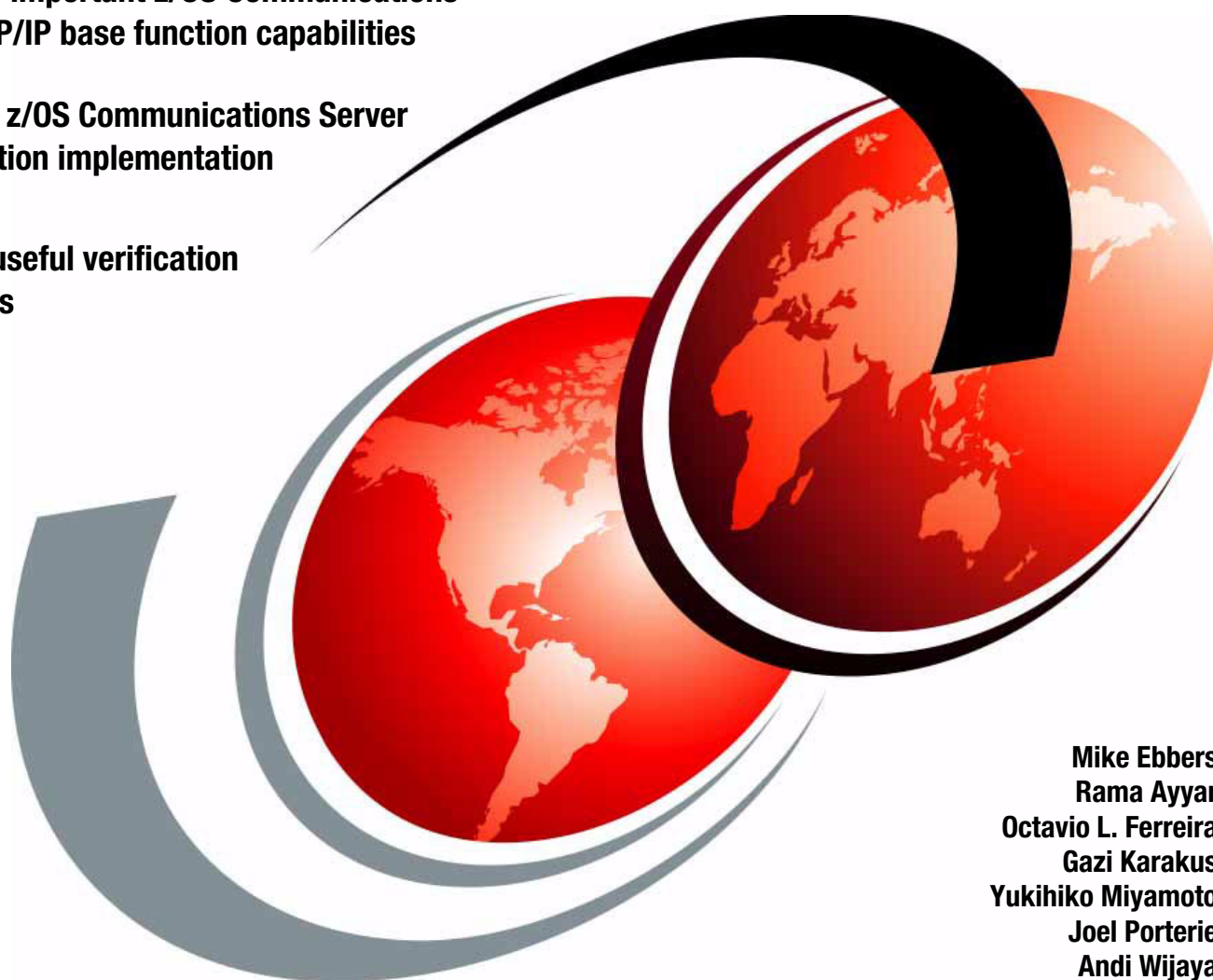


IBM z/OS V1R12 Communications Server TCP/IP Implementation: Volume 1 Base Functions, Connectivity, and Routing

Discusses important z/OS Communications
Server TCP/IP base function capabilities

Describes z/OS Communications Server
base function implementation

Provides useful verification
techniques



Mike Ebbers
Rama Ayyar
Octavio L. Ferreira
Gazi Karakus
Yukihiko Miyamoto
Joel Porterie
Andi Wijaya

Redbooks



International Technical Support Organization

**IBM z/OS V1R12 Communications Server
TCP/IP Implementation: Volume 1 Base
Functions, Connectivity, and Routing**

April 2011

Note: Before using this information and the product it supports, read the information in “Notices” on page ix.

First Edition (April 2011)

This edition applies to Version 1, Release 12 of Communications Server for z/OS (product number 5694-A01).

© Copyright International Business Machines Corporation 2011. All rights reserved.

Note to U.S. Government Users Restricted Rights -- Use, duplication or disclosure restricted by GSA ADP Schedule Contract with IBM Corp.

Contents

Notices	ix
Trademarks	x
 Preface	 xi
The team who wrote this book	xii
Now you can become a published author, too!	xiii
Comments welcome	xiii
Stay connected to IBM Redbooks	xiii
 Chapter 1. Introduction to Communications Server for z/OS IP	 1
1.1 Overview	2
1.1.1 Basic concepts	2
1.2 Featured functions	3
1.3 Communications Server for z/OS IP implementation	4
1.3.1 Functional overview	4
1.3.2 Operating environment	5
1.3.3 Reusable address space ID	6
1.3.4 Protocols and devices	6
1.3.5 Supported routing applications	8
1.3.6 Application programming interfaces	8
1.3.7 z/OS Communications Server applications	10
1.3.8 UNIX System Services	10
1.4 Additional information	17
 Chapter 2. The resolver	 19
2.1 Basic concepts of the resolver	20
2.2 The resolver address space	21
2.2.1 The resolver SETUP data set	22
2.2.2 The resolver configuration file	22
2.2.3 Local hosts file	26
2.2.4 Resolver DNS cache	28
2.2.5 Criteria for indicating an unresponsive DNS name server	33
2.2.6 Unresponsive DNS name servers	33
2.2.7 Affinity servers and generic servers	34
2.2.8 Resolving an IPv6 address	37
2.2.9 Resolver support for EDNS0	39
2.2.10 Considerations	40
2.3 Implementing the resolver	42
2.3.1 Implementation tasks	42
2.3.2 Activation and verification	48
2.4 Problem determination	50
2.5 Additional information	57
 Chapter 3. Base functions	 59
3.1 The base functions	60
3.1.1 Basic concepts	60
3.2 Common design scenarios for base functions	60
3.2.1 Single stack environment	61
3.2.2 Multiple stack environment	62

3.2.3 Recommendation	64
3.2.4 Recommendations for MTU	64
3.3 z/OS UNIX System Services setup for TCP/IP	65
3.3.1 RACF actions for UNIX	65
3.3.2 APF authorization	67
3.3.3 Changes to SYS1.PARMLIB members	68
3.3.4 Changes to SYS1.PROCLIB members	74
3.3.5 Additional z/OS customization for z/OS UNIX	74
3.3.6 TCP/IP server functions	74
3.3.7 TCP/IP client functions	74
3.3.8 UNIX client functions	75
3.3.9 Verification checklist	77
3.4 Configuring z/OS TCP/IP	79
3.4.1 TCP/IP configuration data set names	79
3.4.2 PROFILE.TCPIP	80
3.4.3 VTAM Resource	85
3.4.4 TCPIP.DATA	86
3.4.5 Configuring the local hosts file	86
3.5 Implementing the TCP/IP stack	87
3.5.1 Create TCPIP.DATA file	88
3.5.2 Create the PROFILE.TCPIP file	89
3.5.3 Check BPXPRMxx	91
3.5.4 Create TCP/IP cataloged procedure	91
3.5.5 Add RACF definitions	92
3.5.6 Create a VTAM TRL major node for MPCIPA OSA	92
3.6 Activating the TCP/IP stack	93
3.6.1 UNIX System Services verification	94
3.6.2 Verifying TCP/IP configuration	104
3.7 Reconfiguring the system with z/OS commands	109
3.7.1 Deleting a device and adding or changing a device	110
3.7.2 Modifying a device	110
3.8 Job log versus syslog as diagnosis tool	114
3.9 Message types: Where to find them	114
3.10 Additional information	114
Chapter 4. Connectivity	117
4.1 What is connectivity	118
4.1.1 System z network connectivity	118
4.2 Recommended interfaces	120
4.2.1 High-bandwidth and high-speed networking technologies	120
4.2.2 OSA-Express (MPCIPA)	121
4.2.3 OSA-Express for zEnterprise (z196)	129
4.2.4 HiperSockets (MPCIPA)	130
4.2.5 Dynamic XCF	134
4.3 Connectivity for the z/OS environment	136
4.3.1 IOCP definitions	137
4.3.2 VTAM definitions	139
4.4 OSA-Express QDIO connectivity	139
4.4.1 Dependencies: CHPID, IOCDS, port numbers, portnames, and port sharing	140
4.4.2 Considerations for isolating traffic across a shared OSA port	147
4.4.3 Configuring OSA-Express with VLAN ID	148
4.4.4 Verifying the connectivity status	152
4.5 OSA-Express QDIO connectivity with Connection Isolation	156

4.5.1	Description of Connection Isolation.	158
4.5.2	Dependencies for Connection Isolation	158
4.5.3	Considerations for Connection Isolation	159
4.5.4	Configuring OSA-Express with Connection Isolation	162
4.5.5	Verifying Connection Isolation on OSA2080X.	163
4.5.6	Conclusions and recommendations: best practices for isolating traffic.	180
4.6	HiperSockets connectivity	180
4.6.1	Dependencies	181
4.6.2	Considerations	182
4.6.3	Configuring HiperSockets	182
4.6.4	Verifying the connectivity status	183
4.7	Dynamic XCF connectivity	188
4.7.1	Dependencies	188
4.7.2	Considerations	189
4.7.3	Configuring DYNAMICXCF	189
4.7.4	Verifying connectivity status	190
4.8	Controlling and activating devices.	195
4.8.1	Starting a device	196
4.8.2	Stopping a device	196
4.8.3	Activating modified device definitions	197
4.9	Problem determination	197
4.10	Additional information	203
Chapter 5.	Routing	205
5.1	Basic concepts	206
5.1.1	Terminology	206
5.1.2	Direct routes, indirect routes, and default route	207
5.1.3	Route selection	208
5.1.4	Static routing and dynamic routing	209
5.1.5	Choosing the routing method	211
5.2	Routing in the z/OS environment	212
5.2.1	Static routing	212
5.2.2	Dynamic routing using OMPROUTE.	213
5.2.3	Policy-based routing	217
5.3	Dynamic routing protocols.	217
5.3.1	Open Shortest Path First	217
5.3.2	Routing Information Protocol	222
5.3.3	IPv6 dynamic routing	225
5.4	Implementing static routing in z/OS	227
5.4.1	Dependencies	227
5.4.2	Considerations	228
5.4.3	Implementation tasks	228
5.4.4	Activation and verification	230
5.5	Implementing OSPF routing in z/OS with OMPROUTE	233
5.5.1	Dependencies	234
5.5.2	Considerations	234
5.5.3	Recommendations	235
5.5.4	Implementation tasks	235
5.5.5	Configure routers	243
5.5.6	Activation and verification	244
5.5.7	Managing OMPROUTE.	250
5.6	Problem determination	253
5.6.1	Commands to diagnose networking connectivity problems	254

5.6.2 Diagnosing an OMPROUTE problem	256
5.7 Additional information	264
Chapter 6. VLAN and Virtual MAC support.	265
6.1 Virtual MAC overview	266
6.1.1 Why use virtual MACs.	266
6.1.2 Virtual MAC concept	268
6.1.3 Virtual MAC address assignment	269
6.2 Virtual MAC implementation	269
6.2.1 IP routing when using VMAC	270
6.2.2 Verification	271
6.3 Virtual LAN overview	274
6.3.1 Types of connections	274
6.4 VLAN implementation on z/OS	275
6.4.1 Single VLAN per OSA	275
6.4.2 Multiple VLAN support	276
6.4.3 Multiple VLANs configuration guidelines.	277
6.4.4 Verification	279
6.5 References	281
Chapter 7. Sysplex subplexing	283
7.1 Introduction	284
7.2 Subplex environment	286
7.3 Load Balancing Advisor and subplexing	287
7.4 Subplex implementation	290
7.4.1 XCF group names.	291
7.4.2 TCP/IP structures	292
7.4.3 Subplex 11: Internal subplex.	293
7.4.4 Subplex 22: External subplex	296
7.4.5 Access verifications	297
7.4.6 LBA connected to a subplex	297
7.5 References	298
Chapter 8. Diagnosis.	299
8.1 Debugging a problem in a z/OS TCP/IP environment.	300
8.1.1 An approach to problem analysis	300
8.2 Logs to diagnose CS for z/OS IP problems	302
8.3 Useful commands to diagnose CS for z/OS IP problems	303
8.3.1 The ping command (TSO or z/OS UNIX)	303
8.3.2 traceroute command	306
8.3.3 The netstat command (console, TSO, or z/OS UNIX)	307
8.3.4 NETSTAT Catalog validation	314
8.3.5 Timestamp validation for NETSTAT catalogs	315
8.4 Gathering traces in CS for z/OS IP	315
8.4.1 Taking a component trace	316
8.4.2 Event trace for TCP/IP stacks (SYSTCPIP)	318
8.4.3 Packet trace (SYSTCPDA)	321
8.4.4 OMPROUTE trace (SYSTCPRT)	326
8.4.5 Resolver trace (SYSTCPRE)	329
8.4.6 IKE daemon trace (SYSTCPIK)	330
8.4.7 Intrusion detection services trace (SYSTCPIS)	330
8.4.8 OSAENTA trace (SYSTCPOT)	330
8.4.9 Queued Direct I/O Diagnostic Synchronization.	331
8.4.10 Network security services server trace (SYSTCPNS).	331

8.4.11	Obtaining component trace data with a dump.	332
8.4.12	Analyzing a trace	332
8.4.13	Configuration profile trace	333
8.5	OSA-Express3 Network Traffic Analyzer	333
8.5.1	Determining the microcode level for OSA-Express3.	333
8.5.2	Defining TRLE definitions	334
8.5.3	Checking TCPIP definitions	335
8.5.4	Customizing OSA-Express Network Traffic Analyzer	336
8.5.5	Defining a resource profile in RACF	342
8.5.6	Allocating a VSAM linear data set.	343
8.5.7	Starting the OSAENTA trace	343
8.5.8	Operator command to query and display OSA information.	352
8.5.9	OSM and OSX information	354
8.6	Additional tools for diagnosing CS for z/OS IP problems	355
8.6.1	Network Management Interface API	356
8.6.2	Systems Management Facilities accounting records	357
8.7	MVS console support for selected TCP/IP commands	360
8.7.1	Concept.	360
8.7.2	Commands and environments supported by EZACMD	361
8.7.3	When to use EZACMD	361
8.7.4	How to use the EZACMD command	361
8.7.5	Configuring z/OS for using the EZACMD	362
8.7.6	Using the EZACMD command in the z/OS console	363
8.7.7	Preparing the EZACMD command in z/OS TSO and z/OS NetView	364
8.7.8	Using EZACMD command from z/OS TSO	364
8.7.9	Integrating EZACMD into REXX programs in TSO and NetView	366
8.7.10	Protecting the EZACMD command.	366
8.7.11	Diagnosis: diagnosing the EZACMD command	368
8.8	Additional information	369
Chapter 9.	z/OS in an ensemble.	371
9.1	Basic concepts	372
9.2	Connectivity.	372
9.2.1	Intranode management network (INMN).	372
9.2.2	Intraensemble data network (IEDN)	373
9.3	Enabling z/OS as a member of the ensemble.	373
9.3.1	Enabling z/OS for IPv6	373
9.3.2	Enabling VTAM for the ensemble	375
9.3.3	Validating the ensemble interfaces in z/OS	376
9.3.4	Displaying information about the OSM interfaces.	378
9.4	Defining and activating the z/OS ensemble interfaces	380
9.4.1	Displaying information about the OSX interfaces	382
9.5	References	384
Appendix A.	IPv6 support	385
	Overview of IPv6	386
	Importance of IPv6.	386
	Common design scenarios for IPv6	387
	Tunneling	387
	Dedicated data links	387
	MPLS backbones	388
	Dual-stack backbones.	388
	Dual-mode stack	389

Recommendation	389
How IPv6 is implemented in z/OS Communications Server.	389
IPv6 addressing	389
Stateless address autoconfiguration	390
IPv6 TCP/IP Network part (prefix).	392
IPv6 implementation in z/OS.	395
Verification	405
Appendix B. Additional parameters and functions	411
MVS System symbols	412
Reusable Address Space ID (REUSASID) function examples	415
PROFILE.TCPIP statements	418
IPCONFIG statements	418
GLOBALCONFIG statements.	420
PORT statement	424
IDYNAMICXCF	428
SACONFIG (SNMP subagent)	428
SMFCONFIG	428
Netmonitor	429
INTERFACE statement.	430
DEVICE and LINK statements	436
SRCIP	437
TCP/IP built-in security functions	439
Appendix C. Examples used in our environment.	441
Resolver.	442
TCP/IP stack	444
OMPROUTE dynamic routing	450
Appendix D. Our implementation environment	455
The environment used for all four books	456
Our focus for this book	458
Related publications	459
IBM Redbooks publications	459
Other publications	459
Online resources	460
How to get IBM Redbooks publications	461
Help from IBM	461
Index	463

Notices

This information was developed for products and services offered in the U.S.A.

IBM may not offer the products, services, or features discussed in this document in other countries. Consult your local IBM representative for information on the products and services currently available in your area. Any reference to an IBM product, program, or service is not intended to state or imply that only that IBM product, program, or service may be used. Any functionally equivalent product, program, or service that does not infringe any IBM intellectual property right may be used instead. However, it is the user's responsibility to evaluate and verify the operation of any non-IBM product, program, or service.

IBM may have patents or pending patent applications covering subject matter described in this document. The furnishing of this document does not give you any license to these patents. You can send license inquiries, in writing, to:

IBM Director of Licensing, IBM Corporation, North Castle Drive, Armonk, NY 10504-1785 U.S.A.

The following paragraph does not apply to the United Kingdom or any other country where such provisions are inconsistent with local law: INTERNATIONAL BUSINESS MACHINES CORPORATION PROVIDES THIS PUBLICATION "AS IS" WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESS OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. Some states do not allow disclaimer of express or implied warranties in certain transactions, therefore, this statement may not apply to you.

This information could include technical inaccuracies or typographical errors. Changes are periodically made to the information herein; these changes will be incorporated in new editions of the publication. IBM may make improvements and/or changes in the product(s) and/or the program(s) described in this publication at any time without notice.

Any references in this information to non-IBM Web sites are provided for convenience only and do not in any manner serve as an endorsement of those Web sites. The materials at those Web sites are not part of the materials for this IBM product and use of those Web sites is at your own risk.

IBM may use or distribute any of the information you supply in any way it believes appropriate without incurring any obligation to you.

Information concerning non-IBM products was obtained from the suppliers of those products, their published announcements or other publicly available sources. IBM has not tested those products and cannot confirm the accuracy of performance, compatibility or any other claims related to non-IBM products. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products.

This information contains examples of data and reports used in daily business operations. To illustrate them as completely as possible, the examples include the names of individuals, companies, brands, and products. All of these names are fictitious and any similarity to the names and addresses used by an actual business enterprise is entirely coincidental.


COPYRIGHT LICENSE:

This information contains sample application programs in source language, which illustrate programming techniques on various operating platforms. You may copy, modify, and distribute these sample programs in any form without payment to IBM, for the purposes of developing, using, marketing or distributing application programs conforming to the application programming interface for the operating platform for which the sample programs are written. These examples have not been thoroughly tested under all conditions. IBM, therefore, cannot guarantee or imply reliability, serviceability, or function of these programs.

Trademarks

IBM, the IBM logo, and ibm.com are trademarks or registered trademarks of International Business Machines Corporation in the United States, other countries, or both. These and other IBM trademarked terms are marked on their first occurrence in this information with the appropriate symbol (® or ™), indicating US registered or common law trademarks owned by IBM at the time this information was published. Such trademarks may also be registered or common law trademarks in other countries. A current list of IBM trademarks is available on the Web at <http://www.ibm.com/legal/copytrade.shtml>

The following terms are trademarks of the International Business Machines Corporation in the United States, other countries, or both:

AIX®	NetView®	System z®
BladeCenter®	OMEGAMON®	TDMF®
CICS®	Parallel Sysplex®	Tivoli®
ESCON®	POWER®	VTAM®
FICON®	RACF®	WebSphere®
HiperSockets™	Redbooks®	z/OS®
IBM®	Redbooks (logo)  ®	z/VM®
IMS™	RMF™	z/VSE™
Language Environment®	System p®	z10™
Lotus®	System z10®	z9®
MVS™	System z9®	

The following terms are trademarks of other companies:

Windows, and the Windows logo are trademarks of Microsoft Corporation in the United States, other countries, or both.

UNIX is a registered trademark of The Open Group in the United States and other countries.

Linux is a trademark of Linus Torvalds in the United States, other countries, or both.

Other company, product, or service names may be trademarks or service marks of others.

Preface

For more than 40 years, IBM® mainframes have supported an extraordinary portion of the world's computing work, providing centralized corporate databases and mission-critical enterprise-wide applications. The IBM System z®, the latest generation of the IBM distinguished family of mainframe systems, has come a long way from its IBM System/360 heritage. Likewise, its IBM z/OS® operating system is far superior to its predecessors in providing, among many other capabilities, world class and state-of-the-art support for the TCP/IP Internet protocol suite.

TCP/IP is a large and evolving collection of communication protocols managed by the Internet Engineering Task Force (IETF), an open, volunteer organization. Because of its openness, the TCP/IP protocol suite has become the foundation for the set of technologies that form the basis of the Internet. The convergence of IBM mainframe capabilities with Internet technology, connectivity, and standards (particularly TCP/IP) is dramatically changing the face of information technology and driving requirements for even more secure, scalable, and highly available mainframe TCP/IP implementations.

The *z/OS Communications Server TCP/IP Implementation* series provides understandable, step-by-step guidance about how to enable the most commonly used and important functions of z/OS Communications Server TCP/IP.

In this IBM Redbooks® publication, we provide an introduction to z/OS Communications Server TCP/IP. We then discuss the system resolver, showing the implementation of global and local settings for single and multi-stack environments. Next, we present implementation scenarios for TCP/IP Base functions, Connectivity, Routing, Virtual MAC support, and sysplex subplexing.

For more specific information about z/OS Communications Server standard applications, high availability, and security, refer to the other volumes in the series:

- ▶ *IBM z/OS Communications Server TCP/IP Implementation Volume 2: Standard Applications*, SG24-7897
- ▶ *IBM z/OS V1R11 Communications Server TCP/IP Implementation Volume 3: High Availability, Scalability, and Performance*, SG24-7898
- ▶ *Communications Server for z/OS TCP/IP Implementation Volume 4: Security and Policy-Based Networking*, SG24-7899

For comprehensive descriptions of the individual parameters for setting up and using the functions described in this book, along with step-by-step checklists and supporting examples, refer to the following publications:

- ▶ *z/OS Communications Server: IP Configuration Guide*, SC31-8775
- ▶ *z/OS Communications Server: IP Configuration Reference*, SC31-8776
- ▶ *z/OS Communications Server: IP User's Guide and Commands*, SC31-8780

This book does not duplicate the information in those publications. Instead, it complements them with practical implementation scenarios that can be useful in your environment. To determine at what level a specific function was introduced, refer to *z/OS Communications Server: New Function Summary*, GC31-8771. For complete details, we encourage you to review the documents referred to in the additional resources section at the end of each chapter.

The team who wrote this book

This book was produced by a team of specialists from around the world working at the International Technical Support Organization, Poughkeepsie Center.

Mike Ebberts is a Consulting IT Specialist and Project Leader at the International Technical Support Organization, Poughkeepsie Center. He has worked with IBM mainframe hardware and software products since 1974 in the field, in education, and in the ITSO.

Rama Ayyar is an independent IT Consultant based in Sydney, Australia and a former employee of IBM Australia. He has also worked for CSC Australia in senior technical roles. He is one of the founding members of HCL India. Rama has over 25 years of experience with the MVS™ Operating System. His areas of expertise include TCP/IP, RACF®, DFSMS, z/OS, HCD & Configuration Management, Dump Analysis, and Disaster Recovery. Rama has co-authored seven IBM Redbooks. He holds a Master's Degree in Computer Science from the Indian Institute of Technology, Kanpur. Rama has been in the computer industry for more than 35 years.

Octavio L. Ferreira is a Senior IT Specialist in IBM Brazil. He has 28 years of experience in IBM software support. His areas of expertise include z/OS Communications Server, SNA and TCP/IP, and Communications Server on all platforms. For the last 10 years, Octavio has worked at the Area Program Support Group, providing guidance and support to clients and designing networking solutions such as SNA-TCP/IP Integration, z/OS Connectivity, Enterprise Extender design and implementation, and SNA-to-APPN migration. He has also co-authored other IBM Redbooks publications.

Gazi Karakus is a Network Specialist who has worked for Garanti Technology for four years. He has six years of experience in the networking field. He has a M.Sc. degree in Electronics and Telecommunication Engineering from Istanbul Technical University. His areas of expertise include routing and switching technologies, z/OS Communications Server, SNA and TCP/IP, and Communications Server on other platforms.

Yukihiko Miyamoto is a Senior IT Specialist who has been with IBM Japan for over 14 years. His areas of expertise are WAN, LAN, TCP/IP, and SNA, with a primary focus on router and switch networking technologies. For more than five years, Yukihiko has been working with z/OS Communications Server as a Technical Support member, providing consultation, design, and implementation services for enterprise networking solutions to clients.

Joel Porterie is a Senior IT Specialist who has been with IBM France for over 33 years. He works for Network and Channel Connectivity Services in the PSSC Product Support Group. His areas of expertise include z/OS, TCP/IP, VTAM®, OSA-Express, and Parallel Sysplex®. Joel has taught OSA-Express and FICON® problem determination classes and has provided on site assistance in these areas in numerous countries. He has co-authored many other IBM Redbooks publications.

Andi Wijaya is a Senior Systems Engineer in IBM-JTI Indonesia. His areas of expertise include IT infrastructure management, networking, security, and open source based systems. He is a trainer, consultant, and subject matter expert in Indonesia and also a public quality assurance reviewer for other international books. For more than 10 years, Andi has been working with networking solutions such as fault tolerant infrastructure, high performance enterprise network, and end-to-end integrated security in network infrastructure.

Thanks to the following people for their contributions to this project:

Richard Conway, Robert Haimowitz, Bill White
International Technical Support Organization, Poughkeepsie Center

Roy Brabson
IBM Communications Server Development, Raleigh

Thanks to the authors of the previous editions of this book:

Finally, we want to thank the authors of the previous *z/OS Communications Server TCP/IP Implementation* series for creating the groundwork for this series: Rama Ayyar, Valirio Braga, WenHong Chen, Demerson Cilloti, Sandra Elisa Freitag, Gwen Dente, Gilson Cesar de Oliveira, Mike Ebbers, Octavio Ferreira, Marco Giudici, Adi Horowitz, Michael Jensen, Shizuka Katoh, Sherwin Lake, Bob Loudon, Garth Madella, Yukihiro Miyamoto, Hajime Nagao, Shuo Ni, Carlos Bento Nonato, Yohko Ojima, Roland Peschke, Joel Porterie, Marc Price, Frederick James Rathweg, Micky Reichenberg, Larry Templeton, Rudi van Niekerk, Bill White, Thomas Wienert, Andi Wijaya, and Maulide Xavier.

Now you can become a published author, too!

Here's an opportunity to spotlight your skills, grow your career, and become a published author—all at the same time! Join an ITSO residency project and help write a book in your area of expertise, while honing your experience using leading-edge technologies. Your efforts will help to increase product acceptance and customer satisfaction, as you expand your network of technical contacts and relationships. Residencies run from two to six weeks in length, and you can participate either in person or as a remote resident working from your home base.

Find out more about the residency program, browse the residency index, and apply online at:

ibm.com/redbooks/residencies.html

Comments welcome

Your comments are important to us!

We want our books to be as helpful as possible. Send us your comments about this book or other IBM Redbooks publications in one of the following ways:

- Use the online **Contact us** review Redbooks form found at:

ibm.com/redbooks

- Send your comments in an email to:

redbooks@us.ibm.com

- Mail your comments to:

IBM Corporation, International Technical Support Organization
Dept. HYTD Mail Station P099
2455 South Road
Poughkeepsie, NY 12601-5400

Stay connected to IBM Redbooks

- Find us on Facebook:

<http://www.facebook.com/IBMRedbooks>

- ▶ Follow us on Twitter:
<http://twitter.com/ibmredbooks>
- ▶ Look for us on LinkedIn:
<http://www.linkedin.com/groups?home=&gid=2130806>
- ▶ Explore new Redbooks publications, residencies, and workshops with the IBM Redbooks weekly newsletter:
<https://www.redbooks.ibm.com/Redbooks.nsf/subscribe?OpenForm>
- ▶ Stay current on recent Redbooks publications with RSS Feeds:
<http://www.redbooks.ibm.com/rss.html>



Introduction to Communications Server for z/OS IP

z/OS Communications Server is the IBM implementation of the standard TCP/IP protocol suite on the z/OS platform. TCP/IP is a component product of the z/OS Communications Server, and it provides a multitude of technologies. Collectively, those technologies provide an Open Systems environment for the development, establishment, and maintenance of applications and systems.

The z/OS Communications Server product includes ACF/VTAM, in addition to TCP/IP.

This chapter presents a basic overview of z/OS Communications Server IP as it is implemented in the z/OS environment. You can find a more complete and comprehensive explanation of z/OS Communications Server IP from the publications that are listed in 1.4, “Additional information” on page 17.

This chapter discusses the following topics.

Section	Topic
1.1, “Overview” on page 2	Basic concepts of Communications Server for z/OS IP
1.2, “Featured functions” on page 3	Key characteristics of Communications Server for z/OS IP and
1.3, “Communications Server for z/OS IP implementation” on page 4	Functional overview of how Communications Server for z/OS IP is implemented
1.4, “Additional information” on page 17	Lists IBM publications that provide further details for implementing Communications Server for z/OS IP

1.1 Overview

z/OS Communications Server provides the industry-standard TCP/IP protocol suite, allowing z/OS environments to share data and computing resources with other TCP/IP computing environments, when authorized. Communications Server for z/OS IP enables anyone in a non-z/OS TCP/IP environment to access resources in the z/OS environment and perform tasks and functions provided by the TCP/IP protocol suite.

It provides the computer platform with the freedom desired by organizations to distribute workload to environments best suited to their needs. Communications Server for z/OS IP, therefore, adds the z/OS environment to the list of environments in which an organization can share data and computer processing resources in a TCP/IP network.

Communications Server for z/OS IP supports two environments:

- ▶ It provides a native MVS (z/OS) environment in which users can exploit the TCP/IP protocols in the z/OS applications environment, including batch jobs, started tasks, TSO, CICS® applications, and IMS™ applications.
- ▶ It also provides native TCP/IP support in the UNIX® Systems Services environment in which users can create and use applications that conform to the POSIX or XPG4 standard (a UNIX specification). The UNIX environment and services can also be exploited from the z/OS environment, and vice versa.

1.1.1 Basic concepts

The TCP/IP address space is where the TCP/IP protocol suite is implemented for CS for z/OS IP. The TCP/IP address space is commonly referred to as a *stack*.

Communications Server for z/OS IP has highly efficient direct communication between the UNIX System Services address space (OMVS) and a TCP/IP stack that was integrated in UNIX System Services. This communication path includes the UNIX System Services Physical File System (PFS) component for AF_INET and AF_INET6 (Addressing Family-Internet) sockets communication.

The z/OS Communications Server has the following features:

- ▶ A process model that provides a full multiprocessing capability. It includes full duplex data paths of reduced lengths.
- ▶ An I/O process model that allows VTAM to provide the I/O device drivers. MultiPath Channel (MPC) Data Link Control (DLC) is shared between VTAM and TCP/IP. It executes multiple dispatchable units of work and is tightly integrated with the Common Storage Manager support.
- ▶ A storage management model handles expansion and contraction of storage resources, as well as requests of varying sizes and types of buffers. Common Storage Manager (CSM) manages communication between the Sockets PFS through the transport provider and network protocols to the network interface layer of Communications Server for z/OS IP stack. The data that is placed in the buffers can be accessed by any function all the way down to the protocol stack.

Communications Server for z/OS IP runs as a single stack that serves both the traditional MVS (z/OS) environment and the z/OS UNIX (UNIX System Services) environment, and IP offers two variants of the UNIX shell environment:

- ▶ The OMVS shell, which is much like a native UNIX environment
- ▶ The ISHELL, which is an ISPF interface with access to menu-driven command interfaces

The TCP/IP protocol suite is implemented by an MVS started task within the TCP/IP address space in conjunction with z/OS UNIX (UNIX System Services).

A Communications Server for z/OS IP environment requires a Data Facility Storage Management Subsystem (DFSMS), a z/OS UNIX file system, and a security product such as Resource Access Control Facility (RACF). These resources must be defined and functional before the z/OS Communications Server can be started successfully and establish the TCP/IP environment. We later mention the manner in which these products impact this environment.

1.2 Featured functions

z/OS Communications Server provides a high-performance, highly secure, scalable, and reliable platform on which to build and deploy networking applications.

Communications Server for z/OS IP offers an environment that is accessible to the enterprise IP network and the Internet if so desired. It defines the z/OS environment as a viable platform by making z/OS applications and systems available to the non-z/OS environment, which are typically UNIX/Windows®-centric. Consequently, it eliminates the issues and challenges of many large corporations to migrate or integrate with a more accessible platform and newer technologies.

The following list includes many of the technologies that have been implemented in the z/OS environment to complement TCP/IP.

- ▶ High-speed connectivity, such as:
 - OSA-Express3 10 Gigabit Ethernet in QDIO mode
 - High-speed communication between TCP/IP stacks running in logical partitions using HyperSockets™ in IQDIO mode.
- ▶ High availability for applications using Parallel Sysplex technology in conjunction with:
 - Dynamic Virtual IP Address (VIPA), which provides TCP/IP application availability across z/OS systems in a sysplex and allows participating TCP/IP stacks to provide backup and recovery for each other, for planned and unplanned TCP/IP outages
 - Sysplex Distributor, which provides intelligent load balancing for TCP/IP application servers in a sysplex, and along with Dynamic VIPA provides a single system image for client applications connecting to those servers
 - The Load Balancing Advisor (LBA), which provides z/OS Sysplex server application availability and performance data to outboard load balancers through the Server Application State Protocol (SASP)
- ▶ Enterprise connectivity support is offered through many features, such as:
 - TN3270 Server, which provides workstation connectivity over TCP/IP networks to access z/OS and enterprise SNA applications.
 - Enterprise Extender, which allows SNA Enterprise applications to communicate reliably over an IP network, using SNA HPR over UDP transport layer protocol.

- IPv4 and IPv6 networking functions are provided by the TCP/IP stack operating in a standard dual-mode setup where IPv4 and IPv6 connectivity and applications are supported concurrently by a single TCP/IP stack instance.
- Sockets programming interface support for traditional z/OS workloads provide IP connectivity to applications written in REXX, COBOL, PL/I. Sockets programming interfaces are supported in various environments, such as TSO, batch, CICS, and IMS.
- ▶ Network Security protects sensitive data and the operation of the TCP/IP stack on z/OS by using:
 - IPsec/VPN functions that enable the secure transfer of data over a network using standards for encryption, authentication, and data integrity.
 - Intrusion Detection Services (IDS), which evaluates the stack for attacks that would undermine the integrity of its operation. Events to examine and actions to take (such as logging) at event occurrence are defined by the IDS policy.
 - Transport Layer Security (TLS) enablement ensures data is protected as it flows across the network.
 - Kerberos and GSSAPI support is provided for selected applications.
 - Defensive filtering provides an infrastructure to add, delete and modify short-term TCP/IP filters in real time to counter specific attacks.
 - Network Security Services provides a centralized security infrastructure to extend System z security to NSS clients, such as IKE daemons and XML appliances.
- ▶ Network Management support collects network topology, status, and performance information and makes it available to network management tools, including:
 - Local management applications that can access management data using a specialized high-performing network management programming interface that is known as NMI.
 - Support of remote management applications through the SNMP protocol. Communications Server z/OS supports the latest SNMP standard, SNMPv3. Communications Server z/OS also supports standard TCP/IP-based Management Information Base (MIB) data.
 - Additional MIB support is also provided by Enterprise-specific MIB, which supports management data for Communications Server TCP/IP stack-specific functions.

1.3 Communications Server for z/OS IP implementation

Communications Server for z/OS IP provides TCP/IP support for the native MVS and UNIX System Services environment. It is implemented within a z/OS address space and runs within the native MVS environment, and consequently it has RACF, DFSMS, and z/OS UNIX file system dependencies.

1.3.1 Functional overview

CS for z/OS IP takes advantage of Communications Storage Manager (CSM) and of VTAM's Multipath Channel (MPC) and Queued Direct I/O (QDIO) capabilities in its TCP/IP protocol implementation. This tight coupling with VTAM provides enhanced performance and serviceability.

As illustrated in Figure 1-1, many data link control (DLC) protocols are provided with the z/OS Communications Server by the VTAM component.

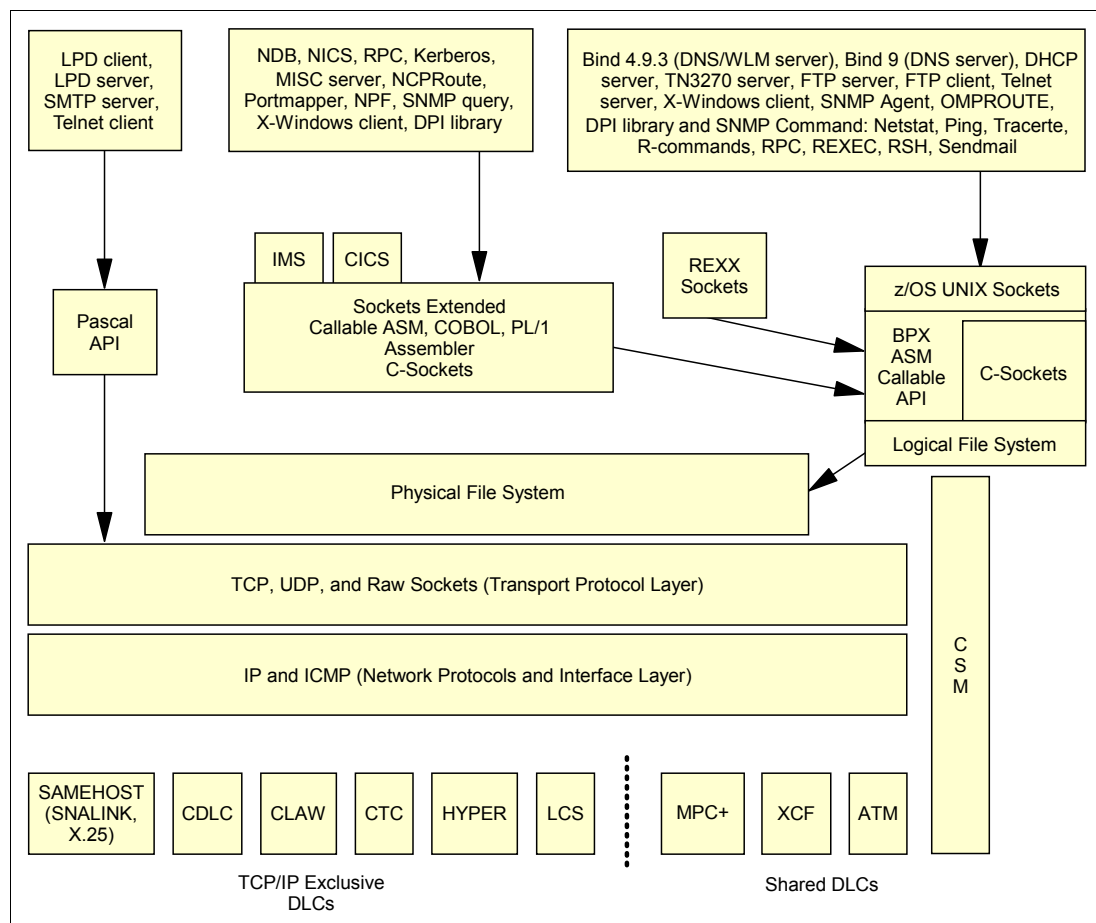


Figure 1-1 Functional overview

With CS for z/OS IP, two worlds converge, providing access to the z/OS UNIX environment and the traditional MVS environment.

1.3.2 Operating environment

Because the z/OS UNIX environment is supported in the MVS environment, there is no need to discuss the creation of an MVS environment here. However, there are customization requirements on the UNIX System Services side of the environment that are needed in order to start Communications Server for z/OS IP successfully. This dependence on UNIX, of course, implies that z/OS UNIX administrators must also be familiar with both traditional MVS commands and interfaces.

I/O flow process

Another feature of the operating environment is the storage and I/O designs. The operating environment design features a tightly integrated storage and I/O model, known as Common Storage Manager (CSM). The CSM facility is used by authorized programs to manage subsystem storage pools. It provides a flat storage model that is accessible by multiple layers of the process model, as Figure 1-1 illustrates. It is also accessible across z/OS address space boundaries, thereby reducing the data moves between processes and tasks that exchange data as they perform work. VTAM and TCP/IP tasks are typical examples.

The CSM facility also manages storage as it automates the addition and subtraction of the different types and sizes of storage requests.

1.3.3 Reusable address space ID

The z/OS system assigns an address space identifier (ASID) to an address space when the address space is created. A limited number of ASIDs are available for the system to assign. When all ASIDs are assigned to existing address spaces, the system cannot start a new address space. This condition can cause lost ASIDs in the system, which are address spaces that have terminated but which the system does not reuse because of the address space's residual cross memory connections.

ASIDs used for the TCP/IP stack, the resolver, VTAM, and TN3270 are non-reusable, because they provide PC-entered services that must be accessible to other address spaces. If these address spaces are terminated enough times, all available ASIDs can be exhausted, preventing the creation of a new address space on the system. That situation might require an IPL.

With the reusable address space support, ASIDs that would otherwise be unusable after termination of the started task are made available for reuse. The reusable ASID function is available for the TCP/IP, resolver, VTAM, and TN3270 started tasks; however, it is not available to TNF, VMCF, and all the products and applications that use their services (such as TSO TELNET).

To enable the reuse ASID function, you must:

- ▶ Specify REUSASID=YES in member DIAGxx of your PARMLIB
- ▶ Specify REUSASID=YES on the **start** command when starting the address space

The REUSASID parameter cannot be coded in the JCL of the started task because the Master Scheduler needs to know this information *before* the JCL is read and the ASID is assigned.

The resolver started task will always use a reusable ASID when started during z/OS UNIX initialization through the BPXRMMxx statement RESOLVER_PROC but will use a non-reusable ASID if stopped and started. You should, therefore, restart resolver with the REUSASID=YES parameter specified on the start command.

The REUSASID parameter is to be used only by address spaces such as TCP/IP, resolver, and TN3270 that are usually non-reusable when terminated, because unnecessary use of REUSASID=YES can reduce the number of ASIDs that are available for satisfying ordinary address space requests.

We include examples of REUSASID coding and its results in Appendix B, "Additional parameters and functions" on page 411.

1.3.4 Protocols and devices

As illustrated in Figure 1-1 on page 5, the DLC is a protocol layer that manages and provides communication between the file I/O subsystem and the I/O device driver of the particular device. The figure also shows two categories of DLCs:

- ▶ TCP/IP exclusive DLCs
- ▶ Shared DLCs

TCP/IP exclusive DLCs

TCP/IP exclusive DLCs are those *only* available for TCP/IP usage and are not shared with ACF/VTAM. Here are examples of TCP/IP exclusive DLCs that are supported by Communications Server for z/OS IP:

- ▶ LAN Channel Station (LCS), which is a protocol used by OSA 1000BASE-T feature, some routers, and the 3746-9x0 MAE.
- ▶ SAMEHOST, which is another TCP/IP exclusive DLC protocol that exists, although it does not make use of System z channels. In the past, this communication was provided by IUCV. Currently, SNALINK LU0, SNALINK LU6.2, and X.25 exploit the SAMEHOST DLC.

Shared DLCs

Shared DLCs are those that can be *simultaneously used* by TCP/IP and ACF/VTAM. Figure 1-1 on page 5 indicates the shared DLCs. The most commonly used DLCs include those that we describe here.

Multipath Channel+

Multipath Channel+ (MPC+) is an enhanced version of VTAM's MPC protocol. The MPC I/O process defines the implementation of the MPC protocols and allows for the efficient use of multiple read and write channels.

MPC handles protocol headers and data separately and executes multiple I/O dispatchable units of work. This process, when used in conjunction with Communication Storage Management, creates efficient I/O throughput. High Performance Data Transfer uses MPC+ together with CSM to decrease the number of data copies that are required to transmit data.

This type of connection can be used in two ways:

- ▶ MPCPTP allows a CS for z/OS IP environment to connect to a peer IP stack in a point-to-point configuration. With MPCPTP, a CS for z/OS IP stack can be connected to:
 - Another CS for z/OS IP stack
 - An IP router with corresponding support
 - A non-z/OS server
 - 3746-9x0 MAE

PTP Samehost (MPCPTP), sometimes referred to IUTSAMEH: This connection type is used to connect two or more CS for z/OS IP stacks running on the same z/OS LPAR. In addition, it can be used to connect these CS for z/OS IP stacks to z/OS VTAM for the use of Enterprise Extender.

- ▶ MPCIPA allows an Open Systems Adapter-Express (OSA-Express) port to act as an extension of the z/OS Communications Server TCP/IP stack and not as a peer TCP/IP stack, as with MPCPTP.
 - OSA-Express provides a mechanism for communication called Queued Direct I/O (QDIO). Although it uses the MPC+ protocol for its control signals, the QDIO interface is quite different from channel protocols. It uses Direct Memory Access (DMA) to avoid the overhead associated with channel programs. A partnership between CS for z/OS IP and the OSA-Express adapter provides compute-intensive functions from the System z server to the adapter.

OSA-Express collaborates with z/OS Communications Server TCP/IP to support up to 10 Gigabit Ethernet, 1000BASE-T, Fast Ethernet, High Speed Token Ring (HSRP), and ATM LAN emulation. TCP/IP hosts supports OSA-Express, OSA-Express2, and OSA-Express3 features.

- HiperSockets (Internal Queued Direct I/O, iQDIO) provides high-speed, low-latency IP message passing between logical partitions (LPARs) within a single System z server. The communication is through processor system memory through Direct Memory Access (DMA). The virtual servers that are connected through HiperSockets form a virtual LAN. HiperSockets uses internal QDIO at memory speeds to pass traffic between virtual servers.

Cross-System Coupling Facility

Cross-System Coupling Facility (XCF) allows communication between multiple CS for z/OS IP stacks within a Parallel Sysplex. The XCF DLC can be defined as with traditional DLCs, but it also supports XCF Dynamics, in which the XCF links are established automatically.

If DYNAMICXCF is coded, z/OS images within the same server will use the HiperSockets DYNAMICXCF connectivity instead of the standard XCF connectivity for data transfer.

For more information about devices and connectivity options, refer to Chapter 4, “Connectivity” on page 117.

1.3.5 Supported routing applications

z/OS Communications Server ships only one routing application, called OMPROUTE. OMPROUTE implements the Open Shortest Path First protocols (OSPF and OSPFv4) and Routing Information Protocols (RIPv1, RIPv2, RIPv3). It enables the Communications Server for z/OS IP to function as an OSPF/RIP-capable router in a TCP/IP network. Either (or both) of these two routing protocols can be used to dynamically maintain the host routing table.

Additionally, Communications Server for z/OS IP provides an OMPROUTE subagent that implements the OSPF MIB variable containing OSPF protocol and state information for SNMP. This MIB variable is defined in RFC 1850. Refer to Chapter 5, “Routing” on page 205, for a detailed discussion about OMPROUTE and its function.

1.3.6 Application programming interfaces

As Figure 1-1 on page 5 illustrates, all of the APIs provided by Communications Server for z/OS IP, with the exception of the Pascal API, interface with the Logical File System (LFS) layer. The APIs are divided into the following categories:

- ▶ Pascal
- ▶ TCP/IP socket APIs
- ▶ z/OS UNIX APIs
- ▶ REXX sockets

We describe these items in more detail in the following sections.

Pascal API

The Pascal application programming interface enables you to develop TCP/IP applications in Pascal language. Supported environments are normal MVS address spaces. Unlike the other APIs, the Pascal API does not interface directly with the LFS. It uses an internal interface to communicate with the TCP/IP protocol stack. The Pascal API only supports AF_INET.

TCP/IP socket APIs

The z/OS Communications Server provides several APIs to access TCP/IP sockets. These APIs can be used in either or both integrated and common INET PFS configurations.

In a common INET PFS configuration, however, they function differently from z/OS UNIX APIs. In this type of configuration, the z/OS Communications Server APIs always bind to a single PFS transport provider, and the transport provider must be the TCP/IP stack provided by the z/OS Communications Server.

The following TCP/IP socket APIs are included in the z/OS Communications Server:

- ▶ The CICS socket interface enables you to write CICS applications that act as clients or servers in a TCP/IP-based network. CICS sockets only support AF_INET.
- ▶ The C sockets interface supports socket function calls that can be invoked from C programs. However, note that for C application development, IBM recommends the use of the UNIX C sockets interface. These programs can be ported between MVS and most UNIX environments relatively easily if the program does not use any other MVS-specific services. C sockets only support AF_INET.
- ▶ The Information Management System (IMS) IPv4 socket interface supports client/server applications in which one part of the application executes on a TCP/IP-connected host and the other part executes as an IMS application program. The IMS sockets API supports AF_INET.
- ▶ The Sockets Extended macro API is a generalized assembler macro-based interface to sockets programming. The Sockets Extended macro API supports AF_INET and AF_INET6.
- ▶ The Sockets Extended Call Instruction API is a generalized call-based, high-level language interface to sockets programming. The Sockets Extended Call Instruction API supports AF_INET and AF_INET6.

z/OS UNIX APIs

The following APIs are provided by the z/OS UNIX element of z/OS and are supported by the TCP/IP stack in the z/OS Communications Server:

- ▶ z/OS UNIX C sockets is used in the z/OS UNIX environment. It is the z/OS UNIX version of the native MVS C sockets programming interface. Programmers use this API to create applications that conform to the POSIX or XPG4 standard (a UNIX specification). The z/OS UNIX C sockets support AF_INET and AF_INET6.
- ▶ z/OS UNIX assembler callable services is a generalized call-based, high-level language interface to z/OS UNIX sockets programming. The z/OS UNIX assembler callable services support AF_INET and AF_INET6.

Refer to *z/OS XL C/C++ Compiler and Run-Time Migration Guide for the Application Programmer*, GC09-4913, for complete documentation of the z/OS UNIX C sockets APIs. You can also find further guidance in *z/OS UNIX System Services Programming Tools*, SA22-7805.

REXX sockets

The REXX sockets programming interface implements facilities for socket communication directly from REXX programs by using an address rxsocket function. REXX socket programs can execute in TSO, online, or batch. The REXX sockets programming interface supports AF_INET and AF_INET6.

Refer to *z/OS Communications Server: IP Sockets Application Programming Interface Guide and Reference*, SC31-8788, for complete documentation of the TCP/IP Services APIs.

1.3.7 z/OS Communications Server applications

z/OS Communications Server TCP/IP provides a number of standard client and server applications, including:

- ▶ SNA 3270 Logon Services (TN3270)
- ▶ z/OS UNIX logging services (syslogd)
- ▶ File Transfer Services (FTP)
- ▶ Network Management Services (SNMP Agents, Subagents, Trap forwarding)
- ▶ IP Printing (LPR, LPD, Infoprint Server)
- ▶ Internet Daemon Listener (INETD)
- ▶ Mail Services (SMTP and sendmail)
- ▶ Client-based mail forwarding Simple Network Management Protocol (CSSMTP)
- ▶ z/OS UNIX logon services (otelnstd)
- ▶ Remote Execution (REXEC, RSHD, REXEC, RSH, orexecd, orshd, orexec, and orsh)
- ▶ Domain Name Services (Caching DNS BIND9 server)

These applications are discussed in detail in *IBM z/OS Communications Server TCP/IP Implementation Volume 2: Standard Applications*, SG24-7897, and *z/OS Communications Server: IP Configuration Guide*, SC31-8775.

z/OS Communications Server also provides a number of specialized services, including:

- ▶ Policy Agent for implementing networking and security policies in a z/OS environment
- ▶ Centralized or Distributed Policy Services
- ▶ Network Security Services (NSS)
- ▶ Defense Manager
- ▶ Integrated Services policies using Resource ReSrvation Protocol (RSVP)
- ▶ Differentiated Services using Quality of Service (QoS) policies.

These applications are discussed in detail in the following publications:

- ▶ *Communications Server for z/OS TCP/IP Implementation Volume 4: Security and Policy-Based Networking*, SG24-7899.
- ▶ *z/OS Communications Server: IP Configuration Guide*, SC31-8775

1.3.8 UNIX System Services

UNIX System Services is the z/OS Communications Server implementation of UNIX as defined by X/Open in XPG 4.2. UNIX System Services coexists with traditional MVS functions and traditional MVS file types (partitioned data sets, sequential files, and so on). It concurrently allows access to z/OS UNIX file system files and to UNIX utilities and commands by means of application programming interfaces and the interactive shell environment.

Communications Server for z/OS IP offers two variants of the UNIX shell environment:

- ▶ The z/OS shell, which is the default shell
- ▶ The tcsh shell (Ishell), which is an enhanced version of the Berkeley UNIX C shell

The Communications Server for z/OS IP requires that UNIX System Services be customized in full-function mode before the TCP/IP stack will successfully initialize. For this reason we present an overview of UNIX System Services to provide an overview of the coding and security considerations that are involved with UNIX System Services.

Customization levels of UNIX System Services

There are two levels of z/OS UNIX services:

- ▶ *Minimum mode*, indicating that although OMVS initializes, it provides few z/OS UNIX services, and there is no support for TCP/IP and the z/OS shell. In this mode there is no need for DFSMS or for a security product such as RACF.
- ▶ *Full-function mode*, indicating that the complete array of z/OS UNIX services is available. In this mode DFSMS, RACF, and the z/OS UNIX file system are required. TCP/IP and z/OS UNIX file system interaction with UNIX System Services is defined within the BPXPRMxx member of SYS1.PARMLIB.

See *z/OS UNIX System Services Planning*, SA22-7800 for a useful description of the UNIX System Services customization process and TCP/IP.

UNIX System Services concepts

z/OS UNIX enables two open systems interfaces on the z/OS operating system:

- ▶ An application program interface (API)
- ▶ An interactive shell interface

With the APIs, programs can run in any environment (including batch jobs, in jobs submitted by TSO/E interactive users, and in most other started tasks) or in any other MVS application task environment. The programs can request:

- ▶ Only MVS services
- ▶ Only z/OS UNIX services
- ▶ Both MVS and z/OS UNIX services

The shell interface is an execution environment similar to TSO/E, with a programming language of shell commands like those in the Restructured Extended Executor (REXX) language. The shell work consists of:

- ▶ Programs that are run interactively by shell users
- ▶ Shell commands and scripts that are run interactively by shell users
- ▶ Shell commands and scripts that are run as batch jobs

In z/OS UNIX Systems Services, address spaces are provided by the `fork()` or `spawn()` functions of the Open Edition callable services.

- ▶ For a `fork()` function, the system copies one process, called the *parent* process, into a new process, called the *child* process, and places the child process in a new address space, the forked address space.
- ▶ A `spawn()` function also starts a new process in a new address space. Unlike a `fork()`, in a `spawn()` call the parent process specifies a name of a program to be run in the child process.

The types of processes can be:

- ▶ User processes, which are associated with a user
- ▶ Daemon processes, which perform continuous or periodic system-wide functions, such as a Web server

Daemons (a UNIX concept) are programs that are typically started when the operating system is initialized and remain active to perform standard services. Some programs are considered daemons that initialize processes for users even though these daemons are not long-running processes. Examples of daemons provided by z/OS UNIX are *cron*, which starts applications at specific times, and *inetd*, which provides service management for a network.

A process can have one or more threads. A *thread* is a single flow of control within a process. Application programmers create multiple threads to structure an application in independent sections that can run in parallel for more efficient use of system resources.

UNIX Hierarchical File System

Data sets and files are comparable terms. If you are familiar with MVS, you probably use the term *data set* to describe a unit of data storage. If you are familiar with AIX® or UNIX, you probably use the term *file* to describe a named set of records stored or processed as a unit. In the UNIX System Services environment, the files are arranged in a z/OS UNIX file system.

The Hierarchical File System allows you to set up a file hierarchy that consists of:

- ▶ Directories, which contain files, other directories, or both. Directories are arranged hierarchically, in a structure that resembles an upside-down tree, with the root directory at the top and branches at the bottom.
- ▶ z/OS UNIX file system files, which contain data or programs. A file containing a load module, shell script, or REXX program is called an *executable file*. Files are kept in directories.
- ▶ Additional local or remote file systems, which are mounted on directories of the root file system or of additional file systems.

To the z/OS system, the UNIX file hierarchy appears as a collection of System z File System data sets. Each z/OS UNIX file system data set is a mountable file system. The root file system is the *first* file system mounted. Subsequent file systems can be mounted logically on a directory within the root file system or on a directory within any mounted file system.

Each mountable file system resides in a z/OS UNIX file system data set on direct access storage. DFSMS/MVS manages the z/OS UNIX file system data sets and the physical files.

For more information about the z/OS UNIX file system, refer to *z/OS CS: IP Migration*, GC31-8773, and *z/OS UNIX System Services Planning*, SA22-7800.

z/OS UNIX file system definitions in BPXPRMxx

To get UNIX System Services active in full-function mode, you need to define the root file system in the BPXPRMxx member of SYS1.PARMLIB. The root file system is usually loaded or copied at z/OS installation time. The BPXPRMxx definition is detailed in *z/OS UNIX System Services Planning*, SA22-7800.

An important part of your z/OS UNIX file system is located in the /etc directory. The /etc directory contains some basic configuration files of UNIX System Services, and most applications keep their configuration files in there as well. To avoid losing all of your configuration when you upgrade your operating system, it is recommended that you put the /etc directory in a separate z/OS UNIX file system data set and mount it at the /etc mountpoint. Refer to *z/OS UNIX System Services Planning*, SA22-7800, for more information about the /etc directory.

z/OS UNIX user identification

All users of an MVS system, including users of z/OS UNIX functions, must have a valid MVS user ID and password. To use standard MVS functions, the user must have the standard MVS identity based on the RACF user ID and group name.

If a unit of work in MVS uses z/OS UNIX functions, this unit of work must have, in addition to a valid MVS identity, a z/OS UNIX identity. A z/OS UNIX identity is based on a UNIX user ID (UID) and a UNIX group ID (GID). Both UID and GID are numeric values ranging from 0 to 2147483647 ($2^{31}-1$).

In a z/OS UNIX system, the UID is defined in the OMVS segment in the user's RACF user profile, and the GID is defined in an OMVS segment in the group's RACF group profile. What we in an MVS environment call the user ID is in a UNIX environment normally termed the user name or the login name. It is the name that users use to present themselves to the operating system. In both a z/OS UNIX system and other UNIX systems, this user name is correlated to a numeric user identification, the UID, which is used to represent this user wherever such information has to be stored in the z/OS UNIX environment. One example of this is in the Hierarchical File System, where the UID of the owning user is stored in the file security portion of each individual file.

Access to resources in the traditional MVS environment is based on the MVS user ID, group ID, and individual resource profiles that are stored in the RACF database.

Access to z/OS UNIX resources is granted *only* if the MVS user ID has a valid OMVS segment with an OMVS UID, or if a default user is configured as explained next. Access to resources in the Hierarchical File System is based on the UID, the GID, and file access permission bits that are stored with each file. The permission bits are three groups of three bits each. The groups describe:

- ▶ The owner of the file itself
- ▶ The users with the same GID as the owner
- ▶ The rest of the world

The three bits are:

- ▶ Read access
- ▶ Write access
- ▶ Search access if it is a directory or if it is a file that is executable

The superuser UID has a special meaning in all UNIX environments, including the z/OS UNIX environment. This user has a UID of zero and can access every resource.

In lieu of or in addition to RACF definitions for individual users, you can define a *default user*. The default user will be used to allow users without an OMVS segment defined to access UNIX System Services. The default user concept should be used with caution, because it could become a security exposure.

You will also find more information about the RACF security aspects of implementing the Communications Server for z/OS IP in *Communications Server for z/OS TCP/IP Implementation Volume 4: Security and Policy-Based Networking*, SG24-7899.

Accessing the z/OS UNIX shells

You can access z/OS UNIX shells in the following ways:

- ▶ The TSO/E OMVS command provides a 3270 interface in the z/OS UNIX shell.
- ▶ The TSO/E ISHELL or ISH command provides a 3270 interface that uses ISPF dialogs.
- ▶ The **rlogin** command provides an ASCII interface.
- ▶ The **Telnet** command provides an ASCII interface. This Telnet is into the UNIX Telnet daemon and not the TN3270 server in the z/OS system space.
- ▶ From a TCP/IP network, the TN3270 command can be used, which provides a full-screen 3270 interface for executing the OMVS or ISHELL commands.

There are two shells, the z/OS shell and the lshell. The login shell is determined by the PROGRAM parameter in the RACF OMVS segment for each user. The default is the z/OS shell.

You can find further information about the z/OS UNIX shells in *z/OS UNIX System Services User's Guide*, SA22-7801.

Operating mode

When a user first logs on to the z/OS UNIX shell, the user is operating in line mode. Depending on the method of accessing the shell, the user can then use utilities that require raw mode (such as vi) or run an X Window System application.

The different workstation operating modes are:

- ▶ Line mode
Input is processed after you press Enter. This is also called *canonical mode*.
- ▶ Raw mode
Each character is processed as it is typed. This is also called *non-canonical mode*.
- ▶ Graphical mode
This is a graphical user interface for X Window System applications.

UNIX System Services communication

A socket is the endpoint of a communication path; it identifies the address of a specific process at a specific computer using a specific transport protocol. The exact syntax of a socket address depends on the protocol being used, that is, on its *addressing family*.

When you obtain a socket using the `socket()` system call, you pass a parameter that tells the socket library to which addressing family the socket should belong. All socket addresses within one addressing family use the same syntax to identify sockets.

Socket addressing families in UNIX System Services

In a z/OS UNIX environment, the most widely used addressing families are AF_INET and AF_UNIX. There is IPv6 support (AF_INET6 addressing family) in Communications Server for z/OS IP in a single transport driver environment configured in Dual-mode. Socket applications written to the IPv6 APIs can use the z/OS TCP/IP stack for IPv6 network connectivity.

Note: Throughout this discussion, information regarding AF_INET (IPv4) also applies to AF_INET6 (IPv6).

The z/OS UNIX Systems Services implements support for a given addressing family through different physical file systems. There is one physical file system for the AF_INET addressing family, and there is another for the AF_UNIX addressing family. A PFS is the part of the z/OS UNIX operating system that handles the storage of data and its manipulation on a storage medium.

AF_UNIX addressing family

The UNIX addressing family is also referred to as the UNIX domain. If two socket applications on the same MVS image want to communicate with each other, they can open a socket as an AF_UNIX family socket. In that case, the z/OS UNIX address space will handle the full communication between the two applications (see Figure 1-2). That is, the AF_UNIX physical file system is self-contained within z/OS UNIX and does not rely on other products to implement the required functions.

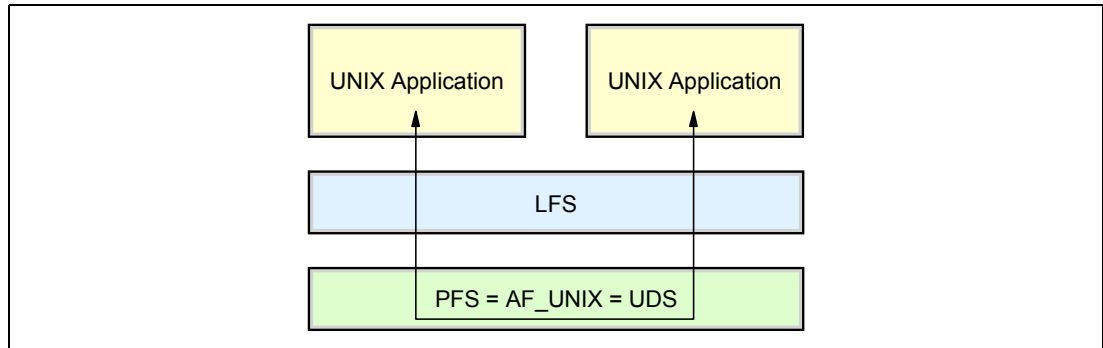


Figure 1-2 AF_UNIX sockets

AF_INET addressing family

This is the Internet addressing family, also referred to as the Internet domain. Socket programs communicate with socket programs on other hosts in the IP network using AF_INET family sockets which, in turn, use the AF_INET physical file system.

You can configure either AF_INET or both AF_INET and AF_INET6. You *cannot* define the stack as IPv6 only. Although coding AF_INET6 alone is not prohibited, TCP/IP will not start because the master socket is AF_INET and the call to open it will fail.

For more on this subject, refer to Chapter 3, “Base functions” on page 59 or *z/OS UNIX System Services Planning*, SA22-7800.

The AF_INET physical file system relies on other products to provide the AF_INET transport services to interact with UNIX System Services and its sockets programs.

For AF_INET/AF_INET6 sockets, the z/OS UNIX address space routes the socket request to the TCP/IP address space directly. As shown in Figure 1-3, the sockets/Physical File System layer is a transform layer between z/OS UNIX and the TCP/IP stack.

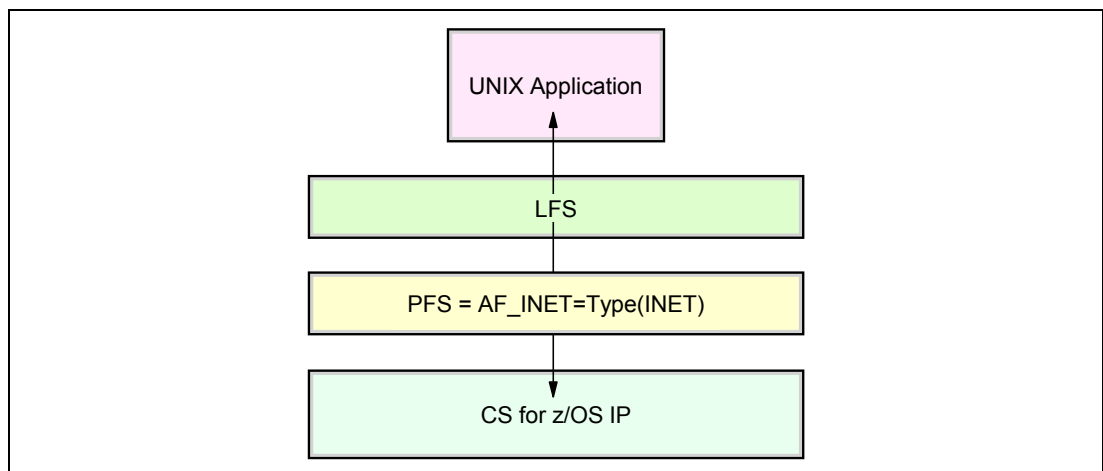


Figure 1-3 AF_INET sockets

The sockets/PFS effectively transforms the sockets calls from the z/OS UNIX interface to the TCP/IP stack regardless of the version of MVS or TCP/IP. The sockets/PFS handles the communication between the TCP/IP address space and the z/OS UNIX address space in much the same manner as High Performance Native Socket (HPNS) handles the communication between the TCP/IP address space and the TCP/IP client and server address spaces.

Physical File System transport providers

TCP/IP requires the use of the Physical File System (AF_INET) configured in two ways:

- ▶ The Integrated Sockets File System type (INET)
- ▶ The Common INET Physical File System type (CINET)

INET is used in a single-stack environment and CINET is used in a multiple-stack environment.

A single Physical File System transport provider

If your background is in a UNIX environment, it might seem strange to have a choice of using INET or CINET, because you are familiar with the TCP/IP protocol stack being an integral part of the UNIX operating system. However, this is not the case in a z/OS environment; it is very versatile. In this environment you can start multiple instances of a TCP/IP protocol stack, each stack running on the same operating system, but each stack having a unique TCP/IP identity in terms of network interfaces, IP addresses, host name, and sockets applications.

A simple example of a situation where you have more TCP/IP stacks running in your z/OS system is if you have two separate IP networks, one production and one test (or one secure and one not). You do not want routing between them, but you do want to give hosts on both IP networks access to your z/OS environment. In this situation you could implement two TCP/IP stacks, one connected to the production IP network and another connected to the test network.

This multi-stack implementation in which you share the UNIX System Services across multiple TCP/IP stacks provides challenges. Sockets applications that need to have an affinity to a particular stack need special considerations, in some cases including the coordination of port number assignments to avoid conflicts. For more information, see Chapter 3, “Base functions” on page 59.

If a single AF_INET(6) transport provider is sufficient, then use the Integrated Sockets physical file system (INET). If you need more than one AF_INET(6) transport provider (multiple TCP/IP stacks), then you must use the Common INET physical file system (CINET).

You can customize z/OS to use the Common INET physical file system with just a single transport provider (AF_INET(6)), but it is generally not recommended due to a slight performance decrease as compared to the Integrated Sockets Physical File System (INET). However, you might consider doing this if you expect to run multiple stacks in the future.

The PFS is also known under the name INET, and this appears in UNIX System Services definitions when a FILESYSTYPE and NETWORK TYPE need to be defined in the BPXPRMxx member of SYS1.PARMLIB.

Common INET Physical File System (CINET)

If you have two or more AF_INET transport providers on an MVS image, such as a production TCP/IP stack together with a test TCP/IP stack, you must use the Common INET Physical File System. Figure 1-4 shows a multiple stack environment with Common INET (CINET).

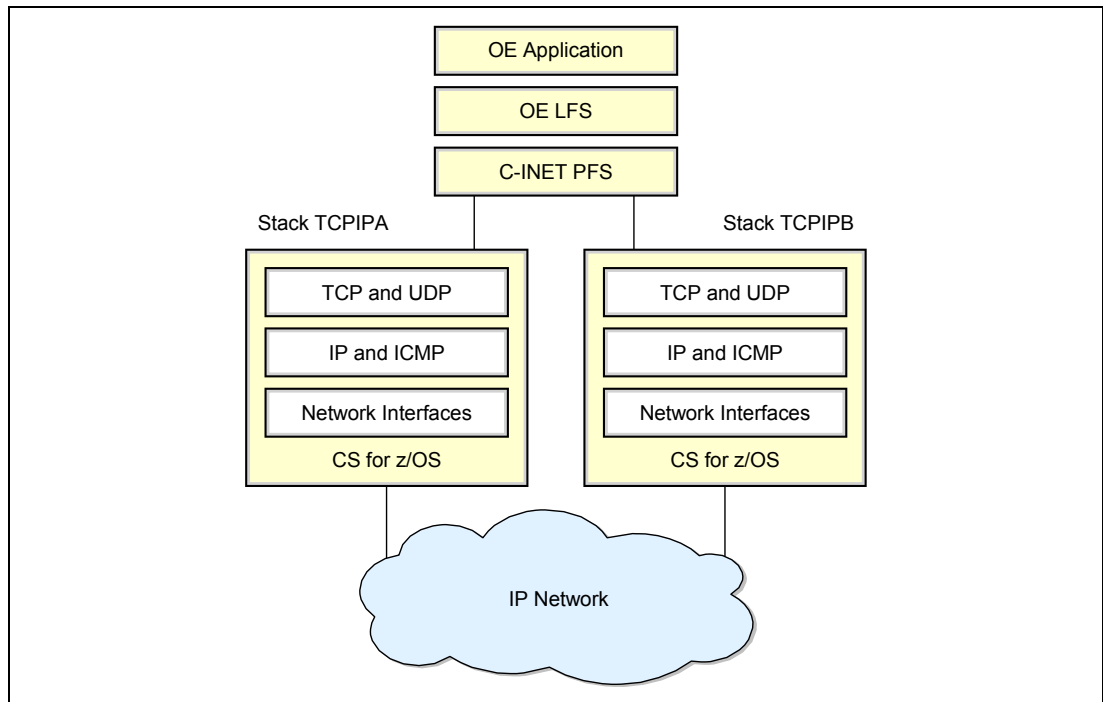


Figure 1-4 Multiple INET transport providers: CINET PFS

Recommendation: Although there are situations where multiple stacks per LPAR can provide value, in general we recommend that you implement only one TCP/IP stack per LPAR for the following reasons:

- ▶ A TCP/IP stack is capable of exploiting all available resources defined to the LPAR in which it is running. Therefore, starting multiple stacks will not yield any increase in throughput.
- ▶ When running multiple TCP/IP stacks, additional system resources, such as memory, CPU cycles, and storage, are required.
- ▶ Multiple TCP/IP stacks add a significant level of complexity to TCP/IP system administration tasks.
- ▶ It is not necessary to start multiple stacks to support multiple instances of an application on a given port number, such as a test HTTP server on port 80 and a production HTTP server also on port 80. This type of support can instead be implemented using BIND-specific support where the two HTTP server instances are each associated to port 80 with their own IP address, using the BIND option on the PORT reservation statement.

1.4 Additional information

The following IBM publications provide further details for implementing a z/OS environment that supports the TCP/IP protocol suite:

- ▶ *z/OS Communications Server: IP Configuration Guide*, SC31-8775

This document explains the major concepts of, and provides implementation guidance for, z/OS Communications Server functions.

- ▶ *z/OS Communications Server: IP Configuration Reference*, SC31-8776
This document details the parameters or statements that can be used to implement z/OS Communications Server functions.
- ▶ *z/OS Communications Server: IP Programmer's Guide and Reference*, SC31-8787
This document provides the guidelines for programming the IP applications on the z/OS.
- ▶ *z/OS Communications Server: IP Sockets Application Programming Interface Guide and Reference*, SC31-8788
This document provides detailed information about the socket API for programming the IP applications on the z/OS.

For migration, the following publications are also helpful:

- ▶ *z/OS Communications Server: New Function Summary*, GC31-8771
This document includes function summary topics to describe all the functional enhancements for the IP and SNA components of Communications Server, including task tables that identify the actions necessary to exploit new function. Use this document as a reference to using all the enhancements of z/OS Communications Server.
- ▶ *z/OS Planning for Installation*, GA22-7504
This document helps you prepare to install z/OS by providing the information you need to write an installation plan.
- ▶ *z/OS Migration*, GA22-7499
This document describes how to migrate (convert) from release to release. Use this document as a reference for keeping all z/OS applications working as they did in previous releases.
- ▶ *z/OS Introduction and Release Guide*, GA22-7502
This document provides an overview of z/OS and lists the enhancements in each release. Use this document to determine whether to obtain a new release, and to decide which new functions to implement.
- ▶ *z/OS Summary of Message and Interface Changes*, SA22-7505
This document describes the changes to interfaces for individual elements and features of z/OS. Use this document as a reference to the new and changed commands, macros, panels, exit routines, data areas, messages, and other interfaces of individual elements and features of z/OS.



The resolver

TCP/IP protocols rely upon IP addressing in order to reach other hosts in a network to communicate. For ease of use, instead of using the IP addresses represented by numbers, they are sometimes mapped to easy-to-remember names. Typically, the names are assigned to the IP address of the servers that many users can access, such as Web servers, FTP servers, and TN3270 servers.

The resolver function allows applications to use names instead of IP addresses to connect to other partners. The mapping of IP addresses and names is managed by name servers or local definitions. The resolver queries those name servers, or searches local definitions, in order to convert the name to an IP address or the IP address to a name. Using the resolver relieves users of having to remember the decimal or hexadecimal IP addresses.

The resolver is important for enabling TCP/IP stacks or TCP/IP applications to establish connections to other hosts.

This chapter discusses the following topics.

Section	Topic
2.1, "Basic concepts of the resolver" on page 20	Basic concepts of the resolver
2.2, "The resolver address space" on page 21	Key characteristics of the resolver address space
2.3, "Implementing the resolver" on page 42	The configuration and verification tasks of resolver implementation
2.4, "Problem determination" on page 50	Problem determination techniques

2.1 Basic concepts of the resolver

A *resolver* is a set of routines that acts as a client on behalf of an application. It reads a local host file or accesses one or more Domain Name System (DNS) servers for name-to-IP address or IP address-to-name resolution.

In most systems, in order for an application to reach a remote partner, it uses two commands to ask the resolver what the IP address is for a host name, or vice versa. The commands are **gethostbyname(nnnnn)** and **gethostbyaddress(aaa.aaa.aaa.aaa)**. The IPv6-enabled equivalent calls are **getaddrinfo(nnnn)** and **getnameinfo(IPaddress)**.

Figure 2-1 illustrates the information request and response flows. The resolver gets a request and, based on its own configuration file, will either look at a local hosts file or send a request to a DNS server. After the relationship between the host name and IP address is established, the resolver returns the response to the application.

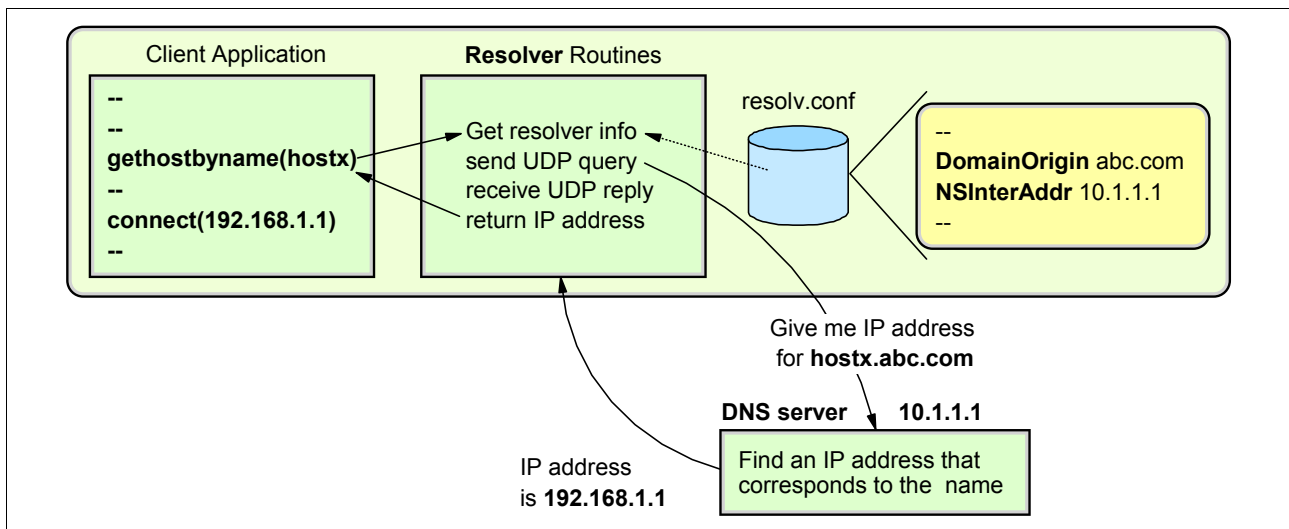


Figure 2-1 How the resolver works

As mentioned, the resolver function allows applications to use names instead of IP addresses to connect to other partners. Although using an IP address might seem to be an easy way to establish such a connection, for applications that need to connect to numerous partners, or for applications that are accessed by thousands of clients, using names is a much easier and more reliable form of establishing access.

Another important reason to use names instead of IP addressing is that a user or an application is not affected by the IP address changes to the underlying network.

Table 2-1 compares the benefits and drawbacks of the use of hard-coded IP addresses and the two name resolution methods:

- The local hosts file
- The name server (DNS)

Table 2-1 Comparing the use of direct addressing with name resolution

	Hard-coded IP addresses	Local hosts file	Domain Name System (DNS)
Technology	None - Use the entered IP address directly on the connect() or sendto() socket call.	Use gethostbyname() and let the resolver find an IP address in the locally configured hosts file.	Use gethostbyname() and let the resolver contact the configured name server for an IP address.
Benefits	Fast (no name resolution). Good in some debugging situations (you know exactly which IP address is being used).	Fast (local name resolution).	IP address changes can be done without any local changes. All host names (in the entire network) can be resolved. A hierarchical name space.
Drawbacks	Difficult to remember IP addresses. Very inconvenient if an IP address change occurs. Just think about IPv6.	If an IP addressing change is needed, all the local hosts files have to be updated. Only locally configured host names can be resolved.	Additional packets (requests) flow to resolve a host name before a destination can be reached. ^a

a. Resolver caching (discussed in 2.2.4, “Resolver DNS cache” on page 28) alleviates some of the need for these flows by saving previously obtained information locally.

2.2 The resolver address space

In z/OS systems, the resolver works as a procedure. The resolver must be started before TCP/IP stacks or any TCP/IP applications issue the resolver calls. It can be started in one of the following ways:

- Default z/OS UNIX resolver

If no customized resolver address space is configured, the z/OS UNIX System Services starts the default resolver. The default resolver is named RESOLVER. To use the default RESOLVER address space, specify the RESOLVER_PROC(DEFAULT) statement or do not specify any RESOLVER_PROC statements in BPXPRMxx.

- Customized resolver address space

The customized resolver address space can specify additional options to control the use of the resolver configuration file. To create the customized resolver address space, create a resolver started procedure and a SETUP data set to specify the additional options. The customized resolver address space can be started automatically with the RESOLVER_PROC(procname) statement in BPXPRMxx.

Although the resolver address space can be started manually, we recommend that you start the resolver address space automatically during initialization of the UNIX System Services by defining the RESOLVER_PROC() statement within BPXPRMxx.

After the resolver address space is activated, the global TCPIP.DATA statements cannot be overridden unless the MODIFY command is issued.

2.2.1 The resolver SETUP data set

The resolver SETUP data set is used by the customized resolver address space. The default z/OS UNIX resolver does not read this file. The SETUP data set can include the following statements:

- ▶ GLOBALTCPIPDATA
- ▶ DEFAULTTCPIPDATA
- ▶ GLOBALIPNODES
- ▶ DEFAULTIPNODES
- ▶ COMMONSEARCH or NOCOMMONSEARCH
- ▶ CACHE or NOCACHE
- ▶ CACHESIZE
- ▶ MAXTTL
- ▶ UNRESPONSIVETHRESHOLD

The use of each statement is discussed in later sections.

2.2.2 The resolver configuration file

The resolver configuration file is often called TCPIP.DATA. In this file you can define how the resolver should query the name-to-address or address-to-name resolution to the name servers or search the local hosts file.

The configuration file can be an MVS data set or a z/OS UNIX Hierarchical File System (HFS) file.

Note: The publication *z/OS Communications Server: IP Configuration Guide*, SC31-8775, contains useful information about the characteristics that are required for the z/OS data sets or file system files that contain resolver SETUP and configuration statements. The guide also points out the security characteristics and file system permission settings that are needed.

TCPIP.DATA configuration statements

The following basic statements should be defined in the TCPIP.DATA file.

- ▶ TCPIPJOBNAME (equivalent to TCPIPUSERID)
The name of the procedure used to start the TCP/IP address space. The default is TCPIP.
- ▶ DOMAIN (equivalent to DOMAINORIGIN)
The domain origin that is appended to the host name to form the fully qualified domain name of a host.
- ▶ HOSTNAME
The TCP host name of the z/OS Communications Server server.
- ▶ LOOKUP
The order in which the DNS or local host files are to be searched for name resolution. By default, DNS is looked up first. If caching is in effect, the resolver cache is considered to be part of the “DNS” lookup step, and resolver will examine its cache data prior to actually

contacting any name server. Then if the resolution is unsuccessful, the local host files are searched.

► **NSINTERADDR** (equivalent to NAMESERVER)

The IP address of a name server the resolver should query to.

► **DATASETPREFIX**

The high-level qualifier for the dynamic allocation of data sets. DATASETPREFIX is referred to as the *hlq* of the TCP/IP stacks.

► **NOCACHE**

You must specify NOCACHE in the TCPIP.DATA data set if you want to prevent applications using this data set from either querying the cache or adding records to the cache.

TCPIP.DATA search order

On z/OS, the configuration file is located based on the search order. You must be mindful of this search order, to ensure that the resolver works in the way you expect.

The TCP/IP applications execute a set of commands in the Sockets API Library to initiate a request to the resolver in z/OS. The Sockets API Library uses one of the following socket environments:

- Native MVS environment
- z/OS UNIX environment

Table 2-2 lists some of the APIs, z/OS applications, and user commands that use the active MVS environment and the z/OS UNIX environment.

Table 2-2 Socket APIs, applications, and commands in Native MVS or z/OS UNIX environment

	Native MVS environment	z/OS UNIX environment
Socket APIs	TCP/IP C Sockets TCP/IP Pascal Sockets TCP/IP REXX Sockets TCP/IP Sockets Extended IMS Sockets CICS Sockets	Language Environment® C Sockets UNIX System Services
z/OS Applications	TN3270 Telnet server SMTP CICS Listener LPD Miscellaneous server PORTMAP RSHD	FTP OMPROUTE CSSMTP SNMP z/OS UNIX OPORTMAP z/OS UNIX OREXECD z/OS UNIX ORSHD

	Native MVS environment	z/OS UNIX environment
User commands	TSO FTP (batch) TSO NETSTAT TSO NSLOOKUP TSO PING TSO TRACERTE TSO DIG TSO LPR TSO REXEC TSO RPCINFO TSO RSH	TSO FTP (command) netstat nslookup ping tracert ftp host hostname dnsdomainname dig rexec rpcinfo sendmail snmp

Each socket environment uses a different search order of the resolver configuration file, as shown in Figure 2-2.

<p>Native MVS environment</p> <ol style="list-style-type: none"> 1. GLOBALTCPIPDATA 2. //SYSTCPD DD statement 3. userid/jobname.TCPIP.DATA 4. SYS1.TCPPARMS(TCPDATA) 5. DEFAULTTCPIPDATA 6. TCPIP.TCPIP.DATA 	<p>z/OS UNIX environment</p> <ol style="list-style-type: none"> 1. GLOBALTCPIPDATA 2. RESOLVER_CONFIG environment variable 3. /etc/resolv.conf 4. //SYSTCPD DD statement 5. userid/jobname.TCPIP.DATA 6. SYS1.TCPPARMS(TCPDATA) 7. DEFAULTTCPIPDATA 8. TCPIP.TCPIP.DATA
--	---

Figure 2-2 The resolver configuration file search order for each socket environment

Note: UNIX System Services Callable sockets use the z/OS UNIX environment search order but cannot use the RESOLVER_CONFIG environment variable.

This provides the flexibility to control the resolver lookup differently, depending on which socket API the application uses. However, because of the difference in search orders, it could sometimes cause an unexpected result in the address resolution.

For example, if you set up /etc/resolv.conf as your resolver configuration file, the FTP server application that uses the z/OS UNIX search order can resolve the name-to-address or address-to-name successfully. However, the TN3270 server, which uses the native MVS search order, would fail because /etc/resolv.conf is not included in its search list.

Using GLOBALTCPIPDATA

In order to deal with the complexity of the different search orders in the environments, the GLOBALTCPIPDATA statement was introduced. Using the GLOBALTCPIPDATA statement, you can use the same resolver configuration file throughout the z/OS system, because it is the first choice in all socket search orders. This consolidation allows for consistent name resolution processing across all TCP/IP applications.

To specify the GLOBALTCPIPDATA statement, you need to create a resolver started procedure and its SETUP data set, instead of using the z/OS UNIX System Services default

RESOLVER address space. The use of the resolver address space and GLOBALTCPIPDATA statement simplifies the resolver configuration on z/OS.

The TCPIP.DATA file specified by the GLOBALTCPIPDATA statement is often called the *global TCPIP.DATA file*. If you define GLOBALTCPIPDATA, the following statements can be included only in the global TCPIP.DATA file:

- ▶ DomainOrigin/Domain or Search
- ▶ NSInterAddr/NameServer
- ▶ NSPortAddr
- ▶ ResolveVia
- ▶ ResolverTimeOut
- ▶ ResolverUDPRetries
- ▶ SortList

Other TCPIP.DATA statements can be optionally included in the global TCPIP.DATA file, and the definition in the global TCPIP.DATA always has precedence. If TCPIPJobname is specified in both the global TCPIP.DATA file and the local (non-global) TCPIP.DATA file, then the one in the global TCPIP.DATA file is used.

If other TCPIP.DATA statements, such as HostName and TCPIPJobname, cannot be found in the global TCPIP.DATA file, then the resolver continues its search according to the search order of the each socket environment. The search stops when the file is found.

If statements such as HostName and TCPIPJobname cannot be found in that file either, the defaults are applied. Note that it does not continue searching in the list. In other words, a maximum of two files can be used (global TCPIP.DATA file and one TCPIP.DATA file in the search order list).

Using GLOBALTCPIPDATA, the administrators can specify which statements should be applied throughout the z/OS image, and decide which statements can be customized by each socket environment by omitting those statements in the global TCPIP.DATA file.

Note: In the Common INET (CINET) multi-stack environment, you should omit the TCPIPJobname statement from the global TCPIP.DATA file so that each TCP/IP stack, or the applications that have affinity to a stack, can specify a local TCP.DATA with its own TCPIPJobname statement.

When using GLOBALTCPIPDATA in the CINET environment, the name server specified by NSInterAddr or NameServer in the global TCPIP.DATA file must be accessible from all TCP/IP stacks that issue resolver calls.

Figure 2-3 on page 26 depicts the relationship between global TCPIP.DATA and local TCPIP.DATA.

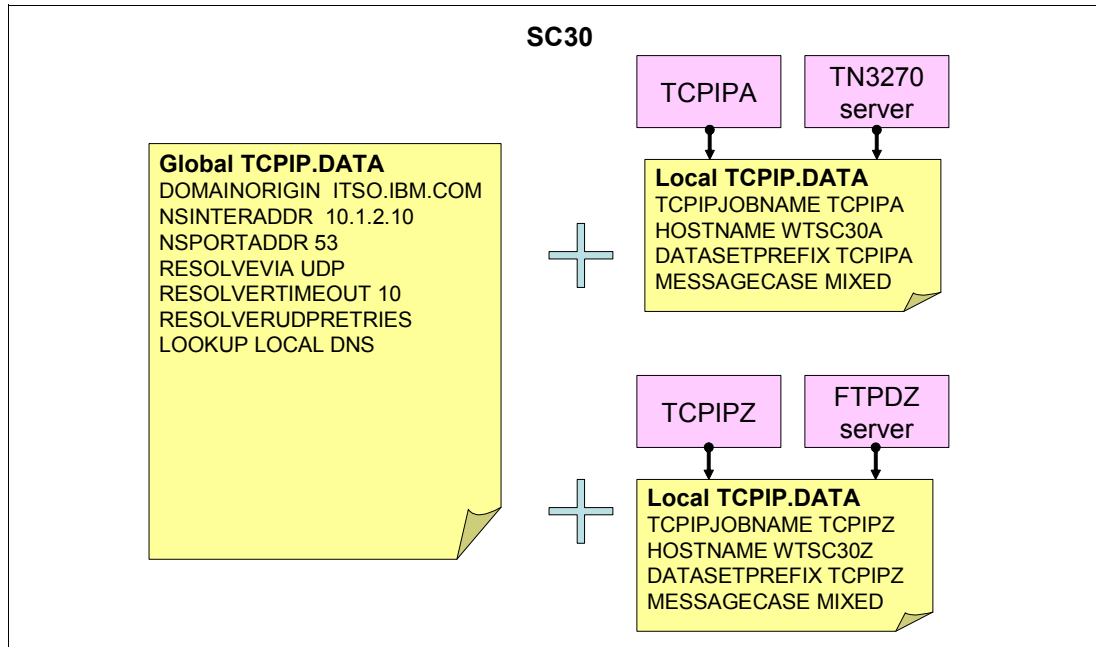


Figure 2-3 Using global TCPIP.DATA and Local TCPIP.DATA

Using DEFAULTTCPIPDATA

DEFAULTTCPIPDATA can be specified in the resolver SETUP data set to define the last choice of the TCPIP.DATA in the search order. The file specified by DEFAULTTCPIPDATA is used when the application does not specify the local (non-global) TCPIP.DATA.

2.2.3 Local hosts file

The local hosts file lists the mapping of the IP addresses and the names just like the name servers, but held locally on the server. The LOOKUP statement in the TCPIP.DATA configuration file defines whether the resolver address space performs the name resolution only in the local files, or using the defined name server (including resolver cache, if active), or both, in any specified order.

Using COMMONSEARCH

When the local hosts file is searched, the search order for the native MVS environment and the z/OS UNIX environment are different. The difference in the search orders adds complexity to configuration tasks and can lead unexpected results of the name resolution.

The simpler approach is to utilize the COMMONSEARCH statement in the resolver SETUP data set. By specifying COMMONSEARCH, native MVS and z/OS UNIX environments use the same search order as shown in Figure 2-4 (except the RESOLVER_IPNODES environment variable, which is only supported by the z/OS UNIX environment). In both environments, the first choice is the file specified by GLOBALIPNODES statement, which is defined in the resolver SETUP data set.

The local hosts files looked up in this search order are typically called ETC.IPNODES files. When the COMMONSEARCH is specified in the resolver SETUP data set, it uses the same search order for both IPv4 and IPv6 queries. You can list both IPv4 and IPv6 addresses in the ETC.IPNODES file.

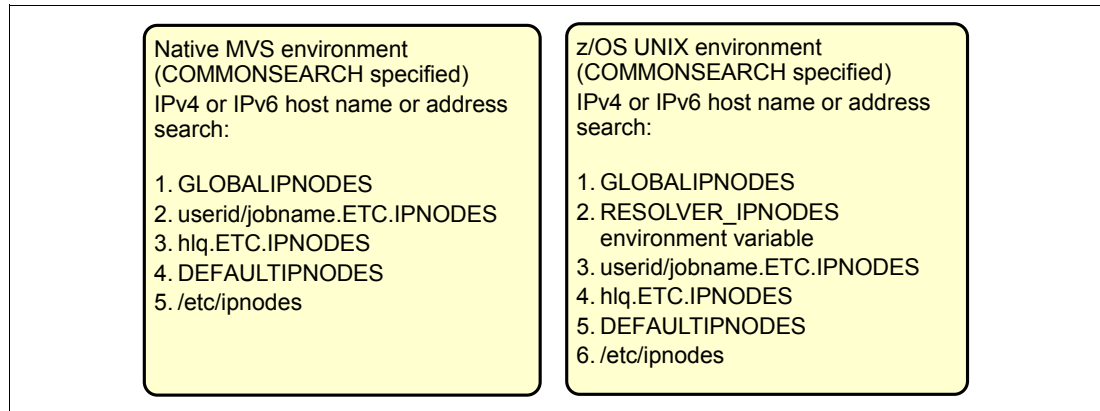


Figure 2-4 Local hosts file search order with COMMONSEARCH specified

To determine which environment is used for a particular socket's APIs, applications, or commands, refer to Table 2-2 on page 23.

If COMMONSEARCH is not specified in the resolver SETUP data set, then the default is NOCOMMONSEARCH and the default search order shown in Figure 2-5 on page 27 is used.

Using GLOBALIPNODES

The GLOBALIPNODES statement specifies the global local host file that is to be used in the entire z/OS image, regardless of which environment (native MVS or z/OS UNIX) that the applications or sockets API use. To put the GLOBALIPNODES statement into effect for the name resolution of IPv4 addresses, also specify COMMONSEARCH in the resolver SETUP data set.

Using DEFAULTIPNODES

The DEFAULTIPNODES statement specifies the last candidate of the local host file search. To put the DEFAULTIPNODES statement into effect for the name resolution of IPv4 addresses, also specify COMMONSEARCH in the resolver SETUP data set.

Default local hosts file search order

If NOCOMMONSEARCH (the default) is specified in the resolver SETUP data set or default z/OS UNIX resolver is used, the default local hosts file search order shown in Figure 2-5 is used. The default local hosts file search order only applies to the query of IPv4 addresses. The query for IPv6 addresses always uses the search order listed in Figure 2-4.

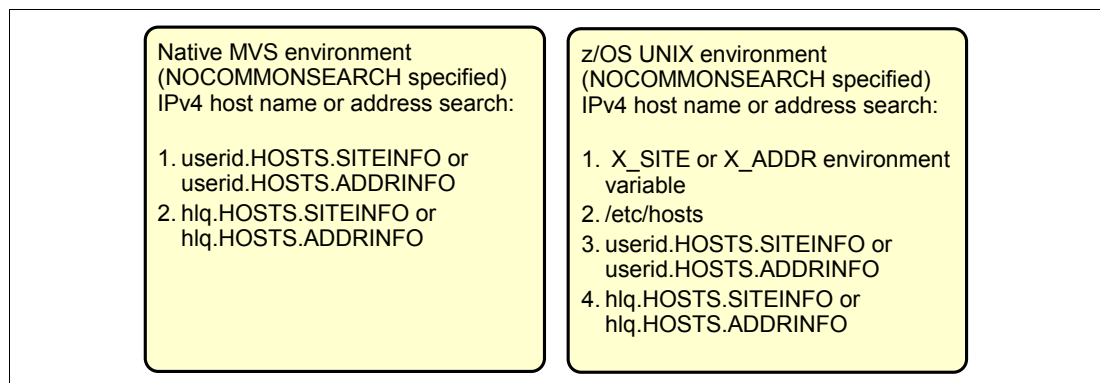


Figure 2-5 Local hosts file search order with NOCOMMONSEARCH specified (default)

2.2.4 Resolver DNS cache

In order to provide better system performance, consider eliminating redundant network flows to DNS servers. You can accomplish this goal by exploiting the resolver DNS cache. This uses the resolver for system-wide caching of Domain Name System (DNS) responses.

Using CACHE or NOCACHE

The resolver cache is enabled by default and is shared across the entire z/OS system image. If you are currently running a caching-only DNS name server, you might be able to use the resolver DNS cache instead; the resolver DNS cache provides the same function with better system performance.

Two of the new resolver setup file statements are **CACHE** and **NOCACHE**. The **CACHE** statement, which is the default, explicitly indicates that resolver caching is active across the entire system. The **NOCACHE** statement explicitly indicates that resolver caching is *not* active across the entire system. You must code **NOCACHE** if you want to maintain the current level of resolver processing. The setting of **CACHE** or **NOCACHE** can be changed dynamically using the **MODIFY RESOLVER,REFRESH,SETUP** command. If you change from a setting of **CACHE** to a setting of **NOCACHE** dynamically, any existing cache records are immediately deleted.

Figure 2-6 shows how the resolver DNS cache works.

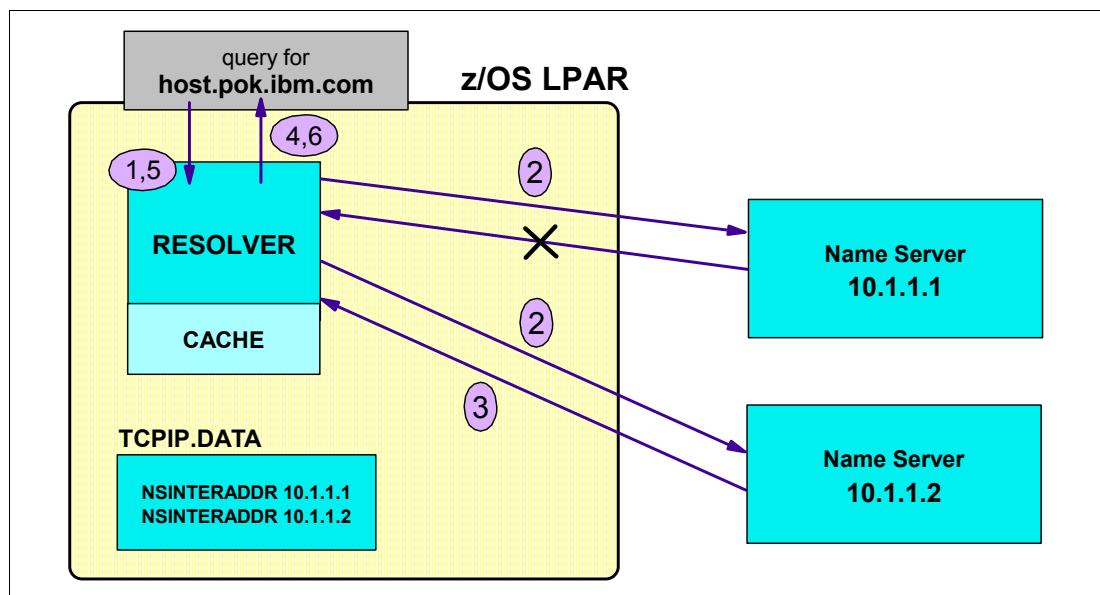


Figure 2-6 How the cache works

In step 1, an application delivers a request to translate the host name `host.raleigh.ibm.com` into an IP address. The resolver, in step 2, forwards the request to the first DNS name server specified in the list of name servers in the `TCPIP.DATA` data set. The response from the name server in step 3 is returned to the application in step 4. If the first DNS name server (10.1.1.1) does not respond in time, the resolver forwards the request to the second name server (10.1.1.2) in the list. In this point, the resolver now saves the information into the local resolver cache. At step 5, when the second request for the host name translation is received, the resolver first queries the local cache for data about the host name. In this example, the information is there, and is still valid, so the resolver returns the response data immediately to the application (step 6).

Using CACHESIZE(*size*)

CACHESIZE indicates how much storage the cache function can use to hold resolver cache information. The valid range for *size* is 1 MB to 999 MB. The default is 200 MB. For planning purposes, assume a megabyte of data holds slightly more than 400 entries and consider coding a CACHESIZE at least 50% greater than your expected needs.

Tip: When CACHESIZE is specified with NOCACHE, the value is ignored.

Important: You can modify the CACHESIZE using the MODIFY REFRESH,SETUP command, but you can only increment the storage amount or keep it the same. To decrease the value of CACHESIZE M, you must stop and restart the resolver.

Using MAXTTL(*time*)

MAXTTL indicates the longest amount of time that the resolver cache can use saved information. The valid range for *time* is 1 to 2147483647 (seconds). The default is 2,147,483,647, which is the largest TTL a name server can return.

You can dynamically change MAXTTL using the MODIFY REFRESH,SETUP command. Changing the value of MAXTTL has no impact on any records currently in the cache. The value can be increased or decreased, but the new value only affects records created after the MODIFY RESOLVER command completes.

Tip: When MAXTTL is specified with NOCACHE, the value is ignored.

Using the cached information

The resolver caching function does not impact the data that is presented to the application across the resolver APIs. The same control block structures are used for returning the information. Applications invoking the resolver should not detect any difference between data supplied from the cache and data that had to be retrieved from a name server.

Furthermore, the cache function is designed to allow resource information to be re-used by compatible API calls. For example, if Getaddrinfo is used to obtain IPv4 addresses for a host name, that same cached information can be retrieved later using Gethostbyname. The same capability exists for Getnameinfo and Gethostbyaddr calls in terms of host names obtained from an IPv4 address. IPv6 processing is only available using Getaddrinfo and Getnameinfo, so IPv6 information cannot be shared in this manner. In addition, the resolver caching function will handle translating the cache information from EBCDIC to ASCII, or vice versa, so cached information is available using either protocol.

One function not provided by the resolver caching function is the ability to return the cached IP addresses in a different, or “round robin”, order than they were received from the name server. The resolver returns the addresses in the same order all the time. It should be noted that Getaddrinfo processing re-orders the list of addresses already; refer to “Default destination address selection” on page 38 for details. Similarly, if you have SORTLIST configuration statements coded, they will also re-order the list of addresses into a more predictable pattern.

Testing the resolver DNS cache

Verify if the resolver is able to cache the expected name-to-address resolution. First, we delete all information from the resolver cache using the `modify resolver,flush,all` command, as shown in Example 2-1.

Example 2-1 MODIFY RESOLVER,FLUSH,ALL command display

```
F RESOLV30,FLUSH,ALL
EZZ9305I 1 CACHE ENTRIES DELETED
EZZ9293I FLUSH COMMAND PROCESSED
```

Display cache entry data using netstat RESCache command

You can use the **netstat RESCache** command in the console or TSO command line to display information regarding the resolver cache. In the UNIX System Services (now called z/OS UNIX) environment the same command is **netstat -q**. Two main types of information can be displayed: statistical information and actual resource information.

You can specify additional modifiers or filters to influence the amount of cache data that is displayed. For statistical information, you can add the DNS modifier to have the overall statistics broken into statistical information on a name server IP address basis. You have even more options for detailed entry information reports. You can filter the information by the IP address of the name server that provided the information. You can filter the information so that only entries related to a specific host name or IP address value are displayed. You can display only negative cache information from the cache, either all entries or subsets of entries based on name server IP address, host name value, or IP address value.

We then verify if the name admin.itso.ibm.com is in the cache by using the **netstat -q** command, as shown in Example 2-2.

Example 2-2 netstat -q DETAIL command display

```
CS02 @ SC30:/u/cs02>netstat -p tcpipa -q DETAIL -H admin.itso.ibm.com
MVS TCP/IP NETSTAT CS V1R12      TCPIP Name: TCPIPA      15:49:22
```

Now we can resolve admin.itso.ibm.com using a **ping** command, and reenter the **netstat -q** command, as shown in Example 2-3 and Example 2-4.

Example 2-3 UNIX ping command display

```
CS02 @ SC30:/u/cs02>ping -p tcpipa admin.itso.ibm.com
CS V1R12: Pinging host admin.itso.ibm.com (10.1.4.11)
Ping #1 response took 0.000 seconds.
```

Example 2-4 netstat -q DETAIL command display

```
CS02 @ SC30:/u/cs02>netstat -p tcpipa -q DETAIL -H admin.itso.ibm.com
MVS TCP/IP NETSTAT CS V1R12      TCPIP Name: TCPIPA      00:12:40
HostName to IPAddress translation
-----
HostName: ADMIN.ITSO.IBM.COM      1
DNS IPAddress: 10.1.2.10          2
DNS Record Type: T_A
Canonical Name: admin.itso.ibm.com
Cache Time: 09/27/2010 04:11:10
Expired Time: 09/27/2010 04:21:10 3
Hits: 1                          4
IPAddress: 10.1.4.11             5
```

In Example 2-4 on page 30, the numbers correspond to the following information:

- 1.** The entry-name is already in the cache.
- 2.** The DNS Server IP address where the resolver found the entry-name admin.itso.ibm.com.
- 3.** The expiration time, which is 600 seconds. This is the time in the MAXTTL statement or in the TTL value for this entry, as supplied by the name server at 10.1.2.10.
- 4.** How many times this entry-name (admin.itso.ibm.com) was used by the resolver while remaining in the cache.
- 5.** The IP address of the entry-name admin.itso.ibm.com.

In Example 2-2 on page 30, the resolver had not yet cached itso.ibm.com. However, after a **ping** command, admin.itso.ibm.com was cached, as shown in Example 2-4 on page 30.

Now we display the cache statistics by using the **netstat -q SUMMARY DNS** command, as shown in Example 2-5.

Example 2-5 netstat -q SUMMARY DNS command display

```
CS02 @ SC30:/u/cs02>netstat -p tcpipa -q SUMMARY DNS
MVS TCP/IP NETSTAT CS V1R12          TCPIP Name: TCPIPA          00:19:21
Storage Usage:
  Maximum: 20M      1
  Current: 19K      MaxUsed: 21K    2
Cache Usage:
  Total number of entries: 1
  Non-NX entries: 1
    A: 1          AAAA: 0          PTR: 0
  NX entries: 0
    A: 0          AAAA: 1          PTR: 0
  Queries: 1              Hits: 0
  SuccessRatio: 0%      3

DNS address: 10.1.2.10    4
  Total number of entries: 1    5
  Non-NX entries: 1
    A: 1          AAAA: 0          PTR: 0
  NX entries: 0
    A: 0          AAAA: 0          PTR: 0
  References: 1              Hits: 0
```

In this example, the numbers correspond to the following information:

- 1.** The maximum amount of storage permitted, or CACHESIZE
- 2.** The current amount of storage in use and maximum amount the resolver has ever used for caching since the resolver was started
- 3.** Percentage of queries satisfied by information in the cache
- 4.** IP address of the name server providing cache data
- 5.** Number of entries in cache, grouped by negative (NX) entries and other (Non-NX) entries

Cache usage statistics include the total number of entries in the cache and the volume of activity involving the cache. The number of entries is differentiated between negative cache entries and non-negative cache entries. Within each of these main categories, the number of

DNS A, AAAA, and PTR records is indicated. These same subsets of entries are displayed for individual name servers.

The number of resolver cache requests and how often usable data was returned by the cache gives you a sense of the efficiency of your cache operations. Note that a single resolver API call can generate multiple cache queries. For example, a Getaddrinfo request for both IPv6 and IPv4 addresses generates two cache queries. On an individual name server level, the “References” value indicates the number of times the set of cache information provided by this name server was examined. Typically, the sum of the name server “References” values is greater than the number of cache queries, because multiple name server information sets can be examined as part of one cache query.

Now, we display the cache statistics by using the **netstat -q DETAIL** command, to display information about a specific cache entry, as shown in Example 2-6 on page 32.

Example 2-6 netstat -q DETAIL command display

```
CS02 @ SC30:/u/cs02>netstat -p tcpipa -q DETAIL -H admin.itso.ibm.com
MVS TCP/IP NETSTAT CS V1R12          TCPIP Name: TCPIPA          00:28:20
HostName to IPAddress translation
-----
HostName: ADMIN.ITSO.IBM.COM
  DNS IPAddress: 10.1.2.10      1
  DNS Record Type: T_A
  Canonical Name: admin.itso.ibm.com
  Cache Time: 09/27/2010 04:26:22 2
  Expired Time: 09/27/2010 04:36:22 3
  Hits: 1      4
  IPAddress: 10.1.4.11      5
```

In this example, the numbers correspond to the following information:

- 1.** The DNS that provided the entry
- 2.** The time and the date of cache entry creation
- 3.** The time and date when the entry will expire
- 4.** The number of times this entry has been re-used
- 5.** IP addresses provided by the specified name server

This is a partial example of a **netstat** report showing a detailed cache entry. The reports are formatted so that DNS A and AAAA records are displayed as one group, and DNS PTR records are displayed as a second group. Negative cache entries can appear in either group, in any order, and are identified using the notation “***NA***”.

For each record, the cache entry key, or the target resource that was searched for to acquire this cache information, is the first line of the entry. After that, the two types of entries are very similar. The IP address of the DNS name server that supplied this particular information is displayed, allowing you to see what values were provided by what name servers. In the case of DNS A and AAAA record entries, the host name used to create the record might really be an alias or nickname for the official name of the resource. For that reason, the display includes the official, or canonical, name, regardless of whether the names match or not. There is no canonical name concept for DNS PTR records.

Two time values are displayed: one is the time and the date of cache entry creation. The other is the time and date when the entry will expire, based on name server TTL or MAXTTL

setting. The **netstat RESCACHE** report will not include any resources that are in the cache that represent expired information. The number of times this entry has been re-used is displayed as the “Hits” value. Finally, for DNS A and AAAA entries, up to 35 IP addresses provided by the specified name server for the host name value are included. For DNS PTR entries, the one host name associated with the input IP address (either IPv4 or IPv6) is included.

2.2.5 Criteria for indicating an unresponsive DNS name server

A DNS name server is considered unresponsive for a specific query when:

- ▶ The resolver sends a UDP or TCP query to a name server and never receives a response.
- ▶ The resolver sends a UDP query to a name server and receives a response after the RESOLVERTIMEOUT value has expired.
- ▶ The resolver attempts to send data to a name server using UDP, but the data cannot be sent to the target IP address (for example, because of an error in the route configuration).
- ▶ The resolver attempts to connect to a name server using TCP, but the connection attempt times out.

The unresponsive DNS notification function is enabled by default. It can be turned off by specifying the UNRESPONSIVETHRESHOLD configuration statement with a value of zero (0).

2.2.6 Unresponsive DNS name servers

Communications Server for z/OS IP has the capability of notifying the operator console when a defined Domain Name System (DNS) name server does not respond to resolver queries during the most recent sliding 5-minute interval.

Resolver also provides statistics for each currently unresponsive name server, regarding the number of queries attempted and the number of queries which received no response during a sliding 5-minute interval.

CS for z/OS IP considers a DNS name server to be unresponsive when the number of unsuccessful queries exceeds a percentage threshold of the total queries sent during a 5-minute interval. By default, the percentage threshold is 25% of the total queries. This percentage can be customized using the UNRESPONSIVETHRESHOLD configuration statement in the resolver setup file.

The percentage threshold value can also be changed while the resolver is active, by changing the UNRESPONSIVETHRESHOLD configuration statement in the resolver setup file and issuing the **MODIFY resolver,REFRESH,SETUP=setup_file_name** command.

Resolver notifications for DNS name server responsiveness

When the resolver detects that a name server is not being responsive, based on the provided failure threshold, network operator messages are issued to report the problem, as shown in Example 2-7.

Example 2-7 Notifying DNS responsiveness

```
*EZZ9308E UNRESPONSIVE NAME SERVER DETECTED AT IP ADDRESS 10.1.2.10
EZZ9310I NAME SERVER 10.1.2.10
      TOTAL NUMBER OF QUERIES SENT 132
      TOTAL NUMBER OF FAILURES    132
      PERCENTAGE                   100%
```

The error message EZZ9308E is issued only once. It remains on the operator console for as long as the resolver considers the name server to be unresponsive. During that time of unresponsiveness, the informational message EZZ9310I is reissued every 5 minutes, giving updated statistics for the unresponsive name server for that sliding 5-minute interval.

If by the end of a subsequent monitor interval, the resolver determines that the name server's failure rate has dropped below the threshold value, the resolver considers this name server to be responsive again, clears message EZZ9308E from the operator console, and issues a message indicating the DNS is responsive again, as shown in Example 2-8.

Example 2-8 Notifying DNS is now responsive

```
EZZ9309I NAME SERVER IS NOW RESPONSIVE AT IP ADDRESS 10.1.2.10
EZZ9310I NAME SERVER 10.1.2.10
      TOTAL NUMBER OF QUERIES SENT 190
      TOTAL NUMBER OF FAILURES    19
      PERCENTAGE                   10%
```

Considerations about UNRESPONSIVETHRESHOLD usage

When you specify the UNRESPONSIVETHRESHOLD value, consider the following factors that have an impact on the network environment:

- ▶ Specifying a small value may generate a large number of console messages.
- ▶ Specifying a value that is too large might result in intermittent problems with the DNS name server or the IP network not being detected.
- ▶ Consider using a higher value for UNRESPONSIVETHRESHOLD if you use a small value for RESOLVERTIMEOUT. If you set a very short timeout value, even temporary problems in the network might generate unnecessary unresponsive name server messages to the operator.
- ▶ The values specified on the RESOLVERUDPRETRIES, SEARCH, and NAMESERVER statements in the TCPIP.DATA file can affect the number of messages generated by the system resolver.

2.2.7 Affinity servers and generic servers

In the multiple stack environment, a TCP/IP application might have an affinity to a specific TCP/IP stack. When designing a multiple stack system, it is important to check each application that will be used and how it will be implemented in the environment.

Affinity server

An *affinity server* is an application that has affinity to a specific TCP/IP stack; it provides service to the clients that are connected through the TCP/IP stack to the applications.

In this case, you need to code a TCP/IPJobname statement that represents the application in order to direct traffic to a specific stack. So, when designing the global definitions in the resolver address space, do not code a TCPIPJobname statement in GLOBALTCPIPDATA. Instead, allow it to be coded in the local TCPIP.DATA.

A native TCP/IP sockets program will always use one stack only, and by default, it will be the stack that is identified in the TCPIPJOBNAME option in the chosen resolver configuration file. However, the stack can also be chosen through the program configuration and API calls to associate the program with a chosen stack, as shown in Figure 2-7 on page 35.

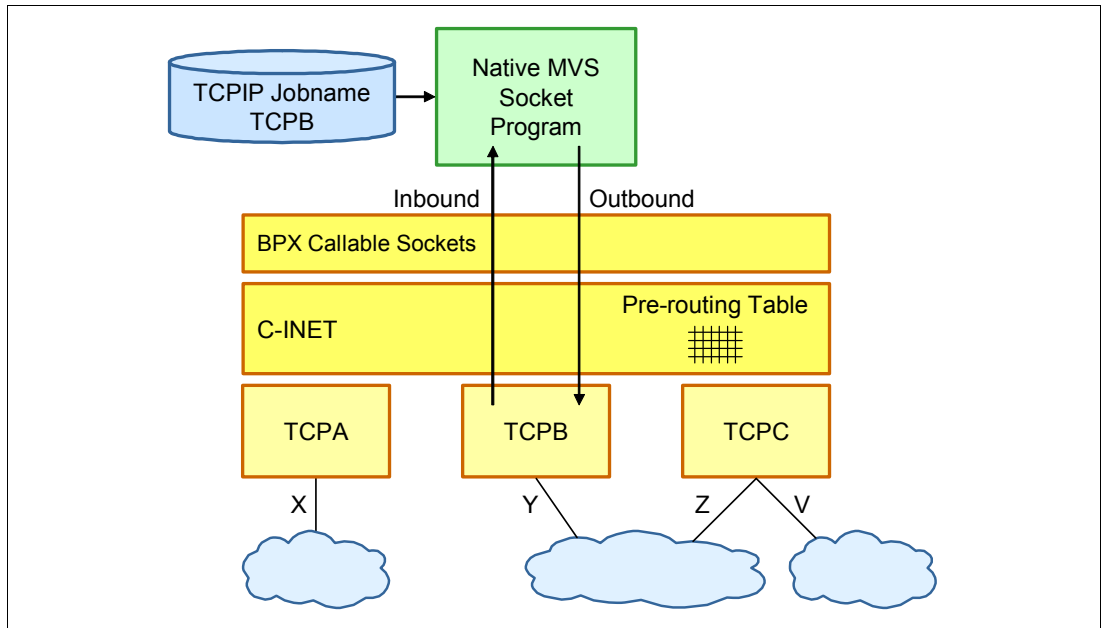


Figure 2-7 Native TCP/IP applications in a multiple stack environment

Applications using UNIX System Services callable APIs or Language Environment C/C++ sockets APIs can also use a specific bind to open a socket. A bind-specific server socket will only receive connections from the stack that owns the IP address to which the socket is bound. Outbound connections or UDP datagrams will be handled by the stack that offers the best route to the destination IP address, as shown in Figure 2-8.

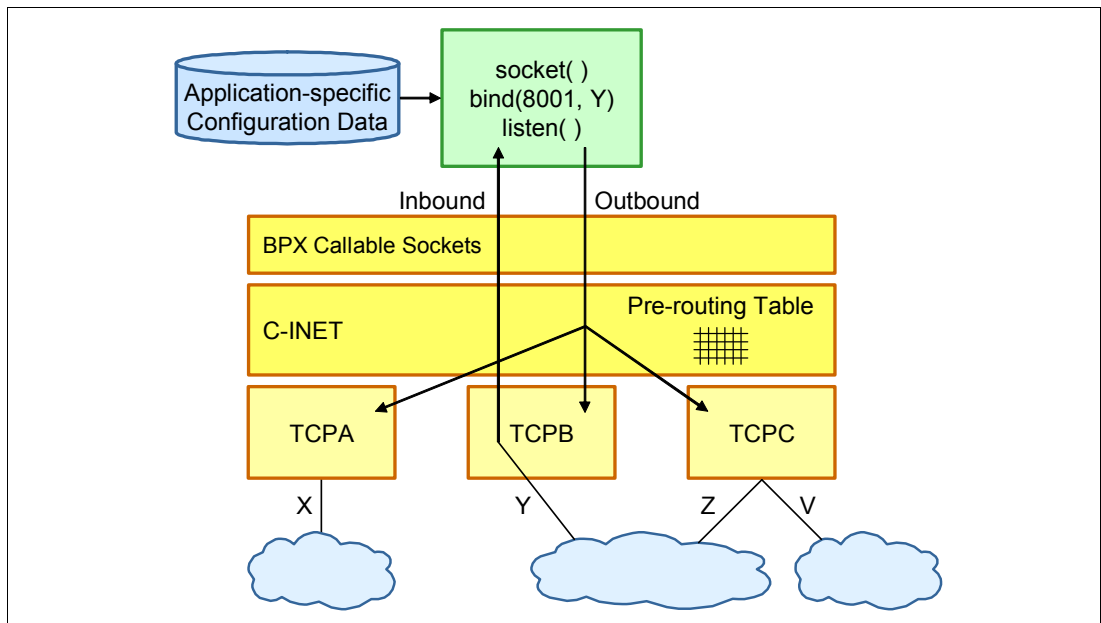


Figure 2-8 APIs, bind-specific

Generic server

A *generic server* is a server without an affinity to a specific stack, and it provides service to any clients that are connected to any TCP/IP stacks on the system.

When using the generic bind, it does not matter if the chosen resolver configuration file has a TCPIPJobname; it is not used when the server is a pure generic server.

Applications using UNIX System Services callable APIs or Language Environment C/C++ sockets APIs can be implemented using a generic bind to open the same port in all TCP/IP stacks. By doing so, the application will accept incoming connections or UDP datagrams over any interface of all connected stacks, as shown in Figure 2-9.

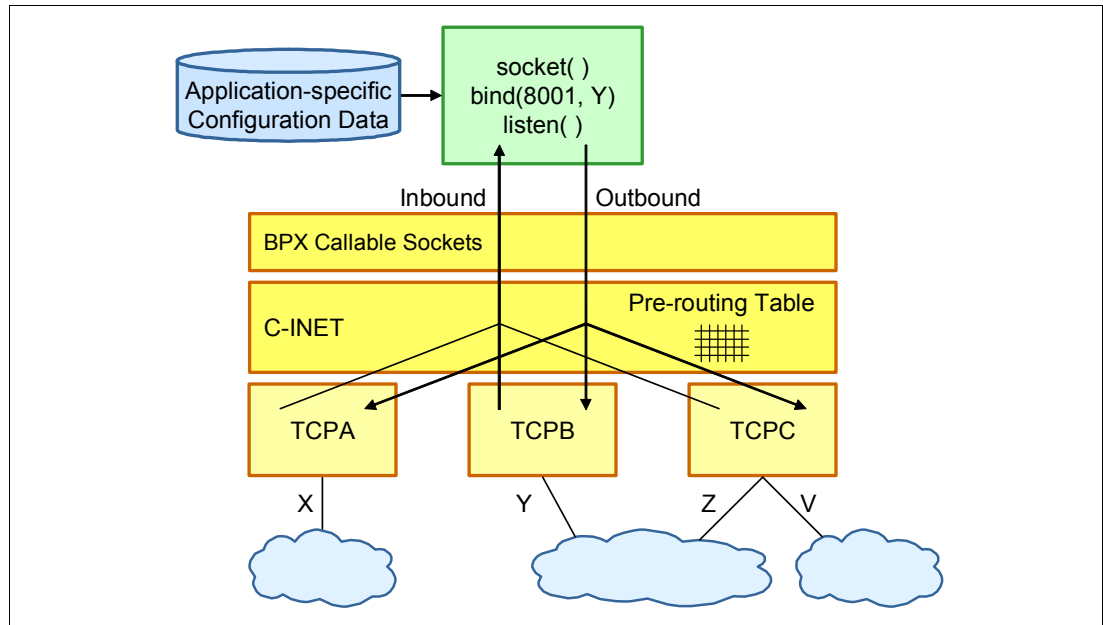


Figure 2-9 APIs, generic bind

Outbound connections or UDP datagrams are processed by the C-INET pre-router, and the stack with the best route to the destination is chosen.

When using a generic bind, the server port number must be reserved in all stacks. If one stack has it reserved to another address space, the `bind()` call fails.

2.2.8 Resolving an IPv6 address

IPv6 support introduces several changes to how host name and IP address resolution is performed. These changes affect several areas of resolver processing, including:

- ▶ Resolver APIs were introduced for IPv6-enabled applications.
- ▶ An algorithm is defined to describe how a resolver needs to sort a list of IP addresses returned for a multihomed host.
- ▶ DNS resource records are defined to represent hosts with IPv6 addresses, and therefore network flows between resolvers and name servers (instead of DNS IPv4 A records).

Resolver support for IPv6 connections to DNS name servers

Communications Server for z/OS IP allows the system resolver to send requests to the Domain Name System (DNS) name servers using IPv6 communication. To implement IPv6 communication with a DNS name server, specify the IPv6 address of the server on the existing `NSINTERADDR` and `NAMESERVER` resolver configuration statements in the `TCPIP.DATA` data set.

Restriction: The `res_state` structure (`nsaddr_list`) contains only the IPv4 addresses coded on the `NSINTERADDR` or `NAMESERVER` statements. Applications that examine or update the `nsaddr_list` cannot manipulate the IPv6 addresses.

IPv6 resolver statements

ETC.IPNODES is a local host file (in the style of `/etc/hosts`), which can contain both IPv4 and IPv6 addresses. IPv6 addresses can only be defined in ETC.IPNODES. This file allows the administration of local host files to more closely resemble that of other TCP/IP platforms, and eliminates the requirement of post-processing the files (specifically, `MAKESITE`).

The IPv6 search order is the same as the `COMMONSEARCH` search order, as shown in Figure 2-4 on page 27. If you do not want to use the `COMMONSEARCH` search order for existing IPv4 local hosts files, you might need to maintain two different local host files (for example, IPv4 addresses in `HOSTS.LOCAL`, and IPv6 and IPv4 addresses in ETC.IPNODES).

Name and address resolution functions

The APIs such as `getaddrinfo`, `getnameinfo`, and `freeaddrinfo` allow applications to resolve host names to IP addresses and vice versa for IPv6. The APIs are designed to work with both IPv4 and IPv6 addressing. The use of these APIs should be considered if an application is being designed for eventual use in an IPv6 environment.

The manner in which host name (`getaddrinfo`) or IP address (`getnameinfo`) resolution is performed is dependent upon resolver specifications contained in the resolver `SETUP` data sets and `TCPIP.DATA` configuration data sets, just like IPv4 address resolution. These specifications determine whether the APIs will query a name server first and then search the local host files, or whether the order will be reversed—or even if one of the steps will be eliminated completely. The specifications also control whether local host files have to be searched, and which local host file will be accessed.

Default destination address selection

Resolver APIs have the capability to return multiple IP addresses as a result of a host name query. However, many applications only use the first address returned to attempt a connection or to send a UDP datagram. Therefore, the sorting of these IP addresses is performed by the default destination address selection algorithm.

Establishing connectivity might depend on whether an IPv6 address or an IPv4 address is selected, thus making this sorting function even more important. Default destination address selection only occurs when the system is enabled for IPv6 and the application is using the `getaddrinfo()` API to retrieve IPv6 or IPv4 addresses.

The default destination address selection algorithm takes a list of destination addresses and sorts them to generate a new list. The algorithm sorts together both IPv6 and IPv4 addresses by a set of rules.

The following rules are applied, in order, to the first and second address, choosing a best address. Rules are then applied to this best address and the third address. This process continues until rules are applied to the entire list of addresses.

Rule 1 Avoid unusable destinations. If one address is reachable (the stack has a route to the particular address) and the other is unreachable, then place the reachable destination address *prior to* the unreachable address.

- Rule 2** Prefer matching scope. If the scope of one address matches the scope of its source address and the other address does not meet this criteria, then the address with the matching scope is placed *before* the other destination address.
- Rule 3** Avoid deprecated addresses. If one address is deprecated and the other is non-deprecated, then the non-deprecated address is placed *prior* to the other address.

Terminology: Deprecated, in this context, means that the state of an IPv6 address has changed from *preferred* state (the address was leased to an interface for a fixed, possibly infinite, length of time) to *deprecated* state. (When a lifetime expires, the binding and address can become invalid, and the address can be reassigned to another interface elsewhere on the Internet.) While in a deprecated state, the use of an address is discouraged but not strictly forbidden.

- Rule 4** Prefer matching address formats. If one address format matches its associated source address format and the other destination does not meet this criteria, then place the destination with the matching format *prior* to the other address.
- Rule 5** Prefer higher precedence. If the precedence of one address is higher than the precedence of the other address, then the address with the higher precedence is placed *before* the other destination address.
- Rule 6** Use the longest matching prefix. If one destination address has a longer CommonPrefixLength with its associated source address than the other destination address has with its source address, then the address with the longer CommonPrefixLength is placed *before* the other address.
- Rule 7** Leave the order unchanged. No rule selected a better address of these two; they are equally good. Choose the first address as the better address of these two and the order is not changed.

2.2.9 Resolver support for EDNS0

An early implementation of DNS, which is discussed in RFC 1035, allows only a maximum of 512 bytes for any DNS packet sent through UDP. This limitation inhibits DNS performance because, when a DNS server or client needs to communicate with a large amount of data, it will have to use the bulky TCP protocol (higher performance cost) instead of the simple UDP protocol (lower performance cost).

Extension Mechanism for DNS (EDNS0) was introduced in RFC 2671 to address the performance improvement limitation imposed by the traditional DNS implementation. The IBM implementation of the EDNS0 standard allows DNS communication of up to 3072 bytes using UDP. This implementation improves DNS's ability to communicate a large amount of data, such as IP version 6 (IPv6).

z/OS Communications Server resolver supports the EDNS0 standard by default. No additional steps are needed to enable this feature. However, the following dependencies are required for the resolver to support EDNS0 properly:

- ▶ The DNS name server must also support EDNS0 protocols in order to use UDP packets larger than 512 bytes.
- ▶ Firewalls that exist between the DNS name server and the z/OS resolver must be configured to accept DNS messages sent as UDP packets of greater than 512 bytes in order to use EDNS0 protocols.

In rare situations where the DNS server was just upgraded to support EDNS0, a refresh of the z/OS resolver is required so that it can relearn the DNS server EDNS0 capabilities. Issue `MODIFY RESOLVER,REFRESH` to the resolver address space to refresh.

2.2.10 Considerations

To implement the resolver address space, it is important to first determine whether your environment requires a single TCP/IP stack or multiple TCP/IP stacks. In both cases the resolver is an independent address space and has to be up and running before the TCP/IP stack is started.

The statements defined in the global `TCPIP.DATA` cannot be overridden by the local `TCPIP.DATA` file of the each TCP/IP stack. The local `TCPIP.DATA` file can only specify the statement if it is not already defined in the global `TCPIP.DATA` file.

Important: In some resolver environments, the use of the trace functions (such as `SockDebug` or `TraceResolver`) might affect performance. Therefore, we recommend using the method that we describe in “`CTTRACE: RESOLVER (SYSTCPRE)`” on page 55.

The resolver in a single stack environment

We recommend that you create a global `TCPIP.DATA` file for a single stack environment. The `TCPIPJobname` statement can be coded in a global `TCPIP.DATA` file or in the local (non-global) `TCPIP.DATA` file, because there is only one stack on the system. If some applications have requirements to specify their own `TCPIP.DATA` statements, then omit them from the global `TCPIP.DATA` so the applications can point to the local `TCPIP.DATA` file to be used.

The resolver in a multiple stack environment

When implementing for a multiple stack environment, each TCP/IP stack should use a local `TCPIP.DATA` file specifying stack-specific statements, such as `TCPIPJobname` and `HostName`. Optionally, you can merge some statements that can be applied to all TCP/IP stacks and all TCP/IP applications to a global `TCPIP.DATA` file. You need to determine which statements should be defined in the global `TCPIP.DATA` and used in the entire z/OS image. This will depend on how much you want to allow each stack or application to define its own definitions.

In the multiple stack environment, we recommend that you create a global `TCPIP.DATA` if all the statements needed in the global `TCPIP.DATA` (see “Using `GLOBALTCPIPDATA`” on page 24) can be applied to all the stacks as shown in Figure 2-3 on page 26. If not, do not use the global `TCPIP.DATA` and only use local `TCPIP.DATA` for each stack.

Figure 2-10 depicts the multiple stack environment without the use of a global TCPIP.DATA.

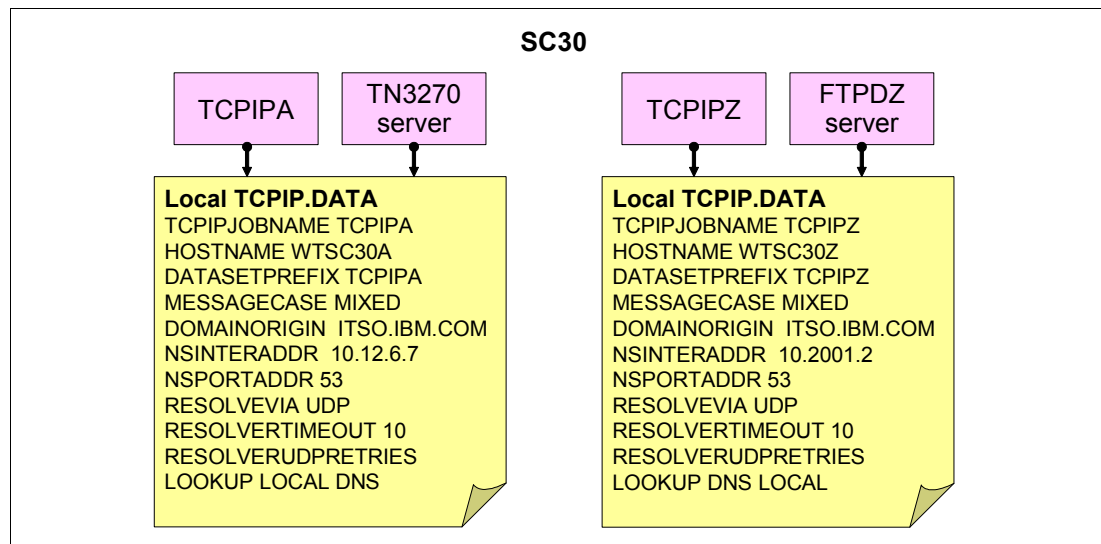


Figure 2-10 The multiple stack environment without global TCPIP.DATA

Recommendation: Although there are specialized cases where multiple stacks per LPAR can provide value, generally we recommend implementing only one TCP/IP stack per LPAR. The reasons for this recommendation are as follows:

- ▶ A TCP/IP stack is capable of exploiting all available resources defined to the LPAR in which it is running. Therefore, starting multiple stacks will not yield any increase in throughput.
- ▶ When running multiple TCP/IP stacks, additional system resources, such as memory, CPU cycles, and storage, are required.
- ▶ Multiple TCP/IP stacks add a significant level of complexity to TCP/IP system administration tasks.
- ▶ It is not necessary to start multiple stacks to support multiple instances of an application on a given port number, such as a test HTTP server on port 80 and a production HTTP server also on port 80. This type of support can instead be implemented using BIND-specific support where the two HTTP server instances are each associated with port 80 with their own IP address, using the BIND option on the PORT reservation statement.

One example where multiple stacks can have value is when an LPAR needs to be connected to multiple isolated security zones in such a way that there is no network level connectivity between the security zones. In this case, a TCP/IP stack per security zone can be used to provide that level of isolation, without any network connectivity between the stacks.

2.3 Implementing the resolver

In this scenario, we use the customized resolver address space and specify GLOBALTCPIPDATA, DEFAULTTCPIPDATA, and GLOBALIPNODES in the resolver SETUP data set. We define a global TCPIP.DATA file and define a common set of parameters for entire z/OS image. We omit some statements in the global TCPIP.DATA file so that the applications or TCP/IP stack can use their own local TCPIP.DATA file for the statements undefined in the global TCPIP.DATA file.

Figure 2-11 depicts the environment that we use for this implementation.

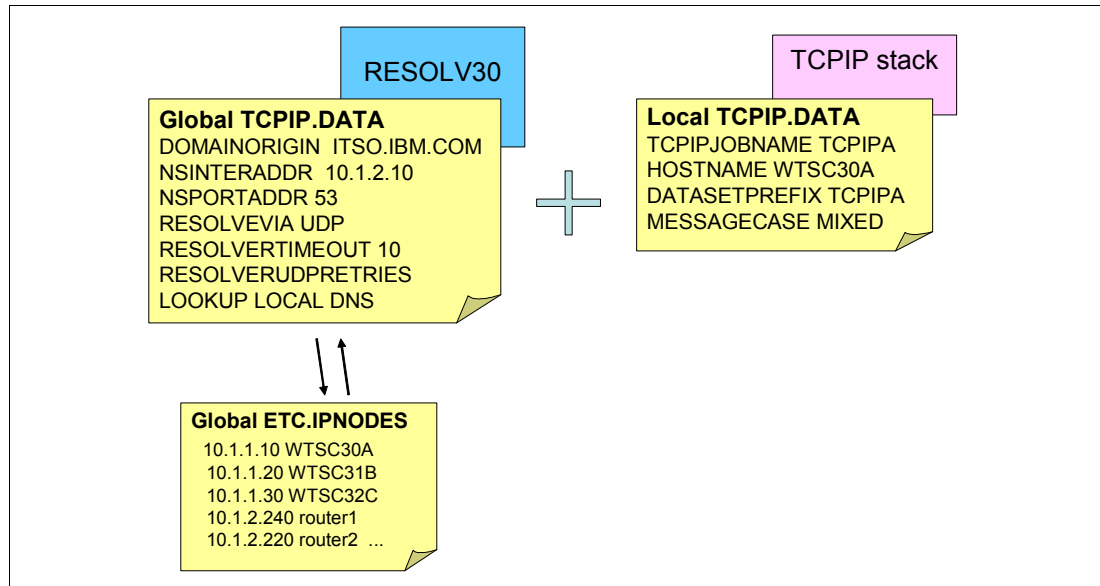


Figure 2-11 Our resolver environment on SC30

2.3.1 Implementation tasks

To implement the resolver address space in our test environment, we performed these steps:

1. Set up the resolver started procedure.
2. Customize BPXPRMxx.
3. Configure the resolver SETUP data set.
4. Create the global TCPIP.DATA file.
5. Create the default TCPIP.DATA file.
6. Create the global IPNODES data set.
7. Create a TCPIP.DATA file for TCPIPA stack.
8. Create the TCPIPA stack started procedure.

We describe these steps in the sections that follow.

Set up the resolver started procedure

We created the resolver procedure so that it would start during the UNIX System Services initialization.

To create the procedure, we copied the sample procedure hlq.SEZAINST(EZBREPRC) and customized it to our environment, as shown in Example 2-9 on page 43. The procedure has only one DD card that must be configured, the SETUP DD card **1**, which describes where the SETUP data set is located.

Example 2-9 The resolver started procedure

```
/* *****  
/* SYS1.PROCLIB(RESOLV30)  
/* *****  
//RESOLV30 PROC PARMS='CTRACE(CTIRES00)' 1  
//EZBREINI EXEC PGM=EZBREINI,REGION=OM,TIME=1440,PARM=&PARMS  
/* SETUP contains resolver setup parameters.  
/* See the section on "Understanding Resolvers" in the  
/* IP Configuration Guide for more information. A sample of  
/* resolver setup parameters is included in member RESSETUP  
/* of the SEZAINST data set.  
/*  
//SETUP DD DSN=TCPIPA.TCPPARMS(RESOLV&SYSCONE),DISP=SHR,FREE=CLOSE 2
```

In this example, the numbers correspond to the following information:

1. The name of the resolver procedure is RESOLV30.
2. Specifies the resolver SETUP data set. The &SYSCONE MVS system symbol value on this system is 30.

Customize BPXPRMxx

We customized the RESOLVER_PROC statement in BPXPRMxx, to specify the procedure name that we used, which causes the resolver to start automatically the next time z/OS UNIX System Services initializes. Example 2-10 shows the partial contents of BPXPRMxx.

Example 2-10 Specifying the resolver procedure to be started

```
/* *****  
/* SYS1.PARMLIB(BPXPRM00)  
/* *****  
/* RESOLVER_PROC is used to specify how the resolver address space */  
/* is processed during Unix System Services initialization.          */  
/* The resolver address space is used by Tcp/Ip applications        */  
/* for name-to-address or address-to-name resolution.              */  
/* In order to create a resolver address space, a system must be    */  
/* configured with an AF_INET or AF_INET6 domain.                  */  
/* RESOLVER_PROC(procname|DEFAULT|NONE)                             */  
/*   procname - The name of the address space for the resolver.     */  
/*               In this case, this is the name of the address      */  
/*               space as well as the procedure member name         */  
/*               in SYS1.PROCLIB. procname is 1 to 8 characters     */  
/*               long.                                              */  
/*   DEFAULT - An address space with the name RESOLVER will         */  
/*               be started. This is the same result that will     */  
/*               occur if the RESOLVER_PROC statement is not        */  
/*               specified in the BPXPRMxx profile.                 */  
/*               @DAA*/  
/*   NONE - Specifies that a RESOLVER address space is             */  
/*               not to be started.                                  */  
/* *****  
RESOLVER_PROC(RESOLV&SYSCONE.) 1
```

In this example, the numbers correspond to the following information:

1. Specifies the name of the resolver procedure we created in the previous step. The &SYSCONE MVS system symbol value on this system is 30.

Important: When the resolver is started by UNIX System Services, you must pay attention to the following information:

- ▶ The resolver address space is started by SUB=MSTR. This means that JES services are not available to the resolver address space. Therefore, no DD cards with SYSOUT can be used.
- ▶ The resolver start procedure needs to reside in a data set that is specified by the MSTJCLxx PARMLIB member's IEFPSI DD card specification. Otherwise, the procedure will not be found and the resolver will not start. SYS1.PROCLIB is usually one of the libraries specified there.

Configure the resolver SETUP data set

We configured the resolver SETUP data set which is specified with the SETUP DD definition in the resolver started procedure. This data set defines the location of the global and default TCPIP.DATA files containing the parameters we wanted to be defined in the z/OS environment.

In our test environment, we copied the SETUP sample data set and changed its contents to meet our requirements, as shown in Example 2-11.

Example 2-11 Resolver address space SETUP data set

```
; *****
; TCPIPA.TCPPARMS(RESOLV30)
; *****
GLOBALTCPIPDATA('TCPIPA.TCPPARMS(GLOBAL)') 1
DEFAULTTCPIPDATA('TCPIPA.TCPPARMS(DEFAULT)') 2
GLOBALIPNODES('TCPIPA.TCPPARMS(IPNODES)') 3
COMMONSEARCH 4
CACHE 5
CACHESIZE(20M) 6
MAXTTL(600) 7
UNRESPONSIVETHRESHOLD(25) 8
```

In this example, the numbers correspond to the following information:

1. Specifies the first choice of the TCPIP.DATA file.
2. Specifies the last choice of the TCPIP.DATA file.
3. Specifies the first choice of the local hosts file.
4. The COMMONSEARCH search order is used. This statement is needed to have GLOBALIPNODES to be applied.
5. Indicates that system-wide caching is enabled for the resolver.
6. Specifies the maximum amount of storage that can be allocated by the resolver to manage cached records. A value of at least 50M should be considered in a production environment.
7. Specifies the maximum amount of time the resolver can use resource information obtained from a DNS server as part of resource resolution.

- 8.** Define the percentage threshold value to be used to calculate when a name server should be declared to be unresponsive to resolver queries.

Important: Be careful when creating these global parameters. The definitions in the resolver SETUP data set is applied to all TCP/IP stacks or applications.

Create the global TCPIP.DATA file

In this step, we provide the global statements that all TCP/IP stacks and applications used in our z/OS environment. To define these statements, we copied the sample TCPIP.DATA file provided in hlq.SEZAINST(TCPDATA) and customized the statements, as shown in Example 2-12.

Example 2-12 Global TCPIP.DATA file

```
; *****
; TCPIPA.TCPPARMS (GLOBAL)
; *****
DOMAINORIGIN  ITS0.IBM.COM      1
NSINTERADDR  10.1.2.10          2
NSPORTADDR   53
RESOLVEVIA   UDP
RESOLVERTIMEOUT 10
RESOLVERUDPRETRIES 1
LOOKUP       DNS LOCAL          3
```

In this example, the numbers correspond to the following information:

- 1.** Specifies the list of domain names appended to the host name when the search is performed.
- 2.** Specifies the IP address of the DNS server.
- 3.** To take advantage of the caching that we enabled in the Global Resolver SETUP file (Example 2-11 on page 44), we changed our previously used LOOKUP sequence to favor the DNS over the LOCAL file. If neither a cache entry or a DNS entry is found for a lookup, then the resolver searches the local file.

Important: If GLOBALTCPIPDATA is specified:

- ▶ Any statements contained in the global TCPIP.DATA file will take precedence over any statements in local TCPIP.DATA file found by way of the appropriate environment's (Native z/OS or z/OS UNIX) search order.
- ▶ The TCPIP.DATA statements in Example 2-12 (with the exception of LOOKUP) can only be specified in GLOBALTCPIPDATA. If the resolver statements are found in any of the other search locations for TCPIP.DATA, they are ignored. If the resolver statements are not found in GLOBALTCPIPDATA, their default value will be used.

Create the default TCPIP.DATA file

We created a default TCPIP.DATA file, as shown in Example 2-13, to be the last choice of the local TCPIP.DATA search order. It is used when the application does not specify the local TCPIP.DATA explicitly.

Example 2-13 Default TCPIP.DATA file

```
; *****  
; TCPIPA.TCPPARMS(DEFAULT)  
; *****  
TCPIPJOBNAME TCPIP 1  
HOSTNAME WTSC30 2
```

In this example, the numbers correspond to the following information:

- 1. Specifies the default TCP/IP procedure name.
- 2. Specifies the default host name.

Important: Applications that use Language Environment services without a TCPIPJOBNAME statement cause applications that issue `__iptcpn()` to receive a job name of NULL, and some of these applications will use INET instead of TCP/IP. Although this presents no problem when running in a single-stack environment, it can potentially cause errors in a multi-stack environment.

Create the global IPNODES data set

We created the global IPNODES data set, which is referred as GLOBALIPNODES in the resolver SETUP data set. It contains name-to-address mappings. This data set is used for the local search to resolve a name into an IP address or vice versa.

We chose to use the COMMONSEARCH, because it allowed us to have a common local search environment with IPv4 or IPv6 hosts. Example 2-14 shows the contents of the GLOBALIPNODES data set. When using COMMONSEARCH, only the IPNODES data set is used.

Example 2-14 GLOBALIPNODES data set

```
; *****  
; TCPIPA.TCPPARMS(IPNODES)  
; *****  
10.1.2.10 OURDNS 1  
10.1.1.10 WTSC30A 1  
10.1.1.20 WTSC31B 1  
10.1.1.30 WTSC32C 1  
10.1.2.240 router1 1  
10.1.3.240 router2 1  
1::2 TESTIPV6ADDRESS1 2  
1:2:3:4:5:6:7:8 TESTIPV6ADDRESS2 2
```

In this example, the numbers correspond to the following information:

- 1. The mapping of a name and a IPv4 address is listed.
- 2. The mapping of a name and a IPv6 address is listed.

Create a TCPIP.DATA file for TCPIPA stack

We created a local TCPIP.DATA file for the TCPIPA stack with file name DATAA30, as shown in Example 2-15.

Example 2-15 TCPIP.DATA file DATAA30

```
; *****
; TCPIPA.TCPPARMS(DATAA30)
; *****
TCPIPJOBNAME TCPIPA           1
HOSTNAME WTSC30A              2
DATASETPREFIX TCPIPA
MESSAGECASE MIXED
```

In this example, the numbers correspond to the following information:

1. Specifies the procedure name of TCPIPA stack.
2. Specifies the host name of the TCPIPA stack.

Create the TCPIPA stack started procedure

We created the TCPIPA stack procedure (RESOLVER_CONFIG) and pointed to TCPIPA.TCPPARMS(DATAA30), using the &sysclone variable to simplify our implementation to allow for a single procedure to be used by any z/OS image in our sysplex environment, as shown in Example 2-16.

Example 2-16 TCPIPA procedure

```
/*****
/* SYS1.PROCLIB(TCPIPA)
/*****
//TCPIPA   PROC  PARM='CTTRACE(CTIEZB00),IDS=00', 1
//          PROFILE=PROFA&SYSCONE.,TCPDATA=DATAA&SYSCONE
//TCPIPA   EXEC  PGM=EZBTCPIP,REGION=OM,TIME=1440,
//          PARM=('&PARMS',
//          'ENVAR("RESOLVER_CONFIG=/' TCPIPA.TCPPARMS(&TCPDATA)'")' ) 2
//SYSPRINT DD SYSOUT=*,DCB=(RECFM=VB,LRECL=132,BLKSIZE=136)
//SYSTCPT DD SYSOUT=*,DCB=(RECFM=VB,LRECL=132,BLKSIZE=136)
//ALGPRINT DD SYSOUT=*,DCB=(RECFM=VB,LRECL=132,BLKSIZE=136)
//CFGPRINT DD SYSOUT=*,DCB=(RECFM=VB,LRECL=132,BLKSIZE=136)
//SYSOUT   DD SYSOUT=*,DCB=(RECFM=VB,LRECL=132,BLKSIZE=136)
//CEEDUMP  DD SYSOUT=*,DCB=(RECFM=VB,LRECL=132,BLKSIZE=136)
//SYSERROR DD SYSOUT=*
//PROFILE  DD DISP=SHR,DSN=TCPIPA.TCPPARMS(&PROFILE.) 3
```

In this example, the numbers correspond to the following information:

1. The TCP/IP procedure name is TCPIPA.
2. The local TCPIP.DATA is specified.
3. The TCP/IP profile is specified (the TCP/IP configuration file example is not shown in this chapter).

2.3.2 Activation and verification

To verify that the resolver address space was working as expected, we performed these steps:

1. Stop the default z/OS UNIX resolver.
2. Start the resolver address space.
3. Display the resolver address space configuration.
4. Use the **ping** command to verify the name resolution.

To implement our resolver address space, we halt the running resolver using the STOP command, as shown in Example 2-17.

Important: Stop and restart the resolver only if you install a new level of the resolver code.

Stop the default z/OS UNIX resolver

In our current environment, the default z/OS UNIX resolver was running. We stopped this default resolver, as shown in Example 2-17, to run the customized resolver.

Example 2-17 Stopping the resolver address space

```
P RESOLV30
EZZ9292I RESOLVER ENDING
IEF352I ADDRESS SPACE UNAVAILABLE
$HASP395 RESOLV30 ENDED
```

Start the resolver address space

We started the customized resolver address space using the procedure we created in the previous step, as shown in Example 2-18.

Example 2-18 Starting a configured resolver address space

```
S RESOLV30
IRR812I PROFILE ** (G) IN THE STARTED CLASS WAS USED 646
      TO START RESOLV30 WITH JOBNAME RESOLV30.
$HASP100 RESOLV30 ON STCINRDR
IEF695I START RESOLV30 WITH JOBNAME RESOLV30 IS ASSIGNED TO USER
IBMUSER , GROUP SYS1
$HASP373 RESOLV30 STARTED
IEE252I MEMBER CTIRES00 FOUND IN SYS1.PARMLIB
EZZ9298I DEFAULTTCPIPDATA - TCPIPA.TCPPARMS(DEFAULT) 1
EZZ9298I GLOBALTCPIPDATA - TCPIPA.TCPPARMS(GLOBAL) 2
EZZ9298I DEFAULTIPNODES - None
EZZ9298I GLOBALIPNODES - TCPIPA.TCPPARMS(IPNODES) 3
EZZ9304I COMMONSEARCH
EZZ9304I CACHE 4
EZZ9298I CACHESIZE - 20M 5
EZZ9298I MAXTTL - 600 6
EZZ9298I UNRESPONSIVETHRESHOLD - 25 7
EZZ9291I RESOLVER INITIALIZATION COMPLETE
```

In this example, the numbers correspond to the following information:

1. The correct DEFAULTTCPIPDATA file is applied.
2. The correct GLOBALTCPIPDATA file is applied.

3. The correct GLOBALIPNODES file is applied.
4. Indicates that system-wide caching is enable for the resolver.
5. Indicates the maximum amount of storage that can be allocated by the resolver.
6. Indicates the maximum amount of time the resolver can use resource information obtained from a Domain Name System (DNS) server as part of resource resolution.
7. Indicates the percentage threshold value to calculate when a name server should be declared to be unresponsive to resolver queries.

Note: If you want to start the default z/OS UNIX resolver, use the following command instead:

```
START IEESYSAS.RESOLVER,PROG=EZBREINI,SUB=MSTR
```

Note: The resolver utilizes non-reusable address spaces. To start resolver using a reusable address space ID (REUSASID), see 1.3.3, “Reusable address space ID” on page 6.

If you want to reload the SETUP data set content changes, use the MODIFY command to refresh the resolver. To show how this is done, we created a new SETUP data set named NEWSETUP, with the same configuration as the RESOLV30 setup file and changed UNRESPONSIVETHRESHOLD statement changed to 35%, and refreshed the resolver to reflect the changes, as shown in Example 2-19.

Example 2-19 Modifying the resolver address space

```
F RESOLV30,REFRESH,SETUP=TCPIPA.TCPPARMS(NEWSETUP)
EZZ9298I DEFAULTTCPIPDATA - TCPIPA.TCPPARMS(DEFAULT)
EZZ9298I GLOBALTCPIPDATA - TCPIPA.TCPPARMS(GLOBAL)
EZZ9298I DEFAULTIPNODES - None
EZZ9298I GLOBALIPNODES - TCPIPA.TCPPARMS(IPNODES)
EZZ9304I COMMONSEARCH
EZZ9304I CACHE
EZZ9298I CACHESIZE - 30M
EZZ9298I MAXTTL - 600
EZZ9298I UNRESPONSIVETHRESHOLD - 35
EZZ9293I REFRESH COMMAND PROCESSED
```

Display the resolver address space configuration

To verify that the correct configuration file is applied to the resolver address space, use the MODIFY command with the display option, as shown in Example 2-20.

Example 2-20 Modify resolver with display option

```
F RESOLV30,DISPLAY
EZZ9298I DEFAULTTCPIPDATA - TCPIPA.TCPPARMS(DEFAULT)
EZZ9298I GLOBALTCPIPDATA - TCPIPA.TCPPARMS(GLOBAL)
EZZ9298I DEFAULTIPNODES - None
EZZ9298I GLOBALIPNODES - TCPIPA.TCPPARMS(IPNODES)
EZZ9304I COMMONSEARCH
EZZ9304I CACHE
EZZ9298I CACHESIZE - 30M
EZZ9298I MAXTTL - 600
```

EZZ9298I UNRESPONSIVETHRESHOLD - 35
EZZ9293I DISPLAY COMMAND PROCESSED

Use the ping command to verify the name resolution

Verify that the resolver is able to perform the expected name-to-address resolution by using the **ping** command, as shown in Example 2-21. As you can see, the name `router1` has resolved to address `10.1.2.240`. Refer to 8.3.1, “The ping command (TSO or z/OS UNIX)” on page 303 for more details about issuing the **ping** command.

Example 2-21 UNIX ping command display

```
CS02 @ SC30:/u/cs02>ping router1
CS V1R12: Pinging host router1 (10.1.2.240)
Ping #1 response took 0.001 seconds.
```

The TSO PING command was also successful, as shown in Example 2-22.

Example 2-22 TSO PING command display

```
TSO PING ROUTER1
CS V1R12: Pinging host router1 (10.1.2.240)
Ping #1 response took 0.001 seconds.
***
```

It is also possible to verify where the resolver is looking by using the `TRACE RESOLVER` parameter in the stack's or application's `TCPIP.DATA` file. For an explanation of how this is done and what the contents of this trace will be, refer to 2.4, “Problem determination” on page 50.

2.4 Problem determination

To diagnose resolver problems, you can use two kinds of trace tools:

- ▶ Trace Resolver

This provides information that can be helpful in debugging problems that an application program could have with using resolver facilities (for example, `GetHostByName` or `GetHostByAddr`).

- ▶ Component Trace RESOLVER (SYSTCPRE)

This is useful for diagnosing resolver problems that cannot be isolated to one particular application.

In this section we provide a brief explanation of when to debug, which trace has to be used, and how to use these trace facilities. For more information about resolver diagnosis, refer to *z/OS Communications Server: IP Diagnosis Guide*, GC31-8782.

Deciding which tool to use to diagnose a resolver problem

The first thing to do when diagnosing a possible resolver problem is to check the symptoms to verify if it is indeed a resolver problem (see Table 2-3).

Table 2-3 What to do if the host name cannot be reached

When we ping a host name, the ping command:	What is the problem?	Solution
Succeeds, but another application fails when resolving the same host name.	The problem is with the resolver configuration for the application in the users environment.	Use the Trace Resolver statement on the local TCPIP.DATA used by the application that has the problem.
Fails, but the host name is converted to an IP address.	The resolution is successful but the host is not reachable or active.	This problem is related to connectivity, not a resolver problem.
Fails to convert the name to an IP address.	The problem might be with the resolver configuration, searching local host files, or using DNS.	Use Trace Resolver to solve the problem.

Tip: If the problem seems to be related to the DNS, use the LOOKUP LOCAL DNS statement to check the local files first.

Trace Resolver

The Trace Resolver informs us what the resolver looked for (the questions) and where it looked (name server's cache and IP addresses or local host file names).

The following situations can be checked in the trace output:

- ▶ Check whether the correct resolver data sets is in use. If an unexpected TCPIP.DATA file is used, check the search orders of the data set.
- ▶ Check whether the data sets defined to be used are authorized by RACF and can be read by the application, TCP/IP stack, or user.
- ▶ Check the TCPIP.DATA parameter values, especially Search, NameServer, NSINTERADDR, and NsPortAddr.
- ▶ Check the questions posed by the resolver to DNS or in searching the local host files. Are these the queries you expected?
- ▶ Look for errors or failures in the trace.
- ▶ Was the information obtained from the resolver cache? If so, use **netstat Rescache/-Q** commands to determine if the cache information is correct.
- ▶ Did DNS respond (if you expected it to)? If not, see whether DNS is active at the IP address you specified for NameServer and NSINTERADDR and what port it is listening on. Also, DNS logs can be helpful, so ask the DNS administrator for help.

Trace Resolver can be activated in the following ways, in its precedence order:

1. The RESOLVER_TRACE environment variable (z/OS UNIX environment only).
2. SYSTCPT DD allocation.
3. TRACE RESOLVER or OPTIONS DEBUG statements (you must allocate STDOUT or SYSPRINT to generate trace data).

4. The resDebug bit set to on in the _res structure option field (you must allocate STDOUT or SYSPRINT to generate trace data).

Next, we illustrate using Trace Resolver in a z/OS UNIX environment, and in a TSO environment.

Using Trace Resolver in z/OS UNIX environment

Example 2-23 shows how to enable and disable the Trace Resolver in z/OS UNIX environment.

Example 2-23 Using Trace Resolver in a z/OS UNIX environment

```
CS02 @ SC30:/u/cs02>export RESOLVER_TRACE=/u/cs02/trace1.txt 1
CS02 @ SC30:/u/cs02>ping admin 2
CS V1R12: Pinging host admin.ITS0.IBM.COM (10.1.4.11)
Ping #1 response took 0.000 seconds.
CS02 @ SC30:/u/cs02>set -A RESOLVER_TRACE 3
CS01 @ SC30:/u/cs01>obrowse /u/cs02/trace1.txt
```

In this example, the numbers correspond to the following information:

1. To enable the Trace Resolver, set the RESOLVER_TRACE environment variable. This command directs the output to the /u/cs06/trace1.txt HFS file. You can also direct the output to STDOUT by specifying RESOLVER_TRACE=STDOUT. If you want to direct it to a new MVS data set, specify the following command:

```
RESOLVER_TRACE="//'SOME.MVS.DATASET' "
```

2. After enabling a Trace Resolver, perform a z/OS UNIX shell command that invokes a resolver call.
3. This command disables the Trace Resolver.

Using Trace Resolver in a TSO environment with SYSTCPT DD

Example 2-24 shows how to enable and disable the Trace Resolver in a TSO environment environment.

Example 2-24 Using Trace Resolver in a TSO environment

```
alloc dd(systcpt) da(*) 1
ping router1 2
free dd(systcpt) 3
```

In this example, the numbers correspond to the following information:

1. To enable the Trace Resolver, allocate a SYSTCPT data set. If you specify da(*), the Trace Resolver output to a TSO terminal. If you want to direct the output to a specific data set, specify da('SOME.DATASET.NAME').
2. After enabling the Trace Resolver, perform a TSO command that invokes a resolver call.
3. To disable the Trace Resolver, free a SYSTCPT data set.

Tip: When directing Trace Resolver output to a TSO terminal, define the screen size to be only 80 columns wide. Otherwise, trace output is difficult to read.

Using Trace Resolver for applications with TCPIP.DATA statements

Allocate STDOUT or SYSPRINT (as a DD statement in the procedure) as an output data set, and define the statement TRACE RESOLVER or OPTIONS DEBUG in the first line of the TCPIP.DATA file that is being used by the application, as shown in Example 2-25. Start the application that invokes a resolver call.

Example 2-25 Using the OPTIONS DEBUG to get a trace of the resolver

```
OPTIONS DEBUG 1
TCPIPJOBNAME TCPIPA
HOSTNAME WTSC30A
DOMAINORIGIN ITS0.IBM.COM
DATASETPREFIX TCPIPA
MESSAGECASE MIXED
NSINTERADDR 10.1.2.10
NSPORTADDR 53
```

In this example, the numbers correspond to the following information:

- 1.** Specify OPTIONS DEBUG or TRACE RESOLVER to enable Trace Resolver.

Displaying the output of the Trace Resolver

Example 2-26 shows the output of the Trace Resolver in the z/OS UNIX environment (which was taken from Example 2-23 on page 52). Note that the Trace Resolver taken in the TSO environment (Example 2-24 on page 52) is almost identical.

Example 2-26 Trace Resolver partial output: z/OS UNIX shell environment

```
Resolver Trace Initialization Complete -> 2010/09/27 15:04:49.709930
res_init Resolver values:
  Global Tcp/Ip Dataset = TCPIPA.TCPPARMS(GLOBAL) 1
  Default Tcp/Ip Dataset = TCPIPA.TCPPARMS(DEFAULT)
  Local Tcp/Ip Dataset = /etc/resolv.conf 2
...
...
  (G) LookUp          = LOCAL DNS 3
  (*) Cache
res_init Succeeded
res_init Started: 2010/09/27 15:04:49.741620
res_init Ended: 2010/09/27 15:04:49.741624
*****
GetAddrInfo Started: 2010/09/27 15:04:49.741646
GetAddrInfo Invoked with following inputs:
  Host Name: admin 4
...
...
GetAddrInfo Only IPv4 Interfaces Exist
GetAddrInfo Searching Local Tables for IPv4 Address
Global IpNodes Dataset = TCPIPA.TCPPARMS(IPNODES) 5
Default IpNodes Dataset = None
Search order           = CommonSearch
...
...
  - Lookup for admin.ITS0.IBM.COM
  - Lookup for admin
res_search(admin, C_IN, T_A)
res_search Host Alias Search found no alias 6
res_querydomain(admin, ITS0.IBM.COM, C_IN, T_A)
res_querydomain resolving name: admin.ITS0.IBM.COM
```

```

res_query(admin.ITSO.IBM.COM, C_IN, T_A)
  Querying resolver cache for admin.ITSO.IBM.COM
...
  No cache information was available 7
res_mkquery(QUERY, admin.ITSO.IBM.COM, C_IN, T_A)
res_mkquery created message:
...
res_send Name Server Capabilities
  Name server 10.1.2.10
...
res_send Sending query to Name Server 10.1.2.10 8
DNS Communication Started: 2010/09/27 15:04:49.752519
  No OPT RR record sent on request to 10.1.2.10
...
BPX1AIO RECVMSG : From 10.1.2.10
UDP Data Length: 86
res_send received data via UDP. Message received:
* * * * * Beginning of Message * * * * *
  Query Id:                62855
...
  Response Code:            NOERROR
  Number of Question RRs:  1
  Question 1:
  admin.ITSO.IBM.COM
...
  Answer 1:
  admin.ITSO.IBM.COM
  Type (0X0001) T_A Class (0X0001) C_IN
  TTL: 86400 (1 days, 0 hours, 0 minutes, 0 seconds)
  10.1.4.11
* * * * * End of Message * * * * *
DNS Communication Ended: 2010/09/27 15:04:49.753095 time used 00:00:00.000576
Name Server Capability Updates
  Name server 10.1.2.10
    Queries sent   = 1
    Failures       = 0
res_send Succeeded
  Attempting to cache results for admin.ITSO.IBM.COM
  EZBRECAB: RetVal = 0, RC = 0, Reason = 0x00000000
  Cache information was saved 9
...
GetAddrInfo Succeeded: IP Address(es) found:
  IP Address(1) is 10.1.4.11 10
GetAddrInfo Ended: 2010/09/27 15:04:49.753194
*****
FreeAddrInfo Started: 2010/09/27 15:04:49.753222
FreeAddrInfo Called to free addrinfo structures
FreeAddrInfo Succeeded, Freed 1 Addrinfos
FreeAddrInfo Ended: 2010/09/27 15:04:49.753229
*****

```

In this example, the numbers correspond to the following information:

1. Informs you that the global TCPIP.DATA is in use.
2. Informs you that the local TCPIP.DATA is in use.
3. The local hosts file is looked up first, followed by the DNS server if it fails.
4. The admin host name is looked up.
5. Informs you that the global ETC.IPNODE is in use.

6. No information was available in ETC.IPNODE.
7. The admin host entry could not be found in the cache.
8. The resolver send a query to name server.
9. The response of name server is cached.
10. The IP Address was found in the name server.

CTRACE: RESOLVER (SYSTCPRE)

Component Trace (CTRACE) is used for the RESOLVER component (SYSTCPRE) to collect debug information. The TRACE RESOLVER traces information about a per-application basis and directs the output to a unique file for each application. The CTRACE shows resolver actions for all applications (although it might be filtered).

The CTRACE support allows for JOBNAME, ASID filtering, or both. The trace buffer is located in the resolver private storage. The trace buffer minimum size is 128 KB. The maximum size is 128 MB. The default size is 16 MB. Trace records can optionally be written to an external writer.

The resolver CTRACE can be started any time needed by using the TRACE CT command, or it can be activated during resolver procedure initialization.

Note: If you suspect that there is an error in the operation of the resolver cache, you must collect CTRACE records, as there are no Trace Resolver trace entries for cache processing.

Using CTRACE for RESOLVER

The resolver CTRACE initialization PARMLIB member can be specified at resolver start time. To activate the resolver CTRACE during resolver initialization, follow these steps:

1. Create a CTWTR procedure in your SYS1.PROCLIB, as shown in Example 2-27.

Example 2-27 CTWTR procedure

```
//CTWTR    PROC
//IEFPROC  EXEC PGM=ITTTRCWR
//TRCOUT01 DD  DSNAME=CS02.CTRACE1,VOL=SER=COMST2,UNIT=3390,
//           SPACE=(CYL,10),DISP=(NEW,KEEP),DSORG=PS
//TRCOUT02 DD  DSNAME=CS02.CTRACE2,VOL=SER=COMST2,UNIT=3390,
//           SPACE=(CYL,10),DISP=(NEW,KEEP),DSORG=PS
//*
```

2. Using the sample resolver procedure shipped with the product, enter the following console command:

```
S RESOLV30,PARMS='CTRACE(CTIRESxx)'
```

Where xx is the suffix of the CTIRESxx PARMLIB member to be used. To customize the parameters used to initialize the trace, you can update the SYS1.PARMLIB member CTIRES00, as shown in Example 2-28.

Example 2-28 Trace options

```
/*****
TRACEOPTS
/* ----- */
/*  Optionally start external writer in this file (use both      */
/*  WTRSTART and WTR with same wtr_procedure)                  */
```

```

                WTRSTART(CTWTR)
/* ----- */
/*   ON OR OFF: PICK 1                               */
/* ----- */
                ON
/*   OFF                                              */
/*   BUFSIZE: A VALUE IN RANGE 128K TO 128M          */
                BUFSIZE(16M)
/*   JOBNAME(jobname1,...)                          */
/*   ASID(Asid1,...)                                */
                WTR(CTWTR)
/* ----- */
/*   OPTIONS: NAMES OF FUNCTIONS TO BE TRACED, OR "ALL" */
/* ----- */
/*   OPTIONS(                                         */
/*       'ALL'                                         */
/*       , 'MINIMUM'                                   */
/*   )                                                 */

```

3. Use the TRACE CT command to define the options, as shown in Example 2-29.

Example 2-29 TRACE CT command flow

```

TRACE CT,ON,COMP=SYSTCPRE,SUB=(RESOLV30)
*189 ITT006A SPECIFY OPERAND(S) FOR TRACE CT COMMAND.
R 189,OPTIONS=(ALL),END
IEE600I REPLY TO 189 IS;OPTIONS=(ALL),END
ITT038I ALL OF THE TRANSACTIONS REQUESTED VIA THE TRACE CT COMMAND
WERE SUCCESSFULLY EXECUTED.
IEE839I ST=(ON,0256K,00512K) AS=ON BR=OFF EX=ON MT=(ON,024K) 497
        ISSUE DISPLAY TRACE CMD FOR SYSTEM AND COMPONENT TRACE STATUS
        ISSUE DISPLAY TRACE,TT CMD FOR TRANSACTION TRACE STATUS

```

4. Reproduce the problem.
5. Save the trace contents into the trace file created by the CTWTR procedure, executing the the commands shown in Example 2-30.

Example 2-30 Saving the trace contents

```

TRACE CT,ON,COMP=SYSTCPRE,SUB=(RESOLV30)
*190 ITT006A SPECIFY OPERAND(S) FOR TRACE CT COMMAND.
R 190,WTR=DISCONNECT,END
IEE600I REPLY TO 190 IS;WTR=DISCONNECT,END
ITT038I ALL OF THE TRANSACTIONS REQUESTED VIA THE TRACE CT COMMAND
WERE SUCCESSFULLY EXECUTED.
IEE839I ST=(ON,0256K,00512K) AS=ON BR=OFF EX=ON MT=(ON,024K) 503
        ISSUE DISPLAY TRACE CMD FOR SYSTEM AND COMPONENT TRACE STATUS
        ISSUE DISPLAY TRACE,TT CMD FOR TRANSACTION TRACE STATUS

```

6. Stop the CTRACE by issuing the command shown in Example 2-31.

Example 2-31 Stopping CTRACE

```

TRACE CT,OFF,COMP=SYSTCPRE,SUB=(RESOLV30)
ITT038I ALL OF THE TRANSACTIONS REQUESTED VIA THE TRACE CT COMMAND
WERE SUCCESSFULLY EXECUTED.
IEE839I ST=(ON,0256K,00512K) AS=ON BR=OFF EX=ON MT=(ON,024K) 506

```

ISSUE DISPLAY TRACE CMD FOR SYSTEM AND COMPONENT TRACE STATUS
ISSUE DISPLAY TRACE,TT CMD FOR TRANSACTION TRACE STATUS

After these steps, we will have a trace file to be formatted using the IPCS command:

CTRACE COMP(SYSTCPRE) TALLY

Displaying the CTRACE result

The resulting display will show the resolver process entries, as shown in Example 2-32.

Example 2-32 Resolver formatted trace entries

COMPONENT TRACE TALLY REPORT

SYSNAME(SC30)

COMP(SYSTCPRE)

TRACE ENTRY COUNTS AND AVERAGE INTERVALS (IN MICROSECONDS)

FMTID	COUNT	Interval	MNEMONIC	DESCRIBE
00000001	0		CTRACE	CTrace Initialized
00000002	0		CTRACE	Status changed or displayed
00000003	0		CTRACE	CTrace Terminated
00000004	0		CTRACE	!CTrace has abended
00000005	0		CTRACE	CTrace Stopped - Buffers Retain
00010001	0		API	GetHostByAddr Entry Parameters
00010002	0		API	GetHostByAddr Stack Affinity
00010003	0		API	GetHostByAddr Failure
00010004	0		API	GetHostByAddr Success
00010005	0		API	GetHostByAddr GetLocalHostName
00010006	0		API	GetHostByName Entry Parameters
00010007	0		API	GetHostByName Stack Affinity
00010008	0		API	GetHostByName Failure
00010009	0		API	GetHostByName Success

2.5 Additional information

For more specific information regarding the resolver address space, refer to *z/OS Communications Server: IP Configuration Guide*, SC31-8775 and *z/OS Communications Server: IP Configuration Reference*, SC31-8776.

For more information about resolver diagnosis, refer to *z/OS Communications Server: IP Diagnosis Guide*, GC31-8782.

Base functions

The term *base functions* in this case implies the minimum configuration required for the proper operation of a z/OS TCP/IP environment. The base functions that we describe in this chapter are considered necessary for any useful deployment of the TCP/IP stack and commonly used applications.

This chapter discusses the following topics.

Section	Topic
3.1, “The base functions” on page 60	Basic concepts of base functions
3.2, “Common design scenarios for base functions” on page 60	Key characteristics of base functions and why they might be important in your environment
3.3, “z/OS UNIX System Services setup for TCP/IP” on page 65	Selected implementation scenarios, tasks, configuration examples, and problem determination suggestions
3.4, “Configuring z/OS TCP/IP” on page 79	Configuration details for the z/OS TCP/IP environment
3.5, “Implementing the TCP/IP stack” on page 87	Implementation tasks for the TCP/IP stack
3.6, “Activating the TCP/IP stack” on page 93	Messages used to verify the accuracy of the current environment customization data sets used in z/OS UNIX and TCP/IP initialization
3.7, “Reconfiguring the system with z/OS commands” on page 109	z/OS commands used to reconfigure the system
3.8, “Job log versus syslog as diagnosis tool” on page 114	Information about using job log versus syslog when diagnosing issues
3.9, “Message types: Where to find them” on page 114	Listing of message types

3.1 The base functions

Base functions are those functions considered to be standard in TCP/IP environments regardless of the implementation. Base functions establish a functional working environment that can be exploited by other features, or upon which many other functions can be implemented or validated. When the base functions are implemented, they exercise the most commonly used features of a TCP/IP environment, providing an effective way to perform integrity tests and validate the TCP/IP environment before embarking on the more complex features, configurations, and implementations of the stack.

Most of these functions are implemented at the lower layers. There are some base functions that are implemented at the application layer (such as Telnet and FTP). The details of the standard applications can be found in *Communications Server for z/OS V1R12 TCP/IP Implementation Volume 3: Standard Applications*, SG24-7897. Here, we discuss the configuration that provides the infrastructure of the TCP/IP protocol suite in the z/OS Communications Server environment.

3.1.1 Basic concepts

The z/OS TCP/IP stack (a TCP/IP instance) is a full functional implementation of the standard RFC protocols that are fully integrated and tightly coupled between z/OS and UNIX System Services. It provides the environment that supports the base functions, as well as the many traditional TCP/IP applications. The two environments that need to be created and customized to support the z/OS Communications Server for TCP/IP are:

- ▶ A native z/OS environment in which users can exploit the TCP/IP protocols in a standard z/OS application environment such as batch jobs (with JES interface), started tasks, TSO, CICS, and IMS applications.
- ▶ A z/OS UNIX System Services environment that lets you develop and use applications and services that conform to the POSIX or XPG4 standards (UNIX specifications). The z/OS UNIX environment also provides some of the base functions to support the z/OS environment and vice versa.

Because the z/OS Communications Server exploits z/OS UNIX services even for traditional z/OS environments and applications, a full-function mode z/OS UNIX environment, including a Data Facility Storage Management Subsystem (DFSMS), a z/OS UNIX file system, and a security product (such as Resource Access Control Facility, or RACF), are required before the z/OS Communications Server can be started successfully and the TCP/IP environment initialized.

3.2 Common design scenarios for base functions

Because base functions are primarily setting up the *primitives* in the TCP/IP environment, we deal with very basic scenarios, which can be built upon at a later time. For the base functions we consider two scenarios:

- ▶ Single stack environment
- ▶ Multiple stack environment

Important: Although there are specialized cases where multiple stacks per LPAR can provide value, in general we recommend implementing only one TCP/IP stack per LPAR.

3.2.1 Single stack environment

A single stack environment refers to the existence of one TCP/IP system address space in a single z/OS image (LPAR) providing support for the functions and features of the TCP/IP protocol suite.

Dependencies

In order to achieve a successful implementation of the z/OS Communications Server - TCP/IP component, we identified certain dependencies, as explained here:

- ▶ Implement a *full-function* UNIX System Services system on z/OS. Detailed information about this topic is available in *z/OS UNIX System Services Planning*, GA22-7800, and in *z/OS MVS Initialization and Tuning Reference*, SA22-7592. Also refer to *z/OS Program Directory*, GI10-0670, which is available at the following address:
<http://publibz.boulder.ibm.com/epubs/pdf/i1006707.pdf>
- ▶ Define a RACF environment for the z/OS Communications Server - TCP/IP component. This includes defining RACF groups to z/OS UNIX groups to manage resources, profiles, user groups, and user IDs.

An OMVS UID must be defined with UID (0) and assigned to the started task name of the CS for z/OS IP system address space. Detailed information is available in *Communications Server for z/OS TCP/IP Implementation Volume 4: Security and Policy-Based Networking*, SG24-7899, *z/OS Security Server RACF Security Administrator's Guide*, SA22-7683, *z/OS Security Server RACF System Programmer's Guide*, SA22-7681, and *z/OS Security Server RACF Command Language Reference*, SA22-7687.
- ▶ Customize SYS1.PARMLIB members with special reference to BPXPRMxx to use the integrated sockets INET with the AF_INET and AF_INET6 physical file system. Detailed information is available in *z/OS MVS Initialization and Tuning Reference*, SA22-7592, *z/OS UNIX System Services Planning*, GA22-7800, and *z/OS V1R7.0 Program Directory* GI10-0670.
- ▶ Customize the TCP/IP configuration data sets:
 - PROFILE.TCPIP
 - TCPIP.DATA
 - Other configuration data sets
- ▶ Fully functional VTAM is required to support the interfaces used by TCP/IP.

Advantages

The advantages of a single stack are:

- ▶ Fewer CPU cycles are spent processing TCP/IP traffic, because there is only one logical instance of each physical interface in a single stack environment versus a multiple stack environment.
- ▶ Servers use fewer CPU cycles when certain periodic updates arrive (OMPROUTE processing routing updates). Multiple stacks mean multiple copies of OMPROUTE.
- ▶ Each stack requires a certain amount of storage, the most significant being virtual storage.
- ▶ Multiple TCP/IP stacks add a significant level of complexity to TCP/IP system administration tasks.

Considerations for a single stack environment

When creating a TCP/IP stack, you need to consider the other requirements upon which the successful initialization of the stack depends. Very often the initial problems encountered are related to the omission of tasks that were not performed by other disciplines such as RACF administration.

CS for z/OS TCP/IP exploits the tightly coupled design of the z/OS Communications Server, the integration of z/OS and UNIX System Services, and the provision of RACF services. Coordination is the key to a successful implementation the TCP/IP stack.

3.2.2 Multiple stack environment

A multiple stack environment consists of more than one stack running concurrently in a single LPAR. These stacks exist independent of each other, with the ability to be uniquely configured. Each stack can support different features and provide different functions. Each stack is configured in its own address space, and can communicate with the other stacks in the LPAR if so desired.

Dependencies

The dependencies for the multiple stack environment are exactly the same as for the single stack environment, as well as:

- ▶ Additional storage, especially virtual storage
- ▶ Additional CPU cycles for processing subsequent interfaces and services performing periodic functions, such as OMPROUTE routing updates

Advantages

There are advantages for running a multiple stack environment, because it provides you with the flexibility to partition your networking environment. Here are advantages to consider:

- ▶ You might want to establish separate stacks to separate workloads based on availability and security. For example, you might have different requirements for a production stack, a system test stack, and a secure stack.

This approach could, for example, be used to establish a test TCP/IP stack, where new socket applications are tested before they are moved into the production system. You might also want to apply maintenance to a non-production stack so it can be tested before you apply it to the production stack.

- ▶ Your strategy might be to separate workload onto multiple stacks based on the functional characteristics of applications, as with UNIX (OpenEdition) applications and non-UNIX (z/OS) applications.
- ▶ You might be running z/OS servers and UNIX (OpenEdition) servers on the same well-known port (TN3270 and otelnet on port 23). An alternative to this is approach is the BIND for INADDR_ANY function.

Whatever the reason, the ability to configure multiple stacks and have them fully functional, independently and concurrently, can be exploited in many different ways.

Considerations for a multiple stack environment

The considerations for a multiple stack environment are primarily the same as they are for a single stack environment. We therefore indicate here only the *differences* and the *additional considerations* regarding the multiple stack environment.

Sharing resolver between multiple stacks

The general recommendation is that you use separate DATASETPREFIX values per stack and create separate copies of configuration data sets or at the very least resolver data sets. Refer to “The resolver in a multiple stack environment” on page 40, for further details.

Selecting the correct configuration data sets

The resolver needs access to all resolver data sets if there are multiple stacks in multiple z/OS LPARs. Refer to Chapter 2, “The resolver” on page 19, for further details.

TSO clients

TSO client functions can be directed against any number of TCP/IP stacks. Keep in mind, though, that the client must be able to find the TCPIP.DATA data set appropriate for the stack of interest. You can modify your TSO logon procedure with a SYSTCPD DD statement, or use a common TSO logon procedure without the SYSTCPD DD statement and allocate the TCPIP.DATA data set to the appropriate stack of interest.

Stack affinity

Any server or client needs to reference the appropriate stack if the desired stack is not the default stack defined in the BPXPRMxx member of SYS1.PARMLIB. Servers can use the BPXK_SETIBMOPT_TRANSPORT environment variable to override the choice of the default stack. There might also be applications that have affinity to the wrong stack and do not have the option of establishing stack affinity. In those instances, you can execute BPXTCAFF prior to the application execution step. For example:

```
//AFFINITY EXEC PGM=BPXTCAFF,PARM='TCPIPA'
```

This assumes TCPIPA is not the default stack.

Port management

When there is a single stack and the relationship of server to stack is 1:1, port management is relatively simple. Using the PORT statement, the port number can be reserved for the server in the PROFILE.TCPIP for that given stack.

Port management becomes more complex, however, in an environment where there are multiple stacks and a potential for multiple combinations of the same server (for example, UNIX System Services TELNET and TN3270 TELNET). With use of VIPA, it is possible to use the same “well-known” port number, in this case 23, for both services. The distinction would be made by different names mapping to different IP addresses (VIPAs). Therefore, in a multiple stack environment, you need to answer some questions based on the following concepts:

- ▶ **Generic server**

A generic server is a server without affinity for a specific stack, and it provides service to any client in the network. FTP is an example, because the stack is merely a connection linking client and server. The service File Transfer is not related to the internal functioning of the stack, and the server can communicate concurrently over any number of stacks.

- ▶ **Servers with an affinity for a specific stack**

There must be an *explicit* binding of the server application to the chosen stack when the service (for example, UNIX System Services DNS, OSNMP, and ONETSTAT) is related to the internal functioning of the stack.

This bind is made using the `setibmopt()` socket call (to specify the chosen stack) or using the C function `_iptcpn()`, which allows applications to search in the TCPIP.DATA file to find the name of a specific stack.

- ▶ Ephemeral ports

In addition to synchronizing PORT reservations for specific applications across all stacks, you have to synchronize reservations for port numbers that will be dynamically assigned across all stacks when running with multiple stacks.

Those ports are called *ephemeral ports*, which are all above 1024, and are assigned by the stack when none is specified on the application `bind()`. Use the `PORTRANGE` statement in the `PROFILE.TCPIP` to reserve a group of ports, and specify the *same* port range for every stack. You also need to let CINET know which ports are guaranteed to be available on every stack, using the `BPXPRMxx` parmlib member through `INADDRANYPORT` and `INADDRANYCOUNT` statements.

CPU resources

Provisions need to be made for additional CPU cycles and storage (especially virtual storage). These increases in resources are just for the existence of the additional stacks running concurrently.

3.2.3 Recommendation

In general, we recommend implementing only one TCP/IP stack per LPAR, for the following reasons:

- ▶ A TCP/IP stack is capable of exploiting all available resources defined to the LPAR in which it is running. Therefore, starting multiple stacks will not yield any increase in throughput.
- ▶ When running multiple TCP/IP stacks additional system resources, such as memory, CPU cycles, and storage, are required.
- ▶ Multiple TCP/IP stacks add a significant level of complexity to TCP/IP system administration tasks.
- ▶ It is not necessary to start multiple stacks to support multiple instances of an application on a given port number, such as a test HTTP server on port 80 and a production HTTP server also on port 80. This type of support can instead be implemented using BIND-specific support where the two HTTP server instances are each associated to port 80 with their own IP address, using the BIND option on the PORT reservation statement.

3.2.4 Recommendations for MTU

The maximum transmission unit (MTU) is the largest packet size that can be sent using this route. If the packet is larger than this size, the packet will have to be fragmented if fragmentation is permitted. If fragmentation is not permitted, the packet is dropped and an ICMP error is returned to the originator of the packet. If a route is inactive, the configured MTU value that was defined using the MTU parameter in the ROUTE statement (or the default MTU value for the specified interface type) is displayed. If a route is active, then the actual MTU value is displayed.

For more information about MTU sizes for OSA-Express and HiperSockets, refer to *Communications Server for z/OS V1R12 TCP/IP Implementation Volume 3: High Availability*, SG24-7898.

3.3 z/OS UNIX System Services setup for TCP/IP

There are several areas that require your attention and action in order to implement a TCP/IP stack successfully. In Chapter 1, “Introduction to Communications Server for z/OS IP” on page 1, we review the UNIX concepts in the z/OS environment. We make specific references to the BPXPRMxx member in SYS1.PARMLIB. However, it is important to first understand the security considerations for the UNIX environment.

3.3.1 RACF actions for UNIX

Security is an important consideration for most z/OS installations, and there are a few features we need to mention here for the base functions of any TCP/IP environment. TCP/IP has some built-in internal security mechanisms, and it relies on the services of a security manager, such as the IBM Resource Access Control Facility (RACF).

A security manager is a requirement in the Communications Server for z/OS IP environment. As an online application, it is important that TCP/IP undergo security checks to eliminate possible security exposures. Some basic security concepts are included in the following sections, but for a more detailed explanation refer to *Communications Server for z/OS V1R12 TCP/IP Implementation Volume 4: Security and Policy-Based Networking*, SG24-7899.

RACF environment

RACF is very flexible and can be set up and tailored to meet almost all security requirements of large enterprises. All RACF implementations are based on the following key elements:

- ▶ User IDs
- ▶ Groups
- ▶ RACF resources
- ▶ RACF profiles
- ▶ RACF facility classes
- ▶ The hierarchical owner principle, which is applicable for all RACF definitions of user IDs, groups, and RACF resources

RACF implementation

Each unit of work in the z/OS system that requires UNIX System Services must be associated with a valid UNIX System Services identity. A valid identity refers to the presence of a valid UNIX user ID (UID) and a valid UNIX group ID (GID) for each such user. The UID and the GID are defined through the OMVS segment in the user's RACF user profile and in the group's RACF group profile.

Each functional RACF access group must be authorized to access a specific TCP/IP RACF resource with a specific access attribute. The details of this process are discussed in *Communications Server for z/OS V1R12 TCP/IP Implementation Volume 4: Security and Policy-Based Networking*, SG24-7899.

Assigning user IDs to started tasks

In some cases, the user ID and started task must be associated with the UNIX superuser. In other cases, you can associate the user ID and started task with the default user.

RACF offers you two techniques to assign user IDs and group IDs to started tasks:

- ▶ The started procedure name table (ICHRIN03)
- ▶ The RACF STARTED resource profiles

By using the STARTED resources, you can add new started tasks to RACF, and immediately make those new definitions active.

```
IEF695I START T03DNS    WITH JOBNAME T03DNS    IS ASSIGNED TO USER TCPIP3, GROUP
OMVSGRP
```

The user ID and default group must be defined in RACF, which then treats the user ID as any other RACF user ID for its resource access checking. RACF allows multiple started procedure names to be assigned to the same RACF user ID. We used this method to assign RACF user IDs to *all* TCP/IP started tasks.

Started task user IDs

The UNIX System Services tasks OMVS and BPXOINIT need to execute in an z/OS system space and have the special user ID OMVSKERN assigned to them. OMVSKERN has to be defined as superuser with UID 0, program /bin/sh, and home directory.

TCP/IP tasks need RACF user IDs with the OMVS segment defined. The user ID associated with the main TCP/IP address space must be defined as a superuser; the requirements for the individual servers vary, but most need to be a superuser as well.

z/OS VARY TCPIP commands

Access to VARY TCPIP commands can be controlled by RACF. This places restrictions on this command, which can be used to alter and disrupt the TCP/IP environment.

NETSTAT command

Access to the TSO NETSTAT command, the UNIX shell command **onetstat**, and command options can be controlled by RACF, by defining NETSTAT resources to the RACF generic class SERVAUTH. This command might also need to be restricted, because it can be used to alter or drop connections or to stop the TN3270 server.

Establish RACF security environment

The notes that follow are merely an overview of the steps in the process. Consult the instructions in *z/OS Security Server RACF Callable Services*, SA22-7691, to accomplish these tasks.

1. Defining commands for CS for z/OS IP in the RACF OPERCMDS class.
2. Establishing a group ID for a default OMVS group segment:

```
ADDGROUP OEDFLTG OMVS(GID(9999))
```

3. Defining a user ID for a default OMVS group segment:

```
RDEFINE FACILITY BPX.DEFAULT.USER APPLDATA('OEDFLTU/OEDFLTG')
ADDUSER OEDFLTU DFLTGRP(OEDFLTG) NAME('OE DEFAULT USER') PASSWORD(xg18ej)
OMVS(UID(999999) HOME('/') PROGRAM('/bin/sh'))
```

4. Activating or refreshing appropriate facility classes:

```
SETROPTS CLASSACT(FACILITY)
SETROPTS RACLIST(FACILITY)
SETROPTS RACLIST(FACILITY) REFRESH
```

5. Defining one or more superuser IDs to be associated with certain UNIX System Services users and TCP/IP started tasks:

```
ADDGROUP OMVSGRP OMVS(GID(1))
ADDUSER TCP3 DFLTGRP(OMVSGRP) OMVS(UID(0) HOME('/') PROGRAM('/bin/sh'))
```

6. Defining other UNIX System Services users.

You might already have defined RACF groups and users. If this is the case, you can set up a z/OS UNIX file system home directory for each user and add an OMVS identity by altering the group to include a GID (ALTGROUP). Then, using the ISHELL utility, add OE segments for UNIX System Services users (associating them with the altered group and giving each user a distinct UID).

Otherwise, you will have to perform these tasks in a more painstaking manner, for example:

```
ADDGROUP usergrp OMVS(GID(10))
ADDUSER user01 DFLTGRP(usergrp) OMVS(UID(20) HOME('/u/user01')
PROGRAM('/bin/sh/'))
```

More information about RACF with z/OS Communications Server TCP/IP

RACF can be used to protect many TCP/IP resources, such as the TCP/IP stack itself and ports. Further information about securing your TCP/IP implementation can be found in *Communications Server for z/OS V1R12 TCP/IP Implementation Volume 4: Security and Policy-Based Networking*, SG24-7899.

3.3.2 APF authorization

The TCP/IP system program libraries must be APF authorized. Authorized Program Facility (APF) means that z/OS built-in security can be bypassed by programs that are executed from such libraries. CS for z/OS IP data sets have to be protected with RACF. Special attention has to be given to the APF authorized libraries defined in PROGxx.

We used the LNKAUTH=LNKLST specification in SYSx.PARMLIB member IEASYSxx, which means that all libraries in the LNKLST concatenation will be APF authorized. If these libraries are accessed through STEPLIB or JOBLIB, they will not be APF authorized unless they have been specifically defined in the IEAAPFxx or PROGxx member.

SEZALOAD is one library that *must* be made part of your LNKLST concatenation. Because of the LNKAUTH=LNKLST specification, it will be APF authorized when it is accessed through the LNKLST concatenation. The SEZALOAD library holds the TCP/IP system code used by both servers and clients.

In addition to the LNKLST libraries, there are libraries that are not accessed through the LNKLST concatenation, but have to be APF authorized. The SEZATCP library holds the TCP/IP system code used by servers. This library is normally placed in the STEPLIB or JOBLIB concatenation, which is part of the server JCL.

The following libraries might have to be APF authorized, depending on the choices that you make during the installation of z/OS:

SEZALPA

This library holds the TCP/IP modules that must be made part of your system's LPA. If you choose to add the library name to your LPALSTxx member in SYSx.PARMLIB, you also have to make sure the library is APF authorized. If you copy the load modules in the library to an existing LPALSTxx data set, you do not need to authorize the SEZALPA data set.

SEZADSIL

This library holds the load modules used by the SNMP command processor running in the NetView® address space. If you choose to concatenate this library to STEPLIB in the NetView address space, you might have to APF authorize it, if other libraries in the concatenation are already APF authorized.

Every APF-authorized online application might have to be reviewed to ensure that it matches the security standards of the installation. A program is a “well-behaved program” if:

- ▶ Logged-on users cannot access or modify system resources for which they are not authorized.
- ▶ The program does not require any special credentials to be able to execute.

Or, in the case of RACF, the program does not need the RACF authorization attribute OPERATIONS for execution.

Note: User IDs with the RACF attribute OPERATIONS have ALTER access to all data sets in the system. The access authority to single data sets can be specifically lowered or excluded.

3.3.3 Changes to SYS1.PARMLIB members

As we noted, the z/OS environment consists of the traditional MVS and UNIX System Services environment. Because the UNIX System Services environment is implemented within a z/OS system space, there are definitions in the z/OS environment upon which the UNIX System Services environment depends.

SYS1.PARMLIB is the single most important data set in the z/OS environment. It contains most of the parameters that define z/OS as well as many other subsystems. The SYS1.PARMLIB data set definition parameters are critical to the proper initialization and functioning of UNIX System Services and, therefore, to the TCP/IP implementation. Some of the members of interest include:

- ▶ IEASYS00
- ▶ BPXPRMxx
- ▶ Integrated Sockets PFS definitions

IEASYS00

Because the z/OS Communications Server exploits z/OS UNIX services even for traditional MVS environments and applications, a full-function mode z/OS UNIX environment, including a Data Facility Storage Management Subsystem (DFSMS) and z/OS File Systems (including z/OS UNIX file system) is required before the z/OS Communications Server can be started and the TCP/IP environment successfully established.

The IEASYS00 parmlib definitions we used that are relevant to TCP/IP are:

OMVS=7A,
SMS=00,

OMVS=7A specifies that BPXPRM7A is used to configure the z/OS UNIX environment at system initialization time. SMS=00 specifies that IGDSMS00 is to be used for definitions of the Data Facility Storage Management Subsystem at z/OS UNIX initialization time.

BPXPRMxx

All the parameters defined in BPXPRMxx should be reviewed and tailored to individual installation specification and resource utilization. *z/OS UNIX System Services Planning*, GA22-7800, and *z/OS MVS Initialization and Tuning Guide*, SA22-7591, explain the details and significance of each parameter in the BPXPRMxx member.

z/OS UNIX System Services Planning, GA22-7800; *z/OS UNIX System Services User's Guide*, SA22-7802; and *z/OS Program Directory*, GI10-0670, detail the structure, design, installation, and implementation of the z/OS UNIX environment.

z/OS Program Directory, GI10-0670, is available at the following address:

<http://publibz.boulder.ibm.com/epubs/pdf/i1006707.pdf>

Concepts such as Logical and Physical File Systems (PFS) are design components of z/OS UNIX and are not discussed here.

Integrated Sockets PFS definitions

We need to define the desired file systems to support the communication provided by the stack. Example 3-1 illustrates how support for IPv4 and IPv6 (dual mode) is defined for a single stack environment.

Specifying NETWORK definitions for both AF_NET and AF_INET6 provides dual support. If IPv6 support is not desired, then you can omit the NETWORK DOMIAINNAME(AF_INET6) statement and subsequent parameters.

Example 3-1 BPXPRMxx definitions for a single stack supporting dual mode

```
FILESYSTYPE TYPE(UDS)
    ENTRYPOINT(BPXTUINT)
NETWORK DOMAINNAME(AF_UNIX)
    DOMAINNUMBER(1)
    MAXSOCKETS(10000)
    TYPE(UDS)

/* IPv4 support
NETWORK DOMAINNAME(AF_INET)      1
    DOMAINNUMBER(2)
    MAXSOCKETS(25000)
    TYPE(INET)                    2
    INADDRANYPORT(10000)
    INADDRANYCOUNT(2000)

FILESYSTYPE TYPE(INET)            2
    ENTRYPOINT(EZBPFINI)          3

/* IPv6 support
NETWORK DOMAINNAME(AF_INET6)      4
    DOMAINNUMBER(19)
    TYPE(INET)
```

INET specifies a single stack with TCP/IP (by default) as the stack name. In this example, the numbers correspond to the following information:

- 1.** AF_INET specifies IPv4 support for the physical file type for socket address used by this stack (TCP/IP).
- 2.** Specify TYPE(INET) for a single stack environment. If you specify INET, you cannot start multiple TCP/IP stacks.
- 3.** EZBPFINI identifies a TCP/IP stack (this is the only valid value).
- 4.** AF_INET6 specifies IPv6 support for the physical file type for socket address used by this stack (TCP/IP).

Example 3-2 shows BPXPRMxx definitions for a multiple stack environment.

Example 3-2 BPXPRMxx definitions for a multiple stack supporting dual mode

```

FILESYSSTYPE TYPE(UDS) ENTRYPOINT(BPXTUINT)
NETWORK DOMAINNAME(AF_UNIX)
        DOMAINNUMBER(1)
        MAXSOCKETS(10000)
        TYPE(UDS)
FILESYSSTYPE TYPE(CINET)
        ENTRYPOINT(BPXTCINT)
NETWORK DOMAINNAME(AF_INET) 1
        DOMAINNUMBER(2)
        MAXSOCKETS(10000)
        TYPE(CINET) 2
        INADDRANYPORT(10000)
        INADDRANYCOUNT(2000)

NETWORK DOMAINNAME(AF_INET6) 3
        DOMAINNUMBER(19)
        MAXSOCKETS(10000)
        TYPE(CINET)

SUBFILESYSSTYPE NAME(TCPIPA) 4
        TYPE(CINET) 2
        ENTRYPOINT(EZBPFINI) 5
        DEFAULT

SUBFILESYSSTYPE NAME(TCPIPB) 4
        TYPE(CINET) 2
        ENTRYPOINT(EZBPFINI) 5

.....

```

In this example, the numbers correspond to the following information:

- 1.** AF_INET specifies IPv4 support for the physical file type for socket address used by this stack (TCP/IP).
- 2.** Specify TYPE(CINET) for a single stack environment. If you specify INET, you cannot start multiple TCP/IP stacks.
- 3.** AF_INET6 specifies IPv6 support for the physical file type for socket address used by this stack (TCP/IP).
- 4.** Specify the name of TCP/IP stack you want to configure.
- 5.** EZBPFINI identifies a TCP/IP stack (this is the only valid value).

Additional SYS1.PARMLIB updates

The updates are:

1. LNKLSTxx

Add the following CS for z/OS IP link libraries to the z/OS system link list:

- hlq.SEZALOAD
- hlq.SEZALNK2

2. LPALSTxx

Add the following CS for z/OS IP LPA modules to the LPA during IPL of z/OS:

- hlq.SEZALPA

Note: hlq.SEZALPA must be cataloged into the MVS master catalog. hlq.SEZALOAD and hlq.SEZALNK2 can be cataloged into the MVS master catalog. You can omit them from the MVS master catalog if you identify them to include a volume specification as in:

```
TCPIP.SEZALOAD(WTLTCP),  
TCPIP.SEZALNK2(WTLTCP)
```

If the three data sets mentioned were renamed during the installation process, then use these names instead.

3. PROGnn or IEAAPFxx

Add the following TCP/IP libraries for APF authorization:

- hlq.SEZATCP
- hlq.SEZADSIL
- hlq.SEZALOAD
- hlq.SEZALNK2
- hlq.SEZALPA
- SYS1.MIGLIB

4. IEFSSNxx

TNF and VMCF may be required for some of the CS for z/OS IP facilities and components you are using. If you need to configure TNF and VMCF, add the subsystem definitions for the MVS address spaces of TNF and VMCF as follows:

- If you choose to use restartable VMCF and TNF, as is recommended:
 - TNF
 - VMCF
- If you will not be using restartable VMCF and TNF:
 - TNF,MVPTSSI
 - VMCF,MVPXSSI,*nodename*

Set the *nodename* to the MVS NJE node name of this MVS system. It is defined in the JES2 parameter member of SYSx.PARMLIB:

```
NJEDEF      ....  
            OWNNODE=03,  
            ....  
  
N03      NAME=SC30,SNA,NETAUTH
```

Before you make this update, make sure that the hlq.SEZALOAD definition has been added to LNKSTxx and the library itself has been APF authorized. z/OS initializes the address spaces of the TNF and VMCF subsystems during IPL as part of the master scheduler initialization.

5. SCHEDxx

You need to specify certain CS for z/OS IP modules as privileged modules in MVS. The following entries are present in the IBM-supplied program properties table (PPT); however, if your installation has a customized version of the PPT, ensure these entries are present:

- For CS for z/OS IP:

```
PPT PGMNAME(EZBTCPIP) KEY(6) NOCANCEL PRIV NOSWAP SYST LPREF SPREF
```

- If you use restartable VMCF and TNF:

```
PPT PGMNAME(MVPTNF) KEY(0) NOCANCEL NOSWAP PRIV SYST
```

```
PPT PGMNAME(MVPXVMCF) KEY(0) NOCANCEL NOSWAP PRIV SYST
```

- For NPF:

```
PPT PGMNAME(EZAPFES) KEY(1) NOSWAP
```

```
PPT PGMNAME(EZAPAAA) NOSWAP
```

- For SNALINK:

```
PPT PGMNAME(SNALINK) KEY(6) NOSWAP SYST
```

6. COMMNDxx

VMCF and TNF might be required for some of the CS for z/OS IP facilities and components you are using. If you use restartable VMCF and TNF, procedure EZAZSSI must be run during your IPL sequence (EZAZSSI starts VMCF and TNF).

Either use your operation's automation software to start EZAZSSI, or add a command to your COMMNDxx member in SYSx.PARMLIB:

```
COM='S EZAZSSI,P=your_node_name'
```

The value of variable P defaults to the value of the MVS symbolic &SYSNAME. If your node name is the same as the value of &SYSNAME, then you can use the following command instead:

```
COM='S EZAZSSI'
```

When the EZAZSSI address space starts, a series of messages is written to the MVS log indicating the status of VMCF and TNF. Then, the EZAZSSI address space terminates. After VMCF and TNF initialize successfully, you can start your TCP/IP system address spaces.

7. IKJTSOxx

You also need to specify CS for z/OS IP modules as authorized for TSO commands. Update the IKJTSOxx member by adding the following to the AUTHCMD section: MVPXDISP, NETSTAT, TRACERTE, RSH, LPQ, LPR, and LPRM.

8. IEASYSxx

Review your CSA and SQA specifications and verify that the numbers allocated are sufficiently large enough to prevent getmain errors.

```
IEASYSxx: CSA(3000,250M)
```

```
IEASYSxx: SQA(8,448)
```


9. IVTPRMxx

Review the computed CSM requirements to reflect ACF/VTAM and CS for z/OS IP usage:

- IVTPRMxx: FIXED MAX(120M)
- IVTPRMxx: ECSA MAX(120M)

10. CTIEZBxx

Copy CTIEZB00 to SYSx.PARMLIB from hlq.SEZAINST for use with CTRACE.

This member can be customized to include a different size buffer. The default buffer size is 8 MB. This should be increased to 32 MB to allow the capture of debugging information. We made a new member, CTIEZB01, with the buffer size change.

For more information about the use of component tracing (CTTRACE), refer to *z/OS CS: IP Diagnosis*, GC31-8782, and *z/OS CS: IP Migration*, GC31-8773. Also see Chapter 8, “Diagnosis” on page 299.

11. BPXPRMxx

In addition to defining the UNIX Physical File Systems, you must ensure that the ports enabled on the system are consistent with what is defined in the PROFILE.TCPIP data set, as shown in Example 3-3.

Example 3-3 BPXPRMxx member with port range provided by a single stack environment

```
/* IPv4 support
NETWORK DOMAINNAME(AF_INET)
        DOMAINNUMBER(2)
        MAXSOCKETS(25000)
        TYPE(INET)
        INADDRANYPORT(10000) 8
        INADDRANYCOUNT(2000) 8
* IPv6 support
NETWORK DOMAINNAME(AF_INET6)
        DOMAINNUMBER(19)
        TYPE(INET)
```

Ensure that the INADDRANYPORT 8 assignment does not conflict with PORT assignments in the TCPIP.PROFILE data set.

Note: The OpenEdition ENTRYPOINT for CS for z/OS IP is EZBPFINI. If you have the incorrect value in BPXPRMxx member, you might see messages such as EZZ4203I or abend codes such as S806.

Review the values specified in BPXPRMxx for MAXPROCSYS, MAXPROCUSER, MAXUIDS, MAXFILEPROC, MAXPTYS, MAXTHREADTASKS, and MAXTHREADS.

12. IFAPRDxx or PROGxx

Use these to add product and feature information in a z/OS environment.

3.3.4 Changes to SYS1.PROCLIB members

This section explains changes that you can make to incorporate the new TCP/IP functions.

TCP/IP JCL procedures

If you choose to use restartable VMCF and TNF, add procedure EZAZSSI:

```
//EZAZSSI PROC P=' '  
//STARTVT EXEC PGM=EZAZSSI,PARM=&P  
//STEPLIB DD DSN=hlq.SEZATCP,DISP=SHR
```

Update your TCP/IP startup JCL procedure. The sample for the CS for z/OS IP procedure is in hlq.SEZAINST(TCPIPROC).

TSO logon procedures

Update your TSO logon procedures by adding the TCP/IP help data set SYS1.HELP to the //SYSHELP DD concatenation. Optionally, add the //SYSTCPD DD statement to your logon procedures.

Add hlq.SEZAMENU to the //ISPMLIB DD concatenation and hlq.SEZAPENU to the //ISPPLIB DD and the //ISPTLIB DD concatenations.

3.3.5 Additional z/OS customization for z/OS UNIX

Updating the MVS system libraries must be done with great care. Follow the instructions in *z/OS Program Directory, Program Number 5694-A01*, GI10-0670, and check the PSP bucket to ensure that all required PTFs and modifications are done as required. You might need to make changes to some or all of the following members, depending on the features you are installing.

3.3.6 TCP/IP server functions

Each CS for z/OS IP server relies on the use of a security manager, such as RACF. Several servers provide some built-in security functions for additional security. These servers are described in *IBM z/OS Communications Server TCP/IP Implementation Volume 2: Standard Applications*, SG24-7897.

3.3.7 TCP/IP client functions

The client functions of Communications Server for z/OS IP are executed in a TSO environment or a UNIX shell environment. Some functions are also available in other environments, such as batch or started task address spaces.

Any TSO user can execute any TCP/IP command and use a TCP/IP client function to access any other TCP/IP server host through the attached TCP/IP network. If these TCP/IP servers have not implemented adequate password protection, then any TSO client user can log on to these servers and access all data.

3.3.8 UNIX client functions

Certain client functions executed from the UNIX shell environment require superuser authority. The user ID accessing the shell must have an OMVS segment associated with it. RACF considerations for UNIX Client functions in CS for z/OS IP are covered in detail in *IBM z/OS Communications Server TCP/IP Implementation Volume 2: Standard Applications*, SG24-7897, and *Communications Server for z/OS TCP/IP Implementation Volume 4: Security and Policy-Based Networking*, SG24-7899.

Common errors implementing UNIX System Services

In this section, we discuss implementation problems frequently encountered.

Superuser mode

Certain commands and operations from OMVS or from the ISHELL are authorized only for superusers. There are two alternatives for running as a superuser:

- ▶ The user ID can have permanent superuser status.
This means that the ID has been created with a UID value of zero (0). TCP/IP started tasks and some of its servers are also defined with a UID of zero.
- ▶ The user ID can have temporary authority for the superuser tasks.
The defined UID will have been set up as a non-zero value in RACF, but the user will have been granted READ access to the RACF facility class of BPX.SUPERUSER. Also, RACF provides superuser granularity enhancements to assign functions to users that need them.

If you need only temporary authority to enter superuser mode, then granting simple READ permission to the BPX.SUPERUSER facility class will allow the user to switch back and forth between superuser mode and standard mode. You can enter `su` from the OMVS shell, or you can select SETUP OPTIONS from the ISHELL and specify Option #7 to obtain superuser mode.

The user is then authorized to enter commands authorized for the superuser function from the ISHELL, or switch to an OMVS shell the user has already signed onto. The basic prompt level, indicated by the dollar sign (\$) prompt, is changed when in superuser mode to a pound sign (#). The `exit` command takes the user out of superuser mode as well as the OMVS (UNIX) shell. Use of the `whoami` command shows the change of user IDs.

Problems with the home directory

In Example 3-4, the TSO user attempted unsuccessfully to enter the OMVS shell interface from ISPF. The user has an OMVS segment defined but another problem occurs. The user entered the TSO OMVS command to enter the UNIX environment and received the response shown in Example 3-4.

Example 3-4 Error executing the TSO OMVS command

```
FSUM2078I No session was started. The home directory for this TS0/E
user does not exist or cannot be accessed, +
FSUM2079I Function = sigprocmask, return value = FFFFFFFF, return code
code = 9C reason code = 0507014D
```

This error occurred because the home directory that is associated with the user is not defined or authorized in to the OMVS segment. You can determine the home directory with the RACF `listuser` command (if you have the RACF authorization to use the command). However, you still have access to the z/OS file, even though the message was displayed.

A similar problem occurs when trying to access the ISHELL environment, as shown in Example 3-5.

Example 3-5 Error executing in the ISHELL

```
Errno=9Cx Process Initialization error; Reason=0507014D The dub failed
due to an error with the initial home directory. Press Enter to continue.
```

In both cases, the user had an OMVS segment defined in RACF. However, the home directory that was associated with the user in the user's OMVS segment was not defined or authorized. (You can determine the home directory with the RACF **listuser** command.) Authorization is provided with the permission bits.

The same symptom shows up for users without an OMVS segment defined if the BPX.DEFAULTUSER facility has been activated with an inaccessible home directory.

UNIX permission bits

You have already read something about setting up appropriate UNIX permission bits. Example 3-6 shows an example of *incorrect* permission bits set for a user.

Example 3-6 Incorrect permission bits set for a user

```
ICH408I USER(CS01 ) GROUP(WTCRES ) NAME(CS RESIDENT      ) 703
/u/CS01 CL(DIRSRCH) FID(01E2D7D3C5E7F34E2B0F000000000003)
INSUFFICIENT AUTHORITY TO LOOKUP
ACCESSINTENT(--X) ACCESS ALLOWED(OTHER ---)
ICH408I USER(CS01 ) GROUP(WTCRES ) NAME(CS RESIDENT      ) 704
```

In this case, although the user has the UNIX permission bit settings of 755 on the /u/cs01/ directory, the permission bits are set at 600 for the /u/ directory. Thus, you must ensure that all directories in the entire path are authorized with suitable permission bits. After the settings are changed to 755 for the /u/ directory, access to the subdirectory is allowed.

You can display UNIX permission bits from the ISHELL environment or by issuing the **ls -a1F** command from the shell.

The **ls -a1F** options indicate that all files should be listed (including hidden files), that the long format should be displayed, and that the flags about the type of file (link, directory, and so on) should be given.

Default search path and symbolic links

The directory search path is specified in the environment variable \$PATH. Normally this environment variable is set system-wide in the /etc/profile and can be further customized for individual users in \$home/.profile. The sample for /etc/profile sets \$PATH to:

```
/bin:.
```

It should be expanded to:

```
/bin:/usr/sbin:.
```

or (depending on whether you want the current directory searched first or last):

```
./:/bin:/usr/sbin
```

The instructions for setting up this user profile are contained in *z/OS UNIX System Services User's Guide*, SA22-7801, and *z/OS UNIX System Services Planning*, SA22-7800.

Note: To view the search path that has been established for you, issue **echo \$PATH** from the shell environment.

A user might attempt to run a simple TCP/IP command, such as **oping**, and receive an error that the command is not found, as shown in Example 3-7.

Example 3-7 Command not found error

```
BROWSE -- /tmp/cs01
Command ==>
***** Top of Data ***
oping: FSUM7351 not found
```

In this case you must preface the command with the directory path necessary to locate it:

`/usr/lpp/tcpip/bin/oping`

If you experience such a problem, check that the symbolic links are correct. Part of the installation is to run the UNIX MKDIR program to set up the symbolic links for the various commands and programs from their real path to `/bin` or `/usr/sbin`, where they can be found using the default search path.

3.3.9 Verification checklist

The following checklist can help ensure that all z/OS and UNIX System Services related setup tasks are complete for the base functions:

1. If you are using TNF and VMCF, have TNF and VMCF initialized successfully?
Check the console log for a successful start of EZAZSSI, TNF, VMCF.
2. Has the TCP/IP feature of z/OS been enabled or registered in IFAPRDxx?
3. Has a *full-function* OMVS (DFSMS, RACF, zFS) started successfully?
 - Is OMVS active when you issue D OMVS?
 - Is SMS active when you issue D SMS?
 - Have z/OS UNIX file systems been mounted? Verify with D OMVS,F.
 - Is RACF enabled on the system?
4. Have the definitions in BPXPRMxx of SYS1.PARMLIB been made to reflect:
 - The correct stack for the stack (or stacks) you will be running?
 - The support for dual-mode is defined to support IPv4 and IPv6 (AF_INET and AF_INET6)?
 - The correct CS for z/OS IP proc names?
 - The correct use of INET versus CINET?
 - The correct ENTRYPOINT name for Communications Server for z/OS IP versus earlier versions of OE function in TCP/IP (z/OS IP ENTRYPOINT = EZBPFINI)?
 - The mounting of file systems for users? (You can verify with D OMVS,F.)
 - Appropriate values for MAXPROCSYS, MAXPROCUSER, MAXUIDS, MAXFILEPROC, MAXPTYS, MAXTHREADTASKS, and MAXTHREADS?
5. Have z/OS UNIX file systems and directories been created and mounted for the users of the system?

6. Have RACF definitions been put in place? For example:
 - OMVS user IDs and group IDs for your CS for z/OS IP procedures
 - OMVS user IDs and group IDs for your other users, for superusers, for a default user, with definitions for appropriate Facility classes, like BPX.SUPERUSER
 - TCP/IP **VARY** commands
 - NETSTAT commands
7. Have you placed the correct definitions in the z/OS data sets? For example:
 - SYSx.LNKLSTxx
 - SYSx.LPALSTxx
 - SYSx.SCHEDxx
 - SYSx.PROGxx
 - SYSx.IEASYSxx
 - SYSx.IEFSSNxx
 - SYSx.IKJTS0xx
 - SYSx.IVTPRMxx
8. Raw sockets require authorization; they run from SEZALOAD and are usually already authorized. If you have moved applications and functions to another library (which is *not* recommended), ensure that this library is authorized.
9. The loopback address is now 127.0.0.1 for IPv4 and ::1 for IPv6. However, If you require 14.0.0.0, have you added this to the HOME list?
10. Have you computed CSA requirements to include not only ACF/VTAM, but also CS for z/OS IP?
 - IEASYSxx: CSA(3000,250M) (need to review)
 - IEASYSxx: SQA(8,448) (need to review)
11. Have you computed CSM requirements to include not only ACF/VTAM, but also CS for z/OS IP?
 - IVTPRMxx: FIXED MAX(120M)
 - IVTPRMxx: ECSA MAX(120M)
12. Have you modified the CTRACE initialization member (CTIEZB00) to reflect 32 MB of buffer storage?
13. Have you created CTRACE Writer procedures for taking traces?
14. Have you updated your TCP/IP procedure?
15. Have you updated your other procedures, for example, the FTP server procedure?
16. Have you revamped your TCP/IP Profile to use the new statements and to comment out the old?
 - Have you made provisions to address device connections that are no longer supported?
 - Have you investigated all your connections to ensure to what extent they are still supported? (In some cases, definitions will have changed.)
17. Have your applications that relied on VMCF and IUCV sockets been converted now that those APIs are no longer supported?
18. If you are migrating from a previous release, have you reviewed the Planning and Migration checklist in *z/OS CS: IP Migration*, GC31-8773, and made appropriate plans to use the sample data sets?
19. Have you reviewed the list and location of configuration data set samples in *z/OS Communications Server: IP Configuration Reference*, SC31-8776?

3.4 Configuring z/OS TCP/IP

A z/OS TCP/IP environment can be very complex. It is controlled using a large variety of settings, including parmlib members, and /etc files for UNIX System Services. Each of these has a different interface and requires special knowledge to configure.

z/OS Communications Server IP continues to be enhanced with new features, enhancements, and defaults. So if you are migrating from a previous release, consult with the migration guide for your particular release from which you are migrating. For further details, refer to *z/OS Communications Server: New Function Summary*, GC31-8771.

3.4.1 TCP/IP configuration data set names

This topic is described in *z/OS CS: IP Configuration Guide*, SC31-8775. We strongly recommend that you read the information about data set names in this book, before you decide on your data set naming conventions.

The purpose here is to give an introduction to the data set naming and allocation techniques that z/OS Communications Server uses. If you choose, you can allocate some of the configuration data sets either implicitly or explicitly. In addition, you need to ensure that both the MVS and the z/OS UNIX functions can find the data sets.

► Implicit allocation

The name of the configuration data set is resolved at runtime based on a set of rules (the search order) implemented in the various components of TCP/IP. When a data set name has been resolved, the TCP/IP component uses the dynamic allocation services of MVS or of UNIX System Services to allocate that configuration data set. See *z/OS CS: IP Configuration Guide*, SC31-8775, for details.

These are some of the data sets (or files) that can only be *implicitly* allocated in an z/OS Communications Server IP:

```
hlq.ETC.PROTO
hlq.ETC.RPC
hlq.HOSTS.ADDRINFO
hlq.HOSTS.SITEINFO
hlq.SRVRFPT.TCPCHBIN
hlq.SRVRFPT.TCPHGBIN
hlq.SRVRFPT.TCPKJBIN
hlq.SRVRFPT.TCPSCBIN
hlq.SRVRFPT.TCPXLBIN
hlq.STANDARD.TCPCHBIN
hlq.STANDARD.TCPHGBIN
hlq.STANDARD.TCPKJBIN
hlq.STANDARD.TCPSCBIN
hlq.STANDARD.TCPXLBIN
```

In these data set names, hlq is determined using the following search sequence:

- User ID or jobname
- DATASETPREFIX value (or its default of TCP/IP), defined in TCPIP.DATA

Dynamically allocated data sets can include a mid-level qualifier (MLQ), for example, a node name, or a function name.

- For data sets containing a PROFILE configuration file:

`xxxx.nodename.zzzz`

- For data sets containing a translate table used by a particular TCP/IP server:

`xxxx.function_name.zzzz` (for the FTP server the function_name is SRVRFTP)

Data set SYS1.TCPPARMS(TCPDATA) can be dynamically allocated if it contains the TCPIP.DATA configuration file.

- Explicit allocation

For some of the configuration files, you can tell TCP/IP which files to use by coding DD statements in JCL procedures, or by setting UNIX environment variables. The various data sets used by TCP/IP functions and their resolution method are described in *z/OS CS: IP Configuration Guide*, SC31-8775.

3.4.2 PROFILE.TCPIP

Before you start your TCP/IP stack, you must configure the operational and address space characteristics. These definitions are defined in the configuration data set which is often called PROFILE.TCPIP. The PROFILE.TCPIP data set is read by the TCP/IP address space during initialization.

The PROFILE data set contains the following major groups of TCP/IP configuration parameters:

- Operating characteristics
- Port number definitions
- Network interface definitions
- Network routing definitions

A sample PROFILE.TCPIP configuration file is provided in hlq.SEZAINST(SAMPPROF).

You can find detailed information about TCP/IP connectivity and routing definitions in Chapter 4, “Connectivity” on page 117, and Chapter 5, “Routing” on page 205.

PROFILE.TCPIP statements

In this section we show some essential statements for configuring TCP/IP stack.

The syntax for the parameters in the PROFILE can be found in *z/OS Communications Server: IP Configuration Reference*, SC31-8776. Additional profile statements and descriptions are available in “PROFILE.TCPIP statements” on page 418.

Most PROFILE parameters required in a basic configuration have default values that will allow the stack to be initialized and ready for operation. There are, however, a few parameters that must be modified or must be unique to the stack.

Appendix D, “Our implementation environment” on page 455, describes the environment we used to create this book.

DEVICE and LINK

Use DEVICE and LINK statements to define the physical or virtual interfaces, such as OSA, HiperSockets, and VIPA. z/OS Communications Server can define multiple interfaces. You need to define a pair of DEVICE and LINK statements for each interface you want to configure for a TCP/IP stack.

Note: You can instead define IPv4 OSA-Express devices (IPQAENET) with the INTERFACE statement. We recommend this approach, as described in “INTERFACE” on page 81.

Each device type has a different set of parameters that you can define. For details on each device type and its definition, refer to Chapter 4, “Connectivity” on page 117.

The following is an example of DEVICE and LINK statements for defining one OSA in QDIO mode.

```
DEVICE OSA20A0    MPCIPA
LINK   OSA20A0I   IPAQENET    OSA20A0
```

The following is an example of DEVICE and LINK statements for defining one VIPA.

```
DEVICE VIPA1      VIRTUAL 0
LINK   VIPA1L     VIRTUAL 0   VIPA1
```

INTERFACE

The INTERFACE statement defines all IPv6 interfaces and is enhanced to define IPv4 IPAQENET devices. This statement combines the definitions of the DEVICE, LINK, and HOME into a single statement for IPv4 and IPv6.

The INTERFACE statement is set to reference the PORTNAME that is defined in the QDIO TRLE definition statement as per DEVICE and LINK definitions and assigns an IP address to it using the IPADDR operand, according to the HOME definition. Optional operands include subnetmask settings using the /subnetmask bit number value in the IPADDR statement and MTU size with the BEGINROTES or BSDROUTINGPARMS and SOURCEVIPAIN statements, which associates a specific VIPA with this INTERFACE only.

Note: If SOURCEVIPAIN is coded, you define the entire INTERFACE definition block in PROFILE *after* the VIPA DEVICE and LINK statements are defined.

You can define the VLANID and VMAC with the LINK statement, with the additional benefit that you can use the INTERFACE statement to set multiple VLANs on the same OSA port. You cannot, however, define multiple VLANs on the same OSA port with the LINK statement.

The devices that are defined through the INTERFACE statement return different displays than devices that are defined through the DEVICE/LINK statements. See examples in “INTERFACE statement” on page 430.

Example 3-8 shows a sample definition of the INTERFACE statement.

Example 3-8 INTERFACE statement in profile TCP/IP for IPv4 IPAQENET devices

```
INTERFACE OSA20A0I
  DEFINE IPAQENET
  PORTNAME OSA20A0
  IPADDR 10.1.2.12/24
  MTU 1492
  VLANID 20
  VMAC
  SOURCEVIPAIN VIPA2L
```

You can delete a previously defined interface from the stack, after you stop it, with the INTERFACE DELETE command using the OBEYFILE command, as shown in Example 3-9.

Example 3-9 INTERFACE delete statement

```
INTERFACE OSA20A0I
DELETE
```

More examples and displays are available in Appendix B, “INTERFACE statement” on page 430.

Refer to *z/OS Communications Server: IP Configuration Guide*, SC31-8775, and *z/OS Communications Server: IP Configuration Reference*, SC31-8776, for further details.

HOME

The HOME statement is used for assigning an IP address for each interface you defined with DEVICE and LINK statements. The following is an example of a HOME statement.

```
HOME
    10.1.1.10    VIPA1L
    10.1.2.12    OSA20A0I
```

Note: The HOME statement (along with DEVICE and LINK) is mutually exclusive from the INTERFACE statement. You must use one or the other. We recommend that you use INTERFACE, as described in “INTERFACE” on page 81.

The TCP/IP stack uses an IP address of 127.0.0.1 for IPv4 and ::1 for IPv6 as the loopback interfaces. If there is a requirement to represent the loopback IP address of 14.0.0.0 for compatibility with earlier TCP/IP versions, you must code an entry in the HOME statement. The link label specified is LOOPBACK and you can define multiple IP addresses with the LOOPBACK interface. For example:

```
HOME
    14.0.0.0    LOOPBACK
```

You can display the HOME IP address defined in a particular TCP/IP stack with a D TCPIP,procname,Netstat HOME command, as shown in Example 3-10. You can also use the z/OS UNIX shell command **onetstat -h**. There is an additional field, called the Flag field, that indicates which interface is the primary interface. The primary interface is the first entry in the HOME list in the PROFILE.TCPIP definitions unless the PRIMARYINTERFACE parameter is specified.

The PRIMARYINTERFACE statement can be used to specify which link is to be designated as the default local host address for the GETHOSTID() function.

Example 3-10 netstat home display

```
D TCPIP,TCPIPA,N,HOME
EZD0101I NETSTAT CS V1R12 TCPIPA
HOME ADDRESS LIST:
LINKNAME:  VIPA3L
ADDRESS:   10.1.30.10
FLAGS:
LINKNAME:  VIPA1L
ADDRESS:   10.1.1.10
FLAGS:     PRIMARY
LINKNAME:  VIPA2L
```

```

ADDRESS: 10.1.2.10
  FLAGS:
LINKNAME: IUTIQDF4L
  ADDRESS: 10.1.4.11
  FLAGS:
LINKNAME: IUTIQDF5L
  ADDRESS: 10.1.5.11
  FLAGS:
LINKNAME: IUTIQDF6L
  ADDRESS: 10.1.6.11
  FLAGS:
LINKNAME: EZASAMEMVS
  ADDRESS: 10.1.7.11
  FLAGS:
LINKNAME: IQDIOLNK0A01070B
  ADDRESS: 10.1.7.11
  FLAGS:
LINKNAME: VIPL0A010817
  ADDRESS: 10.1.8.23
  FLAGS:
LINKNAME: LOOPBACK
  ADDRESS: 127.0.0.1
  FLAGS:
INTFNAME: OSA2080I
  ADDRESS: 10.1.2.11
  FLAGS:
INTFNAME: OSA2081I
  ADDRESS: 10.1.2.14
  FLAGS:
INTFNAME: OSA20A0I
  ADDRESS: 10.1.2.12
  FLAGS:
INTFNAME: OSA20C0I
  ADDRESS: 10.1.3.11
  FLAGS:
INTFNAME: OSA20E0I
  ADDRESS: 10.1.3.12
  FLAGS:
INTFNAME: LOOPBACK6
  ADDRESS: ::1
  TYPE: LOOPBACK
  FLAGS:
16 OF 16 RECORDS DISPLAYED

```

BEGINROUTES

Use this statement to define static routes for TCP/IP routing table. This statement is optional when you use OMPROUTE dynamic routing daemon. However, if you do not configure OMPROUTE dynamic routing daemon, BEGINROUTES is necessary for a TCP/IP stack to communicate with other hosts. For details on static and dynamic routing, refer to Chapter 5, “Routing” on page 205.

VIPADYNAMIC

This statement is not always necessary. Use this statement to define dynamic VIPA or the functions related to dynamic VIPA, such as sysplex distributor and dynamic VIPA takeover.

Refer to *Communications Server for z/OS V1R12 TCP/IP Implementation Volume 3: High Availability*, SG24-7898 for details about high availability and load balancing functions using dynamic VIPA.

AUTOLOG

The procedures specified in AUTOLOG statement are initialized at TCP/IP startup, so you do not have to start the TCP/IP applications manually after the TCP/IP startup. AUTOLOG also monitors procedures started under its auspices, and will restart a procedure that terminates for any reason unless NOAUTOLOG is specified on the PORT statement.

For UNIX servers, some special rules apply;

- ▶ If the procedure name on the AUTOLOG statement is eight characters long, no jobname needs be specified.
- ▶ If the procedure name on the AUTOLOG statement is less than eight characters long and the job spawns listener threads with different names, you might have to specify the JOBNAME parameter and ensure that the jobname matches that coded on the PORT statement. In the following example, jobname FTPDE1 on the PORT statement matches JOBNAME on the AUTOLOG statement:

```
PORT
  20  TCP * NOAUTOLOG ;OMVS
  21  TCP  FTPDA1 ;Contol Port

AUTOLOG 1
  FTPDA JOBNAME FTPDA1 ; FTP Server
ENDAUTOLOG
```

START

Specify a device name on a START statement to initialize the interface at the TCP/IP stack startup. The following is an example of START statement for an OSA and a HiperSockets device. VIPA does not need to be started because it is virtual and always active.

If you do not specify a device name on a START statement, you can initialize the device with the TCPIP,*procname*,START,*devicename* command after the TCP/IP stack startup.

```
START OSA20A0
START IUTIQDF4
```

IPCONFIG

IPv4 features are defined within IPCONFIG. There is a separate configuration section for IPv6 parameters. Refer to “PROFILE.TCPIP statements” on page 418 for commonly used IPCONFIG statements.

TCPCONFIG

TCP features are defined within TCPCONFIG. Refer to “PROFILE.TCPIP statements” on page 418 for commonly used TCPCONFIG statements.

UDPCONFIG

UDP features are defined within UDPCONFIG. Refer to “PROFILE.TCPIP statements” on page 418 for commonly used UDPCONFIG statements.

GLOBALCONFIG

GLOBALCONFIG statement defines the parameters that are affective to the entire TCP/IP stack. Refer to “PROFILE.TCPIP statements” on page 418 for commonly used GLOBALCONFIG statements.

IPCONFIG6

All IPv6 features are defined within IPCONFIG6.

Locating PROFILE.TCPIP

The following search order is used to locate the PROFILE.TCPIP configuration file:

1. //PROFILE DD
2. //PROFILE DD DSN=TCPIPA.TCPPARMS(PROFA30)
3. jobname.nodename.TCPIP
4. hlq.nodename.TCPIP
5. jobname.PROFILE.TCPIP
6. hlq.PROFILE.TCPIP

The PROFILE must exist. Otherwise, the TCP/IP address space will terminate abnormally with the following message:

```
EZZ0332I DD:PROFILE NOT FOUND. CONTINUING PROFILE SEARCH
EZZ0325I INITIAL PROFILE COULD NOT BE FOUND
```

We recommend using the //PROFILE DD statement in the TCP/IP address space JCL procedure to explicitly allocate the PROFILE data set.

3.4.3 VTAM Resource

As mentioned in the introduction, VTAM provides the Data Link Control layer (Layer 2 of the OSI model) for TCP/IP, including support of the Multi-Path Channel (MPC) interfaces. MPC protocols are used to define the DLC layer for OSA-Express devices in QDIO.

OSA-Express QDIO connections are configured through a TRLE definition. All TRLEs are defined as VTAM major nodes. For further information about MPC-related devices/interfaces refer to Chapter 4, “Connectivity” on page 117.

A TRLE definition we used for our OSA-Express in QDIO mode is shown in Example 3-11.

Example 3-11 TRLE VTAM major node definition for device OSA2080

OSA2080	VBUILD TYPE=TRL	
OSA2080T	TRLE LNCTL=MPC,	*
	READ=2080,	*
	WRITE=2081,	*
	DATAPATH=(2082-2087),	*
	PORTNAME=OSA2080, 1	*
	MPCLEVEL=QDIO	

Because VTAM provides the DLC layer for TCP/IP, then VTAM must be started before TCP/IP. The major node (in our case, OSA2080) should be activated when VTAM is initializing. This will ensure the TRLE is active when the TCP/IP stack is started. This is accomplished by placing an entry for OSA2080 in the VTAM startup list ATCCONxx. The portname 1 (Example 3-11) must also be the same as the device name defined in PROFILE.TCPIP data set on the DEVICE and LINK statements.

This definition can be used for OSA-Express, OSA-Express 2 and OSA-Express 3 using only port 0.

With OSA-Express 3, you can use both ports on the same TRL statement as shown in Example 3-12.

Example 3-12 TRL VTAM majnode definition for two ports for device OSA2080

OSA2080	VBUILD	TYPE=TRL	
OSA200T	TRLE	LNCTL=MPC,	*
		READ=2080,	*
		WRITE=2081,	*
		DATAPATH=(2082-2087),	*
		PORTNAME=OSA2080,	*
		PORTNUM=0,	*
		MPCLEVEL=QDIO	
OSA201T	TRLE	LNCTL=MPC,	*
		READ=2088,	*
		WRITE=2089,	*
		DATAPATH=(208A-208D),	*
		PORTNAME=OSA2081,	*
		PORTNUM=1,	*
		MPCLEVEL=QDIO	

3.4.4 TCPIP.DATA

The resolver configuration file is often called TCPIP.DATA. The TCPIP.DATA configuration data set is the anchor configuration data set for the TCP/IP stack and all TCP/IP servers and clients running on that stack.

The TCPIP.DATA configuration data set is read during initialization of *all* TCP/IP server and client functions. TCPIP.DATA contains the configuration for the resolver address space. We define the way name-to-address or address-to-name resolution is performed by the resolver.

TCPIP.DATA is also used by the TCP/IP applications to specify the TCP/IP stack it establishes an affinity with. The associated TCP/IP stack name is specified with TCPIPJOBNAME statement. Other stack-specific statements are HOSTNAME, which is the host name of the TCP/IP stack, and DATASETPREFIX, which is the data set prefix (hlq) to be used for searching a configuration data set.

The syntax for the parameters in the TCPIP.DATA file can be found in *z/OS Communications Server: IP Configuration Guide*, SC31-8775. A sample TCPIP.DATA configuration file is provided in *hlq.SEZAINST(TCPDATA)*. You can define the TCPIP.DATA parameters in an MVS data set or z/OS UNIX file system file.

For further information about TCPIP.DATA file and the resolver address space, refer to Chapter 2, “The resolver” on page 19.

3.4.5 Configuring the local hosts file

You can set up the local hosts file to support local host name resolution. If you use only the local hosts file for this purpose, your sockets applications will only be able to resolve names and IP addresses that appear in your local hosts file.

If you need to resolve host names outside your local area, you can configure the resolver to use a domain name server (see the NSINTERADDR or NAMESERVER statement in the TCPIP.DATA configuration file). A domain name server can be used in conjunction with the local hosts file. If you have configured your resolver to use a name server, it will always try to

do so, unless your applications were written with a RESOLVE_VIA_LOOKUP symbol in the source code.

Refer to Chapter 2, “The resolver” on page 19 for further explanation and details.

3.5 Implementing the TCP/IP stack

In this scenario we create a TCP/IP stack by the name of TCPIPA on the SC30 system (LPAR A11). We define four OSAs and three HiperSockets interfaces and static routing. Figure 3-1 illustrates the OSA2080 and OSA20A0 pair connecting to the same VLAN using two different OSA-Express features. The same applies to the OSA20C0 and OSA20E0 pair. We also defined a dynamic XCF connection, which in our environment can use either a Coupling Facility link or HiperSockets (CHPID F7).

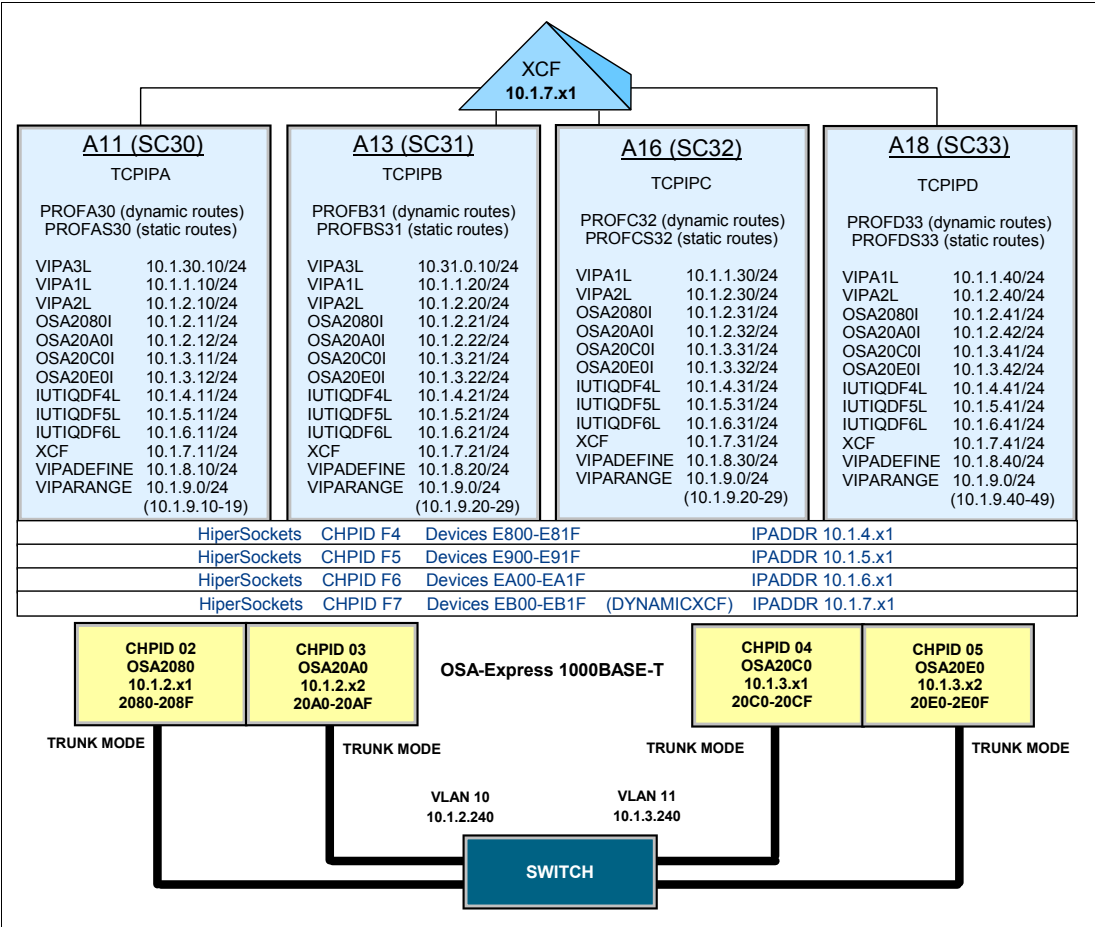


Figure 3-1 Network diagram

To implement the TCP/IP stack to support base functions, perform the following steps:

1. Create a TCPIP.DATA file.
2. Create a PROFILE.TCPIP file.
3. Check BPXPRMxx.
4. Create a TCP/IP cataloged procedure.
5. Add RACF definitions.
6. Create a VTAM TRL major node for MP CIPA OSA.

Allocate the TCPPARMS library to be used for explicitly allocated configuration data sets for the stack, or create a new member in your existing TCPPARMS library. For example, we allocated TCPIPA.TCPPARMS(DATAA30).

3.5.1 Create TCPIP.DATA file

We defined a global TCPIP.DATA and a local TCPIP.DATA for TCPIPA, as shown in Example 3-13 and Example 3-14.

Example 3-13 Global TCPIP.DATA file

```
; *****
; TCPIPA.TCPPARMS (GLOBAL)
; *****\
DOMAINORIGIN  ITS0.IBM.COM
SEARCH  ITS0.IBM.COM IBM.COM
DATASETPREFIX TCPIP
MESSAGECASE MIXED
NSINTERADDR  10.12.6.7
NSPORTADDR 53
RESOLVEVIA UDP
RESOLVERTIMEOUT 10
RESOLVERUDPRETRIES 1
LOOKUP LOCAL
```

We created a local TCPIP.DATA file for the TCPIPA stack.

Example 3-14 Local TCPIP.DATA file

```
; *****
; TCPIPA.TCPPARMS (DATAA30)
; *****
TCPIPJOBNAME TCPIPA
HOSTNAME WTSC30A
DATASETPREFIX TCPIPA
MESSAGECASE MIXED
```

In this example, the numbers correspond to the following information:

- 1.** Specifies the procedure name of TCPIPA stack.
- 2.** Specifies the host name of the TCPIPA stack.

Update the domain name server

If you are using a domain name server, ensure that it is updated with your new host name and address.

Update the local hosts file

If you are not using a domain name server, edit your global ETC.IPNODES file or the local ETC.IPNODES file and add your new host name and address.

3.5.2 Create the PROFILE.TCPIP file

We created a TCP/IP profile and included the statements described in this section.

INTERFACE statement

We configured two OSA-Express3 features, each having four ports. We configured only two ports on each card with the INTERFACE statement. For redundancy we defined two VLANs, with each pair using one port per feature and each pair attached to the same VLAN. This facilitates ARP Takeover.

DEVICE and LINK statement

We defined HiperSockets and VIPA devices with the DEVICE and LINK statements, and the others with the INTERFACE statement.

HOME statement

We assigned a IP address for each interface that was configured with a DEVICE/LINK statement pair.

BEGINROUTES statement

We defined static routes with BEGINROUTES statement to route a traffic to other hosts on a network using the OSA-Express or HiperSockets interfaces.

PORT statement

We reserved TCP ports for some applications with PORT statement

START statement

We defined a START statement to initialize the interfaces at the TCP/IP stack startup.

DYNAMICXCF statement

We defined a DYNAMICXCF statement to dynamically define the device to join the sysplex.

Example 3-15 shows a sample PROFILE.TCPIP file.

Example 3-15 PROFILE.TCPIP file

```
; *****
; TCPIPA.TCPPARMS(PROFA30S)
; *****
ARPAGE 20
;
GLOBALCONFIG NOTCPIPSTATISTICS
;
IPCONFIG DATAGRAMFWD SYSPLEXROUTING
;
DYNAMICXCF 10.1.7.11 255.255.255.0 1
;
SOMAXCONN 10
;
TCPCONFIG TCPSENDBFRSIZE 64K TCPRCVBUFRSIZE 64K SENDGARBAGE FALSE
TCPCONFIG TCPMAXRCVBFRSIZE 256K
TCPCONFIG UNRESTRICTLOWPORTS
;
UDPCONFIG UNRESTRICTLOWPORTS
```

```

;
;INTERFACE OSA20x0I DEFINE IPAQENET (OSA-E) PORTNAME OSA20x0
;TRL MAJ NODE: OSA2080,OSA20A0,OSA20C0,AND OSA20E0
;
INTERFACE OSA2080I
  DEFINE IPAQENET
  PORTNAME OSA2080
  IPADDR 10.1.2.11/24
  VLANID 10
  VMAC
;
INTERFACE OSA20A0I
  DEFINE IPAQENET
  PORTNAME OSA20A0
  IPADDR 10.1.2.12/24
  VLANID 10
  VMAC
;
INTERFACE OSA20C0I
  DEFINE IPAQENET
  PORTNAME OSA20C0
  IPADDR 10.1.3.11/24
  VLANID 11
  VMAC
;
INTERFACE OSA20E0I
  DEFINE IPAQENET
  PORTNAME OSA20E0
  IPADDR 10.1.3.12/24
  VLANID 11
  VMAC
;
;HIPERSOCKETS DEFINITIONS
DEVICE IUTIQDF4 MPCIPA
LINK IUTIQDF4L IPAQIDIO IUTIQDF4
DEVICE IUTIQDF5 MPCIPA
LINK IUTIQDF5L IPAQIDIO IUTIQDF5
DEVICE IUTIQDF6 MPCIPA
LINK IUTIQDF6L IPAQIDIO IUTIQDF6

;
;STATIC VIPA DEFINITIONS
DEVICE VIPA1 VIRTUAL 0
LINK VIPA1L VIRTUAL 0 VIPA1
DEVICE VIPA2 VIRTUAL 0
LINK VIPA2L VIRTUAL 0 VIPA2
;
HOME
  10.1.1.10 VIPA1L
  10.1.2.10 VIPA2L
  10.1.4.11 IUTIQDF4L
  10.1.5.11 IUTIQDF5L
  10.1.6.11 IUTIQDF6L
BEGINRoutes
; Direct Routes - Routes that are directly connected to my interfaces

```

```

;      Destination      Subnet Mask   First Hop Link Name      Packet Size
ROUTE 10.1.2.0/24      =          OSA2080I      MTU 1492
ROUTE 10.1.3.0/24      =          OSA20C0I      MTU 1492
ROUTE 10.1.3.0/24      =          OSA20E0I      MTU 1492
ROUTE 10.1.4.0/24      =          IUTIQDF4L      MTU 8192
ROUTE 10.1.5.0/24      =          IUTIQDF5L      MTU 8192
ROUTE 10.1.6.0/24      =          IUTIQDF6L      MTU 8192
;
PORT
  20  TCP OMVS NOAUTOLOG      ; FTP Server
  21  TCP FTPDA1              ; control port
  23  TCP TN3270A BIND 10.1.9.11 ; OE Telnet Server
  500 UDP IKED                ; @ADI
  520 UDP OMPROUTE NOAUTOLOG   ; OMPROUTE RIPV2 port
  521 UDP OMPROUTE NOAUTOLOG   ; OMPROUTE RIPV2 port
  4500 UDP IKED               ; @AD
  514 UDP OMVS                ; UNIX SyslogD Server  3
;
START OSA2080I
START OSA20C0I
START OSA20E0I
START OSA20A0I
START IUTIQDF4
START IUTIQDF5
START IUTIQDF6

```

3.5.3 Check BPXPRMxx

Refer to SYS1.PARMLIB(BPXPRMxx) and make sure you have your TCP/IP stack name defined in it. If you do not have the stack name in BPXPRMxx, refer to 3.3.3, “Changes to SYS1.PARMLIB members” on page 68.

3.5.4 Create TCP/IP cataloged procedure

We created a cataloged procedure for TCPIPA stack, as shown in Example 3-16.

Example 3-16 Address space JCL procedure (SC30)

```

//TCPIPA  PROC  PARM='CTRACE(CTIEZB00),IDS=00',
//          PROFILE=PROFA&SYSCONE,TCPDATA=DATA&SYSCONE
//TCPIPA  EXEC  PGM=EZBTCPIP,REGION=OM,TIME=1440,
//          PARM=('&PARMS',
//          'ENVAR("RESOLVER_CONFIG=/'TCPIPA.TCPPARMS(&TCPDATA)'")')
//SYSPRINT DD  SYSOUT=*,DCB=(RECFM=VB,LRECL=132,BLKSIZE=136)
//ALGPRINT DD  SYSOUT=*,DCB=(RECFM=VB,LRECL=132,BLKSIZE=136)
//CFGPRINT DD  SYSOUT=*,DCB=(RECFM=VB,LRECL=132,BLKSIZE=136)
//SYSOUT  DD  SYSOUT=*,DCB=(RECFM=VB,LRECL=132,BLKSIZE=136)
//CEEDUMP DD  SYSOUT=*,DCB=(RECFM=VB,LRECL=132,BLKSIZE=136)
//SYSERROR DD  SYSOUT=*
//PROFILE  DD  DISP=SHR,DSN=TCPIPA.TCPPARMS(&PROFILE.)
//SYSTCPD  DD  DSN=TCPIPA.TCPPARMS(&TCPDATA.),DISP=SHR

```

In this example, the numbers correspond to the following information:

- 1.** Illustrates the use of SYSTEM SYMBOLS.
- 2.** SYSTCPD DD statement pointing to TCPIPA.TCPPARMS(DATAA30).

3.5.5 Add RACF definitions

The RACF administrator needs to add RACF definitions to assign started task user IDs to new address spaces, as shown in Example 3-17.

Example 3-17 Defining TCPIP. * procedure to started task*

```

ADDGROUP TCPGRP OMVS(UID(100))
ADDUSER TCPIP DFLTGRP(TCPGRP) OMVS(UID(0) HOME(/) PROGRAM(/bin/sh)) NOPASSWORD
SETOPTS GENERIC(STARTED)
SETOPTS CLASSACT(STARTED) RACLIST(STARTED)
DEFINE STARTED TCPIP*. * STDATA(USER(TCPIP) GROUP(TCPGRP))
SETOPTS RACLIST(STARTED) REFRESH

```

3.5.6 Create a VTAM TRL major node for MPCIPA OSA

We defined our TRLEs in VTAM. Remember to include it in the VTAM startup list in ATCCONxx. Example 3-18 and Example 3-19 are sample TRL major nodes for one OSA device. We then created TRLEs for all OSA devices.

Example 3-18 displays the TRLE VTAM major node definition for device OSA2080.

Example 3-18 TRLE VTAM major node definition for device OSA2080

```

OSA2080  VBUILD TYPE=TRL
OSA200T  TRLE  LNCTL=MPC,                      *
          READ=2080,                            *
          WRITE=2081,                          *
          DATAPATH=(2082-2087),                 *
          PORTNAME=OSA2080,                     *
          PORTNUM=0,                           *
          MPCLEVEL=QDIO
OSA201T  TRLE  LNCTL=MPC,                      *
          READ=2088,                            *
          WRITE=2089,                          *
          DATAPATH=(208A-208D),                 *
          PORTNAME=OSA2081,                     *
          PORTNUM=1,                           *
          MPCLEVEL=QDIO

```

Example 3-19 displays the TRLE's VTAM major node definitions for devices OSA20A0 and OSA20A1.

Example 3-19 TRLE VTAM major node definition for device OSA20A0

```

OSA20A0  VBUILD TYPE=TRL
OSA20A0T TRLE  LNCTL=MPC,                      *
          READ=20A0,                            *
          WRITE=20A1,                          *
          DATAPATH=(20A2-20A7),                 *
          PORTNAME=OSA20A0,                     *

```

```

                                PORTNUM=0,
                                MPCLEVEL=QDIO
*
OSA20A1  VBUILD TYPE=TRL
OSA20A1T TRLE  LNCTL=MPC,
                                READ=20A8,
                                WRITE=20A9,
                                DATAPATH=(20AA-20AE),
                                PORTNAME=OSA20A1,
                                PORTNUM=1,
                                MPCLEVEL=QDIO

```

Note: If server-specific configuration data sets can be explicitly allocated using DD statements, we recommend that you create the configuration data set as a member in the stack-specific TCPPARMS library. If the data set has to be implicitly allocated, remember to create it with the stack-specific data set prefix.

3.6 Activating the TCP/IP stack

If you IPL your z/OS system with PARMLIB definitions similar to our environment, you should get messages similar to those shown in Example 3-20. These are some of the messages that can be used to verify the accuracy of the current environment customization data sets used in z/OS UNIX and TCP/IP initialization.

Note that messages issued by z/OS UNIX begin with the prefix *BPX*.

Example 3-20 IPL and start TCPIPA

```

/* IPL and start of TCPIPA
IEE252I MEMBER BPXPRM1A FOUND IN SYS1.PARMLIB 1
CEE3739I LANGUAGE ENVIRONMENT INITIALIZATION COMPLETE 2
IEE252I MEMBER CTIEZB00 FOUND IN SYS1.IBM.PARMLIB
IEE252I MEMBER CTIIDS00 FOUND IN SYS1.IBM.PARMLIB
IEE252I MEMBER CTINTA00 FOUND IN SYS1.PARMLIB
/*
EZZ4202I Z/OS UNIX - TCP/IP CONNECTION ESTABLISHED FOR TCPIPA
BPXF206I ROUTING INFORMATION FOR TRANSPORT DRIVER TCPIPA HAS BEEN 3
INITIALIZED OR UPDATED.
/*
EZZ4313I INITIALIZATION COMPLETE FOR DEVICE OSA2080
EZZ4313I INITIALIZATION COMPLETE FOR DEVICE OSA2081
EZZ4313I INITIALIZATION COMPLETE FOR DEVICE OSA20C0
EZZ4313I INITIALIZATION COMPLETE FOR DEVICE OSA20E0
EZZ4340I INITIALIZATION COMPLETE FOR INTERFACE OSA20A0I
EZZ4340I INITIALIZATION COMPLETE FOR INTERFACE OSA20A0X
EVB6473I TCP/IP STACK FUNCTIONS INITIALIZATION COMPLETE.
EZZ4313I INITIALIZATION COMPLETE FOR DEVICE IUTIQDF5
EZZ4313I INITIALIZATION COMPLETE FOR DEVICE IUTIQDF6
EZZ4313I INITIALIZATION COMPLETE FOR DEVICE IUTIQDF4
EVB6473I TCP/IP STACK FUNCTIONS INITIALIZATION COMPLETE. 5
EZAIN11I ALL TCPIP SERVICES FOR PROC TCPIPA ARE AVAILABLE.
/*
EVBH006E GLOBALCONFIG SYSPLEXMONITOR RECOVERY was not specified when
IPCONFIG DYNAMICXCF or IPCONFIG6 DYNAMICXCF was configured.
/*
EZD1176I TCPIPA HAS SUCCESSFULLY JOINED THE TCP/IP SYSPLEX GROUP 4
EZBTCPCS

```

```
EZD1214I INITIAL DYNAMIC VIPA PROCESSING HAS COMPLETED FOR TCPIPA
EZZ4313I INITIALIZATION COMPLETE FOR DEVICE IUTIQDIO
/*
```

In this example, the numbers correspond to the following information:

1. The first important message indicates whether the correct UNIX customization data set is used. In our environment it is BPXPRM1A. This contains the *root* system upon which all other file systems are mounted, and it is critical for the current establishment of the correct UNIX System Services environment.

The next set of messages shows the initialization of SMS. SMS is a critical component, because the zFSs are SMS-managed. Note that the file systems are mounted subsequently starting with the root.

2. This message indicates that the Language Environment is available to be exploited by TCP/IP Lotus®, WebSphere®, and parts of the z/OS base, as well as by languages such as C/C++, COBOL, and others.

The next set of messages indicates the successful establishment of the physical file system and availability for socket services for both IPv4 and IPv6. The resolver messages indicate that the resolver process is available to support network resolution, which can be critical to some applications. Note that the initialization of the resolver is completed before TCP/IP.

3. The following two messages indicate the successful initialization of the UNIX System Services environment and TCP/IP services according to the BPXPRMxx definitions.
4. Our environment is defined within a sysplex; therefore, message EZD1176I indicates the connectivity to other active TCP/IP stacks within the sysplex.

Initialization of devices must be completed before they achieve READY status (displayed using the NETSTAT DEVLNKS) and connected to the network.

5. The EZB6473I and EZAIN11I messages are the final initialization messages to complete the successful initialization of the TCP/IP stack.

3.6.1 UNIX System Services verification

A few commands can be used to perform a simple verification of the z/OS UNIX environment after an maintenance IPL; for example, **D SMS** verifies that the system is running a functional SMS environment, as shown in Example 3-21.

Example 3-21 Output of “D SMS” command

```
RESPONSE=SC30
IGD002I 15:43:20 DISPLAY SMS 836
SCDS = SYS1.SMS.SCDS
ACDS = SYS1.SMS.ACDS
COMMDS = SYS1.SMS.COMMDS
DINTERVAL = 150
REVERIFY = NO
ACSDEFAULTS = NO
```

SYSTEM	CONFIGURATION LEVEL	INTERVAL SECONDS
SC30	2010/09/29 15:43:14	10
SC31	2010/09/29 15:43:12	10
SC32	2010/09/29 15:43:14	10
SC33	2010/09/29 15:43:04	10
WTSCPLX5	-----	N/A

Example 3-22 shows the output from the SMS display that supports four different LPARs. The captured information pertains to SC30, SC31, SC32 and SC33.

Example 3-22 Displaying the OMVS system that is running

```

D OMVS,ASID=ALL
BPX0070I 15.48.10 DISPLAY OMVS 878
OMVS      000F ACTIVE              OMVS=(2A) 1
USER      JOBNAME  ASID          PID      PPID STATE   START    CT_SECS
OMVSKERN  BPX0INIT 0025          1        0 MRI----- 07.00.40    1.3 2
  LATCHWAITPID=      0 CMD=BPXPINPR
  SERVER=Init Process
OMVSKERN      0000      65539      1 1L----- 07.00.41    .0
IBMUSER  CEA      0018      16842756  1 1F---P-- 07.00.47    .0
  LATCHWAITPID=      0 CMD=CEAPSRVR
OMVSKERN  SYSLOGDA 0043      50397189  1 HF----- 07.00.41    .6
  LATCHWAITPID=      0 CMD=/usr/sbin/syslogd -c -i -u -f /etc/syslo
NET      NET      001F      65542      1 1F---P-- 07.00.41   28.5
  LATCHWAITPID=      0 CMD=ISTMGCEH
IBMUSER  HZSPROC  002D      33619975  1 1R---B-- 07.00.47    5.3
  LATCHWAITPID=      0 CMD=HZSTKSCH
RMF      RMFGAT   0046      65545      1 1R---P-- 07.00.45   611.9
  LATCHWAITPID=      0 CMD=ERB3GMFC
TCPIP    TCPIP    0045      65546      1 MF---B-- 07.00.45   22.6
  LATCHWAITPID=      0 CMD=EZBTCPIP
IBMUSER  JES2S001 0021      16842765  1 1R----- 07.00.47    .7
  LATCHWAITPID=      0 CMD=IAZNJTCP
TCPIP    TCPIP    0045      16842766  1 1F---B-- 07.00.48   22.6
  LATCHWAITPID=      0 CMD=EZASASUB
TCPIP    TCPIP    0045      33619983  1 1F---B-- 07.00.48   22.6
  LATCHWAITPID=      0 CMD=EZACFALG
TCPIP    FTPMVS1  0042      33619986  1 1FI----- 07.00.59    .0
  LATCHWAITPID=      0 CMD=FTPD
TCPIP    NFSCNT   002C      33619987  1 HR---B-- 07.00.56    .8
  LATCHWAITPID=      0 CMD=GFSCINIT
  SERVER=MVSNFSC
OMVSKERN  INETD1   0040      33619988  1 1FI----- 07.00.55    .0
  LATCHWAITPID=      0 CMD=/usr/sbin/inetd /etc/inetd.conf
TCPIP    REXECD   004B      65557      1 1FI----- 07.00.58    .0
  LATCHWAITPID=      0 CMD=RSMD
TCPIP    TN3270   0049      65558      1 MR---B-- 07.00.58   11.2
  LATCHWAITPID=      0 CMD=EZBTNINI
TCPIP    PORTMAP  004A      33619991  1 1FI----- 07.00.58    .0
  LATCHWAITPID=      0 CMD=PORTMAP
TCPIP    FTPOE1   0041      65560      1 1FI----- 07.00.58    .0
  LATCHWAITPID=      0 CMD=FTPD
IBMUSER  JES2S001 0021      65561      1 1R----- 07.01.04    .7
  LATCHWAITPID=      0 CMD=IAZNJSTK
CS02     CS02     0047      33620005  1 MRI----- 09.09.35    .7
  LATCHWAITPID=      0 CMD=EXEC
CS02     0000      33620006  33620005 1Z----- 09.26.43    .0
TCPIP    TCPIPA   005D      50397230  1 1F---B-- 13.59.13   26.3
  LATCHWAITPID=      0 CMD=EZACFALG
TCPIP    TCPIPD   005F      16842800  1 1R---B-- 14.01.08    1.2
  LATCHWAITPID=      0 CMD=EZACFALG
TCPIP    TCPIPD   005F      33620018  1 MR---B-- 14.01.06    1.2

```

LATCHWAITPID=	0	CMD=EZBTCPIP				
TCPIP TCPIPA	005D	50397245	1	1F---B--	13.59.13	26.3 3
LATCHWAITPID=	0	CMD=EZASASUB				
TCPIP OMPC	0065	50397268	1	HS-----	13.58.29	1.3
LATCHWAITPID=	0	CMD=OMPROUTE				
TCPIP TCPIPB	0056	16842838	1	1F---B--	13.31.54	1.8
LATCHWAITPID=	0	CMD=EZACFALG				
TCPIP TCPIPB	0056	65632	1	1F---B--	13.31.54	1.8
LATCHWAITPID=	0	CMD=EZASASUB				
TCPIP TCPIPB	0056	16842849	1	MF---B--	13.31.51	1.8
LATCHWAITPID=	0	CMD=EZBTCPIP				
TCPIP TRAPFWDB	0064	65637	1	1FI-----	13.32.04	.0
LATCHWAITPID=	0	CMD=EZASNTRA				
TCPIP TNLUNS30	0067	33620070	1	MR---B--	15.22.03	.2
LATCHWAITPID=	0	CMD=EZBTNINI				
TCPIP SNMPQEB	0063	16842855	1	1FI-----	13.32.04	.0
LATCHWAITPID=	0	CMD=SQESERV				
TCPIP TNLUNS30	0067	83951720	1	1F---B--	15.22.03	.2
LATCHWAITPID=	0	CMD=EZBTSSUB				
TCPIP TCPIPC	0060	65646	1	1F---B--	13.58.19	1.6
LATCHWAITPID=	0	CMD=EZASASUB				
TCPIP TCPIPC	0060	65651	1	1F---B--	13.58.19	1.6
LATCHWAITPID=	0	CMD=EZACFALG				
TCPIP TCPIPC	0060	65652	1	MR---B--	13.58.16	1.6
LATCHWAITPID=	0	CMD=EZBTCPIP				
TCPIP TCPIPA	005D	83951733	1	MF---B--	13.59.10	26.3
LATCHWAITPID=	0	CMD=EZBTCPIP				
TCPIP FTPDA1	0050	50397302	1	1FI-----	13.59.23	.0
LATCHWAITPID=	0	CMD=FTPD				
TCPIP OMPA	0053	50397304	1	HS-----	14.17.39	427.9
LATCHWAITPID=	0	CMD=OMPROUTE				
TCPIP TN3270A	005E	67174525	1	1F---B--	15.12.20	.2
LATCHWAITPID=	0	CMD=EZBTSSUB				
TCPIP TN3270A	005E	16842893	1	MR---B--	15.12.20	.2
LATCHWAITPID=	0	CMD=EZBTNINI				
IBMUSER IOASRV	004D	83951763	1	1FI-----	15.44.13	.0
LATCHWAITPID=	0	CMD=IOAXTSRV				
TCPIP SNMPDB	0068	83951764	1	1FI-----	15.46.24	.0
LATCHWAITPID=	0	CMD=EZASNMPD				
CS03 CS03	0052	83951774	1	1RI-----	17.36.25	3.3
LATCHWAITPID=	0	CMD=EXEC				
CS03 CS03	005B	16842915	33620132	1CI-----	17.55.49	.1
LATCHWAITPID=	0	CMD=sh -L				
CS03 CS03	005B	33620132	33619988	1FI-----	17.55.38	.1
LATCHWAITPID=	0	CMD=otelnetsd -Y 9.12.5.202 -p cs03 -A ansi -				

In Example 3-22 on page 95:

- The OMVS member that is running is related to **1** BPXPRM97A.
- The initialization process is running as superuser **2** OMVSKERN, and the PID is 1 (the first process to start).
- There is another TCP/IP started task running **3**.

What is also significant here is that OMVS=DEFAULT is not displayed in the output. In our previous review of the z/OS UNIX environment, we mentioned that the z/OS UNIX System

Services must be customized in *full-function* mode. The display tells you that, at the very least, your system is not running in default mode (*minimal* mode).

Also notice the different TCP/IP stacks and tasks associated with them. There is TCPIPA and TCPIP (the default stack), both executing EZBTCPIP. There are also multiple tasks associated with the same RACF user ID, TCPIP. This offers the advantage of easier maintenance and system definitions. However, this also presents the disadvantage of having no distinguishing features among messages for individual tasks. Many users of TCP/IP and UNIX System Services would assign individual RACF user IDs to each OMVS user for easier problem determination.

For a thorough discussion about the use and implementation of RACF, refer to *Communications Server for z/OS TCP/IP Implementation Volume 4: Security and Policy-Based Networking*, SG24-7899.

Example 3-23 shows the display of available file systems after the initialization of the z/OS UNIX System Services environment. The display should list all of the files defined in the mount statement in the BPXPRMxx member, which in our scenario is BPXPRM9A.

Example 3-23 Output of D OMVS,F

D OMVS,F						
OMVS	000F	ACTIVE	OMVS=(2A)			
TYPENAME	DEVICE	-----STATUS-----	MODE	MOUNTED	LATCHES	
ZFS	784	ACTIVE	RDWR	09/28/2010	L=63	
	NAME=RC30.ZFS			11.48.32	Q=0	
	PATH=/u/rc30					
	OWNER=SC30	AUTOMOVE=Y CLIENT=N				
ZFS	115	ACTIVE	RDWR	09/28/2010	L=57	
	NAME=OMVS.SC33.WEB.HOD			07.13.15	Q=0	
	PATH=/SC33/web/hod					
	OWNER=SC33	AUTOMOVE=U CLIENT=Y				
ZFS	101	ACTIVE	RDWR	09/28/2010	L=55	
	NAME=OMVS.SC33.VAR			07.13.08	Q=0	
	PATH=/SC33/var					
	OWNER=SC33	AUTOMOVE=U CLIENT=Y				
ZFS	100	ACTIVE	RDWR	09/28/2010	L=54	
	NAME=OMVS.SC33.ETC			07.13.08	Q=0	
	PATH=/SC33/etc					
	OWNER=SC33	AUTOMOVE=U CLIENT=Y				
ZFS	99	ACTIVE	RDWR	09/28/2010	L=53	
	NAME=WTSCPLX5.SC33.SYSTEM.ROOT			07.13.08	Q=0	
	PATH=/SC33					
	OWNER=SC33	AUTOMOVE=U CLIENT=Y				
ZFS	72	ACTIVE	RDWR	09/28/2010	L=47	
	NAME=OMVS.SC31.WEB.BW311			07.04.24	Q=0	
	PATH=/SC31/web/bw311					
	OWNER=SC31	AUTOMOVE=U CLIENT=Y				
ZFS	42	ACTIVE	RDWR	09/28/2010	L=40	
	NAME=OMVS.SC30.WEB.BW301			07.00.48	Q=0	
	PATH=/SC30/web/bw301					
	OWNER=SC30	AUTOMOVE=U CLIENT=N				
ZFS	17	ACTIVE	RDWR	09/28/2010	L=33	
	NAME=OMVS.HOM.PRIVATE.HFS			07.00.41	Q=0	
	PATH=/pp/HOD/hostondemand/private					
	OWNER=SC30	AUTOMOVE=Y CLIENT=N				

ZFS	16	ACTIVE	RDWR	09/28/2010	L=32
NAME=OMVS.HOM.HFS				07.00.41	Q=0
PATH=/pp/HOD					
OWNER=SC30 AUTOMOVE=Y CLIENT=N					
ZFS	14	ACTIVE	READ	09/28/2010	L=27
NAME=OMVS.ZOSR1C.Z1CRB1.SIZUR00T				07.00.39	Q=0
PATH=/Z1CRB1/usr/lpp/zosmf/V1R12					
OWNER=SC30 AUTOMOVE=Y CLIENT=N					
ZFS	13	ACTIVE	READ	09/28/2010	L=26
NAME=OMVS.ZOSR1C.Z1CRB1.SBBN7HFS				07.00.39	Q=0
PATH=/Z1CRB1/usr/lpp/zWebSphere0EM/V7R0					
OWNER=SC30 AUTOMOVE=Y CLIENT=N					
ZFS	12	ACTIVE	READ	09/28/2010	L=25
NAME=OMVS.ZOSR1C.Z1CRB1.XML				07.00.39	Q=0
PATH=/Z1CRB1/usr/lpp/ixm					
OWNER=SC30 AUTOMOVE=Y CLIENT=N					
ZFS	11	ACTIVE	READ	09/28/2010	L=24
NAME=OMVS.ZOSR1C.Z1CRB1.JAVA64V6				07.00.39	Q=0
PATH=/Z1CRB1/usr/lpp/java/J6.0_64					
OWNER=SC30 AUTOMOVE=Y CLIENT=N					
ZFS	10	ACTIVE	READ	09/28/2010	L=23
NAME=OMVS.ZOSR1C.Z1CRB1.JAVA31V6				07.00.38	Q=0
PATH=/Z1CRB1/usr/lpp/java/J6.0					
OWNER=SC30 AUTOMOVE=Y CLIENT=N					
ZFS	9	ACTIVE	READ	09/28/2010	L=22
NAME=OMVS.ZOSR1C.Z1CRB1.JAVA64V5				07.00.38	Q=0
PATH=/Z1CRB1/usr/lpp/java/J5.0_64					
OWNER=SC30 AUTOMOVE=Y CLIENT=N					
ZFS	8	ACTIVE	READ	09/28/2010	L=21
NAME=OMVS.ZOSR1C.Z1CRB1.JAVA31V5				07.00.38	Q=0
PATH=/Z1CRB1/usr/lpp/java/J5.0					
OWNER=SC30 AUTOMOVE=Y CLIENT=N					
ZFS	3	ACTIVE	RDWR	09/28/2010	L=16
NAME=OMVS.ZOSR1C.Z1CRB1.ROOT				07.00.37	Q=16
PATH=/Z1CRB1					
OWNER=SC30 AUTOMOVE=Y CLIENT=N					
AUTOMNT	18	ACTIVE	RDWR	09/28/2010	L=34
NAME=*AMD/u				07.00.41	Q=0
PATH=/u					
OWNER=SC30 AUTOMOVE=Y CLIENT=N					
TFS	102	ACTIVE	RDWR	09/28/2010	L=56
NAME=/SC33/TMP				07.13.08	Q=0
PATH=/SC33/tmp					
MOUNT PARM=-s 500					
OWNER=SC33 AUTOMOVE=U CLIENT=Y					
TFS	78	ACTIVE	RDWR	09/28/2010	L=51
NAME=/SC32/TMP				07.05.17	Q=0
PATH=/SC32/tmp					
MOUNT PARM=-s 500					
OWNER=SC32 AUTOMOVE=U CLIENT=Y					
TFS	48	ACTIVE	RDWR	09/28/2010	L=44
NAME=/SC31/TMP				07.04.15	Q=0
PATH=/SC31/tmp					
MOUNT PARM=-s 500					
OWNER=SC31 AUTOMOVE=U CLIENT=Y					

TFS	7	ACTIVE	RDWR	09/28/2010	L=20
NAME=/DEV				07.00.38	Q=0
PATH=/SC30/dev					
MOUNT PARM=-s 10					
OWNER=SC30		AUTOMOVE=U CLIENT=N			
TFS	6	ACTIVE	RDWR	09/28/2010	L=19
NAME=/SC30/TMP				07.00.38	Q=0
PATH=/SC30/tmp					
MOUNT PARM=-s 500					
OWNER=SC30		AUTOMOVE=U CLIENT=N			
HFS	789	ACTIVE	RDWR	09/28/2010	L=62
NAME=CS02.HFS				11.50.30	Q=0
PATH=/u/cs02					
OWNER=SC30		AUTOMOVE=Y CLIENT=N			
HFS	679	ACTIVE	RDWR	09/28/2010	L=61
NAME=CS03.HFS				11.05.22	Q=0
PATH=/u/cs03					
OWNER=SC30		AUTOMOVE=Y CLIENT=N			
HFS	615	ACTIVE	RDWR	09/28/2010	L=60
NAME=CS01.HFS				10.39.56	Q=0
PATH=/u/cs01					
OWNER=SC33		AUTOMOVE=Y CLIENT=Y			
HFS	604	ACTIVE	RDWR	09/28/2010	L=59
NAME=CS06.HFS				10.35.13	Q=0
PATH=/u/cs06					
OWNER=SC32		AUTOMOVE=Y CLIENT=Y			
HFS	602	ACTIVE	RDWR	09/28/2010	L=58
NAME=CS04.HFS				10.35.03	Q=0
PATH=/u/cs04					
OWNER=SC31		AUTOMOVE=Y CLIENT=Y			
HFS	411	ACTIVE	RDWR	09/28/2010	L=52
NAME=CS05.HFS				09.15.21	Q=0
PATH=/u/cs05					
OWNER=SC30		AUTOMOVE=Y CLIENT=N			
HFS	77	ACTIVE	RDWR	09/28/2010	L=50
NAME=OMVS.SC32.VAR				07.05.17	Q=0
PATH=/SC32/var					
OWNER=SC32		AUTOMOVE=U CLIENT=Y			
HFS	76	ACTIVE	RDWR	09/28/2010	L=49
NAME=OMVS.SC32.ETC				07.05.17	Q=0
PATH=/SC32/etc					
OWNER=SC32		AUTOMOVE=U CLIENT=Y			
HFS	75	ACTIVE	RDWR	09/28/2010	L=48
NAME=WTSCPLX5.SC32.SYSTEM.ROOT				07.05.17	Q=0
PATH=/SC32					
OWNER=SC32		AUTOMOVE=U CLIENT=Y			
HFS	67	ACTIVE	RDWR	09/28/2010	L=46
NAME=OMVS.WAS6.BWCELL.BWNODEB.CONFIG.HFS				07.04.23	Q=0
PATH=/SC31/wasbwconfig/bwcell/bwnodeb					
OWNER=SC31		AUTOMOVE=U CLIENT=Y			
HFS	60	ACTIVE	RDWR	09/28/2010	L=45
NAME=BBW6031.SBBOHFS				07.04.22	Q=0
PATH=/SC31/zWebSphereBW					
OWNER=SC31		AUTOMOVE=U CLIENT=Y			
HFS	47	ACTIVE	RDWR	09/28/2010	L=43

NAME=OMVS.SC31.VAR	07.04.15	Q=0
PATH=/SC31/var		
OWNER=SC31 AUTOMOVE=U CLIENT=Y		
HFS 46 ACTIVE	RDWR 09/28/2010	L=42
NAME=OMVS.SC31.ETC	07.04.15	Q=0
PATH=/SC31/etc		
OWNER=SC31 AUTOMOVE=U CLIENT=Y		
HFS 45 ACTIVE	RDWR 09/28/2010	L=39
NAME=WTSCPLX5.SC31.SYSTEM.ROOT	07.04.15	Q=0
PATH=/SC31		
OWNER=SC31 AUTOMOVE=U CLIENT=Y		
HFS 37 ACTIVE	RDWR 09/28/2010	L=38
NAME=OMVS.WAS6.BWCELL.BWNODEA.CONFIG.HFS	07.00.47	Q=0
PATH=/SC30/wasbwconfig/bwcell/bwnodea		
OWNER=SC30 AUTOMOVE=U CLIENT=N		
HFS 30 ACTIVE	RDWR 09/28/2010	L=37
NAME=OMVS.WAS6.BWCELL.BWDMNODE.CONFIG.HFS	07.00.46	Q=0
PATH=/SC30/wasbwconfig/bwcell/bwdmnode		
OWNER=SC30 AUTOMOVE=U CLIENT=N		
HFS 23 ACTIVE	RDWR 09/28/2010	L=36
NAME=BBW6030.SBBOHFS	07.00.45	Q=0
PATH=/SC30/zWebSphereBW		
OWNER=SC30 AUTOMOVE=U CLIENT=N		
HFS 15 ACTIVE	RDWR 09/28/2010	L=31
NAME=OMVS.PP.HFS	07.00.41	Q=0
PATH=/pp		
OWNER=SC30 AUTOMOVE=Y CLIENT=N		
HFS 5 ACTIVE	RDWR 09/28/2010	L=18
NAME=OMVS.SC30.VAR	07.00.37	Q=0
PATH=/SC30/var		
OWNER=SC30 AUTOMOVE=U CLIENT=N		
HFS 4 ACTIVE	RDWR 09/28/2010	L=17
NAME=OMVS.SC30.ETC	07.00.37	Q=0
PATH=/SC30/etc		
OWNER=SC30 AUTOMOVE=U CLIENT=N		
HFS 2 ACTIVE	RDWR 09/28/2010	L=15
NAME=WTSCPLX5.SC30.SYSTEM.ROOT	07.00.37	Q=0
PATH=/SC30		
OWNER=SC30 AUTOMOVE=U CLIENT=N		
HFS 1 ACTIVE	RDWR 09/28/2010	L=14
NAME=WTSCPLX5.SYSPLEX.ROOT	07.00.37	Q=0
PATH=/		
OWNER=SC30 AUTOMOVE=Y CLIENT=N		

Example 3-24 shows some of the files defined in the active BPXPRM9A member for comparative purposes only. We can compare the names defined in the active BPXPRM9A member with the names that are actually active by using the **D OMVS,F** command.

Example 3-24 BPXPRM1A member

```

ROOT  FILESYSTEM('WTSCPLX5.SYSPLEX.ROOT')
      TYPE(HFS)
      AUTOMOVE
      MODE(RDWR)

MOUNT FILESYSTEM('WTSCPLX5.&SYSNAME..SYSTEM.ROOT')
```

```

MOUNTPOINT('/&SYSNAME.')
UNMOUNT
TYPE(HFS)  MODE(RDWR)

MOUNT FILESYSTEM('OMVS.ZOSR1B.&SYSR1..ROOT')
MOUNTPOINT('/$VERSION')
AUTOMOVE
TYPE(HFS)  MODE(RDWR)

MOUNT FILESYSTEM('OMVS.&SYSNAME..ETC')
MOUNTPOINT('/&SYSNAME./etc')
UNMOUNT
TYPE(HFS)  MODE(RDWR)

MOUNT FILESYSTEM('OMVS.&SYSNAME..VAR')
MOUNTPOINT('/&SYSNAME./var')
UNMOUNT
TYPE(HFS)  MODE(RDWR)

MOUNT FILESYSTEM('/&SYSNAME./TMP')
TYPE(TFS)  MODE(RDWR)
MOUNTPOINT('/&SYSNAME./tmp')
PARM('-s 500')
UNMOUNT

MOUNT FILESYSTEM('/DEV')
MOUNTPOINT('/dev')
TYPE(TFS)
PARM('-s 10')
UNMOUNT

MOUNT FILESYSTEM('OMVS.ZOSR1B.&SYSR1..JAVA31V5')
MOUNTPOINT('/usr/lpp/java/J5.0')
TYPE(HFS)  MODE(RDWR)
MOUNT FILESYSTEM('OMVS.ZOSR1B.&SYSR1..JAVA64V5')
MOUNTPOINT('/usr/lpp/java/J5.0_64')
TYPE(HFS)  MODE(RDWR)
MOUNT FILESYSTEM('OMVS.ZOSR1B.&SYSR1..JAVA31V6')
MOUNTPOINT('/usr/lpp/java/J6.0')
TYPE(HFS)  MODE(RDWR)
MOUNT FILESYSTEM('OMVS.ZOSR1B.&SYSR1..JAVA64V6')
MOUNTPOINT('/usr/lpp/java/J6.0_64')
TYPE(HFS)  MODE(RDWR)
MOUNT FILESYSTEM('OMVS.ZOSR1B.&SYSR1..XML')
MOUNTPOINT('/usr/lpp/ixm')
TYPE(HFS)  MODE(RDWR)

```

The OMVS processes can also be displayed within the z/OS UNIX environment, and similar comparisons can be made. Use the shell environment to look at UNIX processes and to execute the UNIX command **ps -ef**. This displays all processes and their environments in forest or family tree format.

Refer to *z/OS UNIX System Services Planning*, GA22-7800 and *z/OS UNIX System Services User's Guide*, SA22-7802 for detailed information about UNIX commands in the z/OS UNIX environment.

Notice that in Example 3-25, the UNIX System Services after this initialization is running with user ID BPXROOT. The reason for this is because RACF cannot map a UNIX System Services UID to an MVS user ID correctly if there are multiple MVS user IDs defined with the same UID. So RACF uses the last referenced MVS user ID.

Example 3-25 UNIX System Services processes display from the shell

```

1 @ SC30:/u/cs01>ps -ef
UID      PID      PPID  C   STIME TTY      TIME CMD
BPXROOT      1          0  -   Oct 25 ?        0:02 BPXPINPR
BPXROOT    33619971      1  -   Oct 25 ?        0:00 CEAPSRVR
BPXROOT    16842756      1  -   Oct 25 ?        0:08 HZSTKSCH
      NET    50397189      1  -   Oct 25 ?        1:29 ISTMGCEH
BPXROOT    50397190      1  -   Oct 25 ?        0:02 BPXVCLNY
BPXROOT      65543      1  -   Oct 25 ?        0:01 /usr/sbin/syslogd -c -i
BPXROOT      65543      1  -   Oct 25 ?        0:01 /usr/sbin/syslogd -c -
      -u -f /etc/sysloga.conf
BPXROOT      65545      1  -   Oct 25 ?        22:55 ERB3GMFC
BPXROOT      65546      1  -   Oct 25 ?        1:14 EZBTCPIP
BPXROOT    16842763      1  -   Oct 25 ?        1:14 EZACFALG
BPXROOT      65548      1  -   Oct 25 ?        1:14 EZASASUB
BPXROOT    33619981      1  -   Oct 25 ?        0:02 BPXVCMT
BPXROOT    16842766      1  -   Oct 25 ?        0:01 IAZNJTCP
BPXROOT    16842768      1  -   Oct 25 ?        0:00 /usr/sbin/inetd
      /etc/inet.conf
BPXROOT      65553      1  -   Oct 25 ?        0:02 GFSCINIT
BPXROOT    33619986      1  -   Oct 25 ?        0:00 PORTMAP
BPXROOT    33619987      1  -   Oct 25 ?        0:01 IAZNJSTK
BPXROOT    16842772      1  -   Oct 25 ?        0:31 EZBTZMST
BPXROOT      65557      1  -   Oct 25 ?        0:00 RSHD
BPXROOT    16842774      1  -   Oct 25 ?        0:00 IOAXTSRV
BPXROOT      65559      1  -   Oct 25 ?        0:00 FTPD
BPXROOT      65560      1  -   Oct 25 ?        0:00 FTPD
BPXROOT    33620011      1  -  13:25:00 ?        0:05 EZBTCPIP
BPXROOT    33620012      1  -  13:25:02 ?        0:05 EZACFALG
BPXROOT      65581      1  -  13:25:02 ?        0:05 EZASASUB
BPXROOT    16842799      1  -  13:25:13 ?        0:04 OMPROUTE
      CS04    33620034      1  -  16:36:32 ?        0:02 OMVS
      CS04      65603    33620034 -  16:36:32 ttyp0000 0:02 sh -L
BPXROOT    83951684      65603 -  16:37:12 ttyp0000 0:00 sh
BPXROOT    16842821    83951684 -  16:37:19 ttyp0000 0:00 ps -ef

```

Here are some typical UNIX commands:

- ▶ The **mkdir /u/cs01** command creates the directory for the user mount point. The permission bits would be set as specified in the `etc/profile` or `$home/.profile`.
- ▶ The **ls -all** command lists the files with their permission bits. From time to time you might need to change the permission bits in the file.
- ▶ The **chmod** command is used to change the permission bits associated with files.
- ▶ The TSO/E interface can be used to work with zOS UNIX files. You can browse files using the ISHELL PDSE interface or you can execute the **obrowse** command from the OMVS shell environment. You can also edit files using the ISHELL tools, or you can use the **oedit** command from the OMVS shell.

Note: Both **obrowse** and **oedit** are TSO commands. If you used telnet or rlogin to get to the UNIX System Services shell, you have to use the **cat** command and the vi editor.

The ISHELL provides an ISPF look and feel. The OMVS shell provides a more UNIX or DOS look and feel, and of course for real UNIX users there is the vi editor.

Starting z/OS Communications Server TCP/IP

Example 3-26 shows the startup of our TCP/IP stack.

Example 3-26 z/OS Communications Server TCP/IP startup

```
S TCPIPA
$HASP373 TCPIPA   STARTED
IEE252I MEMBER CTIEZB00 FOUND IN SYS1.IBM.PARMLIB 1
IEE252I MEMBER CTIIDS00 FOUND IN SYS1.IBM.PARMLIB
IEE252I MEMBER CTINTA00 FOUND IN SYS1.PARMLIB
EZZ0300I OPENED PROFILE FILE DD:PROFILE
EZZ0309I PROFILE PROCESSING BEGINNING FOR DD:PROFILE
EZZ0316I PROFILE PROCESSING COMPLETE FOR FILE DD:PROFILE 2
EZZ0641I IP FORWARDING NOFWMULTIPATH SUPPORT IS ENABLED
EZZ0350I SYSPLEX ROUTING SUPPORT IS ENABLED 3
EZZ0624I DYNAMIC XCF DEFINITIONS ARE ENABLED 4
EZZ0338I TCP PORTS 1 THRU 1023 ARE NOT RESERVED
EZZ0338I UDP PORTS 1 THRU 1023 ARE NOT RESERVED
EZZ0613I TCPIPSTATISTICS IS DISABLED 5
EZZ4202I Z/OS UNIX - TCP/IP CONNECTION ESTABLISHED FOR TCPIPA 6
EZZ4340I INITIALIZATION COMPLETE FOR INTERFACE OSA2080I
EZZ4340I INITIALIZATION COMPLETE FOR INTERFACE OSA20A0I
EZZ4340I INITIALIZATION COMPLETE FOR INTERFACE OSA20C0I
EZZ4340I INITIALIZATION COMPLETE FOR INTERFACE OSA20D0I
EZB6473I TCP/IP STACK FUNCTIONS INITIALIZATION COMPLETE. 7
EZAIN11I ALL TCPIP SERVICES FOR PROC TCPIPA ARE AVAILABLE.
EZD1176I TCPIPA HAS SUCCESSFULLY JOINED THE TCP/IP SYSPLEX GROUP
EZD1214I INITIAL DYNAMIC VIPA PROCESSING HAS COMPLETED FOR TCPIP
```

In this example, the numbers correspond to the following information:

- 1.** Shows how the member that defines CTRACE processing has been found (CTIEZB00). This is discussed in 8.4.1, “Taking a component trace” on page 316.
- 2.** Shows how the PROFILE.TCPIP for the stack has been found and processed.
- 3.** Sysplex routing is enabled, so communication between z/OS TCP/IP is possible.
- 4.** Dynamic XCFs are enabled (DYNAMICXCF parameter).
- 5.** TCPIPSTATISTICS will not be generated.
- 6.** Shows how the stack has been bound to UNIX System Services. It indicates that the Common INET pre-router has successfully obtained a copy of the IP layer routing table from the stack.
- 7.** The stack (TCP/IP) is successfully initialized and READY FOR WORK.

Important: Because TCP/IP shares its Data Link Controls (DLCs) with VTAM, you must restart TCP/IP if you restart VTAM.

Note: See 1.3.3, “Reusable address space ID” on page 6 if you want to run TCPIP in a reusable address space ID.

3.6.2 Verifying TCP/IP configuration

After the configuration files are updated we verified the configuration and we restarted the TCP/IP address space, ensuring that we saw the following message:

```
EZB6473I TCP/IP STACK FUNCTIONS INITIALIZATION COMPLETE
```

If the message is not displayed, the messages issued by the TCP/IP address space should describe why TCP/IP did not start.

Displaying the TCP/IP configuration

To display the enabled features and operating characteristics of a TCP/IP stack, enter any of the following commands:

- ▶ TSO/E command `NETSTAT CONFIG`
- ▶ MVS command `D TCPIP,procname,NETSTAT,CONFIG`
- ▶ UNIX shell command `onetstat -f`

Example 3-27 shows the output from the `NETSTAT CONFIG` display.

Example 3-27 NETSTAT CONFIG display

```
D TCPIP,TCPIPA,NETSTAT,CONFIG
EZD0101I NETSTAT CS V1R12 TCPIPA 442
TCP CONFIGURATION TABLE:
DEFAULTRCVBUFSIZE: 00131072  DEFAULTSNDBUFSIZE: 00131072
DEFLTMAXRCVBUFSIZE: 00262144  SOMAXCONN: 0000000010
MAXRETRANSMITTIME: 120.000  MINRETRANSMITTIME: 0.500
ROUNDTRIPGAIN: 0.125  VARIANCEGAIN: 0.250
VARIANCEMULTIPLIER: 2.000  MAXSEGLIFETIME: 30.000
DEFAULTKEEPALIVE: 00000120  DELAYACK: YES
RESTRICTLOWPORT: NO  SENDGARBAGE: NO
TCPTIMESTAMP: YES  FINWAIT2TIME: 600
TTLS: NO
UDP CONFIGURATION TABLE:
DEFAULTRCVBUFSIZE: 00065535  DEFAULTSNDBUFSIZE: 00065535
CHECKSUM: YES
RESTRICTLOWPORT: NO  UDPQUEUELIMIT: NO
IP CONFIGURATION TABLE:
FORWARDING: YES  TIMETOLIVE: 00064  RSMTIMEOUT: 00060
IPSECURITY: NO
ARPTIMEOUT: 01200  MAXRSMsize: 65535  FORMAT: LONG
IGREDIRECT: YES  SYSPLXROUT: YES  DOUBLENOP: NO
STOPCLAWER: NO  SOURCEVIPA: YES
MULTIPATH: CONN  PATHMTUDSC: YES  DEVRTRYDUR: 0000000090
DYNAMICXCF: YES
  IPADDR: 10.1.7.11  SUBNET: 255.255.255.0  METRIC: 08
  SECCCLASS: 255
QDIOACCEL: YES  QDIOACCELPRIORITY: 1
IQDIOROUTE: N/A
TCPSTACKSRCVIPA: NO
IPV6 CONFIGURATION TABLE:
```



```

FORWARDING: YES HOPLIMIT: 00255 IGREDIRECT: NO
SOURCEVIPA: NO MULTIPATH: NO ICMPERRLIM: 00003
IGRTRHOPLIMIT: NO
IPSECURITY: NO
DYNAMICXCF: NO
TCPSTACKSRCVIPA: NO
TEMPADDRESSES: NO
SMF PARAMETERS:
TYPE 118:
  TCPINIT: 00 TCPTERM: 00 FTPCLIENT: 00
  TN3270CLIENT: 00 TCPIPSTATS: 00
TYPE 119:
  TCPINIT: YES TCPTERM: YES FTPCLIENT: YES
  TCPIPSTATS: YES IFSTATS: NO PORTSTATS: NO
  STACK: NO UDPTERM: NO TN3270CLIENT: YES
  IPSECURITY: NO PROFILE: NO DVIPA: NO
GLOBAL CONFIGURATION INFORMATION:
TCPIPSTATS: NO ECSALIMIT: 0000000K POOLLIMIT: 0000000K
MLSCHKTERM: NO XCFGRPID: 21 IQDVLANID: 21
SEGOFFLOAD: YES SYSPLEXWLMPOLL: 060 MAXRECS: 100
EXPLICITBINDPORTRANGE: 00000-00000 IQDMULTIWRITE: YES
WLMPPRIORITYQ: NO
SYSPLEX MONITOR:
  TIMERSECS: 0060 RECOVERY: YES DELAYJOIN: YES AUTOREJOIN: NO
MONINTF: NO DYNROUTE: NO JOIN: YES
ZIIP:
  IPSECURITY: NO IQDIOMULTIWRITE: YES
NETWORK MONITOR CONFIGURATION INFORMATION:
PKTTRCSRV: NO TCPCNNSRV: NO NTASRV: NO
SMFSRV: YES
  IPSECURITY: YES PROFILE: YES CSSMTP: YES CSMAIL: NO DVIPA: YES
AUTOLOG CONFIGURATION INFORMATION: WAIT TIME: 0300
PROCNAME: FTPDA JOBNAME: FTPDA1
  PARMSTRING:
  DELAYSTART: NO
PROCNAME: OMPA JOBNAME: OMPA
  PARMSTRING:
  DELAYSTART: NO
PROCNAME: IOASRV JOBNAME: IOASRV
  PARMSTRING:
  DELAYSTART: NO
END OF THE REPORT

```

Parameters such as SOURCEVIPA can be either ENABLED or DISABLED. A value of 01 in the NETSTAT CONFIG display means it is ENABLED. In this example, the numbers correspond to the following information:

- 1.** The settings in effect in the TCPCONFIG parameters.
- 2.** The settings for the UDPCONFIG parameters.
- 3.** The settings in effect in the IPCONFIG parameters.
- 4.** The settings in effect for SMFCONFIG.
- 5.** The settings in effect for GLOBALCONFIG.
- 6.** The setting in effect for Network Monitoring Information.

Displaying the status of devices

You can display the status of devices by using the MVS display command TSO NETSTAT, as shown in Example 3-28, or UNIX **onetstat -d**.

Example 3-28 Results of device display

```

D TCPIP,TCPIPA,N,DE
..... Lines deleted
INTFNAME: OSA2080I          INTFTYPE: IPAQENET  INTFSTATUS: READY 1
  PORTNAME: OSA2080    DATAPATH: 2082    DATAPATHSTATUS: READY
  CHPIDTYPE: OSD
  SPEED: 0000001000
  IPBROADCASTCAPABILITY: NO
  VMACADDR: 02000C776873  VMACORIGIN: OSA    VMACROUTER: LOCAL
  ARPOFFLOAD: YES 2          ARPOFFLOADINFO: YES
  CFGMTU: 1492          ACTMTU: 1492
  IPADDR: 10.1.2.11/24
  VLANID: 10 3          VLANPRIORITY: DISABLED
  DYNVLANREGCFG: NO      DYNVLANREGCAP: YES
  READSTORAGE: GLOBAL (4096K)
  INBPERF: BALANCED
  CHECKSUMOFFLOAD: YES   SEGMENTATIONOFFLOAD: YES
  SECCLASS: 255          MONSYSPLEX: NO
  ISOLATE: NO            OPTLATENCYMODE: NO
MULTICAST SPECIFIC:
  MULTICAST CAPABILITY: YES
  GROUP          REFCNT          SRCFLTMD
  -----
  224.0.0.1      0000000001      EXCLUDE
  SRCADDR: NONE
..... Lines deleted

```

In this example, the numbers correspond to the following information:

- 1.** Indicates the overall status of the OSA interface OSA2080I: READY. If this status is not READY, verify that the VTAM Major node is active. You can do this by using the VTAM command D NET,TRL.
- 2.** The OFFLOAD feature is enabled.
- 3.** The VLAN ID defined on the LINK statement in the PROFILE data set.

Example 3-29 Results of the TRLE display

```

D NET,TRL,TRLE=OSA2080T
IST075I NAME = OSA2080T, TYPE = TRLE 337
IST1954I TRL MAJOR NODE = OSA2080
IST486I STATUS= ACTIV, DESIRED STATE= ACTIV 1
IST087I TYPE = LEASED          , CONTROL = MPC , HPDT = YES
IST1715I MPCLEVEL = QDIO      MPCUSAGE = SHARE
IST2263I PORTNAME = OSA2080    PORTNUM = 0    OSA CODE LEVEL = 000C
IST2337I CHPID TYPE = OSD      CHPID = 02
IST1577I HEADER SIZE = 4096 DATA SIZE = 0 STORAGE = ***NA***
IST1221I WRITE DEV = 2081 STATUS = ACTIVE      STATE = ONLINE 2
IST1577I HEADER SIZE = 4092 DATA SIZE = 0 STORAGE = ***NA***
IST1221I READ  DEV = 2080 STATUS = ACTIVE      STATE = ONLINE 2
IST924I -----
IST1221I DATA  DEV = 2082 STATUS = ACTIVE      STATE = N/A 3

```

```

IST1724I I/O TRACE = OFF  TRACE LENGTH = *NA*
IST1717I ULPID = TCPIPA
IST2309I ACCELERATED ROUTING ENABLED
IST2331I QUEUE   QUEUE   READ
IST2332I ID      TYPE    STORAGE
IST2205I -----
IST2333I RD/1    PRIMARY  4.0M(64 SBALS)
IST2305I NUMBER OF DISCARDED INBOUND READ BUFFERS = 0
IST1757I PRIORITY1: UNCONGESTED PRIORITY2: UNCONGESTED
IST1757I PRIORITY3: UNCONGESTED PRIORITY4: UNCONGESTED
IST2190I DEVICEID PARAMETER FOR OSAENTA TRACE COMMAND = 01-01-00-02
IST1801I UNITS OF WORK FOR NCB AT ADDRESS X'0F30D010'
IST1802I P1 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P2 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P3 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P4 CURRENT = 0 AVERAGE = 1 MAXIMUM = 2
IST924I -----
IST1221I DATA  DEV = 2083 STATUS = ACTIVE      STATE = N/A
IST1724I I/O TRACE = OFF  TRACE LENGTH = *NA*
IST1717I ULPID = TCPIPC
IST2309I ACCELERATED ROUTING ENABLED
IST2331I QUEUE   QUEUE   READ
IST2332I ID      TYPE    STORAGE
IST2205I -----
IST2333I RD/1    PRIMARY  4.0M(64 SBALS)
IST2305I NUMBER OF DISCARDED INBOUND READ BUFFERS = 0
IST1757I PRIORITY1: UNCONGESTED PRIORITY2: UNCONGESTED
IST1757I PRIORITY3: UNCONGESTED PRIORITY4: UNCONGESTED
IST2190I DEVICEID PARAMETER FOR OSAENTA TRACE COMMAND = 01-01-00-03
IST1801I UNITS OF WORK FOR NCB AT ADDRESS X'0FCF0010'
IST1802I P1 CURRENT = 0 AVERAGE = 1 MAXIMUM = 2
IST1802I P2 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P3 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P4 CURRENT = 0 AVERAGE = 1 MAXIMUM = 1
IST924I -----
IST1221I DATA  DEV = 2084 STATUS = RESET      STATE = N/A
IST1724I I/O TRACE = OFF  TRACE LENGTH = *NA*
IST924I -----
IST1221I DATA  DEV = 2085 STATUS = RESET      STATE = N/A
IST1724I I/O TRACE = OFF  TRACE LENGTH = *NA*
IST924I -----
IST1221I DATA  DEV = 2086 STATUS = RESET      STATE = N/A
IST1724I I/O TRACE = OFF  TRACE LENGTH = *NA*
IST924I -----
IST1221I DATA  DEV = 2087 STATUS = RESET      STATE = N/A
IST1724I I/O TRACE = OFF  TRACE LENGTH = *NA*
IST924I -----
IST314I END

```

In this example, the numbers correspond to the following information:

- 1.** The Major node is ACTIVE and ONLINE.
- 2.** The READ and WRITE channels are ACTIVE and ONLINE.
- 3.** The data channel is also ACTIVE.

Displaying storage usage

The z/OS Communications Server uses the Common Storage Manager (CSM) to manage storage pools. The recommendation is to increase storage allocations by a minimum of 20 MB for TCP/IP in the CSA definition in IEASYSxx and the FIXED and ECSA definitions in IVTPRMxx.

Check your storage utilization to ensure that you made the correct allocations. Storage usage can also be controlled using the GLOBALCONFIG ECSALIMIT and GLOBALCONFIG POOLLIMIT parameters. ECSALIMIT allows you to specify the maximum amount of extended common service area (ECSA) that TCP/IP can use. POOLLIMIT allows you to specify the maximum amount of authorized private storage that TCP/IP can use within the TCP/IP address space.

The **DISPLAY TCPIP,tcpproc,STOR** command display and the NMI storage statistics report are enhanced to distinguish the common storage that is used by dynamic LPA for load modules from the ECSA storage that is used for control blocks.

You can also use the MVS command **D TCPIP,tcpproc,STOR** to display TCP/IP storage usage, as illustrated in Example 3-30.

Example 3-30 Results of storage display

```
D TCPIP,TCPIPA,STOR
EZZ8453I TCPIP STORAGE
EZZ8454I TCPIPA STORAGE CURRENT MAXIMUM LIMIT
EZZ8455I TCPIPA ECSA 2850K 3270K NOLIMIT
EZZ8455I TCPIPA POOL 9035K 9041K NOLIMIT
EZZ8455I TCPIPA 64-BIT COMMON 1M 1M NOLIMIT
EZZ8455I TCPIPA ECSA MODULES 7451K 7451K NOLIMIT
EZZ8459I DISPLAY TCPIP STOR COMPLETED SUCCESSFULLY
```

Verifying TCPIP.DATA statement values in z/OS

To display which TCPIP.DATA statement values are being used and where they are being obtained from, use trace resolver output. You can obtain trace resolver output at your TSO screen by issuing the following TSO commands:

```
alloc f(sysctcpt) dsn(*)
READY
netstat up
READY
free f(sysctcpt)
READY
```

Tip: When directing trace resolver output to a TSO terminal, define the screen size to be only 80 columns wide. Otherwise, trace output is difficult to read.

Verifying TCPIP.DATA statement values in z/OS UNIX

To display which TCPIP.DATA statement values are being used and where they are being obtained from, use trace resolver output. You can obtain trace resolver output by issuing the following z/OS UNIX shell commands:

```
#
export RESOLVER_TRACE=stdout
#
onetstat -u
#
```

```
set -A RESOLVER_TRACE
```

Verifying PROFILE.TCPIP

Many configuration values specified within the PROFILE.TCPIP file can be verified with the TSO NETSTAT or z/OS UNIX **onetstat** commands. To verify the physical network and hardware definitions, use the D TCPIP,N,DEV, NETSTAT DEVLINKS or **onetstat -d** commands. To see operating characteristics, use z/OS displays, namely NETSTAT CONFIG or **onetstat -f**.

Verifying interfaces with PING and TRACERTE

PING and TRACERTE can be used in the TSO environment to verify adapters or interfaces attached to the z/OS host. In the z/OS UNIX environment, **oping** and **otracer** can be used with identical results. For information about the syntax and output of the commands, refer to *z/OS Communications Server: IP System Administrator's Commands*, SC31-8781. Given that your PROFILE.TCPIP file contains the interfaces of your installation and that the TCPIP.DATA file contains the correct TCPIPJOBNAME, the TCP/IP address space is configured and you can go on to configuring routes, servers, and so on.

3.7 Reconfiguring the system with z/OS commands

The z/OS Communications Server provides the VARY OBEYFILE command to change the running TCP/IP configuration dynamically. This command replaces the OBEYFILE TSO command.

The VARY command is an z/OS Console command. It allows you to add, delete, or completely redefine all devices dynamically, as well as change TN3270 parameters, routing, and almost any TCP/IP parameter in the profile. These changes are in effect until the TCP/IP started task is started again, or another VARY OBEYFILE command overrides them.

Authorization is through the user's RACF profile containing the MVS.VARY.TCPIP.OBEYFILE definition. There is no OBEY statement in the PROFILE.TCPIP, which in earlier MVS TCP/IP implementations provided authorization.

For further details about the VARY OBEYFILE command, see *z/OS CS: IP System Administrator's Commands*, SC31-8781. For more information about RACF definitions, see *IBM z/OS Communications Server TCP/IP Implementation Volume 2: Standard Applications*, SG24-7897.

3.7.1 Deleting a device and adding or changing a device

You can use the OBEYFILE command to reconfigure the devices being used by the stack. Reconfiguration could be the deletion of existing devices, the addition of new devices, or the redefinition of an existing device. The syntax of the statements for OBEYFILE processing is the same as that being used in PROFILE.TCPIP.

Device reconfiguration is a three-step process:

1. Stop the device with an z/OS console command (VARY STOP) or with a VARY OBEYFILE that names a data set in which the STOP command is defined.
2. Activate an OBEYFILE that deletes the links and the devices.
3. Activate an OBEYFILE that adds the new or changed links and devices and then starts them.

Deleting and adding back a device

If you want to delete a device, then the order of the steps that you perform is important. The DELETE statement in PROFILE.TCPIP allows you to remove LINK, DEVICE, and PORT or PORTRANGE definitions. You must delete a resource that is defined using the INTERFACE statement using the DELETE parameter.

The sequence for deleting and adding back a resources that was defined using the INTERFACE statement is as follows:

1. Stop the device.
2. Delete the interface.
3. Add the new or changed interface.
4. Start the device.

Use the following steps to delete and add back a resource that was defined using the DEVICE, LINK, or HOME statements:

1. Stop the device.
2. Remove the HOME address by excluding it from the full stack's HOME list.
3. Delete the link.
4. Delete the device.
5. Add the new or changed device.
6. Add the new or changed link.
7. Add the HOME statements for the full stack.
8. Add the full gateway statements for the stack if you are using static routing.
9. Start the device.

3.7.2 Modifying a device

In this example, we want to change the IP address of our OSA-Express interface/device OSA2080I from 10.1.2.11 to 10.1.2.14. This process involves stopping and deleting the current interface or device, and then redefining and restarting it.

Note: You can delete and redefine OSA-Express resources defined with either the INTERFACE statement or the DEVICE, LINK, or HOME statements by following the same procedure but by creating different OBEYFILE commands. Because the INTERFACE statement is now the preferred way of defining OSA devices, we use that procedure first in the following examples.

Example 3-31 and Example 3-32 show the interface OSA2080I, or link OSA2080L, active with associated IP address 10.1.2.11.

Example 3-31 Displays netstat device before deletion (for INTERFACE defined)

```
D TCPIP,TCPIPA,N,DE
..... Lines deleted
INTFNAME: OSA2080I      INTFTYPE: IPAQENET  INTFSTATUS: READY
  PORTNAME: OSA2080    DATAPATH: 2082      DATAPATHSTATUS: READY
  CHPIDTYPE: OSD
  SPEED: 0000001000
  IPBROADCASTCAPABILITY: NO
  VMACADDR: 02000C776873  VMACORIGIN: OSA  VMACROUTER: LOCAL
  ARPOFFLOAD: YES        ARPOFFLOADINFO: YES
  CFGMTU: 1492           ACTMTU: 1492
  IPADDR: 10.1.2.11/24
  VLANID: 10             VLANPRIORITY: DISABLED
  DYNVLANREGCFG: NO      DYNVLANREGCAP: YES
  READSTORAGE: GLOBAL (4096K)
  INBPERF: BALANCED
  CHECKSUMOFFLOAD: YES   SEGMENTATIONOFFLOAD: YES
  SECCLASS: 255          MONSYSPLEX: NO
  ISOLATE: NO            OPTLATENCYMODE: NO
MULTICAST SPECIFIC:
  MULTICAST CAPABILITY: YES
  GROUP                REFCNT            SRCFLTMD
  -----
  224.0.0.1            0000000001        EXCLUDE
  SRCADDR: NONE
..... Lines deleted
```

Example 3-32 Display netstat home before deletion (for DEVICE/LINK/HOME defined)

```
D TCPIP,TCPIPA,N,HO
..... Lines deleted
INTFNAME:  OSA2080I
  ADDRESS:  10.1.2.11
  FLAGS:
..... Lines deleted
```

Notice the address of OSA2080I (10.1.2.11). We needed to change this in the running system by stopping, deleting, redefining, and adding back the OSA-Express device and link and home address.

Because the STOP command is executed as the last statement within an OBEYFILE regardless of its position within the file, you cannot execute STOP and DELETE in one step. Trying to do so will result in the error messages. You should stop the interface or device with the console command, as shown in Example 3-33.

Example 3-33 Command to stop the interface or device

```
V TCPIP,TCPIPA,STOP,OSA2080I
EZZ0060I PROCESSING COMMAND: VARY TCPIP,TCPIPA,STOP,OSA2080I
EZZ0053I COMMAND VARY STOP COMPLETED SUCCESSFULLY
EZZ4341I DEACTIVATION COMPLETE FOR INTERFACE OSA2080I
EZZ4315I DEACTIVATION COMPLETE FOR DEVICE OSA2080I
```

Then, delete it from the stack, as shown in Example 3-34.

Example 3-34 Command to delete the interface or device

```
V TCPIP,TCPIPA,0,TCPIPA.TCPPARMS(OBDELINT)
EZZ0060I PROCESSING COMMAND: VARY TCPIP,TCPIPA,0,TCPIPA.TCPPARMS(OBDELINT)
EZZ0300I OPENED OBEYFILE FILE 'TCPIPA.TCPPARMS(OBDELINT)'
EZZ0309I PROFILE PROCESSING BEGINNING FOR 'TCPIPA.TCPPARMS(OBDELINT)'
ZZ0316I PROFILE PROCESSING COMPLETE FOR FILE 'TCPIPA.TCPPARMS(OBDELINT)'
EZZ0053I COMMAND VARY OBEY COMPLETED SUCCESSFULLY
```

Enter either the NETSTAT DEV or NETSTAT HOME commands to check that the device you wanted to delete is missing from the list.

Example 3-35 and Example 3-36 show the statements necessary to delete the device.

Example 3-35 OBEYFILE member to delete the device OSA2080I (INTERFACE defined)

```
INTERFACE OSA2080I
DELETE
```

Example 3-36 OBEYFILE member to delete the device OSA2080 (DEVICE/LINK/HOME defined)

```
HOME
  10.1.1.10      VIPA1L
  10.1.2.10      VIPA2L
;;10.1.2.11      OSA2080I
  10.1.3.11      OSA20C0I
  10.1.3.12      OSA20E0I
  10.1.2.12      OSA20A0I
  10.1.4.11      IUTIQDF4L
  10.1.5.11      IUTIQDF5L
  10.1.6.11      IUTIQDF6L
;
DELETE LINK OSA2080I
DELETE DEVICE OSA2080
```

Note: With DEVICE/LINK/HOME defined devices, you have to provide the *complete* HOME definition that excludes the device that you want to delete, because the new HOME statement *replaces* the existing one. This step is *not* necessary with devices defined using the INTERFACE statement.

Then, add either the interface or the device and link back with the changed address definition **3**, as shown in Example 3-37 and Example 3-38.

Example 3-37 OBEYFILE member to add the interface

```
INTERFACE OSA2080I
  DEFINE IPAQENET
  PORTNAME OSA2080
  IPADDR 10.1.2.14/24 3;
  START OSA2080I
```

Example 3-38 OBEYFILE member to add the device (ADDA30)

```
DEVICE OSA2080 MPCIPA
  LINK OSA2080I IPAQENET OSA2080 VLANID 10
;
HOME
  10.1.1.10 VIPA1L
  10.1.2.10 VIPA2L
  10.1.2.14 OSA2080I 3
  10.1.3.11 OSA20C0I
  10.1.3.12 OSA20E0I
  10.1.2.12 OSA20A0I
  10.1.4.11 IUTIQDF4L
  10.1.5.11 IUTIQDF5L
  10.1.6.11 IUTIQDF6L
```

Issue the command shown in Example 3-39 to add the device and link associated with its own IP address.

Example 3-39 Adding the device and link

```
V TCPIP,TCPIPA,0,TCPIPA.TCPPARMS(OBADDINT)
EZZ0060I PROCESSING COMMAND: VARY TCPIP,TCPIPA,0,TCPIPA.TCPPARMS(OBADDINT)
EZZ0300I OPENED OBEYFILE FILE 'TCPIPA.TCPPARMS(OBADDINT)'
EZZ0309I PROFILE PROCESSING BEGINNING FOR 'TCPIPA.TCPPARMS(OBADDINT)'
ZZ0316I PROFILE PROCESSING COMPLETE FOR FILE 'TCPIPA.TCPPARMS(OBADDINT)'
EZZ0053I COMMAND VARY OBEY COMPLETED SUCCESSFULLY
```

Then, follow with a display to verify the addition to the stack, as shown in Example 3-40.

Example 3-40 Display with OSA2080 using a new address

```
D TCPIP,TCPIPE,N,HOME
..... Lines deleted
LINKNAME: OSA2080I
  ADDRESS: 10.1.2.14
  FLAGS:
..... Lines deleted
```

3.8 Job log versus syslog as diagnosis tool

In the past, the TCP/IP job log was used to detect problems. Most procedures now send messages to the syslogd daemon or the MVS console log. Refer to *IBM z/OS Communications Server TCP/IP Implementation Volume 2: Standard Applications*, SG24-7897 for more information about the syslog daemon. Individual server documentation also provides information about diagnosis.

3.9 Message types: Where to find them

For an explanation of z/OS UNIX and TCP/IP messages or SNA sense codes, refer to the publications listed in Table 3-1.

Table 3-1 Messages and return code publications

Message type	Publication
Messages with prefix BPX	<i>z/OS MVS System Messages, Vol 3 (ASB-BPX)</i> , SA22-7633
Messages with prefix EZA	For Communications Server for z/OS IP, refer to <i>z/OS Communications Server: IP Messages Volume 1 (EZA)</i> , SC31-8783
Messages with prefix EZB	For Communications Server for z/OS IP, refer to <i>z/OS Communications Server: IP Messages Volume 2 (EZB, EZD)</i> , SC31-8784
Messages with prefix EZY	For Communications Server for z/OS IP, refer to <i>z/OS Communications Server: IP Messages Volume 3 (EZY)</i> , SC31-8785
Messages with prefix EZZ and SNM	For Communications Server for z/OS IP, refer to <i>z/OS Communications Server: IP Messages Volume 4 (EZZ, SNM)</i> , SC31-8786
Messages with prefix FOMC, FOMM, FOMO, FSUC, and FSUM	<i>z/OS UNIX System Services Messages and Codes</i> , SA22-7807
Eight-digit SNA sense codes and DLC codes	<i>z/OS Communications Server: IP and SNA Codes</i> , SC31-8791
UNIX System Services return codes and reason codes	<i>z/OS UNIX System Services Messages and Codes</i> , SA22-7807

3.10 Additional information

When you install and customize the Communications Server for z/OS IP, it can be very helpful to have the following documentation and product publications available:

- ▶ Implementation and migration plans, fallback plans, and test plans that you have created and customized for your environment
- ▶ Printouts of procedures and data sets that you will be using for the implementation
- ▶ *z/OS Program Directory, Program Number 5694-A01*, GI10-0670
- ▶ *z/OS XL C/C++ Run-Time Library Reference*, SA22-7821
- ▶ *z/OS Migration*, GA22-7499
- ▶ *z/OS Communications Server: IP Configuration Guide*, SC31-8775
- ▶ *z/OS Communications Server: IP Configuration Reference*, SC31-8776

- ▶ *z/OS Communications Server: IP Messages Volume 1 (EZA)*, SC31-8783
- ▶ *z/OS Communications Server: IP Messages Volume 2 (EZB, EZD)*, SC31-8784
- ▶ *z/OS Communications Server: IP Messages Volume 3 (EZY)*, SC31-8785
- ▶ *z/OS Communications Server: IP Messages Volume 4 (EZZ, SNM)*, SC31-8786
- ▶ *z/OS Communications Server: IP and SNA Codes*, SC31-8791
- ▶ *OSA-Express Customer's Guide and Reference*, SA22-7935
- ▶ *z/OS UNIX System Services Planning*, GA22-7800
- ▶ *z/OS UNIX System Services User's Guide*, SA22-7801
- ▶ *z/OS UNIX System Services Messages and Codes*, SA22-7807
- ▶ *z/OS MVS System Messages, Vol 1 (ABA-AOM)*, SA22-7631
- ▶ *z/OS MVS System Messages, Vol 2 (ARC-ASA)*, SA22-7632
- ▶ *z/OS MVS System Messages, Vol 3 (ASB-BPX)*, SA22-7633
- ▶ *z/OS MVS System Messages, Vol 4 (CBD-DMO)*, SA22-7634
- ▶ *z/OS MVS System Messages, Vol 5 (EDG-GFS)*, SA22-7635
- ▶ *z/OS MVS System Messages, Vol 6 (GOS-IEA)*, SA22-7636
- ▶ *z/OS MVS System Messages, Vol 7 (IEB-IEE)*, SA22-7637
- ▶ *z/OS MVS System Messages, Vol 8 (IEF-IGD)*, SA22-7638
- ▶ *z/OS MVS System Messages, Vol 9 (IGF-IWM)*, SA22-7639
- ▶ *z/OS MVS System Messages, Vol 10 (IXC-IZP)*, SA22-7640



Connectivity

In today's networked world, the usability of a computer system is defined by its connectivity. While there are many ways for TCP/IP traffic to reach IBM mainframes, this chapter discusses the most commonly used and the most dynamic types of mainframe connectivity.

Detailed topics regarding these interfaces are provided, including useful implementation information, design scenarios, and setup examples.

This chapter discusses the following topics.

Section	Topic
4.1, "What is connectivity" on page 118	Network connectivity options supported by z/OS and Communications Server TCP/IP, IBM System z servers, and key characteristics of VLAN implementation
4.2, "Recommended interfaces" on page 120	Recommended interfaces supported by System z hardware and z/OS Communications Server
4.3, "Connectivity for the z/OS environment" on page 136	Basic implementation information for z/OS and Communications Server when connecting to the immediate LAN environment
4.4, "OSA-Express QDIO connectivity" on page 139	Configuration examples, with dependencies, considerations, and our recommendations for an OSA-Express interface
4.5, "OSA-Express QDIO connectivity with Connection Isolation" on page 156	Configuration examples, with dependencies, considerations, and our recommendations for isolating traffic across a shared OSA port
4.6, "HiperSockets connectivity" on page 180	Configuration examples, with dependencies, considerations, and our recommendations for a HiperSockets interface
4.7, "Dynamic XCF connectivity" on page 188	Configuration examples, with dependencies, considerations, and our recommendations for a dynamic XCF interface
4.8, "Controlling and activating devices" on page 195	Commands to start and stop devices, as well as activate modified device definitions
4.9, "Problem determination" on page 197	How to determine why certain connectivity options are not working

4.1 What is connectivity

Connectivity is the pipeline through which data is exchanged between clients and servers through physical and logical communication interfaces and the network. IBM System z servers provide a wide range of interface options for connecting your z/OS system to an IP network or to another IP host. Some interfaces offer point-to-point or point-to-multipoint connectivity. Others support Local Area Network (LAN) connectivity.

Figure 4-1 depicts the physical interfaces (and device types) provided by System z servers. The physical network interface is enabled through z/OS Communications Server (TCP/IP) definitions.

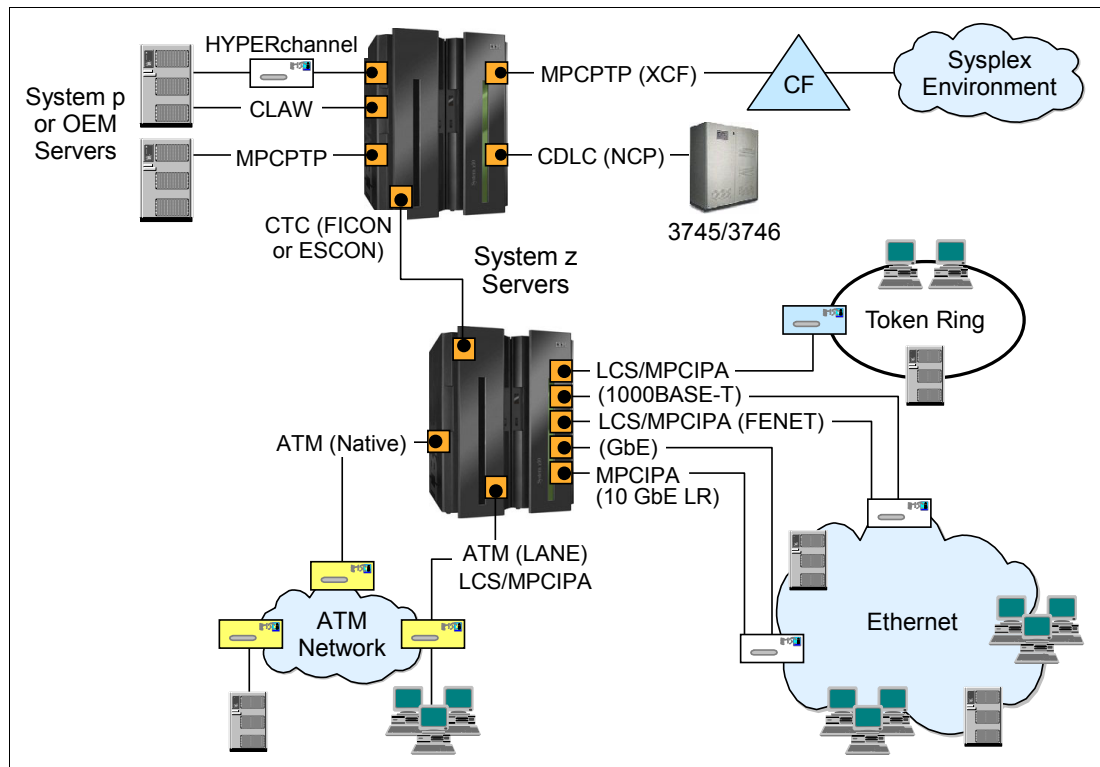


Figure 4-1 System z: physical interfaces

4.1.1 System z network connectivity

System z network connectivity is handled by the physical and logical interfaces to enable the transport of IP datagrams. Using the OSI model as an example, it spans Layer 1 (physical layer) and Layer 2 (data link control layer). The z/OS Communications Server supports several types of interfaces connecting to different networking environments. These environments vary from point-to-point connections (such as MPCPTP, CTC, and CLAW), to LAN connections (such as LCS and MPCIPA).

The supported IPv4 interfaces are listed in Table 4-1.

Table 4-1 System z network interfaces

Interface type	Attachment type	Protocol type	Description
Asynchronous transfer mode (ATM)	ATM Native mode through OSA-Express	ATM network	Enables TCP/IP to send data to an ATM network using an OSA-Express ATM adapter.
Channel data link control (CDLC)	Network Control Program through 3745/3746	Point-to-point	ESCON® attachments can be used to provide native IP transport between the 3746 IP and host systems running the z/OS Communications Server.
Common link access to workstation (CLAW)	IBM System p® Channel attached routers	Point-to-point Point-to-Multipoint	Provides access from IBM System p server directly to a TCP/IP stack over a channel. Can also be used to provide connectivity to other vendor platforms.
Channel-to-channel (CTC)	FICON/ESCON channel	Point-to-point	Provides access to TCP/IP hosts by way of a CTC connection established over a FICON or ESCON channel.
HYPERchannel	Series A devices	Point-to-Multipoint	Provides access to TCP/IP hosts by way of a series A devices and series DX devices that function as series A devices.
LAN Channel Station (LCS)	OSA-Express: <ul style="list-style-type: none"> ▶ 1000BASE-T ▶ Fast Ethernet ▶ Token Ring ▶ ATM Native and LAN Emulation 	LAN: <ul style="list-style-type: none"> ▶ IEEE802.3 ▶ IEEE802.3 ▶ IEEE802.5 ▶ ATM network 	A variety of channel adapters support a protocol called the LCS. The most common are OSA-Express features configured as CHPID type OSE. LCS supports native IP flows on z/OS. CHPID type OSE also supports the Link Station Architecture (LSA) protocol, which supports native SNA flows for VTAM on z/OS.
MultiPath Channel IP Assist (MPCIPA)	HiperSockets ^a OSA-Express: <ul style="list-style-type: none"> ▶ 10 Gigabit Ethernet ▶ Gigabit Ethernet ▶ 1000BASE-T ▶ Fast Ethernet ▶ Token Ring ▶ ATM LAN Emulation 	Internal LAN LAN: <ul style="list-style-type: none"> ▶ IEEE802.3 ▶ IEEE802.3 ▶ IEEE802.3 ▶ IEEE802.3 ▶ IEEE802.5 ▶ ATM network 	Provides access to TCP/IP hosts, using OSA-Express in Queued Direct I/O (QDIO) mode and HiperSockets using the internal Queued Direct I/O (iQDIO).
MultiPath Channel Point-to-Point (MPCPTP)	IUTSAMEH (XCF link)	Point-to-point	Provides access to directly connect z/OS hosts or z/OS LPARs, or by configuring it to utilize Coupling Facility links (if it is part of a sysplex).
SAMEHOST (Data Link Control)	SNALINK LU0 SNALINK LU6.2 X25NPSI	Point-to-point Point-to-point X.25 network	Enables communication between z/OS Communications Server IP and other servers running on the same MVS image.

- a. Can also be used in conjunction with DYNAMICXCF

For further information about these protocols, refer to *z/OS Communications Server: IP Configuration Reference*, SC31-8776.

4.2 Recommended interfaces

This section discusses the recommended interfaces supported by System z hardware and z/OS Communications Server. We highly recommend their use because they deliver the best throughput and performance, as well as offer the most flexibility and highest levels of availability. These interfaces include:

- ▶ OSA-Express
- ▶ HiperSockets
- ▶ Dynamic Cross-system Coupling Facility (dynamic XCF)

4.2.1 High-bandwidth and high-speed networking technologies

z/OS Communications Server supports high-bandwidth and high-speed networking technologies provided by OSA-Express and HiperSockets.

- ▶ The OSA-Express features comply with the most commonly used IEEE standards, used in LAN environments.
- ▶ HiperSockets is used for transporting IP traffic between TCP/IP stacks running in logical partitions (LPARs) within a System z server at memory speed.

Both interfaces use the System z I/O architecture called queued direct input/output (QDIO).

QDIO is a highly efficient data transfer mechanism that satisfies the increasing volume of applications and bandwidth demands. It dramatically reduces system overhead, and improves throughput by using system memory queues and a signaling protocol to directly exchange data between the OSA-Express microprocessor and network software, using data queues in main memory and utilizing Direct Memory Access (DMA).

The components that make up QDIO are Direct Memory Access (DMA), Priority Queueing, dynamic OSA Address Table building, LPAR-to-LPAR communication, and Internet Protocol (IP) Assist functions.

HiperSockets implementation is based on the OSA-Express QDIO protocol, hence the name internal QDIO (iQDIO). The System z microcode for HiperSockets emulates the link control layer of an OSA-Express QDIO interface. The communication is through system memory of the server using I/O queues. IP traffic is transferred at memory speeds between LPARs, eliminating the I/O subsystem overhead and external network delays.

Recommendation: Some OSA-Express features also support LCS (known as non-QDIO mode). However, we recommend the use of QDIO mode in conjunction with the OSA-Express Ethernet features wherever possible.

With QDIO, I/O interrupts and I/O path-lengths are minimized, resulting in significantly improved performance versus non-QDIO mode, reduction of System Assist Processor (SAP) utilization, improved response time, and server cycle reduction.

z/OS Communications Server can only transport IP traffic over OSA-Express in QDIO mode and HiperSockets. However, SNA can be transported over IP connections using encapsulation technologies such as Enterprise Extender (EE) and TN3270.

For more information about EE, refer to *Enterprise Extender Implementation Guide*, SG24-7359. For TN3270 details, refer to *IBM z/OS V1R12 Communications Server TCP/IP Implementation Volume 2: Standard Applications*, SG24-7897.

4.2.2 OSA-Express (MPCIPA)

As mentioned, the OSA-Express can use the I/O architecture called QDIO when defined as channel type (CHPID) OSD. QDIO provides a highly optimized data transfer interface that eliminates the need for channel command words (CCWs) and interrupts during data transmission, resulting in accelerated TCP/IP packet transmission. This is done by providing a data queue between TCP/IP and the OSA-Express. OSA-Express uses a direct memory access (DMA) protocol to transfer the data to and from the TCP/IP stack.

The OSA-Express also provides the offloading of IP processing from the host, which is called IP assist (IPA). With IP assist, the OSA-Express offloads the following processing from the host:

- ▶ All MAC handling is done in the card. The TCP/IP stack no longer has to fully format the datagrams for LAN-specific media.
- ▶ ARP processing for identifying the physical address.
- ▶ Packet filtering, screening, and discarding of LAN packets.

Table 4-2 lists the OSA-Express3, OSA-Express2, and OSA-Express Ethernet features that are available on the System z servers. The mode of operation in which they can run and the necessary TCP/IP and VTAM definition types are included.

Table 4-2 OSA-Express features to support native TCP/IP data flows

OSA-Express feature	Operation mode	TCP/IP device type	TCP/IP link type	VTAM definitions
10 GbE LR	QDIO	MPCIPA	IPAQENET	TRLE
GbE	QDIO	MPCIPA	IPAQENET	TRLE
1000BASE-T	QDIO	MPCIPA	IPAQENET	TRLE
	Non-QDIO	LCS	ETHERNet, 802.3, or ETHEROR802.3	N/A

Note: The 1000Base-T feature can also support native SNA data flows to VTAM when configured in Non-QDIO mode. The VTAM device type protocol is called Link Station Architecture (LSA).

OSA-Express QDIO IPv4 address registration

The Dynamic OSA Address Table (OAT) contains certain active IP addresses displayed in the HOME list of the TCP/IP stack; the addresses are downloaded into the OSA-Express when the interface is started.

z/OS Communications Server registers IPv4 addresses in the OSA OAT for two distinct purposes:

- ▶ Inbound routing
- ▶ ARP offload

Several factors contribute to the types of IPv4 addresses in a TCP/IP stack that are registered in the OAT. These factors are summarized as follows:

- ▶ Which type of definition statement defines the characteristics of the adapter interface in the TCP/IP stack: DEVICE/LINK or INTERFACE?
- ▶ Does the adapter interface definition include a virtual MAC (VMAC) keyword?
- ▶ Has VMAC ROUTEALL been coded or defaulted for the adapter interface?
- ▶ Has VMAC ROUTLCL been coded for the adapter interface?

Depending on these four factors, different addresses are registered in the OSA as described here for the purposes of inbound routing and ARP offload:

- ▶ Inbound routing
 - For INTERFACE statement with VMAC ROUTEALL, we do not register any IP addresses for the purpose of inbound routing. We only register an IP address for the purpose of supporting ARP offload.
 - For INTERFACE without VMAC ROUTEALL or for DEVICE/LINK, we register the entire home list for the purpose of inbound routing. (For DEVICE/LINK with VMAC ROUTEALL, this registration is extraneous and is not used at all.)
- ▶ ARP offload
 - We always register the home IP address for the purpose of ARP offload.
 - If you have multiple OSAs on the same (V)LAN or Physical Network (PNET), and ARP takeover is in effect, then we register the IP address of the interface for which we are taking over connection responsibility.
 - We also register VIPAs for ARP offload purposes as follows:
 - For the INTERFACE statement with subnet mask configured on the statement, we register only the VIPAs that are in the same subnet as the OSA.
 - For the INTERFACE statement without a subnet mask coded on it. For DEV/LINK, we register all the active VIPAs in the Home list.

For both of the above bullets, if there are multiple OSAs on the same (V)LAN or Physical Network (PNET), we register these VIPAs on only one of the OSAs.

Note about displaying registered addresses: OSA/SF has a Get OAT function that retrieves the registered IP addresses in the OAT. However, the displayed table is incomplete, containing only a limited number of the addresses that the stack has registered with the OSA device. When performing problem determination for the OSA, do not assume that OSA/SF is showing you everything you need to know. You may have to solicit the help of Level 2 defect support to see everything that has really been registered in the OSA.

OSA-Express VLAN support

The OSA-Express Ethernet features also support IEEE standards 802.1p/q (priority tagging and VLAN identifier tagging). Deploying VLAN IDs enables a physical LAN to be partitioned or subdivided into discrete virtual LANs. This support is provided by the z/OS TCP/IP stack and OSA-Express in QDIO mode. It allows a TCP/IP stack to register specific single or multiple VLAN IDs for both IPv4 and IPv6 for the same OSA-Express port. Note that the VLAN IDs for IPv4 can be different than the VLAN ID for IPv6.

Note: The INTERFACE statement is required if one stack is going to attach to multiple VLANs through a single OSA port.

When a VLAN ID is configured for an OSA-Express interface in the TCP/IP stack, the following occurs:

- ▶ The TCP/IP stack becomes VLAN-enabled, and the OSA-Express port is considered to be part of a VLAN.
- ▶ During activation, the TCP/IP stack registers the VLAN ID value to the OSA-Express port.
- ▶ A VLAN tag is added to all outbound packets.
- ▶ The OSA-Express port filters all inbound packets based on the configured VLAN ID.

If the TCP/IP stack is also configured with PRIRouter or SECRouter for an OSA-Express port that has a VLAN ID defined, then the stack serves as an IP router for the configured VLAN ID. If OSA-Express ports are shared across multiple TCP/IP routing stacks, consider using virtual MAC support for your environment instead of the PRIRouter and SECRouter options. See Chapter 6, “VLAN and Virtual MAC support” on page 265 for details.

VLAN support of Generic Attribute Registration Protocol (GVRP)

GVRP is defined in the IEEE 802.1p standard for the control of IEEE 802.1q VLANs. It can be used to help simplify networking administration and management of VLANs. With GVRP support, an OSA-Express3 or OSA-Express2 port can register or de-register its VLAN IDs with a GVRP-capable switch and dynamically update its table as the VLANs change. Support of GVRP is exclusive to System z9® or above and is applicable to all of the OSA-Express3 and OSA-Express2 features when in QDIO mode. Defining DYNVLANREG in the LINK statement of the OSA-Express3 and OSA-Express2 port will enable GVRP.

OSA-Express router support

OSA-Express also provides primary (PRIRouter) and secondary (SECRouter) router support. This function allows a single TCP/IP stack, on a per-protocol (IPv4 and IPv6) basis, to register and act as a router stack based on a given OSA-Express port. Secondary routers can also be configured to provide for conditions in which the primary router becomes unavailable and the secondary router takes over for the primary router.

Figure 4-2 shows how the PRIRouter function works in a shared OSA environment.

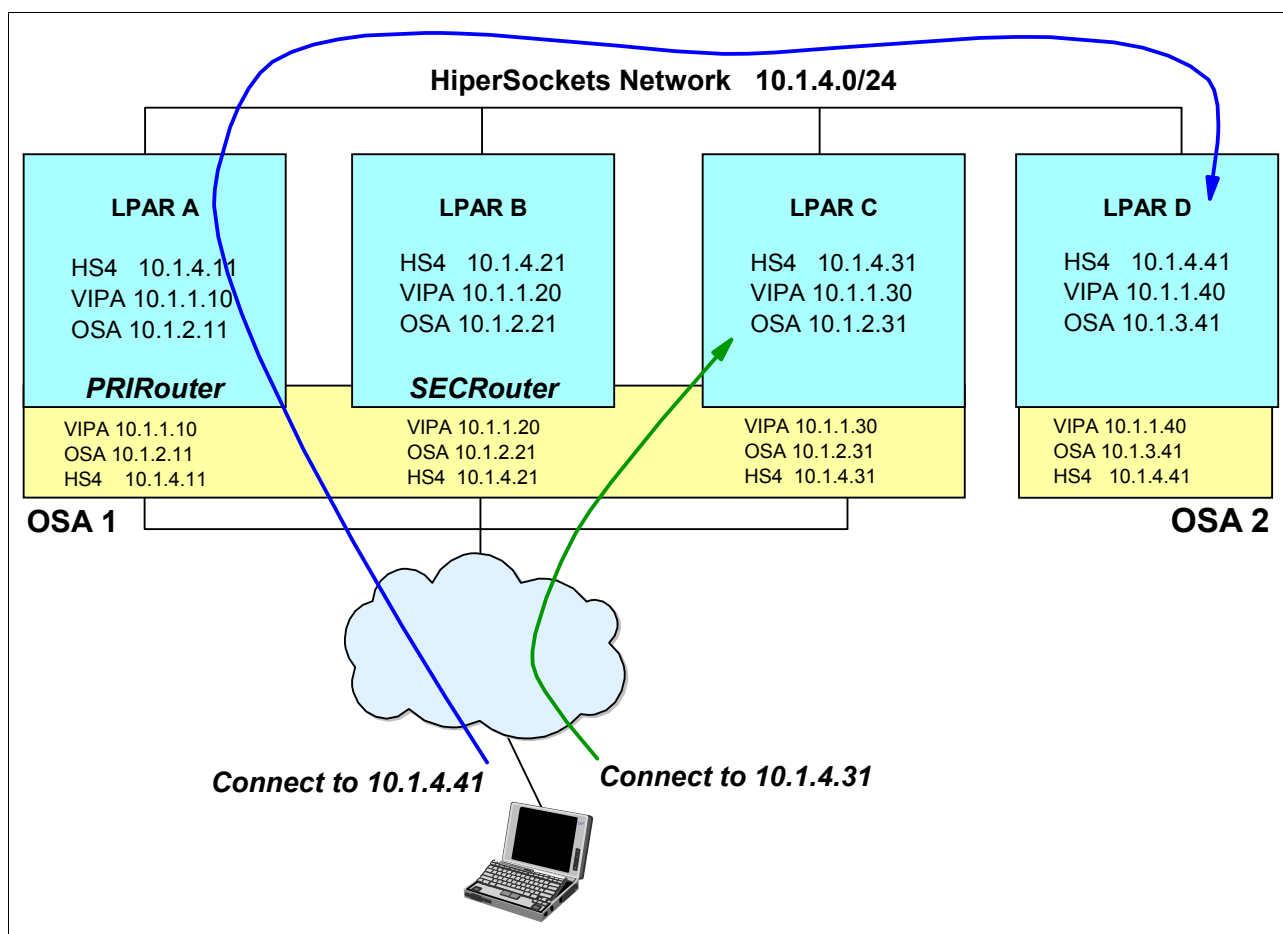


Figure 4-2 How PRIRouter works with a shared OSA

In Figure 4-2, the terminal user connects to 10.1.4.41. Note that each stack sharing OSA1 has registered the IP addresses for VIPAs, OSAs, and the HiperSockets in the OSA Address Table (OAT). However, the address 10.1.4.41 is not represented in OSA1's OAT. Therefore, the packet from the terminal that arrives at OSA1 is sent to the primary routing TCP/IP stack in LPAR A. The TCP/IP stack in LPAR A uses its routing table to forward the packet to LPAR D, where IP address 10.1.4.41 resides.

The connection to IP address 10.1.4.31 is simpler. Because the address is represented in the OAT of OSA1, the OSA can immediately forward the request to the correct TCP/IP stack in LPAR C.

If LPAR A should become unavailable, then the TCP/IP stack in LPAR B or C will take over the routing responsibility for OSA1.

VLAN and primary/secondary router support: VMAC support

The OSA-Express primary router support takes into consideration VLAN ID support (VLAN ID registration and tagging) and interacts with it. OSA-Express supports a primary and secondary router on a per-VLAN basis (per registered VLAN ID).

Therefore, if an OSA interface is configured with a specific VLAN ID and also configured as a primary or secondary router, that stack serves as a router for just that specific VLAN. This

allows each OSA-Express (CHPID) to have a primary router per VLAN. Configuring primary routers (one per VLAN) has many advantages and preserves traffic isolation for each VLAN.

If OSA-Express ports are shared across multiple TCP/IP routing stacks, consider using virtual MAC support for your environment instead of the PRIRouter and SECRouter options. See Chapter 6, “VLAN and Virtual MAC support” on page 265 for details.

High latency network

Streaming a workload over large bandwidth and high latency networks (such as satellite links) is, in general, constrained by the TCP window size. The problem is that it takes time to send data over such a network. At any given point in time, data filling the full window size is in transit and cannot be acknowledged until it starts arriving at the receiver side. The sender can send up to the window size and then must wait for an ACK to advance the window size before the next chunk can be sent.

The left hand side of Figure 4-3 depicts a high-latency network where the TCP window size is too small. The round trip time (RTT) is relatively long and the window size is relatively small. Therefore, the sender fills the window before it receives an ACK for the data at the start of the window. This forces the sender to delay sending additional data until it receives an ACK or a window update. Over a long distance connection, this can cause transmission stalls and suboptimal performance.

The right hand side demonstrates a situation where the window size is large enough for the high-latency network. The sender has not yet sent the last bit of the window size before it receives an ACK for the first bit of the current window. The z/OS TCP maximum windows size is 512K (defined in the TCPMAXRCVBUFRSIZE in the TCPCONFIG section). However, a window size of 512K may not always be enough to achieve this behavior.

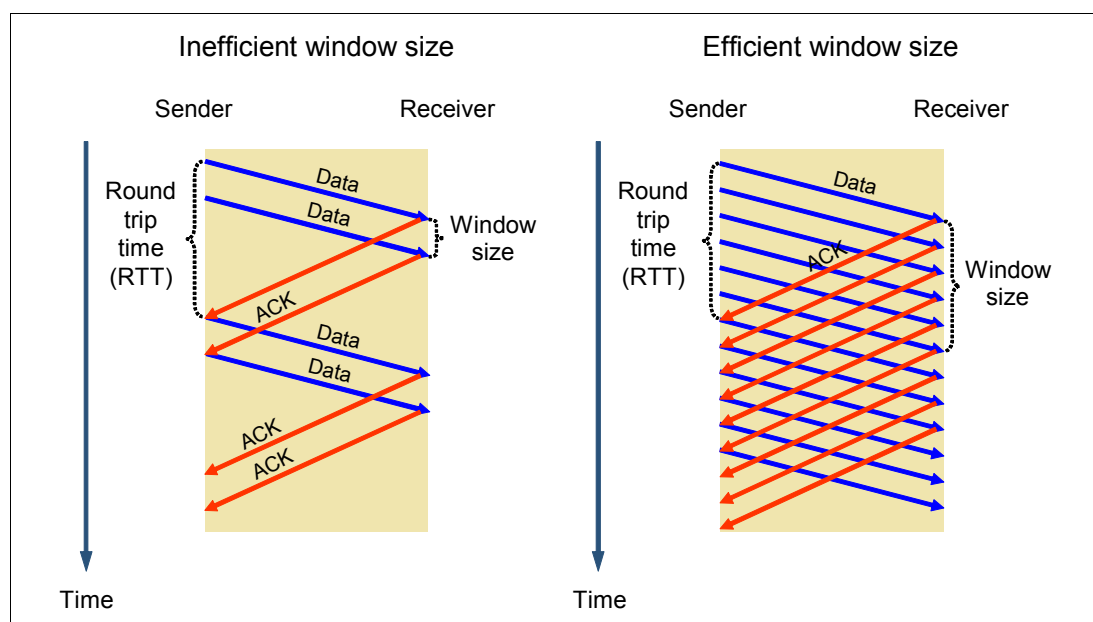


Figure 4-3 High latency network and window size

The solution for high latency networks

z/OS Communications Server implements Dynamic Right Sizing (DRS) to address the problem related to high latency networks. DRS is described in a paper published by the Los Alamos National Laboratory (LANL), which can be found at the following address:

<http://public.lanl.gov/radiant/software/drs.html>

The goal of the DRS function is to keep the pipe full for inbound streaming TCP connections over networks with large capacity and high latency and prevent the sender from being constrained by the receiver's advertised window size.

If a TCP connection uses a receive buffer size larger than 64K, the stack detects a high latency inbound streaming TCP connection and dynamically increases the receive buffer size for the connection (in an attempt to not constrain the sender). This in turn adjusts the advertised receive window and allows window size to grow as high as 2M. In other words, The TCP receive buffer size can grow as high as 2M for certain TCP connections irrespective of the TCP_MAXRCVBUFSIZE value. The stack disables the function for a connection if the application is not keeping up with the pace.

DRS does not take effect for applications that set a value less than 64K on the SO_RCVBUF socket option on SETSOCKOPT().

If TCPRCVBUFSIZE is less than 64K, then DRS does not take effect for applications that do not use the SO_RCVBUF socket option.

Implementation

To configure an OSA-Express3 feature to operate in optimized latency mode, use the INTERFACE statement with the OLM parameter. Because optimized latency mode affects both inbound and outbound interrupts, it supersedes other inbound performance settings set by the INBPERF parameter.

Optimized latency mode is limited to the OSA-Express3 Ethernet feature in QDIO mode running with an IBM System z10®.

Restrictions

You must observe the following restrictions:

- ▶ Traffic that is either inbound over or being forwarded to an OSA-Express3 feature configured to operate in optimized latency mode is not eligible for the accelerated routing provided by HiperSockets Accelerator and QDIO Accelerator.
- ▶ For an OSA-Express3 configured to operate in optimized latency mode, the stack ignores the value coded on the INBPERF parameter. The value assigned to the INBPERF is DYNAMIC.

Guidelines

Because of the operating characteristics of optimized latency mode, other configuration changes might be required:

- ▶ For outbound traffic to gain the benefit of optimized latency mode, direct traffic to priority queues 1, 2, or 3 using the WLM_PRIORITYQ parameter in the GLOBALCONFIG statement or using Policy Agent and configuring a policy with the SetSubnetPrioToMask statement.
- ▶ Although an OSA-Express feature supports multiple outbound write priority queues, outbound optimized latency mode is performed only for traffic on priority queue 1 (priority level 1). The TCP/IP stack combines all the traffic directed to priority queues 1, 2, and 3 into priority queue 1 for any OSA-Express3 feature operating in optimized latency mode.
- ▶ Configure the WLM_PRIORITYQ parameter with no subparameters, which assigns a default mapping of service class importance levels to OSA-Express outbound priority queues. This default mapping directs traffic assigned to the higher priority service class importance levels 1–4 to queues that operate in optimized latency mode, and enables the appropriate types of traffic to benefit from optimized latency mode.

- Ensure that there are no more than four concurrent users of an OSA-Express3 feature that is configured with optimized latency mode.
- When enabling multipath routing using the PERPACKET option, do not configure a multipath group that contains an OSA-Express3 feature configured with optimized latency mode and any other type of device.

For more information regarding OSA-Express features and capabilities, refer to *OSA-Express Implementation Guide*, SG24-5948.

OSA multiple inbound queue support

Outbound traffic separation (assignment to specific priority queue) on the multiple write queues can be accomplished by using Policy Agent and configuring a policy with the SetSubnetPrioTosMask statement, and by using the WLMRIORITYQ parameter on the GLOBALCONFIG statement. Each priority queue is processed independently of the others.

The left side of Figure 4-4 depicts OSA single inbound queue support. All inbound QDIO traffic is received on a single read queue regardless of the data type. This includes both batch and interactive traffic and both traffic destined for this TCP/IP stack and traffic to be forwarded by this TCP/IP stack. The maximum amount of storage available for inbound traffic is limited to the read buffer size (64K read SBALs) times the maximum number of read buffers (126). Multiple processes only run for inbound traffic when data is accumulating on the read queue – typically during burst periods when z/OS Communications Server is not keeping up with the OSA. This can cause bulk data packets for a single TCP connection to arrive at the TCP layer out of order. Each time the TCP layer on the receiving side sees out of order data, it transmits a duplicate ACK. A single process is used to package the data, queue it, and schedule the TCP/IP stack to process it. This same process also performs acceleration functions, such as Sysplex Distributor connection routing accelerator. The TCP/IP stack separates the traffic types to be forwarded to the appropriate stack component that will process them.

For these reasons, z/OS Communications Server is becoming the bottleneck as OSA-Express3 10 GbE nears line speed. z/OS Communications Server is injecting latency and increasing processor utilization.

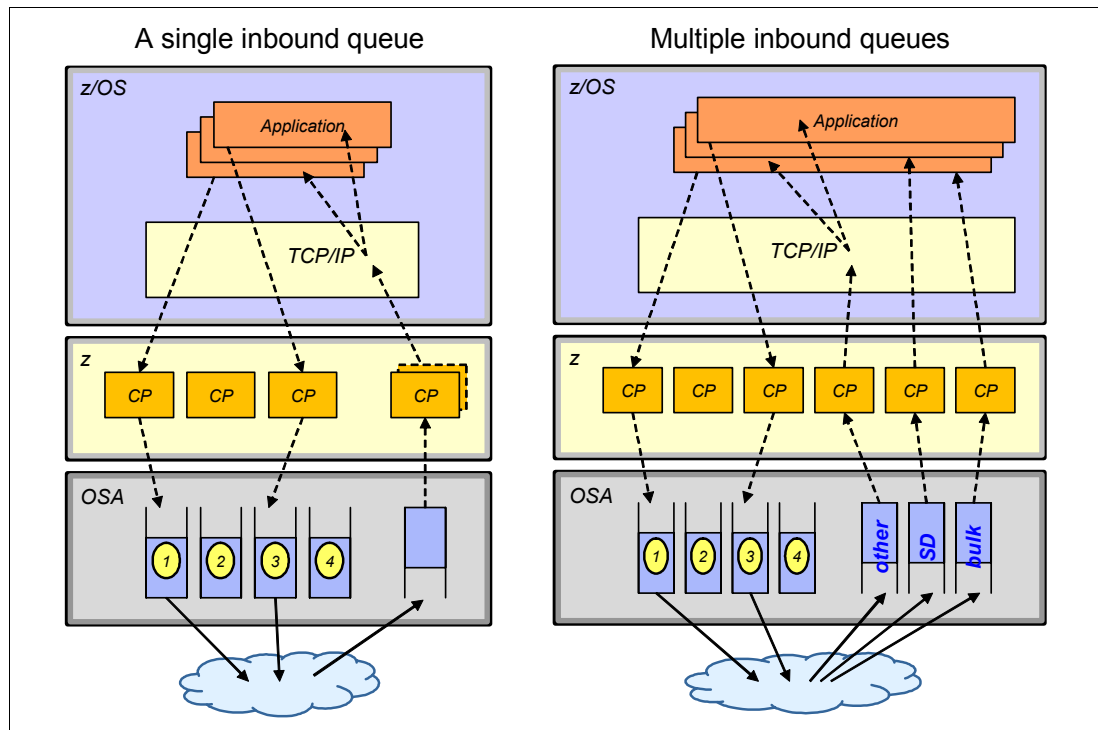


Figure 4-4 QDIO inbound queueing

The solution for the bottleneck on a single read queue

z/OS Communications Server supports inbound traffic separation using multiple read queues. The right hand side of Figure 4-4 on page 128 depicts OSA multiple inbound queue support. TCP/IP will register with OSA which traffic is to be received on each read queue. The OSA-Express Data Router function routes traffic to the correct queue. Each read queue can be serviced by a separate process. The primary input queue is used for general traffic. One or more ancillary input queues (AIQs) are used for specific traffic types.

The supported traffic types are streaming bulk data and sysplex distributor. Examples of bulk data traffic are FTP, TSM, NFS, and TDMF®. Both IP versions are supported for all types of traffic.

With bulk data traffic separated onto its own read queue, TCP/IP will service the bulk data queue from a single processor. This solves the out of order delivery issue. With sysplex distributor traffic separated onto its own read queue, it can be efficiently accelerated or presented to the target application. All other traffic is processed simultaneously with the bulk data and sysplex distributor traffic.

The dynamic LAN idle timer is updated independently for each read queue. This ensures the most efficient processing of inbound traffic based on the traffic type.

Implementation

The QDIO inbound workload queueing function is enabled with the INBPERF DYNAMIC WORKLOADQ setting on IPAQENET and IPAQENET6 INTERFACE statements. WORKLOADQ is not supported for INBPERF DYNAMIC on IPAQENET LINK statements. The VMAC parameter can be specified with or without macaddr.

For more information, refer to the IPAQENET INTERFACE and IPAQENET6 INTERFACE statements in *z/OS Communications Server: IP Configuration Reference*, SC31-8776.

Verification

Refer to a WorkloadQueueing field in the Netstat DEvlinks/-d report to see whether the QDIO inbound workload queueing function is enabled. This information can also be returned by the Getlfs callable NMI.

Moreover, you can use other commands to obtain more information about the QDIO inbound workload queueing function for the QDIO interface. The output for the Display ID=trlename and Display TRL,TRLE=trlename commands shows whether this function is in use for the QDIO interface as follows:

- ▶ For each input queue, it includes the queue ID and queue type in addition to the read storage. The queue type is PRIMARY for the primary input queue, BULKDATA for the bulk data AIQ, and SYSDIST for the sysplex distributor connection routing AIQ.
- ▶ The queue type value N/A indicates that the queue is initialized but is not currently in use by the TCP/IP stack.

In addition, the queue ID and queue type can be used to correlate with VTAM tuning statistics, packet trace, and OSA-Express Network Traffic Analyzer (OSAENTA) trace output for the QDIO interface. The Netstat ALL/-A report includes the interface name for bulk data TCP connections that are using this function. This information can also be returned by the GetConnectionDetail callable NMI. The Netstat STATS/-S report includes the total number of segments received for all connections from the bulk data AIQ of this function. This information can also be returned by the GetGlobalStats callable NMI.

For more information, refer to *z/OS Communications Server: IP System Administrator's Commands*, SC31-8781, and the DISPLAY ID and DISPLAY TRL commands in *z/OS Communications Server: SNA Operation*, SC31-8779.

4.2.3 OSA-Express for zEnterprise (z196)

The IBM zEnterprise 196 (z196) offers communications access to two new internal networks through OSA-Express3 adapters that are configured with an appropriate channel path ID (CHPID) type. The following list describes the two new internal networks:

- ▶ The *intranode management network* provides connectivity between network management applications within the z196 node and it can be accessed through 1000BASE-T Ethernet OSA-Express3 adapters that are configured with a CHPID type of OSM.
- ▶ The *intraensemble data network* provides access to other images that are connected to the intraensemble data network and to applications and appliances that are running in an IBM zEnterprise BladeCenter® Extension (zBX). This internal network can be accessed through 10 GbE OSA-Express3 adapters that are configured with a CHPID type of OSX.

Restrictions:

- ▶ Access to the intranode management network is restricted to authorized management applications, and is only available through Port 0 of any OSA-Express3 CHPID configured with type OSM. Port 1 is not available for these communications.
- ▶ Connectivity to the intranode management network is restricted to stacks that are enabled for IPv6.
- ▶ Connectivity to the intranode management network and to the intraensemble data network is allowed only when the central processor complex (CPC) is a member of an ensemble.

See Figure 4-5 on page 130 for zEnterprise design.

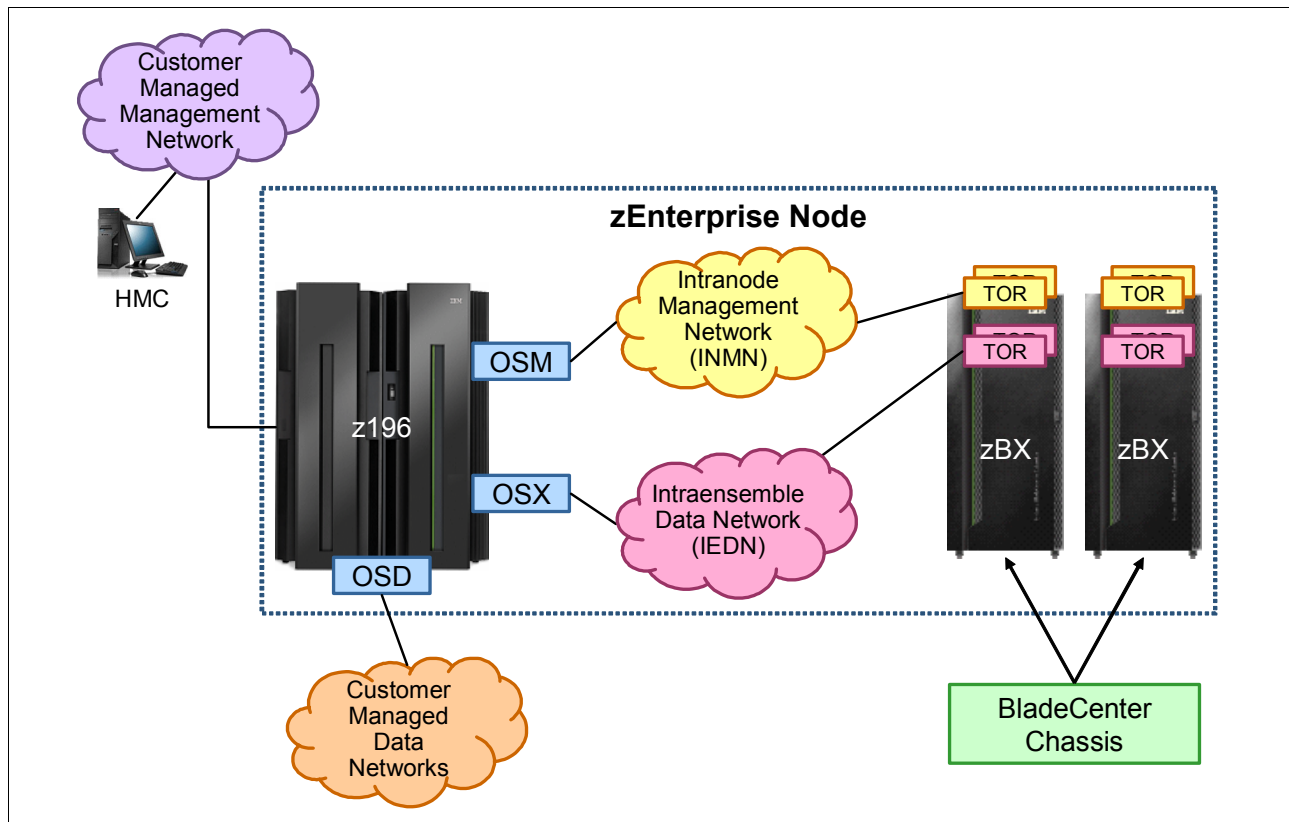


Figure 4-5 zEnterprise Node

Tip: OSA/SF cannot manage these CHPIDs. You can see the related information about these CHPIDs by using the `DISPLAY TCPIP, ,OSAINFO` command.

4.2.4 HiperSockets (MPCIPA)

HiperSockets, also known as *internal Queued Direct I/O (iQDIO)*, is a hardware feature that provides high-speed LPAR-to-LPAR communications within the same server (through memory). It also provides secure data flows between LPARs and high availability, if there is no network attachment dependency or exposure to adapter failures.

HiperSockets can be used to communicate among consolidated servers within a single System z server. All the hardware boxes running these separate servers can be eliminated, along with the cost, complexity, and maintenance of the networking components that interconnect them.

Consolidated servers that have to access corporate data residing on the System z server can do so at memory speeds, bypassing all the network overhead and delays.

HiperSockets can be customized to accommodate varying traffic sizes. With HiperSockets, a maximum frame size can be defined according to the traffic characteristics transported for each HiperSockets.

Because there is no server-to-server traffic outside the System z server, a much higher level of network availability, security, simplicity, performance, and cost effectiveness is achieved as compared with servers communicating across a LAN, such as:

- ▶ HiperSockets has no external components. It provides a very secure connection. For security purposes, servers can be connected to different HiperSockets or VLANs within the same HiperSockets. All security features, like IPSec or IP filtering, are available for HiperSockets interfaces as they are with other TCP/IP network interfaces.
- ▶ HiperSockets looks like any other TCP/IP interface; therefore, it is transparent to applications and supported operating systems.
- ▶ HiperSockets can also improve TCP/IP communications within a sysplex environment when the DYNAMICXCF is used (for example, in cases where Sysplex Distributor uses HiperSockets within the same System z server to transfer IP packets to the target systems).

The HiperSockets device is represented by the IQD channel ID (CHPID) and its associated subchannel devices. All LPARs that are configured in HCD/IOCP to use the same IQD CHPID have internal connectivity and, therefore, have the capability to communicate using HiperSockets.

VTAM will build a single HiperSockets MPC group using the subchannel devices associated with a single IQD CHPID. VTAM will use two subchannel devices for the read and write control devices, and 1 to 8 devices for data devices. Each TCP/IP stack will be assigned a single data device.

Therefore, in order to build the MPC group, there must be a minimum of three subchannel devices defined (within HCD) and associated with the same IQD CHPID. The maximum number of subchannel devices that VTAM will use is 10 (supporting 8 data devices or 8 TCP/IP stacks) per LPAR or MVS image.

When the server that supports HiperSockets and the CHPIDs has been configured in HCD (IOCP), TCP/IP connectivity is provided if:

- ▶ DYNAMICXCF is configured on the IPCONFIG (IPv4) or the IPCONFIG6 (IPv6) statements.
- ▶ A user-defined HiperSockets (MPCIPA) DEVICE and LINK for IPv4 or (IPAQIDIO) INTERFACE for IPv6 is configured and started.

IQD CHPID can be viewed as a *logical* LAN within the server. System z servers allow up to 16 separate IQD CHPIDs, creating the capability of having up to 16 separate logical LANs within the same server.

Each IQD CHPID can be assigned to a set of LPARs (configured in HCD), making it possible to isolate these LPARs in separate logical LANs, as shown in of Figure 4-6.

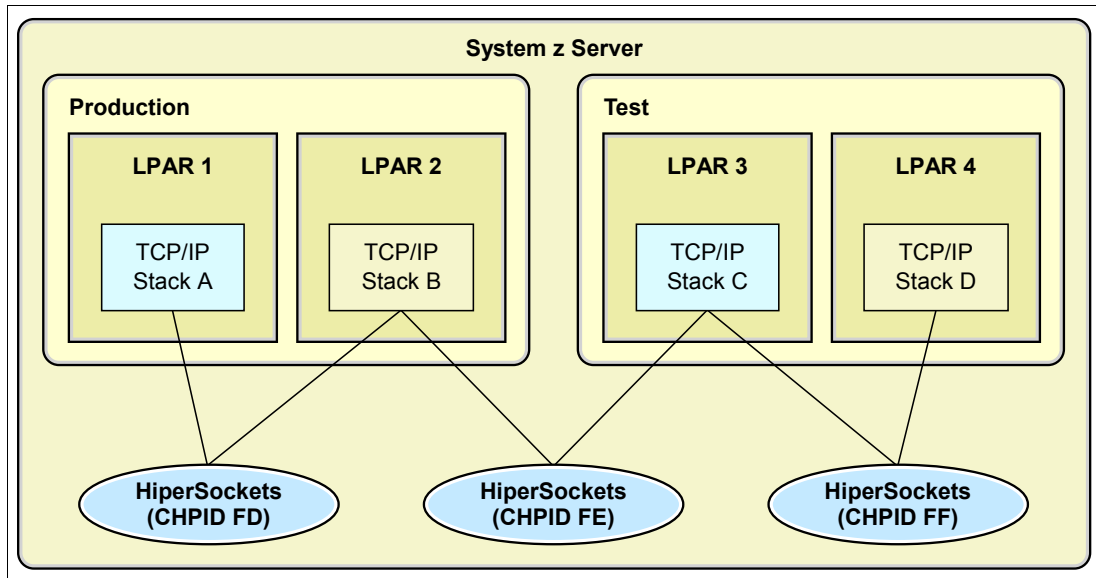


Figure 4-6 HiperSockets: multiple logical LANs

HiperSockets multiple write

The HiperSockets multiple write facility moves multiple buffers of data with a single write operation. This facility was added to reduce CPU utilization and to improve performance for large outbound messages over HiperSockets.

Restriction: HiperSockets multiple write is effective only on an IBM System z10 and when z/OS is not running as a guest in a z/VM® environment.

To enable the HiperSockets multiple write facility on all HiperSockets interfaces, including interfaces created for dynamic XCF, add the IQDMULTIWRITE parameter to the GLOBALCONFIG statement.

For more information, see Appendix B, “Additional parameters and functions” on page 411.

For a review of the scenarios we used to test HiperSockets multiple write, see Appendix A, “HiperSockets Multiple Write”, in *IBM z/OS V1R11 Communications Server TCP/IP Implementation Volume 3: High Availability, Scalability, and Performance*, SG24-7800.

HiperSockets multiple write assist with IBM zIIP

On an IBM System z10, an additional assist for HiperSockets data that is using the multiple write facility is available through the IBM System z10 Integrated Information Processor (zIIP).

To enable HiperSockets traffic that is using the multiple write facility to be processed on available zIIPs, specify the ZIIP IQDIOMULTIWRITE parameter on the GLOBALCONFIG statement.

HiperSockets VLAN support

HiperSockets connections defined through DYNAMICXCF coding or through individual DEVICE and LINK statement coding also support VLAN tagging. This allows you to split the internal LAN represented by a single HiperSockets CHPID into multiple virtual LANs, providing isolation for security or administrative purposes. Only stacks attached to the same HiperSockets VLAN can communicate with each other. Stacks attached to a different

HiperSockets VLAN on the same CHPID cannot use the HiperSockets path to communicate with the stacks on a different VLAN.

Note: The VLAN ID assigned to a HiperSockets device applies to both IPv4 and IPv6 connections over that CHPID.

HiperSockets accelerator

The Communications Server takes advantage of the technological advances and high-performing nature of the I/O processing offered by HiperSockets with the IBM System z servers and OSA-Express, using the QDIO architecture. This is achieved by optimizing IP packet forwarding processing that occurs across these two types of technologies. This function is referred to as HiperSockets Accelerator. It is a configurable option, and is activated by defining the IQDIORouting option on the IPCONFIG statement.

When the TCP/IP stack is configured with HiperSockets Accelerator, it allows IP packets received from HiperSockets to be forwarded to an OSA-Express port (or vice versa) without the need for those IP packets to be processed by the TCP/IP stack.

When using this function, one or more LPARs contain the *routing* stack, which manages connectivity through OSA-Express ports to the LAN, while the other LPARs connect to the routing stack using the HiperSockets, as shown in Figure 4-7.

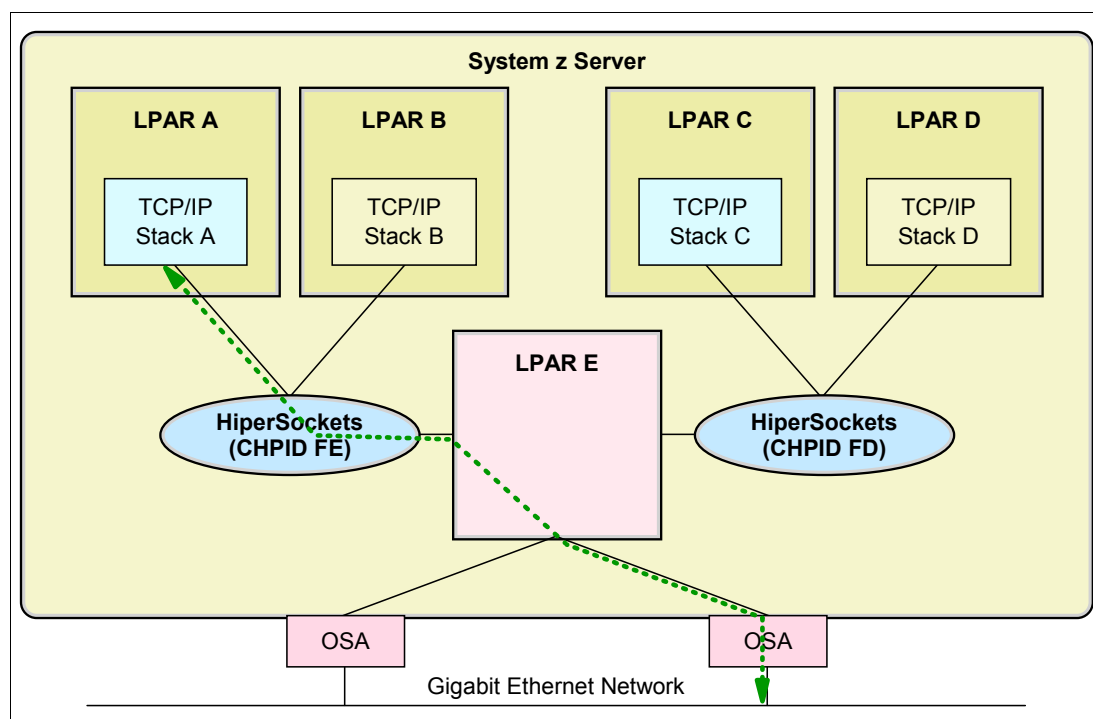


Figure 4-7 HiperSockets Accelerator

Note: This example is intended purely to demonstrate IP traffic flow. We do not recommend implementing HiperSockets Accelerator using a single LPAR.

Detailed information about the subject of HiperSockets is available in *HiperSockets Implementation Guide*, SG24-6816.

4.2.5 Dynamic XCF

You have a choice of defining the XCF connectivity to other TCP/IP stacks individually or using the dynamic XCF definition facility. Dynamic XCF significantly reduces the number of definitions that you need to create whenever a new system joins the sysplex or when you need to start up a new TCP/IP stack. These changes become more numerous as the number of stacks and systems in the sysplex grows. This could lead to configuration errors. With dynamic XCF, you do not need to change the definitions of the existing stacks in order to accommodate the new stack.

From an IP topology perspective, DYNAMICXCF establishes fully meshed IP connectivity to all other z/OS TCP/IP stacks in the sysplex. You only need one end-point specification in each stack for fully meshed connectivity to all other stacks in the sysplex. When a new stack gets started, Dynamic XCF connectivity is automatically established.

Note: Only one dynamic XCF network is supported per sysplex.

Dynamic XCF is required to support Sysplex Distributor and nondisruptive dynamic VIPA movement (discussed in detail in *IBM z/OS V1R11 Communications Server TCP/IP Implementation Volume 3: High Availability, Scalability, and Performance*, SG24-7800).

Dynamic XCF uses Sysplex Sockets support, allowing the stacks to communicate with each other and exchange information such as VTAM CPNAMES, MVS SYSCONE value, and IP addresses. The dynamic XCF definition is activated by coding the IPCONFIG DYNAMICXCF parameter in the TCP/IP profile.

Dynamic XCF creates definitions for DEVICE, LINK, HOME, and BSDROUTINGPARMS statements and the START statement dynamically. When activated, the dynamic XCF devices and links appear to the stack as though they had been defined in the TCP/IP profile. They can be displayed using standard commands, and they can be stopped and started.

During TCP/IP initialization the stack joins the XCF group, ISTXCF, through VTAM. When other stacks in the group discover the new stack, the definitions are created automatically, the links are activated, and the remote IP address for each link is added to the routing table. After the remote IP address has been added, IP traffic can flow across one of the following interfaces:

- ▶ IUTSAMEH (within the same LPAR)
- ▶ HiperSockets (within the same server)
- ▶ XCF signaling (different server, either using the Coupling Facility link or a CTC connection)

Dynamic XCF support is illustrated in Figure 4-8.

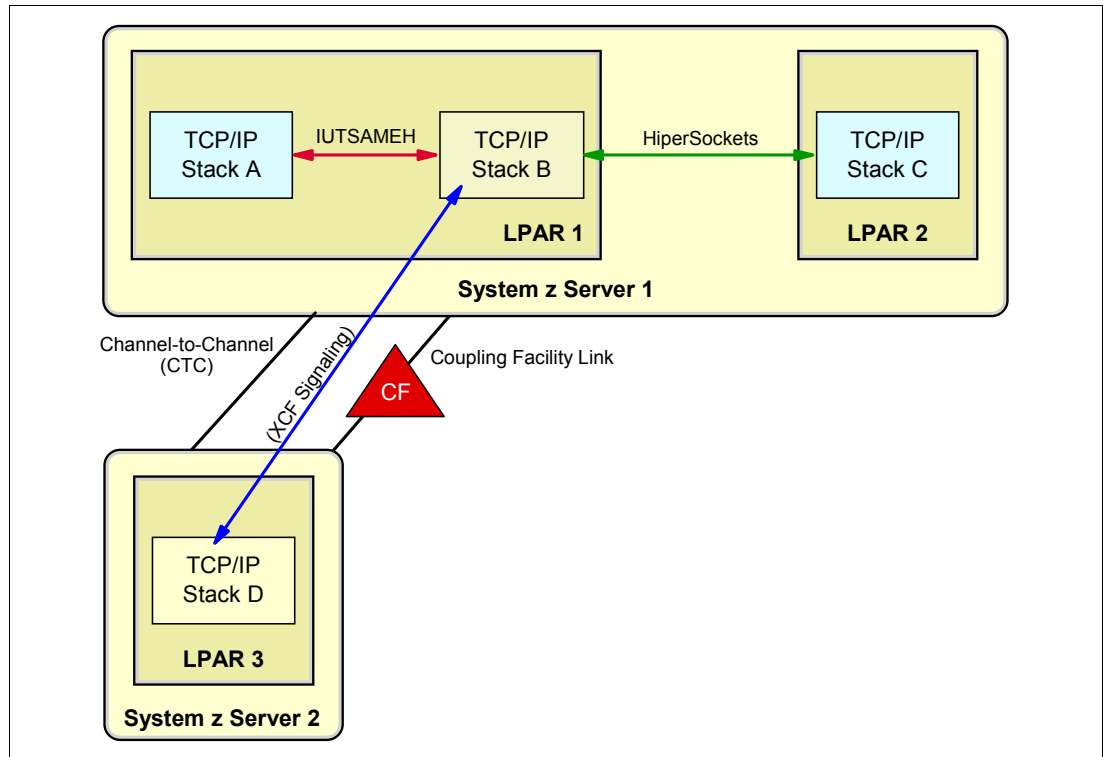


Figure 4-8 Dynamic XCF support

HiperSockets DYNAMICXCF connectivity

z/OS images within the same server with DYNAMICXCF coded will use HiperSockets DYNAMICXCF connectivity instead of standard XCF connectivity, under these conditions:

- ▶ The TCP/IP stacks must be on the same server.
- ▶ For the DYNAMICXCF HiperSockets device (IUTIQDIO), the stacks must be using the same IQD CHPID, even with different channel subsystems (spanning).
- ▶ The stacks must be configured through HCD or the IOCDS to use HiperSockets.
- ▶ For IPv6 HiperSockets connectivity, both stacks must be at z/OS V1R7 or higher.
- ▶ The initial HiperSockets activation must complete successfully.

When an IPv4 DYNAMICXCF HiperSockets device and link are created and successfully activated, a subnetwork route is created across the HiperSockets link. The subnetwork is created by using the DYNAMICXCF IP address and mask. This allows any LPAR within the same server to be reached, even ones that are not within the sysplex. To do that, the LPAR that is outside of the sysplex environment must define at least one IP address for the HiperSockets endpoint that is within the subnetwork defined by the DYNAMICXCF IP address and mask.

When multiple stacks reside within the same LPAR that supports HiperSockets, both IUTSAMEH and HiperSockets links or interfaces will coexist. In this case, it is possible to transfer data across either link. Because IUTSAMEH links have better performance, it is always better to use them for intra-stack communication. A host route will be created by DYNAMICXCF processing across the IUTSAMEH link, but not across the HiperSockets link.

For additional information about dynamic XCF, Sysplex Distributor, and nondisruptive dynamic VIPA movement refer to *IBM z/OS V1R12 Communications Server TCP/IP Implementation Volume 3: High Availability*, SG24-7898.

4.3 Connectivity for the z/OS environment

The subsequent sections focus on the interface implementation only, which means establishing Layer 2 and a subset of Layer 3 (IP addressing) connectivity. For details beyond the basic implementation of the immediate LAN environment, also refer to:

- ▶ Chapter 5, “Routing” on page 205 for IP routing details
- ▶ Chapter 6, “VLAN and Virtual MAC support” on page 265 for use of virtual MAC addresses
- ▶ Chapter 7, “Sysplex subplexing” on page 283 for isolating TCP/IP stack in a sysplex

To design connectivity in a z/OS environment, you must take the following considerations into account:

- ▶ As a server environment, network connectivity to the external corporate network should be carefully designed to provide a high-availability environment, avoiding single points of failures.
- ▶ If a z/OS LPAR is seen as a stand-alone server environment on the corporate network, it should be designed as an endpoint.
- ▶ If a z/OS LPAR will be used as a front-end concentrator (for example, making use of HiperSockets Accelerator), it should be designed as an intermediate network or node.

Recommendation: Although there are specialized cases where multiple stacks per LPAR can provide value, in general we recommend implementing only one TCP/IP stack per LPAR. The reasons for this recommendation are as follows:

- ▶ A TCP/IP stack is capable of exploiting all available resources defined to the LPAR in which it is running. Therefore, starting multiple stacks will not yield any increase in throughput.
- ▶ When running multiple TCP/IP stacks, additional system resources, such as memory, CPU cycles, and storage, are required.
- ▶ Multiple TCP/IP stacks add a significant level of complexity to TCP/IP system administration tasks.
- ▶ It is not necessary to start multiple stacks to support multiple instances of an application on a given port number, such as a test HTTP server on port 80 and a production HTTP server also on port 80. This type of support can instead be implemented using BIND-specific support where the two HTTP server instances are each associated to port 80 with their own IP address, using the BIND option on the PORT reservation statement.

One example where multiple stacks can have value is when an LPAR needs to be connected to multiple isolated security zones in such a way that there is no network level connectivity between the security zones. In this case, a TCP/IP stack per security zone can be used to provide that level of isolation, without any network connectivity between the stacks.

Based on these considerations, in the following sections we present best practice scenarios for building a z/OS Communications Server TCP/IP configuration, using OSA-Express (QDIO), HiperSockets (iQDIO), and dynamic XCF.

We built our connectivity scenarios with two OSA-Express3 1000BASE-T features (four ports each) that are connected to the LAN environment (one layer3 switch). We also implemented a

HiperSockets internal LAN to interconnect all LPARs within the same System z10. Finally, we used dynamic XCF connectivity for the Sysplex environment.

Note: In our environment we connected all the OSA ports to one switch, but in a production implementation it is best to connect your OSAs to at least two switches

The scenarios we discuss are as follows:

- ▶ 4.4.3, “Configuring OSA-Express with VLAN ID” on page 148
- ▶ 4.6.3, “Configuring HiperSockets” on page 182
- ▶ 4.7.3, “Configuring DYNAMICXCF” on page 189

Note: In this chapter, we define only our LPARs as end points.

For a complete picture of our implementation environment, refer to Appendix D, “Our implementation environment” on page 455.

4.3.1 IOCP definitions

Example 4-1 on page 137 is an excerpt of the IOCP statements we used in our System z environment (only showing OSA-Express CHPID 02 and HiperSockets CHPID F4). These statements are required by the input/output subsystem and the operating system. Because all of our OSA-Express and HiperSockets connectivity will be used across all four LPARs, we defined the CHPIDs as shared.

Example 4-1 IOCP statements

```
ID      MSG2='SYS6.IODF64 - 2010-09-23 11:18',SYSTEM=(2817,1),  *
        LSYSTEM=SCZP301,                                         *
        TOK=('SCZP301',00800006991E2094111808480110266F00000000,*
        00000000,'10-09-23','11:18:08','SYS6','IODF64')
RESOURCE PARTITION=((CSS(0),(A0A,A),(A0B,B),(A0C,C),(A0D,D),(A*
        0E,E),(A0F,F),(A01,1),(A02,2),(A03,3),(A04,4),(A05,5),(A*
        06,6),(A07,7),(A08,8),(A09,9)),(CSS(1),(A1B,B),(A1E,E),(*
        A1F,F),(A11,1),(A12,2),(A13,3),(A14,4),(A15,5),(A16,6),(*
        A17,7),(A18,8),(*,9),(*,A),(*,C),(*,D)),(CSS(2),(A2F,F),*
        (A21,1),(A22,2),(*,3),(*,4),(*,5),(*,6),(*,7),(*,8),(*,9*
        ),(*,A),(*,B),(*,C),(*,D),(*,E)),(CSS(3),(A31,1),(*,2),(*
        *,3),(*,4),(*,5),(*,6),(*,7),(*,8),(*,9),(*,A),(*,B),(*,*
        C),(*,D),(*,E),(*,F)))

CHPID  PATH=(CSS(1),0A),SHARED,
        PARTITION=((A11,A13,A16,A18),('=')),CHPARM=02,PCHID=531,
        TYPE=OSM
CHPID  PATH=(CSS(1),0B),SHARED,
        PARTITION=((A11,A13,A16,A18),('=')),CHPARM=02,PCHID=101,
        TYPE=OSM

CNTLUNIT CUNUMBR=2340,PATH=((CSS(1),0A)),UNIT=OSM
IODEVICE ADDRESS=(2340,015),MODEL=M,UNITADD=00,CUNUMBR=(2340),
        UNIT=OSA,MODEL=M,DYNAMIC=YES,LOCANY=YES
CNTLUNIT CUNUMBR=2360,PATH=((CSS(1),0B)),UNIT=OSM
IODEVICE ADDRESS=(2360,015),MODEL=M,UNITADD=00,CUNUMBR=(2360),
        UNIT=OSA,MODEL=M,DYNAMIC=YES,LOCANY=YES
```

```

CHPID PATH=(CSS(1),18),SHARED,
      PARTITION=((A11,A13,A16,A18),(-)),PCHID=590,TYPE=OSX 1
CHPID PATH=(CSS(1),19),SHARED,
      PARTITION=((A11,A13,A16,A18),(-)),CHPARM=02,PCHID=510,
      TYPE=OSX

CNTLUNIT CUNUMBR=2300,PATH=((CSS(1),18)),UNIT=OSX
IODEVICE ADDRESS=(2300,015),MODEL=X,CUNUMBR=(2300),UNIT=OSA,
      MODEL=X,DYNAMIC=YES,LOCANY=YES
CNTLUNIT CUNUMBR=2320,PATH=((CSS(1),19)),UNIT=OSX
IODEVICE ADDRESS=(2320,015),MODEL=X,UNITADD=00,CUNUMBR=(2320),
      UNIT=OSA,MODEL=X,DYNAMIC=YES,LOCANY=YES


CHPID PATH=(CSS(1),02),SHARED,
      PARTITION=((A11,A13,A16,A18),(-)),PCHID=530,TYPE=OSD
CHPID PATH=(CSS(1),03),SHARED,
      PARTITION=((A11,A13,A16,A18),(-)),PCHID=100,TYPE=OSD
CHPID PATH=(CSS(1),04),SHARED,
      PARTITION=((A11,A13,A16,A18),(-)),PCHID=181,TYPE=OSD
CHPID PATH=(CSS(1),05),SHARED,
      PARTITION=((A11,A13,A16,A18),(-)),PCHID=291,TYPE=OSD

CNTLUNIT CUNUMBR=2080,PATH=((CSS(1),02)),UNIT=OSA
IODEVICE ADDRESS=(2080,015),UNITADD=00,CUNUMBR=(2080),UNIT=OSA
IODEVICE ADDRESS=208F,UNITADD=FE,CUNUMBR=(2080),UNIT=OSAD
CNTLUNIT CUNUMBR=20A0,PATH=((CSS(1),03)),UNIT=OSA
IODEVICE ADDRESS=(20A0,015),UNITADD=00,CUNUMBR=(20A0),UNIT=OSA
IODEVICE ADDRESS=20AF,UNITADD=FE,CUNUMBR=(20A0),UNIT=OSAD
CNTLUNIT CUNUMBR=20C0,PATH=((CSS(1),04)),UNIT=OSA
IODEVICE ADDRESS=(20C0,015),UNITADD=00,CUNUMBR=(20C0),UNIT=OSA
IODEVICE ADDRESS=20CF,UNITADD=FE,CUNUMBR=(20C0),UNIT=OSAD
CNTLUNIT CUNUMBR=20E0,PATH=((CSS(1),05)),UNIT=OSA
IODEVICE ADDRESS=(20E0,015),UNITADD=00,CUNUMBR=(20E0),UNIT=OSA
IODEVICE ADDRESS=20EF,UNITADD=FE,CUNUMBR=(20E0),UNIT=OSAD

CHPID PATH=(CSS(1),F4),SHARED,
      PARTITION=((A11,A13,A16,A18),(-)),TYPE=IQD
CHPID PATH=(CSS(1),F5),SHARED,
      PARTITION=((A11,A13,A16,A18),(-)),TYPE=IQD
CHPID PATH=(CSS(1),F6),SHARED,
      PARTITION=((A11,A13,A16,A18),(-)),TYPE=IQD
CHPID PATH=(CSS(1),F7),SHARED,
      PARTITION=((A11,A13,A16,A18),(-)),TYPE=IQD

CNTLUNIT CUNUMBR=E800,PATH=((CSS(1),F4)),UNIT=IQD
IODEVICE ADDRESS=(E800,032),CUNUMBR=(E800),UNIT=IQD
CNTLUNIT CUNUMBR=E900,PATH=((CSS(1),F5)),UNIT=IQD
IODEVICE ADDRESS=(E900,032),CUNUMBR=(E900),UNIT=IQD
CNTLUNIT CUNUMBR=EA00,PATH=((CSS(1),F6)),UNIT=IQD
IODEVICE ADDRESS=(EA00,032),CUNUMBR=(EA00),UNIT=IQD
CNTLUNIT CUNUMBR=EB00,PATH=((CSS(1),F7)),UNIT=IQD
IODEVICE ADDRESS=(EB00,032),CUNUMBR=(EB00),UNIT=IQD

```

Attention:  The CHPIDs type OSM and OSX are only used if you are connected to a zBX (zEnterprise Blade Center).

There are other ways to build the IOCDS for an OSA-Express adapter than the one depicted in Example 4-1 on page 137. This applies particularly to an OSA-Express3, which can contain more than a single port on the same CHPID. However, in our labs, we used the method shown in Example 4-1 on page 137. Consult 4.4.1, “Dependencies: CHPID, IOCDS, port numbers, portnames, and port sharing” on page 140 to see other alternatives to define the IOCDS and to review our recommendations.

4.3.2 VTAM definitions

Before getting started with configuring the scenarios in the following sections, it is important to understand the role of VTAM in the TCP/IP configuration.

z/OS Communications Server provides a set of High Performance Data Transfer (HPDT) services that includes MultiPath Channel (MPC), a high-speed channel interface designed for network protocol use (for example, APPN or TCP/IP).

Multiple protocols can either share or have exclusive use of a set of channel paths to an attached platform. MPC provides the ability to have multiple device paths, defined as a single logical connection.

The term MPC group is used to define a single MPC connection that can contain multiple read and write paths. The number of read and write paths does not have to be equal, but there must be at least one read and write path defined within each MPC group.

MPC groups are defined using the Transport Resource List (TRL), where each defined MPC group becomes an entry (that is, a TRLE) in the TRL table. The configuration and control of the MultiPath Channel (MPC) interfaces are provided by VTAM. They are enabled in VTAM as TRLE minor nodes.

You must define the channel paths that are a part of the group in the TRLE. Each TRLE is identified by a resource_name. For OSA-Express, the TRLE also has a port_name to identify the association between VTAM and TCP/IP, allowing connectivity to the OSA-Express port. OSA-Express3 Gigabit Ethernet and 1000Base-T also defines port_num to identify which port the TRLE definition applies to.

For HiperSockets, the TRLE is generated dynamically by VTAM.

For details about defining a TRLE, refer to *z/OS Communications Server: SNA Resource Definition*, SC31-8778.

4.4 OSA-Express QDIO connectivity

Configuring an OSA-Express (QDIO mode) in a single stack scenario is the simplest way to integrate your z/OS TCP/IP stack into a LAN environment. This scenario, however, still needs to be planned to avoid any single points of failure. Therefore, we must have at least two OSA-Express features connecting to two different switches in the network.

Because we are dealing with multiple LPARs in our server, for redundancy purposes we have shared the OSA-Express ports (CHPID type OSD) across all LPARs.

In this scenario, we have two OSA-Express3 1000BASE-T features, each with four ports, two ports per channel. One port of each channel was used unless the second port was needed for testing of new functions. This allowed us to have four CHPIDs (02, 03, 04, and 05), shared by our four LPARs (SC30, SC31, SC32 and SC33), as shown in Figure 4-9 on page 140.

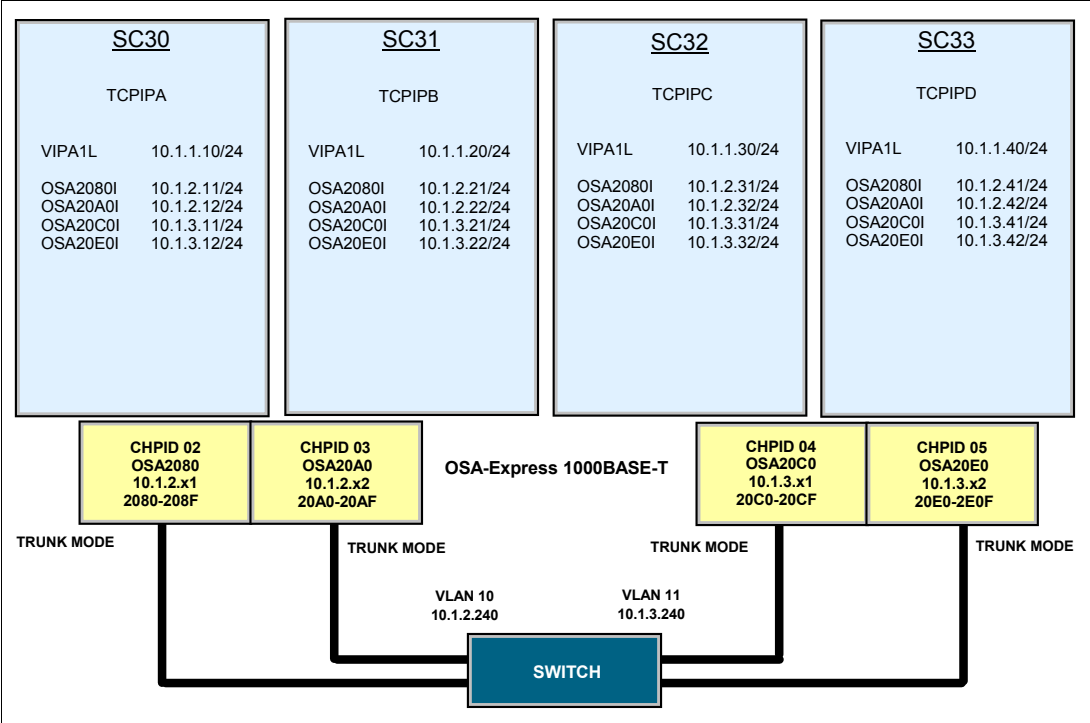


Figure 4-9 OSA-Express (QDIO) implementation

To make better use of our OSA-Express ports and to control data traffic patterns, we defined one port on each OSA-Express feature with a separate VLAN ID, creating two subnetworks to be used by all LPARs. In a high availability configuration, these OSA-Express ports will be the path to all of our IP addresses for the LAN environment.

4.4.1 Dependencies: CHPID, IOCDS, port numbers, portnames, and port sharing

To implement this scenario, we have the following dependencies:

- ▶ The OSA-Express port must be defined as CHPID type OSD to the server using HCD or IOCP to enable QDIO. This CHPID must be defined as shared to all LPARs that will use the OSA-Express port (see Example 4-1 on page 137).
- ▶ To define an OSA-Express port in QDIO mode, we use the MPCIPA DEVICE statement, specifying the PORTNAME value from the TRLE definition as the device_name. The TRLE must be defined as MPCLEVEL=QDIO.
- ▶ The Virtual LAN identifiers (VLAN IDs) defined to each OSA-Express port must be recognized by the switch.
- ▶ The switch ports where the OSA-Express ports are connected must be configured in trunk mode.

OSA-Express2 and OSA-Express3 Adapter and port layouts

While an OSA-Express2 adapter (with the exception of the 10 Gbe adapter) contains two ports, the OSA-Express3 (with the exception of the 10 Gbe adapter) houses four ports. Compare and contrast the layout of an OSA-Express2 and OSA-Express3 in Figure 4-10.

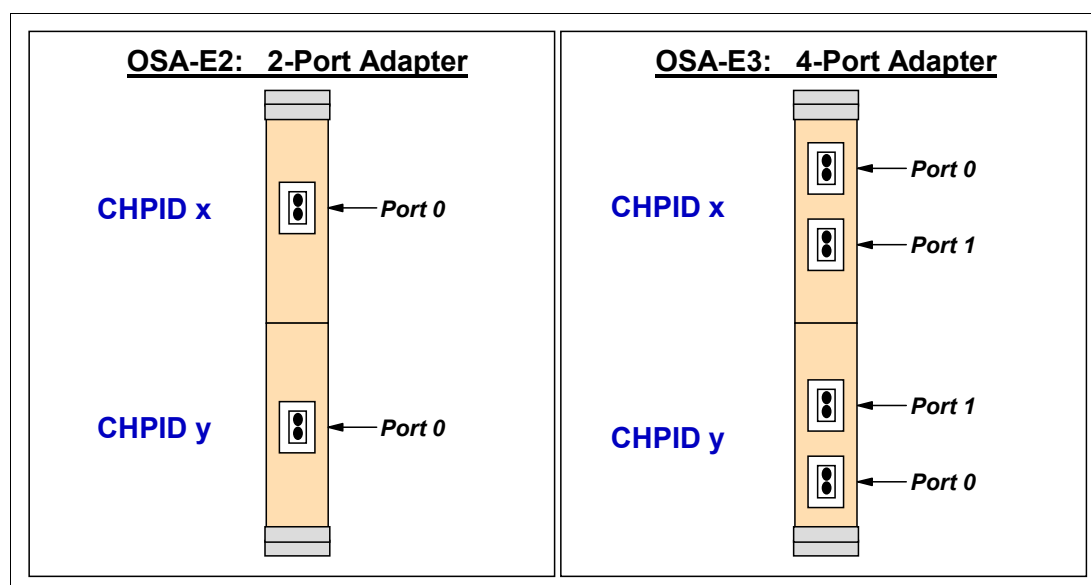


Figure 4-10 Comparison: OSA-E2 2-port adapter and OSA-E3 4-port adapter

Each port of the OSA-Express2 adapter depicted in Figure 4-10 is on a separate CHPID: CHPID x and CHPID y. Each port on each CHPID is defined with a separate port name and resides at port number 0.

The OSA-Express3 is engineered with two ports on each CHPID: CHPID x and CHPID y. The two ports on each CHPID are numbered port 0 and port 1. Note how the top half of the OSA-E3 is the mirror image of the bottom half with regard to the port number assignments; reading from top to bottom, you see Port 0, Port 1, Port 1, Port 0. As with any OSA port, the portnames on the multi-port OSA-E3 must be unique to a CHPID. An explanation of this portname assignment is provided in “Considerations for assigning the OSA portname” on page 147.

Considerations for the IOCP or IOCDS Definitions for an OSA-Express3

We have shown you in Example 4-1 on page 137 an IOCDS that was originally built for an OSA-Express2 configuration. When migrating to an OSA-Express3, you can choose to use a similar IOCDS and spread the assigned addresses from a single address range across two ports of the same CHPID that originally connected to only one port. You can also choose to change your IOCDS to reflect separate address ranges or even separate *logical* control units, despite the presence of only a single *physical* control unit on the CHPID. We now want to show you a couple of different ways to implement an IOCDS for an OSA-Express3 implementation.

Alternative 1: Single IODEVICE range for two E3 ports on single CHPID

In our scenarios where we used an OSA-Express3, we used exactly the same IOCDS definitions as those we deployed for an OSA-Express2. You have seen this IOCDS in Example 4-1 on page 137. We use as an example the IOCDS definitions for the devices on OSA port OSA2080. In Example 4-2, you see at **1** that we have allocated fifteen addresses (2080-208E) to QDIO connections starting with device address 2080.

Example 4-2 Sample CNTLUNIT and IODEVICE for an OSA on CHPID Type OSD (QDIO)

```
CNTLUNIT CUNUMBR=2080,PATH=((CSS(2),02)),UNIT=OSA
IODEVICE ADDRESS=(2080,015),CUNUMBR=(2080),UNIT=OSA 1
```

Example 4-2 corresponds to what you must code in a VTAM TRLE definition in order to support a QDIO connection of a TCP/IP stack. Look at Example 4-3, where you see that the VTAM TRLE that defines port number 0 (**A**) (the only port number on an OSA-Express2) utilizes only the first nine addresses (2080-2088) of the allocated fifteen addresses (2080-208E) on this CNTLUNIT.

Example 4-3 TRLE definition for PORTNUM=0 (Portname of OSA2080)

```
OSA2080 VBUILD TYPE=TRL
OSA2080P TRLE LNCTL=MPC, *
                READ=2080, *
                WRITE=2081, *
                DATAPATH=(2082-2088), *
                PORTNAME=OSA2080, *
                PORTNUM=0, A *
                MPCLEVEL=QDIO
```

To add the OSA-Express3 port that resides at port number 1 of the same CHPID, we use the same IOCDS as before, but now we add a new TRLE definition for PORTNUM=1 (**B**). See the TRLE example in Example 4-4.

Example 4-4 TRLE definition for PORTNUM=1 (Portname of OSA2081)

```
OSA2081 VBUILD TYPE=TRL
OSA2081P TRLE LNCTL=MPC,
                READ=2089, C
                WRITE=208A,
                DATAPATH=(208B-208D),
                PORTNAME=OSA2081,
                PORTNUM=1, B
                MPCLEVEL=QDIO
```

In Example 4-4, we have simply started the addresses for PORTNUM=1 at 2089 of the IOCDS **C**.

Figure 4-11 shows the allocation of all the device addresses across the two ports of an OSA-Express 3 card.

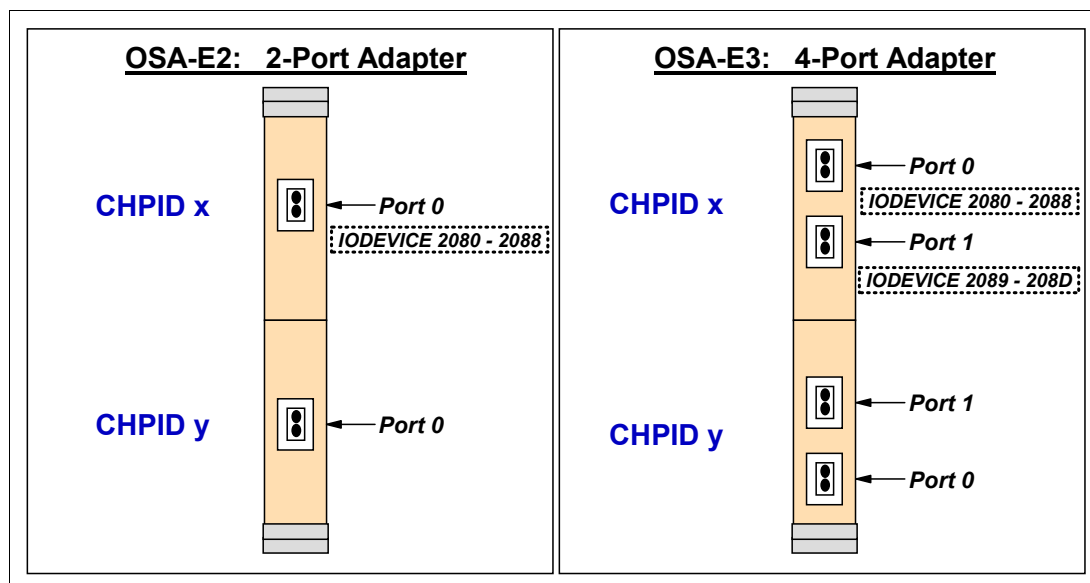


Figure 4-11 Allocation of device addresses across two ports of an OSA-Express3

As you saw in Example 4-2 on page 142, the IOCP definitions have no awareness of the OSA adapter's two ports and simply assign device addresses; the VTAM definition for z/OS does care about the port numbers and maps the number to the addresses (Example 4-3 on page 142 and Example 4-4 on page 142). This address allocation scheme worked well for us because we did not have to reconfigure the IOCP for our test. Other schemes may work better for you, particularly if you are consolidating OSA ports from separate CHPIDs onto the same CHPID of a new OSA-Express3.

Note: Our examples show you how to point to the two different ports with the PORTNUM parameter in a z/OS example. Other System z operating systems, such as z/VM, Linux® on z, z/VSE™, or TPF, have similar coding parameters to allocate addresses to port number 0 versus port number 1. See the appropriate operating system documentation for those definitions.

Bear in mind that a migration to OSA-Express 3 can affect more than just the IOCDs. You also have other types of definitions in the operating system and potentially in access methods (like VTAM) to migrate. The more you can keep the definitions the same across migrations, the easier and more efficient the migration to a new platform or release becomes. This is where the next two alternatives can make a difference for you.

Alternative 2: Two IODEVICE ranges for two E3 ports on a single CHPID

An alternative to the coding scheme we just showed you in Example 4-2 on page 142, Example 4-3 on page 142, and Example 4-4 on page 142 is to use a different address for each of the two ports. Such a scheme might make problem determination and operator procedures easier for you, as message displays very clearly show the distinction between the two ports, although they reside on the same CHPID.

In Example 4-5, you see a range of addresses starting with 1000 (A) and another range starting with 2000 (B), and the VTAM definitions in Example 4-6 show that these addresses are used for OSA-E3 port numbers 0 and 1. (Compare with 1 and 2 in Example 4-6.)

Example 4-5 Separate device ranges for separate OSA-Express3 ports

```
CNTLUNIT CUNUMBR=1000,PATH=((CSS(0),10)),UNIT=OSA
IODEVICE ADDRESS=(1000,032),CUNUMBR=(1000),UNIT=OSA (A)
IODEVICE ADDRESS=(10FE,001),CUNUMBR=(1000),UNIT=OSAD
IODEVICE ADDRESS=(2000,032),UNITADD=20,CUNUMBR=(1000),UNIT=OSA (B)
```

Example 4-6 VTAM definitions for OSA-E3 port numbers 0 and 1 (two device ranges)

```
OSA1000 VBUILD TYPE=TRL
OSA1000P TRLE LNCTL=MPC, *
READ=1000, *
WRITE=1001, *
DATAPATH=(1002), *
PORTNAME=OSA1000, *
PORTNUM=0, 1 *
MPCLEVEL=QDIO

OSA2000 VBUILD TYPE=TRL
OSA2000P TRLE LNCTL=MPC, *
READ=2000, *
WRITE=2001, *
DATAPATH=(2002), *
PORTNAME=OSA2000, *
PORTNUM=1, 2 *
MPCLEVEL=QDIO
```

The diagram in Figure 4-12 shows you how the device addresses are allocated for this example.

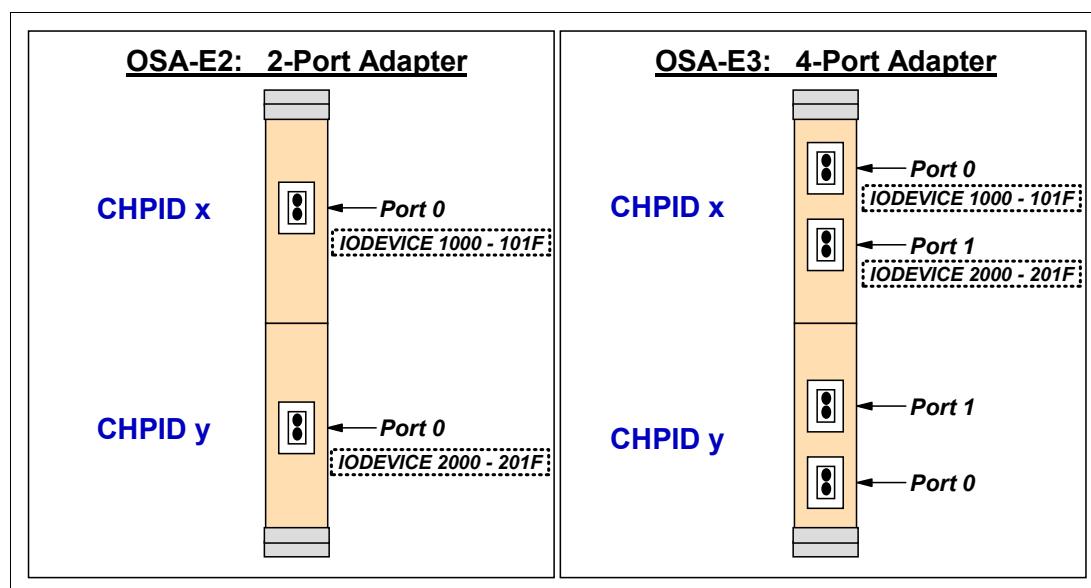


Figure 4-12 Consolidating two OSA ports from OSA-E2 onto a single CHPID of OSA-E3

With this alternative, you can preserve the device addresses in your VTAM definitions and simply deal with a few changes in the IOCDs. This might represent a simple migration scenario for you if you have many VTAM definitions to change.

Alternative 3: Two Logical Control Units on Physical CU for two E3 ports

The following examples show you how to make a logical distinction within the IOCP between ports 0 and 1 of an OSA-E3 CHPID. You can specify separate logical control units with the CUADD parameter. While in the past defining multiple Control Units had value only if you were defining many devices, it appears customers migrating from OSA-Express2 channels to multi-port OSA-Express3s are finding it easier to combine two OSA-Express CHPIDs into the two ports of an OSA-Express3 CHPID.

Refer to Example 4-7, where you find the device range for port number 0 under CUADD=0 (A) and the device range for port number 1 under CUADD=1 (B).

Example 4-7 Separate logical control unit for each OSA-E3 port

```

CNTLUNIT CUNUMBR=3000,CUADD=0 A,PATH=((CSS(0),02),(CSS(1),02)),UNIT=OSA
IODEVICE ADDRESS=(3000,032),UNITADD=00,CUNUMBR=(3000),UNIT=OSA
IODEVICE ADDRESS=3020,UNITADD=FE,CUNUMBR=(3000),UNIT=OSAD
CNTLUNIT CUNUMBR=3500,CUADD=1 B,PATH=((CSS(0),02),(CSS(1),02)),UNIT=OSA
IODEVICE ADDRESS=(3500,032),UNITADD=00,CUNUMBR=(3500),UNIT=OSA

```

The VTAM definitions look similar to what you have seen before. Examine the coding in Example 4-8.

Example 4-8 VTAM TRLEs for two logical control units and port numbers of an OSA-E3

```

OSA3000 VBUILD TYPE=TRL
OSA3000P TRLE  LNCTL=MPC,                                *
                READ=3000,                                *
                WRITE=3001,                                *
                DATAPATH=(3002),                            *
                PORTNAME=OSA3000,                            *
                PORTNUM=0, 1                                *
                MPCLEVEL=QDIO

OSA3500 VBUILD TYPE=TRL
OSA3500P TRLE  LNCTL=MPC,                                *
                READ=3500,                                *
                WRITE=3501,                                *
                DATAPATH=(3502),                            *
                PORTNAME=OSA3500,                            *
                PORTNUM=1, 2                                *
                MPCLEVEL=QDIO

```

The device range beginning with 3000 has been assigned to port number 0 (1); the device range starting with 3500 has been assigned to port number 1 (2). Refer to Figure 4-13.

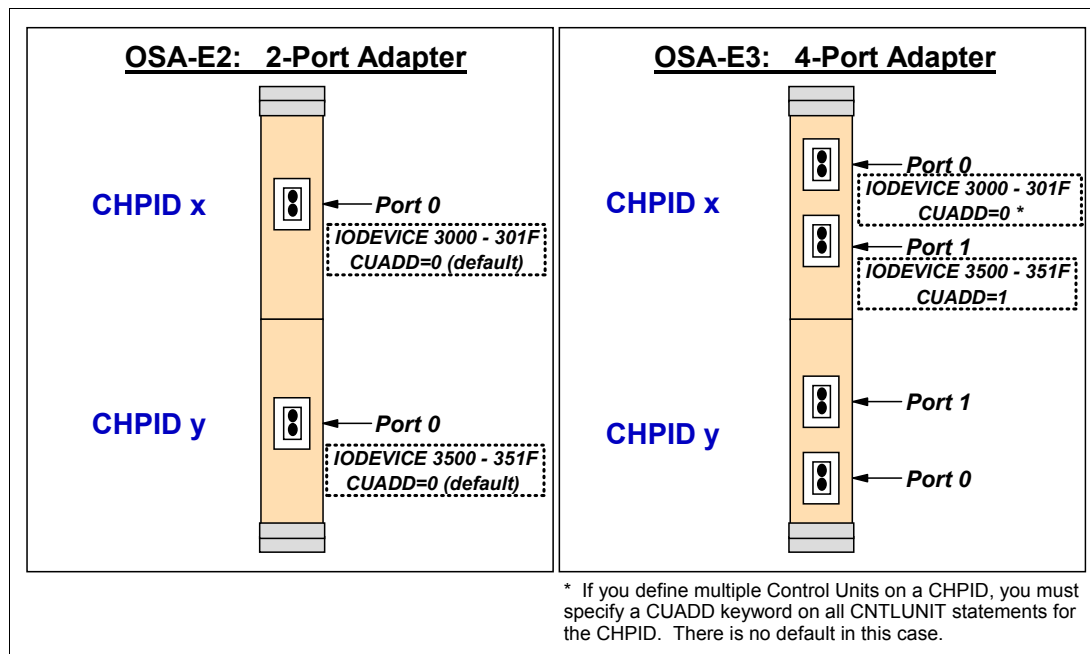


Figure 4-13 Distinguishing OSA-E3 port numbers in the IOCDS with a CUADD parameter

Just as with the second alternative, you might find it easier to merge what were OSA connections on two separate CHPIDs into a single CHPID and distinguish them not only with separate address ranges, but also with separate logical control unit numbers.

Notes:

1. In all the IOCDS definitions we have shown you, we have coded the OSA/SF device on CUADD=0, either by default or through explicit coding. The OSA/SF device must reside on CUADD=0.
2. OSA supports Outbound Priority Queueing (multiple Outbound Queues) as long as no more than 480 valid subchannels are defined for all LPARs sharing a CHPID. Each logical partition sharing a CHPID gets a subchannel for every device defined on that CHPID. Therefore, if you define a CHPID shared by 15 logical partitions and define 32 devices (either on one port or across two ports), you have used 480 valid subchannels ($15 * 32 = 480$). If your definition requires more than 480 valid subchannels (with a maximum of 1920), then the user must explicitly turn off Outbound Priority Queueing on the CHPID definition by specifying CHPARM=02 in the IOCP or by specifying it in HCD. HCD will prevent a device definition that will cause the 480 subchannel limit to be broken. IOCP will issue an error message and not create an IOCDS if the limit is broken.
3. If you need to define more than 254 devices for an unshared OSD channel path, multiple control units must be defined. Specify a unique logical address for each control unit using the CUADD keyword.

Considerations for assigning the OSA portname

OSA Portname assignment for a QDIO implementation (CHPID Type of OSD) is important in the z/OS operating system. The rule for assigning a portname is the same regardless of the type of OSA adapter being implemented:

RULE: The portname of an OSA port must be unique on a CHPID.

This rule seems obvious, but you may find yourself confused when you contemplate a migration from certain configurations of the OSA-Express2 to an implementation of a new OSA-Express3. Consider Figure 4-14.

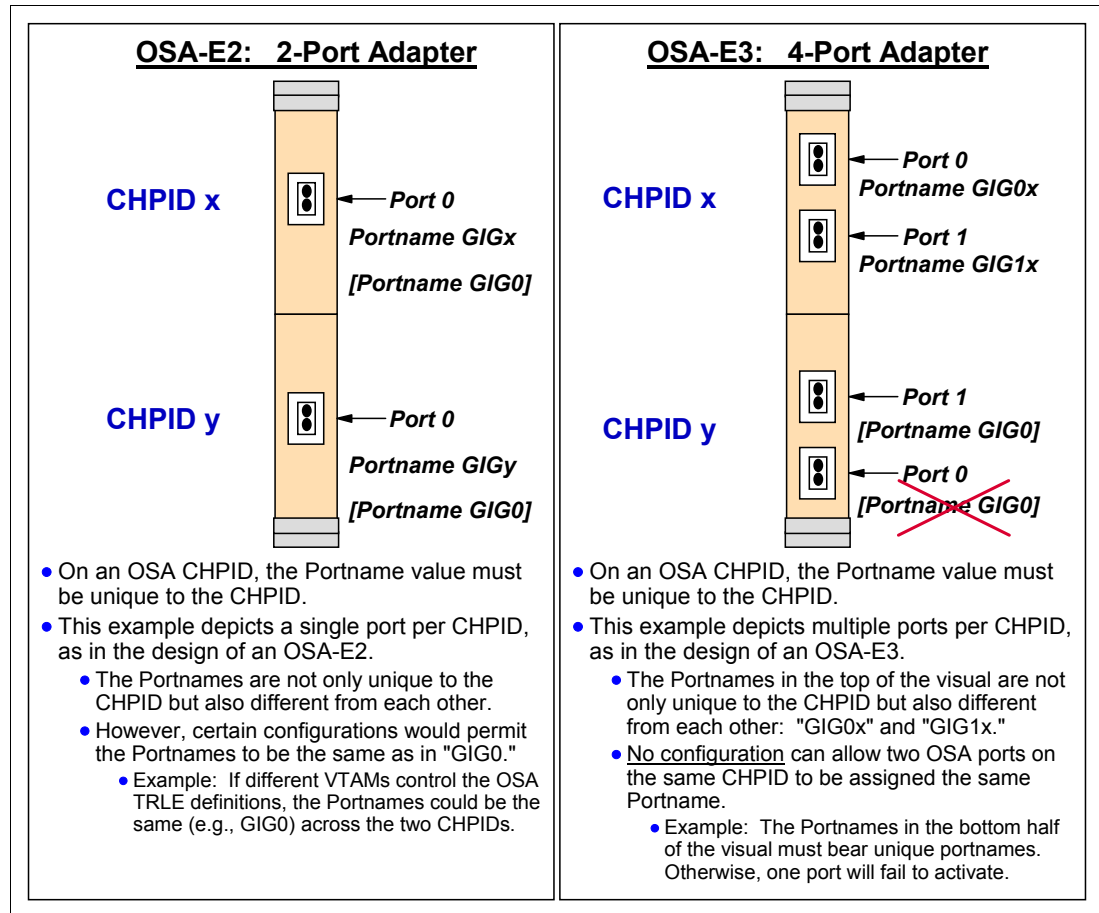


Figure 4-14 Providing unique portnames for OSA-Express ports

The figure shows that if you had attempted to move both ports named GIG0 to CHPIDy of the OSA-E3, one port would not activate because the names are no longer unique to the CHPID. The presence of duplicate names on the same CHPID would generate an SNA sense code of 8010311B.

4.4.2 Considerations for isolating traffic across a shared OSA port

VLANs, when properly implemented, can isolate traffic over a shared network and shared OSA port. The isolation is complete if all TCP/IP stacks that share an OSA port implement VLAN ID tagging and assign separate VLANIDs. For more information about this subject, consult Chapter 6, "VLAN and Virtual MAC support" on page 265.

Another method that is available to isolate traffic across a shared OSA port is *OSA Connection Isolation*. This method can be deployed with or without out assigning a VLAN ID or a VMAC to the OSA port. You can read more about this method in 4.5, “OSA-Express QDIO connectivity with Connection Isolation” on page 156.

When planning connectivity for a LAN environment, there might not be a requirement to isolate data traffic or services for certain servers or clients as we have shown in this scenario. In such cases, VLAN IDs can be omitted.

If there is a requirement for VLANs, however, we recommend adding the VLAN IDs to your IP addressing scheme to aid in the mapping of IP addresses to VLANs based on data traffic patterns or access to resources.

Also, to simplify administration and management of VLANs, consider using Generic Attribute VLAN Registration Protocol (GVRP) wherever possible. For details, refer to “VLAN support of Generic Attribute Registration Protocol (GVRP)” on page 123.

4.4.3 Configuring OSA-Express with VLAN ID

To implement OSA-Express (QDIO) in our environment, we performed these tasks:

1. Verify the switch port configuration.
2. Define a TRLE in VTAM to represent each OSA-Express port.

In the TCP/IP profile:

1. Create DEVICE and LINK or INTERFACE statements for each OSA-Express port.
2. Create a HOME address to each defined LINK.
3. Define the characteristics of each LINK statement using BSDROUTINGPARMS. You can code the BSDROUTINGPARMS statement even if you define the LINK characteristics in OMROUTE.

We explain these tasks in more detail in the following sections.

Verify the switch port configuration

It is important to be aware of the switch configuration and definitions to which the OSA-Express ports will be connected. You will need confirm the following information:

- The switch ports to which the OSA-Express ports are connected.

Table 4-3 shows the OSA-Express and switch port assignment with VLAN IDs and mode type in our configuration.

Table 4-3 OSA-Express and switch port assignment with VLAN IDs

OSA-Express port	Connects to switch	Switch port	VLAN ID (mode)
CHPID 02 (2080)	Switch 1	Interface GIGA 1/8	10 (Trunk mode)
CHPID 03 (20A0)	Switch 1	interface GIGA 1/41	10 (Trunk mode)
CHPID 04 (20C0)	Switch 1	Interface GIGA 1/43	11 (Trunk mode)
CHPID 05 (20E0)	Switch 1	Interface GIGA 1/19	11 (Trunk mode)

- The IP subnetwork and mask.
We used the following:
 - Subnetwork 10.1.2.0, mask 255.255.255.0 for VLAN 10
 - Subnetwork 10.1.3.0, mask 255.255.255.0 for VLAN 11
- The appropriate switch ports should be defined in trunk mode, as shown in Example 4-9.

Example 4-9 Switch port definition from Switch 1 port 1/41

```
interface GigabitEthernet1/41
  switchport
  switchport trunk encapsulation dot1q
  switchport mode trunk
  no ip address
```

Define a TRLE in VTAM to represent each OSA-Express port

Each OSA-Express port must have a TRLE definition defined; see Example 4-10. The PORTNAME **1** must match the device name of the DEVICE definition or the portname in INTERFACE definition in the TCP/IP profile. The PORTNUM **2** operand is optional (default 0), but required when defining the second port of an OSA-Express3. The statement MPCLEVEL **3** must be specified as QDIO.

Example 4-10 TRLE definition

```
OSA2080  VBUILD TYPE=TRL
OSA2080P  TRLE  LNCTL=MPC,                      *
              READ=2080,                         *
              WRITE=2081,                        *
              DATAPATH=(2082-2088),              *
1 PORTNAME=OSA2080,                             *
2 PORTNUM=0,                                     *
3 MPCLEVEL=QDIO
```

For all OSA-Express ports in our scenarios, we used the following PORTNAMES:

- OSA2080
- OSA20A0
- OSA20C0
- OSA20E0

Create DEVICE and LINK or INTERFACE statements for each OSA-Express port

The next step is to create the device and link or interface statements for each OSA-Express port, as shown in Example 4-11 on page 150.

Note: We encourage and recommend the use of the INTERFACE statement, as it groups all the definitions required in one spot, while DEVICE and LINK require that the IP address be assigned in the HOME list.

The device definition of an OSA-Express port must be set as an MPCIPA device type **1**. The link definition describes the type of transport used (in our case, QDIO Ethernet, defined as IPAQENET **2**). VLAN ID **3** defines the VLAN number the packets will be tagged with as they are being sent out to the switch.

Note: You can only define a single VLAN per each OSA port with device and link statement. If you want to define multiple VLANs on a single OSA port, you need to define it with the interface statement.

Example 4-11 OSA-Express device and link definitions

```

;OSA DEFINITIONS
;TRL MAJ NODE: OSA2080,OSA20A0,OSA20C0,AND OSA20E0
DEVICE OSA2080    MPCIPA 1
LINK  OSA2080L    IPAQENET 2 OSA2080 VLANID 10 3
DEVICE OSA20C0    MPCIPA
LINK  OSA20C0L    IPAQENET    OSA20C0 VLANID 11
DEVICE OSA20E0    MPCIPA
LINK  OSA20E0L    IPAQENET    OSA20E0
DEVICE OSA20A0    MPCIPA
LINK  OSA20A0L    IPAQENET    OSA20A0

```

The alternative interface statement of OSA-Express ports combines the definitions otherwise coded in the device, link, home, beginroutes and bsdroutingparms statements, and as such requires a label **1**, the type of transport used (QDIO Ethernet, defined as IPAQENET **2** is the only type allowed for IPv4 devices), a portname **3** matching the TRLE portname, an IP address and optional subnetmask **4**, optional MTU size **5**, VLANID **6**, VMAC **7** (required when setting multiple VLANs on the same physical OSA port) and SOURCEVIPAIN **8** which associates a specific VIPA with *this* interface.

Note: If SOURCEVIPAIN is coded, the whole INTERFACE definition block must be defined in PROFILE *after* the VIPA DEVICE and LINK statements are defined.

Example 4-12 OSA-Express interface definition

```

INTERFACE OSA20A0I 1
DEFINE IPAQENET 2
PORTNAME OSA20A0 3
IPADDR 10.1.2.12/24 4
MTU 1492 5
VLANID 20 6
VMAC 7
SOURCEVIPAIN VIPA2L 8

```

Create a HOME address to each defined LINK

If you are not implementing the connection with the INTERFACE statement, you must assign an IP address to the LINK of each DEVICE/LINK pair. Each link configured must have its own IP address configured on the HOME statement of the TCP/IP profile. Our OSA-Express ports are defined with the IP addresses shown in Example 4-13.

Note: This step is not required when defining OSA ports through the INTERFACE statement.

Example 4-13 OSA-Express HOME addresses

```
HOME
  10.1.2.11      OSA2080L
  10.1.3.11      OSA20C0L
  10.1.3.12      OSA20E0L
  10.1.2.12      OSA20A0L
```

Define the characteristics of each LINK statement using BSDROUTINGPARMS

To define the link characteristics, such as MTU size **1** and subnet mask **2**, we used the BSDROUTINGPARMS statements (see Example 4-14).

Note: This step is not required when defining OSA ports through the INTERFACE statement.

If not supplied, defaults are used from static routing definitions in BEGINROUTES or the OMROUTE configuration (dynamic routing definitions), if implemented.

If the link characteristics, BEGINROUTES statements, or the OMROUTE configuration are not defined, then the stack's interface layer (based on hardware capabilities) and the characteristics of devices and links are used. This, however, might not provide the performance or function desired.

Example 4-14 BSDRoutingparms statements

```
BSDROUTINGPARMS TRUE
; Link name      MTU      Cost metric  Subnet Mask  Dest address
   VIPA1L 1492 1        0          255.255.255.252  0
   OSA2080L 1492        0          255.255.255.0   2 0
   OSA20A0L 1492        0          255.255.255.0   0
   OSA20C0L 1492        0          255.255.255.0   0
   OSA20E0L 1492        0          255.255.255.0   0
ENDBSDROUTINGPARMS
```

Note: Static and dynamic routing definitions will override or replace the link characteristics defined through the BSDROUTINGPARMS statements. Refer to Chapter 5, "Routing" on page 205 for more information about static and dynamic routing.

4.4.4 Verifying the connectivity status

In this section, we verify the status of the OSA devices defined to the TCP/IP stack and VTAM.

Verifying the device status in TCP/IP stack

To verify the status of all devices being activated in the TCP/IP stack we use the NETSTAT command with the DEVLIST option, as shown in Example 4-15.

Example 4-15 Using command D TCPIP,TCPIPA,N,DEV to verify the Device status

```
D TCPIP,TCPIPA,N,DEV
..... Lines
deleted
INTFNAME: OSA2080I          INTFTYPE: IPAQENET  INTFSTATUS: READY
PORTNAME: OSA2080          DATAPATH: 2082      DATAPATHSTATUS: READY
CHPIDTYPE: OSD
SPEED: 0000001000
IPBROADCASTCAPABILITY: NO
VMACADDR: 02000C776873     VMACORIGIN: OSA    VMACROUTER: LOCAL
ARPOFFLOAD: YES            ARPOFFLOADINFO: YES
CFGMTU: 1492               ACTMTU: 1492
IPADDR: 10.1.2.11/24
VLANID: 10                 VLANPRIORITY: DISABLED
DYNVLANREGCFG: NO          DYNVLANREGCAP: YES
READSTORAGE: GLOBAL (4096K)
INBPERF: BALANCED
CHECKSUMOFFLOAD: YES       SEGMENTATIONOFFLOAD: YES
SECCLASS: 255              MONSYSPLEX: NO
ISOLATE: NO                OPTLATENCYMODE: NO
MULTICAST SPECIFIC:
MULTICAST CAPABILITY: YES
GROUP          REFCNT          SRCFLTMD
-----
224.0.0.1      0000000001      EXCLUDE
INTERFACE STATISTICS:
BYTESIN                = 0
INBOUND PACKETS        = 0
INBOUND PACKETS IN ERROR = 0
INBOUND PACKETS DISCARDED = 0
INBOUND PACKETS WITH NO PROTOCOL = 0
BYTESOUT               = 168
OUTBOUND PACKETS       = 2
OUTBOUND PACKETS IN ERROR = 0
OUTBOUND PACKETS DISCARDED = 0
..... Lines
deleted
```

Displaying TCP/IP device resources in VTAM

The device drivers for TCP/IP are provided by VTAM. When CS for z/OS IP devices are activated, there must be an equivalent Transport Resource List Element (TRLE) defined to VTAM. The devices that are exclusively used by z/OS Communications Server IP have TRLEs that are automatically generated for them.

Because the device driver resources are provided by VTAM, you have the ability to display the resources using VTAM display commands. To display a list of all TRLEs active in VTAM, use the command D NET,TRL, as shown in Example 4-16.

Example 4-16 D NET,TRL command output

```

D NET,TRL
IST350I DISPLAY TYPE = TRL 135
IST924I -----
IST1954I TRL MAJOR NODE = ISTTRL
IST1314I TRLE = IUTIQDF6 STATUS = ACTIV CONTROL = MPC
IST1314I TRLE = IUTIQDF5 STATUS = ACTIV CONTROL = MPC
IST1314I TRLE = IUTIQDF4 STATUS = ACTIV CONTROL = MPC
IST1314I TRLE = ISTT3033 STATUS = ACTIV CONTROL = XCF
IST1314I TRLE = ISTT3032 STATUS = ACTIV CONTROL = XCF
IST1314I TRLE = ISTT3031 STATUS = ACTIV CONTROL = XCF
IST1314I TRLE = IUTIQDIO STATUS = ACTIV CONTROL = MPC
IST1314I TRLE = IUTSAMEH STATUS = ACTIV CONTROL = MPC
IST1454I 8 TRLE(S) DISPLAYED
IST924I -----
IST1954I TRL MAJOR NODE = TRLTNET
IST1314I TRLE = MPCNET STATUS = ACTIV CONTROL = MPC
IST1454I 1 TRLE(S) DISPLAYED
IST924I -----
IST1954I TRL MAJOR NODE = OSA2000
IST1314I TRLE = OSA2000P STATUS = NEVAC CONTROL = MPC
IST1454I 1 TRLE(S) DISPLAYED
IST924I -----
IST1954I TRL MAJOR NODE = OSA2020
IST1314I TRLE = OSA2020P STATUS = NEVAC CONTROL = MPC
IST1454I 1 TRLE(S) DISPLAYED
IST924I -----
IST1954I TRL MAJOR NODE = OSA2080
IST1314I TRLE = OSA2080T STATUS = ACTIV CONTROL = MPC
IST1454I 1 TRLE(S) DISPLAYED
IST924I -----
IST1954I TRL MAJOR NODE = OSA20A0
IST1314I TRLE = OSA20A0P STATUS = ACTIV CONTROL = MPC
IST1454I 1 TRLE(S) DISPLAYED
IST924I -----
IST1954I TRL MAJOR NODE = OSA20A1
IST1314I TRLE = OSA20A1P STATUS = NEVAC CONTROL = MPC
IST1454I 1 TRLE(S) DISPLAYED
IST924I -----
IST1954I TRL MAJOR NODE = OSA20C0
IST1314I TRLE = OSA20C0P STATUS = ACTIV CONTROL = MPC
IST1454I 1 TRLE(S) DISPLAYED
IST924I -----
IST1954I TRL MAJOR NODE = OSA20E0
IST1314I TRLE = OSA20E0P STATUS = ACTIV CONTROL = MPC
IST1454I 1 TRLE(S) DISPLAYED
IST924I -----
IST1954I TRL MAJOR NODE = OSA2100
IST1314I TRLE = OSA2100T STATUS = ACTIV CONTROL = MPC
IST1454I 1 TRLE(S) DISPLAYED
IST924I -----

```

```

IST1954I TRL MAJOR NODE = OSA2120
IST1314I TRLE = OSA2120T STATUS = ACTIV          CONTROL = MPC
IST1454I 1 TRLE(S) DISPLAYED
IST314I END

```

You can also display information of TRLEs grouped by control type, such as MPC or XCF devices, as shown in Example 4-17.

Example 4-17 D NET,TRL,CONTROL=MPC

D NET,TRL,CONTROL=MPC

```

IST350I DISPLAY TYPE = TRL 276
IST924I -----
IST1954I TRL MAJOR NODE = ISTTRL
IST1314I TRLE = IUTIQDF6 STATUS = ACTIV          CONTROL = MPC
IST1314I TRLE = IUTIQDF5 STATUS = ACTIV          CONTROL = MPC
IST1314I TRLE = IUTIQDF4 STATUS = ACTIV          CONTROL = MPC
IST1314I TRLE = IUTIQDIO STATUS = ACTIV          CONTROL = MPC
IST1314I TRLE = IUTSAMEH STATUS = ACTIV          CONTROL = MPC
IST1454I 5 TRLE(S) DISPLAYED
IST924I -----
IST1954I TRL MAJOR NODE = TRLTNET
IST1314I TRLE = MPCNET STATUS = ACTIV          CONTROL = MPC
IST1454I 1 TRLE(S) DISPLAYED
IST924I -----
IST1954I TRL MAJOR NODE = OSA2000
IST1314I TRLE = OSA2000P STATUS = NEVAC          CONTROL = MPC
IST1454I 1 TRLE(S) DISPLAYED
IST924I -----
IST1954I TRL MAJOR NODE = OSA2020
IST1314I TRLE = OSA2020P STATUS = NEVAC          CONTROL = MPC
IST1454I 1 TRLE(S) DISPLAYED
IST924I -----
IST1954I TRL MAJOR NODE = OSA2080
IST1314I TRLE = OSA2080T STATUS = ACTIV          CONTROL = MPC
IST1454I 1 TRLE(S) DISPLAYED
IST924I -----
IST1954I TRL MAJOR NODE = OSA20A0
IST1314I TRLE = OSA20A0P STATUS = ACTIV          CONTROL = MPC
IST1454I 1 TRLE(S) DISPLAYED
IST924I -----
IST1954I TRL MAJOR NODE = OSA20A1
IST1314I TRLE = OSA20A1P STATUS = NEVAC          CONTROL = MPC
IST1454I 1 TRLE(S) DISPLAYED
IST924I -----
IST1954I TRL MAJOR NODE = OSA20C0
IST1314I TRLE = OSA20C0P STATUS = ACTIV          CONTROL = MPC
IST1454I 1 TRLE(S) DISPLAYED
IST924I -----
IST1954I TRL MAJOR NODE = OSA20E0
IST1314I TRLE = OSA20E0P STATUS = ACTIV          CONTROL = MPC
IST1454I 1 TRLE(S) DISPLAYED
IST924I -----
IST1954I TRL MAJOR NODE = OSA2100
IST1314I TRLE = OSA2100T STATUS = ACTIV          CONTROL = MPC
IST1454I 1 TRLE(S) DISPLAYED

```

```

IST924I -----
IST1954I TRL MAJOR NODE = OSA2120
IST1314I TRLE = OSA2120T STATUS = ACTIV          CONTROL = MPC
IST1454I 1 TRLE(S) DISPLAYED
IST314I END

```

We can also get specific information about a single TRLE, using the TRLE name as shown in Example 4-18, for an OSA-Express device.

Example 4-18 D NET,TRL,TRLE=OSA2080T

D NET,TRL,TRLE=OSA2080T

```

IST075I NAME = OSA2080T, TYPE = TRLE 336
IST1954I TRL MAJOR NODE = OSA2080
IST486I STATUS= ACTIV, DESIRED STATE= ACTIV
IST087I TYPE = LEASED          , CONTROL = MPC , HPDT = YES
IST1715I MPCLEVEL = QDIO      MPCUSAGE = SHARE
IST2263I PORTNAME = OSA2080   PORTNUM = 0   OSA CODE LEVEL = 000C
IST2337I CHPID TYPE = OSD     CHPID = 02
IST1577I HEADER SIZE = 4096 DATA SIZE = 0 STORAGE = ***NA***
IST1221I WRITE DEV = 2081 STATUS = ACTIVE      STATE = ONLINE
IST1577I HEADER SIZE = 4092 DATA SIZE = 0 STORAGE = ***NA***
IST1221I READ  DEV = 2080 STATUS = ACTIVE      STATE = ONLINE
IST924I -----
IST1221I DATA DEV = 2082 STATUS = ACTIVE      STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST1717I ULPID = TCPIPA
IST2309I ACCELERATED ROUTING ENABLED
IST2331I QUEUE  QUEUE      READ
IST2332I ID     TYPE       STORAGE
IST2205I -----
IST2333I RD/1    PRIMARY   4.0M(64 SBALS)
IST2305I NUMBER OF DISCARDED INBOUND READ BUFFERS = 0
IST1757I PRIORITY1: UNCONGESTED PRIORITY2: UNCONGESTED
IST1757I PRIORITY3: UNCONGESTED PRIORITY4: UNCONGESTED
IST2190I DEVICEID PARAMETER FOR OSAENTA TRACE COMMAND = 01-01-00-02
IST1801I UNITS OF WORK FOR NCB AT ADDRESS X'0F512010'
IST1802I P1 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P2 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P3 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P4 CURRENT = 0 AVERAGE = 2 MAXIMUM = 2
IST924I -----
IST1221I DATA DEV = 2083 STATUS = ACTIVE      STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST1717I ULPID = TCPIPC
IST2309I ACCELERATED ROUTING ENABLED
IST2331I QUEUE  QUEUE      READ
IST2332I ID     TYPE       STORAGE
IST2205I -----
IST2333I RD/1    PRIMARY   4.0M(64 SBALS)
IST2305I NUMBER OF DISCARDED INBOUND READ BUFFERS = 0
IST1757I PRIORITY1: UNCONGESTED PRIORITY2: UNCONGESTED
IST1757I PRIORITY3: UNCONGESTED PRIORITY4: UNCONGESTED
IST2190I DEVICEID PARAMETER FOR OSAENTA TRACE COMMAND = 01-01-00-03
IST1801I UNITS OF WORK FOR NCB AT ADDRESS X'0FCF0010'
IST1802I P1 CURRENT = 0 AVERAGE = 1 MAXIMUM = 2

```

```

IST1802I P2 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P3 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P4 CURRENT = 0 AVERAGE = 1 MAXIMUM = 1
IST924I -----
IST1221I DATA DEV = 2084 STATUS = ACTIVE STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST1717I ULPID = TCPIPB
IST2310I ACCELERATED ROUTING DISABLED
IST2331I QUEUE QUEUE READ
IST2332I ID TYPE STORAGE
IST2205I -----
IST2333I RD/1 PRIMARY 4.0M(64 SBALS)
IST2305I NUMBER OF DISCARDED INBOUND READ BUFFERS = 0
IST1757I PRIORITY1: UNCONGESTED PRIORITY2: UNCONGESTED
IST1757I PRIORITY3: UNCONGESTED PRIORITY4: UNCONGESTED
IST2190I DEVICEID PARAMETER FOR OSAENTA TRACE COMMAND = 01-01-00-04
IST1801I UNITS OF WORK FOR NCB AT ADDRESS X'0F03F010'
IST1802I P1 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P2 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P3 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P4 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST924I -----
IST1221I DATA DEV = 2085 STATUS = RESET STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST924I -----
IST1221I DATA DEV = 2086 STATUS = RESET STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST924I -----
IST1221I DATA DEV = 2087 STATUS = RESET STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST924I -----
IST314I END

```

4.5 OSA-Express QDIO connectivity with Connection Isolation

Many customers share OSA-Express ports across logical partitions, especially if capacity is not an issue. Each stack sharing the OSA port registers certain IP addresses and multicast groups with the OSA.

Note: You may wish to revisit the discussion of IPv4 address registration in “OSA-Express QDIO IPv4 address registration” on page 121.

For performance reasons, the OSA-Express bypasses the LAN and routes packets directly between the stacks when possible. Examine Figure 4-15, where you see two TCP/IP stacks, TCPIPA and TCPIPB, which share the same OSA port connected to subnet 10.1.2.0/24.

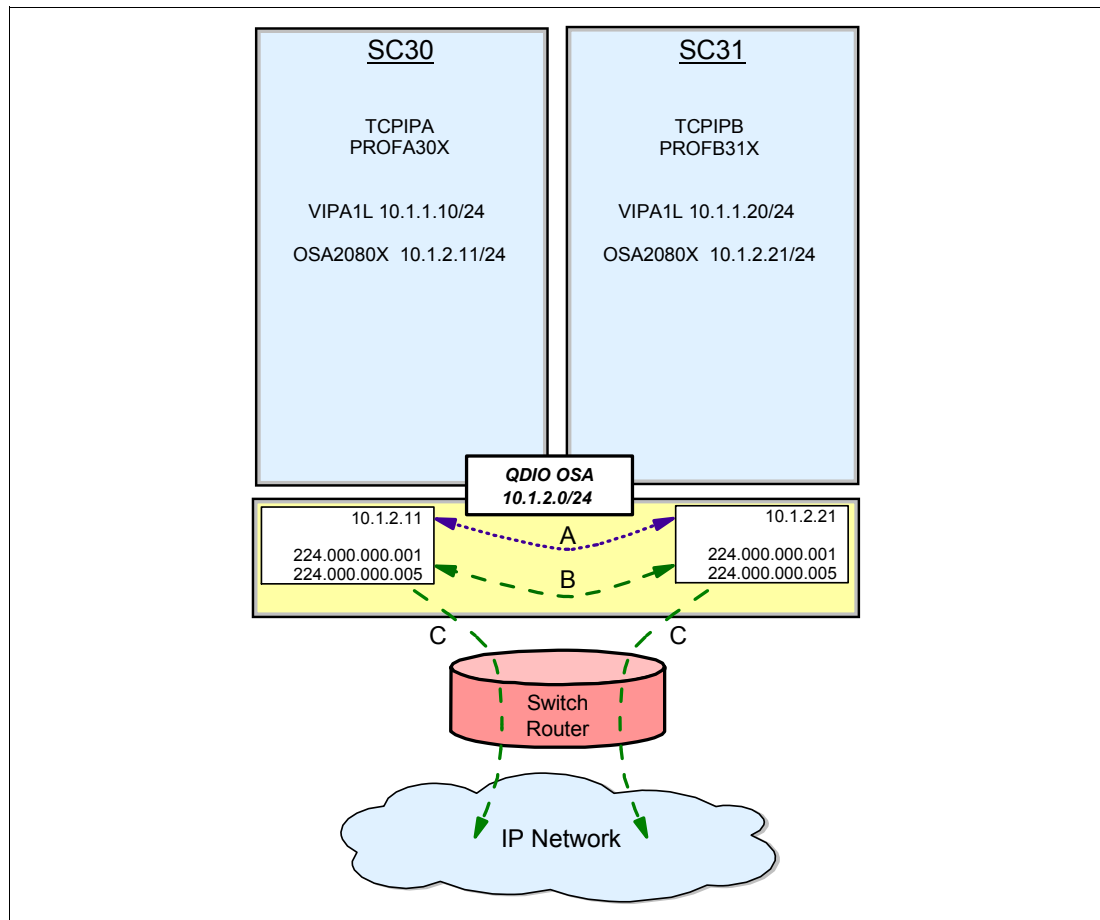


Figure 4-15 Routing paths over an OSA port

For performance reasons, the OSA-Express bypasses the LAN and routes packets directly between the stacks when possible. For unicast packets, OSA internally routes the packet when the next-hop IP address is registered on the same LAN or VLAN by another stack sharing the OSA port. Figure 4-15 illustrates examples of this action, where:

- ▶ **A**: You see how TCPIPA routes a packet to 10.1.2.21 in TCPIP over the OSA port without exiting out onto the LAN because the next hop to reach the destination is registered in the OSA Address Table (OAT); the TCPIPA routing table indicates that the destination can be reached by hopping through the direct connection to the 10.1.2.0/24 network.
- ▶ **B**: For multicast (for example, OSPF protocol packets), OSA internally routes the packet to all sharing stacks on the same LAN or VLAN that registered the multicast group. Note how TCPIPA and TCPIP have each registered multicast addresses for OSP (224.000.000.00n) in the OSA port.
- ▶ **C**: OSA also sends the multicast/broadcast packet to the LAN. For broadcast (not depicted), OSA internally routes the packet to all sharing stacks on the same LAN or VLAN.

Thus, you see that stacks sharing an OSA-Express port can communicate over the OSA. Some customers may express concerns about this efficient communication path and wish to disable it because traffic flowing internally through the OSA adapter bypasses any security features implemented on the external LAN. For example, the customer may have exploited the virtualization features of the System z and of 10 Gigabit OSA adapters to build a demilitarized zone (DMZ) on several LPARs of a System z as well as several production

LPARs on the same System z footprint. Although they can implement firewall and Intrusion Detection technologies within the LPARs to isolate the two zones (DMZ and Production) from each other, they may have already invested in external security mechanisms on the LAN. If traffic through a shared OSA bypasses the security on the LAN, they need to find a way to prevent the internal routing across the shared OSA path.

Several network designs are available to provide isolation and force the traffic to bypass the shared OSA path or to be prevented from using it:

- ▶ Implement IP filtering on the stacks in the adjacent zones by exploiting z/OS Policy Agent with IP filtering and Intrusion Detection Services (IDS).
- ▶ Implement routing filters that block the advertisement of certain routing zones to parts of the network from which they should remain concealed. Examples of such features are OSPF range checking, RIP, or EIGRP routing filters.
- ▶ Implement Policy Based Routing (PBR) to eliminate the internal OSA path where it is not desired.
- ▶ Define static routes so that paths to a stack sharing the OSA are forced to hop through a router on the LAN.
- ▶ Configure the TCP/IP stacks in separate zones (IP subnets) with separate VLANs that extend into the stacks themselves.
- ▶ Implement *OSA Connection Isolation*.

4.5.1 Description of Connection Isolation

Some environments require strict controls for routing data traffic between servers or nodes. In certain cases, the LPAR-to-LPAR capability of a shared OSA port can prevent such controls from being enforced. For example, you may need to ensure that traffic flowing through the OSA adapter does not bypass firewalls or intrusion detection systems implemented on the external LAN. We have described several ways to isolate traffic from different LPARs on a shared OSA port, with one of these methods being *OSA Connection Isolation*.

The feature is called *OSA Connection Isolation* in z/OS, but it is also available in z/VM, where it is called *QDIO data connection isolation* or *VSWITCH port isolation*. It allows you to disable the internal routing on a QDIO connection basis, providing a means for creating security zones and preventing network traffic between the zones. It also provides extra insurance against a misconfiguration that might otherwise allow such traffic to flow as in the case of an incorrectly defined IP filter. With interface isolation, internal routing can be controlled on an LPAR basis. When interface isolation is enabled, the OSA will discard any packets destined for a z/OS LPAR that is registered in the OAT as isolated.

4.5.2 Dependencies for Connection Isolation

QDIO interface isolation is supported by Communications Server for z/OS V1R12 and all OSA-Express3 and OSA-Express2 features on System z10, and by all OSA-Express2 features on System z9, with an MCL update. Refer to the appropriate Preventive Service Planning bucket for details regarding your System z server.

Coding ISOLATE on your INTERFACE statement enables the function. It tells the OSA-Express not to allow communications to this stack other than over the LAN. OSA-Express requires that both stacks sharing the port be non-isolated for direct routing to occur.

Because this function is specific to security, an OSA-Express interface that does not support the Connection Isolation function cannot be activated. Examine the messages at **1** and **2** in Example 4-19 which show an unsuccessful activation attempt for a QDIO interface whose OSA does not support the ISOLATE function that was coded on it.

Example 4-19 Failure to activate an OSA interface that does not support the ISOLATE feature

```
V TCPIP,TCPIPF,START,OSA2080X
EZZ0060I PROCESSING COMMAND: VARY TCPIP,,START,OSA2080X
EZZ0053I COMMAND VARY START COMPLETED SUCCESSFULLY
EZD0022I INTERFACE OSA2080X DOES NOT SUPPORT THE ISOLATE FUNCTION      1
EZZ4341I DEACTIVATION COMPLETE FOR INTERFACE OSA2080X      2
```

To eliminate the ISOLATE specification on the device so that you can successfully activate it, you must first STOP the interface before using the V TCPIP,,OBEYFILE command to modify the ISOLATE parameter.

4.5.3 Considerations for Connection Isolation

When Connection Isolation is in effect on either or both endpoints of a connection on a shared OSA port, OSA-Express will discard any packets when the next-hop address is registered in the OSA by a sharing stack, that is, OSA discards unicast packets which previously qualified for internal routing. It also ceases to route internal multicast or broadcast packets. It does, however, continue to send the multicast or broadcast packets to the LAN. OSA-Express requires that both stacks sharing the port be non-isolated for direct routing to occur.

If you have implemented static routing where Connection Isolation is in effect, it is simple to code the appropriate routing statements to bypass the direct path through the OSA. If you are running a dynamic routing protocol, you may see routing errors when the routing protocol attempts to send packets over the ISOLATED OSA port. Such errors are “working as designed” when ISOLATION has been introduced into the configuration.

Dynamic routing protocol implementations with RIP or OSPF require careful planning on LANs where OSA-Express connection isolation is in effect; the dynamic routing protocol learns of the existence of the direct path but is unaware of the isolated configuration, which renders the direct path across the OSA port to the registered target unusable. If the direct path that is operating as ISOLATED is selected, you will experience routing failures.

If the visibility of such errors is undesirable, you can take other measures to avoid the failure messages. If you are simply attempting to bypass the direct route in favor of another, indirect route, you can accomplish this as well with some thoughtful design.

For example, you might purposely bypass the direct path by using Policy Based Routing (PBR) or by coding static routes that supersede the routes learned by the dynamic routing protocol. You might adjust the weights of connections to favor alternate interfaces over the interfaces that have been coded with ISOLATE.

Static routes to override the direct OSA path

Examine the sample network diagram in Figure 4-16.

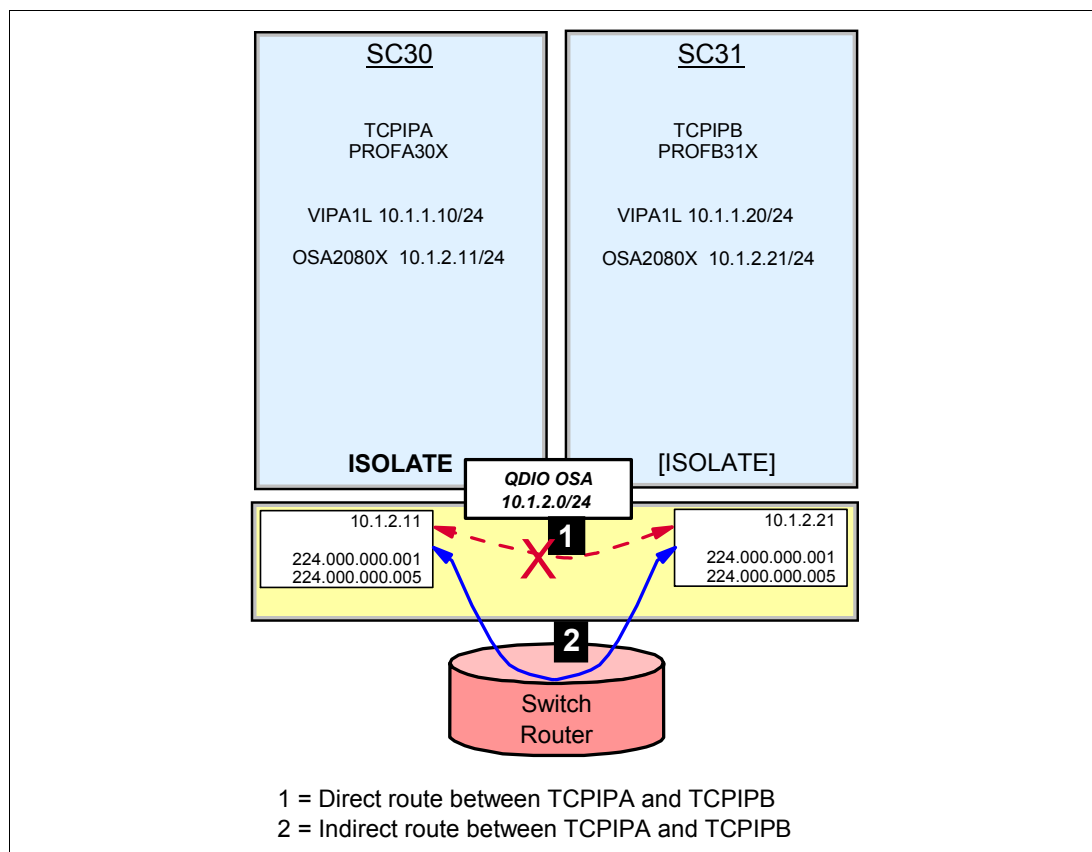


Figure 4-16 Routing TCPIPA and TCPIPB: block direct OSA path and multi-hop static route

In Figure 4-16 you see two TCP/IP stacks: TCPIPA and TCPIPB. They share an OSA port on network 10.1.2.0/24. Both stacks are running a dynamic routing protocol that informs them that there is a direct path (1) through the OSA port between each other. The routing protocol knows nothing of the ISOLATE function that was introduced to prevent packets from using the direct route. (ISOLATE must be coded on only one of the two TCP/IP stacks, although you can code it on both in this diagram.)

Another path between the two TCP/IP stacks is available through an external, next-hop router (2). However, the dynamic routing protocol does not apprise the TCP/IP stack of this route's existence because it is not the shortest path. Therefore, when a packet is sent from TCPIPA to TCPIPB, the stack's routing table will always try to send that packet through the shortest path; the send will not be successful because the stacks have been ISOLATED from each other over the OSA port.

If TCPIPA and TCPIPB are not to communicate with each other at all, then there is no need to alter the appearance of the existing routing table. A route failure in this instance could be desirable. In order to produce a message that explains that the two endpoints are ineligible for routing to each other at all, you would probably introduce an IP filter. (Note that the routing failure itself has no failure message that indicates that ISOLATE is at fault.)

If, however, TCPIPA and TCPIPB do need to exchange information, you will need to deploy an effective route that bypasses the direct route between them. Therefore, at TCPIPA, you might add a non-replaceable static route to an IP address in TCPIPB; the static route in the

BEGINROUTES block points to the next-hop router on the path indicated with (2) in Figure 4-16 on page 160.

The effect of ICMP redirect packets

To avoid the override of the ICMP redirect packets that would most likely occur from the router to the originating host, you need to disable the receipt of ICMP redirects in the IP stacks or disable ICMP redirects at the router. If you are using OMPROUTE, ICMP redirects are automatically disabled, as evidenced by the message that appears during OMPROUTE initialization:

```
EZZ7475I ICMP WILL IGNORE REDIRECTS DUE TO ROUTING APPLICATION BEING ACTIVE
```

An alternate path that is more desirable than the direct path

If you do choose to have two TCP/IP stacks that are sharing an OSA port communicate with each other, and if you do not wish to introduce the static routes just described, another alternative is to provide another path that has, for example, a lower weighted path. See the diagram in Figure 4-17.

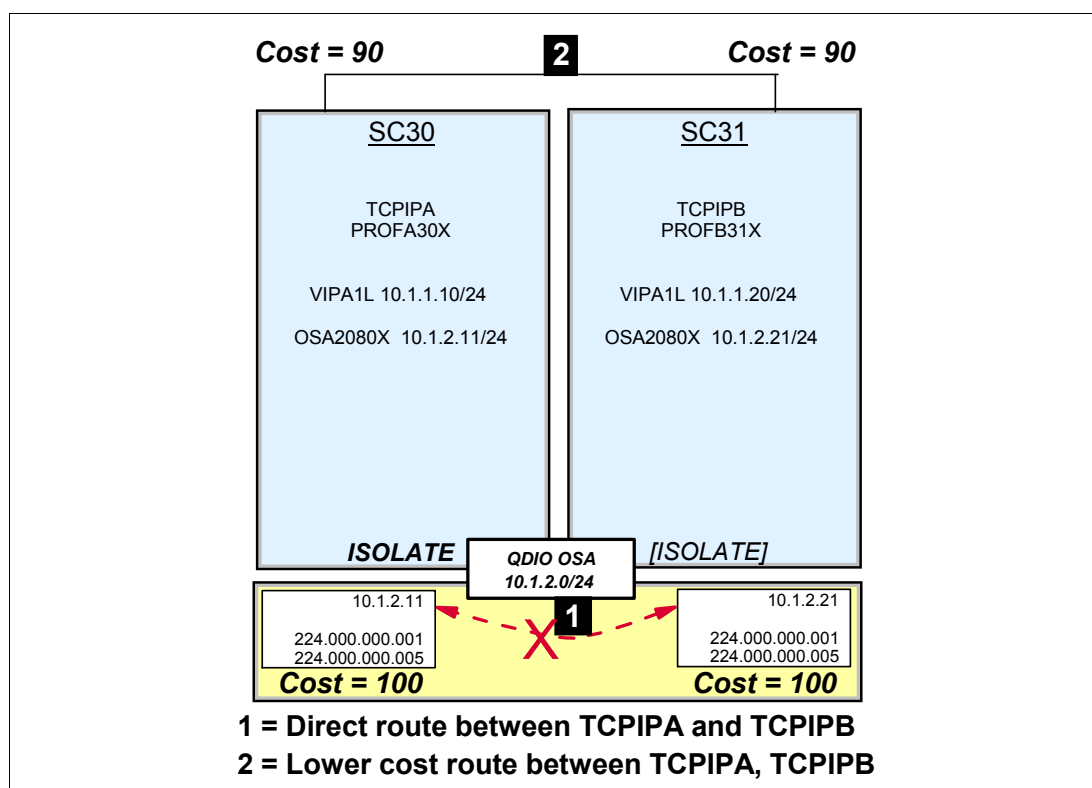


Figure 4-17 Routing TCPIPA and TCPIPB: block direct path provides a lower-cost alternate path

In Figure 4-17 on page 161 you see a lower-cost route at (2). The dynamic routing protocol continues to run, but now the favored route is the one over HiperSockets, XCF, CTC, or over an alternate LAN connection. Although the dynamic routing protocol continues its awareness of the direct OSA path, it prefers the path at (2).

Altering the routing table with Policy-Based Routing (PBR)

You may also wish to deploy Policy-Based Routing in order to bypass direct routes. Refer to Chapter 4, “Policy Agent”, in *IBM z/OS V1R12 Communications Server TCP/IP Implementation Volume 4: Security*, SG24-7899 for more information about how to accomplish this task.

4.5.4 Configuring OSA-Express with Connection Isolation

In Figure 4-18, you see the test network that will be the basis of our testing for OSA Express Connection Isolation. This diagram shows you how we have depicted the shared OSA environment in other manuals in this series.

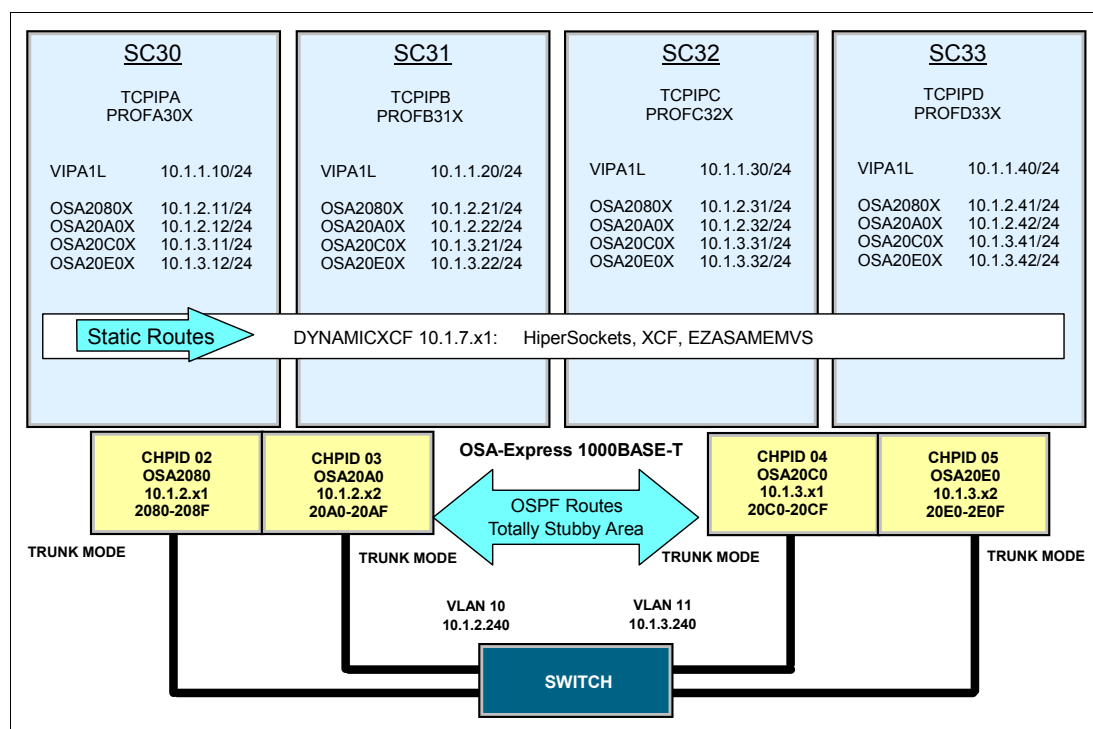


Figure 4-18 Stacks started with test profiles PROFA30X, B31X, C32X, and D33X

All of the System z TCP/IP stacks are members of an OSPF Totally Stubby Network. Note that the TCP/IP profiles at each stack are named PROFA30X, PROFB31X, PROFC32X, and PROFD33X. Each stack shares each of the four OSA ports depicted. In VLAN 10 and on Subnet 10.1.2.0/24, you see two OSA ports on each stack: OSA2080 and OSA20A0. In VLAN 11 and on Subnet 10.1.3.0/24, you see two OSA ports on each stack: OSA20C0 and OSA20E0. Each stack also has a static VIPA in subnet 10.1.1.0/24. The OSA and VIPA interfaces are all advertised with OSPF protocols. However, the connections implemented with the DYNAMICXCF keyword use only static routing.

Examine the revised visual in Figure 4-19. It attempts to depict in a clearer fashion than Figure 4-18 on page 162 how the OSA ports are shared across the four LPARs. Each TCP/IP stack has two connections into subnet 10.1.2.0/24 and two into subnet 10.1.3.0/24.

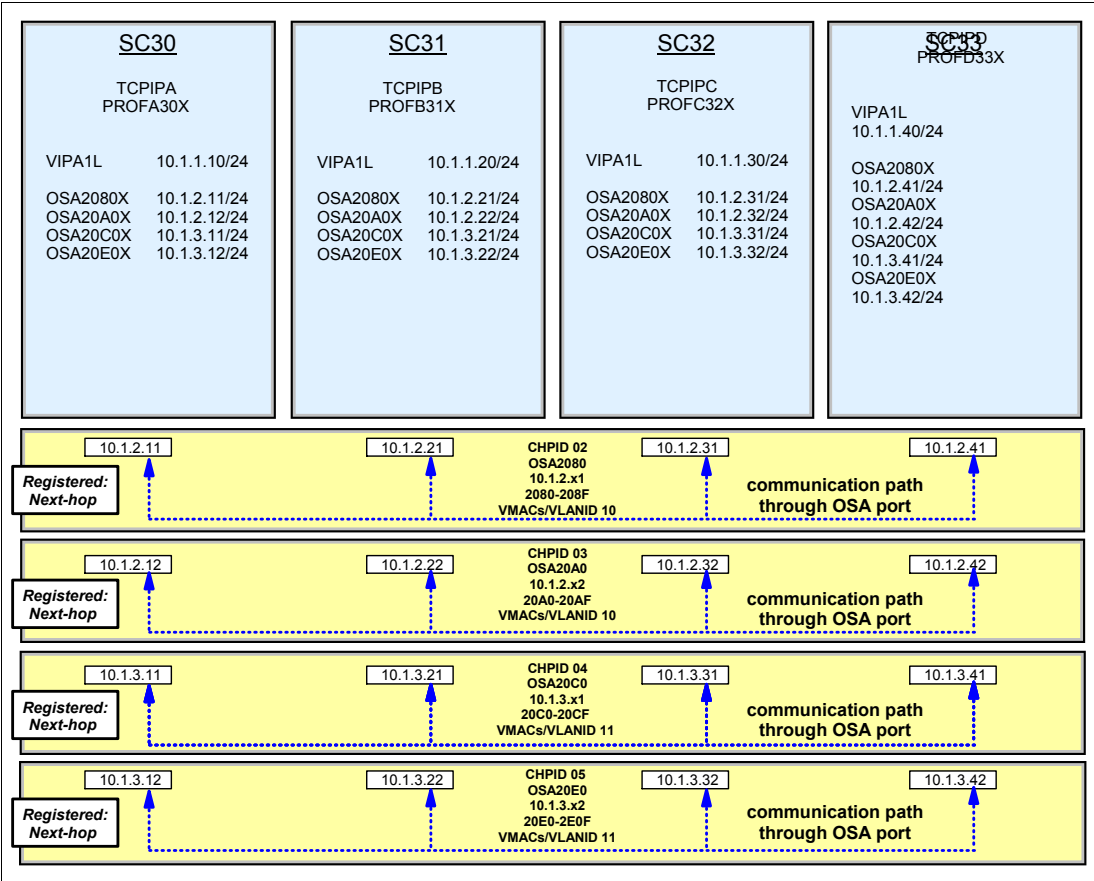


Figure 4-19 OAT entries for the stacks sharing the four OSA ports

The revised diagram shows you how stacks communicate with each other over the shared OSA ports when the next-hop router IP address is registered in the OSA. For performance reasons, the OSA-Express bypasses the LAN and routes packets directly between the stacks when possible.

4.5.5 Verifying Connection Isolation on OSA2080X

In this section, we discuss verifying Connection Isolation on OSA2080X.

Scenario for testing

To simplify the testing of the connection isolation scenario, we started only the connections to the OSA on CHPID 2. We then modified the existing configuration to implement OSA Connection Isolation only on TCPIPA and TCPIPB on CHPID 2. Connection Isolation was not exploited on CHPID 2 for TCPIPC and TCPIPD.

The new configuration is depicted in Figure 4-20.

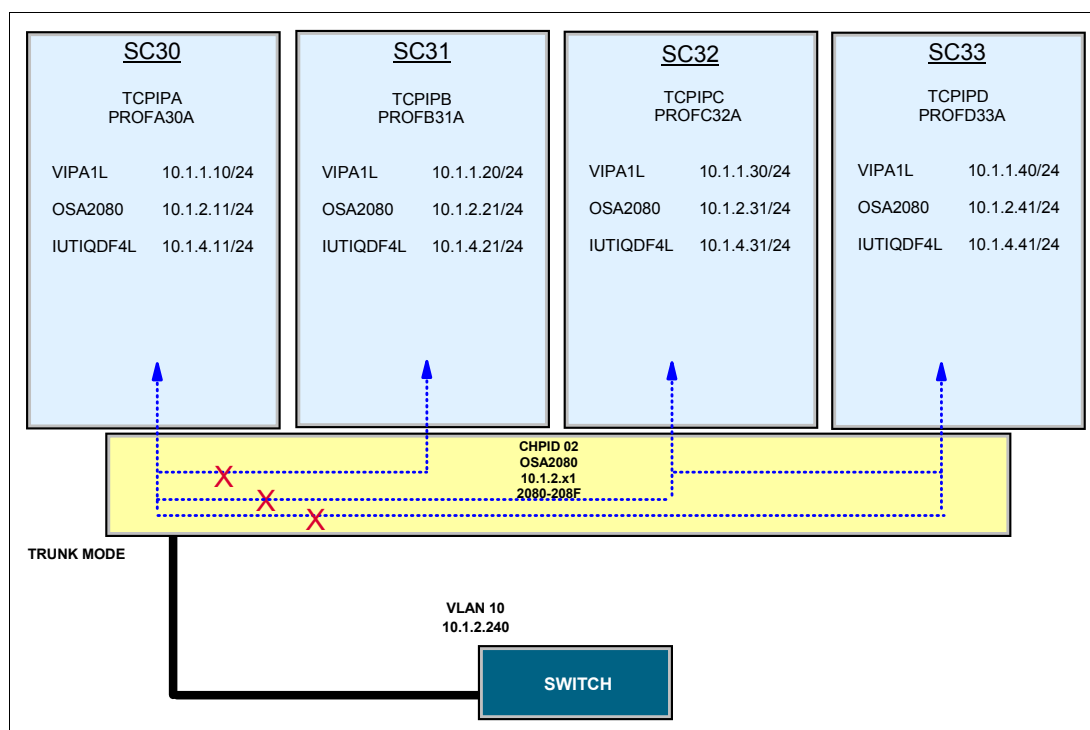


Figure 4-20 OSA2080 shared across four TCP/IP Images

Notice in Figure 4-20 the communication paths that we have indicated with an X. In our testing, we will not permit TCPIPA or TCPIPB to be reached directly over the shared OSA port. Using the ISOLATE function, we prevented direct communication between TCPIPA and TCPIPB by way of this port; we also prevented direct communication between either TCPIPA or TCPIPB and either of the two remaining stacks in our configuration: TCPIPC and TCPIPD.

We continued to permit TCPIPC and TCPIPD to share the OSA path between each other.

Note: You might choose to design your OSA ISOLATE function so that no sharing TCP/IP stack may use the direct path through the OSA. On the other hand, if you have abundant bandwidth on the OSA port, you might choose to implement ISOLATE on only selected sharing TCP/IP stacks, as we have done in our test.

Coding ISOLATE on the INTERFACE statements

The ISOLATE keyword can be coded only on an INTERFACE statement. The coding for TCPIPA and for TCPIPB is displayed at **1** and **2** in Example 4-20.

Example 4-20 ISOLATE coding on CHPID2 (OSA2080X) for PROFA30X and PROFB31X

```
INTERFACE OSA2080X
  DEFINE IPAQENET
  PORTNAME OSA2080
  IPADDR 10.1.2.11/24
  MTU 1492
  VLANID 10
  VMAC ROUTEALL
  ISOLATE 1
```

```
INTERFACE OSA2080X
  DEFINE IPAQENET
  PORTNAME OSA2080
  IPADDR 10.1.2.21/24
  MTU 1492
  VLANID 10
  VMAC ROUTEALL
  ISOLATE 2
```

The definitions for the interface in stacks TCPIPC and TCPIPD contain NOISOLATE, which is also the default. See **3** and **4** in Example 4-21.

Example 4-21 NOISOLATE coding on CHPID2 (OSA2080X) for PROFC32X and PROFD33X

```
INTERFACE OSA2080X
  DEFINE IPAQENET
  PORTNAME OSA2080
  IPADDR 10.1.2.31/24
  MTU 1492
  VLANID 10
  VMAC ROUTEALL
  NOISOLATE 3
```

```
INTERFACE OSA2080X
  DEFINE IPAQENET
  PORTNAME OSA2080
  IPADDR 10.1.2.41/24
  MTU 1492
  VLANID 10
  VMAC ROUTEALL
  NOISOLATE 4
```

Displaying the DEVICE to verify that ISOLATE is enabled

A display of all the INTERFACES on which we coded ISOLATE shows that ISOLATE is in force (**A**), as shown in Example 4-22 on page 166.

Example 4-22 ISOLATE coding on the Interface definition

```

INTFNAME: OSA2080X          INTFTYPE: IPAQENET  INTFSTATUS: READY
PORTNAME: OSA2080          DATAPATH: 2082      DATAPATHSTATUS: READY
SPEED: 0000001000
IPBROADCASTCAPABILITY: NO
VMACADDR: 020004749925    VMACORIGIN: OSA     VMACROUTER: ALL
ARPOFFLOAD: YES           ARPOFFLOADINFO: YES
CFGMTU: 1492              ACTMTU: 1492
IPADDR: 10.1.2.21/24
VLANID: 10                VLANPRIORITY: DISABLED
DYNVLANREGCFG: NO         DYNVLANREGCAP: YES
READSTORAGE: GLOBAL (4096K) INBPERF: BALANCED
CHECKSUMOFFLOAD: YES
SECCLASS: 255             MONSYSPLEX: NO
ISOLATE: YES A            OPTLATENCYMODE: NO

```

Viewing OSA/SF for ISOLATE enablement

The OSA/SF display shows where ISOLATE is enabled, as you can see in Example 4-23. Look for the entries marked with **X**.

Example 4-23 OSA/SF and ISOLATE

```

*****
*** OSA/SF Get OAT output created 18:28:38 on 11/02/2009          ***
*** IOACMD APAR level - OA26486                                   ***
*** Host   APAR level - OA27643                                   ***
*****
***                      Start of OSA address table for CHPID 02      ***
*****
* UA(Dev) Mode   Port   Entry specific information   Entry Valid
*****
                                Image 2.3 (A11 ) CULA 0
80(2080)* MPC      N/A    OSA2080 (QDIO control)          SIU ALL
82(2082) MPC 00 No4 No6 OSA2080 (QDIO data)   Isolated X SIU ALL
                                VLAN 10 (IPv4)

                                Group Address   Multicast Address
                                01005E000001    224.000.000.001
                                01005E000005    224.000.000.005

                                VMAC             IP address
                                HOME 020005749925 010.001.002.011
...
*****
                                Image 2.4 (A24 ) CULA 0
80(2080)* MPC      N/A    OSA2080 (QDIO control)          SIU ALL
82(2082) MPC 00 No4 No6 OSA2080 (QDIO data)   Isolated X SIU ALL
                                VLAN 10 (IPv4)

                                Group Address   Multicast Address
                                01005E000001    224.000.000.001
                                01005E000005    224.000.000.005

                                VMAC             IP address
                                HOME 020004749925 010.001.002.021
...
*****

```

Viewing routes with ISOLATE in effect

We recycled all four stacks that were sharing the OSA port named *OSA2080*. We are going to test connectivity between 10.1.2.11 at TCPIPA and 10.1.2.21, 31, and 33 at the other three stacks. We also want to test connectivity between 10.1.2.11 and the VIPAs in 10.1.1.0 at the other three stacks.

First, we examine the routing table at TCPIPA to determine if we have routes that will take us to those destinations, as shown in Example 4-24.

Example 4-24 View of routing at TCPIPA prior to adding static routes

D TCPIP,TCPIPA,N,ROUTE				
IPV4 DESTINATIONS				
DESTINATION	GATEWAY	FLAGS	REFCNT	INTERFACE
DEFAULT	10.1.2.240	UGO	0000000000	OSA2080I
DEFAULT	10.1.2.240	UGO	0000000000	OSA20A0I
10.1.1.10/32	0.0.0.0	UH	0000000000	VIPA1L
10.1.1.12/32	10.1.4.12	UGHO	0000000000	IUTIQDF4L
10.1.1.12/32	10.1.5.12	UGHO	0000000000	IUTIQDF5L
10.1.1.20/32	10.1.4.21	UGHO	0000000001	IUTIQDF4L
10.1.1.20/32	10.1.5.21	UGHO	0000000000	IUTIQDF5L
10.1.1.25/32	10.1.4.25	UGHO	0000000000	IUTIQDF4L
10.1.1.25/32	10.1.5.25	UGHO	0000000000	IUTIQDF5L
10.1.1.30/32	10.1.4.31	UGHO	0000000000	IUTIQDF4L
10.1.1.30/32	10.1.5.31	UGHO	0000000000	IUTIQDF5L
10.1.1.40/32	10.1.4.41	UGHO	0000000000	IUTIQDF4L
10.1.1.40/32	10.1.5.41	UGHO	0000000000	IUTIQDF5L
10.1.1.50/32	10.1.2.52	UGHO	0000000000	OSA2080I
10.1.1.50/32	10.1.2.51	UGHO	0000000000	OSA2080I
10.1.1.50/32	10.1.2.52	UGHO	0000000000	OSA20A0I
10.1.1.50/32	10.1.2.51	UGHO	0000000000	OSA20A0I
10.1.2.0/24	0.0.0.0	UO	0000000000	OSA20A0I
10.1.2.0/24	0.0.0.0	UO	0000000000	OSA2080I
10.1.2.10/32	0.0.0.0	UH	0000000000	VIPA2L
10.1.2.11/32	0.0.0.0	UH	0000000000	OSA2080I
10.1.2.12/32	0.0.0.0	UH	0000000000	OSA20A0I
10.1.2.14/32	0.0.0.0	H	0000000000	OSA2081I
10.1.2.17/32	10.1.4.12	UGHO	0000000000	IUTIQDF4L
10.1.2.17/32	10.1.5.12	UGHO	0000000000	IUTIQDF5L
10.1.2.20/32	10.1.4.21	UGHO	0000000000	IUTIQDF4L
10.1.2.20/32	10.1.5.21	UGHO	0000000000	IUTIQDF5L
10.1.2.25/32	10.1.4.25	UGHO	0000000000	IUTIQDF4L
10.1.2.25/32	10.1.5.25	UGHO	0000000000	IUTIQDF5L
10.1.2.30/32	10.1.4.31	UGHO	0000000000	IUTIQDF4L
10.1.2.30/32	10.1.5.31	UGHO	0000000000	IUTIQDF5L
10.1.3.0/24	0.0.0.0	UO	0000000000	OSA20E0I
10.1.3.11/32	0.0.0.0	UH	0000000000	OSA20C0I
10.1.3.12/32	0.0.0.0	UH	0000000000	OSA20E0I
10.1.4.0/24	0.0.0.0	UO	0000000000	IUTIQDF4L
10.1.4.11/32	0.0.0.0	UH	0000000000	IUTIQDF4L
10.1.5.0/24	0.0.0.0	UO	0000000000	IUTIQDF5L
10.1.5.11/32	0.0.0.0	UH	0000000000	IUTIQDF5L
10.1.6.0/24	10.1.4.41	UGO	0000000000	IUTIQDF4L
10.1.6.0/24	10.1.5.41	UGO	0000000000	IUTIQDF5L
10.1.6.11/32	0.0.0.0	UH	0000000000	IUTIQDF6L
10.1.7.0/24	0.0.0.0	US	0000000000	IQDIOLNK0A01070

B			
10.1.7.11/32	0.0.0.0	UH	0000000000 EZASAMEMVS
10.1.7.11/32	0.0.0.0	UH	0000000000 IQDIOLNK0A01070
B			
10.1.7.21/32	0.0.0.0	UHS	0000000000 IQDIOLNK0A01070
B			
10.1.7.25/32	0.0.0.0	UHS	0000000000 EZASAMEMVS
10.1.8.10/32	10.1.4.41	UGHO	0000000000 IUTIQDF4L
10.1.8.10/32	10.1.5.41	UGHO	0000000000 IUTIQDF5L
10.1.8.20/32	10.1.4.41	UGHO	0000000000 IUTIQDF4L
10.1.8.20/32	10.1.5.41	UGHO	0000000000 IUTIQDF5L
10.1.8.21/32	10.1.4.21	UGHO	0000000000 IUTIQDF4L
10.1.8.21/32	10.1.5.21	UGHO	0000000000 IUTIQDF5L
10.1.8.23/32	0.0.0.0	UH	0000000000 VIPL0A010817
10.1.8.28/32	10.1.4.12	UGHO	0000000000 IUTIQDF4L
10.1.8.28/32	10.1.5.12	UGHO	0000000000 IUTIQDF5L
10.1.8.30/32	10.1.4.31	UGHO	0000000000 IUTIQDF4L
10.1.8.30/32	10.1.4.41	UGHO	0000000000 IUTIQDF4L
10.1.8.30/32	10.1.5.31	UGHO	0000000000 IUTIQDF5L
10.1.8.30/32	10.1.5.41	UGHO	0000000000 IUTIQDF5L
10.1.8.40/32	10.1.4.41	UGHO	0000000000 IUTIQDF4L
10.1.8.40/32	10.1.5.41	UGHO	0000000000 IUTIQDF5L
10.1.8.41/32	10.1.4.41	UGHO	0000000000 IUTIQDF4L
10.1.8.41/32	10.1.5.41	UGHO	0000000000 IUTIQDF5L
10.1.8.42/32	10.1.4.41	UGHO	0000000000 IUTIQDF4L
10.1.8.42/32	10.1.5.41	UGHO	0000000000 IUTIQDF5L
10.1.8.43/32	10.1.4.41	UGHO	0000000000 IUTIQDF4L
10.1.8.43/32	10.1.5.41	UGHO	0000000000 IUTIQDF5L
10.1.8.44/32	10.1.4.41	UGHO	0000000000 IUTIQDF4L
10.1.8.44/32	10.1.5.41	UGHO	0000000000 IUTIQDF5L
10.1.8.50/32	10.1.2.52	UGHO	0000000000 OSA2080I
10.1.8.50/32	10.1.2.51	UGHO	0000000000 OSA2080I
10.1.8.50/32	10.1.2.52	UGHO	0000000000 OSA20A0I
10.1.8.50/32	10.1.2.51	UGHO	0000000000 OSA20A0I
10.1.30.10/32	0.0.0.0	UH	0000000000 VIPA3L
10.1.100.0/24	10.1.2.240	UGO	0000000002 OSA2080I
10.1.100.0/24	10.1.2.240	UGO	0000000001 OSA20A0I
127.0.0.1/32	0.0.0.0	UH	0000000005 LOOPBACK
192.168.1.0/24	10.1.2.240	UGO	0000000000 OSA2080I
192.168.1.0/24	10.1.2.240	UGO	0000000000 OSA20A0I
192.168.1.40/32	10.1.2.45	UGHO	0000000000 OSA2080I
192.168.1.40/32	10.1.2.45	UGHO	0000000000 OSA20A0I
192.168.2.0/24	10.1.2.240	UGO	0000000000 OSA2080I
192.168.2.0/24	10.1.2.240	UGO	0000000000 OSA20A0I
192.168.3.0/24	10.1.2.240	UGO	0000000000 OSA2080I
192.168.3.0/24	10.1.2.240	UGO	0000000000 OSA20A0I
IPV6 DESTINATIONS			
DESTIP: ::1/128			
GW: ::			
INTF: LOOPBACK6		REFCNT: 0000000000	
FLGS: UH		MTU: 65535	
86 OF 86 RECORDS DISPLAYED			
END OF THE REPORT			

Test 1: effect of ISOLATE with basic dynamic routing table

Our first test uses a limited routing table, where there is only one path available among the four TCP/IP stacks, and that path is blocked with the ISOLATE keyword at TCPIPA and TCPIPB. We execute all the following commands at TCPIPA, where ISOLATE is coded.

First, we use **tracert** against the three target addresses in network 10.1.2.0/24, as shown in Example 4-25.

Example 4-25 *tracert from TCPIPA to native OSA Home address of TCPIPB*

```
====> tracert 10.1.2.21 (tcp tcpipa V srcip 10.1.2.11 Intf OSA2080X
```

```
CS V1R12: Traceroute to 10.1.2.21 (10.1.2.21):
```

```
1 * * *
```

The results are the same when trying to reach TCPIPC and TCPIPD from either TCPIPA or TCPIPB: Because the route table indicates a direct path through the OSA, the stack attempts to send the packet over the direct route and experiences a failure. This is what we expected because we coded ISOLATE on OSA2080X in TCPIPA (and TCPIPB). Can we reach the VIPAs over the OSA port that is indicated as a route in Example 4-24 on page 167? We issue a **tracert** to the VIPAs and discover that the available routes will not allow us to reach them. See Example 4-26.

Example 4-26 *tracert from TCPIPA to VIPA address of TCPIPB*

```
====> tracert 10.1.1.20 (tcp tcpipa V srcip 10.1.1.10
```

```
CS V1R12: Traceroute to 10.1.1.20 (10.1.1.20):
```

```
1 * * *
```

The results are the same when trying to reach the VIPAs at TCPIPC and TCPIPD from either TCPIPA or TCPIPB: Because the route table indicates a direct path through the OSA, the stack attempts to send the packet over the direct route and experiences a failure. This is what we expected, because we coded ISOLATE on OSA2080X in TCPIPA (and TCPIPB).

ARP takeover with ISOLATE

As part of this test, we also activated a second interface on the same subnet in order to test the ARP takeover function. The second interface was also coded with ISOLATE on TCPIPA and TCPIPB. TCPIPC and TCPIPD were not defined with ISOLATE. When we deactivated the second interface, the address from the adapter port we were taking over moved to the remaining interface on that subnet, as you see at **1** in Example 4-27.

Example 4-27 *Effect of ARP takeover with ISOLATE*

```
/******  
/* OSA/SF Query created 14:31:35 on 10/12/2010 */  
/* IOACMD APAR level - 0A26486 */  
/* Host APAR level - 0A31645 */  
/******  
*** Start of OSA address table for CHPID 02 ***  
*****  
* UA(Dev) Mode Port Entry specific information Entry Valid  
*****  
Image 1.1 (A11 ) CULA 0  
00(2080)* MPC N/A OSA2080 (QDIO control) SIU ALL
```

```
02(2082) MPC 00 No4 No6 OSA2080 (QDIO data) Isolated SIU ALL
          VLAN 10 (IPv4)
```

	Group Address	Multicast Address	
	01005E000001	224.000.000.001	
	VMAC	IP address	
HOME	02004F776872	010.001.002.010	
HOME	02004F776872	010.001.002.011	Y

Note: The ARP takeover function still works as expected if we start a second device on the same subnet in the same stack. ISOLATE does not alter this function.

Test 2: effect of ISOLATE and NOISOLATE for multiple stacks

We next test to see if the shared OSA on CHPID2 will allow TCPIPC and TCPIP2D to communicate directly or not with each other and with TCPIPA and TCPIPB. We also test to confirm that the stack routes still allow us to reach TCP/IP stacks in the external network across the OSA ports, even if they have been coded with ISOLATE. Examine Figure 4-21.

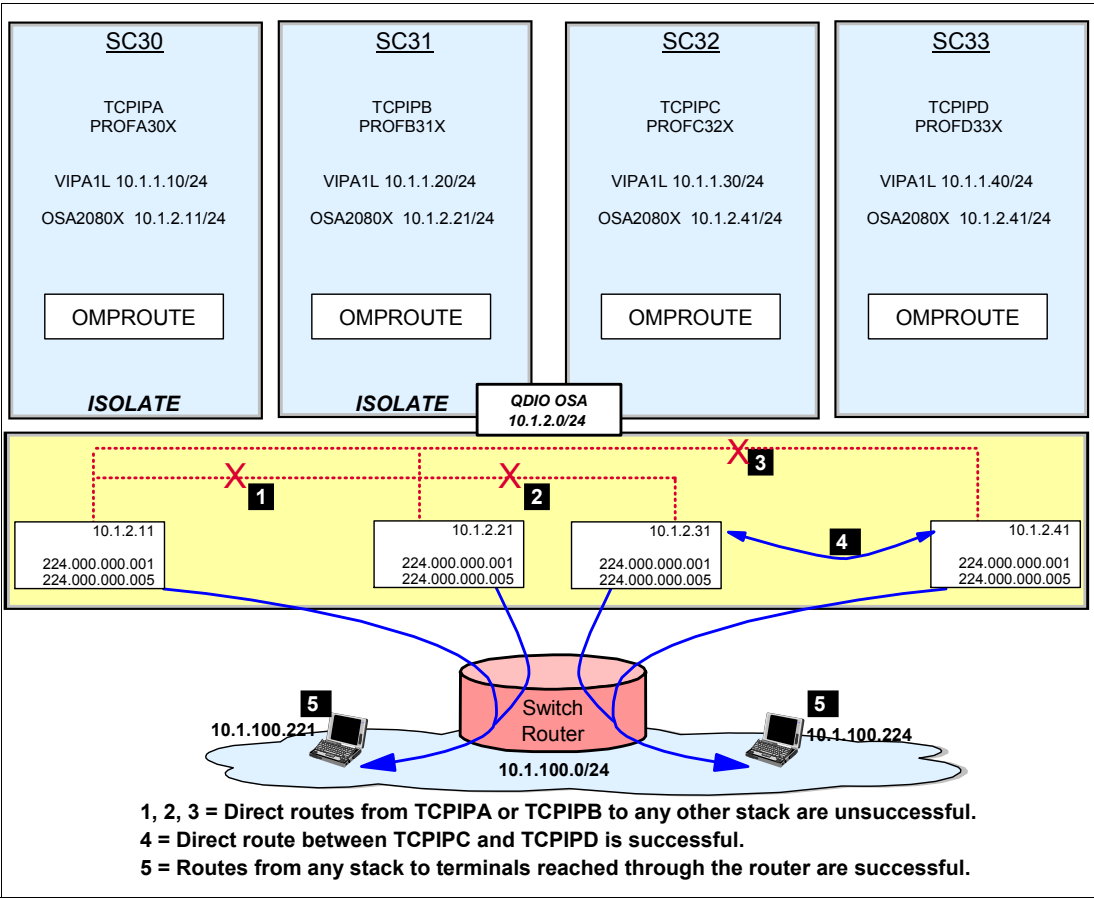


Figure 4-21 Available paths when ISOLATE has been defined and dynamic routing is enabled

Those tests show that the existing basic routing table at each of the stacks allows us to communicate with TCP/IP networks reached through the external router (5).

The routing tables also permit TCPIPC and TCIPD to communicate with each other (4);

Note: ISOLATE prevents only direct communication across the OSA port to any stack that has coded the ISOLATE keyword. However, it does *not* prevent communication between the stack defined with ISOLATE and any destinations beyond the OSA port. Therefore, we continue to be able to use the stack's routing table to TELNET or FTP from our workstations into any of the TCP/IP stacks, regardless of the coding of ISOLATE or NOISOLATE.

NOISOLATE is either coded or defaulted on the INTERFACE in the two stacks. However, TCPIPC and TCIPD cannot communicate at all with either TCPIPA or TCIPB, and TCPIPA and TCIPB cannot communicate with each other over the internal OSA path (1, 2, 3). In Example 4-28, we see the typical responses when a target cannot be reached.

Example 4-28 Unsuccessful attempts to reach TCIPB from TCPIPC

```
==> traceroute 10.1.2.21 (tcp tcpip V srcip 10.1.2.31
CS V1R12: Traceroute to 10.1.2.21 (10.1.2.21):
1 * * *
2 * * *
```

Unfortunately, the only path that TCPIPC and TCIPD for reaching TCPIPA and TCIPB is the direct route through the OSA port, but this port prevents internal routing because the parameter ISOLATE has been coded at TCPIPA and TCIPB. Examine the routing table in Example 4-24 on page 167, where you see that the table points to a network route for 10.1.2.0/24 that happens to be reached by way of a directly attached next-hop router (0.0.0.0):

10.1.2.0/24	0.0.0.0	U0	0000000000 OSA2080X
-------------	---------	----	---------------------

There is no route for any of the stacks to reach each other over the external router.

Again, the issue is that the dynamic routing table knows nothing of the ISOLATE feature because ISOLATE is not a Layer 3 function. The dynamic routing protocol is working per the protocol standards. So, how do we rectify this situation if we really want the stacks to communicate with each other, but just not directly over the OSAs? It will be a matter of adjusting the routing table by adding some non-replaceable static routes.

Recall that we earlier gave you options for dealing with this situation:

- Bypass the direct path by using Policy Based Routing (PBR).
- Bypass the direct path by coding static routes that supersede the routes learned by the dynamic routing protocol (see Figure 4-16 on page 160).
- Adjust the costs or weights of connections to favor alternate interfaces over the interfaces that have been coded with ISOLATE (see Figure 4-17 on page 161).

We tested only one of these options: coding static routes to supersede the dynamically learned routes.

Test 3: effect of ISOLATE: alternate path is present in routing table

In this test, we assume that we want to establish connectivity between TCPIPA and the other stacks on System z, but not over the direct OSA path. We also want to establish connectivity between TCPIPB and the other stacks, but not over the direct OSA path. In short, we want to leave our dynamic routing with OSPF in place, but we need to ensure that the learned routes over the direct OSA path are not always used. However, in doing so, we need to ensure that the direct paths between TCPIPC and TCPIPD stay in place. We accomplish all of this by adding non-replaceable static routes to the TCP/IP stacks, as shown in Figure 4-22.

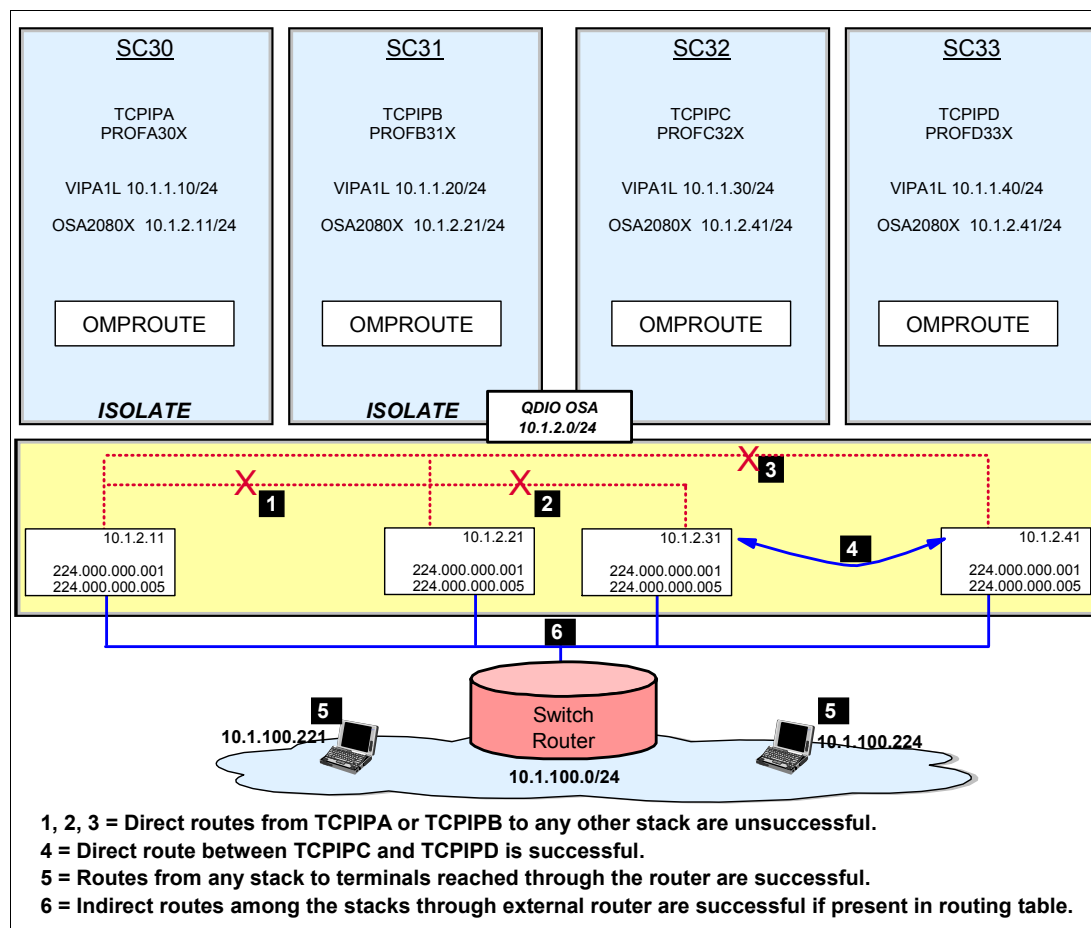


Figure 4-22 Making paths available through an external next-hop router

The routes to the remote TCP/IP nodes (5) through the OSA ports continue to be successful in our scenario; no changes are necessary here. The routing table between TCPIPC and TCPIPD continues to function as expected to permit direct routing between the two stacks (4); changes to the routing table are unnecessary here as well.

The routing paths indicated with 1, 2, and 3 in Figure 4-22 continue to be unsuccessful in this test because we want to enforce ISOLATE. However, we can make the two-hop paths through the external router (6) available if we code non-replaceable static routes. These routes will supersede the dynamically learned routes in the stack's routing table.

Coding non-replaceable static routes with the BEGINROUTES statement block

We add the static routing statements shown in Example 4-29 to TCPIPA. Note that TCPIPA's VIPA should not be present in the table.

Example 4-29 Static non-replaceable routes at TCPIPA to override direct route through OSA port

```
;TCPIPA.TCPPARMS(ROUTA30X)
BEGINRoutes
; Direct Routes - Routes that are directly connected to my interfaces
; Destination Subnet Mask First Hop Link Name Packet Size
;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;BELOW IS FOR TESTING ISOLATION;;;;;;;;;;;;;;;;;
ROUTE 10.1.2.0/24 10.1.2.240 OSA2080X mtu 1492 1
ROUTE 10.1.1.0/24 10.1.2.240 OSA2080X mtu 1492 2
ROUTE 10.1.1.20/32 10.1.2.240 OSA2080X mtu 1492 3
ROUTE 10.1.1.30/32 10.1.2.240 OSA2080X mtu 1492 3
ROUTE 10.1.1.40/32 10.1.2.240 OSA2080X mtu 1492 3
;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;ABOVE IS FOR TESTING ISOLATION;;;;;;;;;;;;;;;;;
ENDRoutes
```

In Example 4-29, we see, at **1** and **2**, the indirect route to both the native OSA port IP subnet and the VIPA IP subnet. In our scenario, these two statements do not suffice, because our OSPF configuration indicates that we are advertising HOST routes for the VIPAs. As a result, we also need the statements you see at **3**, that is, the statements that point to a route over the external router in order to reach the specific host VIPA addresses. If we do not code these statements, OSPF will advertise HOST routes and our stack will always try unsuccessfully to reach the target VIPAs over the OSA port.

We add the static routing statements shown in Example 4-30 to TCPIPB. The only difference to the statements at TCPIPA is the absence of TCPIPB's VIPA and the presence of TCPIPA's VIPA address.

Example 4-30 Static non-replaceable routes at TCPIPB to override direct route through OSA port

```
;TCPIPB.TCPPARMS(ROUTB31X)
BEGINRoutes
; Direct Routes - Routes that are directly connected to my interfaces
; Destination Subnet Mask First Hop Link Name Packet Size
;
;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;BELOW IS FOR TESTING ISOLATION;;;;;;;;;;;;;;;;;
ROUTE 10.1.2.0/24 10.1.2.240 OSA2080X mtu 1492
ROUTE 10.1.1.0/24 10.1.2.240 OSA2080X mtu 1492
ROUTE 10.1.1.10/32 10.1.2.240 OSA2080X mtu 1492
ROUTE 10.1.1.30/32 10.1.2.240 OSA2080X mtu 1492
ROUTE 10.1.1.40/32 10.1.2.240 OSA2080X mtu 1492
;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;ABOVE IS FOR TESTING ISOLATION;;;;;;;;;;;;;;;;;
ENDRoutes
```

We are going to test only a subset of all addresses available at the four stacks, that is, the connectivity with the VIPAs and the native OSA port addresses. Therefore, we have limited our BEGINROUTES coding only to these two address types.

Note: If you also need connectivity to other addresses, such as CTC or HiperSockets, then you might have to add more routes to your list of non-replaceable routes.

To simplify the test, we stop all interfaces but OSA2080X and take a snapshot of the current routing table (shown in Example 4-31).

Example 4-31 Routing table built by OMPROUTE (OSPF) at TCPIPA

```

D TCPIP,TCPIPA,N,ROUTE
IPV4 DESTINATIONS
DESTINATION      GATEWAY      FLAGS      REFCNT      INTERFACE
DEFAULT          10.1.2.240   UGO        0000000000  OSA2080X
10.1.1.10/32     0.0.0.0      UH         0000000000  VIPA1L
10.1.1.20/32     10.1.2.21    UGHO       0000000000  OSA2080X
10.1.1.30/32     10.1.2.31    UGHO       0000000000  OSA2080X
10.1.1.40/32     10.1.2.41    UGHO       0000000000  OSA2080X
10.1.2.0/24      0.0.0.0      UO         0000000000  OSA2080X
10.1.2.11/32     0.0.0.0      UH         0000000000  OSA2080X
10.1.2.12/32     0.0.0.0      H          0000000000  OSA20A0X
10.1.3.0/24      10.1.2.21    UGO        0000000000  OSA2080X
10.1.3.0/24      10.1.2.31    UGO        0000000000  OSA2080X
10.1.3.11/32     0.0.0.0      H          0000000000  OSA20C0X
10.1.3.12/32     0.0.0.0      H          0000000000  OSA20E0X
10.1.7.0/24      0.0.0.0      US         0000000000  IQDIOLNK0A01070B
10.1.7.11/32     0.0.0.0      H          0000000000  EZASAMEMVS
10.1.7.11/32     0.0.0.0      UH         0000000000  IQDIOLNK0A01070B
10.1.7.31/32     0.0.0.0      UHS        0000000000  IQDIOLNK0A01070B
10.1.7.41/32     0.0.0.0      UHS        0000000000  IQDIOLNK0A01070B
10.1.100.0/24    10.1.2.240   UGO        0000000000  OSA2080X
127.0.0.1/32    0.0.0.0      UH         0000000002  LOOPBACK
192.168.1.0/24   10.1.2.240   UGO        0000000000  OSA2080X
192.168.2.0/24   10.1.2.240   UGO        0000000000  OSA2080X
192.168.3.0/24   10.1.2.240   UGO        0000000000  OSA2080X
IPV6 DESTINATIONS
DESTIP:  ::1/128
  GW:    ::
  INTF:  LOOPBACK6      REFCNT:  0000000000
  FLGS:  UH             MTU:      65535
23 OF 23 RECORDS DISPLAYED
END OF THE REPORT

```

At **A** in Example 4-31, we see that OSPF reaches the VIPAs in subnet 10.1.1.0/24 over the OSA port. At **B**, we see that OSPF has informed the stack that the network 10.1.2.0/24 is directly attached.

We place OBEYFILE in the BEGINROUTES block. The new routing table is depicted in Example 4-32.

Example 4-32 Routing table at TCPIPA with static routes inserted

```

D TCPIP,TCPIPA,N,ROUTE
IPV4 DESTINATIONS
DESTINATION      GATEWAY      FLAGS      REFCNT      INTERFACE
DEFAULT          10.1.2.240   UGO        0000000000  OSA2080X
10.1.1.0/24      10.1.2.240   UGS        0000000000  OSA2080X
10.1.1.10/32     0.0.0.0      UH         0000000000  VIPA1L
10.1.1.20/32     10.1.2.240   UGHS       0000000000  OSA2080X
10.1.1.30/32     10.1.2.240   UGHS       0000000000  OSA2080X
10.1.1.40/32     10.1.2.240   UGHS       0000000000  OSA2080X
10.1.2.0/24      10.1.2.240   UGS        0000000000  OSA2080X
10.1.2.11/32     0.0.0.0      UH         0000000000  OSA2080X
10.1.2.12/32     0.0.0.0      H          0000000000  OSA20A0X
10.1.3.0/24      10.1.2.21    UGO        0000000000  OSA2080X
10.1.3.0/24      10.1.2.31    UGO        0000000000  OSA2080X
10.1.3.11/32     0.0.0.0      H          0000000000  OSA20C0X
10.1.3.12/32     0.0.0.0      H          0000000000  OSA20E0X
10.1.7.0/24      0.0.0.0      US         0000000000  IQDIOLNK0A01070B
10.1.7.11/32     0.0.0.0      H          0000000000  EZASAMEMVS
10.1.7.11/32     0.0.0.0      UH         0000000000  IQDIOLNK0A01070B
10.1.7.31/32     0.0.0.0      UHS        0000000000  IQDIOLNK0A01070B
10.1.7.41/32     0.0.0.0      UHS        0000000000  IQDIOLNK0A01070B
10.1.100.0/24    10.1.2.240   UGO        0000000000  OSA2080X
127.0.0.1/32    0.0.0.0      UH         0000000003  LOOPBACK
192.168.1.0/24   10.1.2.240   UGO        0000000000  OSA2080X
192.168.2.0/24   10.1.2.240   UGO        0000000000  OSA2080X
192.168.3.0/24   10.1.2.240   UGO        0000000000  OSA2080X
IPV6 DESTINATIONS
DESTIP:  ::1/128
  GW:    ::
  INTF:  LOOPBACK6      REFCNT:  0000000000
  FLGS:  UH              MTU:      65535
24 OF 24 RECORDS DISPLAYED
END OF THE REPORT

```

TCPIPB's routing table looks similar after the changes are made. Both tables now have HOST routes that point directly to the VIPAs and to the native OSA port addresses; however, the route statement now sends any packets destined for those addresses through the router with an IP address of 10.1.2.240. See the lines marked with **A** and **B** in Example 4-32.

We also need to make routing changes at TCPIPC. See the statements we have added to this stack in Example 4-33.

Example 4-33 TCPIPC: non-replaceable static routes to other TCP/IP nodes on System z

```

;TCPIPC.TCPPARMS(ROUTC32X)
BEGINRoutes
; Direct Routes - Routes that are directly connected to my interfaces
; Destination Subnet Mask First Hop Link Name Packet Size
;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;BELOW IS FOR TESTING ISOLATION;;;;;;;;;;;;;;;;;
ROUTE 10.1.2.11/32 10.1.2.240 OSA2080X mtu 1492 1
ROUTE 10.1.2.21/32 10.1.2.240 OSA2080X mtu 1492 1
ROUTE 10.1.1.10/32 10.1.2.240 OSA2080X mtu 1492 2
ROUTE 10.1.1.20/32 10.1.2.240 OSA2080X mtu 1492 2
;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;ABOVE IS FOR TESTING ISOLATION;;;;;;;;;;;;;;;;;
ENDRoutes

```

At TCPIPC and TCPIPD, we need to override the routes learned from OSPF that point to the addresses at TCPIPA and TCPIPB. In Example 4-33 at **1**, we have defined host routes to the native OSA port IP addresses at TCPIPA and TCPIPB that point to the external router. Note how we have not explicitly coded any static routes for the TCPIPD stack. At **2**, we add routes to the host VIPAs that are in TCPIPA and TCPIPB, but not in TCPIPD.

Of course, we need to make the same types of routing changes at TCPIPD. See the statements we have added to this stack in Example 4-34.

Example 4-34 TCPIPD: non-replaceable static routes to other TCP/IP nodes on System z

```

;TCPIPD.TCPPARMS(ROUTD33X)
BEGINRoutes
; Direct Routes - Routes that are directly connected to my interfaces
; Destination Subnet Mask First Hop Link Name Packet Size
;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;BELOW IS FOR TESTING ISOLATION;;;;;;;;;;;;;;;;;
ROUTE 10.1.2.11/32 10.1.2.240 OSA2080X mtu 1492 1
ROUTE 10.1.2.21/32 10.1.2.240 OSA2080X mtu 1492 1
ROUTE 10.1.1.10/32 10.1.2.240 OSA2080X mtu 1492 2
ROUTE 10.1.1.20/32 10.1.2.240 OSA2080X mtu 1492 2
;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;ABOVE IS FOR TESTING ISOLATION;;;;;;;;;;;;;;;;;
ENDRoutes

```

At **1**, we have defined host routes to the native OSA port IP addresses that point to the external router. Note how we have not explicitly coded any static routes for the TCPIPC stack. At **2**, we add routes to the host VIPAs that are in TCPIPA and TCPIPB, but not in TCPIPC.

In Example 4-35, we show you the new routing table structure at TCPIPC. (The routing table at TCPIPD resembles the one at TCPIPC, and so we do not illustrate it here.)

Example 4-35 Routing table at TCPIPC with entries provided by OSPF and by static routes

D TCPIP,TCPIPC,N,ROUTE					
IPV4 DESTINATIONS					
DESTINATION	GATEWAY	FLAGS	REFCNT	INTERFACE	
DEFAULT	10.1.2.240	UGO	0000000000	OSA2080X	
10.1.1.10/32	10.1.2.240	UGHS	0000000000	OSA2080X	A
10.1.1.20/32	10.1.2.240	UGHS	0000000000	OSA2080X	A
10.1.1.30/32	0.0.0.0	UH	0000000000	VIPA1L	
10.1.1.40/32	10.1.2.41	UGHO	0000000000	OSA2080X	B
10.1.2.0/24	0.0.0.0	UO	0000000000	OSA2080X	C
10.1.2.11/32	10.1.2.240	UGHS	0000000000	OSA2080X	D
10.1.2.21/32	10.1.2.240	UGHS	0000000000	OSA2080X	D
10.1.2.31/32	0.0.0.0	UH	0000000000	OSA2080X	
10.1.2.32/32	0.0.0.0	H	0000000000	OSA20A0X	
10.1.3.0/24	10.1.2.41	UGO	0000000000	OSA2080X	
10.1.3.0/24	10.1.2.240	UGO	0000000000	OSA2080X	
10.1.3.31/32	0.0.0.0	H	0000000000	OSA20C0X	
10.1.3.32/32	0.0.0.0	H	0000000000	OSA20E0X	
10.1.7.0/24	0.0.0.0	US	0000000000	IQDIOLNK0A01071F	
10.1.7.11/32	0.0.0.0	UHS	0000000000	IQDIOLNK0A01071F	
10.1.7.31/32	0.0.0.0	H	0000000000	EZASAMEMVS	
10.1.7.31/32	0.0.0.0	UH	0000000000	IQDIOLNK0A01071F	
10.1.7.41/32	0.0.0.0	UHS	0000000000	IQDIOLNK0A01071F	
10.1.100.0/24	10.1.2.240	UGO	0000000000	OSA2080X	
127.0.0.1/32	0.0.0.0	UH	0000000002	LOOPBACK	
192.168.1.0/24	10.1.2.240	UGO	0000000000	OSA2080X	
192.168.2.0/24	10.1.2.240	UGO	0000000000	OSA2080X	
192.168.3.0/24	10.1.2.240	UGO	0000000000	OSA2080X	
IPV6 DESTINATIONS					
DESTIP: ::1/128					
GW: ::					
INTF: LOOPBACK6		REFCNT: 0000000000			
FLGS: UH		MTU: 65535			
25 OF 25 RECORDS DISPLAYED					
END OF THE REPORT					

Look more closely at Example 4-35. The entries marked with **A** were statically added to override learned routes from OSPF. The entries at **B** and **C** remain as OSPF originally advertised them. These are for addresses in TCPIPD or for other 10.1.2.0/24 addresses that are not to be found in TCPIPA or TCPIPB. The entries marked with **D** were statically added to override learned routes from OSPF.

Testing with the non-replaceable static routes and OSPF

We use **tracert** to determine whether we are now taking a one-hop or two-hop route through the router. See the output in Example 4-36.

Example 4-36 tracert tests from TCPIPA to TCPIPB

TO NATIVE OSA PORT ADDRESS AT TCPIPB:

```
===> tracert 10.1.2.21 (tcp tcpipa V srcip 10.1.2.11

CS V1R12: Traceroute to 10.1.2.21 (10.1.2.21):
 1 router1 (10.1.2.240) 36 bytes to 10.1.2.11  1 ms  0 ms  0 ms A
 2 10.1.2.21 (10.1.2.21) 36 bytes to 10.1.2.11  0 ms  0 ms  0 ms
***
```

TO VIPA AT TCPIPB from NATIVE OSA PORT on TCPIPA:

```
===> tracert 10.1.1.20 (tcp tcpipa V srcip 10.1.2.11

CS V1R12: Traceroute to 10.1.1.20 (10.1.1.20):
 1 router1 (10.1.2.240) 36 bytes to 10.1.2.11  0 ms  0 ms  0 ms A
 2 WTSC31B (10.1.1.20) 36 bytes to 10.1.2.11  2 ms  0 ms  0 ms
***
```

TO VIPA AT TCPIPB from VIPA on TCPIPA:

```
===> tracert 10.1.1.20 (tcp tcpipa V srcip 10.1.1.10

CS V1R12: Traceroute to 10.1.1.20 (10.1.1.20):
 1 router1 (10.1.2.240) 36 bytes to 10.1.1.10  0 ms  0 ms  0 ms A
 2 WTSC31B (10.1.1.20) 36 bytes to 10.1.1.10  0 ms  0 ms  0 ms
***
```

As you can see in Example 4-36, our command executions are successful and point to a two-hop route across the router (A) between the two isolated TCPIP stacks (TCPIPA and TCPIPB).

Our tests to the external terminals from TCPIPA are also successful. (See Figure 4-22 on page 172 for a diagram of where the terminals reside.) Our test in Example 4-37 shows a verbose **ping** to the terminal at address 10.1.100.221:

Example 4-37 Connectivity through the ISOLATED OSA to the remote network

```
===> ping 10.1.100.221 (tcp tcpipa V srcip 10.1.1.10
CS V1R12: Pinging host 10.1.100.221
with 256 bytes of ICMP data
Ping #1 from 10.1.100.221: bytes=264 seq=1 ttl=127 time=1.28 ms
***
Ping #2 from 10.1.100.221: bytes=264 seq=2 ttl=127 time=0.37 ms
Ping #3 from 10.1.100.221: bytes=264 seq=3 ttl=127 time=0.91 ms
Ping statistics for 10.1.100.221
    Packets: Sent=3, Received=3, Lost=0 (0% loss)
    Approximate round trip times in milliseconds:
    Minimum=0.37 ms, Maximum=1.28 ms, Average=0.85 ms, StdDev=0.46 ms
***
```

Again, notice how our test is successful. The ISOLATE parameter did not inhibit us from reaching our external network over the OSA port.

We must now test our connectivity from TCPIPA to TCPIPC and TCPIPD to see if the two-hop route is successful now that we have updated the routing tables at all four stacks. See the indications of a two-hop route (2) in Example 4-38.

Example 4-38 traceroute tests from TCPIPA to TCPIPC

TO NATIVE OSA PORT ADDRESS AT TCPIPC from NATIVE OSA PORT at TCPIPA:

```
====> tracerte 10.1.2.31 (tcp tcpipa V srcip 10.1.2.11
```

```
CS V1R12: Traceroute to 10.1.2.31 (10.1.2.31):
 1 router1 (10.1.2.240) 36 bytes to 10.1.2.11 0 ms 0 ms 0 ms
 2 10.1.2.31 (10.1.2.31) 36 bytes to 10.1.2.11 0 ms 0 ms 0 ms 2
***
```

TO NATIVE OSA PORT ADDRESS AT TCPIPC from VIPA at TCPIPA:

```
====> tracerte 10.1.2.31 (tcp tcpipa V srcip 10.1.1.10
```

```
CS V1R12: Traceroute to 10.1.2.31 (10.1.2.31):
 1 router1 (10.1.2.240) 36 bytes to 10.1.1.10 0 ms 0 ms 0 ms
 2 10.1.2.31 (10.1.2.31) 36 bytes to 10.1.1.10 0 ms 0 ms 0 ms 2
***
```

TO VIPA AT TCPIPC from VIPA on TCPIPA:

```
====> tracerte 10.1.1.30 (tcp tcpipa V srcip 10.1.1.10
```

```
CS V1R12: Traceroute to 10.1.1.30 (10.1.1.30):
 1 router1 (10.1.2.240) 36 bytes to 10.1.1.10 0 ms 0 ms 0 ms
 2 WTSC32C (10.1.1.30) 36 bytes to 10.1.1.10 0 ms 0 ms 0 ms 2
***
```

Finally, we test the connectivity between TCPIPC and TCIPD to ensure that we are still taking the direct path through the OSA port despite the addition of our static routes. See Example 4-39.

Example 4-39 Static and dynamic routes at TCPIPC and TCIPD

TO NATIVE OSA PORT ADDRESS AT TCIPD from NATIVE OSA PORT at TCIPIC:

====> `tracerte 10.1.2.41 (tcp tcpipc srcip 10.1.2.31`

CS V1R12: Traceroute to 10.1.2.41 (10.1.2.41):

1 10.1.2.41 (10.1.2.41) 0 ms 0 ms 0 ms A

TO VIPA AT TCIPD from VIPA on TCIPIC:

====> `tracerte 10.1.1.40 (tcp tcpipc srcip 10.1.1.30`

CS V1R12: Traceroute to 10.1.1.40 (10.1.1.40):

1 10.1.1.40 (10.1.1.40) 0 ms 0 ms 0 ms A

You see in Example 4-39 that we are indeed taking the one-hop route (A).

4.5.6 Conclusions and recommendations: best practices for isolating traffic

Our experience shows that the ISOLATE function is not necessary in order to segregate traffic that is flowing across a shared OSA port. We have also seen that the ISOLATE function requires careful consideration, especially when you have implemented dynamic routing protocols in order to simplify the maintenance of valid routes in your network.

If you are using static routing protocols at z/OS and must isolate traffic over shared OSA ports, then either deploy a VLAN implementation with separate VLAN IDs assigned to separate IP subnets or exploit the ISOLATE feature and remember to disable ICMP redirects.

If you are using a dynamic routing protocol at z/OS and must isolate traffic over shared OSA ports, use a VLAN implementation with separate VLAN IDs assigned to separate IP subnets for each of the sharing TCP/IP stacks.

If you are using a dynamic routing protocol at z/OS and must isolate traffic over shared OSA ports but are reluctant to deploy VLANs in the System z TCP/IP stacks, use the OSA Connection Isolation feature. When doing so, plan a strategy to include some non-replaceable static routes in the TCP/IP stack's routing table that will force a hop over an external router. Create a robust testing plan to ensure that you are permitting only the type of routing that you desire.

4.6 HiperSockets connectivity

HiperSockets provides very fast TCP/IP communications between different logical partitions (LPARs) through the system memory of the System z server. The LPARs that are connected this way form an *internal* LAN, passing data between the LPARs at memory speeds, and thereby totally eliminating the I/O subsystem overhead and external network delays.

To create this scenario, we define the HiperSockets, which is represented by the IQD CHPID and its associated devices. All LPARs that are configured to use the shared IQD CHPID have internal connectivity, and therefore have the capability to communicate using HiperSockets.

In our environment we use three IQD CHPIDs (F4, F5, and F6). Each will create a separate logical LAN with its own subnetwork. Figure 4-23 depicts these interfaces to our scenario.

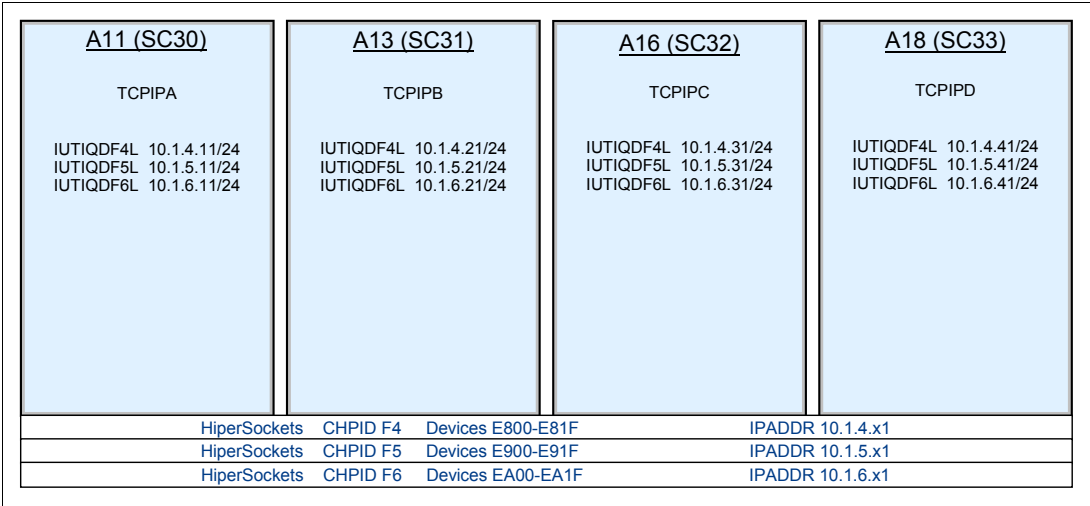


Figure 4-23 HiperSockets implementation scenario

4.6.1 Dependencies

The dependencies are:

- ▶ The HiperSockets must be defined as CHPID type IQD to the server using HCD or IOCP. This CHPID must be defined as shared to all LPARs that will be part of the HiperSockets internal LAN (see Example 4-1 on page 137).
- ▶ When explicitly defined, a correspondent TRLE must be created in VTAM using a port name IUTIQDxx, where xx is the CHPID number.
- ▶ When more than one IQD CHPID is configured to a specific LPAR, the VTAM start option IQDCHPID must be used to specify which specific IQD CHPID this LPAR should use.

Note: In both cases, the TRLE is dynamically built by VTAM. The IQDCHPID VTAM start option controls the VTAM selection of which IQD CHPID (and related devices) to include in the HiperSockets MPC group (IUTIQDIO) when it is dynamically built for DYNAMICXCF connectivity.

For additional details regarding how to configure a user-defined HiperSockets device or interface, refer to *z/OS Communications Server: IP Configuration Reference*, SC31-8776.

4.6.2 Considerations

For isolation of IP traffic between LPARs through HiperSockets, consider using VLANs, which means that you can logically subdivide the internal LAN for a HiperSockets CHPID into multiple virtual LANs. Therefore, stacks that configure the same VLAN ID for the same CHPID can communicate over that given HiperSockets, while stacks that have no VLAN ID or a different VLAN ID configured cannot.

For HiperSockets, the VLAN ID applies to IPv4 and IPv6 connections. HiperSockets VLAN IDs can be defined using the VLANID parameter on a LINK or INTERFACE statement. Valid VLAN IDs are in the range of 1 to 4094.

4.6.3 Configuring HiperSockets

The steps to implement HiperSockets are basically the same as with an OSA-Express interface. What changes is that there is no external configuration to be done, and the TRLE is created dynamically by VTAM.

The steps in the TCP/IP profile are as follows:

1. Create a DEVICE and LINK statements for each HiperSockets CHPID.
2. Create a HOME address to each defined LINK.
3. Define the characteristics of each LINK statement using BSDROUTINGPARMS.

Create a DEVICE and LINK statements for each HiperSockets CHPID

When defining an MPCIPA HiperSockets, use the DEVICE statement to specify the IQD CHPID hexadecimal value. The reserved device name prefix IUTIQDxx must be specified. The suffix xx indicates the hexadecimal value of the corresponding IQD CHPID that was configured with HCD or IOCP.

Define the device and link statements for each HiperSockets CHPID being implemented, as shown in Example 4-40. A HiperSockets CHPID must be defined as an MPCIPA type of device **1**.

The link definition describes the type of transport being used. A HiperSockets link is defined as IPAQIDIO **2**.

Example 4-40 HiperSockets device and link definitions

```
;HiperSockets definition. The TRLE is dynamically created on VTAMs
DEVICE IUTIQDF4 MPCIPA 1
LINK IUTIQDF4L IPAQIDIO 2 IUTIQDF4
DEVICE IUTIQDF5 MPCIPA 1
LINK IUTIQDF5L IPAQIDIO 2 IUTIQDF5
DEVICE IUTIQDF6 MPCIPA 1
LINK IUTIQDF6L IPAQIDIO 2 IUTIQDF6
```

Important: The hexadecimal value specified here represents the CHPID, and it cannot be the same value as that used for the dynamic XCF HiperSockets interface.

Create a HOME address to each defined LINK

Each link configured must have its own IP address. Our HiperSockets links are defined with the IP addresses, as shown in Example 4-41.

Example 4-41 HiperSockets HOME addresses

```
HOME
  10.1.4.11      IUTIQDF4L
  10.1.5.11      IUTIQDF5L
  10.1.6.11      IUTIQDF6L
```

Define the characteristics of each LINK statement using BSDROUTINGPARMS

To define the link characteristics, such as MTU size (1) and subnet mask (2), we used the BSDROUTINGPARMS statements (see Example 4-42). If not supplied, defaults will be used from static routing definitions in BEGINROUTES or the OMPROUTE configuration (dynamic routing definitions), if implemented.

If the link characteristics, BEGINROUTES statements, or the OMPROUTE configuration are not defined, then the stack's interface layer (based on hardware capabilities) and the characteristics of devices and links are used. This, however, might not provide the performance or function desired.

Example 4-42 BSDRoutingparms statements

```
BSDROUTINGPARMS TRUE
; Link name      MTU      Cost metric  Subnet Mask  Dest address
  VIPA1L 1492          0      255.255.255.252  0
  OSA2080L 1492        0      255.255.255.0    0
  OSA20A0L 1492        0      255.255.255.0    0
  OSA20C0L 1492        0      255.255.255.0    0
  OSA20E0L 1492        0      255.255.255.0    0
  IUTIQDF4L 8192 1      0      255.255.255.0 2      0
  IUTIQDF5L 8192        0      255.255.255.0    0
  IUTIQDF6L 8192        0      255.255.255.0    0
ENDBSDROUTINGPARMS
```

Note: Static and dynamic routing definitions will override or replace the link characteristics defined through the BSDROUTINGPARMS statements. Refer to Chapter 5, "Routing" on page 205 for more information about static and dynamic routing.

4.6.4 Verifying the connectivity status

In this section, we verify the status of all devices defined to the TCP/IP stack or VTAM.

Verifying the device status in TCP/IP stack

To verify the status of all devices being activated in the TCP/IP stack we use the NETSTAT command with the DEVLIST option, as shown in Example 4-43.

Example 4-43 Using command D TCPIP,TCPIPA,N,DEV to verify the HiperSockets connection

```
DEVNAME: IUTIQDF4          DEVTYPE: MPCIPA
DEVSTATUS: READY
LNKNAME: IUTIQDF4L        LNKTYPE: IPAQIDIO  LNKSTATUS: READY
IPBROADCASTCAPABILITY: NO
CFGROUTER: NON              ACTROUTER: NON
ARPOFFLOAD: YES             ARPOFFLOADINFO: YES
ACTMTU: 8192
VLANID: NONE
READSTORAGE: GLOBAL (2048K)
SECCLASS: 255               MONSYSPLEX: NO
IQDMULTIWRITE: ENABLED (ZIIP)
ROUTING PARAMETERS:
MTU SIZE: 8192              METRIC: 80
DESTADDR: 0.0.0.0           SUBNETMASK: 255.255.255.0
MULTICAST SPECIFIC:
MULTICAST CAPABILITY: YES
GROUP      REFCNT           SRCFLTMD
-----
224.0.0.5   0000000001      EXCLUDE
SRCADDR: NONE
224.0.0.1   0000000001      EXCLUDE
SRCADDR: NONE
LINK STATISTICS:
BYTESIN                      = 196650
INBOUND PACKETS              = 1647
INBOUND PACKETS IN ERROR     = 0
INBOUND PACKETS DISCARDED    = 0
INBOUND PACKETS WITH NO PROTOCOL = 0
BYTESOUT                     = 82841
OUTBOUND PACKETS             = 670
OUTBOUND PACKETS IN ERROR    = 0
OUTBOUND PACKETS DISCARDED   = 0
```

Displaying TCP/IP device resources in VTAM

The device drivers for TCP/IP are provided by VTAM. When CS for z/OS IP devices are activated, there must be an equivalent Transport Resource List Element (TRLE) defined to VTAM. The devices that are exclusively used by z/OS Communications Server IP have TRLEs that are automatically generated for them.

Because the device driver resources are provided by VTAM, you have the ability to display the resources using VTAM display commands.

For TRLEs that are generated dynamically, the device type and address can be decoded from the generated TRLE name. The format of the TRLE name is *IUTtaaaa*:

IUT Fixed for all TRLEs that are generated dynamically.

t Shows the device type, which indicates the following:

- C** Indicates this is a CDLC device.
- H** Indicates this is a HYPERCHANNEL device.
- I** Indicates this is a QDIO device.
- L** Indicates this is a LCS device.
- S** Indicates this is a SAMEHOST device.
- W** Indicates this is a CLAW device.
- X** Indicates this is a CTC device.

aaaa The read device number. For SAMEHOST connections, this is a sequence number.

To display a list of all TRLEs active in VTAM, use the D NET,TRL command, as shown in Example 4-44.

Example 4-44 D NET,TRL command output

```

D NET,TRL
IST350I DISPLAY TYPE = TRL 468
IST924I -----
IST1954I TRL MAJOR NODE = ISTTRL
IST1314I TRLE = IUTIQDF6 STATUS = ACTIV CONTROL = MPC
IST1314I TRLE = IUTIQDF5 STATUS = ACTIV CONTROL = MPC
IST1314I TRLE = IUTIQDF4 STATUS = ACTIV CONTROL = MPC
IST1314I TRLE = ISTT3033 STATUS = ACTIV CONTROL = XCF
IST1314I TRLE = ISTT3032 STATUS = ACTIV CONTROL = XCF
IST1314I TRLE = ISTT3031 STATUS = ACTIV CONTROL = XCF
IST1314I TRLE = IUTIQDIO STATUS = ACTIV CONTROL = MPC
IST1314I TRLE = IUTSAMEH STATUS = ACTIV CONTROL = MPC
IST1454I 8 TRLE(S) DISPLAYED

```

The D NET,TRL,TRLE command that is used to obtain information about a HiperSockets device is shown in Example 4-45.

Example 4-45 D NET,TRL,TRLE=IUTIQDF6

```

D NET,TRL,TRLE=IUTIQDF6
IST075I NAME = IUTIQDF6, TYPE = TRLE 512
IST1954I TRL MAJOR NODE = ISTTRL
IST486I STATUS= ACTIV, DESIRED STATE= ACTIV
IST087I TYPE = LEASED , CONTROL = MPC , HPDT = YES
IST1715I MPCLEVEL = QDIO MPCUSAGE = SHARE
IST2263I PORTNAME = PORTNUM = 0 OSA CODE LEVEL = 4752
IST2337I CHPID TYPE = IQD CHPID = F6
IST2319I IQD NETWORK ID = 0720
IST1577I HEADER SIZE = 4096 DATA SIZE = 16384 STORAGE = ***NA***
IST1221I WRITE DEV = EA01 STATUS = ACTIVE STATE = ONLINE
IST1577I HEADER SIZE = 4092 DATA SIZE = 0 STORAGE = ***NA***
IST1221I READ DEV = EA00 STATUS = ACTIVE STATE = ONLINE
IST924I -----
IST1221I DATA DEV = EA02 STATUS = ACTIVE STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST1717I ULPID = TCPIPA
IST2309I ACCELERATED ROUTING ENABLED
IST2331I QUEUE QUEUE READ

```

```

IST2332I ID      TYPE      STORAGE
IST2205I -----
IST2333I RD/1    PRIMARY   2.0M(126 SBALS)
IST2305I NUMBER OF DISCARDED INBOUND READ BUFFERS = 0
IST1757I PRIORITY1: UNCONGESTED PRIORITY2: UNCONGESTED
IST1757I PRIORITY3: UNCONGESTED PRIORITY4: UNCONGESTED
IST1801I UNITS OF WORK FOR NCB AT ADDRESS X'0F191010'
IST1802I P1 CURRENT = 0 AVERAGE = 1 MAXIMUM = 1
IST1802I P2 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P3 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P4 CURRENT = 0 AVERAGE = 1 MAXIMUM = 2
IST924I -----
IST1221I DATA DEV = EA03 STATUS = ACTIVE      STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST1717I ULPID = TCPIPC
IST2309I ACCELERATED ROUTING ENABLED
IST2331I QUEUE   QUEUE   READ
IST2332I ID      TYPE      STORAGE
IST2205I -----
IST2333I RD/1    PRIMARY   2.0M(126 SBALS)
IST2305I NUMBER OF DISCARDED INBOUND READ BUFFERS = 0
IST1757I PRIORITY1: UNCONGESTED PRIORITY2: UNCONGESTED
IST1757I PRIORITY3: UNCONGESTED PRIORITY4: UNCONGESTED
IST1801I UNITS OF WORK FOR NCB AT ADDRESS X'0FC87010'
IST1802I P1 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P2 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P3 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P4 CURRENT = 0 AVERAGE = 1 MAXIMUM = 2
IST924I -----
IST1221I DATA DEV = EA04 STATUS = ACTIVE      STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST1717I ULPID = TCPIPB
IST2310I ACCELERATED ROUTING DISABLED
IST2331I QUEUE   QUEUE   READ
IST2332I ID      TYPE      STORAGE
IST2205I -----
IST2333I RD/1    PRIMARY   2.0M(126 SBALS)
IST2305I NUMBER OF DISCARDED INBOUND READ BUFFERS = 0
IST1757I PRIORITY1: UNCONGESTED PRIORITY2: UNCONGESTED
IST1757I PRIORITY3: UNCONGESTED PRIORITY4: UNCONGESTED
IST1801I UNITS OF WORK FOR NCB AT ADDRESS X'0F530010'
IST1802I P1 CURRENT = 0 AVERAGE = 1 MAXIMUM = 1
IST1802I P2 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P3 CURRENT = 0 AVERAGE = 1 MAXIMUM = 1
IST1802I P4 CURRENT = 0 AVERAGE = 1 MAXIMUM = 2
IST924I -----
IST1221I DATA DEV = EA05 STATUS = RESET      STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST924I -----
IST1221I DATA DEV = EA06 STATUS = RESET      STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST924I -----
IST1221I DATA DEV = EA07 STATUS = RESET      STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST924I -----

```

```
IST1221I DATA DEV = EA08 STATUS = RESET      STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST924I -----
IST1221I DATA DEV = EA09 STATUS = RESET      STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST924I -----
IST314I END
```

4.7 Dynamic XCF connectivity

The last connectivity scenario that we add to our environment connects all images within the same sysplex environment through a dynamic XCF connection that is created by the DYNAMICXCF definition in the TCP/IP profile.

After being defined, DYNAMICXCF provides connectivity between stacks under the same LPAR by using the IUTSAMEH device (SAMEHOST) and between LPARs through HiperSockets using a IUTiQDIO device. To connect other z/OS images or other servers, an XCF Coupling Facility link is created.

In our scenario, we use DYNAMICXCF through HiperSockets with IQD CHPID F7. So, by defining the DYNAMICXCF statement, we create the XCF subnetwork through HiperSockets, as shown in Figure 4-24.

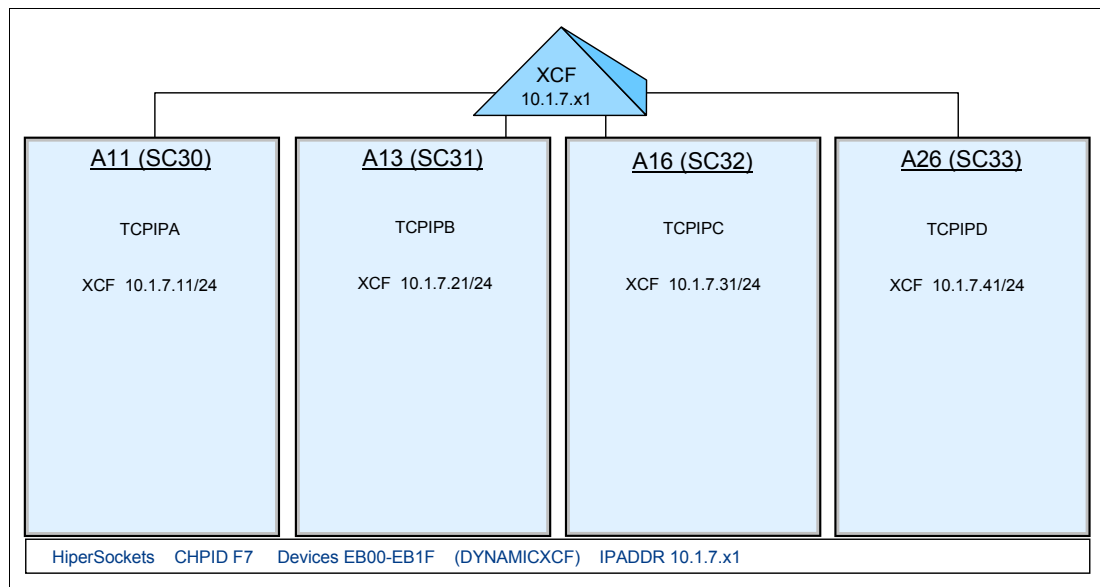


Figure 4-24 Dynamic XCF environment

4.7.1 Dependencies

The dependencies are as follows:

- ▶ All z/OS hosts must belong to the same sysplex.
- ▶ VTAM must have XCF communications enabled by specifying XCFINIT=YES or XCFINIT=DEFINE as a startup parameter or by activating the VTAM XCF local SNA major node, ISTLSXCF. For details about configuration, refer to *z/OS Communications Server: SNA Network Implementation*, SC31-8777.
- ▶ DYNAMICXCF must be specified in the TCP/IP profile of each stack.
- ▶ The IQD CHPID being used for the DYNAMICXCF device cannot be the user-defined HiperSockets device (IQD CHPID). To avoid this, a VTAM start option, IQDCHPID, can be used to identify which IQD CHPID will be used by DYNAMICXCF.

4.7.2 Considerations

z/OS Communications Server has improved and optimized Sysplex IP routing. In a sysplex environment, you might prefer to use a connection other than a Coupling Facility link for cross-server connectivity because XCF is heavily used by other workloads (in particular, for distributed application data sharing).

This option can be configured with the VIPAROUTE statement in the VIPADYNAMIC statement. It allows for the use of OSA-Express features such as 1000BASE-T Ethernet, Gigabit Ethernet, and 10 Gigabit Ethernet. For details, refer to *IBM z/OS V1R12 Communications Server TCP/IP Implementation Volume 3: High Availability*, SG24-7898.

z/OS Communications Server supports sysplex subplexing in conjunction with HiperSockets and DYNAMICXCF. Refer to Chapter 7, “Sysplex subplexing” on page 283, for details about restricting data traffic flow among certain TCP/IP stacks in a sysplex environment.

4.7.3 Configuring DYNAMICXCF

To implement XCF connections in our environment, you can use three different types of devices:

- ▶ DynamicXCF HiperSockets device (IUTIQDIO) for connections between z/OS LPARs within the same server
- ▶ DynamicXCF SAMEHOST device (IUTSAMEH) for stacks within the same LPAR
- ▶ VTAM dynamically created ISTLSXCF to connect z/OS LPARs in other servers within the same sysplex

Figure 4-24 on page 188 shows the DynamicXCF implementation in our environment using HiperSockets CHPID F7.

When you use dynamic XCF for sysplex configuration, make sure that XCFINIT=YES or XCFINIT=DEFINE is coded in the VTAM start options.

If XCFINIT=NO was specified, issue the VARY ACTIVATE command for the ISTLSXCF major node. This ensures that XCF connections between TCP stacks on different VTAM nodes in the sysplex can be established.

The VTAM ISTLSXCF major node must be active for DYNAMICXCF work, except for the following scenarios:

- ▶ Multiple TCP/IP stacks on the same LPAR; a dynamic SAMEHOST definition is generated whether or not ISTLSXCF is active.
- ▶ HiperSockets is configured and enabled across multiple z/OS LPARs that are in the same sysplex and the same server. If this is the case, a dynamic IUTIQDIO link is created whether ISTLSXCF is active or not.

To implement DYNAMICXCF in our environment, we coded the IPCONFIG definitions in the TCP/IP profile, as shown in Example 4-46. To control the IP subnetwork used to connect all z/OS images, we define the XCF IP address, the IP mask, and the link cost in the DYNAMICXCF statement **1**.

Example 4-46 IPCONFIG DYNAMICXCF configuration

```
IPCONFIG DATAGRAMFWD SYSPLEXROUTING IPSECURITY
DYNAMICXCF 10.1.7.11 255.255.255.0 1 1
```

4.7.4 Verifying connectivity status

In this section, we verify the status of all devices that are defined to the TCP/IP stack or VTAM.

Verifying the device status in the TCP/IP stack

To verify the status of all devices being activated in the TCP/IP stack, use the NETSTAT command with the DEVLIST option, as shown in Example 4-47.

Example 4-47 Using command D TCPIP,TCPIPA,N,DEV to verify the device status

```
D TCPIP,TCPIPA,N,DEV
..... Lines
deleted
DEVNAME: TUTSAMEH          DEVTYPE: MPCPTP
  DEVSTATUS: READY
  LNKNNAME: EZASAMEMVS      LNKNTYPE: MPCPTP      LNKNSTATUS: READY
    ACTMTU: 65535
    SECCCLASS: 255
  ROUTING PARAMETERS:
    MTU SIZE: 65535          METRIC: 00
    DESTADDR: 0.0.0.0        SUBNETMASK: 255.255.255.0
  MULTICAST SPECIFIC:
    MULTICAST CAPABILITY: YES
    GROUP          REFCNT          SRCFLTMD
    -----
    224.0.0.1      0000000001      EXCLUDE
    SRCADDR: NONE
  LINK STATISTICS:
    BYTESIN                      = 0
    INBOUND PACKETS              = 0
    INBOUND PACKETS IN ERROR     = 0
    INBOUND PACKETS DISCARDED    = 0
    INBOUND PACKETS WITH NO PROTOCOL = 0
    BYTESOUT                     = 96
    OUTBOUND PACKETS            = 4
    OUTBOUND PACKETS IN ERROR    = 0
    OUTBOUND PACKETS DISCARDED   = 0
DEVNAME: TUTIQDIO          DEVTYPE: MPCIPA
  DEVSTATUS: READY
  LNKNNAME: IQDIOLNKOA01070B LNKNTYPE: IPAQIDIO      LNKNSTATUS: READY
    IPBROADCASTCAPABILITY: NO
    CFGROUTER: NON              ACTROUTER: NON
    ARPOFFLOAD: YES             ARPOFFLOADINFO: YES
    ACTMTU: 8192
```

```

VLANID: 21
READSTORAGE: GLOBAL (2048K)
SECCLASS: 255
IQDMULTIWRITE: ENABLED (ZIIP)
ROUTING PARAMETERS:
  MTU SIZE: 65535          METRIC: 00
  DESTADDR: 0.0.0.0       SUBNETMASK: 255.255.255.0
MULTICAST SPECIFIC:
  MULTICAST CAPABILITY: YES
  GROUP              REFCNT          SRCFLTMD
  -----
  224.0.0.1          0000000001      EXCLUDE
  SRCADDR: NONE
LINK STATISTICS:
  BYTESIN                      = 0
  INBOUND PACKETS              = 0
  INBOUND PACKETS IN ERROR     = 0
  INBOUND PACKETS DISCARDED    = 0
  INBOUND PACKETS WITH NO PROTOCOL = 0
  BYTESOUT                    = 0
  OUTBOUND PACKETS            = 0
  OUTBOUND PACKETS IN ERROR   = 0
  OUTBOUND PACKETS DISCARDED  = 0

```

Note: The link name for device IUTIQDIO is defined dynamically as IQDIOLNK0A01070B. In the link name, 0A01070B is the hexadecimal value of the assigned IP address (10.1.7.11).

Displaying TCP/IP device resources in VTAM

The device drivers for TCP/IP are provided by VTAM. When CS for z/OS IP devices are activated, there must be an equivalent Transport Resource List Element (TRLE) defined to VTAM. The devices that are exclusively used by z/OS Communications Server IP have TRLEs that are automatically generated for them.

Because the device driver resources are provided by VTAM, you have the ability to display the resources using VTAM display commands.

For TRLEs that are generated dynamically, the device type and address can be decoded from the generated TRLE name. The format of the TRLE name is *IUTtaaaa*:

IUT Fixed for all TRLEs that are generated dynamically.

t Shows the device type, which indicates the following:

- C** Indicates this is a CDLC device.
- H** Indicates this is a HYPERCHANNEL device.
- I** Indicates this is a QDIO device.
- L** Indicates this is a LCS device.
- S** Indicates this is a SAMEHOST device.
- W** Indicates this is a CLAW device.
- X** Indicates this is a CTC device.

aaaa The read device number. For SAMEHOST connections, this is a sequence number.

For XCF links, the format of the TRLE name is ISTTxyy. ISTT is fixed, xx is the SYSCONE value of the originating VTAM, and yy is the SYSCONE value of the destination VTAM.

To display a list of all TRLEs active in VTAM use the D NET,TRL command, as shown in Example 4-48.

Example 4-48 D NET,TRL command output

```

D NET,TRL
IST097I DISPLAY ACCEPTED
IST350I DISPLAY TYPE = TRL 605
IST924I -----
IST1954I TRL MAJOR NODE = ISTTRL
IST1314I TRLE = ISTT3032 STATUS = ACTIV          CONTROL = XCF
IST1314I TRLE = ISTT3031 STATUS = ACTIV          CONTROL = XCF
IST1314I TRLE = ISTT3033 STATUS = ACTIV          CONTROL = XCF
IST1314I TRLE = IUTIQDF6 STATUS = ACTIV          CONTROL = MPC
IST1314I TRLE = IUTIQDF5 STATUS = ACTIV          CONTROL = MPC
IST1314I TRLE = IUTIQDF4 STATUS = ACTIV          CONTROL = MPC
IST1314I TRLE = IUTIQDIO STATUS = ACTIV          CONTROL = MPC
IST1314I TRLE = IUTSAMEH STATUS = ACTIV          CONTROL = MPC
IST1454I 8 TRLE(S) DISPLAYED
IST924I -----
IST314I END

```

You can display information of TRLEs grouped by control type, such as MPC or XCF devices, as shown in Example 4-49.

Example 4-49 D NET,TRL,CONTROL=XCF

```

D NET,TRL,CONTROL=XCF
IST350I DISPLAY TYPE = TRL 911
IST924I -----
IST1954I TRL MAJOR NODE = ISTTRL
IST1314I TRLE = ISTT3033 STATUS = ACTIV          CONTROL = XCF
IST1314I TRLE = ISTT3032 STATUS = ACTIV          CONTROL = XCF
IST1314I TRLE = ISTT3031 STATUS = ACTIV          CONTROL = XCF
IST1454I 3 TRLE(S) DISPLAYED
IST1454I 3 TRLE(S) DISPLAYED

```

You can also display XCF TRLE-specific information, as shown in Example 4-50.

Example 4-50 D NET,TRL,TRLE=ISTT3031

D NET,TRL,TRLE=ISTT3031

```
IST075I NAME = ISTT3031, TYPE = TRLE 925
IST1954I TRL MAJOR NODE = ISTTRL
IST486I STATUS= ACTIV, DESIRED STATE= ACTIV
IST087I TYPE = LEASED           , CONTROL = XCF , HPDT = *NA*
IST1715I MPCLEVEL = HPDT        MPCUSAGE = SHARE
IST1717I ULPID = ISTP3031
IST1503I XCF TOKEN = 02000058001F0002    STATUS = ACTIVE
IST1502I ADJACENT CP = USIBMSC.SC31M
IST314I END
```

The DYNAMICXCF configuration created a HiperSockets TRLE named IUTIQDIO. The related TRLE status can also be displayed, as shown in Example 4-51.

Example 4-51 D NET,TRL,TRLE=IUTIQDIO

D NET,TRL,TRLE=IUTIQDIO

```
IST075I NAME = IUTIQDIO, TYPE = TRLE 933
IST1954I TRL MAJOR NODE = ISTTRL
IST486I STATUS= ACTIV, DESIRED STATE= ACTIV
IST087I TYPE = LEASED           , CONTROL = MPC , HPDT = YES
IST1715I MPCLEVEL = QDIO        MPCUSAGE = SHARE
IST2263I PORTNAME = IUTIQDF3    PORTNUM = 0    OSA CODE LEVEL = 4752
IST2337I CHPID TYPE = IQD       CHPID = F3
IST2319I IQD NETWORK ID = 0713
IST1577I HEADER SIZE = 4096 DATA SIZE = 16384 STORAGE = ***NA***
IST1221I WRITE DEV = 7301 STATUS = ACTIVE    STATE = ONLINE
IST1577I HEADER SIZE = 4092 DATA SIZE = 0 STORAGE = ***NA***
IST1221I READ DEV = 7300 STATUS = ACTIVE     STATE = ONLINE
IST924I -----
IST1221I DATA DEV = 7302 STATUS = ACTIVE     STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST1717I ULPID = TCPIPA
IST2309I ACCELERATED ROUTING ENABLED
IST2331I QUEUE    QUEUE    READ
IST2332I ID       TYPE      STORAGE
IST2205I -----
IST2333I RD/1     PRIMARY   2.0M(126 SBALS)
IST2305I NUMBER OF DISCARDED INBOUND READ BUFFERS = 0
IST1757I PRIORITY1: UNCONGESTED PRIORITY2: UNCONGESTED
IST1757I PRIORITY3: UNCONGESTED PRIORITY4: UNCONGESTED
IST1801I UNITS OF WORK FOR NCB AT ADDRESS X'0F30D010'
IST1802I P1 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P2 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P3 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P4 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST924I -----
IST1221I DATA DEV = 7303 STATUS = ACTIVE     STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST1717I ULPID = TCPIPC
IST2309I ACCELERATED ROUTING ENABLED
IST2331I QUEUE    QUEUE    READ
IST2332I ID       TYPE      STORAGE
IST2205I -----
```

```

IST2333I RD/1    PRIMARY    2.0M(126 SBALS)
IST2305I NUMBER OF DISCARDED INBOUND READ BUFFERS = 0
IST1757I PRIORITY1: UNCONGESTED PRIORITY2: UNCONGESTED
IST1757I PRIORITY3: UNCONGESTED PRIORITY4: UNCONGESTED
IST1801I UNITS OF WORK FOR NCB AT ADDRESS X'0F731010'
IST1802I P1 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P2 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P3 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P4 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST924I -----
IST1221I DATA DEV = 7304 STATUS = ACTIVE      STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST1717I ULPID = TCPIPD
IST2310I ACCELERATED ROUTING DISABLED
IST2331I QUEUE    QUEUE    READ
IST2332I ID      TYPE      STORAGE
IST2205I -----
IST2333I RD/1    PRIMARY    2.0M(126 SBALS)
IST2305I NUMBER OF DISCARDED INBOUND READ BUFFERS = 0
IST1757I PRIORITY1: UNCONGESTED PRIORITY2: UNCONGESTED
IST1757I PRIORITY3: UNCONGESTED PRIORITY4: UNCONGESTED
IST1801I UNITS OF WORK FOR NCB AT ADDRESS X'0CDDA010'
IST1802I P1 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P2 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P3 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P4 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST924I -----
IST1221I DATA DEV = 7305 STATUS = ACTIVE      STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST1717I ULPID = TCPIPB
IST2310I ACCELERATED ROUTING DISABLED
IST2331I QUEUE    QUEUE    READ
IST2332I ID      TYPE      STORAGE
IST2205I -----
IST2333I RD/1    PRIMARY    2.0M(126 SBALS)
IST2305I NUMBER OF DISCARDED INBOUND READ BUFFERS = 0
IST1757I PRIORITY1: UNCONGESTED PRIORITY2: UNCONGESTED
IST1757I PRIORITY3: UNCONGESTED PRIORITY4: UNCONGESTED
IST1801I UNITS OF WORK FOR NCB AT ADDRESS X'0F529010'
IST1802I P1 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P2 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P3 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P4 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST924I -----
IST1221I DATA DEV = 7306 STATUS = RESET      STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST924I -----
IST1221I DATA DEV = 7307 STATUS = RESET      STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST924I -----
IST1221I DATA DEV = 7308 STATUS = RESET      STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST924I -----
IST1221I DATA DEV = 7309 STATUS = RESET      STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*

```

```
IST924I -----  
IST314I END
```

The DYNAMICXCF configuration created a SAMEHOST TRLE named IUTSAMEH. The related TRLE status can be displayed, as shown in Example 4-52.

Example 4-52 D NET,TRL,TRLE=IUTSAMEH

```
D NET,TRL,TRLE=IUTSAMEH  
IST075I NAME = IUTSAMEH, TYPE = TRLE 970  
IST1954I TRL MAJOR NODE = ISTTRL  
IST486I STATUS= ACTIV, DESIRED STATE= ACTIV  
IST087I TYPE = LEASED , CONTROL = MPC , HPDT = YES  
IST1715I MPCLEVEL = HPDT MPCUSAGE = SHARE  
IST1717I ULPID = TCPIPB  
IST1717I ULPID = TCPIPA  
IST1717I ULPID = TCPIPD  
IST1717I ULPID = TCPIPC  
IST314I END
```

The DYNAMICXCF statement dynamically generates the DEVICE, LINK, and HOME statements. It also starts the device when the TCP/IP stack is activated, as we can see in the messages shown in Example 4-53.

Example 4-53 DYNAMICXCF messages

```
$HASP373 TCPIPA STARTED  
  
EZZ0350I SYSPLEX ROUTING SUPPORT IS ENABLED  
EZZ0624I DYNAMIC XCF DEFINITIONS ARE ENABLED  
EZD1176I TCPIPA HAS SUCCESSFULLY JOINED THE TCP/IP SYSPLEX GROUP EZBTCPCS  
EZZ4324I CONNECTION TO 10.1.7.51 ACTIVE FOR DEVICE IUTSAMEH 1  
EZZ4313I INITIALIZATION COMPLETE FOR DEVICE IUTSAMEH  
EZZ4324I CONNECTION TO 10.1.7.31 ACTIVE FOR DEVICE IUTSAMEH 1  
EZZ4313I INITIALIZATION COMPLETE FOR DEVICE IUTIQDIO 2
```

In this example, the numbers correspond to the following information:

- 1.** Indicates that the TCPIPA stack has been connected to the other stacks through XCF using a SAMEHOST device.
- 2.** Indicates that XCF will also use HiperSockets to connect other TCP/IP stacks within the same server, using a IUTIQDIO device.

4.8 Controlling and activating devices

After all required connectivity definitions are defined in the TCP/IP profile and the stack is started, you have the option to start and stop devices, as well as activate modified device definitions. In this section we show the commands used to perform these tasks.

4.8.1 Starting a device

A device can be started by any of the following methods:

- ▶ Defining the START statement in the TCP/IP profile, as shown in Example 4-54.

Example 4-54 Start statements in TCP/IP profile

```
START OSA2080
START OSA20C0
START OSA20E0
START OSA20A0
START IUTIQDF4
START IUTIQDF5
START IUTIQDF6
```

- ▶ Using the z/OS console command **VARY TCPIP,tcpipproc,start,devicename**.
- ▶ Creating a file with a start statement and using the z/OS console command **Vary TCPIP,tcpipproc,OBEYFILE,datasetname**. The file defined by the file name has the START statement to activate the desired device or devices.

Using any of the starting methods will result in a series of messages, as shown in Example 4-55.

Example 4-55 Starting a TCP/IP device

```
V TCPIP,TCPIPA,START,OSA2080
EZZ0060I PROCESSING COMMAND: VARY TCPIP,TCPIPA,START,OSA2080
EZZ0053I COMMAND VARY START COMPLETED SUCCESSFULLY
EZZ4313I INITIALIZATION COMPLETE FOR DEVICE OSA2080
```

4.8.2 Stopping a device

You can stop a device using any of the following methods:

- ▶ Using the z/OS console command **Vary TCPIP,tcpipproc,STOP,devicename**.
- ▶ Creating a file with the stop statement to the desired device or devices and using the z/OS console command **Vary TCPIP,tcpipproc,OBEYFILE,datasetname**.

When you stop a device, you see messages as shown in Example 4-56.

Example 4-56 Stop command resulting messages

```
V TCPIP,TCPIPA,STOP,OSA2080
EZZ0060I PROCESSING COMMAND: VARY TCPIP,TCPIPA,STOP,OSA2080
EZZ0053I COMMAND VARY STOP COMPLETED SUCCESSFULLY
EZZ4329I LINK OSA20A0L HAS TAKEN OVER ARP RESPONSIBILITY FOR INACTIVE LINK
OSA2080L
EZZ4315I DEACTIVATION COMPLETE FOR DEVICE OSA2080
```

Note: When an OSA-Express device is stopped or loses its connection to the switch, another OSA-Express device defined to the TCP/IP stack will take over the IP address. A gratuitous ARP is broadcasted on the LAN to advertise the new MAC address related to the IP address being taken over. Message EZZ4329I in Example 4-56 indicates this action.

4.8.3 Activating modified device definitions

You can activate modified device definitions by issuing the OBEY command:

Vary TCP/IP, *tcpiplib*, OBEYFILE, *datasetname*

Authorization to use this command is through the user's RACF profile. The *datasetname* variable cannot be a z/OS UNIX file system file. The data set contains the modified TCP/IP configuration statements. See Example 4-57.

Example 4-57 OBEYFILE example

;Original BSDROUTINGPARMS statement for link OSA2080

```
;BSDROUTINGPARMS TRUE
; Link name      MTU      Cost metric  Subnet Mask  Dest address
;OSA2080L        1492      0           255.255.255.0  0
;ENDBSDROUTINGPARMS
```

;Modified BSDROUTINGPARMS statement for link OSA20C0

```
BSDROUTINGPARMS TRUE
; Link name      MTU      Cost metric  Subnet Mask  Dest address
OSA2080L        1024      0           255.255.255.0  0
ENDBSDROUTINGPARMS
```

Important: Dynamic XCF cannot be changed by using the OBEYFILE command. If you want to change the IPCONFIG DYNAMICXCF parameters, stop TCP/IP, code a new IPCONFIG DYNAMICXCF statement in the initial profile, and restart TCP/IP.

4.9 Problem determination

Isolating network problems is an essential step to verify a connectivity problem in your environment. This section introduces commands and techniques that can use to diagnose network connectivity problems related to a specific interface. The diagnostic commands that we discuss in this section are available for either the z/OS UNIX environment or the TSO environment.

The ping command

The **ping** command can be very useful for determining if a destination address can be reached in the network. Based on the results, it is possible to define whether the problem is related to the interface being tested, or whether it is a network-related problem.

Using **ping**, you can verify the following information:

- ▶ The directly-attached network is defined correctly.
- ▶ The device is properly connected to the network.
- ▶ The device is able to send and receive packets in the network.
- ▶ The remote host is able to receive and send packets.

When you issue a **ping** command, you can receive any of the responses listed in Table 4-4. See 8.3.1, “The ping command (TSO or z/OS UNIX)” on page 303 for more details about issuing the **ping** command.

Table 4-4 Using the ping command as a debugging tool

ping command (direct network)	ping response	Possible cause and actions
ping 10.1.2.11 (intf osa20801)	CS V1R12: Pinging host 10.1.2.11 sendMessage(): EDC8130I Host cannot be reached.	The interface being tested has a problem. Use the netstat command to verify the interface status.
ping 10.1.2.11 (intf osa20801)	CS V1R12: Pinging host 10.1.2.11 Ping #1 timed out.	The ICMP packet has been sent to the network, but the destination address is either invalid or it is not able to answer. Correct the destination address or verify the destination host status. This problem should be verified in the network.
ping 10.1.2.11 (intf osa20801)	CS V1R12: Pinging host 10.1.2.11 Ping #1 response took 0.000 seconds.	This is the expected response. The interface is working.

The netstat command

You can use the **netstat** command to verify the TCP/IP configuration. You need to check the information that is provided in the output from the **netstat** command against the values in the configuration data sets for the TCP/IP stack. To verify connectivity status from an interface perspective, use the following **netstat** options:

► netstat HOME/-h

Displays all defined interfaces and their IP addresses, even those interfaces that are created dynamically, as shown in Example 4-58.

Example 4-58 NETSTAT HOME command results

```
D TCPIP,TCPIPA,N,HOME
EZD0101I NETSTAT CS V1R12 TCPIPA 973
HOME ADDRESS LIST:
LINKNAME:  VIPA3L
ADDRESS:   10.1.30.10
FLAGS:
LINKNAME:  VIPA1L
ADDRESS:   10.1.1.10
FLAGS:    PRIMARY
LINKNAME:  VIPA2L
ADDRESS:   10.1.2.10
FLAGS:
LINKNAME:  IUTIQDF4L
ADDRESS:   10.1.4.11
FLAGS:
LINKNAME:  IUTIQDF5L
ADDRESS:   10.1.5.11
FLAGS:
LINKNAME:  IUTIQDF6L
ADDRESS:   10.1.6.11
FLAGS:
LINKNAME:  EZASAMEMVS
ADDRESS:   10.1.7.11
FLAGS:
```

```

LINKNAME:  IQDIOLNKOAO1070B
  ADDRESS:  10.1.7.11
  FLAGS:
LINKNAME:  VIPL0A010817
  ADDRESS:  10.1.8.23
  FLAGS:
LINKNAME:  VIPL0A010815
  ADDRESS:  10.1.8.21
  FLAGS:  INTERNAL
LINKNAME:  LOOPBACK
  ADDRESS:  127.0.0.1
  FLAGS:
INTFNAME:  OSA2080I
  ADDRESS:  10.1.2.11
  FLAGS:
INTFNAME:  OSA2081I
  ADDRESS:  10.1.2.14
  FLAGS:
INTFNAME:  OSA20A0I
  ADDRESS:  10.1.2.12
  FLAGS:
INTFNAME:  OSA20C0I
  ADDRESS:  10.1.3.11
  FLAGS:
INTFNAME:  OSA20E0I
  ADDRESS:  10.1.3.12
  FLAGS:
INTFNAME:  LOOPBACK6
  ADDRESS:  ::1
  TYPE:    LOOPBACK
  FLAGS:
17 OF 17 RECORDS
DISPLAYED

```

netstat DEVLINKS/-d

Displays the status of each interface, physical and logical, that is defined in the TCP/IP stack, as illustrated in Example 4-59 (only one interface shown as a sample).

Example 4-59 NETSTAT DEVLINKS command results

```

D TCP/IP,TCPIPA,N,DEV,INTFN=OSA2080I
EZD0101I NETSTAT CS V1R12 TCPIPA 996
INTFNAME: OSA2080I          INTFTYPE: IPAQENET   INTFSTATUS: READY
  PORTNAME: OSA2080    DATAPATH: 2082    DATAPATHSTATUS: READY
  CHPIDTYPE: OSD
  SPEED: 0000001000
  IPBROADCASTCAPABILITY: NO
  VMACADDR: 020010776873  VMACORIGIN: OSA    VMACROUTER: LOCAL
  ARPOFFLOAD: YES          ARPOFFLOADINFO: YES
  CFGMTU: 1492             ACTMTU: 1492
  IPADDR: 10.1.2.11/24
  VLANID: 10               VLANPRIORITY: DISABLED
  DYNVLANREGCFG: NO        DYNVLANREGCAP: YES
  READSTORAGE: GLOBAL (4096K)
  INBPERF: BALANCED
  CHECKSUMOFFLOAD: YES     SEGMENTATIONOFFLOAD: YES

```

```

SECCLASS: 255                                MONSYSPLEX: NO
ISOLATE: NO                                OPTLATENCYMODE: NO
MULTICAST SPECIFIC:
MULTICAST CAPABILITY: YES
GROUP          REFCNT          SRCFLTMD
-----
224.0.0.1      0000000001      EXCLUDE
SRCADDR: NONE
INTERFACE STATISTICS:
BYTESIN                      = 2672
INBOUND PACKETS              = 25
INBOUND PACKETS IN ERROR     = 0
INBOUND PACKETS DISCARDED    = 0
INBOUND PACKETS WITH NO PROTOCOL = 0
BYTESOUT                     = 3576
OUTBOUND PACKETS             = 39
OUTBOUND PACKETS IN ERROR    = 0
OUTBOUND PACKETS DISCARDED   = 0
IPV4 LAN GROUP SUMMARY
LANGROUP: 00001
NAME          STATUS        ARPOWNER          VIPAOWNER
-----
OSA20E0I      ACTIVE        OSA20E0I      YES
OSA2080I      ACTIVE        OSA2080I      NO
OSA20A0I      ACTIVE        OSA20A0I      NO
OSA20C0I      ACTIVE        OSA20C0I      NO
1 OF 1 RECORDS DISPLAYED
END OF THE REPORT

```

► **netstat ARP/-R** (for OSA-Express devices)

Used to query the ARP cache for a given address. Use this command when the remote host does not answer as expected, to check whether an ARP entry has been created for the remote host. It also allows you to check if the relationship between the IP and MAC address is the expected one. The resulting display is shown in Example 4-60 on page 201

DISPLAY TCPIP,TCPIPA,N,ARP

```
EZD0101I NETSTAT CS V1R12 TCPIPA 009
QUERYING ARP CACHE FOR ADDRESS 10.1.2.23
INTERFACE: OSA2080I          ETHERNET: 020011776873
QUERYING ARP CACHE FOR ADDRESS 10.1.2.61
INTERFACE: OSA2080I          ETHERNET: 020007776873
QUERYING ARP CACHE FOR ADDRESS 10.1.2.41
INTERFACE: OSA2080I          ETHERNET: 00145E776872
QUERYING ARP CACHE FOR ADDRESS 10.1.2.11
INTERFACE: OSA2080I          ETHERNET: 020010776873
QUERYING ARP CACHE FOR ADDRESS 10.1.2.13
INTERFACE: OSA2080I          ETHERNET: 020002776873
QUERYING ARP CACHE FOR ADDRESS 10.1.2.51
INTERFACE: OSA2080I          ETHERNET: 020003776873
QUERYING ARP CACHE FOR ADDRESS 10.1.2.21
INTERFACE: OSA2080I          ETHERNET: 020014776873
QUERYING ARP CACHE FOR ADDRESS 10.1.2.240
INTERFACE: OSA2080I          ETHERNET: 0014F1464600
QUERYING ARP CACHE FOR ADDRESS 10.1.2.24
INTERFACE: OSA20A0I          ETHERNET: 02001277688D
QUERYING ARP CACHE FOR ADDRESS 10.1.2.62
INTERFACE: OSA20A0I          ETHERNET: 02000677688D
QUERYING ARP CACHE FOR ADDRESS 10.1.2.60
INTERFACE: OSA20A0I          ETHERNET: 02000677688D
QUERYING ARP CACHE FOR ADDRESS 10.1.2.42
INTERFACE: OSA20A0I          ETHERNET: 00145E77688C
QUERYING ARP CACHE FOR ADDRESS 10.1.2.45
INTERFACE: OSA20A0I          ETHERNET: 00145E77688C
QUERYING ARP CACHE FOR ADDRESS 10.1.2.52
INTERFACE: OSA20A0I          ETHERNET: 02000277688D
QUERYING ARP CACHE FOR ADDRESS 10.1.2.12
INTERFACE: OSA20A0I          ETHERNET: 02001177688D
QUERYING ARP CACHE FOR ADDRESS 10.1.2.32
INTERFACE: OSA20A0I          ETHERNET: 02000B77688D
QUERYING ARP CACHE FOR ADDRESS 10.1.2.22
INTERFACE: OSA20A0I          ETHERNET: 02001577688D
QUERYING ARP CACHE FOR ADDRESS 10.1.2.240
INTERFACE: OSA20A0I          ETHERNET: 0014F1464600
QUERYING ARP CACHE FOR ADDRESS 10.1.3.11
INTERFACE: OSA20C0I          ETHERNET: 02000E776C05
QUERYING ARP CACHE FOR ADDRESS 10.1.3.14
INTERFACE: OSA20E0I          ETHERNET: 02000177855F
QUERYING ARP CACHE FOR ADDRESS 10.1.3.42
INTERFACE: OSA20E0I          ETHERNET: 00145E77855F
QUERYING ARP CACHE FOR ADDRESS 10.1.3.52
INTERFACE: OSA20E0I          ETHERNET: 02000277855F
QUERYING ARP CACHE FOR ADDRESS 10.1.3.12
INTERFACE: OSA20E0I          ETHERNET: 00145E77855F
QUERYING ARP CACHE FOR ADDRESS 10.1.3.62
INTERFACE: OSA20E0I          ETHERNET: 02000577855F
QUERYING ARP CACHE FOR ADDRESS 10.1.3.22
INTERFACE: OSA20E0I          ETHERNET: 02000C77855E
QUERYING ARP CACHE FOR ADDRESS 10.1.3.240
INTERFACE: OSA20E0I          ETHERNET: 0014F1464600
```

```
QUERYING ARP CACHE FOR ADDRESS 10.1.4.61
INTERFACE: IUTIQDF4L
QUERYING ARP CACHE FOR ADDRESS 10.1.4.41
INTERFACE: IUTIQDF4L
QUERYING ARP CACHE FOR ADDRESS 10.1.4.31
INTERFACE: IUTIQDF4L
QUERYING ARP CACHE FOR ADDRESS 10.1.4.25
INTERFACE: IUTIQDF4L
QUERYING ARP CACHE FOR ADDRESS 10.1.4.21
INTERFACE: IUTIQDF4L
QUERYING ARP CACHE FOR ADDRESS 10.1.4.12
INTERFACE: IUTIQDF4L
QUERYING ARP CACHE FOR ADDRESS 10.1.4.11
INTERFACE: IUTIQDF4L
QUERYING ARP CACHE FOR ADDRESS 10.1.5.61
INTERFACE: IUTIQDF5L
QUERYING ARP CACHE FOR ADDRESS 10.1.5.41
INTERFACE: IUTIQDF5L
QUERYING ARP CACHE FOR ADDRESS 10.1.5.31
INTERFACE: IUTIQDF5L
QUERYING ARP CACHE FOR ADDRESS 10.1.5.25
INTERFACE: IUTIQDF5L
QUERYING ARP CACHE FOR ADDRESS 10.1.5.21
INTERFACE: IUTIQDF5L
QUERYING ARP CACHE FOR ADDRESS 10.1.5.12
INTERFACE: IUTIQDF5L
QUERYING ARP CACHE FOR ADDRESS 10.1.5.11
INTERFACE: IUTIQDF5L
QUERYING ARP CACHE FOR ADDRESS 10.1.6.11
INTERFACE: IUTIQDF6L
QUERYING ARP CACHE FOR ADDRESS 10.1.7.61
INTERFACE: IQDIOLNKOAO1070B
QUERYING ARP CACHE FOR ADDRESS 10.1.7.51
INTERFACE: IQDIOLNKOAO1070B
QUERYING ARP CACHE FOR ADDRESS 10.1.7.41
INTERFACE: IQDIOLNKOAO1070B
QUERYING ARP CACHE FOR ADDRESS 10.1.7.31
INTERFACE: IQDIOLNKOAO1070B
QUERYING ARP CACHE FOR ADDRESS 10.1.7.21
INTERFACE: IQDIOLNKOAO1070B
QUERYING ARP CACHE FOR ADDRESS 10.1.7.12
INTERFACE: IQDIOLNKOAO1070B
QUERYING ARP CACHE FOR ADDRESS 10.1.7.11
INTERFACE: IQDIOLNKOAO1070B
48 OF 48 RECORDS DISPLAYED
END OF THE REPORT
```

These commands can help you to locate connectivity problems. If they do not, the next step in debugging a direct attached network problem is to gather documentation that shows more detailed information about traffic problems related to the interface and network.

To get this detailed information, the z/OS Communications Server typically uses the component trace to capture event data and save it to an internal buffer—or write the internal buffer to an external writer, if requested. You can later format these trace records using the Interactive Problem Control System (IPCS) subcommand CTRACE.

To debug a network connectivity problem you can use the Component trace with either of the two specific components, as follows:

- ▶ SYSTCPIP component trace with options
 - VTAM, which shows all of the nondata-path signaling occurring between the devices and VTAM
 - VTAMDATA, which shows data-path signaling between the devices and VTAM, including a snapshot of media headers and some data

Important: Using this option slows performance considerably; therefore, use it with caution.

- ▶ SYSTCPDA component trace, used with the VARY TCPIP,PKTTRACE command. You can use the PKTTRACE statement to copy IP packets as they enter or leave TCP/IP, and then examine the contents of the copied packets.

For more information about how to set up and activate a CTRACE, refer to Chapter 8, “Diagnosis” on page 299.

- ▶ OSA-Express network traffic analyzer (OSAENTA) trace

This trace provides a way to trace inbound and outbound frames for an OSA-Express3 and OSA-Express2 feature in QDIO mode.

- The function allows the z/OS Communications Server to control and format the tracing of frames collected in the OSA-Express3 and OSA-Express2 feature at the network port.
- It also provides the capability to trace frames discarded by the OSA-Express3 and OSA-Express2 feature.

SYSTCPOT is a new CTRACE component for collecting NTA trace data. The trace records can be formatted using the IPCS CTRACE command, specifying a component name of SYSTCPOT.

For more information about how to set up and enable the network traffic analyzer, refer to Chapter 8, “Diagnosis” on page 299.

4.10 Additional information

For additional information, refer to:

- ▶ *HiperSockets Implementation Guide*, SG24-6816
- ▶ *OSA-Express Implementation Guide*, SG24-5948
- ▶ *z/OS Communications Server: IP Configuration Reference*, SC31-8776
- ▶ *z/OS Communications Server: SNA Resource Definition*, SC31-8778



Routing

One of the major functions of a network protocol such as TCP/IP is to efficiently interconnect a number of disparate networks. These networks can include LANs and WANs, fast and slow, reliable and unreliable, inexpensive and expensive connections.

To interconnect these networks, some level of intelligence is needed at the boundaries to look at the data packets as they pass, and make rational decisions as to where and how they should be forwarded. This is known as IP routing. In this chapter we look at the various types of IP routing supported in a z/OS Communications Server environment.

This chapter includes the following topics.

Section	Topic
5.1, "Basic concepts" on page 206	The basic concepts of IP routing
5.2, "Routing in the z/OS environment" on page 212	Key characteristics of IP routing in z/OS Communications Server and performance considerations
5.3, "Dynamic routing protocols" on page 217	Detailed characteristics of dynamic routing protocols
5.4, "Implementing static routing in z/OS" on page 227	The implementation tasks and configuration examples for static routing
5.5, "Implementing OSPF routing in z/OS with OMPROUTE" on page 233	The implementation tasks and configuration examples for OSPF dynamic routing
5.6, "Problem determination" on page 253	Techniques for problem determination

5.1 Basic concepts

When we talk about networks, one of the key issues is how to transport data across the network. Based on the OSI reference model, the act of moving data traffic across a network from a source to a destination can be accomplished either by bridging or routing this data between the endpoints.

Bridging is often compared with routing, which might seem to accomplish precisely the same thing. However, note the primary difference between these functions:

- ▶ Bridging occurs at Layer 2 (the data link control layer) of the OSI reference model.
- ▶ Routing occurs at Layer 3 (the network layer).

This distinction provides bridging and routing with different information to use in the process of moving information from source to destination, so the two functions accomplish their tasks in different ways.

5.1.1 Terminology

To help understand the concepts described in this section, Table 5-1 lists some of the common IP routing-related terms. Most of the functions or protocols listed are supported by the z/OS Communications Server.

Table 5-1 IP routing terms

Term	Definition
Routing	The process used in an IP network to deliver a datagram to the correct destination.
Routing daemon	A server process that manages the IP routing table. OMPROUTE is the z/OS Communications Server component that acts as the routing daemon.
Replaceable <i>static routes</i>	Static routes that can be replaced by OMPROUTE.
Dynamic routing	Routing that is dynamically managed by a routing daemon and automatically changes in response to network topology changes.
Static routing	Routing that is manually configured and does not change automatically in response to network topology changes.
Autonomous system (AS)	A group of routers exchanging routing information through a common routing protocol. A single AS can represent a large number of IP networks.
Router	A device or host that interprets protocols at the Internet Protocol (IP) layer and forwards datagrams on a path toward their correct destination.
Gateway	A router that is placed between networks or subnetworks. The term is used to represent routers between autonomous systems.
Interior gateway protocols (IGP)	Dynamic route update protocol used between dynamic routers running on TCP/IP hosts within a single autonomous system.
Exterior gateway protocols (EGP)	Dynamic route update protocols used between routers that are placed between two or more autonomous systems.

In order to route packets in the network, each network interface must have a unique IP address assigned. Whenever a packet is sent, the destination and source IP addresses are included in the packet's header information. The network layer (Layer 3) of the TCP/IP stack examines the destination IP address to determine how the packet should be forwarded. The packet is either sent to its destination on the same network (direct routing) or, based on a routing table entry, to another network using a router (indirect routing).

5.1.2 Direct routes, indirect routes, and default route

Every IP host is capable of routing IP datagrams and maintaining an IP routing table. There are three types of entries in an IP routing table:

- Direct routes

The networks to which the host is directly attached are called direct routes. If the destination host is attached to the same IP network as the source host, IP datagrams can be exchanged directly.

- Indirect routes

When the destination host is not connected to the same IP network as the source host, the only way to reach the destination host is through one or more IP routers. The routing entry with the destination IP address and the IP address of the first router (the next hop) is called an indirect route in the IP routing algorithm.

The IP address of the first router is the only information required by the source host to send a packet to the destination host. If the source and destination hosts are on the same physical network, but belong to different subnetworks, indirect routing is used to communicate between the endpoints. A router is needed to forward packets between subnetworks.

- The default route

The default route entry contains the IP address of the first router (the next hop) to be used when the destination IP address or network is not found in any of the direct or indirect routes.

Figure 5-1 illustrates the concept of IP routing.

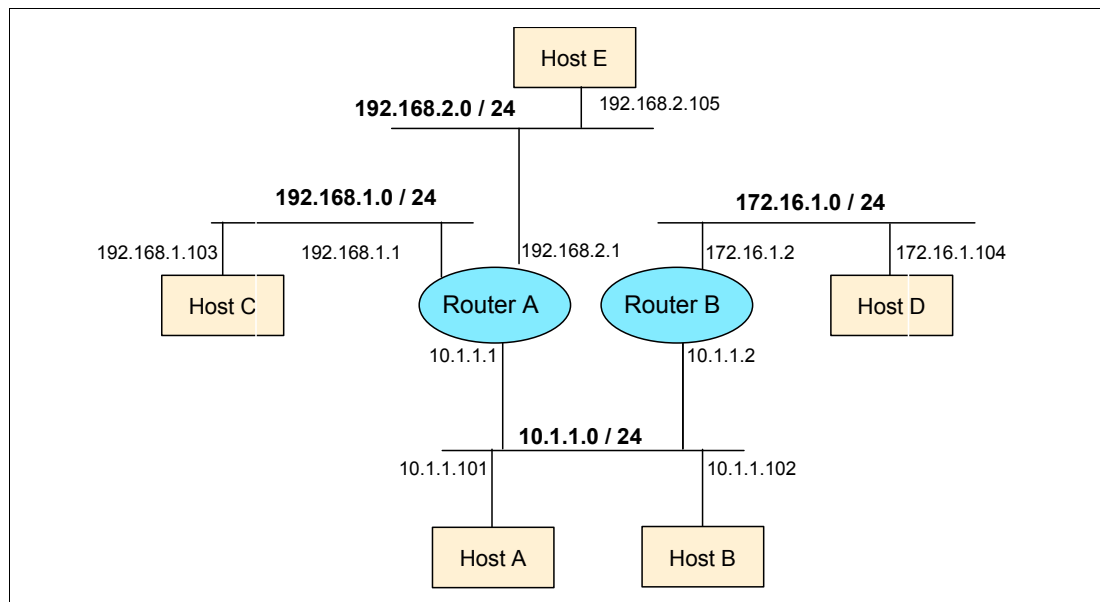


Figure 5-1 Sample network with multiple subnetworks

This example has hosts and routers located in multiple networks, and to achieve connectivity between these hosts, the routers are connected to multiple networks, creating a path between them.

In this scenario, if host A wants to connect to host D, both resources must create and maintain a routing table to define which path has to be used to reach its destination. Host A in this example might contain the (symbolic) entries as listed in Table 5-2.

Table 5-2 IP routing table for Host A

Destination	IP address of Next Hop Router
10.1.1.0/24	Directly connected
192.168.1.0/24	10.1.1.1 (Router A)
172.16.1.0/24	10.1.1.2 (Router B)
Default	10.1.1.1 (Router A)
127.0.0.1	Loopback

The routing table contains routes to different routers in this network. When host A has an IP datagram to forward, it determines which IP address to forward it to using the IP routing algorithm and the routing table.

Note: The /24 (a *prefix*), represents the length of subnet mask (a 24-bit mask, in this case).

Because Host A is directly attached to network 10.1.1.0/24, it maintains a direct route for this network. To reach other networks such as 192.168.1.0/24 and 172.16.1.0/24, it must have an indirect route through router A and router B, respectively, because these networks are not directly attached to it. Another option is to define a default route. If the indirect route to the network is not defined explicitly, the default route is used.

In this example, Host A reaches to Host B using the direct route. To reach Host C (192.168.1.103), it uses the indirect route to 192.168.1.0/24 and forwards the packet to Router A (10.1.1.1).

Likewise, to reach Host D, it uses the indirect route to 172.16.1.0/24 and forwards the packet to Router B (10.1.1.2). The indirect route to Host E (192.168.2.105) is not explicitly defined in the Host A. So, the default route is used and the Host A forwards the packet to Router A (10.1.1.1).

To reach any given IP network address, each host or router in the network needs to know only the next hop's IP address and not the full network topology.

5.1.3 Route selection

IP uses a unique algorithm to route an IP datagram. In a network without subnetworks, each host in the path from source host to destination host will:

1. Inspect the destination address of the packet.
2. Divide the destination address into network and host addresses.
3. Determine whether the network is directly attached:
 - If so, send the IP datagram directly to the destination.
 - If not, send the IP datagram to the next router, as defined by the routing tables.

In a subnetted network, each host in the path from source host to destination host will:

1. Inspect the destination address of the packet.
2. Divide the destination address into subnetwork and host addresses.
3. Determine whether the subnetwork is directly attached:
 - If so, forward the packet directly to the destination.
 - If not, forward the packet to the next router as defined in the routing tables.

If two or more indirect routes are defined for the same destination, the route selection depends on the implementation of the routers or hosts. Some implementation always uses the top entry in the list, and some implementation uses all routes to distribute the packets. In some cases it is configurable with the provided parameters.

If two or more indirect routes are defined for the same destination but with different subnet mask length, the route with longest mask length is selected. This method is called *the longest match*.

5.1.4 Static routing and dynamic routing

In this section we explain the two ways to set up the necessary routing table in a system: using static routing, or dynamic routing.

Static routing

Static routing requires you to *manually* configure the routing tables yourself. This task is part of the configuration steps you follow when customizing TCP/IP. It implies that you know the address of every network you want to communicate with and how to get there. That is, you must know the address of the first router on the way.

The task of statically defining all necessary routes can be simple for a small network. It offers the advantage of avoiding the network traffic overhead of a dynamic route update protocol. It also allows you to enforce rigid control on the allocation of addresses and resource access. However, it will require manual reconfiguration if you move or add a resource.

The another drawback of static routing is that, even if the network failure occurs in the intermediate path to the destination, the routing table remains unchanged and keeps sending the packet according to the statically defined next hop routers. Sometimes it might cause the network to be unreachable. Also, if you fail to define the right next hop router in the route entry, the routers keep forwarding the packet using that entry. Even if there is a better route, the router does not change its next hop router until the changes are made to the static route entry.

If your network environment is small and manageable, with few to no network changes anticipated, then using static routes is an option (keeping in mind that your z/OS system is basically an application server environment). A good practice is to define only the default gateways to the exterior networks, and let the routers do the exterior routing. You can implement the static routing between the z/OS system and external router, and still let the external routers use the dynamic routing protocol to exchange route information.

Dynamic routing

Dynamic routing removes the need for static definition of the routing table. The network routing table is built dynamically, automatically exchanging route information among the routers in the network. This sharing of the routing information enables the routers to always calculate the best path through the network to any destination. When a network outage occurs in the intermediate route to the destination, the routers exchange the information about the outage and the best path is recalculated.

If your routing tables are complex due to network growth, or if the system must act as a gateway, it is far easier to let the system do the work for you by using dynamic routing.

The drawback of dynamic routing is the burden of route information exchange. There are some configuration techniques you can use to reduce this burden, as explained in 5.2, “Routing in the z/OS environment” on page 212.

Dynamic routing protocols can be divided into two types: interior gateway protocols, and exterior gateway protocols.

Interior gateway protocols (IGPs) are dynamic route update protocols used between dynamic routers running on TCP/IP hosts within a single autonomous system. These protocols are used by the routers to exchange information about which IP routes the IP hosts that they have knowledge of. By exchanging IP routing information with each other, the routers are able to maintain a complete picture of all available routes inside an autonomous system.

Exterior gateway protocols (EGPs) are dynamic route update protocols that are used between routers that are placed between two or more Autonomous Systems.

OSPF and RIP

In this section we discuss the interior gateway protocols OSPF, RIP version1, and RIP version2, which are supported by z/OS Communications Server.

- ▶ Open Shortest Path First (OSPF)

OSPF uses a link state or shortest path first algorithm. OSPF's most significant advantage compared to RIP is the reduced time needed to converge after a network change. In general, OSPF is more complicated to configure than RIP and might not be suitable for small networks.

- ▶ Routing Information Protocol (RIP)

RIP uses a distance vector algorithm to calculate the best path to a destination based on the number of hops in the path. RIP has several limitations. Some of the limitations that exist in RIP version 1 are resolved by RIP version 2.

RIP version 2 expands RIP version 1. Among the improvements are support for multicasting and variable subnetting. Variable subnetting allows the division of networks into variable size subnets.

- ▶ IPv6 OSPF

IPv6 OSPF (OSPFv3) uses a link state or shortest path first algorithm to calculate the best path to a destination. IPv6 OSPF has the same advantages and more complicated configuration compared to IPv6 RIP (as with OSPF compared to RIP).

- ▶ IPv6 RIP

IPv6 RIP uses the same distance vector algorithm used by RIP to calculate the best path to a destination. It is intended to allow routers to exchange information for computing routes through an IPv6-based network.

Table 5-3 lists the main characteristics of the routing protocols supported by the z/OS Communications Server.

Table 5-3 Interior Gateway Protocol characteristics

	RIP V1	RIP V2	IPv6 RIP	OSPF	IPv6 OSPF
Algorithm	Distance vector	Distance vector	Distance vector	Shortest path first	Shortest path first
Network load ^a	High	High	High	Low	Low
CPU processing requirements ^a	Low	Low	Low	High	High
IP network design restrictions	Many	Some	Some	Virtually none	Virtually none
Convergence time	Up to 180 seconds	Up to 180 seconds	Up to 180 seconds	Low	Low
Multicast support ^b	No	Yes	Yes	Yes	Yes
Multiple equal-cost routes	No ^c	No ^c	No ^c	Yes	Yes

a. Depends on network size and stability.

b. Multicast saves CPU cycles on hosts that do not require certain periodic updates, such as OSPF link state advertisements or RIP-2 routing table updates. Multicast frames are filtered out, either in the device driver or directly on the interface card, if this host has not joined the specific multicast group.

c. RIP in OMROUTE allows multiple equal-cost routes only for directly connected destination over redundant interfaces.

5.1.5 Choosing the routing method

The choice of a routing protocol is a major decision for the network administrator, and has a major impact on overall network performance. The selection depends on the network complexity, size, and administrative policies. The protocol chosen for one type of network might be inappropriate for other types of networks. Each unique environment must be evaluated against a number of fundamental design requirements, as explained here.

► Scalability to large environments

The potential growth of the network dictates the importance of this requirement. If support is needed for large, highly redundant networks, then link state or hybrid algorithms should be considered. Distance vector algorithms do not scale into these environments. Static routing also does not usually scale into large environments.

► Stability during outages

Distance vector algorithms can introduce network instability during outage periods. The counting to infinity problems might cause routing loops or other non-optimal routing paths. Link state or hybrid algorithms reduce the potential for these problems. Static routing can provide stability if the platform implements protocols like Virtual Router Redundancy Protocol (VRRP), Hot Standby Router Protocol (HSRP), or if redirected routes are accepted.

On a z platform, OSAs can provide stability in a static routing environment through a feature called ARP Takeover. Refer to *IBM z/OS V1R11 Communications Server TCP/IP Implementation Volume 3: High Availability, Scalability, and Performance*, SG24-7800 for more detailed information about ARP Takeover.

- Speed of convergence

Triggered updates provide the ability to immediately initiate convergence when a failure is detected. All three types of protocols support this feature.

One contributing factor to convergence is the time required to detect a failure. In OSPF networks, a series of “hello” packets must be missed before convergence begins.

In RIP environments, subsequent route advertisements must be missed before convergence is initiated.

These detection times increase the time required to restore communication. In static routing environments, convergence is a factor limited by the time it takes to update static routing tables manually.

- Metrics

Metrics provide the ability to groom appropriate routing paths through the network. Link state algorithms consider bandwidth when calculating routes.

- Vendor interoperability

The types of devices deployed in a network indicate the importance of this requirement. If the network contains equipment from a number of vendors, then standard routing protocols should be used. The IETF has dictated the operating policies for the distance vector and link state algorithms described in this book. Implementing these algorithms avoids any interoperability problems encountered with nonstandard protocols.

The administrator must assess the importance of each of these requirements when determining the appropriate routing protocol for an environment.

5.2 Routing in the z/OS environment

This section discusses the two IP routing methods provided by z/OS Communications Server. It also describes OMROUTE and what to consider when implementing dynamic and static routes.

5.2.1 Static routing

In z/OS Communications Server, the static routes are defined with the BEGINROUTES statement block in the TCP/IP profile. The defined static routes are installed into the routing table of the TCP/IP stack. The GATEWAY statement also can be used in the TCP/IP profile to define static routes, but it is an obsolete statement. Instead, we recommend the use of the BEGINROUTES statement block.

Static routing can be combined with dynamic routing through the use of the OMROUTE routing daemon. If the ROUTE statement in the BEGINROUTES statement block is coded with NOREPLACEABLE, then the static route is always preferred over the dynamically learned route for the same destination with the same subnet mask length.

If two or more routes to the same destination with same subnet mask length are defined in the z/OS Communications Server routing table, then the TCP/IP stack always uses the first *active* entry, by default. If you specify a IPCONFIG MULTIPATH statement in the TCP/IP profile, all routes for the same destination are used by per connection or per packet, depending on which option you specify for MULTIPATH.

5.2.2 Dynamic routing using OMPROUTE

In z/OS Communications Server IP, there is a multiprotocol routing daemon for dynamic routing called OMPROUTE. (The term *daemon* is used in UNIX to refer to a background server process.) It provides an alternative to the static TCP/IP routing definitions. The z/OS host running with OMPROUTE becomes an active OSPF or RIP router in a TCP/IP network. Either or both of these routing protocols can be used to dynamically maintain the routing table.

Supported dynamic routing protocols

OMPROUTE supports the OSPF, RIP version 1, and RIP version 2 routing protocols.

For IPv4, OMPROUTE implements the OSPF protocol described in RFC 1583 (OSPF version 2), the OSPF subagent protocol described in RFC 1850 (OSPF version 2 Management Information Base), and the RIP protocols described in RFC 1058 (Routing Information Protocol) and in RFC 1723 (RIP version 2 - Carrying Additional Information).

For IPv6, OMPROUTE implements the IPv6 RIP protocol described in RFC 2080 (RIPng for IPv6) and the IPv6 OSPF protocol described in RFC 2740 (OSPF for IPv6).

How OMPROUTE works

OMPROUTE manages an OMPROUTE routing table. OMPROUTE installs the routes that are learned dynamically through other routers with routing protocol (OSPF or RIP) to the TCP/IP stack's routing table. When routing a packet to its destination, the TCP/IP stack makes decisions for route selection based on TCP/IP stack's routing table, not the OMPROUTE routing table.

A one-to-one relationship exists between an OMPROUTE and a TCP/IP stack. OSPF/RIP support for multiple TCP/IP stacks requires multiple instances of OMPROUTE. The affinity to the TCP/IP stack is made by specifying the TCPIPJobname statement with the TCP/IP stack name in TCPIP.DATA file that OMPROUTE uses.

OMPROUTE supports Virtual IP Addressing (VIPA) to handle network interface failures by switching to alternate paths. VIPA routes are included in the OSPF and RIP advertisements to adjacent routers. Adjacent routers learn about VIPA routes from advertisements and can use them to reach destinations at the z/OS.

OMPROUTE does not make use of the BSDROUTINGPARMS statement. Instead, its parameters are defined in the OMPROUTE configuration file. The OMPROUTE configuration file is used to define both OSPF and RIP environments.

Note: If the INTERFACE statement is used in the TCP/IP stack to define an interface, the subnet mask and MTU coded in OMPROUTE need to agree, or OMPROUTE will issue an error message and use the values you configure to OMPROUTE.

For IPv4, the OSPF and RIP protocols are communicated over interfaces defined with the OSPF_INTERFACE and RIP_INTERFACE configuration statements. Interfaces that are not involved in the communication of the RIP or OSPF protocol are configured with the INTERFACE configuration statement (unless it is a non-point-to-point interface and all default values specified on the INTERFACE statement are acceptable).

If both OSPF and RIP protocols are used in an OMPROUTE environment, then OSPF takes precedence over RIP. OSPF routes will be preferred over RIP routes to the same destination.

OMPROUTE allows the generation of multiple, equal-cost routes to a destination (with OSPF, not RIP). If there are multiple routes for same destination with the same subnet mask length, the stack uses the first active route for all traffic. If you specify an IPCONFIG MULTIPATH statement in the TCP/IP profile, the stack uses all routes for the same destination per connection or per packet, depending on which option you specify for MULTIPATH.

Considerations: combining OMPROUTE with BEGINROUTES

Note that when coding static routes in BEGINROUTES statements, in conjunction with the OMPROUTE configuration, you have the following options for static routes:

- ▶ NOREPLACEABLE (the default)
- ▶ REPLACEABLE

OMPROUTE does *not* replace a NOREPLACEABLE static route, even if it has detected a dynamic route to the same destination, and the TCP/IP stack uses a NOREPLACEABLE static route to forward the packet. OMPROUTE replaces a REPLACEABLE static route if it detects a dynamic route to the same destination. The REPLACEABLE option enables the last resort to the destination in cases where OMPROUTE has not detected a dynamic route to the destination.

Also, take care to ensure that the z/OS Communications Server host is not overly burdened with routing work. Unlike routers or other network boxes whose sole purpose is routing, an application host z/OS Communications Server will be doing many things other than routing, and it is not desirable for a large percentage of machine resources (memory and CPU) to be used for routing tasks, as can happen in very complex or unstable networks.

The most common and recommended way to use dynamic routing in the z/OS environment is to define the stack as a OSPF Stub Area or, even better, as a Totally Stubby Area. Stub and Totally Stubby Areas minimize the amount of routing work that z/OS must perform.

Effect of storage shortages on OMPROUTE

Dynamic routing protocols depend on the timely exchange of routing updates with neighbor routers in the network. Responsiveness of the dynamic routing nodes and the network is essential to maintaining valid routing tables. If OMPROUTE fails to receive routing updates from neighbors, the dynamic routes learned from these neighbors are deleted from the stack route table. If OMPROUTE fails to send updates to neighbors, the dynamic routes affected by the missing updates are deleted from the stack route table at these neighbors. If OMPROUTE exits for any reason, the dynamic routes in the stack route table are not deleted, but they become stale because they no longer reflect an accurate network status.

Given the need for a responsive OMPROUTE node, a storage shortage in the node can lead to lost connectivity in the network. For example, OMPROUTE might exit if the stack is unable to allocate storage for OMPROUTE dispatchable unit control blocks or for sending routing updates to neighbor routers. Messages that advise you about storage shortages are depicted in Example 5-1.

Example 5-1 Messages that indicate storage shortages

```
EZZ4360I jobname ECSA CONSTRAINED
EZZ4361I jobname ECSA CRITICAL
EZZ4364I jobname POOL CONSTRAINED
EZZ4365I jobname POOL CRITICAL
IVT5591I CSM ECSA STORAGE AT CONSTRAINED LEVEL
IVT5562I CSM ECSA STORAGE AT CRITICAL LEVEL
IVT5592I CSM FIXED STORAGE AT CONSTRAINED LEVEL
IVT5563I CSM FIXED STORAGE AT CRITICAL LEVEL
```

When storage shortages are relieved, other console messages advise you of this fact, as shown in Example 5-2.

Example 5-2 Messages that indicate storage shortage relief

```
EZZ4363I jobname ECSA SHORTAGE RELIEVED
EZZ4367I jobname POOL SHORTAGE RELIEVED
IVT5564I CSM ECSA STORAGE SHORTAGE RELIEVED
IVT5565I CSM FIXED STORAGE SHORTAGE RELIEVED
```

Proper design of the dynamic routing environment can eliminate or reduce the likelihood of storage shortages that affect OMPROUTE. For example, the most common and recommended way to use dynamic routing in the z/OS environment is to define the stack as an OSPF Stub Area or, even better, as a Totally Stubby Area.

Stub Areas minimize storage and CPU processing at the nodes that are part of the Stub Area because they maintain less knowledge about the topology of the Autonomous System (AS) than do other types of non-backbone routers. They maintain knowledge only of intra-area destinations and summaries of inter-area destinations and default routes within the AS in order to reach external destinations.

A Totally Stubby Area receives even less routing information than a Stub Area. It only knows of intra-area destinations and default routes within the Stub Area to reach external destinations. Thus, its storage and CPU processing requirements are even less than what is required for a Stub Area.

Providing tolerance for storage shortage conditions affecting OMPROUTE

OMPROUTE and the TCP/IP stack work together to provide tolerance for storage shortage conditions. Notifications are sent to OMPROUTE by the TCP/IP stack to inform OMPROUTE when the stack enters or exits a storage shortage condition. During a storage shortage, OMPROUTE uses these notifications to temporarily suspend the requirement that it receive periodic routing updates from neighbor routers.

The TCP/IP stack ensures that there are always control blocks available for dispatchable units doing work for OMPROUTE. In addition, the stack satisfies requests for stack storage made on behalf of OMPROUTE as long as storage remains available. Requests made on behalf of other applications are not satisfied during a storage shortage.

These actions temporarily keep OMPROUTE from deleting routes during a storage shortage when OMPROUTE fails to receive the usual periodic routing updates from neighboring routers. In addition, they decrease the likelihood that OMPROUTE will exit, time out routes, or fail to send routing updates to neighbor routers during a storage shortage. This temporary reprieve lasts for five minutes, at which time OMPROUTE automatically resumes the requirement for periodic routing table updates.

Learning of OMPROUTE's tolerance to a storage shortage

OMPROUTE displays can reveal that OMPROUTE is responding to a storage shortage condition. For example, the detailed information about an OSPF neighbor could show that the time interval since receipt of the last HELLO packet is longer than the configured Dead Router Interval. Example 5-3 shows an example of this display, where the Dead Router Interval is 40 seconds, but the HELLO packet was last received 60 seconds ago **1**.

Example 5-3 OSPF neighbor display

```
D TCPIP,TCPIPA,OMPROUTE,OSPF,NEIGHBOR,IPADDR=10.1.2.240
EZZ7852I NEIGHBOR DETAILS 968
      NEIGHBOR IP ADDRESS:    10.1.2.240
      OSPF ROUTER ID:         10.1.3.240
      NEIGHBOR STATE:         128
      PHYSICAL INTERFACE:     OSA2080I
      DR CHOICE:               10.1.2.240
      BACKUP CHOICE:           0.0.0.0
      DR PRIORITY:             100
      NBR OPTIONS:             (0X50)
DB SUMM QLEN:      0  LS RXMT QLEN:      0  LS REQ QLEN:      0
LAST HELLO:        1 60 NO HELLO:         OFF
# LS RXMITS:        1  # DIRECT ACKS:      0  # DUP LS RCVD:    11
# OLD LS RCVD:      1  # DUP ACKS RCVD:    0  # NBR LOSSES:     0
# ADJ. RESETS:      0
```

With RIP routes, you might discover that OMPROUTE is responding to the shortage event when several route displays reveal that the age of RIP routes ceases to increase. See an example of such a display at **2** in Example 5-4. Several iterations of the OMPROUTE command showed that the age of the route never increased beyond 10.

Example 5-4 Display of OMPROUTE RTTABLE

```
D TCPIP,,OMPROUTE,RTTABLE
EZZ7847I ROUTING TABLE 796
TYPE DEST      NET MASK COST AGE NEXT HOP(S)
RIP 30.1.1.0 FFFFFFF00 2    10 2 9.67.103.6
```

A trace of OMPROUTE activity using a trace level of -t2 and a debug level of -d1 also provides information about OMPROUTE's automatic tolerance of a storage shortage condition. Messages shown in Example 5-5 advise you that OMPROUTE is reacting as designed to a storage shortage.

Example 5-5 OMPROUTE trace messages for toleration of storage shortage

```
EZZ8166I Received type storage shortage notification for ip_version
EZZ8167I OSPF dead router checking is resumed for ip_version
EZZ8168I OSPF dead router checking is suspended for ip_version
EZZ8169I RIP route aging is resumed for ip_version
EZZ8170I RIP route aging is suspended for ip_version

IPv4 route aging bypassed - in stack storage shortage
IPv6 route aging bypassed - in stack storage shortage
IPv4 dead router checks bypassed - in stack storage shortage
IPv6 dead router checks bypassed - in stack storage shortage
```

In Example 5-5, the value of the type field can be *begin* or *end* and the *ip_version* field can be *IPv4* or *IPv6*.

Despite OMPROUTE's built-in tolerance to storage shortages, problems can still occur. First, the already mentioned relief from storage shortage conditions lasts only five minutes. If the storage shortage lasts longer, the local routes begin to be deleted from the stack route table if updates from neighbors are still not reaching the stack. Second, OMPROUTE will still exit if stack storage becomes totally exhausted so that OMPROUTE can no longer send data. In such a case, messages other than those depicted in Example 5-5 on page 216 or in addition to these messages could appear on the console or in a trace to advise you of these further problems.

5.2.3 Policy-based routing

In a TCP/IP environment, the route is selected based on the destination IP address of the packet. The TCP/IP routing table is looked up for the matching entry for the destination IP address. This means that all types of packets destined to the same destination IP address, including interactive traffic (TSO, for example) and bulk traffic (FTP, for example), are forwarded to the same next hop router. In some cases, the bulk traffic might cause traffic congestion and can lead to a performance problem for interactive traffic.

The policy-based routing determines the destination based on the defined policy. Traffic descriptors such as TCP/UDP port numbers, application name, and source IP addresses can be used to define the policy to enable the optimized route selection.

Policy-based routing can use both static routes and dynamic routes, which are obtained with the OMPROUTE routing daemon.

For detailed information about policy-based routing, refer to *IBM z/OS V1R11 Communications Server TCP/IP Implementation Volume 4: Security and Policy-Based Networking*, SG24-7801.

5.3 Dynamic routing protocols

z/OS Communications Server supports two different type of dynamic routing:

- ▶ Open Shortest Path First
- ▶ Routing Information Protocol

5.3.1 Open Shortest Path First

This section provides a brief overview of the Open Shortest Path First (OSPF) routing protocol.

The OSPF protocol is based on link-state or shortest path first technology. In other words, OSPF routing tables contain details of the connections between routers, their status (active or inactive), their cost (desirability for routing), and so on.

Updates are broadcast whenever a link changes status, and consist merely of a description of the changed status. OSPF can divide its network into topology subsections, known as *areas*, within which broadcasts are confined. OSPF has been designed for the TCP/IP Internet environment. In CS for z/OS IP, OSPF is configured using the UNIX daemon OMPROUTE.

OSPF features include the following:

- ▶ OSPF supports variable length subnetting.
- ▶ OSPF can be configured so that all its protocol exchanges are authenticated.
- ▶ Only trusted routers can participate in an AS that has been configured with authentication.
- ▶ Least-cost routing allows you to configure path costs based on any combination of network parameters. Bandwidth, delay, and metric cost are some examples.
- ▶ There are no limitations to the routing metric. Although RIP restricts the routing metric to 16 hops, OSPF has virtually no restrictions.
- ▶ Multipath routing is allowed. OSPF supports multiple paths of equal cost that connect the same points. These paths are then used for network load distribution, resulting in more use of the network bandwidth.
- ▶ OSPF's area routing capability provides an additional level of routing protection and a reduction in routing protocol traffic.

OSPF terminology

This section describes some of the more common IP routing-related terms and concepts used in OSPF.

- ▶ Router ID

This is a 32-bit number allocated to each router in the OSPF network protocol. This number is unique in the autonomous system. It represents the IP address of an interface defined on the OSPF node.

For the z/OS implementation of the Router ID in OSPF, we recommend the use of a static VIPA address. Do not use a Dynamic VIPA as the Router ID, because the movement of the Router ID causes confusion in the OSPF routing protocol exchanges.

- ▶ Areas

OSPF networks can be divided into areas. An area consists of networks and routers that are logically grouped together. All routers within an area maintain the same topology database.

All OSPF networks consist of at least one area, typically the backbone area. If you define more than one area, one of the areas must be the backbone area and the other area or areas are defined as non-backbone areas.

- ▶ Backbone area

All OSPF networks should have a backbone area. The area identifier of the backbone area is always 0.0.0.0. The backbone area is special in that it distributes routing information to all areas connected to it.

- ▶ Area border routers

These are routers that connect two or more areas. The area border router maintains a topology database of each area to which it is attached. All area border routers must have at least one interface in the backbone area. A virtual link can be used to satisfy this requirement.

- ▶ AS boundary routers

These are the routers that connect the OSPF internetwork and exchange reachability information with other routers in other Autonomous Systems. They can use the exterior gateway protocols. The AS boundary routers are used to import static routes, RIP routes into the OSPF network (and vice versa).

- ▶ Virtual link

This is a logical link that connects an area that does not have a physical link to a backbone area. The link is treated as a point-to-point link.

- ▶ Neighboring routers

Routers that have interfaces to the same connection are called neighboring routers. To become neighbors, routers must belong to the same OSPF area, use the same security scheme, and have the same Hello and Dead intervals.

- ▶ Adjacency

Neighboring routers are considered adjacent after they have exchanged link state information and synchronized their topology database.

- ▶ Link State Advertisement

Link State Advertisement (LSA) is the unit of data describing the topology of the network and its adjacent routers. LSAs are flooded to other routers after the Hello protocol has established connection.

- ▶ Link state database

Also called the topology database, the link state database contains the link state advertisements that describe the OSPF area. Each router within the OSPF area maintains an identical copy of the link state database.

- ▶ Flooding

Flooding is the OSPF function that distributes link state advertisements and synchronizes the link state database between routers after the network topology has changed.

- ▶ OSPF Hello protocol

The OSPF Hello protocol is used to detect and establish contact with neighboring routers. It dynamically maintains the relationship by periodically sending a Hello packet to all adjacent routers.

- ▶ Non-backbone area

There are several types of non-backbone areas. A non-backbone area is identified by a four-octet area number that is not 0.0.0.0. There is a standard non-backbone area. There are also two special types of non-backbone areas: the Stub area and the Totally Stubby Area.

- ▶ Stub Area

A Stub Area is a non-backbone area that is connected to the backbone area through an Area Border Router. The Stub Area does not receive advertisements about destinations that are in other Autonomous Systems. Such advertisements are called “external link state advertisements” because they refer to Autonomous Systems external to this Autonomous System.

The Stub Area knows only about intra-area destinations within the Stub Area. It knows about the Totally Stubby Area destinations that exist outside the Stub Area. It reaches external destinations through default routes sent to it by the ABR. With smaller link-state databases and smaller routing tables, Stub Areas consume less CPU storage and fewer CPU cycles.

- ▶ Totally Stubby Area

Nodes in a Totally Stubby Area consume even less CPU storage and fewer CPU cycles for OSPF processing, because they maintain knowledge only of the intra-area destinations and the default routes to reach inter-area and external destinations.

Note: We recommend that, if possible, you define z/OS OSPF nodes as members of a Totally Stubby Area in order to reduce the size of the link state database and reduce the CPU cycles required to produce a routing table. If Totally Stubby is not an option, then we recommend that you find other ways to minimize storage and CPU.

For example, you might integrate a mainframe network running OSPF with a router network running Enhanced Interior Gateway Routing Protocol (EIGRP) to take advantage of the filtering capabilities of EIGRP, thus reducing the amount of protocol traffic between the OSPF network and the EIGRP network.

► Designated router

A designated router (DR) is a router on a shared multi-access medium such as a LAN or ATM network. A DR performs most of the OSPF protocol activities for that network, like synchronizing database information and informing members of the broadcast network of changes to the network. The DR must be adjacent to all other routers on the broadcast medium. Every network or subnetwork on a broadcast network must have a DR and preferably a backup designated router (BDR).

Note: We recommend that you define non-z/OS routers attached to z/OS OSPF LAN broadcast networks as the DRs. z/OS CPU utilization is reduced if a non-z/OS router performs the work of the DR.

There is one exception to this rule when dealing with a HiperSockets network. A HiperSockets network is also a broadcast network; however, only z/OS, z/VM, or Linux on System z nodes participate in a HiperSockets network. Therefore, at least some node inside the mainframe must be a DR on a HiperSockets LAN.

Complications can occur if the z/OS node is the DR on a LAN network when parallel interfaces into the LAN over a shared OSA exist. Shared OSAs can route over the shared OSA port without entering the network.

If the packet arrives over the backup interface instead of the primary parallel interface, the recipient discards the packet. The databases at the nodes become corrupted due to missing information, and lost adjacencies can result.

Therefore, we recommend that you not allow z/OS nodes with parallel interfaces and shared LANs to be the DR. If a z/OS node must be the DR, it should be connected to the broadcast medium through a non-shared OSA port.

► Backup designated router (BDR)

The BDR is also adjacent to all other routers on the medium. It listens to DR conversations, and takes over if the DR fails. After the DR fails, the BDR becomes the DR and a new BDR is elected according to the router priority value. The router priority value is between 0 and 127. If you do not want a router to be elected a DR, configure it with a router priority of zero.

► Transit Area

A Transit Area is an area through which the virtual link ends. Remember that virtual links behave like point-to-point links.

Link-state routing

Link-state routing is a concept used in the routing of packet-switched networks. The routers tell every router in the network about its closest neighbors. The entire routing table is not distributed from any router, only the part of the table containing its neighbors. Basically, implementing link-state routing by OSPF uses the following process:

- ▶ Routers identify other routing devices on directly connected networks, and exchange identification information with them.
- ▶ Routers advertise the details of directly connected network links and the cost of those links by exchanging link state advertisements (LSAs) with other routers in the network.
- ▶ Each router creates a link state database based on the link state advertisements, and the database describes the network topology for the OSPF area.
- ▶ All routers in an area maintain an identical link state database.
- ▶ A routing table is constructed from the link state database.

Link state advertisements are normally sent under the following circumstances:

- ▶ When a router discovers a new neighbor has been added to the area network
- ▶ When a connection to a neighbor is unavailable
- ▶ When the cost of a link changes
- ▶ When basic LSA refreshes are transmitted every 30 minutes

Each area has its own topology and has a gateway that connects it to the rest of the network. It dynamically detects and establishes contacts with its neighboring routers by periodically sending Hello packets.

Link-state advertisements (LSAs)

As mentioned previously, OSPF routers exchange one or more link state advertisements with adjacent routers. LSAs describe the state and cost of an individual router's interfaces that are within a specific area, and the status of an individual network component.

There are five types of LSAs:

- ▶ Router LSAs (Type-1) describe the state and cost of the routers' interfaces within the area. They are generated by *every* OSPF router and are flooded throughout the area.
- ▶ Network LSAs (Type-2) describe all routers attached to the network. They are generated by the *designated* router and are flooded through the area.
- ▶ Summary LSAs (Type-3) describe routes to destinations in other areas in the OSPF network. They are generated by an *area border* router.
- ▶ Summary LSAs (Type-4) are also generated by an *area border* router and describe routes to an AS boundary router.
- ▶ AS External LSAs (Type-5) describe routes to destinations outside the OSPF network. They are generated by an *AS boundary* router.

Link-state database

The link-state database is a collection of OSPF Link State Advertisements. OSPF, being a dynamic IP routing protocol, does not need to have routes defined to it. It dynamically discovers all the routes and the attached routers through its Hello part of the protocol. The OSPF Hello part of the protocol transmits Hello packets to all its router neighbors to establish connection. After the neighbors have been discovered, the connection is made.

Before the link state databases are exchanged, however, the OSPF routers transmit only their LSA headers. After receiving the LSA headers, they are examined for any corruptions. If everything is fine, the request for the most recent LSAs is made. This process is bidirectional between routers.

After the Hello protocol has concluded that all the connections have been established, the link state databases are synchronized. This exchange is performed starting with the most recently updated LSAs. The link state databases are synchronized until all router LSAs in the network (within an area) have the same information. The link state protocol maintains a loop-free routing because of the synchronization of the link state databases.

Physical network types

OSPF supports a combination of different physical networks. In this section, we give a brief description of each physical network and how OSPF supports them.

- Point-to-point

This refers to a network that connects two routers together. A PPP serial line that connects two routers is an example of a point-to-point network.

- Point-to-multipoint

This refers to networks that support more than two attached routers with no broadcast capabilities. These networks are treated as a collection of point-to-point links. OSPF does not use designated routers on point-to-multipoint networks. The Hello protocol is used to detect the status of the neighbors.

- Broadcast multiaccess

This refers to networks that support more than two attached routers and are capable of addressing a single message to all the attached routers. OSPF's Hello Protocol discovers the adjacent routers by periodically sending and receiving Hello packets. This is a typical example of how OSPF exploits a broadcast network. OSPF utilizes multicast in a broadcast network if implemented.

- Nonbroadcast multiaccess (NBMA)

This refers to networks that support more than two attached routers, but have no broadcast capabilities. Because NBMA does not support multicasting, the OSPF Hello packets must be specifically addressed to each router. And because OSPF cannot discover its neighbors through broadcasting, more configuration is required: all routers attached to the NBMA network must be configured. These routers must be configured whether or not they are eligible to become designated routers.

5.3.2 Routing Information Protocol

This section provides an overview of the Routing Information Protocol (RIP) protocol. RIP is designed to manage relatively small networks.

RIP uses a hop count (distance vector) to determine the best possible route to a network or host. The hop count is also known as the *routing metric*, or *the cost of the route*. A router is defined as being zero hops away from its directly connected networks, one hop away from networks that can be reached through one gateway, and so on. The fewer hops, the better.

The route that has the fewest hops will be the preferred path to a destination. A hop count of 16 means infinity, or that the destination cannot be reached. Thus, very large networks with more than 15 hops between potential partners cannot make use of RIP.

The information is kept in a distance vector table, which is periodically advertised to each neighboring router. The router also receives updates from neighboring gateways and uses these to update its routing tables. If an update is not received for three minutes, a gateway is assumed to be down, and all routes through that gateway are set to a metric of 16 (infinity).

Basic distance vector algorithm

The following procedure is carried out by every entity that participates in the RIP routing protocol. This must include all of the gateways in the system. Hosts that are not gateways can participate as well.

- ▶ Keep a table with an entry for every possible destination in the system. The entry contains the distance D to the destination, and the first gateway G on the route to the network.
- ▶ Periodically, send a routing update to every neighbor. The update is a set of messages that contains all the information from the routing table. It contains an entry for each destination, with the distance shown to that destination.
- ▶ When a routing update arrives from the neighbor G' , add the metric associated with the network that is shared with G' . Call the resulting distance D' . Compare the resulting distance with the current routing table entries.

If the new distance D' for N is smaller than the existing value D , then adopt the new route. That is, change the table entry for N to have metric D' and gateway G' . If G' is the gateway from which the existing route came, $G' = G$, then use the new metric, even if it is larger than the old one.

RIP Version 1

RIP is a protocol that manages IP routing table entries dynamically. The gateways using RIP exchange their routing information in order to allow the neighbors to learn of topology changes. The RIP server updates the local routing tables dynamically, resulting in current and accurate routing tables. The protocol is based on the exchange of protocol data units (PDUs) between RIP servers (such as OMPROUTE).

There are various types of PDUs, but the two most important PDUs are:

- | | |
|---------------------|---|
| REQUEST PDU | This PDU is sent from a RIP server as a request to other RIP servers to transmit their routing tables immediately. |
| RESPONSE PDU | This PDU is sent from a RIP server to other RIP servers either as a response to a REQUEST PDU or as a result of expiration of the broadcast timer (every 30 seconds). |

RIP V1 limitations

Because RIP is designed for a specific network environment, it has some limitations as described here. Consider these limitations before implementing RIP in your network.

- ▶ RIP V1 declares a route invalid if it passes through 16 or more gateways. Therefore, RIP V1 places a limitation of 15 hops on the size of a large network.
- ▶ RIP V1 uses fixed metrics to compare alternative routes versus actual parameters, such as measured delay, reliability, and load. This means that the number of hops is the only parameter that differentiates a preferred route from non-preferred routes.
- ▶ The routing tables can take a relatively long time to converge or stabilize.
- ▶ RIP V1 does not support variable subnet masks or variable subnetting because it does not pass the subnet mask in its routing advertisements. *Variable subnet masking* refers to the capability of assigning different subnet masks to interfaces that belong to the same Class A, B, or C network.

- ▶ RIP V1 does not support discontinuous subnets. Discontinuous subnets are built when interfaces belong to the same Class A, B, or C network, but to different subnets that are not adjacent to each other. Rather, they are separated from each other by interfaces that belong to a different network.

With RIP version 1, discontinuous subnets represent unreachable networks. If you find it necessary to build discontinuous subnets, you must use one of the following techniques:

- An OSPF implementation
- RIP version 2 protocol
- Static routing

RIP Version 2

Rather than being another protocol, RIP V2 is an extension to the functions provided by RIP V1. To use these new functions, RIP V2 routers exchange the same RIP V1 messages. The version field in the message will specify version number 2 for RIP messages that use authentication or carry information in any of the newly defined fields.

RIP V2 protocol extensions provide features such as:

- ▶ Route tags to provide EGP-RIP and BGP-RIP implementation
Route tags are used to separate *internal* RIP routes (routes for networks within the RIP routing domain) from *external* RIP routes, which might have been imported from an EGP (external gateway protocol) or another IGP. OMROUTE does not generate route tags, but preserves them in received routes and readvertises them when necessary.
- ▶ Variable subnetting support
Variable length subnet masks are included in routing information so that dynamically added routes to destinations outside subnetworks or networks can be reached.
- ▶ Immediate next hop for shorter paths
Next hop IP addresses, whenever applicable, are included in the routing information. Their purpose is to eliminate packets being routed through extra hops in the network. OMROUTE will not generate immediate next hops, but will preserve them if they are included in RIP packets.
- ▶ Multicasting to reduce load on hosts
An IP multicast address 224.0.0.9, reserved for RIP version 2 packets, is used to reduce unnecessary load on hosts that are not listening to RIP version 2 messages. RIP version 2 multicasting is dependent on interfaces that are multicast-capable.
- ▶ Authentication for routing update security
Authentication keys can be included in outgoing RIP version 2 packets for authentication by adjacent routers as a routing update security protection. Likewise, incoming RIP version 2 packets are checked against local authentication keys. The authentication keys are configurable on a router-wide or per-interface basis.
- ▶ Configuration switches for RIP V1 and RIP V2 packets
Configuration switches are provided to selectively control which versions of RIP packets are to be sent and received over network interfaces. You can configure them router-wide or per-interface.
- ▶ Supernetting support
The supernetting feature is part of the Classless InterDomain Routing (CIDR) function. Supernetting provides a way to combine multiple network routes into fewer supernet routes. Therefore, the number of network routes in the routing tables becomes smaller for advertisements. Supernet routes are received and sent in RIP V2 messages.

RIP V2 packets are backward compatible with existing RIP V1 implementations. A RIP V1 system will process RIP V2 packets but without the RIP V2 extensions, and broadcast them as RIP V1 packets to other routers. Note that routing problems might occur when variable subnet masks are used in mixed RIP V1 and RIP V2 systems. RIP V2 is based on a distance vector algorithm, just as RIP V1 is.

5.3.3 IPv6 dynamic routing

Dynamic routing in a IPv6 network can be implemented in a z/OS Communications Server in two different ways:

- ▶ IPv6 dynamic routing using router discovery
- ▶ IPv6 dynamic routing using OMPROUTE

IPv6 dynamic routing using router discovery

Enabling IPv6 router discovery in the z/OS Communications Server requires no additional z/OS Communications Server configuration. All that is needed is at least one IPv6 interface that is defined and started, and at least one adjacent router through that interface that is configured for IPv6 router discovery. If these things exist, then the z/OS Communications Server begins receiving router advertisements from the adjacent routers.

Depending on the configuration in the adjacent routers, the following types of routes can be learned from the received router advertisements:

- ▶ Default route, for which the originator of the router advertisement is the next hop
- ▶ Direct routes (no next hop) to prefixes that reside on the link shared by the z/OS Communications Server and the originator of the router advertisement

IPv6 dynamic routing using OMPROUTE

For IPv6, OMPROUTE implements the IPv6 RIP protocol described in RFC 2080 (RIPng for IPv6) and the IPv6 OSPF protocol described in RFC 2740 (OSPF for IPv6). It provides an alternative to the static TCP/IP gateway definitions.

The z/OS host running with OMPROUTE becomes an active OSPF or RIP router in a TCP/IP network. Either or both of these routing protocols can be used to dynamically maintain the host IPv6 routing table. For example, OMPROUTE can detect when a route is created, is temporarily unavailable, or if a more efficient route exists. If both IPv6 OSPF and IPv6 RIP protocols are used simultaneously, then IPv6 OSPF routes will be preferred over IPv6 RIP routes to the same destination.

RIPng or RIP next generation

RIP Next Generation (RIPng) is a distance vector routing protocol for IPv6 that is defined in RFC 2080. RIPng for IPv6 is an adaptation of the RIP V2 protocol to advertise IPv6 network prefixes. RIPng for IPv6 uses UDP port 521 to periodically advertise its routes, respond to requests for routes, and advertise route changes.

RIPng for IPv6, like other distance vector protocols, has a maximum distance of 15, in which 15 is the accumulated cost (hop count). Locations that are a distance of 16 or further are considered unreachable. RIPng for IPv6 is a simple routing protocol with a periodic route-advertising mechanism designed for use in small to medium-sized IPv6 networks. RIPng for IPv6 does not scale well to a large or very large IPv6 network.

Differences between RIPng and RIP-2

There are two important distinctions between RIP-2 and RIPng:

- ▶ Support for authentication

The RIP-2 standard includes support for authenticating a node transmitting routing information. RIPng does not include any native authentication support. Rather, RIPng uses the security features inherent in IPv6.

In addition to authentication, these security features provide the ability to encrypt each RIPng packet. This can control the set of devices that receive the routing information.

One consequence of using IPv6 security features is that the AFI field within the RIPng packet is eliminated. There is no longer a need to distinguish between authentication entries and routing entries within an advertisement.

- ▶ Support for IPv6 addressing formats

The fields contained in RIPng packets were updated to support the longer IPv6 address format.

OSPF for IPv6

OSPF for IPv6 is a link state routing protocol defined in RFC 2740 and designed for routing table maintenance within a single autonomous system. OSPF for IPv6 is an adaptation of the OSPF routing protocol version 2 for IPv4 defined in RFC 2328.

IPv6 OSPF is classified as an Interior Gateway Protocol (IGP). This means that it distributes routing information between routers belonging to a single autonomous system (AS), a group of routers all using a common routing protocol. The IPv6 OSPF protocol is based on link-state or shortest path first (SPF) technology.

At a glance, the OSPF implementation is basically the same as it is for IPv4, except for some primary differences.

Primary differences between IPv6 OSPF and IPv4 OSPFv2

IP addressing and topology semantics have been separated where possible (many LSAs do not carry IP addresses at all, only abstract topology information). Removing IP addressing from the topology description makes OSPFv3 more protocol-independent.

New LSA types are added (to carry addressing and link-local information). Because IP addressing has been removed from some of the basic LSA types, new LSA types are provided to communicate IP addresses, which routers then correlate to topology information in other LSA types.

The Concept of Flooding Scope is added (scopes are link, area, autonomous system). It indicates how far an advertisement can be flooded. For example, link scope means an LSA can only be flooded on the originating link.

Support for Unknown LSA types is added (this makes the protocol more extensible). Unknown LSA types can be ignored, or they can be stored and forwarded by the router, depending on the settings of bits in the LSA type field. This vastly improves interoperability between routers running different versions of the protocol. For example, a designated router could conceivably have a lower level of support than another router on the same link; because the designated router floods on behalf of the other routers on the link, it could store and forward unknown LSA types received from its peers.

Multiple OSPF instances are supported on a link. An "instance id" field is added to OSPF headers, and OSPF processes only process packets whose instance ID matches their own.

This opens up the possibility of one link belonging to completely different autonomous systems.

Subnet loses its importance, replaced by *link* (because multiple IPv6 prefixes per link are allowed and expected, routing by subnet/prefix makes less sense). In OSPFv2, most routing is done by subnet. In OSPFv3 it is done by link. This is because in IPv6 a subnet (prefix) does not always uniquely identify a link, and a link can have more than one prefix assigned.

5.4 Implementing static routing in z/OS

In this section we implement a static routing scenario, as illustrated in Figure 5-2. We only provide definition examples for the TCPIPA stack on SC30 because the examples for the TCPIPB stack on SC31 are similar. On TCPIPA, we define direct routes for interfaces such as OSA and HiperSockets and indirect routes for TCPIPB VIPAs. We also define default routes through our switches (layer 3 switch).

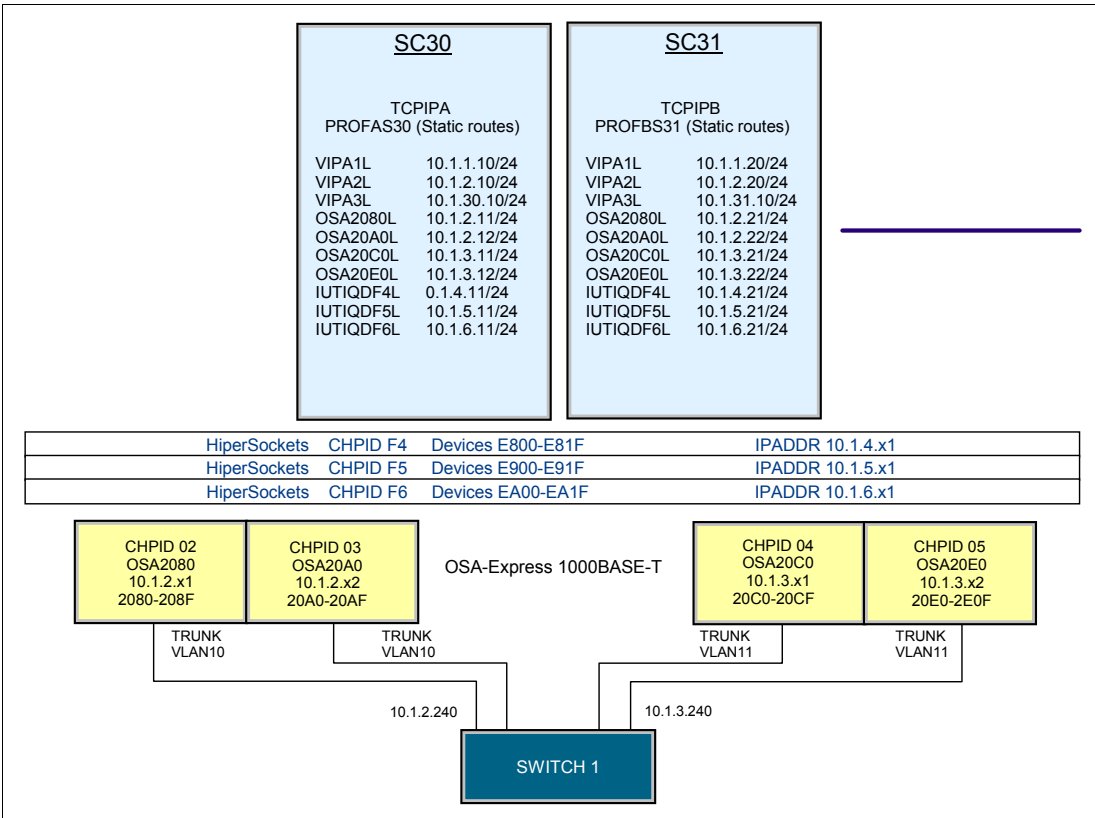


Figure 5-2 Static routing scenario

5.4.1 Dependencies

All subnetworks defined in the TCP/IP stack that are used by the application servers, including static and dynamic VIPAs, must also have static routing definitions in the routers. In our case, the layer 3 switches (routers) do not need static route definitions for direct routes. We defined indirect routes for TCPIPA and TCPIPB VIPAs in the routers.

5.4.2 Considerations

When planning to design a static routing environment on a z/OS Communications Server, you need to address some issues. Keep in mind that static routes are configured manually for each router by the system administrator. Available network interfaces and routes through the network must be determined before the routes are configured. Except for potential ICMP router redirections, routers do not communicate with each other about the topology of the network.

The routing table's management is manual, thus increasing the possibility of outages caused by definition errors. If a destination (sub)network becomes unreachable, then the static routes for that (sub)network will remain in the routing table, and packets will still be forwarded to the destination. The only way to remove static routes from the routing table is for the network administrator to update the routing table.

We recommend that you define as few static routing definitions as possible when implementing a static routing environment, keeping in mind that our z/OS system is basically an application server environment. It is good practice to define only the default gateways to the exterior networks, and let the routers do the exterior routing. You can implement the static routing between the z/OS system and external router, and still let the external router use the dynamic routing protocol.

In the router, we recommend that you define only the route definitions to the VIPA subnetworks. The interior subnetworks, such as XCF and HiperSockets, do not usually need to be reached by the corporate network, so they do not need to be defined.

Attention: If you choose to implement the OSA Connection Isolation feature together with dynamic routing and yet still need to communicate between two or more nodes sharing the same OSA adapter port, you will need to override the dynamically generated subnet or host route between the two TCP/IP stacks with a non-replaceable static route that indicates a next-hop address of an external router. Refer to the discussion of OSA connection isolation in "Considerations for assigning the OSA portname" on page 147.

5.4.3 Implementation tasks

To implement the static routing scenario, follow these steps:

1. Update the TCP/IP profile.
2. Configure the router.

Update the TCP/IP profile

In the TCP/IP profile, use the BEGINROUTES block and ROUTE statement to define the following routes:

- ▶ A direct route to all local interfaces (except static VIPAs, dynamic VIPAs, or XCF)
To define a direct route, specify = for its First Hop.
- ▶ An indirect route to the subnetwork
To define a direct route, specify the IP address of the next hop router for its first hop.
- ▶ Default gateway statements to route all packets being sent to unknown destinations

Example 5-6 shows our definition example.

When multiple default routes are defined, the traffic will be sent to the first default route defined. If the MULTIPATH parameter is specified on the IPCONFIG statement, then all default routes will be used.

Example 5-6 Direct routes configuration

```

; *****
; TCPIPA.TCPPARMS(PROFA30S)
; *****
.....
BEGINRoutes
; Direct Routes - Routes that are directly connected to my interfaces
;   Destination      Subnet Mask  First Hop Link Name      Packet Size
ROUTE 10.1.2.0        255.255.255.0 = 1      OSA2080L      MTU 1492
ROUTE 10.1.2.0/24      = 1      OSA20A0L      MTU 1492
ROUTE 10.1.3.0/24      = 1      OSA20C0L      MTU 1492
ROUTE 10.1.3.0/24      = 1      OSA20E0L      MTU 1492
ROUTE 10.1.4.0/24      = 2      IUTIQDF4L      MTU 8192
ROUTE 10.1.5.0/24      = 2      IUTIQDF5L      MTU 8192
ROUTE 10.1.6.0/24      = 2      IUTIQDF6L      MTU 8192
;
; Indirect Routes - Routes that are not directly connected to my interfaces
;   Destination      Subnet Mask  First Hop  Link Name  Packet Size ;
ROUTE 10.1.1.20/32      10.1.4.21  IUTIQDF4L  MTU 8192
ROUTE 10.1.2.20/32      10.1.4.21  IUTIQDF4L  MTU 8192
ROUTE 10.1.31.10/32     10.1.4.21  IUTIQDF4L  MTU 8192
ROUTE 10.1.100.0/24     10.1.2.240 3 OSA2080L  MTU 1492
ROUTE 10.1.100.0/24     10.1.2.240 3 OSA20A0L  MTU 1492
ROUTE 10.1.100.0/24     10.1.3.240 3 OSA20C0L  MTU 1492
ROUTE 10.1.100.0/24     10.1.3.240 3 OSA20E0L  MTU 1492
;
; Default Routes - Routes directly connected to my interfaces
;   Destination      Subnet Mask  First Hop  Link Name  Packet Size ;
ROUTE DEFAULT          10.1.2.240 4 OSA2080L  MTU 1492
ROUTE DEFAULT          10.1.2.240 4 OSA20A0L  MTU 1492
ROUTE DEFAULT          10.1.3.240 4 OSA20C0L  MTU 1492
ROUTE DEFAULT          10.1.3.240 4 OSA20E0L  MTU 1492
;
ENDROUTES

```

In this example, the numbers correspond to the following information:

- 1.** Define the direct routes for OSA interfaces. Specify the subnet mask with decimal format (such as 255.255.255.0) or the *prefix* length.
- 2.** Define the direct routes for HiperSockets interfaces.
 Note that the first hop parameter is defined as an equal sign (=) **1** to identify this as a direct route.
- 3.** Define the indirect routes to reach the external network. The next hop is router 1 (10.1.2.240 and 10.1.3.240).
- 4.** Define the default routes, to reach the external network, which are not explicitly defined as indirect routes. The next hop is router 1 (10.1.2.240 and 10.1.3.240).

Configure the router

Define the static routes to the VIPA or the physical interfaces which are not on the subnet that the routers are directly connected to (HiperSockets, for example). In our example, 10.1.2.0/24 and 10.1.3.0/24 are direct routes of the router, and we do not need to define the static routes for those subnets.

Example 5-7 shows the example of router (layer 3 switch) configuration.

Example 5-7 Static route definition in router

```
.....
ip route 10.1.1.10 255.255.255.255 10.1.2.11 1
ip route 10.1.2.10 255.255.255.255 10.1.2.11 1
ip route 10.1.30.10 255.255.255.255 10.1.2.11 1
.....
```

In this example, the number corresponds to the following information:

1. Define the static route to the static VIPA in TCPIPA. The next hop address is the IP address of the OSA physical interface. In our example we define this static route with 32-bit mask (255.255.255.255), but you can use a mask length shorter than 32-bit.

5.4.4 Activation and verification

To activate and verify the static routing scenario, follow these steps:

1. Apply changes to TCP/IP profile.
2. Verify the connectivity.

Apply changes to TCP/IP profile

To apply the changes to static routes, do *one* of the following:

- ▶ Restart the TCP/IP stack.
- ▶ Modify the TCP/IP definition with VARY TCPIP,*procname*,OBEYFILE command.

After you perform one of these tasks, then all static routes are listed in the TCP/IP routing table.

Example 5-8 illustrates applying changes by using the OBEYFILE command.

Important: When using the OBEYFILE command, include all static routes that you want to define. The OBEYFILE command replaces the entire BEGINROUTES block.

Example 5-8 Applying changes with the OBEYFILE command

```
V TCPIP,TCPIPA,0,DSN=TCPIPA.TCPPARMS(PROFA30S)
EZZ0060I PROCESSING COMMAND: VARY TCPIP,TCPIPA,0,DSN=TCPIPA.TCPPARMS(P
ROFA30S)
EZZ0300I OPENED OBEYFILE FILE 'TCPIPA.TCPPARMS(PROFA30S)'
EZZ0309I PROFILE PROCESSING BEGINNING FOR 'TCPIPA.TCPPARMS(PROFA30S)'
.....
EZZ0316I PROFILE PROCESSING COMPLETE FOR FILE 'TCPIPA.TCPPARMS(PROFA30S)'
.....
```

Verify the connectivity

To verify if the static routing table is built as expected, the following commands are useful.

Note: `netstat` commands can be executed as TSO commands, z/OS UNIX shell commands, or Display commands on the system console. Our examples are the result of Display commands on the system console, but their output is identical to the TSO and z/OS UNIX shell output.

Display the device status

Use the `D TCPIP,TCPIPA,Netstat,DEVlink` command to review the status of all devices defined in the TCP/IP environment. If a device is not ready, there will be no routing through this device. Example 5-9 shows the resulting display of this command.

Example 5-9 netstat DEVlink command display

```
D TCPIP,TCPIPA,N,DEV
DEVNAME: OSA2080          DEVTYPE: MPCIPA
DEVSTATUS: READY 1
LNKNAME: OSA2080L         LNKTYPE: IPAQENET  LNKSTATUS: READY 1
NETNUM: N/A  QUESIZE: N/A  SPEED: 0000001000
IPBROADCASTCAPABILITY: NO
CFGROUTER: NON           ACTROUTER: NON
ARPOFFLOAD: YES          ARPOFFLOADINFO: YES
ACTMTU: 8992
VLANID: 10               VLANPRIORITY: DISABLED
DYNVLANREGCFG: NO        DYNVLANREGCAP: YES
READSTORAGE: GLOBAL (4096K)  INBPERF: BALANCED
CHECKSUMOFFLOAD: YES
SECCLASS: 255            MONSYSPLEX: NO
BSD ROUTING PARAMETERS:
MTU SIZE: 1492           METRIC: 100
DESTADDR: 0.0.0.0        SUBNETMASK: 255.255.255.0
MULTICAST SPECIFIC:
MULTICAST CAPABILITY: YES
GROUP          REFCNT          SRCFLTMD
-----
224.0.0.1      0000000001    EXCLUDE
SRCADDR: NONE
LINK STATISTICS:
BYTESIN                    = 2492
INBOUND PACKETS            = 14
INBOUND PACKETS IN ERROR   = 0
INBOUND PACKETS DISCARDED  = 0
INBOUND PACKETS WITH NO PROTOCOL = 0
BYTESOUT                   = 536
OUTBOUND PACKETS           = 6
OUTBOUND PACKETS IN ERROR  = 0
OUTBOUND PACKETS DISCARDED = 0
.....
```

In this example, the number corresponds to the following information:

- 1.** Make sure the DEVSTATUS and LNKSTATUS are both READY.

Display routing table

Use the **netstat ROUTe** command to display the routing table in a TCP/IP stack. A sample of the command is shown in Example 5-10.

Example 5-10 *netstat ROUTe* resulting display

```
D TCPIP,TCPIPA,N,ROUTE
IPV4 DESTINATIONS
DESTINATION      GATEWAY      FLAGS      REFCNT      INTERFACE
DEFAULT          10.1.2.240   UGS 4      000000     OSA2080L 1
DEFAULT          10.1.2.240   UGS        000001     OSA20A0L
DEFAULT          10.1.3.240   UGS        000000     OSA20C0L
DEFAULT          10.1.3.240   UGS        000000     OSA20E0L
10.1.1.10/32      0.0.0.0      UH         000000     VIPA1L
10.1.1.20/32      10.1.4.21    UGHS       000000     IUTIQDF4L 2
10.1.2.0/24       0.0.0.0      US         000000     OSA2080L 3
10.1.2.0/24       0.0.0.0      US         000000     OSA20A0L
10.1.2.10/32      0.0.0.0      UH         000000     VIPA2L
10.1.2.11/32      0.0.0.0      UH         000000     OSA2080L
10.1.2.12/32      0.0.0.0      UH         000000     OSA20A0L
10.1.3.0/24       0.0.0.0      US         000000     OSA20C0L
10.1.3.0/24       0.0.0.0      US         000000     OSA20E0L
10.1.3.11/32      0.0.0.0      UH         000000     OSA20C0L
10.1.3.12/32      0.0.0.0      UH         000000     OSA20E0L
10.1.100.0/24     10.1.2.240   UGS        000000     OSA2080L
10.1.100.0/24     10.1.2.240   UGS        000000     OSA20A0L
10.1.100.0/24     10.1.3.240   UGS        000000     OSA20C0L
10.1.100.0/24     10.1.3.240   UGS        000000     OSA20E0L
127.0.0.1/32      0.0.0.0      UH         000001     LOOPBACK
IPV6 DESTINATIONS
DESTIP:  ::1/128
GW:      ::
      INTF:  LOOPBACK6      REFCNT:  000000
      FLGS:  UH             MTU:  65535
41 OF 41 RECORDS DISPLAYED
END OF THE REPORT
```

In this example, the numbers correspond to the following information:

- 1.** The default route is defined. If there are multiple default route entries as shown in the example, only the first active entry (interface OSA2080L) is used.
- 2.** The indirect route to VIPA in TCPIPB is defined.
- 3.** The direct route for OSA physical interface is defined.
- 4.** *S* in the FLAG field stands for non-replaceable static route entry. For replaceable static route entries, FLAG *Z* would be displayed.

Check the connectivity using PING command

The PING command can be executed using the TSO PING command or the z/OS UNIX **ping** command. Example 5-11 on page 233 shows the display of the TSO PING command; the ping is successful.

In a CINET environment where multiple TCP/IP stacks are configured, use the TCP option for the TSO PING command and the **-p** option for the z/OS UNIX **ping** command to specify the TCP/IP stack name from which you want to issue the **ping** command.

You do not need to specify these options if the user issuing this command is already associated to the TCP/IP stack (with SYSTCPD DD, for example).

You do not need to specify these options if your environment is an INET environment where only one TCP/IP stack is configured.

Example 5-11 TSO PING command display

```
TSO PING 10.1.1.20 (TCP TCPIPA
CS V1R12: Pinging host 10.1.1.20
Ping #1 response took 0.000 seconds.
***
```

Example 5-12 shows the display of z/OS UNIX **ping** command.

Example 5-12 z/OS UNIX ping command display

```
CS02 @ SC30:/u/cs02>ping -p tcpipa 10.1.1.20
CS V1R12: Pinging host 10.1.1.20
Ping #1 response took 0.000 seconds.
```

Verify the selected route with TRACEROUTE command

TRACEROUTE can be invoked by either the TSO TRACERTE command or the z/OS UNIX shell **tracert**/**otracert** command. Example 5-13 shows the example of the display. We see the router 1 (10.1.2.240) is the next hop router to reach the destination IP address 10.1.100.221.

In a CINET environment where multiple TCP/IP stacks are configured, use the TCP option for TSO TRACERTE command and the **-a** option for the z/OS UNIX **tracert** command to specify the TCP/IP stack name you want to issue the TRACEROUTE command from.

You do not need to specify these options if the user issuing this command is already associated to the TCP/IP stack (with SYSTCPD DD, for example).

You do not need to specify these options if your environment is an INET environment where only one TCP/IP stack is configured.

Example 5-13 tracert command results

```
CS02 @ SC30:/u/cs02>otracert 10.1.100.221
CS V1R12: Traceroute to 10.1.100.221 (10.1.100.221)
1 router1 (10.1.2.240) 0 ms 0 ms 0 ms
2 10.1.100.221 (10.1.100.221) 0 ms 0 ms 0 ms
***
```

5.5 Implementing OSPF routing in z/OS with OMPROUTE

In this scenario we show a dynamic routing implementation. In our example we configure OSPF for lesser network load, more IP network design flexibility, and lower convergence time compared to RIP v1 and RIP v2 (see Table 5-3 on page 211). Although OSPF requires higher CPU processing, we can reduce that requirement by making the z/OS Communications Server a part of the OSPF Stub Area or Totally Stubby Area.

Figure 5-3 depicts the environment we use for the OSPF scenario. The TCPIPA stack is running on SC30. We create the OMPROUTE procedure OMPA to establish affinity to TCPIPA. Likewise, we create OMPB for TCPIPB on SC31.

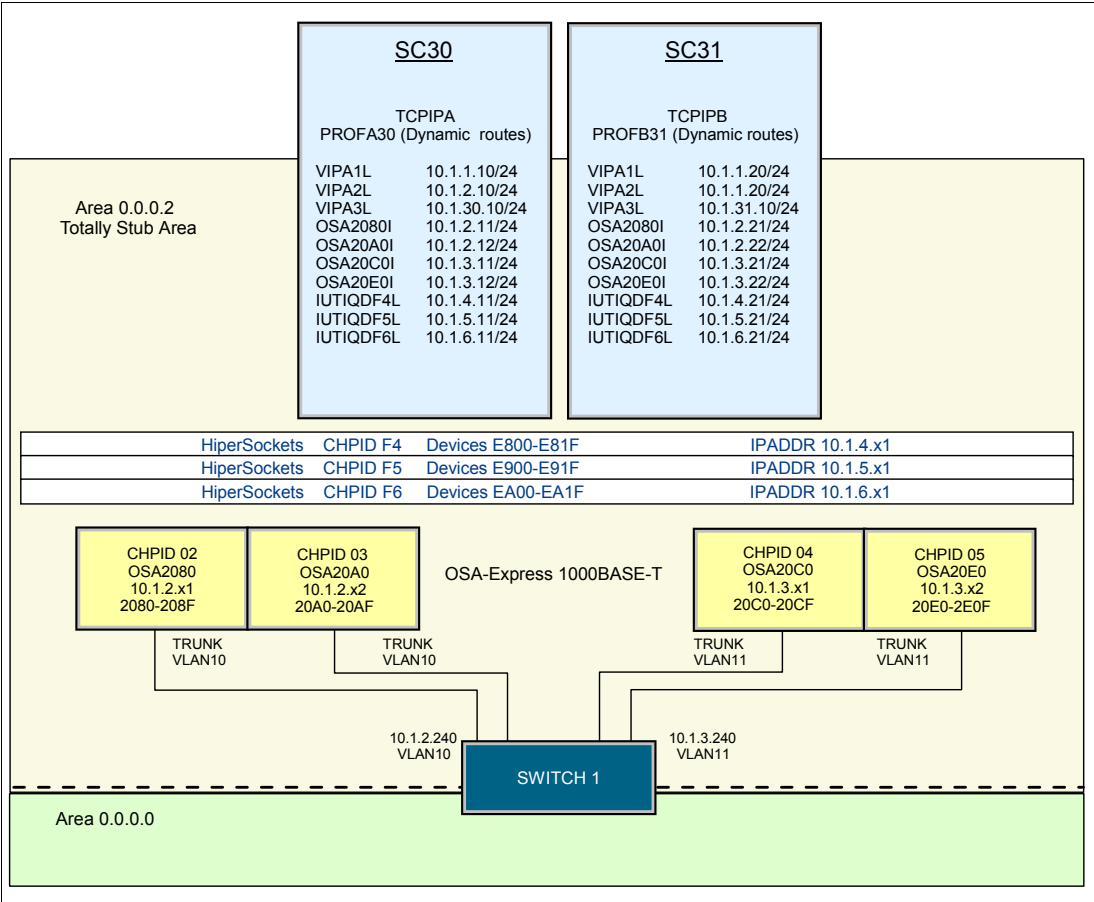


Figure 5-3 Dynamic routing scenario using OSPF

We define a z/OS TCP/IP to be a member of OSPF Totally Stubby Area. The external routers (layer 3 switches) represent the ABRs between the Totally Stubby Area and the backbone area. We made the external routers to be DR or BDR to reduce the routing workloads required in the z/OS.

Because the configuration examples for TCPIPB and OMPB on SC31 are similar to those examples for TCPIPA and OMPA on SC30, we only show the configuration examples on SC30.

5.5.1 Dependencies

The IP routers that will be involved in establishing access to the external network must support OSPF, and the configuration parameters set in OMPROUTE must be consistent with those defined to the IP routers.

5.5.2 Considerations

A z/OS Communications Server host is usually used as an application server and the routing daemon is running primarily to provide access to network resources and vice versa. With this

in mind, take care to ensure that the z/OS Communications Server host is not overly burdened with routing work.

The z/OS Communications Server should *not* be configured as a backbone router, either intentionally or inadvertently. Careful network design can minimize the routing burdens on the z/OS Communications Server (application host), without compromising accessibility.

5.5.3 Recommendations

We recommend that you define the z/OS Communications Server environment as an OSPF Stub Area to reduce the CPU process needed for managing the routing table. A Stub Area can be configured so that route summaries from other areas are not flooded into the Stub Area by the area border routers. When this is done, only routes to destinations within the Stub Area are shared among the hosts. Default routes are used to represent all destinations outside the Stub Area. The Stub Area's resources are still advertised to the network at large by the area-border routers. You can use this optimization, sometimes referred to as a Totally Stubby Area.

We also recommend that you make the external routers be DR or BDR, and do not allow z/OS systems to be DR or BDR, in order to reduce the routing burden for z/OS systems. DR or BDR is selected in each LAN segment or VLAN. However, on HiperSockets links, z/OS systems are the only participants. One of the z/OS on the HiperSockets network has to take the role of DR (optionally, another one can take the role of BDR).

Note: Recall our earlier warning in this chapter about the use of OSA Connection Isolation: It is generally incompatible with a dynamic routing protocol like OSPF. If implemented, you may need to introduce non-replaceable static routes pointing to external next-hop routers. Refer to the discussion of OSA connection isolation in "Considerations for assigning the OSA portname" on page 147.

5.5.4 Implementation tasks

To implement and configure OMPROUTE in the z/OS Communications Server, follow these steps:

1. Create the OMPROUTE cataloged procedure.
2. Define the OMPROUTE environment variables.
3. Update the TCPIP.DATA file.
4. RACF-authorize user IDs for starting OMPROUTE.
5. Start syslogd.
6. Change port 520 and 521 definitions to NOAUTOLOG.
7. Create the OMPROUTE configuration file.
8. Configure routers.

In the sections that follow, we show only the configuration examples for the TCPIPA stack and omit the examples for the TCPIPB stack. We do not define any static routes in TCP/IP profile in conjunction with the dynamic routing.

Create the OMPROUTE cataloged procedure

We create the OMPROUTE cataloged procedure by copying the sample in hlq.SEZAINST(OMPROUTE) to our PROCLIB. We specify the STDENV file name and OMPCFG file name, as shown in Example 5-14.

Example 5-14 OMPROUTE cataloged procedure

```
//OMPA30 PROC STDENV=OMPEN&SYSCONE 1  
//OMPA30 EXEC PGM=OMPROUTE,REGION=OM,TIME=NOLIMIT,  
//          PARM=(' POSIX(ON) ALL31(ON) ',  
//          'ENVAR("_BPXK_SETIBMOPT_TRANSPORT=TCPIPA"',  
//          '" _CEE_ENVFILE=DD:STDENV")/') 2  
//STDENV DD DISP=SHR,DSN=TCPIP.SC&SYSCONE..STDENV(&STDENV)  
//SYSOUT DD SYSOUT=*  
//OMPCFG DD DSN=TCPIPA.TCPPARMS(OMPA&SYSCONE.),DISP=SHR 3  
//CEEDUMP DD SYSOUT=*,DCB=(RECFM=FB,LRECL=132,BLKSIZE=132)
```

In this example, the numbers correspond to the following information:

1. Specifies the STDENV variable. We can use a common procedure for all images within the same server environment by specifying the &SYSCONE variable. The &SYSCONE value for this LPAR is 30.
2. Each OMPROUTE procedure in the same server will have its own environment variables based on this DD.
3. The OMPCFG DD card permits you to specify the OMPROUTE configuration file within the JCL. The DD card enables the use of an MVS system symbol that can make the procedure shareable across TCP/IP stacks. If you specify the configuration file here, you can omit the statement OMPROUTE_FILE from the STDENV file.

Tip: OMPROUTE can be started as a z/OS procedure, or from the z/OS shell, or from AUTOLOG.

Define the OMPROUTE environment variables

To define our OMPROUTE environment variables we use a STDENV file, pointed to by the STDENV DD statement in our OMPROUTE procedure. Example 5-15 shows the STDENV file we use in our example.

Example 5-15 OMPROUTE environment variables

```
; *****  
; TCPIP.SC30.STDENV(OMPEN&SYSCONE)  
; *****  
RESOLVER_CONFIG=//'TCPIPA.TCPPARMS(DATA&SYSCONE.)' 1  
;OMPROUTE_FILE=//'TCPIPA.TCPPARMS(OMPA30)' 2  
OMPROUTE_DEBUG_FILE=/etc/omproute/debuga30  
OMPROUTE_DEBUG_FILE_CONTROL=100000,5
```

In this example, the numbers correspond to the following information:

1. Specify the TCPIP.DATA file. The &SYSCONE value for this LPAR is 30. If you do not want to use MVS system symbols, you can define the hard-coded member name as shown:

RESOLVER_CONFIG=//'TCPIPA.TCPPARMS(DATA&SYSCONE30)'
2. We can omit the OMPROUTE_FILE statement if we have coded the OMPCFG DD statement in the OMPROUTE started procedure.

With the appropriate naming conventions, we can make both the OMPROUTE environment variable file and the OMPROUTE started procedure shareable across multiple TCP/IP stacks.

Important: When defining the STDENV (_CEE_ENVFILE) file with a z/OS data set, the data set must be allocated with RECFM=V. Using RECFM=F or FB is not recommended, because the fixed setting enables padding with blanks for the environment variables.

Although you can include a UNIX time zone variable (TZ=...) in either the JCL or the environment variable file, the recommended procedure is to insert the appropriate time zone for all applications into the z/OS SYS1.PARMLIB(CEEPRMxx) member, as shown in Example 5-16. You should define the TZ environment variable for all three LE option sets (CEEDOPT, CEECOPT, and CELQDOPT).

Example 5-16 Setting the time zone variable for all applications

```
CEECDPT(ALL31(ON), ENVAR('TZ=EST5EDT')) )
CEEDOPT(ALL31(ON), ENVAR('TZ=EST5EDT')) )
CELQDOPT(ALL31(ON), ENVAR('TZ=EST5EDT')) )
```

Update the TCPIP.DATA file

Our test environment is running under CINET. With CINET there is often a global TCPIP.DATA file and a stack-specific local TCPIP.DATA file. The keywords specified in the global TCPIP.DATA cannot be overridden with parameters in any local TCPIP.DATA files.

In the CINET environment, the Global Resolver configuration file contains keywords that are shared with all TCP/IP stacks on the z/OS image, and should omit the stack-specific keywords such as TCPIPJobname and Hostname. Those parameters should be specified in the local TCPIP.DATA file. If a specific parameter is not found in the global TCPIP.DATA, the local TCPIP.DATA file is searched according to the search order. You can read more about the resolver in Chapter 2, “The resolver” on page 19.

Example 5-17 shows the global TCPIP.DATA file used in our example.

Example 5-17 Global TCPIP.DATA file

```
; *****
; TCPIPA.TCPPARMS(GLOBAL)
; *****
DOMAINORIGIN  ITS0.IBM.COM
NSINTERADDR   10.12.6.7
NSPORTADDR    53
RESOLVEVIA    UDP
RESOLVERTIMEOUT 10
RESOLVERUDPRETRIES 1
LOOKUP        LOCAL DNS
```

Then each stack has a stack-specific local TCPIP.DATA file identifying stack-specific parameters such as TCPIPJobname and Hostname. Example 5-18 on page 237 shows the local TCPIP.DATA file used in our example.

Example 5-18 Local TCPIP.DATA file

```
; *****
; TCPIPA.TCPPARMS(DATAA30)
; *****
```

```

TCPIPJOBNAME TCPIPA      1
HOSTNAME WTSC30A
DATASETPREFIX TCPIPA     2
MESSAGECASE MIXED

```

In this example, the numbers correspond to the following information:

- 1.** Specify the TCP/IP stack name that OMPROUTE should establish affinity to, using the TCPIPJobname statement.
- 2.** Specify the data set prefix (hlq) that OMPROUTE should use.

In an INET environment, usually only a global TCPIP.DATA file is used. It should contain the keywords (TCPIPJobname and DATASETPREFIX) used by OMPROUTE. The TCPIPJobname parameter specifies the name of TCP/IP stack with which OMPROUTE establishes an affinity.

RACF-authorize user IDs for starting OMPROUTE

To reduce the risk of an unauthorized user starting OMPROUTE and affecting the contents of the routing table, users who start OMPROUTE must be RACF-authorized to the entity MVS.ROUTE MGR.OMPROUTE and require a UID of zero (0). In our test environment, we executed the command shown in Example 5-19.

Example 5-19 RACF commands to authorize the starting of OMPROUTE

```

RDEFINE OPERCMDS (MVS.ROUTE MGR.OMPROUTE) UACC(NONE)
PERMIT MVS.ROUTE MGR.OMPROUTE ACCESS(CONTROL) CLASS(OPERCMDS) ID(OMPA) 1
SETROPTS RACLIST(OPERCMDS) REFRESH

```

In this example, the numbers correspond to the following information:

- 1.** Specify the OMPROUTE cataloged procedure name for ID parameter.

Important: OMPROUTE must be started by a RACF-authorized user ID.

Start syslogd

Syslogd can and should be used to receive the specified messages from OMPROUTE. It can be configured to receive all OMPROUTE non-critical messages. Update the syslogd.conf file to isolate all OMPROUTE messages to a specific output destination file. Example 5-20 shows how we configured the syslogd to receive all error, warning, info, and notice messages in a syslog file.

Example 5-20 Syslogd configuration file

```

##*****
#*
#* syslog.conf - Defines the actions to be taken for the specified
#*               facilities/priorities by the syslogd daemon.
#*
*.OMPA*.*.err /tmp/syslog/ompa.err.log 1

```

In this example, the numbers correspond to the following information:

1. Specify the syslog output destination file name for the OMPA-related messages.

Change port 520 and 521 definitions to NOAUTOLOG

If OMPROUTE is started with AUTOLOG and only the OSPF protocol is used, it is important to do *one* of the following tasks:

- ▶ Ensure that the RIP UDP port (520) and the IPv6 RIP UDP port (521) are *not* reserved by the PORT statement in the PROFILE.TCPIP.
- ▶ Add the NOAUTOLOG parameter to the PORT statement, as shown in Example 5-21.

Example 5-21 Ports 520 and 521 defined as NOAUTOLOG

```
; *****  
; TCPIPA.TCPPARMS (PROFA30)  
; *****  
.....  
    520 UDP OMPA NOAUTOLOG ; OMPROUTE IPv4 RIPV2  
    521 UDP OMPA NOAUTOLOG ; OMPROUTE IPv6 RIPV2  
.....
```

Important: If you fail to take one of these actions, OMPROUTE will be periodically canceled and restarted by TCP/IP.

Create the OMPROUTE configuration file

We defined the parameters for OSPF implementation in the OMPROUTER configuration file. Example 5-22 shows the configuration we used in our example. We defined a z/OS TCP/IP to be a Totally Stubby Area, the interfaces as part of the Stub Area, and other parameters.

The search order for OMPROUTE configuration file is as follows:

1. OMPCFG DD statement in the OMPROUTE started procedure
2. OMPROUTE_FILE environment variable
3. /etc/omproute.conf
4. hlq.ETC.OMPROUTE.CONF

Example 5-22 OMPROUTE configuration file

```
Area Area_Number=0.0.0.2      1  
    Stub_Area=YES             2  
    Authentication_type=None  
    Import_Summaries=Yes;      3  
OSPF  
    RouterID=10.1.1.10         4  
    Comparison=Type2  
    DR_Max_Adj_Attempt = 10    5  
    Demand_Circuit=YES;  
Global_Options  
    Ignore_Undefined_Interfaces=YES 6  
;  
Routesa_Config Enabled=No;  
; Static vipa  
OSPF_interface ip_address=10.1.1.10  
    name=VIPAI1  
    subnet_mask=255.255.255.0  
    Advertise_VIPA_Routes=HOST_ONLY
```

```

        attaches_to_area=0.0.0.2
        cost0=10
        mtu=65535
OSPF_interface ip_address=10.1.2.10
        name=VIPA2L
        subnet_mask=255.255.255.0
        Advertise_VIPA_Routes=HOST_ONLY
        attaches_to_area=0.0.0.2
        cost0=10
        mtu=65535
; OSA Qdio VLAN10
OSPF_Interface IP_address=10.1.2.*
        Subnet_mask=255.255.255.0
        Router_Priority=0
        Attaches_To_Area=0.0.0.2
        Cost0=100
        MTU=1492;
; OSA Qdio VLAN11
OSPF_Interface IP_address=10.1.3.*
        Subnet_mask=255.255.255.0
        Router_Priority=0
        Attaches_To_Area=0.0.0.2
        Cost0=100
        MTU=1492;
; Hipersockets 10.1.4.x
OSPF_Interface IP_address=10.1.4.*
        Subnet_mask=255.255.255.0
        Router_Priority=1
        Attaches_To_Area=0.0.0.2
        Cost0=90
        MTU=8192;
; Hipersockets 10.1.5.x
OSPF_Interface IP_address=10.1.5.*
        Subnet_mask=255.255.255.0
        Router_Priority=1
        Attaches_To_Area=0.0.0.2
        Cost0=90
        MTU=8192;
; Hipersockets 10.1.6.x
OSPF_Interface IP_address=10.1.6.*
        Subnet_mask=255.255.255.0
        Router_Priority=1
        Attaches_To_Area=0.0.0.2
        Cost0=90
        MTU=8192;
;
; Dynamic vipa VIPADEFINE
ospf_interface ip_address=10.1.8.*
        subnet_mask=255.255.255.0
        Advertise_VIPA_Routes=HOST_ONLY
-         attaches_to_area=0.0.0.2
        cost0=10
        mtu=65535
;
; Dynamic vipa VIPADEFINE

```

12

13


```

ospf_interface ip_address=10.1.9.*
    subnet_mask=255.255.255.0
    Advertise_VIPA_Routes=HOST_ONLY
    attaches_to_area=0.0.0.2
    cost0=10
    mtu=65535
;
;AS_Boundary_routing
; Import_Direct_Routes=yes;

```

14

In this example, the numbers correspond to the following information:

1. Define the OSPF Area (Area 2).
2. Indicates Area 2 is a Stub Area.
3. Import_Summaries has meaning only if coded in an ABR. It makes the area connected to the ABR a Totally Stubby Area. If you coded the parameter in OMPROUTE on a Stub Area, it is ignored but it functions as a reminder that the ABR -- the layer 3 switch -- is defining our Stub Area as a Totally Stubby Area.
4. Defines the router internal IP address to be represented as the router ID.

Note: The RouterID *must be unique* on each OMPROUTE configuration. Otherwise, routing problems, timeouts, or poor performance will occur; the constant flooding of LSAs that contradict previous LSAs congests the network and consumes CPU as OSPF attempts to update its routing tables with the frequent changes in the topology database. We highly recommend that you code the RouterID statement (4), either with the static VIPA address or with an interface IP address, because the dynamic VIPAs (DVIPAs) can move between z/OS hosts within a sysplex.

OMPROUTE will issue message EZZ8165I when OSPF packets are received from a adjacent router with the same router ID that OMPROUTE is using. EZZ8165I is issued to the console once every 10 minutes per OSPF version. So, if a router is using the same router ID for both IPv4 and IPv6 OSPF, message EZZ8165I is issued twice. Automation can be put in place to monitor for the EZZ8165I message:

```

EZZ8165I DUPLICATE ip_version OSPF ROUTER ID router_number DETECTED

```

OMPROUTE cycles through all the OSPF interfaces until it finds a non-DVIPA interface if no RouterID is coded. A message (EZZ8134I) will be issued if a Dynamic VIPA is explicitly coded or chosen as a RouterID, because this is not recommended.

The interface with an IP address that represents RouterID (5) must be explicitly coded, as shown in Example 5-22, and not by using the wild card (*). Therefore, if the RouterID is hardcoded, as it is in this example, only portions of the configuration file can be shared across systems.

5. Defines the DR_MAX_ADJ_ATTEMPT parameter on the OSPF to enable this function. OMPROUTE will then report and control futile neighbor state loops during the adjacency formation process. Futile neighbor state loops are automatically detected and reported using message EZZ8157I. If a parallel OSPF interface is not available, adjacency formation attempts continue to be retried over the same interface. If parallel OSPF are available, an interface change is reported using message EZZ8158I.
6. Tells OSPF not to build an INTERFACE statement automatically for interfaces not defined in the OMPROUTE configuration file.

- 7.** Defines a specific interface to be an OSPF interface. When the specific IP address is coded (not the wildcard as in **12**) the NAME statement must also be configured (see **7**).
- 8.** The NAME statement identifies the link as an OSPF interface. The NAME statement must match the name specified in the LINK statement in the TCP/IP profile.
- 9.** Defines the interface should belong to Area 2.
- 10.** If the OSPF_Interface is a VIPA link, you can use this parameter to tell OMPROUTE how you want the VIPA address to be advertised. By default both the host and the subnet routes are advertised. Only the VIPA host route is advertised when this option is set to HOST_ONLY. We recommend the use of HOST_ONLY for VIPAs unless you have a compelling reason to advertise the subnet route.
- 11.** We define the interface cost of VIPA to be 10, HiperSockets to be 90, and OSA to be 100. We made the cost of HiperSockets smaller than OSA so that the HiperSockets route is preferred for the mainframe-to-mainframe communication.
- 12.** When defining OSPF_Interface IP_address with a wild card (*), all interfaces within the defined range will be seen as OSPF interfaces. Individual definitions with wild cards can be used for seeding other OMPROUTE configuration files. Remember that the unique RouterID of each file makes the entire file unshareable unless MVS system symbolics are employed.
- 13.** The z/OS Communications Server should be prevented from becoming the designated router in the LAN environment, when routers that are present can perform this function. To do this, define statement Router_Priority with value 0.
- 14.** Stub Areas do not permit importation of OSPF external, direct, or static routes; although the z/OS Communications Server on this node can learn about them, they will not be advertised. Therefore, the AS_Boundary_Routing statement is useless in this configuration and it is commented out. If it is not commented out, it is ignored when the node belongs to a Stub Area or Totally Stubby Area.

Our next example of an OMPROUTE configuration file can be shared across multiple stacks using MVS system symbols and too we can to use the statement INCLUDE. It can group OMPROUTE configuration statements that are common to several OMPROUTE instances into a single file. You do not need to repeat the configuration information in multiple places, we need only put the INCLUDE.

The use of MVS system symbols and statement INCLUDE in the OMPROUTE configuration file are introduced in Communications Server. We show it in Example 5-23.

Example 5-23 Shareable OMPROUTE configuration file using MVS system symbols and INCLUDE

```
OSPF
  RouterID=10.1.&SYSCclone..10 1
  Comparison=Type2
  Demand_Circuit=YES;
Global_Options
  Ignore_Undefined_Interfaces=YES

Routesa_Config Enabled=No;
; Static vipa
OSPF_Interface IP_address=10.1.&SYSCclone..10 2
  Subnet_mask=255.255.255.0
  Name=VIPa3L
  Attaches_To_Area=0.0.0.2
  Advertise_VIPa_Routes=HOST_ONLY
  Cost0=10
  MTU=65535;
INCLUDE //'TCPIPA.TCPPARMS(OMPA30IN)' 3
```

This OMPROUTE configuration file is now completely shareable. We have fully exploited wildcards, the MVS system symbolics and the statement INCLUDE.

In this example, the numbers correspond to the following information:

- 1.** We used a MVS system symbol &SYSCclone value to express the OSPF_Interface that is also used to represent our RouterID of 10.1.0.10 (the &SYSCclone value resolves to 30 on this LPAR; we used only one digit starting with the second digit of the &SYSCclone value).
- 2.** We used the same SYSCclone value to define the OSPF_Interface.
- 3.** We used the INCLUDE TCPIPA.TCPPARMS(OMPA30IN), and put in this file the statements of common interfaces Dynamic XCF. Example 5-24 shows this data set.

Example 5-24 Data set OMPA30IN with include configuration

```
;Dynamic XCF
interface ip_address=10.1.7.*
  subnet_mask=255.255.255.0
  mtu=65535;
```

5.5.5 Configure routers

In our router, we created a router service 100, which is analogous to an OMPROUTE procedure (or OSPF process). We defined the OSPF configuration for a Totally Stubby Area under this process.

To configure router 1, we used the configuration statements shown in Example 5-25.

Example 5-25 Router A configuration statements

```
interface Loopback1
  ip address 10.1.200.1 255.255.255.0
!
interface Vlan10
  ip address 10.1.2.240 255.255.255.0
  ip ospf cost 100
  ip ospf priority 100
!
interface Vlan11
  ip address 10.1.3.240 255.255.255.0
  ip ospf cost 100
  ip ospf priority 100
!
router ospf 100 1
  router-id 10.1.3.240 2
  log-adjacency-changes
  area 2 stub no-summary 3
  network 10.1.2.0 0.0.0.255 area 2 4
  network 10.1.3.0 0.0.0.255 area 2 4
  network 10.1.100.0 0.0.0.255 area 2 4
  network 10.200.1.0 0.0.0.255 area 0 5
  default-information originate always metric-type 1
```

In this example, the numbers correspond to the following information:

- 1.** Designates the process 100 to be an OSPF routing service.
- 2.** Defines the router ID of this process (100).
- 3.** Creates an OSPF area to this process (Area 2) and defines it to be a Stub Area.
- 4.** Defines the network range designated to Area 2. All interfaces within this IP address range (10.10.0.0) will belong to Area 2.
- 5.** Defines the network range designated to the backbone Area 0. All interfaces within this IP address range (10.200.0.0) will belong to Area 0.

5.5.6 Activation and verification

To activate and verify the OMPROUTE configuration, perform the following steps:

1. Start OMPROUTE.
2. Verify the configuration.

Start OMPROUTE

OMPROUTE can be started from an z/OS procedure, from the z/OS shell, or from AUTOLOG.

In our test environment, we started the OMPROUTE from the z/OS procedure, as shown in Example 5-26.

Example 5-26 OMPROUTE initialization

```
S OMPA
$HASP100 OMPA      ON STCINRDR
IEF695I START OMPA      WITH JOBNAME OMPA      IS ASSIGNED TO USER
TCP/IP      , GROUP TCPGRP
$HASP373 OMPA      STARTED
IEE252I MEMBER CTIORA00 FOUND IN SYS1.PARMLIB
EZZ7800I OMPA STARTING 1
EZZ7975I OMPA IGNORING UNDEFINED INTERFACE EZASAMEMVS 2
EZZ7475I ICMP WILL IGNORE REDIRECTS DUE TO ROUTING APPLICATION BEING
ACTIVE
EZZ8100I OMPA SUBAGENT STARTING
IEA989I SLIP TRAP ID=X13E MATCHED.  JOBNAME=RMFGAT      , ASID=0043.
EZZ7898I OMPA INITIALIZATION COMPLETE
```

In this example, the numbers correspond to the following information:

- 1.** The procedure name OMPA appears in several of the informational messages: EZZ7871I, EZZ7975I, and EZZ8100I. This facilitates problem determination in a SYSPLEX or CINET environment with messages flowing to a single console.
- 2.** Message EZZ7975I shows the effect of “Ignore_Undefined_Interfaces=YES” coding in the OMPROUTE configuration file shown in Example 5-22 on page 239.

You can use the AUTOLOG statement to start OMPROUTE automatically during TCP/IP initialization. Insert the name of the OMPROUTE start procedure in the AUTOLOG statement of the PROFILE.TCPIP data set (see Example 5-27).

Example 5-27 AUTOLOG statements

```
; *****
;
; TCPIPA.TCPPARMS(PROFA30)
; *****
.....
AUTOLOG 5
      OMPA ; OMPROUTE procedure
ENDAUTOLOG
.....
```

Verify the configuration

To verify that OMPROUTE is configured correctly as we defined, we can use either the DISPLAY or the MODIFY command.

In this section, we show some of the most useful DISPLAY commands and outputs. To see other display command options and to get more detailed information about specific commands, refer to *z/OS Communications Server: IP System Administrator's Commands*, SC31-8781.

To display OSPF configuration information, use the Display OMPROUTE,OSPF,LIST,ALL command. The sample display is shown in Example 5-28.

Example 5-28 -----+ command display

```

D TCPIP,TCPIPA,OMP,OSPF,LIST,ALL
EZZ7831I GLOBAL CONFIGURATION 413
  TRACE: 0, DEBUG: 1, SADEBUG LEVEL: 0
  STACK AFFINITY:          TCPIPA
  OSPF PROTOCOL:           ENABLED
  EXTERNAL COMPARISON:     TYPE 2
  AS BOUNDARY CAPABILITY:  DISABLED
  DEMAND CIRCUITS:         ENABLED
  DR MAX ADJ. ATTEMPT:     10

EZZ7832I AREA CONFIGURATION
AREA ID      AUTYPE      STUB?  DEFAULT-COST  IMPORT-SUMMARIES?
0.0.0.2      0=NONE      YES      1              YES
0.0.0.0      0=NONE      NO       N/A            N/A

```

1

```

EZZ7833I INTERFACE CONFIGURATION
IP ADDRESS   AREA      COST  RTRNS  TRDLY  PRI  HELLO  DEAD  DB_E*
10.1.8.10    0.0.0.2    10    N/A    N/A    N/A  N/A    N/A  N/A
10.1.6.11    0.0.0.2    90     5      1      1    10     40   40
10.1.5.11    0.0.0.2    90     5      1      1    10     40   40
10.1.4.11    0.0.0.2    90     5      1      1    10     40   40
10.1.2.12    0.0.0.2   100     5      1      0    10     40   40
10.1.3.12    0.0.0.2   100     5      1      0    10     40   40
10.1.3.11    0.0.0.2   100     5      1      0    10     40   40
10.1.2.11    0.0.0.2   100     5      1      0    10     40   40
10.1.2.10    0.0.0.2    10    N/A    N/A    N/A  N/A    N/A  N/A
10.1.1.10    0.0.0.2    10    N/A    N/A    N/A  N/A    N/A  N/A
10.1.30.10   0.0.0.2    10    N/A    N/A    N/A  N/A    N/A  N/A

```

2

3

```

ADVERTISED VIPA ROUTES
10.1.8.10    /255.255.255.255  10.1.2.10    /255.255.255.255
10.1.1.10    /255.255.255.255  10.1.30.10   /255.255.255.255

```

In this example, the numbers correspond to the following information:

- 1.** The Area 2 is defined as a Totally Stubby Area.
- 2.** The OSA interface has Router_Priority=1, Hello_Interval=10, Dead_Interval=40 specified to establish neighbors with other routers.
- 3.** The VIPA interface does not have Router_Priority, Hello_Interval, or Dead_Interval specified because they do not establish neighbors.

Display OSPF interfaces

Use the Display OMROUTE,OSPF,INTERFACE command to display the defined OSPF interfaces and their current status. Our display example is shown in Example 5-29.

Example 5-29 D TCPIP,TCPIPA,OMPR,OSPF,INTERFACE command display

D TCPIP,TCPIPA,OMPR,OSPF,INTERFACE							
EZZ7849I INTERFACES 400							
IFC	ADDRESS	PHYS	ASSOC. AREA	TYPE	STATE	#NBRS	#ADJS
10.1.8.10	VIPL0A01080A	0.0.0.2	VIPA	N/A	N/A	N/A	
10.1.6.11	IUTIQDF6L	0.0.0.2	BRDCST	32	4	2	
10.1.5.11	IUTIQDF5L	0.0.0.2	BRDCST	32	4	2	
10.1.4.11	IUTIQDF4L	0.0.0.2	BRDCST	32	4	2	
10.1.3.12	OSA20E0I	0.0.0.2	BRDCST	32	4	1	1
10.1.3.11	OSA20C0I	0.0.0.2	BRDCST	2	0	0	
10.1.2.12	OSA20A0I	0.0.0.2	BRDCST	32	3	1	
10.1.2.14	OSA2081I	0.0.0.2	BRDCST	2	0	0	
10.1.2.11	OSA2080I	0.0.0.2	BRDCST	2	0	0	
10.1.2.10	VIPA2L	0.0.0.2	VIPA	N/A	N/A	N/A	
10.1.1.10	VIPA1L	0.0.0.2	VIPA	N/A	N/A	N/A	2
10.1.30.10	VIPA3L	0.0.0.2	VIPA	N/A	N/A	N/A	
* -- LINK NAME TRUNCATED							

In this example, the numbers correspond to the following information:

- 1.** The OSA interface is attached to Area 2 and has four neighbors established.
- 2.** The VIPA interface is attached to Area 2 but does not establish neighbors.

Display OSPF neighbors

Use the Display OMROUTE,OSPF,NBRS command to display the OSPF neighbors and their current status. Our display example is shown in Example 5-30.

Example 5-30 D TCPIP,TCPIPA,OMPR,OSPF,NBRS command display

D TCPIP,TCPIPA,OMPR,OSPF,NBRS							
EZZ7851I NEIGHBOR SUMMARY 392							
NEIGHBOR ADDR	NEIGHBOR ID	STATE	LSRXL	DBSUM	LSREQ	HSUP	IFC
10.1.3.240	10.1.3.240	128	0	0	0	OFF	OSA20E0I 1
10.1.3.41	10.1.3.10	8	0	0	0	OFF	OSA20E0I 2
10.1.2.22	10.1.1.20	8	0	0	0	OFF	OSA20A0I 2
10.1.2.240	10.1.3.240	128	0	0	0	OFF	OSA20A0I 1
10.1.2.22	10.1.31.10	8	0	0	0	OFF	OSA20A0I 2
10.1.4.21	10.1.31.10	128	0	0	0	OFF	IUTIQDF4L 3
10.1.5.21	10.1.31.10	128	0	0	0	OFF	IUTIQDF5L 3
10.1.6.21	10.1.31.10	128	0	0	0	OFF	IUTIQDF6L 3
* -- LINK NAME TRUNCATED							

In this example, the numbers correspond to the following information:

- 1.** The neighbor with router 1 is established in each VLAN that OSA belongs to. The state is 128 (Full).
- 2.** The neighbor with TCPIPB stack is established on the OSA interface. The state is 8 (2-way), because router 1 and router 2 are DR/BDR, so TCPIPA and TCPIPB are both DR other.
- 3.** The neighbor with TCPIPB stack is established on the HiperSockets interface. The state is 128 (Full), because TCPIPB is the DR on the HiperSockets subnet.

Display OSPF routers

Use the Display OMROUTE,OSPF,ROUTERS command to display the OSPF routes to ABRs and ASBRs. Our display example is shown in Example 5-30 on page 247.

Example 5-31 D TCPIP,TCPIPA,OMPR,OSPF,ROUTERS command display

```
D TCPIP,TCPIPA,OMPR,OSPF,ROUTERS
EZZ7855I OSPF ROUTERS 402
DTYPE RTYPE DESTINATION      AREA      COST      NEXT HOP(S)
BR  SPF   10.1.3.240          0.0.0.2    100        10.1.2.240 * 1
```

In this example, the number corresponds to the following information:

- 1.** Router 1 is the ABRs.

Display OMROUTE routing table

Use the Display OMROUTE,RTTABLE command to display the OMROUTE routing table. Our display example is shown in Example 5-32.

Example 5-32 D TCPIP,TCPIPA,OMPR,RTTABLE command display

```
D TCPIP,TCPIPA,OMPR,RTTABLE
EZZ7847I ROUTING TABLE 404
TYPE  DEST NET      MASK      COST    AGE    NEXT HOP(S)
SPIA  0.0.0.0        0 101    35859   10.1.2.240 (5) 1
DIR*  10.1.1.0        FFFFFFF0 1    37110   10.1.1.10 2
DIR*  10.1.1.10       FFFFFFFF 1    37110   VIPA1L 2
SPF   10.1.1.12       FFFFFFFF 90    37098   10.1.4.12 (2) 3
SPF   10.1.1.20       FFFFFFFF 90    37098   10.1.4.21 (2)
SPF   10.1.1.40       FFFFFFFF 110    35859   10.1.2.42 (5)
SPF*  10.1.2.0         FFFFFFF0 100    35859   OSA2080I (3)
DIR*  10.1.2.10       FFFFFFFF 1    37110   VIPA2L 2
SPF   10.1.2.17       FFFFFFFF 90    37098   10.1.4.12 (2)
SPF*  10.1.3.0         FFFFFFF0 100    37103   OSA20C0I (2)
SPF*  10.1.4.0         FFFFFFF0 80    37100   IUTIQDF4L
SPF*  10.1.5.0         FFFFFFF0 80    37098   IUTIQDF5L
SPF*  10.1.6.0         FFFFFFF0 190    37103   IUTIQDF6L
SPF   10.1.8.10       FFFFFFFF 90    35859   10.1.4.12 (4)
SPF   10.1.8.20       FFFFFFFF 90    37098   10.1.4.21 (2)
```

DEFAULT GATEWAY IN USE.

```
TYPE COST    AGE    NEXT HOP
SPIA 101      13248   10.1.2.240 (3)
0 NETS DELETED, 1 NETS INACTIVE
```


In this example, the numbers correspond to the following information:

- 1.** Only the default route is advertised from ABR.
- 2.** Direct routes to the subnet to which local interfaces belong are listed.
- 3.** Indirect routes to the VIPA in TCPIPB are listed.

Display TCP/IP routing table

Use the **netstat** command to display the OSPF routes to ABRs and ASBRs. Our display example is shown in Example 5-33.

Example 5-33 D TCPIP,TCPIPA,N,ROUTE,MAX= command display*

```

D TCPIP,TCPIPA,N,ROUTE,MAX=*
IPV4 DESTINATIONS
DESTINATION      GATEWAY      FLAGS      REFCNT      INTERFACE
DEFAULT          10.1.2.240   UGO        000000      OSA2080I 1
DEFAULT          10.1.2.240   UGO        000000      OSA20A0I
DEFAULT          10.1.3.240   UGO        000000      OSA20C0I
DEFAULT          10.1.3.240   UGO        000000      OSA20E0I
10.1.1.10/32     0.0.0.0      UH         000000      VIPA1L
10.1.1.12/32     10.1.4.12    UGHO       000000      IUTIQDF4 2
10.1.1.12/32     10.1.5.12    UGHO       000000      IUTIQDF5
10.1.1.20/32     10.1.4.21    UGHO       000000      IUTIQDF4
10.1.1.20/32     10.1.5.21    UGHO       000000      IUTIQDF5
10.1.1.40/32     10.1.2.42    UGHO       000000      OSA2080I
10.1.1.40/32     10.1.2.42    UGHO       000000      OSA20A0I
10.1.1.40/32     10.1.3.41    UGHO       000000      OSA20C0I
10.1.1.40/32     10.1.3.41    UGHO       000000      OSA20E0I
10.1.2.0/24      0.0.0.0      UO         000000      OSA2080I
127.0.0.1/32    0.0.0.0      UH         000004      LOOPBACK
IPV6 DESTINATIONS
DESTIP:  ::1/128
  GW:    ::
  INTF:  LOOPBACK6      REFCNT:  000000
  FLGS:  UH             MTU:  65535
END OF THE REPORT

```

In this example, the numbers correspond to the following information:

- 1.** The default routes to the router 1 is listed. We do not have IPCONFIG MULTIPATH specified, so the first active default route entry (interface OSA2080I and gateway address 10.1.2.240) is always used for a destination that does not have an explicit route entry.
- 2.** Indirect routes to the VIPA in TCPIPB stack are listed. We do not have IPCONFIG MULTIPATH specified, so the first active route entry (interface IUTIQDF4L) is always used for the VIPA destination.

Check the connectivity using PING command

The PING command can be executed with the TSO PING command or the z/OS UNIX **ping** command. Example 5-34 on page 250 shows the display of the TSO PING command. We see the ping is successful.

In a CINET environment where multiple TCP/IP stacks are configured, use the TCP option for the TSO PING command and the -p option for the z/OS UNIX **ping** command to specify the TCP/IP stack name from which you want to issue the **ping** command.

You do not need to specify those options if the user issuing this command is already associated to the TCP/IP stack (with SYSTCPD DD, for example). There is no need to specify these options if your environment is an INET environment where only one TCP/IP stack is configured.

Example 5-34 TSO PING command display

```
TSO PING 10.1.1.20 (TCP TCPIPA
CS V1R12: Pinging host 10.1.1.20
Ping #1 response took 0.000 seconds.
***
```

Example 5-35 shows the display of z/OS UNIX **ping** command.

Example 5-35 z/OS UNIX ping command display

```
CS02 @ SC30:/u/cs02>ping -p TCPIPA 10.1.1.20
CS V1R12: Pinging host 10.1.1.20
Ping #1 response took 0.000 seconds.
```

Verify the selected route with the TRACEROUTE command

TRACEROUTE can be invoked by either the TSO TRACERTE command or the z/OS UNIX shell **traceroute**/**otracer** command. Example 5-36 shows an example of the display. We see that the router 1 (10.1.2.240) is the next hop router to reach destination IP address 10.1.100.221.

In a CINET environment where multiple TCP/IP stacks are configured, use the TCP option for the TSO TRACERTE command and the **-a** option for the z/OS UNIX **traceroute** command to specify the TCP/IP stack name from which you want to issue the TRACEROUTE command.

You do not need to specify those options if the user issuing this command is already associated to the TCP/IP stack (with SYSTCPD DD, for example). There is no need to specify those options if your environment is an INET environment where only one TCP/IP stack is configured.

Example 5-36 TRACERTE command results

```
TSO TRACERTE 10.1.100.221 (TCP TCPIPA
CS V1R12: Traceroute to 10.1.100.221 (10.1.100.221):
  1 router1 (10.1.2.240)  0 ms  0 ms  0 ms
  2 10.1.100.221 (10.1.100.221)  0 ms  0 ms  0 ms
***
```

5.5.7 Managing OMPROUTE

You can manage OMPROUTE from a z/OS operator console. Commands are available to perform the following:

- ▶ Stop OMPROUTE.
- ▶ Modify OMPROUTE.
- ▶ Display OMPROUTE.

Stop OMPROUTE from z/OS console

OMPROUTE can be stopped from the z/OS console by issuing **STOP <procname>** or **MODIFY <procname>, KILL**.

Example 5-37 shows the display.

Example 5-37 Stopping OMPROUTE from z/OS console

```
P OMPA
EZZ7804I OMPA EXITING
ITT120I SOME CTRACE DATA LOST, LAST 5 BUFFER(S) NOT WRITTEN
$HASP395 OMPA      ENDED
```

Stop OMPROUTE from z/OS UNIX shell

You can also stop OMPROUTE from a z/OS UNIX shell superuser ID by issuing the **kill** command to the process ID (PID) associated with OMPROUTE. To determine the PID, use one of the following methods:

- From the z/OS console, issue **D OMVS,U=userid** (where *userid* is the user ID that started omproute from the shell). From the resulting display, look at the PID number related to OMPROUTE, as shown in Example 5-38 (1).

Example 5-38 Stopping OMPROUTE from z/OS UNIX

```
D OMVS,U=TCPIP
BPX0040I 14.56.39 DISPLAY OMVS 617
OMVS      000E ACTIVE          OMVS=(7A)
USER      JOBNAME  ASID      PID      PPID STATE   START      CT_SECS
TCPIP     OMPA     00EF      50397483 1 1  HS----  14.43.13  243.843
LATCHWAITPID=      0 CMD=OMPRROUTE
```

- From the z/OS UNIX shell, issue the **ps -ef** command, as shown in Example 5-39 (1).
- Using the PID number, stop omproute using the **kill pidnumber** command as seen in Example 5-39 (2)

Example 5-39 kill command in z/OS UNIX shell

```
CS03 @ SC30:/u/cs03>ps -ef
      UID      PID      PPID  C   STIME TTY      TIME CMD
BPXR00T      1      0   - Oct 11 ?      0:14 BPXPINPR
BPXR00T 16842754      1   - Oct 11 ?      1:07 BPXVCMT
BPXR00T 50397483 1      1   - 14:43:14 ?      4:10 OMPROUTE 1
CS03 @ SC30:/u/cs03>kill 50397483 2
```

Modify OMPROUTE configuration

We can use the **MODIFY (F)** command to change some configuration statements and to start, stop, or change the level of OMPROUTE tracing and debugging, as follows:

- **F procname,RECONFIG** command: Used to reread the OMPROUTE configuration file, adding new OSPF_interfaces
- **F procname,ROUTESA=ENABLE/DISABLE** command: Used to enable or disable the OMPROUTE subagent
- **F procname,OSPF,WEIGHT,NAME=<if_name>,COST=<cost>** command: Changes dynamically the cost of an OSPF interface

In OSPF environments in which there might be a problem with some remote hardware (for example, router, switch, or network cable) that is beyond detection by z/OS hardware or software, OMPROUTE can get into an infinite neighbor state loop over one of its interfaces with a neighbor. This loop might contribute to increased workload. In LAN configurations in which there are parallel OSPF interfaces that can reach the same neighbor for adjacency formation, unless you are using OMPROUTE futile neighbor state loop detection or unless you manually fix the problem, the backup interfaces are not used until after an outage occurs for the OSPF interface that was initially involved in an adjacency formation attempt with a designated router

We can use the MODIFY (F) command to suspend and, after fixing the problem, activate an OSPF interface using the **F *procname*,OSPF,INTERFACES,NAME=*interfname*,SUSPEND or ACTIVATE** command, which suspends or activates the OMPROUTE interface. In Example 5-40, we shows these commands and the status of the interface. First, we suspend the interface OSA20A01 **1**, then we issue the display and see the state **1*** (suspend) **2**. Finally, we reactivate the interface to normal status **3**.

Example 5-40 MODIFY SUSPEND and ACTIVATE commands

F OMPA,OSPF,INTERFACES,NAME=OSA20A0I,SUSPEND 1							
EZZ7866I OMPA MODIFY COMMAND ACCEPTED							
EZZ8159I OMPA MODIFY SUSPEND COMMAND FOR OSPF IPV4 INTERFACE OSA20A0I IS SUCCESSFUL							
D TCPIP,TCPIPA,OMP,OSPF,INTERFACES							
EZZ7849I INTERFACES 803							
IFC	ADDRESS	PHYS	ASSOC.	AREA	TYPE	STATE	#NBRS
#ADJS							
10.1.6.11		IUTIQDF6L	0.0.0.2		BRDCST	64	1
10.1.5.11		IUTIQDF5L	0.0.0.2		BRDCST	32	2
10.1.4.11		IUTIQDF4L	0.0.0.2		BRDCST	32	2
10.1.3.12		OSA20E0I	0.0.0.2		BRDCST	32	4
10.1.3.11		OSA20C0I	0.0.0.2		BRDCST	2	0
10.1.2.12		OSA20A0I	0.0.0.2		BRDCST	1*	0 2
10.1.2.14		OSA2081I	0.0.0.2		BRDCST	1	0
F OMPA,OSPF,INTERFACES,NAME=OSA20A0I,ACTIVATE 3							
EZZ7866I OMPA MODIFY COMMAND ACCEPTED							
EZZ8160I OMPA MODIFY ACTIVATE COMMAND FOR OSPF IPV4 INTERFACE OSA20A0I IS SUCCESSFUL							

Note: Use the MODIFY SUSPEND command to stop OSPF traffic on an OSPF interface, rather than using the VARY TCPIP command to deactivate the corresponding physical interface in TCPIP. This will allow existing sessions using static routes on the affected interface not to be disrupted.

Display OMPROUTE information

We can use the MODIFY (F) command instead of the DISPLAY TCPIP command to display information for OMPROUTE. Both commands provide the same information and use the same statements, as shown in the following samples:

- **F *procname*,RTTABLE** command: The resulting display provides the same contents as though we were using **D TCPIP,*procnamename*,OMP,RTTABLE**.
- **F *procname*,OSPF,LIST ALL** command: The resulting display provides us the same contents as though we were using **D TCPIP,*procnamename*,OMP,OSPF,LIST ALL**.

Start, stop, or change the level of OMPROUTE tracing and debugging

We can use the MODIFY (F) command to start, stop, or change the level of tracing and debugging.

- **F *procname*,TRACE=n:** For OMPROUTE tracing for initialization and IPv4 routing protocols; n can be 0–2.
- **F *procname*,DEBUG=n:** For OMPROUTE debugging for initialization and IPv4 routing protocols; n can be 0–4.
- **F *procname*,SADEBUG=n:** For OMPROUTE subagent debugging; n can be 0 or 1.

For further information about these commands and options, refer to *z/OS Communications Server: IP System Administrator's Commands*, SC31-8781.

5.6 Problem determination

When implementing a network environment with indirect access to external hosts or networks using routing definitions, it is important to understand how to isolate networking problems. This means that using the correct diagnostic tools and techniques is essential. In this section we describe the tools and techniques needed to debug routing problems in a static routing environment and in a dynamic OSPF routing environment. To debug a network problem in a z/OS environment, we suggest following the flow shown in Figure 5-4.

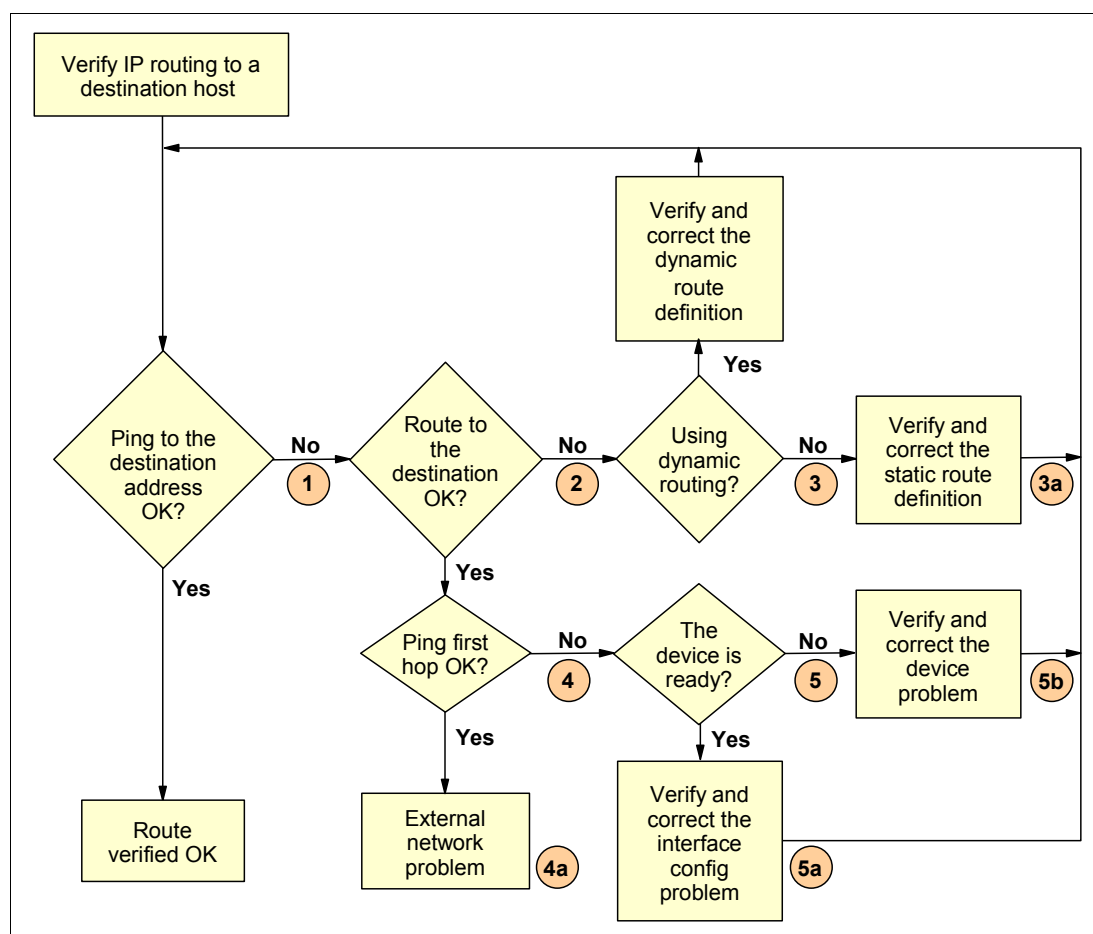


Figure 5-4 Routing problem determination flow

The descriptions for the tags shown in Figure 5-4 on page 253 are as follows:

1. Use the **ping** command to determine whether there is connectivity to the destination IP address. More information about the **ping** command can be found in “PING command (TSO or z/OS UNIX)” on page 255.
2. If the **ping** command fails immediately, there might not be a route to the destination host or subnet. Use the **netstat ROUTE/ -r** command to display routes to the network, as shown in Example 5-10 on page 232. Verify that TCP/IP has a route to the destination address.
If there is no route, proceed to step 3. If a route does exist, then proceed to step 4.
3. If there is no route to the destination, problem resolution depends on whether static or dynamic routing is being used. In either case, do the following:
 - a. If TCP/IP is configured using Static Routing, review and correct the configuration.
 - b. If OMPROUTE is being used to generate dynamic routes, verify and correct the configuration. If it seems correct, then diagnose the problem using the debugging tools described in “Diagnosing an OMPROUTE problem” on page 256.
4. If a route exists, verify that the route is correct for the destination. Determine whether the gateway identified for the route to the destination is reachable. To verify this, use the PING command to confirm connectivity to the gateway.

Do *one* of the following:

- a. If the gateway responds to a ping, it means there is a network problem at the gateway or beyond. To get further debug information, use the **tracert** command with the final destination address to determine which hop in the route is failing.
 - b. If the gateway does not respond to a ping, proceed to step 5.
5. Determine which network interface is associated with the route to the destination. Verify that it is operational by issuing the **netstat Devlink** command, as shown in Example 5-9 on page 231.

Based on the resulting display, do *one* of the following:

- a. If the device is ready, the problem might be in the interface configuration. Check the network configuration (VLAN ID, IP address, subnet mask, and so on). Correct this and resume testing.

Otherwise, a packet trace should be taken to verify that the packets are being sent to the network. A LAN Analyzer could also be used to verify the network traffic in the switch port where the OSA-Express port is connected.

- b. If the device is not ready, the problem might be that the device is not varied online to z/OS, or that there is an error in the device configuration. Also verify the VTAM TRLE definitions, HCD/IOCP configuration, as well as the physical connection, cable, and switch port.

5.6.1 Commands to diagnose networking connectivity problems

In this section we describe briefly the commands that can be used to diagnose connectivity problems. For additional help and information about diagnosing problems, refer to *z/OS Communications Server: IP System Administrator's Commands*, SC31-8781.

PING command (TSO or z/OS UNIX)

We used the **ping** command to verify:

- ▶ The route to the remote network is defined correctly.
- ▶ The router is able to forward packets to the remote network.
- ▶ The remote host is able to send and receive packets in the network.
- ▶ The remote host has a route back to the local host.

The **ping** command can be executed with the TSO PING command or the z/OS UNIX **ping** command. Example 5-41 shows the display of TSO PING command. We see that the ping is successful.

In a CINET environment where multiple TCP/IP stacks are configured, use the TCP option for the TSO PING command and the **-p** option for the /OS UNIX **ping** command to specify the TCP/IP stack name from which you want to issue the **ping** command. You do not need to specify those options if you are issuing this command in the associated TCP/IP stack (with SYSTCPD DD, for example). There is no need to specify this option if your environment is an INET environment where only one TCP/IP stack is configured.

Example 5-41 TSO PING command display

```
TSO PING 10.1.1.20 (TCP TCPIPA
CS V1R12: Pinging host 10.1.1.20
Ping #1 response took 0.000 seconds.
***
```

Example 5-42 shows the display of the z/OS UNIX **ping** command.

Example 5-42 z/OS UNIX ping command display

```
CS02 @ SC31:/u/cs02>ping -p TCPIPA 10.1.1.20
CS V1R12: Pinging host 10.1.1.20
Ping #1 response took 0.000 seconds.
```

TRACEROUTE command

Traceroute can be invoked by either the TSO TRACERTE command or the z/OS UNIX shell **traceroute**/**otracert** command.

Traceroute displays the route that a packet takes to reach the requested target. Traceroute starts at the first router and uses a series of UDP probe packets with increasing IP time-to-live (TTL) or hop count values to determine the sequence of routers that must be traversed to reach the target host. The output generated by this command can be seen in Example 5-43.

Example 5-43 TSO TRACERTE command results

```
TSO TRACERTE 10.1.100.221 (TCP TCPIPA
CS V1R12: Traceroute to 10.1.100.221 (10.1.100.221):
  1 10.1.2.240 (10.1.2.240)  0 ms  0 ms  0 ms
  2 10.1.100.221 (10.1.100.221)  0 ms  0 ms  0 ms
***
```

In a CINET environment where multiple TCP/IP stacks are configured, use the TCP option for the TSO TRACERTE command and the **-a** option for the z/OS UNIX **traceroute** command to specify the TCP/IP stack name from which you want to issue the TRACEROUTE command.

Note that you do not need to specify those options if the user issuing this command is already associated to the TCP/IP stack (with SYSTCPD DD, for example). There is no need to specify those options if your environment is an INET environment where only one TCP/IP stack is configured.

Example 5-44 shows the display of the z/OS UNIX **tracert** command.

Example 5-44 z/OS UNIX tracert command result

```
CS02 @ SC31:/u/cs02>tracert -a TCPIPA 10.1.100.221
CS V1R12: Tracert to 10.1.100.221 (10.1.100.221)
Enter ESC character plus C or c to interrupt
1 10.1.2.240 (10.1.2.240)  0 ms  0 ms  0 ms
2 10.1.100.221 (10.1.100.221)  0 ms  0 ms  0 ms
```

Tip: Using a name instead of IP address would need the resolver or DNS to do the translation. This adds more variables to the problem determination, and should be avoided when you are diagnosing network problems. Use the host IP address instead.

NETSTAT,DEVLINK command (console or z/OS UNIX)

Use the D TCPIP,*procname*,NETSTAT, DEVLINK command to display the status and associated configuration values for a device and its defined interfaces. From the z/OS UNIX shell, use the **netstat -d -p *procname*** command. The results are identical in the console or the z/OS UNIX shell. Example 5-9 on page 231 shows a sample display.

NETSTAT,ROUTE command (console or z/OS UNIX)

Use the D TCPIP,*procname*,NETSTAT,ROUTE command to display the current routing tables for TCP/IP. From z/OS UNIX shell, use the **netstat -r -p *procname*** command. Example 5-33 on page 249 shows a sample display.

5.6.2 Diagnosing an OMPROUTE problem

This section describe methods that you can use to diagnose an OMPROUTE problem.

Useful commands

In addition to the commands that we show in 5.6.1, “Commands to diagnose networking connectivity problems” on page 254, you can use additional commands to diagnose OMPROUTE problems, as described here.

D TCPIP,TCPIPA,OMP,OSPF,NBR command

This command displays all the OSPF neighbors. Make sure you have established the neighbor with other routers as you expected. Example 5-30 on page 247 shows a sample display.

D TCPIP,TCPIPA,OMP,RTTABLE command

This command displays the OMPROUTE routing table. Make sure you have the expected route listed in the table. If you have multiple routes for the destination, with different costs, only the best route (least cost route) is added to the OMPROUTE and TCP/IP routing tables. Example 5-32 on page 248 shows a sample display.

D TCPIP,TCPIPA,OMPR,RTTABLE,DELETED command

This command displays all of the route destinations that have been deleted from the OMPROUTE routing table since the initialization of OMPROUTE at this node. The routes that have changed the next hop are *not* considered deleted, and are therefore *not displayed* with this command. Example 5-45 shows the results of this display after OMPROUTE is terminated at SC31 (OMPB), another member of the SYSPLEX.

Example 5-45 Deleted OMPROUTE destinations

```
D TCPIP,TCPIPA,OMPR,RTTABLE,DELETED
EZZ8137I  IPV4 DELETED ROUTES 182
TYPE      DEST NET          MASK      COST      AGE      NEXT HOP(S)

DEL    10.1.1.20          FFFFFFFF  16        66        NONE
DEL    10.1.2.20          FFFFFFFF  16        66        NONE
DEL    10.1.8.20          FFFFFFFF  16        66        NONE
      3 NETS DELETED, 2 NETS INACTIVE
```

Observing initialization messages and taking traces

If these commands do not help, use traces for further diagnosis and verify whether you have any error messages related to OMPROUTE during the startup process. To do so, examine SYSLOGD, JES MSG output and the console log for errors.

If there is no apparent error message that could help you to solve the problem, then prepare OMPROUTE to generate more detailed information by using the debug tools available in OMPROUTE. This can be activated by coding the Debug and Trace options in the startup procedure, or by using the MODIFY command to implement these options.

Using OMPROUTE trace and debug for initialization

A trace from startup is ideal because some information is only shown in the trace at startup, and because the time for problem determination and resolution is faster when the trace captures the entire flow of events rather than just a small subset of events.

An OMPROUTE trace from startup can be enabled by coding the trace options after the forward slash (/) in the PARM field of the OMPROUTE catalogued procedure, as shown in Example 5-46.

Example 5-46 Trace options defined in the OMPROUTE startup procedure

```
//OMPR30A PROC STDENV=STDENV&SYSLCLONE
//OMPR30A EXEC PGM=OMPROUTE,REGION=4096K,TIME=NOLIMIT,
//          PARM=(' POSIX(ON) ALL31(ON) ',
//          'ENVAR("_BPXK_SETIBMOPT_TRANSPORT=TCPIPA"',
//          '" _CEE_ENVFILE=DD:STDENV")/-t2 -d1')
//          /*
//STDENV DD DISP=SHR,DSN=TCPIPA.OMPROUTE.&STDENV
```

If a trace cannot be enabled from startup, then the following commands can dynamically enable and disable tracing:

- ▶ To enable:
 - MODIFY omproute,TRACE=2 (TRACE6=2 for IPv6)
 - MODIFY omproute,DEBUG=1 (DEBUG6=1 for IPv6)
- ▶ To disable:
 - MODIFY omproute,TRACE=0 (TRACE6=0 for IPv6)
 - MODIFY omproute,DEBUG=0 (DEBUG6=0 for IPv6)

Trace output is sent to one of the following locations:

- ▶ A destination referenced by the `OMPROUTE_DEBUG_FILE` environment variable (which is coded in the `STDENV DD` data set).
- ▶ `STDOUT DD`, but trace output is only output to this location if `OMPROUTE_DEBUG_FILE` is *not* defined, and the trace is started at initialization.
- ▶ `/tmp/omproute_debug` (`TMPDIR` is usually `/tmp`.)

By default, `OMPROUTE` will create five debug files, each 200 KB in size, for a total of 1 MB of trace data. The size and number of trace files can be controlled with the `OMPROUTE_DEBUG_CONTROL` environment variable. For further information about the `TRACE` and `DEBUG` options, refer to *z/OS Communications Server: IP Diagnosis Guide*, GC31-8782.

Important: Using the `OMPROUTE` Trace and Debug options and directing the output to z/OS UNIX file system files generates additional overhead that might cause OSPF adjacency failures or other routing problems. To prevent that, change the output destination to the `CTRACE` Facility.

Using OMPROUTE CTRACE to get debugging information

As mentioned, the overhead problems that can occur when using z/OS UNIX file system files to save the trace and debug output data can be resolved by using the `CTRACE` facility. To use this facility, we strongly recommend using the `OMPROUTE` option (`DEBUGTRC`) in the startup procedure, which changes the output destination of the `OMPROUTE` trace. In this section we briefly describe how to define and use `CTRACE` to debug `OMPROUTE` problems.

You can start the `OMPROUTE CTRACE` anytime by using the command `TRACE CT`. Or it can be activated during `OMPROUTE` initialization. If not defined, `OMPROUTE` component trace is started with a buffer size of 1 MB and the `MINIMUM` tracing option.

A parmlib member can be used to customize the parameters and to initialize the trace. The default `OMPROUTE` Component Trace parmlib member is the `SYS1.PARMLIB` member `CTIORA00`. The parmlib member name can be changed by using the `OMPROUTE_CTRACE_MEMBER` environment variable.

In addition to specifying the trace options, you can also change the `OMPROUTE` trace buffer size. (Note that the buffer size can be changed *only* at `OMPROUTE` initialization.) The maximum `OMPROUTE` trace buffer size is 100 MB. The `OMPROUTE REGION` size in the `OMPROUTE` catalog procedure must be large enough to accommodate a large buffer size.

When `OMPROUTE` is initialized using the `DEBUGTRC` option, we recommend that you use a larger internal `CTRACE` buffer or an external writer. When using the internal `CTRACE` buffer, we must get a `DUMP` of `OMPROUTE` in order to see the trace output.

In this section we illustrate the steps needed to start the CTRACE for OMPROUTE and direct the trace output to an external writer:

1. Create a CTWTR procedure in your SYS1.PROCLIB, as shown in Example 5-47.

Example 5-47 CTWTR procedure

```
//CTWTR    PROC
//IEFPROC  EXEC PGM=ITTTTCWR
//TRCOUT01 DD  DSN=SYS1.&SYSNAME..OMPA.CTRACE,
//          VOL=SER=COMST2,UNIT=3390,
//          SPACE=(CYL,10),DISP=(NEW,CATLG),DSORG=PS
//*
```

2. Prepare the SYS1.PARMLIB member CTIORA00 to get the desired output data.
Example 5-48 shows a sample of CTIORA00 contents.

Example 5-48 CTIORA00 sample

```
/******
/*
/* DESCRIPTION = This parmlib member causes component trace for
/*               the TCP/IP OMPROUTE application to be initialized
/*               with a trace buffer size of 1M
/*
/*               This parmlib member only lists those TRACEOPTS
/*               values specific to OMPROUTE. For a complete list
/*               of TRACEOPTS keywords and their values see
/*               z/OS MVS INITIALIZATION AND TUNING REFERENCE.
/*
/*
/* $MAC(CTIORA00),COMP(OSPF ),PROD(TCPIP ): Component Trace
/*                                     SYS1.PARMLIB member
/*
/******
TRACEOPTS
/* ----- */
/*  Optionally start external writer in this file (use both
/*  WTRSTART and WTR with same wtr_procedure)
/* ----- */
/*      WTRSTART(CTWTR)                a
/* ----- */
/*  ON OR OFF: PICK 1
/* ----- */
/*      ON
/*      OFF
/* ----- */
/*  BUFSIZE: A VALUE IN RANGE 128K TO 100M
/*          CTRACE buffers reside in OMPROUTE Private storage
/*          which is in the regions address space.
/* ----- */
/*      BUFSIZE(50M)                   b
/*      WTR(CTWTR)                     a
/* ----- */
/*  OPTIONS: NAMES OF FUNCTIONS TO BE TRACED, OR "ALL"
/* ----- */
/*      OPTIONS(                        c
/*
```

```

/*          'ALL          '          */
/*          , 'MINIMUM '    */
/*          , 'ROUTE  '    */
/*          , 'PACKET  '    */
/*          , 'OPACKET '    */
/*          , 'RPACKET '    */
/*          , 'IPACKET '    */
/*          , 'SPACKET '    */
/*          , 'DEBUGTRC'    d
/*                                     )          */

```

In this example, the letters correspond to the following information:

- **a** Define whether we are going to use an external writer to save the output trace data.
- **b** Define the CTRACE buffer size allocated in the OMPROUTE private storage.
- **c** Define the trace options to be used to get specific debug information. MINIMUM is the default option.
- **d** This option indicates we will send to CTRACE the trace and debug level options defined in the OMPROUTE startup procedure.

3. Start the OMPROUTE procedure using the desired Debug and Trace options, as shown in Example 5-49.

Example 5-49 OMPROUTE procedure

```

//OMPA  PROC  STDENV=OMPENA&SYSCONE
//OMPA  EXEC  PGM=OMPROUTE,REGION=OM,TIME=NOLIMIT,
//        PARM=(' POSIX(ON) ALL31(ON) ',
//              ' ENVAR("_BPXK_SETIBMOPT_TRANSPORT=TCPIPA" ',
//              ' "_CEE_ENVFILE=DD:STDENV")/-t2 -d1' ) a
//STDENV DD  DISP=SHR,DSN=TCPIP.SC&SYSCONE..STDENV(&STDENV)

```

In this example, the letters correspond to the following information:

- **a** The parameters -t (trace) and -d (debug) define how detailed we want the output data to be. We recommend using -t2 and -d1.

To verify whether CTRACE has been started as expected, display the CTRACE status, issuing the console command shown in Example 5-50.

Example 5-50 Displaying OMPROUTE CTRACE status

```

D TRACE,COMP=SYSTCPRT,SUB=(OMPA)
IEE843I 16.23.40 TRACE DISPLAY 677
      SYSTEM STATUS INFORMATION
ST=(ON,0001M,00004M) AS=ON BR=OFF EX=ON MO=OFF MT=(ON,024K)
TRACENAME
=====
SYSTCPRT          MODE BUFFER HEAD SUBS
=====
                   OFF          HEAD    2
      NO HEAD OPTIONS
SUBTRACE          MODE BUFFER HEAD SUBS
-----
OMPA              ON    0010M
  ASIDS           *NONE*
  JOBNAMES        *NONE*
  OPTIONS         MINIMUM ,DEBUGTRC
  WRITER          CTWTR

```

We can also use TRACE CT command to define the options we want after OMPROUTE has been initialized, and send the trace to an external writer, following these steps:

1. Start the the ctrace external writer, as seen in Example 5-51.

Example 5-51 Starting the ctrace external writer, CTWTR, partial console output

```

TRACE CT,WTRSTART=CTWTR
ITT038I ALL OF THE TRANSACTIONS REQUESTED VIA THE TRACE CT COMMAND
WERE SUCCESSFULLY EXECUTED.
...
IRR812I PROFILE ** (G) IN THE STARTED CLASS WAS USED
      TO START CTWTR WITH JOBNAME CTWTR.
...
IEF196I DSNAME=SYS1.SC30.OMPA.CTRACE,VOL=SER=COMST2,UNIT=3390,
IEF196I SPACE=(CYL,10),DISP=(NEW,
IEF196I          CATLG),DSORG=PS
...
ITT110I INITIALIZATION OF CTRACE WRITER CTWTR COMPLETE.

```

2. Activate ctrace with the omproute options, as shown in Example 5-52.

Example 5-52 TRACE CT command flow

```

TRACE CT,ON,COMP=SYSTCPRT,SUB=(OMPA)
*011 ITT006A SPECIFY OPERAND(S) FOR TRACE CT COMMAND.
R 11,OPTIONS=(ALL),END
IEE600I REPLY TO 011 IS;OPTIONS=(ALL),END
ITT038I ALL OF THE TRANSACTIONS REQUESTED VIA THE TRACE CT COMMAND
WERE SUCCESSFULLY EXECUTED.
IEE839I ST=(ON,0001M,00004M) AS=ON BR=OFF EX=ON MO=OFF MT=(ON,024K)
      ISSUE DISPLAY TRACE CMD FOR SYSTEM AND COMPONENT TRACE STATUS
      ISSUE DISPLAY TRACE,TT CMD FOR TRANSACTION TRACE STATUS

```

3. Modify the trace or debug trace levels as needed, issuing one or both the following commands, as shown in Example 5-53:

- a. **modify omp_proc,trace=x**
- b. **modify omp_proc,debug=x**

Example 5-53 Modify the omproute to use the desired trace and debug levels

```
F OMPA,TRACE=1
EZZ7866I OMPA MODIFY COMMAND ACCEPTED
F OMPA,DEBUG=2
EZZ7866I OMPA MODIFY COMMAND ACCEPTE
```

4. Reproduce the problem.
5. Stop the CTRACE by issuing the command in Example 5-54.

Example 5-54 Stopping CTRACE

```
TRACE CT,OFF,COMP=SYSTCPRT,SUB=(OMPA)
ITT038I ALL OF THE TRANSACTIONS REQUESTED VIA THE TRACE CT COMMAND
WERE SUCCESSFULLY EXECUTED.
IEE839I ST=(ON,0001M,00004M) AS=ON BR=OFF EX=ON MO=OFF MT=(ON,024K)
        ISSUE DISPLAY TRACE CMD FOR SYSTEM AND COMPONENT TRACE STATUS
        ISSUE DISPLAY TRACE,TT CMD FOR TRANSACTION TRACE STATUS
```

6. Save the trace contents into the trace file created by the CTWTR procedure, by executing the command in Example 5-55.

Example 5-55 Saving the trace contents

```
TRACE CT,ON,COMP=SYSTCPRT,SUB=(OMPA)
R 12,WTR=DISCONNECT,END
IEE600I REPLY TO 012 IS;WTR=DISCONNECT,END
ITT038I ALL OF THE TRANSACTIONS REQUESTED VIA THE TRACE CT COMMAND
WERE SUCCESSFULLY EXECUTED.
IEE839I ST=(ON,0001M,00004M) AS=ON BR=OFF EX=ON MO=OFF MT=(ON,024K)
        ISSUE DISPLAY TRACE CMD FOR SYSTEM AND COMPONENT TRACE STATUS
        ISSUE DISPLAY TRACE,TT CMD FOR TRANSACTION TRACE STATUS
```

7. Stop the external writer procedure CTWTR by issuing the command shown in Example 5-56.

Example 5-56 Stopping CTWTR

```
TRACE CT,WTRSTOP=CTWTR
ITT038I ALL OF THE TRANSACTIONS REQUESTED VIA THE TRACE CT COMMAND
WERE SUCCESSFULLY EXECUTED.
IEF196I AHL904I THE FOLLOWING TRACE DATASETS CONTAIN TRACE DATA :
IEF196I          SYS1.SC30.OMPA.CTRACE
AHL904I THE FOLLOWING TRACE DATASETS CONTAIN TRACE DATA : 404
          SYS1.SC30.OMPA.CTRACE
IEE839I ST=(ON,0001M,00004M) AS=ON BR=OFF EX=ON MO=OFF MT=(ON,024K)
        ISSUE DISPLAY TRACE CMD FOR SYSTEM AND COMPONENT TRACE STATUS
        ISSUE DISPLAY TRACE,TT CMD FOR TRANSACTION TRACE STATUS
ITT111I CTRACE WRITER CTWTR TERMINATED BECAUSE OF A WTRSTOP REQUEST.
IEF196I IEF142I CTWTR CTWTR - STEP WAS EXECUTED - COND CODE 0000
```

8. Change the OMPROUTE Debug and Trace level, as shown in Example 5-57, to avoid performance problems using the MODIFY command.

Example 5-57 Modifying Debug and Trace level

```
F OMPA,TRACE=0
EZZ7866I OMPROUTE MODIFY COMMAND ACCEPTED
F OMPA,DEBUG=0
EZZ7866I OMPROUTE MODIFY COMMAND ACCEPTED
```

After these steps, the trace file must be formatted, using the following IPCS command, into the IPCS Subcommand screen (option 6), as shown in Example 5-58.

Example 5-58 Formatting the OMPROUTE CTRACE

```
CTRACE COMP(SYSTCPRT) FULL
```

The next display shows the omproute debug entries entries, as shown in Example 5-59.

Example 5-59 Sample of formatted OMPROUTE CTRACE

```
COMPONENT TRACE FULL FORMAT
SYSNAME(SC30)
COMP(SYSTCPRT)
**** 09/28/2010
```

SYSNAME	MNEMONIC	ENTRY ID	TIME STAMP	DESCRIPTION
SC30	DEBUGTRC	00060001	21:02:16.210715	Trace Message
09/28 17:02:16 EZZ7878I OSPF Version: 2				Packet Length: 56
=====00089D17				
SC30	DEBUGTRC	00060001	21:02:16.210775	Trace Message
09/28 17:02:16 EZZ7908I Received packet type 1 from 10.1.5.12				
=====00089D21				
SC30	OPACKET	00020001	21:02:16.212164	OSPF RECVFROM PEEK
ASID...0053	TCB...007E60D0		JOBN...OMPA	
MODID...EZAORORT	CID...00000009		REG14..13F74178	
ADDR...14BF1E18	LEN...00000014		OSPF PEEK Packet Buffer	
000000 4500004C 54D50000 01597672	0A01040C	E0000005		
...<.N.....\...				
ADDR...14BF1E30	LEN...00000010		OSPF PEEK RECVFROM from address	
000000 00020000 0A01040C 0A01040B	00000006			
.....				
ADDR...14BF1E08	LEN...00000004		OSPF PEEK RECVFROM size	
000000 00000014				
....				
=====00089D22				

For more information about OMPROUTE diagnosis, refer to *z/OS Communications Server: IP Diagnosis Guide*, GC31-8782.

5.7 Additional information

For more details on these topics, refer to:

- ▶ *z/OS Communications Server: IP Configuration Guide*, SC31-8775
- ▶ *z/OS Communications Server: IP Configuration Reference*, SC31-8776
- ▶ *z/OS Communications Server: IP Diagnosis Guide*, GC31-8782

VLAN and Virtual MAC support

Virtual LAN (VLAN) technology is becoming more important in network planning for many customers. A VLAN is a configured logical grouping of nodes using switches. Nodes on a VLAN can communicate as though they were on the same LAN.

You need a switch to communicate across VLANs, but typically separate VLANs are in separate IP subnets; therefore, you often need a router to communicate across VLANs.

Virtual Medium Access Control (VMAC) support for z/OS Communications Server is a function that affects the operation of an OSA interface at the OSI layer 2 level. This is the Data Link Control (DLC) layer with its sub-layer Medium Access Control (MAC) layer.

This chapter discusses the following topics.

Section	Topic
6.1, "Virtual MAC overview" on page 266	The VMAC concept, and the environment in which it can be used
6.2, "Virtual MAC implementation" on page 269	An implementation example of VMACs
6.3, "Virtual LAN overview" on page 274	VLAN basics
6.4, "VLAN implementation on z/OS" on page 275	Single VLAN and multiple VLAN implementation scenarios on z/OS

6.1 Virtual MAC overview

Prior to the introduction of the *virtual* MAC function, an OSA interface only had one MAC address. This restriction caused problems when using load balancing technologies in conjunction with TCP/IP stacks that share OSA interfaces. The single MAC address of the OSA also causes a problem when using TCP/IP stacks as a forwarding router for packets destined to unregistered IP addresses.

VMAC support enables an OSA interface to have not only a physical MAC address, but also many distinct virtual MAC addresses for each device or interface in a stack. That is, each stack can define up to eight VMACs per protocol (IPv4 or IPv6) for each OSA interface.

Using VMACs, forwarding decisions in the OSA can be made without having to involve the OSI layer 3 level (network layer / IP layer). From a LAN perspective, the OSA interface with a VMAC appears as a dedicated device or interface to a TCP/IP stack. Packets destined for a TCP/IP stack are identified by an assigned VMAC address and packets sent to the LAN from the stack use the VMAC address as the source MAC address. This means that all IP addresses associated with a TCP/IP stack are accessible using their own VMAC address, instead of sharing a single physical MAC address of an OSA interface.

6.1.1 Why use virtual MACs

A shared OSA environment can be a challenge in certain network designs, and it requires careful planning when selecting the correct TCP/IP stacks to act as routers.

As mentioned in “OSA-Express router support” on page 123, the PRIRouter and SECRouter functions enable routing through a TCP/IP stack to IP addresses that are not registered in the OSA. The stack that has the OSA interface defined with PRIRouter will receive packets destined for IP addresses that do not reside in the given stack. The stack will then forwards the packets to the next hop.

Only one PRIRouter can be defined per OSA interface, although multiple SECRouters can be defined to an OSA interface for other TCP/IP routing stacks. However, only one SECRouter function can take over services if the PRIRouter is not available. If the first SECRouter function is not available, then the next defined SECRouter will forward IP packets to the associated stack. This means the OSA interface cannot serve multiple TCP/IP routing stacks concurrently even with the use of the PRIRouter and SECRouter functions.

Another challenge with shared OSA interfaces is one that requires load balancing of traffic across multiple TCP/IP stacks and IP addresses. For example, certain load balancing technologies use a concept of distributing packets to the appropriate adjacent systems based on knowledge of the MAC address.

We use load balancing (LB) with Sysplex Distributor to illustrate this challenge. If there is a shared OSA environment, the MAC address is attached to the Sysplex Distributor and to the selected target system. However, the target IP address can reside on a system other than the Sysplex Distributor.

As a result, the LB forwarding agent sends the packets to be distributed to the OSA’s physical MAC address, but the OSA only knows to send the information to the system that has registered the target address; it does not know to forward the information to the actual target stack. Mechanisms that are in place to overcome this challenge are Generic Resource Encapsulation (GRE) and Network Address Translation (NAT).

VMAC is a solution for both these problems and we recommend defining VMAC whenever multiple TCP/IP stacks share an OSA interface. VMAC support can provide the following:

- ▶ Allow for multiple concurrent TCP/IP routing stacks sharing an OSA interface
- ▶ Simplify the LAN infrastructure
- ▶ Eliminate the need for PRIRouter/SECRouter
- ▶ Improve outbound routing
- ▶ Improve IP workload balancing
- ▶ Remove the dependency on GRE and NAT

Note that there are two modes that can be used with load balancing technologies:

- ▶ *Directed mode* is where the load balancer converts the destination IP address (cluster IP address) to an IP address owned by the target system, using NAT. When IP packets from the target system are sent back to the clients, the load balancer converts the source IP address back to the cluster IP address. Therefore, the packets must return through the same load balancer that will recognize the changes and do the reverse mapping to ensure packets can flow from the original destination to the original source.
- ▶ *Dispatch mode* does not convert IP addresses, therefore eliminating the need for performing NAT. This mode requires VMAC support if the target stacks share the OSAs. In addition, all target applications must bind to the IP address specified by INADDR_ANY (or in6addr_any for IPv6), and the cluster IP address must be defined to the stack. The cluster IP address must not be advertised through a dynamic routing protocol. Otherwise, some systems might not have work routed to them. This can be done by defining the cluster IP address in the HOME list as a loopback address.

For more information regarding load balancing modes (directed and dispatch), refer to *z/OS Communications Server: IP Configuration Guide*, SC31-8775.

6.1.2 Virtual MAC concept

Figure 6-1 depicts how the definition of VMACs in the TCP/IP stacks gives the appearance of having a dedicated OSA interface on each stack. When packets arrive at the shared OSA interface, the individual VMAC assignments allow the packets to be forwarded directly to the correct stack. In the example shown, no individual stack needs to be defined as a primary or secondary router, thus offloading this function from a TCP/IP stack.

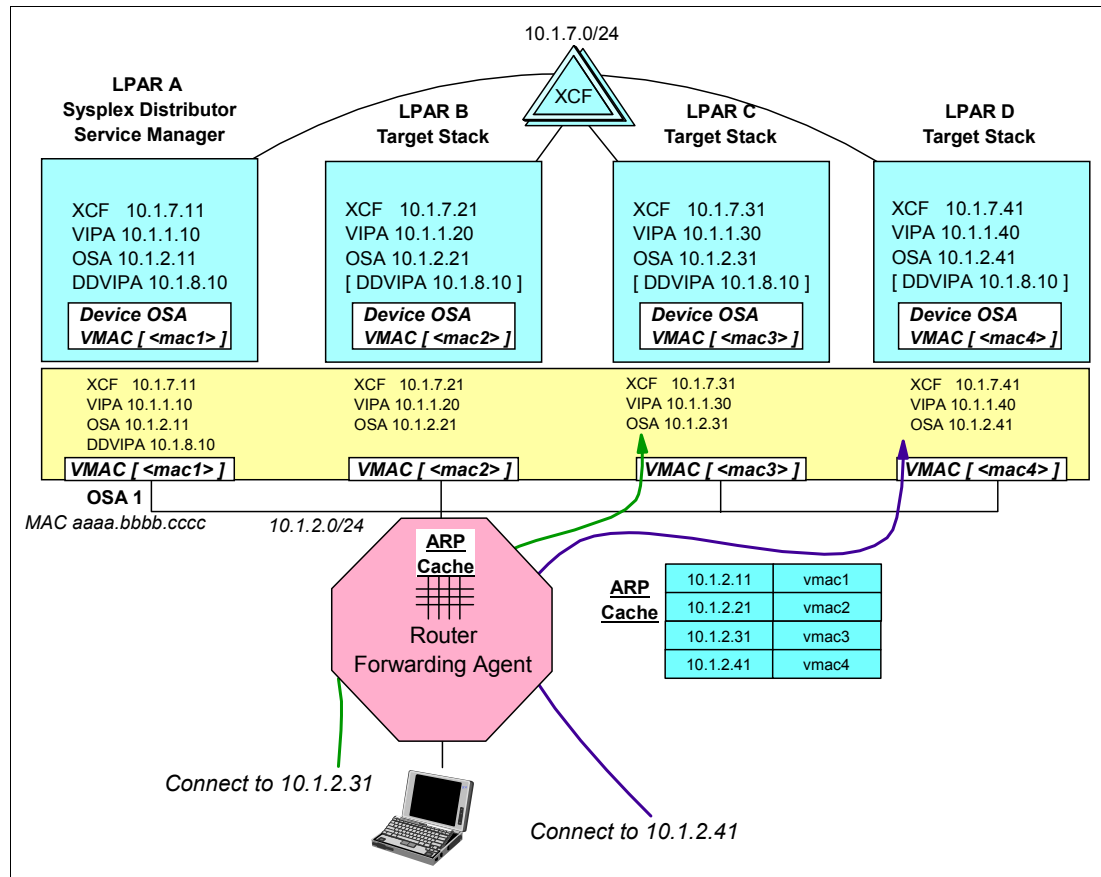


Figure 6-1 Forwarding packets to VMAC targets

This simplifies a shared OSA configuration significantly. Defining VMACs has very little administrative overhead. It is also an alternative to GRE or NAT when load balancing technologies are used. In Figure 6-1, the Dynamic VIPA targets are found without the use of GRE and without routing through the Sysplex Distributor. One of the options for defining VMACs permits the OSA to bypass IP address lookup. As a result, when the packet arrives at the correct VMAC, it is routed to the stack even though the DDVIPA is not registered in the OAT.

For IPV6, TCP/IP uses the VMAC address for all neighbor discovery address resolution flows for that stack's IP addresses, and likewise uses the VMAC as the source MAC address for all IPv6 packets sent from that stack. Again, from a LAN perspective, the OSA interface with a VMAC appears as a dedicated device to that stack.

Note: VMAC definitions on a device in a TCP/IP stack override any NONRouter, PRIRouter, or SECRouter parameters on devices in a TCP/IP stack. If necessary, selected stacks on a shared OSA can define the device with VMAC and others can define the device with PRIRouter and SECRouter capability.

6.1.3 Virtual MAC address assignment

The VMAC address can be defined in the stack, or generated by the OSA. If generated by the OSA, it is guaranteed to be unique from all other physical MAC addresses and from all other VMAC addresses generated by any OSA-Express feature.

Note: We recommend letting the OSA generate the VMACs instead of assigning an address in the TCP/IP profile. If VMACs are defined in the LINK statement, they must be defined as locally administered MAC addresses, and should be unique to the LAN on which they reside.

The same VMAC can be defined for both IPv4 and IPv6 usage, or a stack can use one VMAC for IPv4 and one for IPv6. Also, a VLAN ID can be associated with an OSA-Express device or interface defined with a VMAC.

6.2 Virtual MAC implementation

In this section, we show a scenario using VMAC as a replacement for PRIRouter and SECRouter. However the same implementation would apply to an environment using load balancing technologies. For details regarding load balancing technologies, refer to *Communications Server for z/OS V1R12 TCP/IP Implementation Volume 3: High Availability, Scalability, and Performance*, SG24-7898.

When implementing VMAC support, keep in mind the following points:

- ▶ The VMAC function is only available for OSA interfaces configured in QDIO mode.
- ▶ Each stack can define one VMAC per protocol (IPv4 or IPv6) for each OSA interface.
- ▶ If a VMAC is defined, the stack will not receive any packets destined to the physical MAC.
- ▶ VLAN IDs also apply to VMACs such as physical MACs.
- ▶ Allow the OSA to generate VMAC addresses.
- ▶ When configuring VMACs to solve load balancing issues, remember to:
 - Remove GRE tunnels as appropriate.
 - Change external load balancer configurations (such as directed mode to dispatch mode).

Note: VMAC support is only available with the IBM System z9, z10 and z196 Enterprise Class servers.

6.2.1 IP routing when using VMAC

In our scenario, as illustrated in Figure 6-2, we define VMACs to make two TCP/IP stacks act as forwarding stacks to route unregistered IP addresses, using OMPRoute. TCPIPA and TCPIPC share an OSA interface. We configured TCPIPA to forward packets to TCPIPB, and we configured TCPIPC to forward packets to TCIPID.

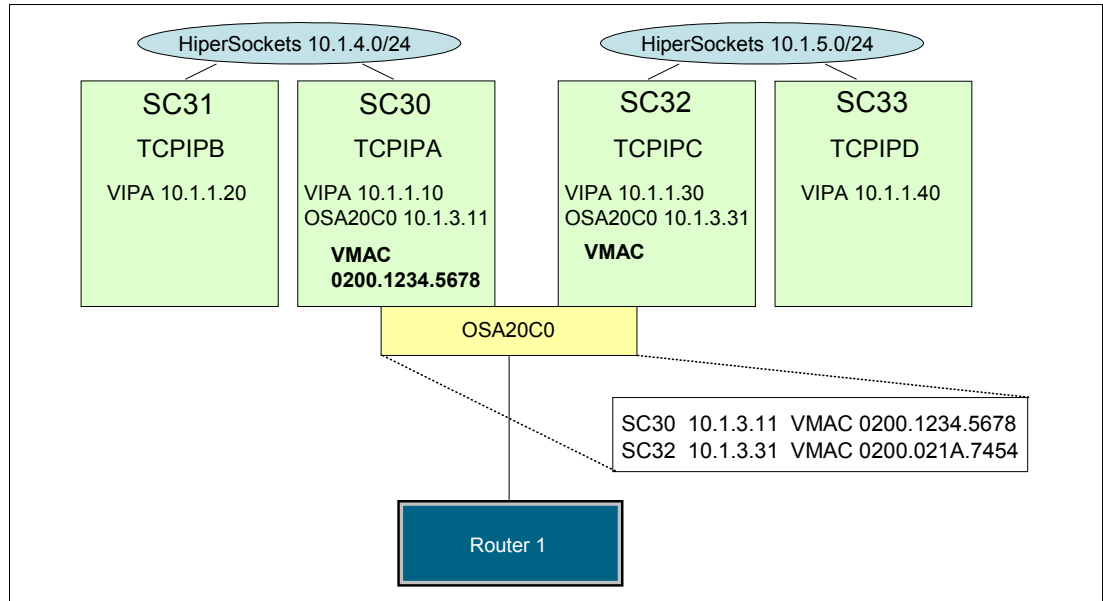


Figure 6-2 IP Routing using VMAC

We omitted the DEVICE, LINK, and HOME statements for OSA20C0 on TCPIPB and TCIPID, and modified the IP routing definitions on all stacks.

Note that Figure 6-2 is used only for demonstration purposes. We do *not* recommend implementing any configuration with single-points-of-failure.

Configuring the VMAC

The VMAC is defined on the LINK statement in the TCP/IP profile. Example 6-1 and Example 6-2 on page 271 show the VMAC definitions for TCPIPA and TCPIPC. In our example, we defined VMAC for OSA with VLAN ID. However, VLAN ID is not a prerequisite.

Example 6-1 Device and link statements: VMAC definition for TCPIPA

```

DEVICE OSA20C0    MPCIPA
LINK   OSA20C0L   IPAQENET    OSA20C0 VLANID 11 VMAC 020012345678 1
DEVICE IUTIQDF4   MPCIPA
LINK   IUTIQDF4L IPAQIDIO     IUTIQDF4
DEVICE VIPA1      VIRTUAL 0
LINK   VIPA1L     VIRTUAL 0    VIPA1
    
```

If VMAC is defined without a MAC address **2**, then OSA generates a VMAC using a part of the “burned-in” MAC address of the OSA. You can also specify the MAC address for VMAC **1**. If you decide to specify a MAC address, it must be a locally administered address, which means bit 6 of the first byte is 1 and bit 7 of the first byte is 0.

Example 6-2 Device and link statements: VMAC definition for TCPIPC

```

DEVICE OSA20C0      MPCIPA
LINK   OSA20COL     IPAQENET      OSA20C0 VLANID 11 VMAC 2
DEVICE IUTIQDF5     MPCIPA
LINK   IUTIQDF5L    IPAQIDIO      IUTIQDF5
DEVICE VIPA1        VIRTUAL 0
LINK   VIPA1L       VIRTUAL 0     VIPA1

```

There is no need to define PRIRouter or SECRouter on the DEVICE statement. When VMAC is specified on LINK statement, PRIRouter or SECRouter is ignored.

Note: z/OS Communications Server has been enhanced and IPV4 interfaces VLANs can be defined using the INTERFACE statement. More details are available in “INTERFACE” on page 81.

6.2.2 Verification

We verified that VMAC was correctly defined in TCPIPA (see Example 6-3). We specified a MAC address 1 for the OSA in TCPIPA, so VMACORIGIN is CFG 2.

Example 6-3 Display VMAC on TCPIPA

```

D TCPIP,TCPIPA,N,DEV
INTFNAME: OSA20C0I          INTFTYPE: IPAQENET  INTFSTATUS: READY
PORTNAME: OSA20C0          DATAPATH: 20C2      DATAPATHSTATUS: READY
CHPIDTYPE: OSD
SPEED: 0000001000
IPBROADCASTCAPABILITY: NO
VMACADDR: 020012345678 1 VMACORIGIN: CFG 2 VMACROUTER: ALL
ARPOFFLOAD: YES           ARPOFFLOADINFO: YES
CFGMTU: 1492              ACTMTU: 1492
IPADDR: 10.1.3.11/24
VLANID: 11                VLANPRIORITY: DISABLED
DYNVLANREGCFG: NO         DYNVLANREGCAP: YES
READSTORAGE: GLOBAL (4096K)
INBPERF: BALANCED
CHECKSUMOFFLOAD: YES      SEGMENTATIONOFFLOAD: YES
SECCLASS: 255             MONSYSPLEX: NO
ISOLATE: NO               OPTLATENCYMODE: NO
MULTICAST SPECIFIC:
MULTICAST CAPABILITY: YES
GROUP                     REFCNT             SRCFLTMD
-----
224.0.0.1                0000000001    EXCLUDE
SRCADDR: NONE
INTERFACE STATISTICS:
BYTESIN                   = 0
INBOUND PACKETS           = 0
INBOUND PACKETS IN ERROR  = 0
INBOUND PACKETS DISCARDED = 0
INBOUND PACKETS WITH NO PROTOCOL = 0
OUTBOUND PACKETS          = 1
OUTBOUND PACKETS IN ERROR = 0

```

We verified that VMAC was correctly defined in TCPIPC (see Example 6-4). Because we did not specify a MAC address for the OSA in TCPIPC, the OSA generated the MAC address **3**. Because this is a OSA-generated MAC address, VMACORIGIN is OSA **4**.

Example 6-4 Display VMAC on TCPIPC

```

D TCPIP,TCPIP,N,DEV
INTFNAME: OSA20C0I          INTFTYPE: IPAQENET  INTFSTATUS: READY
PORTNAME: OSA20C0          DATAPATH: 20C2      DATAPATHSTATUS: READY
CHIPIDTYPE: OSD
SPEED: 0000001000
IPBROADCASTCAPABILITY: NO
VMACADDR: 02000E776C05 1 VMACORIGIN: OSA 2 VMACROUTER: ALL
ARPOFFLOAD: YES            ARPOFFLOADINFO: YES
CFGMTU: 1492               ACTMTU: 1492
IPADDR: 10.1.3.11/24
VLANID: 11                 VLANPRIORITY: DISABLED
DYNVLANREGCFG: NO          DYNVLANREGCAP: YES

DEVNAME: OSA20C0           DEVTYPE: MPCIPA
DEVSTATUS: READY
LNKNAME: OSA20C0L          LNKTYPE: IPAQENET  LNKSTATUS: READY
NETNUM: N/A QUESIZE: N/A   SPEED: 0000000100
IPBROADCASTCAPABILITY: NO
VMACADDR: 0200021A7454 3 VMACORIGIN: OSA 4 VMACROUTER: ALL

```

We can also see the VMAC in the OSA Address Table (OAT) queried by OSA/SF (see Example 6-5). OSA registers all IP addresses (including VIPA) in the TCP/IP stack, and maps them to the VMAC address.

Example 6-5 Display OAT queried with OSA/SF

```

Local MAC address -----> 00145E776C05 5
Universal MAC address -----> 00145E776C05
*****
Image 1.1 (A11      ) CULA 0
00(20C0)* MPC      N/A      OSA20C0 (QDIO control)          SIU ALL
02(20C2) MPC 00 No4 No6 OSA20C0 (QDIO data)          SIU ALL
                   VLAN 11 (IPv4)

                   Group Address  Multicast Address
                   01005E000001    224.000.000.001

                   VMAC            IP address
HOME      02000E776C05          010.001.003.011 7

03(20C3) MPC 00 No4 No6 OSA20C0 (QDIO data)          SIU ALL
                   VLAN 11 (IPv4)

                   Group Address  Multicast Address
                   01005E000001    224.000.000.001
                   01005E000005    224.000.000.005

                   VMAC            IP address
HOME      02000F776C05          010.001.003.023 7

```

```

                                Image 1.6 (A16      ) CULA 0
00(20C0)* MPC          N/A      OSA20C0 (QDIO control)      SIU ALL
02(20C2) MPC 00 No4 No6 OSA20C0 (QDIO data)      SIU ALL
                                VLAN 10 (IPv4)

                                Group Address  Multicast Address
                                01005E000001   224.000.000.001
                                01005E000005   224.000.000.005

                                VMAC          IP address
HOME 020007776C05 7 010.001.002.030
HOME 020007776C05 010.001.002.033
```

Note that the last three bytes of the OSA-generated VMAC 7 are identical to that of the universal MAC address (“burned-in” address) of the OSA 5. The first byte of the OSA-generated VMAC is always 02, in order to make the VMAC a locally administered address. To make the VMAC unique among all TCP/IP stacks, the second and third bytes are used as a counter that is incremented each time OSA generates a MAC address.

Example 6-6 shows the ARP cache of the router. IP address 10.1.3.11 in TCPIPA is mapped to the VMAC defined in TCPIPA 8, and IP address 10.1.3.31 in TCPIPC is mapped to the VMAC defined in TCPIPC 9.

Example 6-6 Display ARP cache in Router 1

```
Router1#sh arp
Internet 10.1.3.11          10 0200.1234.5678 8 ARPA  Vlan11
Internet 10.1.3.31          20 0200.021a.7454 9 ARPA  Vlan11
```

Note that each IP address is mapped to a different MAC address, even if these stacks share the same OSA interface. OSA responds to ARP requests for all registered IP addresses by using a VMAC instead of a “burned-in” MAC address.

According to the routing table, the router chooses 10.1.3.11 as the next hop for destination address 10.1.1.20, and chooses 10.1.3.31 as the next hop for destination address 10.1.1.40. The router forwards the packet with the destination IP address 10.1.1.20 to the destination MAC address 0200.1234.5678. When the packet reaches the OSA interface, OSA forwards the packet to TCPIPA, because OSA knows the VMAC 200.1234.5678 is mapped to TCPIPA. The same can be said for the TCPIPC VMAC.

Example 6-7 shows that the two stacks (TCPIPA and TCPIPC) sharing one OSA interface are able to route packets correctly.

Example 6-7 Display traceroute from Router 1

```
Router1#traceroute 10.1.1.20
 1 10.1.3.11 0 msec 0 msec 0 msec
 2 10.1.1.20 0 msec 0 msec 0 msec

Router1#traceroute 10.1.1.40
 1 10.1.3.31 4 msec 0 msec 0 msec
 2 10.1.1.40 0 msec 0 msec 0 msec
```

6.3 Virtual LAN overview

A virtual LAN (VLAN) is the grouping of workstations, independent of physical location, that have a common set of requirements. VLANs have the same attributes as physical LANs, although they might not be located physically on the same LAN segment.

A VLAN configuration provides several benefits:

- ▶ VLANs can improve network performance by reducing traffic on a physical LAN. VLANs can enhance security by isolating traffic.
- ▶ VLANs provide more flexibility in configuring networks.
- ▶ VLANs can be used to increase link optimization by allowing networks to be organized for optimum traffic flow through implementation of network segregation and a quality of service policy.
- ▶ VLANs can be used to increase bandwidth and reduce overhead.

6.3.1 Types of connections

VLANs operate by defining switch ports as members of virtual LANs. Devices on a VLAN can use three types of connections, based on whether the connected devices are VLAN-aware or VLAN-unaware. VLAN-aware devices understand VLAN memberships (which users belong to a particular VLAN) and VLAN formats.

Ports used to attach VLAN-unaware equipment are called *access ports*, while ports used to connect to other switches or VLAN-aware servers are known as *trunk ports*. Network frames generated by VLAN-aware equipment are marked with a *tag*, which identifies the frame to the VLAN.

The types of connections include:

- ▶ Trunk mode

Trunk mode indicates that the switch should allow all VLAN ID tagged packets to pass through the switch port without altering the VLAN ID. This mode is intended for servers that are VLAN capable. It filters and processes all VLAN ID tagged packets. In trunk mode, the switch expects to see VLAN ID tagged packets inbound to the switch port.

- ▶ Access mode

Access mode indicates that the switch should filter on specific VLAN IDs and only allow packets that match the configured VLAN IDs to pass through the switch port. The VLAN ID is then removed from the packet before it is sent to the server. That is, VLAN ID filtering is controlled by the switch. In access mode, the switch expects to see packets without VLAN ID tags inbound to the switch port.

- ▶ Hybrid mode

Hybrid mode is a combination of the previous two modes. This mode defines a port where both VLAN-aware and VLAN-unaware devices are attached. A hybrid port can have both tagged and untagged frames.

6.4 VLAN implementation on z/OS

In this section, we discuss single and multiple VLANs and their configuration.

6.4.1 Single VLAN per OSA

Figure 6-3 shows one physical LAN subdivided into two VLANs, VLAN A and VLAN B. To configure VLAN for an OSA-Express in QDIO mode, you specify a VLAN ID in the TCP/IP profile. In earlier releases, the only way to access multiple VLANs from a given z/OS stack was to use multiple OSAs.

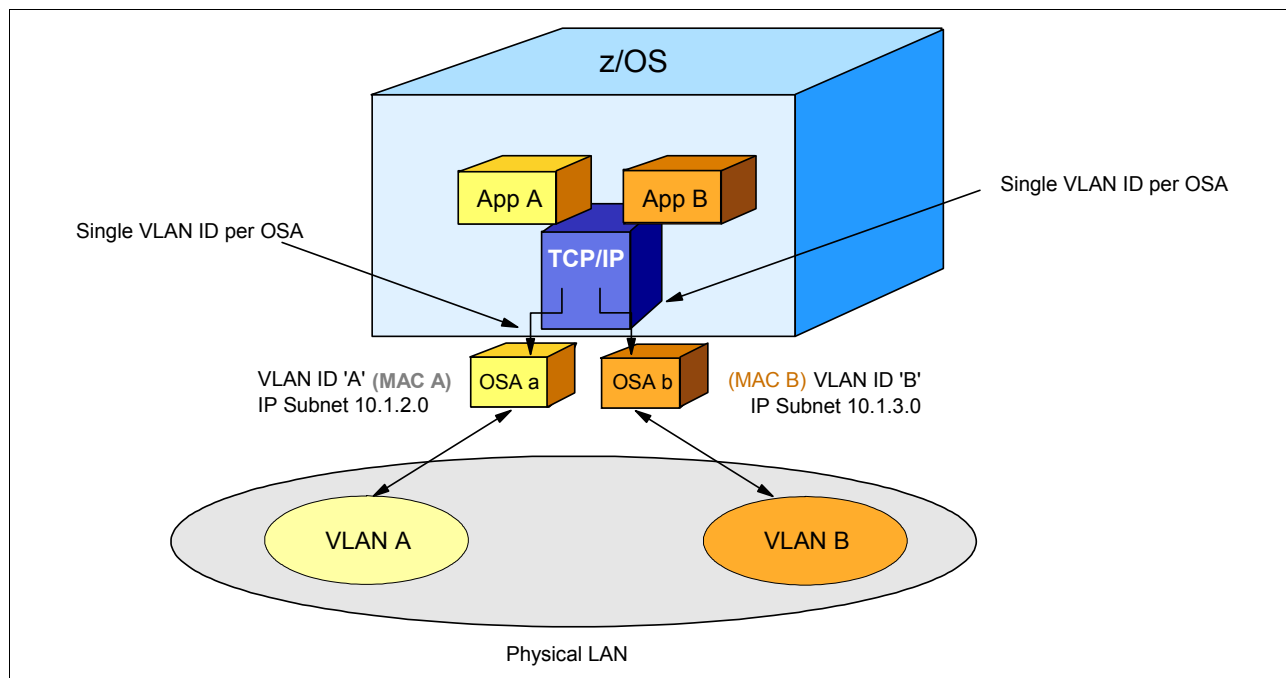


Figure 6-3 Single VLAN per OSA

The z/OS stack registers the VLAN ID to OSA, which means that the OSA:

- Appends a Layer 2 VLAN tag with this VLAN ID on all outbound packets (For IPv6 unicast packets, the stack, not the OSA, appends the VLAN tags.)
- Filters out any inbound packets that have a VLAN tag containing a different VLAN ID

VLANs on a single footprint, as shown in Figure 6-3 on page 275, typically map to separate IP subnets. This one-to-one mapping is not a requirement, because the same IP subnet (a Layer 3 construct) can be subdivided into separate VLANs. Likewise, separate IP subnets on the same footprint can be mapped to the same VLAN. Nevertheless, it is more common to assign a separate IP subnet to separate VLAN IDs, as shown in Figure 6-3 on page 275. The latter type of network design simplifies network topology and the planning of a Layer 3 routing infrastructure.

6.4.2 Multiple VLAN support

Figure 6-4 the multiple VLAN functions that are supported. Multiple VLANs can be configured for the same OSA-Express feature (up to eight for IPv4 and eight for IPv6) from the same z/OS stack. This is done by defining multiple interfaces to the same OSA-Express (one for each VLAN ID).

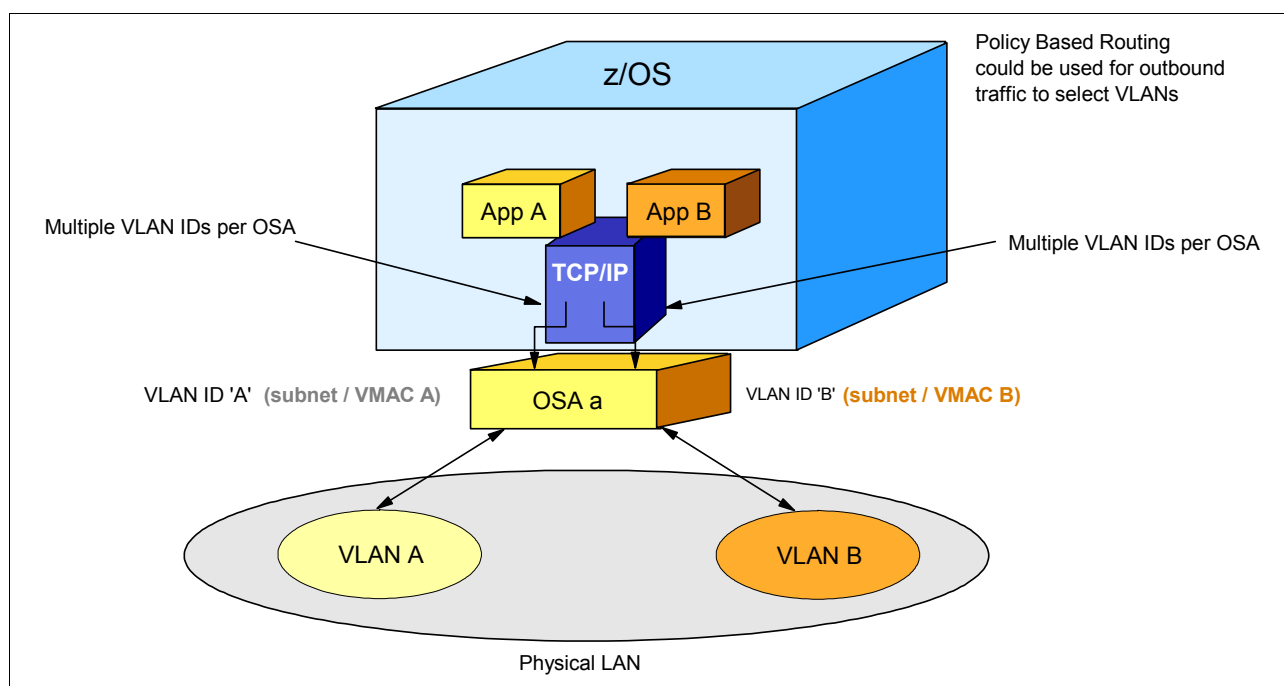


Figure 6-4 Multiple VLANs in z/OS

Multiple VLAN support provides:

- OSA port consolidation: The multiple VLAN function allows a customer to consolidate multiple OSAs (for example, three 1 Gb OSA ports) into a single OSA (for example, one 10 Gb OSA port) serving multiple VLANs.
- Server consolidation: The multiple VLAN function allows a customer to consolidate multiple application servers across multiple stacks into a single z/OS image where the traffic related to these servers is on unique VLANs.

- Improved QoS with policy-based routing: The policy-based routing function allows a z/OS stack to make routing decisions for IPv4 traffic that take into account additional criteria, such as job name, source port, destination port, protocol type (TCP or UDP), source IP address, NetAccess security zone, and a multilevel secure environment security label. It enables routing of traffic that meets certain criteria to one VLAN and traffic that meets different criteria to another VLAN.

Defining multiple VLANs

The INTERFACE statement in the TCPIP profile that was used to define IPV6 OSA interfaces, is extended to use the IPV4 QDIO OSA devices, as shown in Example 6-8.

Note: Non QDIO IPV4 network interfaces are defined using the old syntax of DEVICE and LINK.

Example 6-8 INTERFACE statement

```
INTERFACE OSAQDI0024
  DEFINE IPAQENET
  PORTNAME OSAQDI002
  SOURCEVIPAINTE VIPAV4      1
  IPADDR 10.1.1.1/24
  MTU 8992
  VLANID 200                  2
  VMAC ROUTEALL               3
```

For a detailed description of the INTERFACE statement refer to “INTERFACE” on page 81. In this example, the numbers correspond to the following information:

- 1.** For IPV4 source VIPA specify the VIPA LINKNAME. For IPV6 source VIPA specify the interface name specified on the VIRTUAL6 interface statement.
- 2.** This is the VLAN ID of the VLAN.
- 3.** Specifies that all IP traffic destined to the virtual MAC is forwarded by the OSA-Express device to the TCP/IP stack.

6.4.3 Multiple VLANs configuration guidelines

To define multiple interfaces to the same OSA express or define multiple VLANs on the same OSA express port or more then one OSA express port, follow these rules:

- Configure each IPv4 interface for the OSA-Express feature in the TCP/IP profile using the INTERFACE statement for IPAQENET rather than DEVICE/LINK/HOME. Configure each IPv6 interface using the INTERFACE statement for IPAQENET6.
- Configure a VLANID value on each IPv4 and each IPv6 INTERFACE statement for this OSA. Within each IP version, VLANID values must be unique.
- Configure the VMAC parameter on each of these INTERFACE statements with the default ROUTEALL attribute. The VMAC address can either be specified or OSA-generated. If you specify a VMAC address, it must be unique for each INTERFACE statement.

Note: By using the ROUTEALL attribute, you allow the interface to forward IP packets. You can use the ROUTELOCAL attribute if you do not want the interface to forward IP packets.

- ▶ Configure a unique subnet for each IPv4 interface for this OSA-Express feature using the subnet mask specification on the IPADDR parameter on the INTERFACE statement.
- ▶ To use multiple VLANs for an OSA port, you need to configure a separate interface to the OSA port for each VLAN. Each of these interfaces requires a separate DATAPATH device in the TRLE definition. Furthermore, each DATAPATH device requires a certain amount of fixed storage. See “VTAM considerations” on page 279.
- ▶ VLAN IDs must be unique on a single OSA port within a single stack. If you code multiple INTERFACES from one stack to the same OSA and do not configure a VLAN ID for one INTERFACE, the INTERFACE definition will be rejected.
- ▶ If one INTERFACE within a stack that is connecting to an OSA port is implemented with VLAN/VMAC, then *all* INTERFACES connecting to the same OSA Port within that stack must specify VLAN/VMAC.
- ▶ If more than one INTERFACE is defined for a particular IP version for a single OSA port within a stack, then the VLANID, VMAC and IP subnet values must be unique on each of the INTERFACE statements. If parallel interfaces are desired with the same IP subnet and same VLANID, then the parallel INTERFACE statements must be coded on different OSA ports.
- ▶ When a z/OS TCP/IP stack has access to multiple OSAs that are on the same physical LAN and when a VLAN ID is configured on any of the OSAs, it is recommended that this stack configure a VLAN ID for all OSAs on the same physical LAN. That is, do not mix interfaces configured with and without VLAN on the same physical network when a stack has access to the same LAN through multiple OSAs. Doing otherwise can cause problems with various ARP takeover scenarios.
- ▶ The multi-VLAN configuration rules apply only within a stack, that is, each stack on a shared OSA port is completely independent of any other stacks sharing the OSA port. Therefore, if you have one TCP/IP stack (at an earlier release) sharing an OSA port with a second TCP/IP stack (at a current release), the first stack can be configured to use the DEVICE/LINK statement for a single connection to a shared OSA port and the second stack can be configured to use the INTERFACE statement for any connections to the shared OSA port.
- ▶ A network switch can establish VLANIDs on some connections of a trunk port to a single OSA port. The switch can also configure other connections with what is called a *Native VLANID* or a *Default VLANID* on the same trunk port to the shared OSA port. If a single TCP/IP stack has configured multiple VLAN connections to the switch and one of those connections is to the Native VLAN, then the z/OS TCP/IP stack must set a VLANID for the Native VLAN on that connection. Do not mix Native VLAN and other VLANIDs on the same OSA port to the same TCP/IP stack.

Note: Some switch vendors use VLAN ID 1 as the default value when a VLAN ID value is not explicitly configured. It is recommended that you avoid the value of 1 when configuring a VLAN ID value. By convention the “Native VLANID” is often coded as “1”.

Source VIPA

Use the following guidelines when selecting a source VIPA:

- ▶ In earlier CS releases, for IPv4, when source VIPA is in effect, the stack selects a source VIPA based on the order of the home list (from the ordering of IP addresses in the HOME statement in the profile). So, for IPv4, the user controls source VIPA selection using the HOME statement.
- ▶ For IPv6, there is no HOME statement. The user controls source VIPA selection using the SOURCEVIPAINTERFACE parameter on the INTERFACE statement.

- ▶ The source VIPA selection for interfaces defined with the IPv4 INTERFACE statement works the same way as IPv6 (using the SOURCEVIPAINTERFACE parameter, which must point to the link name of an IPv4 static VIPA).
- ▶ For IPv4 interfaces defined using DEVICE/LINK, source VIPA selection continues to work based on the ordering of the home list.
- ▶ You can specify SOURCEVIPAINTERFACE for every VLAN you define. The VIPA IP address can be in the same or different subnet from the IP address of the OSA interface.

ARP processing

In QDIO mode, the OSA performs all Address Resolution Protocol (ARP) processing for IPv4. The z/OS stack informs the OSA of the IP addresses for which it should perform ARP processing. Because the z/OS stack also supports configurations where ARPs flow for VIPAs (which one might see on some flat network configurations using static routing), the stack also informs the OSA of the VIPAs for which it should perform ARP processing. OSA sends gratuitous ARPs for these IP addresses during interface takeover scenarios to provide fault tolerance.

If the OSA is defined using DEVICE/LINK statements, then the stack will inform OSA to perform ARP processing for all VIPAs in the home list, which can result in numerous unnecessary gratuitous ARPs for VIPAs in an interface takeover scenario. However, if using the IPv4 INTERFACE statement for IPAQENET, and a subnet mask is configured (non-0 *num_mask_bits*) on the IPADDR parameter of the INTERFACE statement, then the stack will inform OSA to perform ARP processing for a VIPA only if the VIPA is configured in the same subnet as the OSA.

VTAM considerations

The QDIOSTG VTAM start parameter specifies how much storage VTAM keeps available for all OSA QDIO devices. Each OSA express QDIO DATAPATH device consumes large amount of fixed storage. The QDIOSTG value can be overridden by using the READSTORAGE parameter on the IPAQENET LINK or the INTERFACE statement on the TCPIP profile. As every VLAN adds another OSA device (DATAPATH) and environment it is recommended in a multi-VLAN to use VTAM tuning statistics and evaluate the needs and storage.

6.4.4 Verification

We performed TCPIP device displays and retrieved the OSA address Table (OAT) to present how multiple VLANs are recognized by the system. Example 6-9 on page 280 shows the output of the TCPIP device display. We defined two VLANs and a source VIPA on the interface statement.

D TCPIP,TCPIPA,N,DEV,INTFN=OSA2080I **VLAN 10**
EZD0101I NETSTAT CS V1R12 TCPIPA 219
INTFNAME: OSA2080I INTFTYPE: IPAQENET INTFSTATUS: READY
PORTNAME: OSA2080 DATAPATH: 2082 DATAPATHSTATUS: READY
CHPIDTYPE: OSD
SPEED: 0000001000
IPBROADCASTCAPABILITY: NO
VMACADDR: 020010776873 VMACORIGIN: OSA VMACROUTER: LOCAL
ARPOFFLOAD: YES ARPOFFLOADINFO: YES
CFGMTU: 1492 ACTMTU: 1492
IPADDR: 10.1.2.11/24
VLANID: **10** VLANPRIORITY: DISABLED
DYNVLANREGCFG: NO DYNVLANREGCAP: YES
READSTORAGE: GLOBAL (4096K)
INBPERF: BALANCED
CHECKSUMOFFLOAD: YES SEGMENTATIONOFFLOAD: YES
SECCLASS: 255 MONSYSPLEX: NO
ISOLATE: NO OPTLATENCYMODE: NO
MULTICAST SPECIFIC:
MULTICAST CAPABILITY: YES

D TCPIP,TCPIPA,N,DEV,INTFN=OSA20C0I **VLAN 11**
EZD0101I NETSTAT CS V1R12 TCPIPA 678
INTFNAME: OSA20C0I INTFTYPE: IPAQENET INTFSTATUS: READY
PORTNAME: OSA20C0 DATAPATH: 20C2 DATAPATHSTATUS: READY
CHPIDTYPE: OSD
SPEED: 0000001000
IPBROADCASTCAPABILITY: NO
VMACADDR: 02000E776C05 VMACORIGIN: OSA VMACROUTER: ALL
ARPOFFLOAD: YES ARPOFFLOADINFO: YES
CFGMTU: 1492 ACTMTU: 1492
IPADDR: 10.1.3.11/24
VLANID: **11** VLANPRIORITY: DISABLED
DYNVLANREGCFG: NO DYNVLANREGCAP: YES
READSTORAGE: GLOBAL (4096K)
INBPERF: BALANCED
CHECKSUMOFFLOAD: YES SEGMENTATIONOFFLOAD: YES
SECCLASS: 255 MONSYSPLEX: NO
ISOLATE: NO OPTLATENCYMODE: NO
MULTICAST SPECIFIC:
MULTICAST CAPABILITY: YES

- The IP address and the subnet mask assigned to this VLAN **1**
- The VLAN ID **2**

Example 6-10 shows the OAT of a CHPID defined as multiple VLANs and source VIPA

Example 6-10 OAT of a CHPID defined as multiple VLANs and source VIPA

Image 1.1 (A11) CULA 0									
00(20C0)* MPC			N/A		OSA20C0	(QDIO control)		SIU ALL	
02(20C2) MPC	00	No4	No6		OSA20C0	(QDIO data)		SIU ALL	
		VLAN	11		(IPv4)				
		Group Address			Multicast Address				
		01005E000001			224.000.000.001				
		VMAC			IP address				
HOME		02000E776C05			010.001.003.011			3	
Image 1.1 (A11) CULA 0									
04(2084) MPC	00	No4	No6		OSA2080	(QDIO data)		SIU ALL	
		VLAN	10		(IPv4)				
		Group Address			Multicast Address				
		01005E000001			224.000.000.001				
		VMAC			IP address				
HOME		020011776873			010.001.002.023			1	
HOME		020011776873			010.001.002.025			2	

In this example, the numbers correspond to the following information:

- 1. Source VIPA address as defined on the Interface statement
- 2. VLAN 10 IP address and the assigned VMAC
- 3. VLAN 11 IP address and the assigned VMAC

Note: The same VMAC is assigned for the VLAN IP address and the source VIPA IP address.

Because VLAN 11 belongs to a different IP subnet mask from the source VIPA, the source VIPA is not displayed on this VLAN.

6.5 References

For more information about the VMAC function, refer to the following documentation:

- *z/OS Communications Server: IP Configuration Guide*, SC31-8775
- *z/OS Communications Server: IP Configuration Reference*, SC31-8776
- *Communications Server for z/OS V1R12 TCP/IP Implementation Volume 3: High Availability, Scalability, and Performance*, SG24-7898



Sysplex subplexing

In large sysplex environments, there can be strict security requirements to isolate access to certain VTAM nodes or TCP/IP stacks within the sysplex. A z/OS Communications Server function, called *subplexing*, provides this type of support. It enables the user to implement automatically controlled access to subplex groups.

As mentioned, the subplexing support is also for VTAM nodes. However, this chapter only describes subplexing for TCP/IP stacks. For information about VTAM subplexing, refer to *SNA Network Implementation Guide*, SC31-8777.

This chapter discusses the following topics.

Section	Topic
7.1, “Introduction” on page 284	The subplexing concept, and the environment in which it can be used
7.2, “Subplex environment” on page 286	Our TCP/IP subplexing environment
7.3, “Load Balancing Advisor and subplexing” on page 287	The Load Balancing Advisor allows any external load balancing solution to become sysplex aware
7.4, “Subplex implementation” on page 290	Implementation examples of TCP/IP subplexing

7.1 Introduction

Prior to subplexing, VTAM and TCP/IP sysplex functions were deployed sysplex-wide and users had to implement complex resource controls and disable many of the dynamic XCF and routing functions to support multiple security zones. For example, as shown in Figure 7-1, TCP/IP stacks access different networks with diverse security requirements within the same sysplex:

- ▶ In the upper configuration, two TCP/IP stacks in the left LPARs access an internal network. The TCP/IP stacks in the right two LPARs access the external network. Presumably, the security requirements would include isolating external traffic from the internal network. However, all TCP/IP stacks in the sysplex can dynamically establish connectivity with all the other TCP/IP stacks in the sysplex.
- ▶ In the lower configuration, TCP/IP stacks in the same LPAR have different security requirements. The first stack in each LPAR connects to the internal network, and the second stack connects to the external network. Through the IUTSAMEH connection, the two stacks in each LPAR can establish connectivity with each other dynamically and possibly violating security policies.

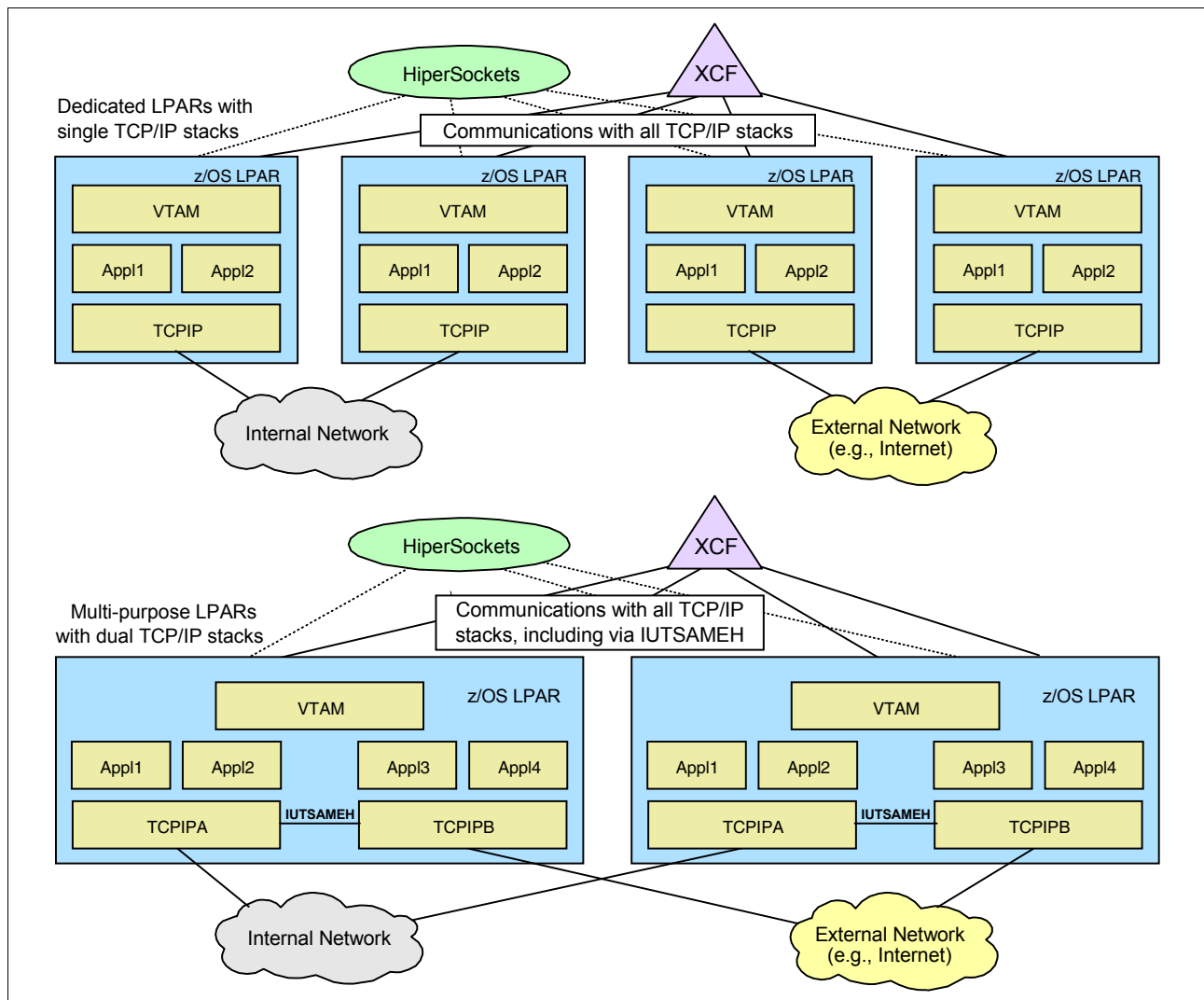


Figure 7-1 Sysplex connectivity: examples

With subplexing, you can build *security zones*. Only members within the same security zone can communicate with each other. Subplex members are VTAM nodes and TCP/IP stacks that are grouped in security zones to isolate communication.

Concept of subplexing

A *subplex* is a subset of a sysplex that consists of selected members. Those members are connected and they communicate through dynamic cross-system coupling facility (XCF) groups to each other, using the following methods:

- ▶ XCF links (for cross-system IP and VTAM connections)
- ▶ IUTSAMEH (for IP connections within an LPAR)
- ▶ HiperSockets (IP connections cross-LPAR in the same server)

Subplexes do not communicate with members outside the subset of the sysplex. For example, in Figure 7-2, TCP/IP stacks with connectivity to the internal network can be isolated from TCP/IP stacks connected to the external network using subplexing.

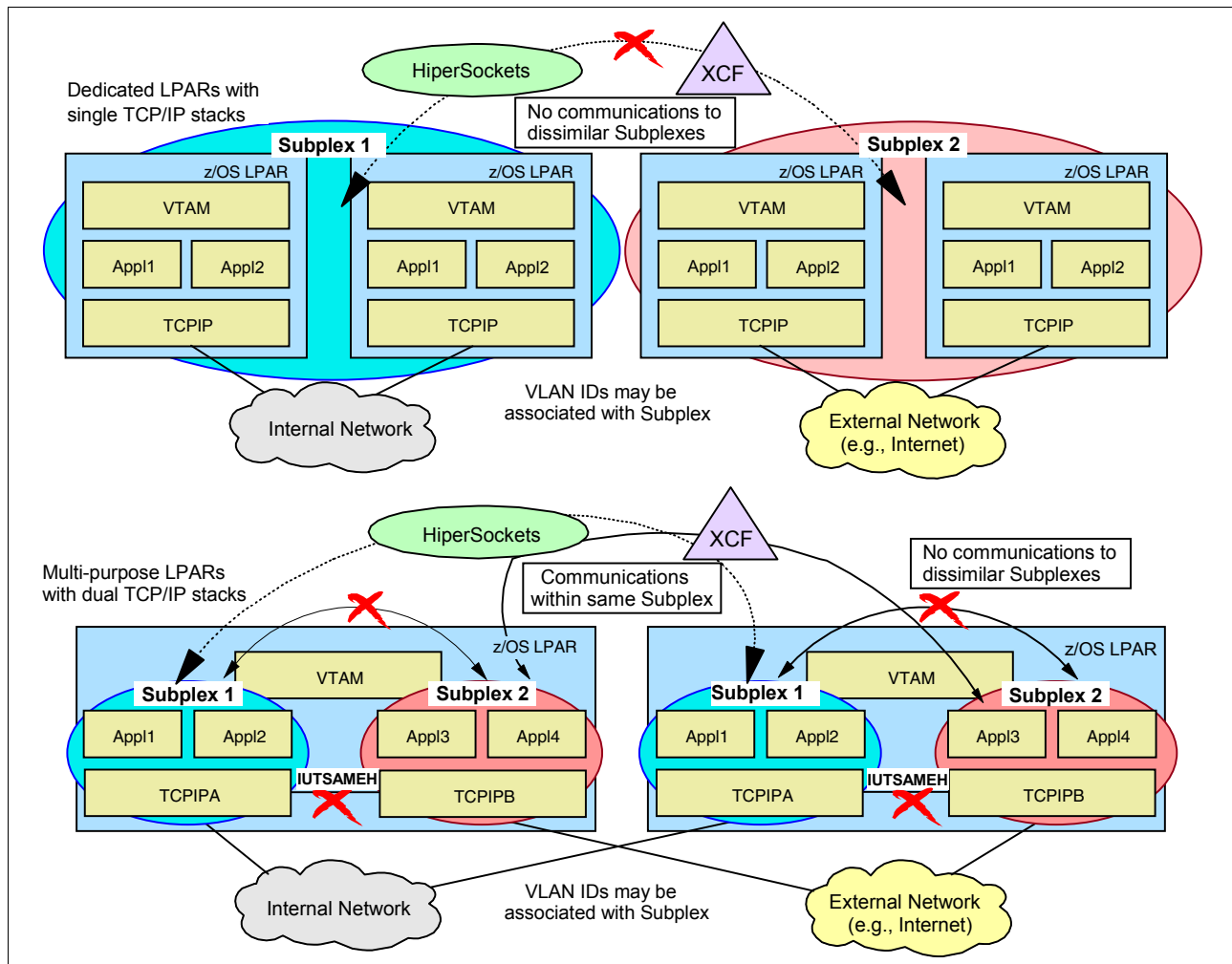


Figure 7-2 Subplexing multiple security zones

TCP/IP stacks are defined as members of a subplex group with a defined group ID. For example, in Figure 7-2 TCP/IP stacks within subplex 1 are able to communicate only with stacks within the same subplex group. They are not able to communicate with stacks in subplex 2.

In an environment where a single LPAR has access to internal and external networks through two TCP/IP stacks, those stacks are assigned to two different subplex group IDs. Even though IUTSAMEH is the communication method, it is controlled automatically through the association of subplex group IDs, thus creating two separate security zones within the LPAR.

Recommendation: Network connectivity provided through an OSA port in a multiple security zone environment should *not* be shared across different subplex groups. The OSA ports and HiperSockets connections should be physically isolated or logically separated using firewall and VLAN technologies.

7.2 Subplex environment

In this section we describe the environment used to demonstrate subplexing in a multiple security zone, based on Figure 7-2 on page 285. All LPARs in our scenarios were configured in a single server with multiple stacks for demonstration purposes only.

Note: Although there are specialized cases where multiple stacks per LPAR can provide value, we generally recommend implementing only one TCP/IP stack per LPAR whenever possible.

Figure 7-3 illustrates our TCP/IP subplexing environment with the following attributes:

- The first subplex is a VTAM subplex, which is not within the scope of this book. However, when defining only a TCP/IP subplex, a default VTAM subplex is defined automatically.

Note: A TCP/IP subplex uses VTAM XCF support for DYNAMICXCF connectivity. Therefore, a TCP/IP stack cannot span different VTAM subplexes.

- The second subplex was configured with TCP/IP C stacks running in LPARs A11 and A13, representing the internal subplex.
- The third subplex was configured with TCP/IP D stacks running in LPARs A13 and A16, representing the external subplex.

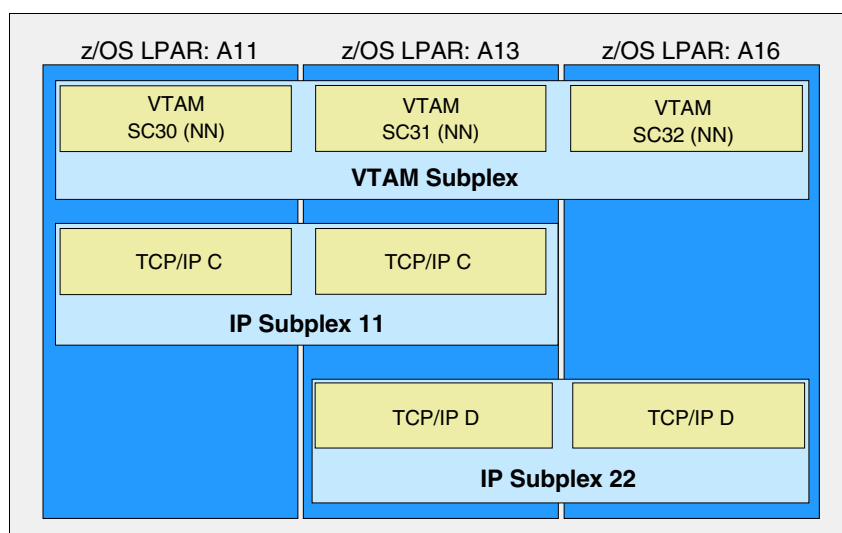


Figure 7-3 Our TCP/IP subplexing environment

We do not describe or discuss OSA connectivity in this chapter. For details regarding OSA functions and configuration information refer to Chapter 4, “Connectivity” on page 117.

7.3 Load Balancing Advisor and subplexing

The Load Balancing Advisor is a z/OS Communications Server component that allows any external load balancing solution to become sysplex aware. Subplex support for Load Balancing Advisor enhances the Load Balancing Advisor and the Load Balancing Agent function, so that they can participate in a sysplex subplexing. Prior to this support, only one Load Balancing Advisor was implemented in an LPAR. In a multiple TCP/IP stack environment, one Load Balancing Agent reported on all servers on all stacks, not just those stacks in a subplex.

With subplex support for Load Balancing Advisor, more than one Advisor can be active in the sysplex at any given time. In fact, there should be one Advisor active for each subplex in the sysplex that participates in load balancing through the Load Balancing Advisor. Each Advisor reads configuration data from a file, which can exist as a z/OS UNIX file, a PDS or PDSE member, or a sequential data set.

In the configuration file for each Advisor, the `sysplex_group_name` statement specifies the TCP/IP sysplex group name in the form of `EZBTvvtt`, where `vv` is the VTAM subplex group ID that is specified on the VTAM XCFGRPID start option and `tt` is the TCP/IP subplex group ID specified by the XCFGRPID parameter on the GLOBALCONFIG statement in the TCP/IP profile. If no VTAM subplex ID is specified when VTAM is started, then `vv` is CP. If no TCP/IP subplex ID is specified in the TCP/IP profile, then `tt` is CS. If you have a default subplex in your sysplex (that is, a subplex in which both the VTAM and TCP/IP subplex IDs are not specified), configure the Load Balancing Advisor for that subplex with a sysplex group name of `EZBTCPCS`.

Note: XCFGRPID is explained in 7.4.1, “XCF group names” on page 291.

In Figure 7-4, a Load Balancing Advisor application is configured to allow an external LBA to connect to the Internet subplex and the intranet production subplex.

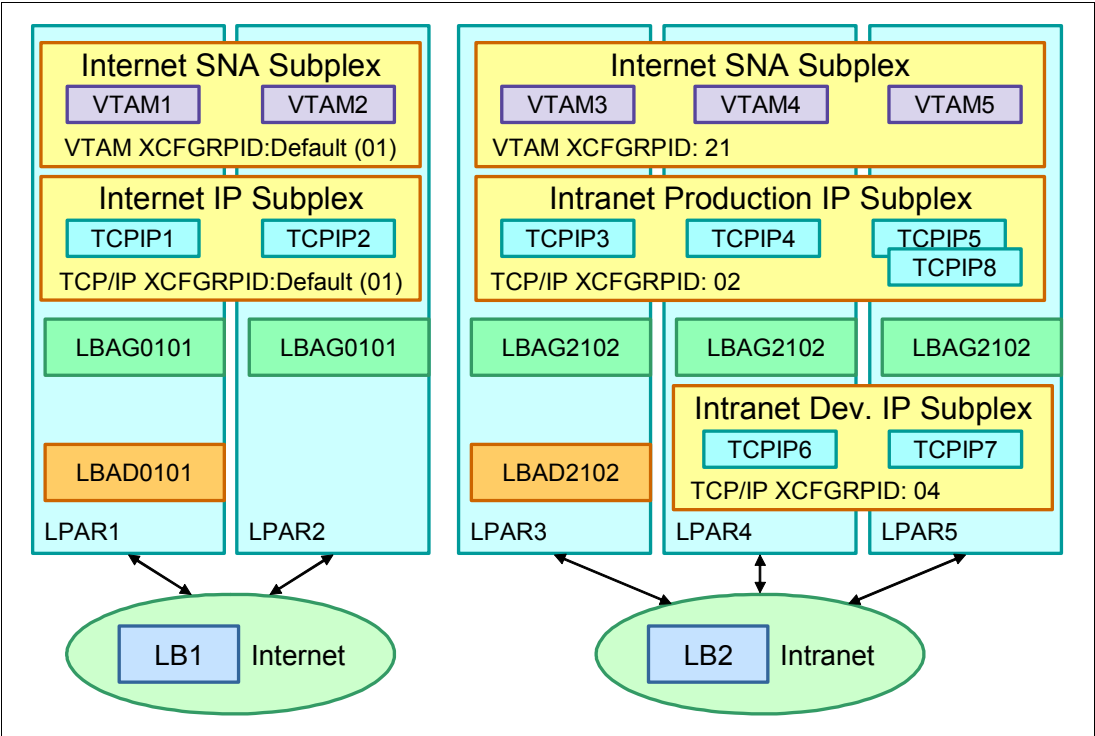


Figure 7-4 Load Balancing Advisor in subplexed sysplex

LB1 is balancing connections to applications running on TCP/IP stacks in the Internet IP subplex on LPAR1 and LPAR2. The TCP/IP sysplex group name is EZBTCPCS (VTAM XCFGRP ID 01 and TCP/IP XCFGRP ID 01). This is the default TCP/IP sysplex group name when the TCP/IP subplex ID is 0101 (the default VTAM and TCP/IP XCFGRP ID). LB1 connects to the Load Balancing Advisor in this subplex. The Advisor job, LBAD0101, is configured to use stacks that are members of TCP/IP subplex ID of 0101. A single instance of this Advisor can run in LPAR1 or LPAR2. It is currently running in LPAR1. Two Agents are configured to use the stacks that are members of TCP/IP subplex ID of 0101. The Agent job names are LBAG0101 on LPAR1 and LBAG0101 on LPAR2.

LB2 is balancing connections to applications running TCP/IP stacks in the intranet production IP subplex on LPAR3, LPAR4, and LPAR5. The TCP/IP sysplex group name is EZBT2102 (VTAM XCFGRP ID 21 & TCP/IP XCFGRP ID 02). The TCP/IP subplex ID is 2102. LB2 connects to a Load Balancing Advisor in this subplex. The Advisor, LBAD2102, is configured to use stacks that are members of the TCP/IP subplex ID of 2102. A single instance of this Advisor can run in LPAR3, LPAR4, or LPAR5. It is currently running in LPAR3. Three agents are configured to use stacks that are members of TCP/IP subplex ID of 2102. The three agent job names are as follows:

- ▶ LBAG02102 on LPAR3
- ▶ LBAG2102 on LPAR4
- ▶ LBAG2102 on LPAR5

Note: Although there are two TCP/IP stacks in LPAR5 in subplex 2102, there is only one Load Balancing Agent for that subplex on that LPAR. The one agent reports on all servers in that LPAR in that subplex.

There is no Load Balancing for applications that are running in the intranet development IP subplex. Therefore, no Advisor and no Agents need to run in this subplex. If you want to load balance in the intranet development IP subplex, configure an Advisor instance to run on either LPAR4 or LPAR5. Also, configure an Agent instance to run on both LPAR4 and LPAR5, and configure the Advisor and Agent applications to use stacks that are members of TCP/IP subplex ID 2104 (TCPIP6 and TCPIP7).

There are two subplexes in the three LPARs on the right side of the figure. The production IP subplex has TCP/IP subplex ID 2102 because the VTAM XCF group ID is 21 and the TCP/IP XCF group ID is 02. Subplex 2102 spans LPAR3, LPAR4, and LPAR5. The TCP/IP sysplex group name is EZBT2102. This subplex includes the following stacks:

- ▶ Stack TCPIP3 on LPAR3
- ▶ Stack TCPIP4 on LPAR4
- ▶ Stacks TCPIP5 and TCPIP8 on LPAR5

The Development IP subplex spans only LPAR4 and LPAR5. This subplex has a TCP/IP subplex ID of 2104 which is VTAM XCF group ID 21 and TCP/IP XCF group ID 04. The TCP/IP sysplex group name is EZBT2104. This subplex includes the following stacks:

- ▶ Stack TCPIP6 on LPAR4
- ▶ Stack TCPIP7 on LPAR5

Note: A TCP/IP subplex cannot span multiple VTAM subplexes, because all TCP/IP stacks on an LPAR use the same VTAM for their dynamic XCF communication.

7.4 Subplex implementation

TCP/IP stacks in the sysplex must be at a current release under the following conditions:

- ▶ Complete isolation between TCP/IP stacks in different subplexes is required.
- ▶ HiperSockets are used in support of dynamic XCF connectivity for TCP/IP stacks in a subplex.
- ▶ TCP/IP stacks in different subplexes accessing HiperSockets with the same CHPID.

An IP subplex is built through association of selected TCP/IP stack members to an XCF group. This is done by defining the XCFGRPID parameter in the GLOBALCONFIG statement of the TCP/IP profile. The subplex is created automatically at the start of the first stack member using this XCFGRPID definition plus the dynamic XCF IP address taken from the IPCONFIG statement DYNAMICXCF.

If the IP traffic for a defined subplex uses HiperSockets, which is the recommended method for cross-LPAR connectivity within the same server, then an additional parameter (IQDVLANID) in the GLOBALCONFIG is needed for the HiperSockets VLAN ID of the HiperSockets connection built with the DYNAMICXCF definition. Values from 2 to 31 are valid for XCFGRPID, while IQDVLANID allows values from 1 to 4094. If defining HiperSockets with DEVICE and LINK statements, the parameter VLANID on the LINK statement is required for assigning the VLAN for the subplex.

Requirement: A z890 or z990 at GA2 hardware level, or a z9 EC or z9 BC, is required to support VLAN IDs on HiperSockets.

Figure 7-5 depicts our subplexing environment. It shows three LPARs with a VTAM subplex, and two IP subplexes 11 and 22. Because we did not define the VTAM subplex, the XCFGRPID value for the VTAM subplex automatically defaults to CP.

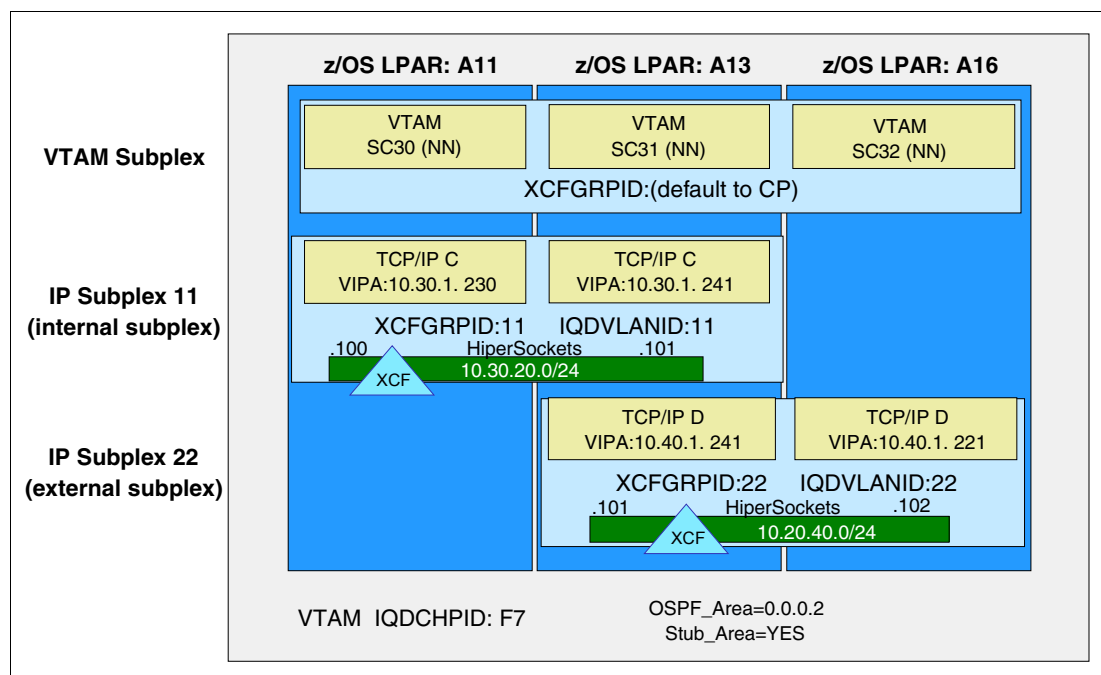


Figure 7-5 Subplex configuration

7.4.1 XCF group names

Basically, XCF group names for subplexes are created through the XCFGRPID parameter for the VTAM and TCP/IP environment; for example:

- ▶ For defining a VTAM subplex, use the XCFGRPID parameter in the VTAM start option. For detailed information about group and structure names, refer to *SNA Network Implementation Guide*, SC31-8777.
- ▶ For defining a TCP/IP subplex, use the XCFGRPID parameter on the GLOBALCONFIG statement in the TCP/IP profile.

For TCP/IP, both the VTAM group ID suffix and the TCP group ID suffix will be used to build the TCP/IP group name. This group name is also used to join the sysplex. Remember, when starting TCP/IP under Sysplex Autonomics control in previous z/OS releases, the stack joined the sysplex group with the name EZBTCPCS. You can verify this using the D XCF,GROUP command.

EZBTCPCS is the default TCP/IP group name. The format of this group name is *EZBTvvtt*, where *vv* is a 2-digit VTAM group ID suffix, specified on the VTAM XCFGRPID start option (the default is CP if not specified) and *tt* is a 2-digit TCP group ID suffix, specified on the XCFGRPID parameter of the GLOBALCONFIG statement (the default is CS if not specified).

In our scenario (see Example 7-3 on page 294 [5](#)), we defined XCFGRPID 11 for TCP/IP, and we did not define an XCFGRPID for VTAM. The result was an XCF group name of EZBTCP11 (Example 7-4 on page 294 [6](#)).

You might recognize that both XCFGRPIDs are important in creating the subplex group name. Be aware that changing the VTAM XCFGRPID will change the XCF group name for the TCP/IP stack. Thus, the stack is no longer a member of the previous TCP/IP subplex group.

For example, in our environment no VTAM XCFGRPID was defined and XCFGRPID 11 was specified for TCP/IP. Therefore, the XCF group name was dynamically built as EZBTCP11. If we add XCFGRPID=02 to the VTAM start option, then the new XCF group name will be EZBT0211.

Although nothing was changed in the TCP/IP profile definitions in this example, the TCP/IP stack with the new subplex group name is no longer a member of the previous subplex (EZBTCP11). Thus, the TCP/IP stack will lose the connectivity to the subplex.

Important: If VTAM is brought down and restarted with a different XCFGRPID, the TCP/IP stacks must be stopped and restarted to pick up the new VTAM subplex group ID. Otherwise, the TCP/IP stacks will continue to act as though there were in the original sysplex group, resulting in unpredictable connectivity.

7.4.2 TCP/IP structures

This section is intended for TCP/IP implementations using functions for Sysplex-wide Security Associations (SWSA) or for SYSPLEXPORTS, which is needed for sysplex-wide source VIPA to use one source VIPA for all outbound TCP connections within the sysplex.

- ▶ SWSA list structure (EZBDVIPA)

This stores information about IPsec tunnels addressed to distributed DVIPAs within the sysplex or subplex. This information is used to renegotiate IPsec tunnels in case of distributed DVIPA takeover. SWSA is enabled through definitions in the IPSEC statement.

- ▶ SYSPLEXPORTS list structure (EZBEPORT)

This contains all the ephemeral ports allocated in support of the sysplex-wide source VIPA function. Ephemeral ports that establish connections with external servers and use the sysplex-wide source VIPA function are allocated as participating clients from TCP/IP stacks within the sysplex or subplex.

This function is defined using TCPSTACKSOURCEVIPA on the IPCONFIG statement and SYSPLEXPORTS on the VIPADISTRIBUTE statement.

If TCP and VTAM Coupling Facility structures are used, names must also be unique for each subplex in order to preserve separation between the subplexes. This means that the TCP structures EZBDVIPA and EZBEPORT must be appended with the VTAM and TCP XCF group ID suffixes to the end of the structure names (for example, EZBDVIPAvvtt and EZBEPORTvvtt, where vv is the 2-digit VTAM group ID suffix specified on the XCFGRPID start option and tt is the 2-digit TCP group ID specified in the TCP/IP profile).

The default suffixes are as follows:

- ▶ If no VTAM XCFGRPID is specified, then the structure names will be EZBDVIPA01tt and EZBEPORT01tt.
- ▶ If no TCP/IP XCFGRPID is specified, then a null value is used for tt when the structure names are built.
- ▶ If no VTAM XCFGRPID and no TCP/IP XCFGRPID are specified, then vv and tt are both null.

The TCP structure names, including the suffixes, must be defined in the sysplex CFRM policy (see Example 7-1).

Example 7-1 TCP/IP structure example for SYSPLEXPORTS: subplex 11

```
STRUCTURE NAME(EZBEPORT0111)  
    INITSIZE(4096)  
    SIZE(8192)  
    PREFLIST(FACIL02,FACIL01)
```

Note: Example 7-1 is only a sample. The size depends on the number of source DVIPAs and concurrently established TCP outbound connections from all TCPSTACKSOURCEVIPA of the participating stacks within the sysplex. The ephemeral port number for each connection is stored to avoid duplicate source port numbers.

For more information regarding TCP and VTAM structures, refer to *z/OS MVS Setting Up a Sysplex*, SA22-7625.

The following sections describe the implementation for each subplex in detail.

7.4.3 Subplex 11: Internal subplex

Figure 7-6 depicts the configuration for IP Subplex 11 (the internal subplex).

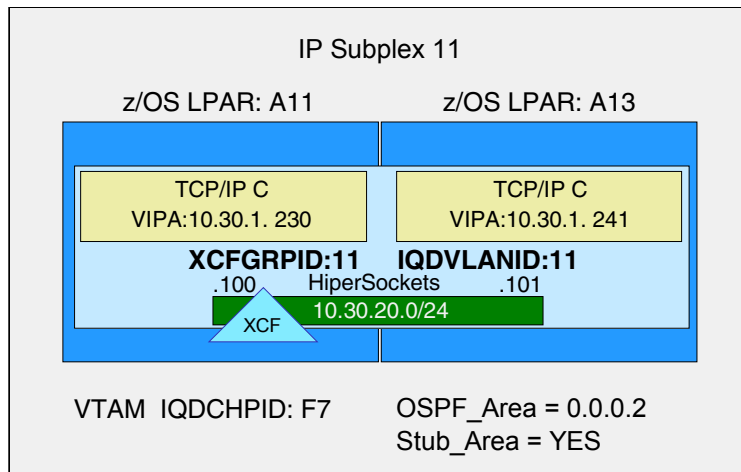


Figure 7-6 Subplex 11: internal subplex

TCP/IP profile definitions for subplex 11 in LPAR A13 stack C

Because we used automatically defined HiperSockets for the IP traffic within the subplex, we only had to add the VTAM start option IQDCHPID = F7 ¹. This CHPID is used when HiperSockets are implemented under z/OS.

The VTAM start option is needed by VTAM to automatically create the Transmission Resource List Element (TRLE) for the HiperSockets interface of the stack. The TRLE points to its IUTQDIO name, which is defined to the TCP/IP profile DEVICE name. The PORTNAME created by VTAM is IUTQDxx, where xx is the used Channel Path ID (CHPID).

The DYNAMICXCF function also requires the VTAM start option XCFINIT=YES ², which creates the XCF major node dynamically.

Tip: You can check your VTAM start options by using the D NET,VTAMOPTS command.

Example 7-2 ATCSTRxx definitions needed for DYNAMICXCF and HiperSockets interface

```
SYS1.VTAMLST(ATCSTR31)
IQDCHPID=F7,
XCFINIT=YES
```

¹
²

Example 7-3 shows the TCP/IP profile definitions needed for assigning stack C in LPAR A13 to subplex 11. Based on the parameters XCFGRPID **3** and IQDVLANID **4**, stack C belongs to subplex 11. The group interface is defined using the IPCONFIG parameter DYNAMICXCF with its IP address 10.30.20.101 **5**.

Example 7-3 TCP/IP profile: subplex definitions for stack C in LPAR A13

```
GLOBALCONFIG
  XCFGRPID 11 3
  IQDVLANID 11 4
;
IPCONFIG
DYNAMICXCF 10.30.20.101 255.255.255.0 8 5
```

Note: We used the same value for XCFGRPID and IQDVLANID. These values do not have to match. XCFGRPID allows values from 2 to 31, while IQDVLANID allows values from 1 to 4094.

The definitions for LPAR A11 are not shown because the XCFGRPID is the same. Only the DYNAMICXCF IP address **5** is different (10.30.20.100).

Verification of the subplex 11

The group name used is in the form EZBTvvtt, where vv is the 2-digit VTAM group ID suffix specified on the XCFGRPID start option or default (CP) and tt is the TCP group.

In our scenarios we did not define the VTAM start option XCFGRPID. A display from LPAR A13 TCP/IP stack C (see Example 7-4) shows that the stack is a member of the VTAM subplex group ID CP and TCP/IP subplex group 11, with the name EZBTCP11 **6**.

In the same LPAR there is another stack member of subplex group 22 with the name EZBTCP22 **7** (see definitions in 7.4.4, “Subplex 22: External subplex” on page 296).

Example 7-4 Displays of XCF groups

```
D XCF, GROUP
IXC331I 12.13.08 DISPLAY XCF 229
  GROUPS(SIZE):  ATRRRS(3)      COFVLFN0(3)      DBCDU(3)
                  EZBTCPCS(5)    EZBTCP11(2) 6
                  EZBTCP22(2) 7    IDAVQUI0(3)      IGWXSGIS(6)
                  IOEZFS(3)      IRRXCF00(3)      ISTCFS01(3)
                  ISTXCF(3)      IXCL000A(3)      IXCL000B(3)
                  IXCL0006(3)     SYSBPX(3)      SYSCNZMG(3)
```

The number between the parenthesis is related to the number of stacks active in the XCF group.

Example 7-5 displays that the stack in LPAR A13 is located in subplex 11 with its name EZBTCP11 **9**. The definitions for the subplex 22 (EZBTCP22) are described in 7.4.4, “Subplex 22: External subplex” on page 296.

Example 7-5 Display of specific stacks that belong to an XCF group

```
D TCPIP,TCPIPC,SYSPLEX,GROUP
EZZ8270I SYSPLEX GROUP FOR TCPIPC    AT SC31      IS EZBTCP11  9

D TCPIP,TCPIPD,SYSPLEX,GROUP
EZZ8270I SYSPLEX GROUP FOR TCIPD     AT SC31      IS EZBTCP22
```

The NETSTAT CONFIG display shows the XCFGRPID **10** and the IQDVLANID **11**.

Example 7-6 NETSTAT CONFIG with XCFGRPID and IQDVLANID for stack C

```
D TCPIP,TCPIPC,NETSTAT,CONFIG
GLOBAL CONFIGURATION INFORMATION:
TCPIPSTATS: NO  ECSALIMIT: 0000000K  POOLLIMIT: 0000000K
MLSCHKTERM: NO  XCFGRPID: 11 10 IQDVLANID: 11 11
SEGOFFLOAD: NO  SYSPLEXWLMPOLL: 060  MAXRECS: 100
EXPLICITBINDPORTRANGE: 00000-00000  IQDMULTIWRITE: NO
WLMPPRIORITYQ: NO
SYSPLEX MONITOR:
  TIMERSECS: 0060  RECOVERY: NO  DELAYJOIN: NO  AUTOREJOIN: NO
  MONINTF:  NO  DYNROUTE: NO  JOIN:  YES
```

The command NETSTAT DEV also shows the HiperSockets connection with VLANID **12**, which is the same value as IQDVLANID, as shown in Example 7-7.

Example 7-7 NETSTAT Device showing the HiperSockets VLAN ID

```
D TCPIP,TCPIPC,NETSTAT,DEV
DEVNAME: IUTIQDIO          DEVTYPE: MPCIPA
DEVSTATUS: READY
LNKNAME: IQDIOLNKOA1E1465  LNKTYPE: IPAQIDIO  LNKSTATUS: READY
IPBROADCASTCAPABILITY: NO
CFGROUTER: NON              ACTROUTER: NON
ARPOFFLOAD: YES             ARPOFFLOADINFO: YES
ACTMTU: 8192
VLANID: 11 12
READSTORAGE: GLOBAL (2048K)
SECCLASS: 255
IQDMULTIWRITE: DISABLED
ROUTING PARAMETERS:
  MTU SIZE: 8192              METRIC: 00
  DESTADDR: 0.0.0.0          SUBNETMASK: 255.255.255.0
MULTICAST SPECIFIC:
  MULTICAST CAPABILITY: YES
  GROUP              REFCNT              SRCFLTMD
  -----
  224.0.0.1          0000000001          EXCLUDE
  SRCADDR: NONE
LINK STATISTICS:
  BYTESIN              = 57156
  INBOUND PACKETS      = 548
  INBOUND PACKETS IN ERROR = 0
```

INBOUND PACKETS DISCARDED	= 0
INBOUND PACKETS WITH NO PROTOCOL	= 0
BYTESOUT	= 18296
OUTBOUND PACKETS	= 168
OUTBOUND PACKETS IN ERROR	= 0
OUTBOUND PACKETS DISCARDED	= 0

7.4.4 Subplex 22: External subplex

Figure 7-7 depicts the configuration for IP Subplex 22 (the external subplex). Note that both subplexes are using the same HiperSockets (CHPID F7).

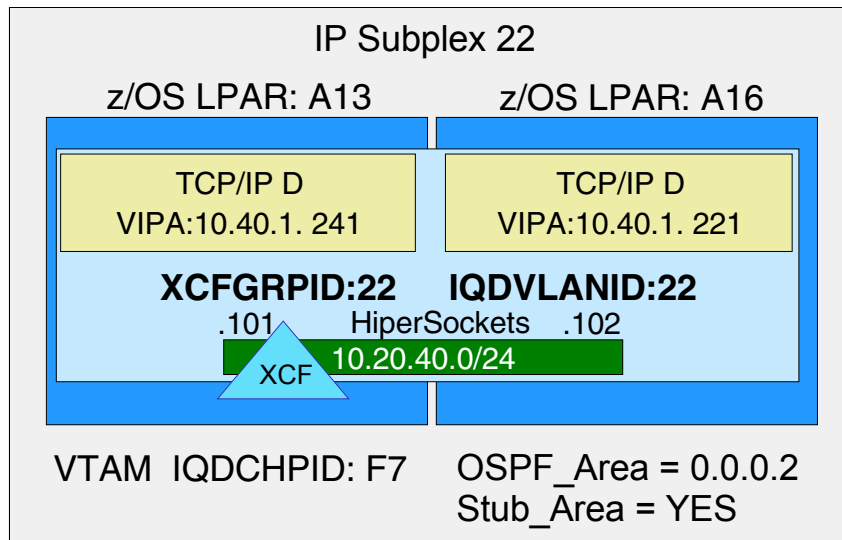


Figure 7-7 Subplex 22: external subplex

TCP/IP profile definitions for subplex 22 in LPAR A13 stack D

If you compare the definitions for stack D (shown in Example 7-8) with stack C (shown in Example 7-3 on page 294), you discover that only XCFGRPID **1**, IQDVLANID **2**, and DYNAMICXCF IP address **3** values are different.

Example 7-8 TCP/IP profile: subplex definitions for stack D in LPAR A13

```

GLOBALCONFIG
XCFGRPID 22 1
IQDVLANID 22 2
;
IPCONFIG
DYNAMICXCF 10.20.40.101 255.255.255.0 8 3

```

TCP/IP profile definitions for subplex 22 in LPAR A16 stack D

If you compare the definitions for stack D in LPAR A16 with stack D in LPAR A13 (see Example 7-8 on page 296), you will discover that XCFGRPID **4** and IQDVLANID **5** have the same values. Only the DYNAMICXCF IP address value **6** is different.

Example 7-9 TCP/IP profile: subplex definitions for stack D in LPAR A16

```
GLOBALCONFIG
XCFGRPID 22 4
IQDVLANID 22 5
;
IPCONFIG
DYNAMICXCF 10.20.40.102 255.255.255.0 8 6
```

Verification of the subplex 22

The command NETSTAT CONFIG shows the definitions used by the stack.

Example 7-10 NETSTAT CONFIG from LPAR A13 stack D

```
D TCPIP,TCPIP,NETSTAT,CONFIG
GLOBAL CONFIGURATION INFORMATION:
TCPIPSTATS: NO ECSALIMIT: 0000000K POOLLIMIT: 0000000K
MLSCHKTERM: NO XCFGRPID: 22 IQDVLANID: 22
SEGOFFLOAD: NO SYSPLEXWLPOLL: 060 MAXRECS: 100
EXPLICITBINDPORTRANGE: 00000-00000 IQDMULTIWRITE: NO
WLMPPRIORITYQ: NO
SYSPLEX MONITOR:
  TIMERSECS: 0060 RECOVERY: NO DELAYJOIN: NO AUTOREJOIN: NO
  MONINTF: NO DYNROUTE: NO JOIN: YES
```

7.4.5 Access verifications

We executed **ping** commands from all TCP/IP stacks in all LPARs. Example 7-11 shows a ping to IP address 10.30.20.101 (XCF and HiperSockets interface) from outside Subplex 11, which failed. All ping requests within each subplex were successful. Requests from other subplex groups or non-subplex groups were rejected.

Example 7-11 ping test

```
==> ping 10.30.20.101
CS V1R12: Pinging host 10.30.20.101
Ping #1 timed out
```

7.4.6 LBA connected to a subplex

Ensure that the Advisor and Agent are configured for a subplexed environment.

- ▶ There should be one Agent on each LPAR in the subplex.
- ▶ The Agents report to the Advisor only about applications within their subplex.
- ▶ If the Agent is configured with `sysplex_group_name` EZBTvvtt, the Agent will report only applications that are on the VTAM subplex vv and TCP/IP stacks with subplex tt.
- ▶ When configured for subplexing, the Agents will not report on other applications in the same LPAR.

- ▶ There can be more than one Agent in an LPAR if they are in different subplexes.
- ▶ IP addresses used as the source IP address for outbound Agent connection to the Advisor should be configured/owned by the proper stacks.
- ▶ The DVIPA for the Advisor needs to be defined in all the stacks associated with the subplex (and where a restart of the Advisor can occur).

7.5 References

For more information about subplexing, refer to the following publications:

- ▶ *z/OS Communications Server: IP Configuration Reference*, SC31-8776
- ▶ *z/OS Communications Server: IP Configuration Guide*, SC31-8775
- ▶ *HiperSockets Implementation Guide*, SG24-6816

Diagnosis

A key topic in any TCP/IP network infrastructure is documenting and analyzing problems. In this chapter we describe tools available in z/OS Communications Server as well as techniques to gather and diagnose problems related to the TCP/IP environment.

This chapter discusses the following topics.

Section	Topic
8.1, "Debugging a problem in a z/OS TCP/IP environment" on page 300	Problem determination techniques and the tools available to debug a problem in z/OS Communications Server - TCP/IP component
8.2, "Logs to diagnose CS for z/OS IP problems" on page 302	Why logs are important in problem analysis
8.3, "Useful commands to diagnose CS for z/OS IP problems" on page 303	Commands used to debug network problems
8.4, "Gathering traces in CS for z/OS IP" on page 315	Using z/OS Component Trace Service to capture trace data for the main z/OS Communications Server - TCP/IP component
8.5, "OSA-Express3 Network Traffic Analyzer" on page 333	Using OSAENTA to diagnose OSA problems
8.6, "Additional tools for diagnosing CS for z/OS IP problems" on page 355	Other tools that can be used to diagnose network problems
8.7, "MVS console support for selected TCP/IP commands" on page 360	Using EZACMD to run z/OS CS UNIX commands from the MVS console, in NetView and in TSO
8.8, "Additional information" on page 369	More information regarding use of logs, standard commands, tools, and utilities

8.1 Debugging a problem in a z/OS TCP/IP environment

In a TCP/IP network, several different types of problems can arise. Therefore, the support staff needs to develop debugging techniques that can help them better understand, define, and debug such problems. In this section we discuss a problem determination approach that uses logs, standard commands, tools, and utilities.

When problems arise in a TCP/IP environment, they can sometimes be very challenging to isolate. Without the proper tools, techniques, and knowledge of the environment, it can be very difficult to debug any problem. The culprit could be any one of the many components between the affected endpoints.

Therefore, we suggest categorizing the problem. Problems in TCP/IP networks can usually be classified into three major categories:

- ▶ Network connectivity problems

These occur when a z/OS server cannot establish a connection with another server or client because the node is unreachable (for example, it does not respond to the **ping** command).

- ▶ Application-related problems

These occur when a host is reachable, but communication with the desired application fails.

- ▶ Stack-related problems

These occur when the z/OS TCP/IP stack does not work as implemented, or ends with a dump.

Most problems can easily be associated to one of these categories and the information needed to debug them can be retrieved from logs, commands, or utilities.

Logs are the first and most important tool to help you understand the nature of the problem. In logs you will find messages that might explain what happened or even lead you to the actions needed to solve the problem.

Sometimes, however, problems like connectivity or routing do not provide messages that clearly show what went wrong. Therefore, you need further information, which can be obtained by using commands such as **netstat**, **ping**, or **tracert**. If the commands do not provide enough information to solve or isolate the problem, then you can invoke the z/OS Communications Server trace utilities that gather data as it passes through the devices and the stack.

Many problems related to the TCP/IP stack are due to configuration errors. Here you can use logs to find useful messages that indicate where the error is located.

If the TCP/IP stack happens to abend, a dump is generated. In such a situation, the dump and related information can be sent to IBM Support for further analysis.

8.1.1 An approach to problem analysis

When performing problem analysis, it is essential that you have readily available current, accurate documentation describing the physical and logical network environment. This documentation should include network diagrams, naming conventions, addressing schemes, and system configuration information.

When a problem occurs, the first step is to verify that the operating environment is behaving as expected. After this is confirmed, you can then focus on other areas. To help isolate the problem, a useful approach is to answer such basic questions including:

- Is the TCP/IP stack running correctly?

This generic question can help determine whether the problem is stack-related. It can be answered by verifying the behavior of the entire Communications Server for z/OS IP environment.

Usually the tools used to answer this question are the logs where messages related to the problem can be found (see 8.2, “Logs to diagnose CS for z/OS IP problems” on page 302) and tools that receive information using the Network Management Interface (see 8.6, “Additional tools for diagnosing CS for z/OS IP problems” on page 355).

If the problem is an abend, save the generated dump for analysis. The configuration should also be checked for inconsistencies. If you conclude it is not a stack-related problem, then the next step will be based on your findings to determine whether it is a network- or application-related problem.

- Has this ever worked before? If so, what has changed?

These two basic questions might seem obvious, but they are in fact the most common reasons for problems encountered in a Communications Server for z/OS IP environment.

If the problem is with a production and stable environment, you must first check whether any changes have been made. In some cases, changes do not take effect until a system or stack recycle is done. The only useful approach in this case is to keep track of any changes and always use change management processes.

If you are dealing with a new implementation, was a step-by-step approach being used? If so, you will probably know in which step the problem occurred and can adapt your problem determination procedure based on the step being implemented.

- Are the physical connections and interfaces active and working properly?

This question is related to a connectivity problem, and it leads to checking interface definitions and status. You also need to look at the log files, and use commands to determine whether the interfaces are operational. The **netstat** command can be used to verify this, as discussed in 8.3, “Useful commands to diagnose CS for z/OS IP problems” on page 303.

If it is an intermittent problem, or if you cannot find the cause of the problem, Communications Server for z/OS IP provides a set of trace tools that you can use to gather more information. See 8.4, “Gathering traces in CS for z/OS IP” on page 315.

- Can the destination host be reached?

In cases where the physical connections are up and running, but a specific host cannot be reached, the problem is probably related to routing. In this case, you need to look at the logs files for related error messages. You can also use commands such as **ping**, **tracert**, and **netstat** to discover why you are not able to reach this host.

If these steps do not provide you with enough information to isolate the problem, you will need to use the packet trace utility as described in 8.4, “Gathering traces in CS for z/OS IP” on page 315. A packet trace allows you to check if there is any data going to or coming from the host you are trying to reach.

- Is the problem affecting multiple connections?

There might be a problem with the proper configuration of a firewall policy, an incorrect interface configuration, or an application problem.

The approach in this case is to review the configuration files, looking for inconsistencies. Also examine the log files, which might contain error messages about this problem. If necessary, you can also debug this problem using the component trace for event and packet tracing; see 8.4, “Gathering traces in CS for z/OS IP” on page 315.

- Is this problem related to a single application?

To analyze application problems, you need general knowledge of the application protocol. You should know what transport protocol is used, which port numbers are used, how a connection is established, and the application protocol semantics.

Mainly, the following tools are used to diagnose application problems:

- Debugging commands
- Specific application traces
- Packet trace
- Component trace

You can check whether an application is running by using display commands. With TELNET, for example, the D TCPIP,,TELNET and V TCPIP,,TELNET commands can be used to verify and control Telnet connections.

Specific application traces are useful for following the execution of an application (either client or server), and checking whether there are error messages. The application trace might not be sufficient to diagnose some problems because it shows the commands (rather than the data) exchanged during a connection.

For an in-depth investigation you need to use a packet trace, which can be interpreted relatively easily for standard applications (see 8.4, “Gathering traces in CS for z/OS IP” on page 315), or a CTRACE (when requested by IBM).

8.2 Logs to diagnose CS for z/OS IP problems

To start the problem determination process, the most important step is to pull together reliable information to verify, classify, and define possible lines of action to resolve a problem. Examining logs is an excellent starting point in the problem determination process. Logs contain different types of messages (informational, error-based, and warning-based) that provide very useful information. Logs are very important in the problem determination process because they can be the only source of information about a problem.

For example, in a production environment, problems are often business- or service-related, and end users are the first to notice there is a problem (usually, because they are unable to access applications or execute services). The operational response taken when a problem is discovered is often based on business or service recovery; it is usually only after these actions are taken that support personnel are called upon to evaluate and determine why the outage occurred in the first place. In some cases, there is no information given other than a problem description based on the business or service point of view, with no technical perspective.

In many situations, the information obtained during the problem determination process comes from separate logs (system, application, and stack logs). To be able to build a clear picture of the problem or outage, all significant information must be correlated.

Because of this, we recommend that you implement syslogd to control where all messages are sent. This way, you will have a single place to refer to when debugging a problem. The syslogd process is a UNIX process that logs UNIX application messages to one or more files.

TCP/IP services that run as UNIX processes log application messages using syslogd can consolidate logging information from several systems to one system through UDP communications.

For further information about setting up syslogd, refer to *Communications Server for z/OS V1R9 TCP/IP Implementation Volume 2: Standard Applications*, SG24-7533.

8.3 Useful commands to diagnose CS for z/OS IP problems

To solve problems it is important to know what tools are available and how to make best use of them. Some commands or utilities can be used to review configuration options or settings; others can be used to test connectivity.

In this section we describe briefly the main commands that you can use to diagnose problems in a Communications Server for z/OS IP environment. For additional help and detailed information about the commands described and other commands that can be used for problem determination, refer to *z/OS Communications Server: IP System Administrator's Commands*, SC31-8781, and *z/OS Communications Server: IP Diagnosis Guide*, GC31-8782.

In this section, we describe these commands:

- ▶ **ping** command (TSO or z/OS UNIX)
- ▶ **tracert** command
- ▶ **netstat** command (console, TSO, or z/OS UNIX)

8.3.1 The ping command (TSO or z/OS UNIX)

The **ping** command is relatively simple, but it is one of the best tools you can use to check basic connectivity. It sends an ICMP echo request message to the target system and waits for an ICMP echo reply message. Because this command uses only two ICMP messages (echo request and echo response), it cannot be used to test transport or application protocols. In order for PING to work, the sending system and all intermediate systems must be correctly set up for both the outbound and inbound journeys.

Typically, **ping** is used to verify:

- ▶ The route to a network is defined correctly.
- ▶ The router is able to forward packets to the network.
- ▶ The remote host is able to send and receive packets in the network
- ▶ The remote host has a route back to the local host.

Tip: Using names instead of IP address needs the resolver or DNS to do the translation, thus adding more variables to the problem determination task. This should be avoided when diagnosing network problems. Use the host IP address instead.

In most cases, the default options of **ping** are used. However, in a z/OS Communications Server environment, using the default options might lead to a false conclusion, given the number of interfaces that can be used to transport the ICMP request.

This command provides several options that can be used to analyze a problem in more detail. For example, the **intf** option of the TSO **ping** command (or the **-i** option, if you use the z/OS UNIX **ping** command) specifies the local interface over which the packets are sent. This can be useful to determine if a remote host is reachable through the desired path.

Table 8-1 shows the available options that can be used with the **ping** command in TSO and z/OS UNIX environments.

Table 8-1 Options available with ping

Options	TSO	z/OS UNIX
Address type (ipv4 /ipv6)	Addrtype	-A
Number of ping interactions	Count	-c
Interface to be used as path	Intf	-i
Amount of data being sent	Length	-l
Do not resolve IP addresses to host names (used with Pmtu)	NOName	-n
Determine the path MTU size of a host in the network.	PMTU	-P
Source IP address	Srcip	-s
TCP/IP stack to be used.	TCP	-p
How long it waits for a response	Timeout	-t
Display details of the echo reply packets	Verbose	-v
Help information	HELP or ?	-h or -?

Figure 8-1 illustrates the use of the **ping** command for problem determination.

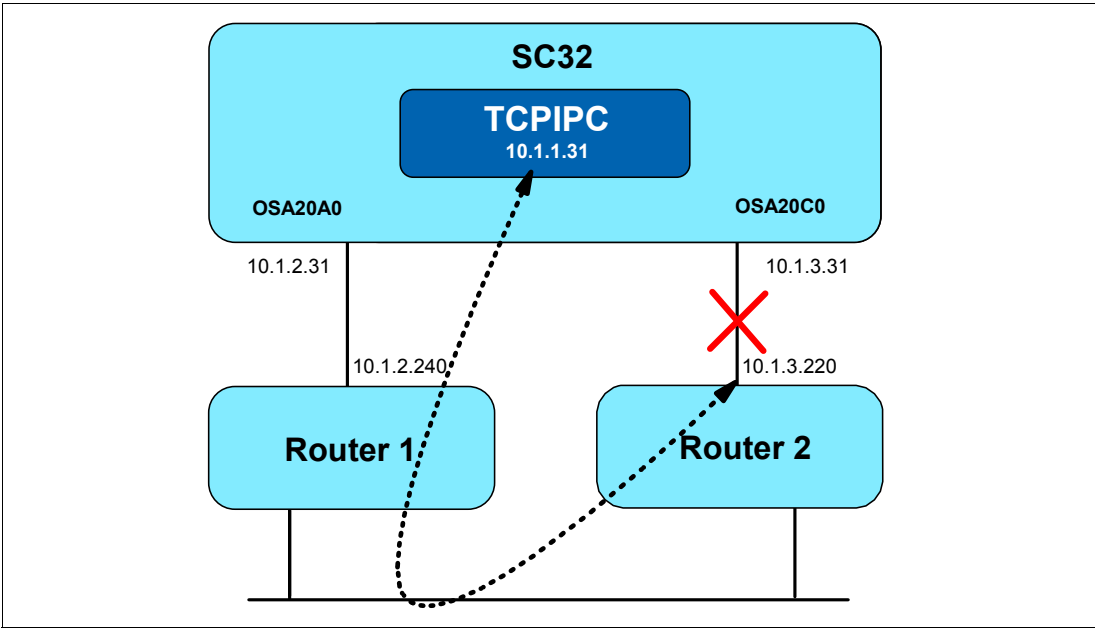


Figure 8-1 Using the ping command with the interface option

This is a situation where a ping might work even when the expected route to the endpoint is down. In this case the endpoint was accessed using an alternate route. However, using the **ping** command *without* the correct option would hide the problem, as shown in Example 8-1.

Example 8-1 ping command without the intf option

```
ping 10.1.3.220 (tcp tcpipc
CS V1R12: Pinging host 10.1.3.220
Ping #1 response took 0.001 seconds.
***
```

To avoid such confusion, indicate which path to verify by using the interface (intf) option, as shown in Example 8-2.

Example 8-2 ping command with the intf option

```
ping 10.1.3.220 (tcp tcpipc intf osa20c0l
CS V1R12: Pinging host 10.1.3.220
sendMessage(): EDC8130I Host cannot be reached. (errno2=0x74420291)
***
```

After using the correct command, you can see there is a problem using interface OSA20C0L, which is the direct connection to the 10.1.3.0 subnetwork.

Using the pmtu option

Use the pmtu option on the **ping** command to determine where fragmentation is necessary in the network. The pmtu yes option differs from the pmtu ignore option in the way **ping** completes its echo function. Refer to *z/OS Communications Server: IP System Administrator's Commands*, SC31-8781, for a detailed discussion about the pmtu option.

Example 8-3 shows a ping with a very large packet size, with no pmtu option specified. We used the noname option to avoid a reverse DNS lookup on the IP address.

Example 8-3 Using ping without the pmtu option

```
ping 10.1.1.10 (noname tcp tcpipb l 25000 c 1 t 1
CS V1R12: Pinging host 10.1.1.10
Ping #1 response took 0.000 seconds.
***
```

Example 8-4 shows that by adding the pmtu yes option **1** to the **ping** command, you can determine at which hop **2** the fragmentation is necessary, and the MTU size **3**.

Example 8-4 Using ping with the pmtu option

```
ping 10.1.1.10 (noname tcp tcpipb l 25000 c 1 t 1 pmtu yes 1
CS V1R12: Pinging host 10.1.1.10
Ping #1 needs fragmentation at: 10.1.7.21 2
Next-hop MTU size is 8192 3
***
```

By varying the size of the ping packet, you can work your way through the path to the hop requiring fragmentation, as shown in Example 8-5.

Example 8-5 Varying the ping packet size

```
ping 10.1.1.10 (noname tcp tcpihb l 6000 c 1 t 1 pmtu yes
CS V1R12: Pinging host 10.1.1.10
Ping #1 response took 0.000 seconds.
***

ping 10.1.1.10 (noname tcp tcpihb l 8164 c 1 t 1 pmtu yes
CS V1R12: Pinging host 10.1.1.10
Ping #1 response took 0.000 seconds.
***

ping 10.1.1.10 (noname tcp tcpihb l 8165 c 1 t 1 pmtu yes
CS V1R12: Pinging host 10.1.1.10
Ping #1 needs fragmentation at: 10.1.5.21 (10.1.5.21)
Next-hop MTU size is 8192
***
```

Example 8-6 illustrates the use of the pmtu ignore option.

Example 8-6 ping with the pmtu ignore option

```
ping 10.1.1.10 (noname tcp tcpihb l 25000 c 1 t 1 pmtu ignore
CS V1R12: Pinging host 10.1.1.10
Ping #1 needs fragmentation at: 10.1.7.21
Next-hop MTU size is 8192
***
```

8.3.2 traceroute command

The **traceroute** (z/OS UNIX) or **tracerte** (TSO) command is used to determine the route that IP datagrams follow through the network. **traceroute** is based on ICMP and UDP. It sends an IP datagram with a time-to-live (TTL) of 1 to the destination host. The first router decrements the TTL to 0, discards the datagram, and returns an ICMP Time Exceeded message to the source. In this way, the first router in the path is identified. This process is repeated with successively larger TTL values to identify the exact series of routers in the path to the destination host.

On most platforms, the **traceroute** command sends UDP datagrams to the destination host. These datagrams reference a port number outside the standard range. The source knows when it has reached the destination host when it receives an ICMP “Port Unreachable” message.

traceroute displays the route that a packet takes to reach the requested target. The output generated by this command can be seen in Example 8-7.

Example 8-7 tracerte command results

```
CS V1R12: Traceroute to 10.1.100.222 (10.1.100.222)
 1 10.1.2.240 (10.1.2.240) 0 ms 0 ms 0 ms
 2 10.1.100.222 (10.1.100.222) 0 ms 0 ms 0 ms
***
```

8.3.3 The netstat command (console, TSO, or z/OS UNIX)

The **netstat** command provides information about the status of the local host, including information about TCP/IP configuration, connections, network clients, gateways, and devices. It also has options to drop connections for users who have the MVS.VARY.TCPIP.DROP statement defined in their RACF profile.

As shown in Figure 8-2, there is a variety of **netstat** options, and these can be further qualified by filter criteria, depending on the option you choose. The output can be displayed to the terminal (default), to a data set (report), or to the REXX data stack. The Output Format (short or long) supports IPv6 addresses.

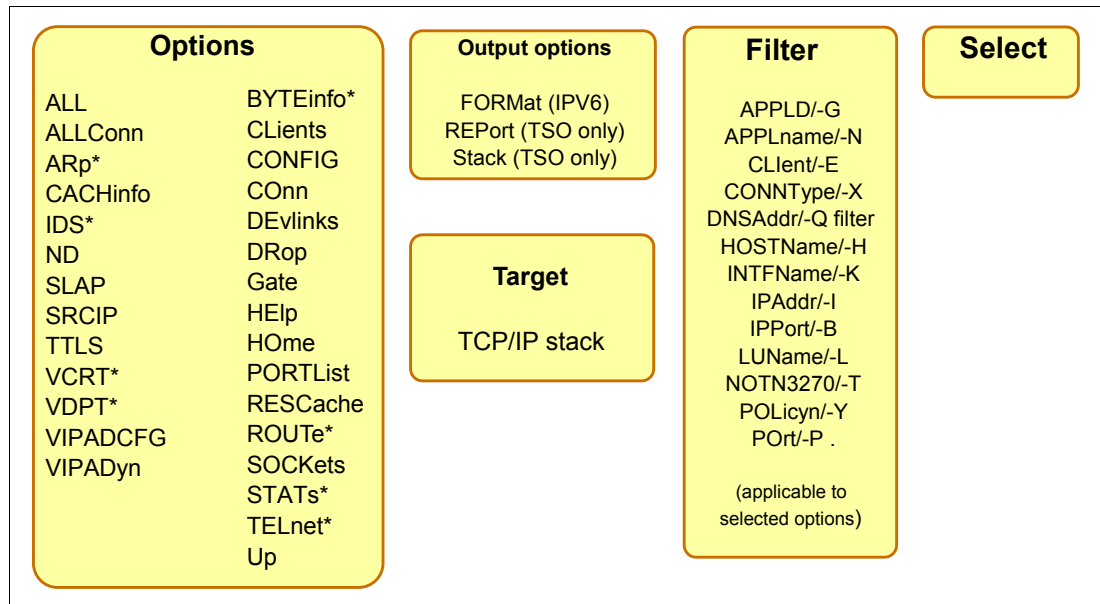


Figure 8-2 netstat command options: target output (filter select)

The remainder of this section shows examples of **netstat** commands used for diagnostic purposes, and their outputs.

Example 8-8 displays the results of the D TCPIP,TCPIPA,NETSTAT,CONN command.

Example 8-8 D TCPIP,TCPIPA,NETSTAT,CONN command

```
EZD0101I NETSTAT CS V1R12 TCPIPA 337
USER ID  CONN      STATE
FTPDA1   00000021  LISTEN
  LOCAL SOCKET:   ::FFFF:10.1.1.10..21
  FOREIGN SOCKET: ::FFFF:0.0.0.0..0
JES2S001 00000013  LISTEN
  LOCAL SOCKET:   ::..175
  FOREIGN SOCKET: ::..0
OMPA      0000007A  ESTBLSH
  LOCAL SOCKET:   127.0.0.1..1027
  FOREIGN SOCKET: 127.0.0.1..1028
TCPIPA    00000078  ESTBLSH
  LOCAL SOCKET:   127.0.0.1..1028
  FOREIGN SOCKET: 127.0.0.1..1027
TN3270A   00001278  LISTEN
  LOCAL SOCKET:   ::..4992
  FOREIGN SOCKET: ::..0
TN3270A   00001279  LISTEN
  LOCAL SOCKET:   ::..992
  FOREIGN SOCKET: ::..0
TCPIPA    0000142E  UDP
  LOCAL SOCKET:   ::..3271
  FOREIGN SOCKET: *.*
7 OF 7 RECORDS DISPLAYED
END OF THE REPORT
```

Use the D TCPIP,tcpipproc,NETSTAT, DEVLINK command to display the status and associated configuration values for a device and its defined interfaces. This command can be filtered to display only the interface you want, as shown in Example 8-9.

Example 8-9 D TCPIP,TCPIPA,N,DEV,INTFN=OSA2080I

```
D TCPIP,TCPIPA,N,DEV,INTFN=OSA2080I
INTFNAME: OSA2080I          INTFTYPE: IPAQENET  INTFSTATUS: READY
PORTNAME: OSA2080          DATAPATH: 2082      DATAPATHSTATUS: READY
CHPIDTYPE: OSD
SPEED: 0000001000
IPBROADCASTCAPABILITY: NO
VMACADDR: 02000A776873     VMACORIGIN: OSA     VMACROUTER: LOCAL
ARPOFFLOAD: YES            ARPOFFLOADINFO: YES
CFGMTU: 1492               ACTMTU: 1492
IPADDR: 10.1.2.11/24
VLANID: 10                 VLANPRIORITY: DISABLED
DYNVLANREGCFG: NO          DYNVLANREGCAP: YES
READSTORAGE: GLOBAL (4096K)
INBPERF: BALANCED
CHECKSUMOFFLOAD: YES       SEGMENTATIONOFFLOAD: YES
SECCLASS: 255              MONSYSPLEX: NO
ISOLATE: NO                OPTLATENCYMODE: NO
MULTICAST SPECIFIC:
MULTICAST CAPABILITY: YES
GROUP          REFCNT      SRCFLTMD
-----
224.0.0.1      0000000001  EXCLUDE
SRCADDR: NONE
INTERFACE STATISTICS:
BYTESIN                = 9548
INBOUND PACKETS        = 128
INBOUND PACKETS IN ERROR = 0
```

```

INBOUND PACKETS DISCARDED          = 0
INBOUND PACKETS WITH NO PROTOCOL  = 0
BYTESOUT                           = 4891
OUTBOUND PACKETS                    = 66
OUTBOUND PACKETS IN ERROR          = 0
OUTBOUND PACKETS DISCARDED         = 0
IPV4 LAN GROUP SUMMARY
LANGROUP: 00001
NAME          STATUS      ARPOWNER      VIPAOWNER
----          -
OSA2080I      ACTIVE      OSA2080I      NO
OSA20A0I      ACTIVE      OSA20A0I      YES
1 OF 1 RECORDS DISPLAYED Outbound Packets Discarded      = 0

```

The `D TCPIP,tcpipproc,NETSTAT,ROUTE` command displays the current routing tables for TCP/IP, and it can be filtered to show specific routes as shown in Example 8-10.

Example 8-10 D TCPIP,TCPIPA,N,ROUTE,IPADDR=10.1.100.0/24

```
D TCPIP,TCPIPA,N,ROUTE,IPADDR=10.1.100.0/24
IPV4 DESTINATIONS
DESTINATION      GATEWAY      FLAGS      REFCNT  INTERFACE
10.1.100.0/24    10.1.2.240   UGO        000000  OSA2080L
10.1.100.0/24    10.1.2.240   UGO        000000  OSA20A0L
10.1.100.0/24    10.1.3.240   UGO        000000  OSA20C0L
10.1.100.0/24    10.1.3.240   UGO        000000  OSA20E0L
4 OF 4 RECORDS DISPLAYED
END OF THE REPORT
```

You can optionally display additional application connection data by using the `APPLDATA` parameter on the `NETSTAT CONN` and `NETSTAT ALLCONN` commands. Example 8-11 contrasts the output of two `NETSTAT CONN` commands: one without the `APPLDATA` parameter **1**, the other with the `APPLDATA` parameter **2**. The TN3270 server populates the `APPLDATA` field with connection data, as documented in *z/OS Communications Server: IP Configuration Reference*, SC31-8776. The TN3270 `appldata` fields shown for the connection are the component ID, LU name, the SNA application name, connection mode, client type, security method, security level, and security cipher **3**.

Example 8-11 NETSTAT CONN without and with the APPLDATA option

```
D TCPIP,TCPIPB,N,conn 1
EZD0101I NETSTAT CS V1R12 TCPIPB
USER ID  CONN      STATE
JES2S001 00000013 LISTEN
  LOCAL SOCKET:  ::..175
  FOREIGN SOCKET: ::..0
SNMPQEB 00000023 LISTEN
  LOCAL SOCKET:  0.0.0.0..1025
  FOREIGN SOCKET: 0.0.0.0..0
SNMPQEB 00000021 UDP
  LOCAL SOCKET:  0.0.0.0..1026
  FOREIGN SOCKET: *.*
SNMPQEB 00000022 UDP
  LOCAL SOCKET:  0.0.0.0..162
  FOREIGN SOCKET: *.*
TCPIPB 00000026 UDP
  LOCAL SOCKET:  ::..1027
  FOREIGN SOCKET: *.*
5 OF 5 RECORDS DISPLAYED

D TCPIP,TCPIPB,N,conn,appldata 2
```

```

EZD0101I NETSTAT CS V1R12 TCPIPA 759
USER ID  CONN      STATE
FTPDA1   00000021 LISTEN
  LOCAL SOCKET:  ::FFFF:10.1.1.10..21
  FOREIGN SOCKET: ::FFFF:0.0.0.0..0
  APPLICATION DATA: EZAFTPD 3
JES2S001 0000145F LISTEN
  LOCAL SOCKET:  ::..175
  FOREIGN SOCKET: ::..0
OMPA      0000007A ESTBLSH
  LOCAL SOCKET:  127.0.0.1..1027
  FOREIGN SOCKET: 127.0.0.1..1028
SNMPDB    000014CF LISTEN
  LOCAL SOCKET:  ::..1035
  FOREIGN SOCKET: ::..0
TCPIPA    00000078 ESTBLSH
  LOCAL SOCKET:  127.0.0.1..1028
  FOREIGN SOCKET: 127.0.0.1..1027
TN3270A   000014CA LISTEN
  LOCAL SOCKET:  ::..992
  FOREIGN SOCKET: ::..0
  APPLICATION DATA: EZBTNSRV LISTENER 3
TN3270A   000014CB LISTEN
  LOCAL SOCKET:  ::..23
  FOREIGN SOCKET: ::..0
  APPLICATION DATA: EZBTNSRV LISTENER 3

```

You can optionally display the report provided by the netstat ALL/-A, that is now available when using the DISPLAY TCPIP,NETSTAT command, in addition to being available using the TSO or z/OS UNIX shell environment. You can filter this command to display only the client IPADDR that you want, and receive a complete details of this session, such as the maximum segment size in use, as shown in Example 8-12.

Example 8-12 D TCPIP,TCPIPB,N,ALL,IPADDR=10.1.100.222

```

D TCPIP,TCPIPB,N,ALL,IPADDR=10.1.100.222
EZD0101I NETSTAT CS V1R12 TCPIPB 382
CLIENT NAME: TN3270B                CLIENT ID: 00000DA3
  LOCAL SOCKET: ::FFFF:10.1.1.20..23
  FOREIGN SOCKET: ::FFFF:10.1.100.222..1401
    BYTESIN:                00000000000000000035
    BYTESOUT:               000000000000000000321
    SEGMENTSIN:            000000000000000000008
    SEGMENTSOUT:           000000000000000000012
    LAST TOUCHED:          18:45:49          STATE:          ESTABLSH
    RCVNXT:                2832645553      SNDNXT:
2175448364
  CLIENTRCVNXT:            2832645553      CLIENTSNDNXT:
2175448364
    INITRCVSEQNUM:         2832645517      INITSDNSEQNUM:
2175448042
    CONGESTIONWINDOW:      0000014520      SLOWSTARTTHRESHOLD:
0000065535
    INCOMINGWINDOWNUM:     2833169838      OUTGOINGWINDOWNUM:
2175513578
    SNDWL1:                2832645550      SNDWL2:
2175448364
    SNDWND:                0000065214      MAXSNDWND:
0000065535
    SNDUNA:                2175448364      RTT_SEQ:
2175448361
    MAXIMUMSEGMENTSIZE: 0000001452      DSFIELD:          00
    ROUND-TRIP INFORMATION:
SMOOTH TRIP TIME: 15.000          SMOOTHTRIPVARIANCE: 229.00
    REXMT:                0000000000      REXMTCOUNT:
0000000000
    DUPACKS:              0000000000      RCVWND:
0000524285
    SOCKOPT:              C000          TCPTIMER:          00
    TCPSIG:               01          TCPSEL:            00
    TCPDET:               E0          TCPPOL:            00
    TCPPRF:               00
    QOSPOLICY:            NO
    ROUTINGPOLICY:        NO
    RECEIVEBUFFERSIZE:    0000262144      SENDBUFFERSIZE:
0000262144
    RECEIVEDATAQUEUED:    0000000000
    SENDDATAQUEUED:       0000000000
    ANCILLARY INPUT QUEUE: N/A
    APPLICATION DATA:    EZBTNSRV SHLU02          ET B
-----
1 OF 1 RECORDS DISPLAYED

```

You can filter the output of the NETSTAT CONN,APPLDATA command by adding the APPLD filter option and specifying the filter criteria. The APPLDATA field is a total of 40 bytes. By using an asterisk (*) in the filter criteria, you can filter on any part of the 40 bytes. Example 8-13 shows several filter criteria strings being used.

Example 8-13 NETSTAT CONN APPLDATA with APPLD filter

```

D TCPIP,TCPIPB,N,CONN,APPLDATA,APPLD=*TNSRV*
EZD0101I NETSTAT CS V1R12 TCPIPB 480
USER ID  CONN      STATE
TN3270B  00000DA3  ESTBLSH
    LOCAL SOCKET:  ::FFFF:10.1.1.20..23
    FOREIGN SOCKET: ::FFFF:10.1.100.222..1401
    APPLICATION DATA: EZBTNSRV SHLU02                ET B
TN3270B  00000D7B  LISTEN
    LOCAL SOCKET:  ::..23
    FOREIGN SOCKET: ::..0
    APPLICATION DATA: EZBTNSRV LISTENER
TN3270B  00000DC0  ESTBLSH
    LOCAL SOCKET:  ::FFFF:10.1.1.20..23
    FOREIGN SOCKET: ::FFFF:10.1.100.221..4908
    APPLICATION DATA: EZBTNSRV SH99LU02              ET B
TN3270B  00000D7A  LISTEN
    LOCAL SOCKET:  ::..992
    FOREIGN SOCKET: ::..0
    APPLICATION DATA: EZBTNSRV LISTENER
4 OF 4 RECORDS DISPLAYED
END OF THE REPORT

D TCPIP,TCPIPB,N,CONN,APPLDATA,APPLD=*SC31*
USER ID  CONN      STATE
TN3270B  00000111  ESTBLSH
    LOCAL SOCKET:  ::FFFF:10.1.1.20..23
    FOREIGN SOCKET: ::FFFF:10.1.100.222..1028
    APPLICATION DATA: EZBTNSRV SC31BB05 SC31TS03 3T B
1 OF 1 RECORDS DISPLAYED
END OF THE REPORT

```

Note: The MAXRECS parameter is now available on the GLOBALCONFIG TCP/IP profile statement for configuring a default value for the DISPLAY TCPIP,NETSTAT command's MAX parameter. The default value is 100.

You can drop (reset) each connection using the Netstat DROP/-D command with a connection ID obtained by issuing the Netstat CONN/-c command. If you want to move workload from one server application to another you can quiesce the creation of new connections to the old server, but all persistent connections need to be ended using the Netstat DROP/-D command. And also, you might want to drop dozens of connections with an unexpected state, such as CLOSWT, for the purpose of solving any problems using the Netstat DROP/-D command.

The VARY TCPIP, DROP command allows all TCP connections associated with a server matching the specified filter to be reset. If more than one server application is found to match the input filter values, the command is failed. Existing TCP connections are reset by this command, but new connection requests are not quiesced. If necessary, you might quiesce new connection requests to the server application before issuing this command.

Example 8-14 shows the output for two VARY TCPIP,,DROP commands: one with the PORT parameter and the optional JOBNAME parameter **1**, the other with the JOBNAME parameter **2**. You can optionally specify the address space ID (ASID). You can see EZD2013I message **3**, which includes the number of connections that were reset. The following messages depend on the server application that the command was issued on.

Example 8-14 V TCPIP,,DROP with PORT filter and JOBNAME filter

```

V TCPIP,TCPIPB,DROP,PORT=23,JOBNAME=TN3270B 1
EZZ0060I PROCESSING COMMAND: VARY TCPIP,TCPIPB,DROP,PORT=23,JOBNAME=TN
3270B
EZD2013I    3 CONNECTIONS WERE SUCCESSFULLY DROPPED 3
IKT100I USERID          CANCELED DUE TO UNCONDITIONAL LOGOFF
IKT122I IPADDR..PORT 10.1.100.221..4214
IKT100I USERID          CANCELED DUE TO UNCONDITIONAL LOGOFF
IKT100I USERID          CANCELED DUE TO UNCONDITIONAL LOGOFF
IKT122I IPADDR..PORT 10.1.100.221..4213
IKT122I IPADDR..PORT 10.1.100.221..4212
EZZ6034I TN3270B CONN 000000AE LU SC31BB26 CONN DROP  ERR 1010 141
      IP..PORT: ::FFFF:10.1.100.221..4212          EZBTTRCV
EZZ6034I TN3270B CONN 000000B0 LU MULTIPLE CONN DROP  ERR 1010 143
      IP..PORT: ::FFFF:10.1.100.221..4213          EZBTTRCV

V TCPIP,TCPIPB,DROP,JOBNAME=TN3270B 2
EZZ0060I PROCESSING COMMAND: VARY TCPIP,TCPIPB,DROP,JOBNAME=TN3270B
IKT100I USERID          CANCELED DUE TO UNCONDITIONAL LOGOFF
IKT100I USERID          CANCELED DUE TO UNCONDITIONAL LOGOFF
EZD2013I    3 CONNECTIONS WERE SUCCESSFULLY DROPPED 3
IKT122I IPADDR..PORT 10.1.100.221..4217
IKT100I USERID          CANCELED DUE TO UNCONDITIONAL LOGOFF
IKT122I IPADDR..PORT 10.1.100.221..4218
IKT122I IPADDR..PORT 10.1.100.221..4216
EZZ6034I TN3270B CONN 000000BE LU SC31BB29 CONN DROP  ERR 1010 158
      IP..PORT: ::FFFF:10.1.100.221..4216          EZBTTRCV
EZZ6034I TN3270B CONN 000000C0 LU MULTIPLE CONN DROP  ERR 1010 159
      IP..PORT: ::FFFF:10.1.100.221..4217          EZBTTRCV

```

Note: The VARY TCPIP,,DROP command will drop all connections for a server. The Netstat DROP/-D command supports dropping only one connection per command invocation.

8.3.4 NETSTAT Catalog validation

You can invoke NETSTAT in three ways:

- ▶ TSO
- ▶ MVS console
- ▶ z/OS UNIX command

NETSTAT opens the message catalog in this way:

- ▶ netmsg.cat (IP v4)
- ▶ netmsg6.cat (IP v6)
- ▶ default message text

If the catalogs and command processor are not in synch an ABEND0C4 can occur in the module ONETSTAT.

The following message can be issued:

```
EZZ0157I CONFIGURATION: THE CONFIGURATION COMPONENT HAS TERMINATED
```

This error can occur during z/OS migration when the new z/OS version is pointing to the old load library (TCPIP.SEZALOAD).

8.3.5 Timestamp validation for NETSTAT catalogs

NETSTAT checks the message catalog timestamp against the timestamp expected by the command.

If there is a mismatch, the following message is displayed:

```
EZZ2394I Netstat was expecting netmsg.cat to be at service level HIP61C0 and  
2010 091 20:03 UTC - Netstat is using default messages.
```

When we use a previous z/OS catalog version, we get the message 2008 100 19:39 UTC.ØIBM-1047 instead of 2010 091 20:03 UTC.yIBM-1047.

Maintenance level for the catalog must be at least EZASERVICE Service Level HIP61C0.

8.4 Gathering traces in CS for z/OS IP

Using trace tools is helpful when you have a concern about what is happening in the *flow* of data. Communications Server for z/OS IP provides general trace and application-specific trace facilities. Included in the general traces are packet trace and event trace. Both use the Component Trace (CTRACE) facilities of z/OS.

- ▶ A *packet trace* captures data packets that flow in or out of the IP stack.
- ▶ An *event trace* can capture data flows within the stack, through the application socket interfaces as well as other network flows, such as the ARP process.

This section deals with the trace facilities that are available to analyze TCP/IP problems on z/OS servers and clients. It also discusses how to process those traces.

The MVS component trace can be used to diagnose most TCP/IP problems. Some components of TCP/IP continue to maintain their own tracing mechanisms, for example, the FTP server. Consult *z/OS Communications Server: IP Diagnosis Guide*, GC31-8782, and *z/OS MVS Diagnosis: Tools and Service Aids*, GA22-7589, for more information about the various trace options.

The following TCP/IP traces are available using the component trace:

- ▶ Event trace for TCP/IP stacks (SYSTCPIP)
- ▶ Packet trace (SYSTCPDA)
- ▶ Socket data trace
- ▶ OMPROUTE trace (SYSTCPRT)
- ▶ Resolver trace (SYSTCPRE)
- ▶ Intrusion detection services trace (SYSTCPIS)
- ▶ IKE daemon trace (SYSTCPIK)
- ▶ OSAENTA trace (SYSTCPOT)
- ▶ Network security services server trace (SYSTCPNS)

► Configuration profile trace

Figure 8-3 shows the traces that can be used for debugging. Some applications have their own internal trace functions. The output from those traces can be to the window, a file, or to the syslogd logging function. The data from the z/OS Component Trace is written to either an external writer or the TCP/IP data space TCPIPDS1.

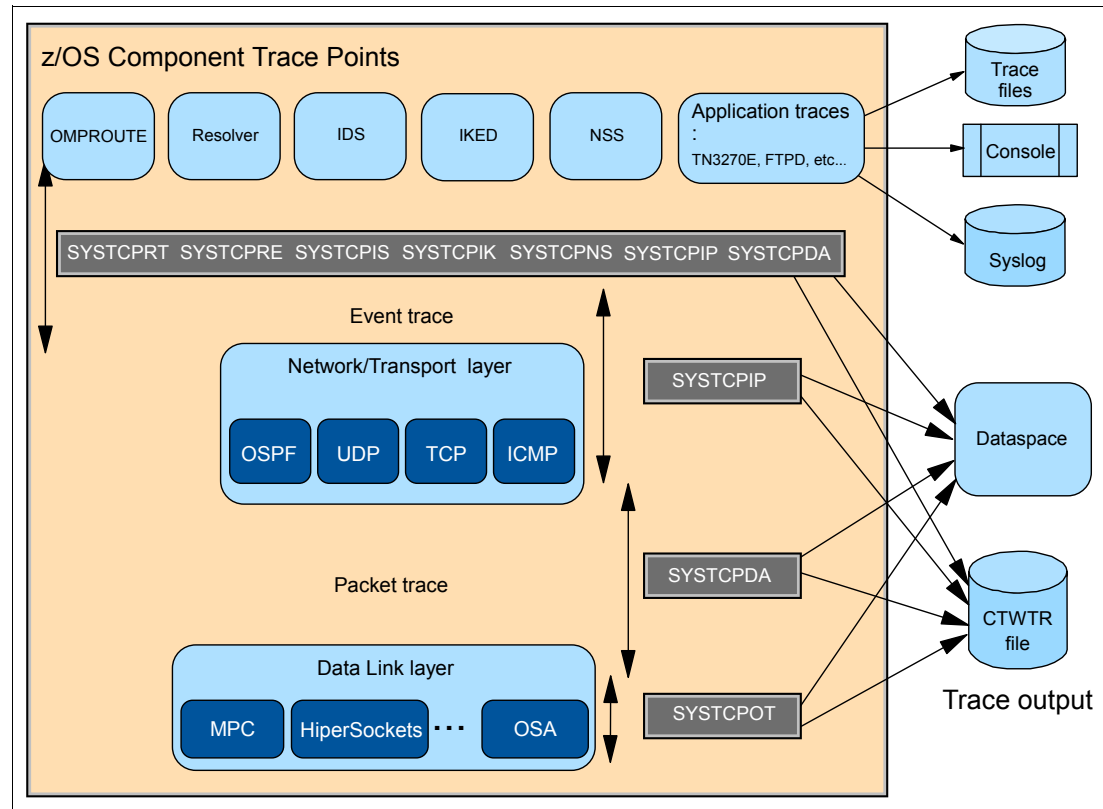


Figure 8-3 Trace points

Information APAR II12014 is a useful source of information about the TCP/IP component and packet trace. For general information about the MVS component trace, see *z/OS MVS Diagnosis: Tools and Service Aids*, GA22-7589.

8.4.1 Taking a component trace

Component trace data is written to either an external writer or the TCP/IP data space TCPIPDS1 (the default is to write trace data to the data space). In this section we show the necessary steps to start a component trace that uses the external writer; this allows you to store trace data in data sets, which can later be used as input to IPCS.

Before starting the traces, create the external write procedure in the SYS1.PROCLIB library, which allocates the trace data set. This procedure is activated using the **trace** command. A sample procedure named CTWTR is shown in Figure 8-4.

```
//CTWTR    PROC
//IEFPROC  EXEC PGM=ITTTTCWR
//TRCOUT01 DD DSN=SYS1.&SYSNAME..CTRACE,
//          VOL=SER=COMST2,UNIT=3390,
//          SPACE=(CYL,10),DISP=(NEW,CATLG),DSORG=PS
//*
```

Figure 8-4 Sample External Write procedure

Next, follow these steps using the **trace** command to activate, capture data, and stop the trace process:

1. Start the external writer (CTRACE writer).

```
TRACE CT,WTRSTART=ctwtr
```

Where *ctwtr* is the name of the procedure created to allocate the trace data set.

2. Start the CTRACE and connect to the external writer.

```
TRACE CT,ON,COMP=component,SUB=(proc_name)
R xx,OPTION=(valid_options),WTR=ctwtr,END
```

Where:

- *component* is the component name of the trace being started and can be any of these:

```
SYSTCPIP (Event trace)
SYSTCPDA (Packet trace)
SYSTCPDA (Data trace)
SYSTCPIS (Intrusion Detection Services trace)
SYSTCPIK (IKE daemon trace)
SYSTCPOT (OSAENTA trace)
SYSTCPNS (Network security services (NSS) server trace)
SYSTCPRT (OMPROUTE trace)
SYSTCPRE (RESOLVER trace)
Configuration profile trace
```

- *proc_name* is the procedure related to the component trace being started, and can be any of these:

```
tcpip_proc
iked_proc
nss_proc
omp_proc
```

- *ctwtr* is the started procedure name of the external writer.

3. To verify the trace is started, use the **display trace** command.

```
DISPLAY TRACE,COMP=component,SUB=(proc_name)
```

4. Perform the operation you want to trace.

5. Disconnect the external writer.

```
TRACE CT,ON,COMP=component,SUB=(proc_name)
R xx,WTR=DISCONNECT,END
```

6. Stop the component trace.

```
TRACE CT,OFF,COMP=component,SUB=(proc_name)
```

7. Stop the external writer.

```
TRACE CT,WTRSTOP=ctwtr
```

The next sections describe each component trace used by z/OS Communications Server - TCP/IP component for documenting problems. For a detailed explanation for each component trace, refer to *z/OS Communications Server: IP Diagnosis Guide*, GC31-8782.

8.4.2 Event trace for TCP/IP stacks (SYSTCPIP)

The TCP/IP event trace, SYSTCPIP, traces TCP/IP stack components such as IP, ARP, TCP, UDP, TELNET, VTAM, and Socket API (SOCKAPI). It is automatically started at TCP/IP initialization using the CTRACE parm option in the parms statement of the TCP/IP sack startup procedure.

z/OS Communications Server provides a default trace options set in the SYS1.PARMLIB member (CTIEZB00 for SYSTCPIP, and CTIEZBTN for the TN3270 Telnet server). The options provided can be changed using an alternate member with the desired options (for example, CTIEZBXX), and then changing the value in the parm CTRACE keyword in your TCP/IP procedure; see Figure 8-5.

Note: The buffer size option is defined during TCP/IP startup only, so any change needs to be done using the CTIEZBxx parmlib member and cannot be reset without restarting the TCP/IP address space. The default is 8 MB.

```
//TCPIP  PROC  PARMS='CTRACE(CTIEZBXX)'  
//*  
//TCPIP  EXEC  PGM=EZBTCPIP,REGION=0M,TIME=1440,  
//        PARM='&PARMS'
```

Figure 8-5 Overriding CTIEZB00 with CTIEZBXX

If you want to specify different trace options after TCP/IP initialization, you can execute the TRACE CT command and either specify the new component trace options file or respond to prompts from the command.

Figure 8-6 shows the status of the component trace for TCP/IP procedure TCPIPA as it has been initialized using SYS1.PARMLIB member CTIEZB01. Note that we have changed the default value for BUFSIZE to 4 M.

```

RESPONSE=SC30
IEE843I 17.09.02 TRACE DISPLAY 466
      SYSTEM STATUS INFORMATION
ST=(ON,0256K,00512K) AS=ON BR=OFF EX=ON MO=OFF MT=(ON,024K)
TRACENAME
=====
SYSTCPIP
                                MODE BUFFER HEAD SUBS
                                =====
                                OFF          HEAD    4
      NO HEAD OPTIONS
SUBTRACE      MODE BUFFER HEAD SUBS
-----
TCPIPA              ON    0004M
  ASIDS             *NONE*
  JOBNAMES           *NONE*
  OPTIONS            MINIMUM
  WRITER             *NONE*
```

Figure 8-6 *DISPLAY TRACE,COMP=SYSTCPIP,SUB=(TCPIPC) output*

The MINIMUM trace option is always active. During minimum tracing, certain exceptional conditions are being traced so the trace records for these events will be available for easier debugging in case the TCP/IP address space should encounter an abend condition.

Socket API trace

The SOCKAPI option for the TCP/IP CTRACE component SYSTCPIP is intended to be used for application programmers to debug problems in their applications. The SOCKAPI option captures trace information related to the socket API calls that an application might issue.

When you need to trace application-related problems using the SOCKAPI option, it is recommended that you follow these guidelines:

- ▶ Trace only one application. Use the job name or ASID option when capturing the trace to limit the trace data to one application.
- ▶ Trace only the SOCKAPI option. To get the maximum number of SOCKAPI trace records, specify only the SOCKAPI option.
- ▶ Use an external writer. The external writer is recommended to save more trace data.
- ▶ Trace only one TCP/IP stack.
- ▶ Activate the data trace only if more data is required. The SOCKAPI trace contains the first 96 bytes of data sent or received, which is usually sufficient.

However, the SOCKET option is primarily intended for use by TCP/IP Service and provides information meant to be used to debug problems in the TCP/IP socket layer, UNIX System Services, or the TCP/IP stack. Refer to *z/OS CS: IP Diagnosis*, GC31-8782, for further details on the SOCKAPI option.

Sample SYSTCPIP trace

In this section we follow the steps described in 8.4.1, “Taking a component trace” on page 316 to start, get data, and stop a CTRACE for component SYSTCPIP. The TN3270 server address space also uses the SYSTCPIP event trace. Therefore, all discussions that follow here where TCP/IP is used also pertain to the Telnet server, with the following exceptions:

- ▶ The Telnet server does not use a data space for trace data collection; it uses its own private storage.
- ▶ A subset of the trace commands are used by Telnet so a new default member, CTIEZBTN, is created, which provides an indication of the trace options available. This member can also be overwritten in the same manner as the TCP/IP parmlib member can be overwritten.
- ▶ A subset of IPCS commands are used by Telnet.

Note: If using the Telnet option, do not specify the JOBNAME parm when starting CTRACE.

The resulting messages are shown after each command, as follows:

1. Start the external writer (CTRACE writer).

```
TRACE CT,WTRSTART=CTWTR
ITT038I ALL OF THE TRANSACTIONS REQUESTED VIA THE TRACE CT COMMAND
WERE SUCCESSFULLY EXECUTED.
IEF196I IGD100I 8623 ALLOCATED TO DDNAME TRCOUT01 DATACLAS (      )
ITT110I INITIALIZATION OF CTRACE WRITER CTWTR COMPLETE.
```

2. Connect to the CTRACE external writer and specify trace options.

```
TRACE CT,ON,COMP=SYSTCPIP,SUB=(TCPIPC)
*060 ITT006A SPECIFY OPERAND(S) FOR TRACE CT COMMAND.
R 60,JOBNAME=(FTPDC),OPTIONS=(SOCKAPI),WTR=CTWTR,END
IEE600I REPLY TO 060 IS;JOBNAME=(FTPDC),OPTIONS=(SOCKAPI),WTR=CTWTR
ITT038I ALL OF THE TRANSACTIONS REQUESTED VIA THE TRACE CT COMMAND
WERE SUCCESSFULLY EXECUTED.
```

Note: You can use the parmlib member CTIEZBxx to provide the same options:

```
TRACE CT,ON,COMP=SYSTCPIP,SUB=(TCPIPC),PARM=(CTIEZBXX)
```


3. Display the active component trace options to verify they are correct.

```

DISPLAY TRACE,COMP=SYSTCPIP,SUB=(TCPIPC)
IEE843I 12.12.22 TRACE DISPLAY 206
      SYSTEM STATUS INFORMATION
ST=(ON,0256K,00512K) AS=ON BR=OFF EX=ON MO=OFF MT=(ON,024K)
TRACENAME
=====
SYSTCPIP
                                MODE BUFFER HEAD SUBS
                                =====
                                OFF          HEAD    3
      NO HEAD OPTIONS
SUBTRACE      MODE BUFFER HEAD SUBS
-----
TCPIPC        ON    0008M
  ASIDS      *NONE*
  JOBNAME    FTPDC
  OPTIONS    SOCKAPI
  WRITER     CTWTR

```

4. Reproduce the failure that you want to trace.

5. Disconnect the external writer.

```

TRACE CT,ON,COMP=SYSTCPIP,SUB=(TCPIPC)
*061 ITT006A SPECIFY OPERAND(S) FOR TRACE CT COMMAND.
R 61,WTR=DISCONNECT,END
IEE600I REPLY TO 061 IS;WTR=DISCONNECT,END
ITT038I ALL OF THE TRANSACTIONS REQUESTED VIA THE TRACE CT COMMAND
WERE SUCCESSFULLY EXECUTED.

```

6. Stop the component trace.

```

TRACE CT,OFF,COMP=SYSTCPIP,SUB=(TCPIPC)
ITT038I ALL OF THE TRANSACTIONS REQUESTED VIA THE TRACE CT COMMAND
WERE SUCCESSFULLY EXECUTED.

```

7. Stop the external writer.

```

TRACE CT,WTRSTOP=CTWTR
ITT038I ALL OF THE TRANSACTIONS REQUESTED VIA THE TRACE CT COMMAND
WERE SUCCESSFULLY EXECUTED.
ITT111I CTRACE WRITER CTWTR TERMINATED BECAUSE OF A WTRSTOP REQUEST.
IEF196I IEF142I CTWTR CTWTR - STEP WAS EXECUTED - COND CODE 0000
IEF196I IEF285I   SYS1.SC32.CTRACE CATALOGED

```

After your events trace data is captured, the trace data set created by the external writer procedure is saved and IPCS is used to format and analyze its contents. Refer to *z/OS CS: IP Diagnosis*, GC31-8782, for further details about SYSTCPIP events trace.

8.4.3 Packet trace (SYSTCPDA)

Packet tracing captures IP packets as they enter or leave TCP/IP. You select what you want to trace using the PKTTRACE statement within the PROFILE.TCPIP, or using the VARY PKTTRACE command entered from the MVS console. RACF authorization is required to execute this command.

You can also use the packet trace to capture data traffic going through fast local socket (local traffic).

With the VARY PKTTRACE command or PKTTRACE statement in PROFILE.TCPIP, you can specify options such as IP address, port number, discard, and protocol type. If you are planning to gather a trace for relatively long hours, or if your system experiences heavy traffic, it is recommended that you specify these filtering options so that TCP/IP does not have to gather unnecessary packets.

To run a packet trace, follow the steps described in 8.4.1, "Taking a component trace" on page 316 to activate the component trace for component SYSTCPDA. Then, activate the packet trace with the desired options and filters. All data in this sample is written to an external writer:

1. Start the CTRACE external writer.

```
TRACE CT,WTRSTART=CTWTR
ITT038I ALL OF THE TRANSACTIONS REQUESTED VIA THE TRACE CT COMMAND
WERE SUCCESSFULLY EXECUTED.
IEF196I IGD100I 8623 ALLOCATED TO DDNAME TRCOUT01 DATACLAS (      )
ITT110I INITIALIZATION OF CTRACE WRITER CTWTR COMPLETE.
```

2. Start the CTRACE and connect the external writer to the TCP/IP address space.

```
TRACE CT,ON,COMP=SYSTCPDA,SUB=(TCPIPA)
063 ITT006A SPECIFY OPERAND(S) FOR TRACE CT COMMAND.
R 63,WTR=CTWTR,END
IEE600I REPLY TO 063 IS;WTR=CTWTR,END
ITT038I ALL OF THE TRANSACTIONS REQUESTED VIA THE TRACE CT COMMAND
WERE SUCCESSFULLY EXECUTED.
```

3. Check that the trace started successfully.

```
D TRACE,COMP=SYSTCPDA,SUB=(TCPIPA)
IEE843I 14.00.29 TRACE DISPLAY 388
      SYSTEM STATUS INFORMATION
TRACENAME
=====
SYSTCPDA
                                MODE BUFFER HEAD SUBS
                                =====
                                OFF          HEAD    2
      NO HEAD OPTIONS
SUBTRACE                      MODE BUFFER HEAD SUBS
NO HEAD OPTIONS
SUBTRACE                      MODE BUFFER HEAD SUBS
-----
TCPIPA                        MIN  0016M
ASIDS      *NONE*
JOBNAMES    *NONE*
OPTIONS     MINIMUM
WRITER      CTWTR
```

4. Start the trace through the PROFILE.TCPIP statement and the VARY OBEYFILE command, or through the V TCPIP,,PKT command.

```
VARY TCPIP,TCPIPA,PKT,ON
EZZ0060I PROCESSING COMMAND: VARY TCPIP,TCPIPA,PKT,ON
EZZ0053I COMMAND VARY PKTTRACE COMPLETED SUCCESSFULLY
```

5. Optionally, modify the trace options to filter the data that is captured using the VARY command. If both options IPAddr and PORTNUM are specified in the same command, an AND condition is created so data is only captured if *both* conditions are met.

For example, issuing the following VARY command records only the packets with both IPAddr=10.1.8.21 *and* PORTNUM=23. Example 8-15 shows the output generated by this command.

Example 8-15 Command to modify the trace options

```
V TCPIP,TCPIPA,PKTTRACE,IP=10.1.8.21,PORTNUM=23
EZZ0053I COMMAND VARY PKTTRACE COMPLETED SUCCESSFULLY
```

It can also create an OR condition issuing multiple VARY commands to apply filters together. For example, if you want to record all packets with destination ports *xx* OR source ports *yy*, use the following commands:

```
VARY TCPIP,tcpprocname,PKT,DEST=xx
VARY TCPIP,tcpprocname,PKT,SRCP=yy
```

When VIPAROUTE statements are defined to a sysplex distributor to select routes, the sysplex distributor encapsulates the packet with a new header before sending it to the target stack. The IPAddr option can allow filtering to be performed on not only the outer packet but the *inner* packet.

Additionally, z/OS Communications Server provides the DISCARD option which allows you to filter inbound packets that are discarded by the stack. You can also filter packet trace collection and formatting by using discard reason codes. For example, if you want to record *all* packets that are discarded or filter the packets with reason code such as 4136, use the commands:

```
VARY TCPIP,TCPIPA,PKT,DISCARD=*
EZZ0060I PROCESSING COMMAND: VARY TCPIP,TCPIPA,PKT,DISCARD=*
EZZ0053I COMMAND VARY PKTTRACE COMPLETED SUCCESSFULLY

VARY TCPIP,TCPIPA,PKT,DISCARD=4136
EZZ0060I PROCESSING COMMAND: VARY TCPIP,TCPIPA,PKT,DISCARD=4136
EZZ0053I COMMAND VARY PKTTRACE COMPLETED SUCCESSFULLY
```

6. Check whether the packet trace options set are correct, using the **netstat dev (-d)** command. We can see the PORTNUM = 23 option **1** , the IPADDR = 10.1.8.21 option **2** and the Discard Code = 4136 option **3**. Example 8-16 shows a sample packet trace setting.

Example 8-16 Packet trace options setting

```
D TCPIP,TCPIPA,N,DEV
...
DEVNAME: LOOPBACK          DEVTYP: LOOPBACK
DEVSTATUS: READY
LNKNAME: LOOPBACK  LNKTYPE: LOOPBACK  LNKSTATUS: READY
ACTMTU: 65535
BSD ROUTING PARAMETERS:
  MTU SIZE: N/A              METRIC: 00
  DESTADDR: 0.0.0.0          SUBNETMASK: 0.0.0.0
PACKET TRACE SETTING:
  PROTOCOL: *                TRRECCNT: 00000000  PCKLENGTH: FULL
  DISCARD:  NONE
  SRCPORT:  *                DESTPORT: *          PORTNUM: 23
  IPADDR:   10.1.8.21       SUBNET:  *
1
2
```

```

        PROTOCOL: *                TRRECCNT: 00000000  PCKLENGTH: FULL
        DISCARD: 4136
        SRCPORT: *                DESTPORT: *          PORTNUM: *
        IPADDR: *                SUBNET: *
        MULTICAST SPECIFIC:
        MULTICAST CAPABILITY: NO
        ...

```

3

7. Perform the operation that you want to trace.

8. Stop the trace.

```

VARY TCPIP,TCPIPA,PKT,OFF
EZZ0060I PROCESSING COMMAND: VARY TCPIP,TCPIPA,PKT,OFF
EZZ0053I COMMAND VARY PKTTRACE COMPLETED SUCCESSFULLY

```

9. Disconnect the external writer from TCP/IP.

```

TRACE CT,ON,COMP=SYSTCPDA,SUB=(TCPIPA)
064 ITT006A SPECIFY OPERAND(S) FOR TRACE CT COMMAND.
R 64,WTR=DISCONNECT,END
IEE600I REPLY TO 064 IS;WTR=DISCONNECT,END
ITT120I SOME CTRACE DATA LOST, LAST 9 BUFFER(S) NOT WRITTEN
ITT038I ALL OF THE TRANSACTIONS REQUESTED VIA THE TRACE CT COMMAND
WERE SUCCESSFULLY EXECUTED.

```

10. Stop the CTRACE.

```

TRACE CT,OFF,COMP=SYSTCPDA,SUB=(TCPIPA)
ITT038I ALL OF THE TRANSACTIONS REQUESTED VIA THE TRACE CT COMMAND
WERE SUCCESSFULLY EXECUTED

```

11. Stop the external writer.

```
TRACE CT,WTRSTOP=CTWTR
ITTO38I ALL OF THE TRANSACTIONS REQUESTED VIA THE TRACE CT COMMAND
WERE SUCCESSFULLY EXECUTED.
ITT111I CTRACE WRITER CTWTR TERMINATED BECAUSE OF A WTRSTOP REQUEST.
IEF196I IEF142I CTWTR CTWTR - STEP WAS EXECUTED - COND CODE 0000
IEF196I IEF285I   SYS1.SC30.CTRACE CATALOGED
```

After the packet trace or the socket data is captured, the trace data set that is created by the external writer procedure is saved. Use IPCS to format and analyze the saved contents. Refer to *z/OS CS: IP Diagnosis*, GC31-8782 for further details about these traces.

Note: The next hop IP address is provided for all outbound packets. This information is only viewable if the packet trace is formatted with the “FULL” option, and also available externally by way of the real-time packet trace NMI. Additionally, CTRACE with `OPTIONS((LAST IPADDR(ipaddress)))` can select packets for the inner IP address.

Socket data trace

Using the SYSTCPDA component CS for z/OS IP provides a way to capture socket data into and out of the Physical File System (PFS). It helps to diagnose application data-related problems. To activate this trace, we follow the same steps we used to activate a packet trace and change only the command to start and stop the socket data trace:

```
V TCPIP,tcpproc,DATTRACE,ON
V TCPIP,tcpproc,DATTRACE,OFF
```

A `PORTNUM` parameter is supported on the `VARY TCPIP,,DATTRACE` command that you can use to trace only packets that have a source or destination port that matches a specific port number.

The Socket data trace options can modify the data being captured using the `VARY` command. For example, issuing the `VARY` command records only the data packets with both `IPaddr (10.1.9.11)` and `Portnum= 21`, as shown in Example 8-17.

Example 8-17 V TCPIP,TCPIPB,DAT options

```
V TCPIP,TCPIPB,DAT,JOBNAME=*,IP=10.1.9.11/32,PORTNUM=21
EZZ0060I PROCESSING COMMAND: VARY TCPIP,TCPIPB,DAT,JOBNAME=*,
IP=10.1.9.11/32,PORTNUM=21
EZZ0053I COMMAND VARY DATTRACE COMPLETED SUCCESSFULLY
```

To verify if the data trace options setting are correct, use the `NETSTAT CONFIG` command. See Example 8-18.

Example 8-18 NETSTAT CONFIG command shows the options

D TCPIP,TCPIPB,N,CONFIG

```
EZD0101I NETSTAT CS V1R12 TCPIPB 541
TCP CONFIGURATION TABLE:
DEFAULTRCVBUFSIZE: 00262144  DEFAULTSNDBUFSIZE: 00262144
DEFLTMAXRCVBUFSIZE: 00524288  SOMAXCONN: 000000010
MAXRETRANSMITTIME: 120.000  MINRETRANSMITTIME: 0.500
ROUNDTRIPGAIN: 0.125  VARIANCEGAIN: 0.250
VARIANCEMULTIPLIER: 2.000  MAXSEGLIFETIME: 30.000
DEFAULTKEEPALIVE: 00000015  DELAYACK: NO
RESTRICTLOWPORT: NO  SENDGARBAGE: NO
```

```

TCPTIMESTAMP:      YES      FINWAIT2TIME:      600
TTLS:              NO
DATA TRACE SETTING:
JOBNAME: *          TRRECCNT: 00000000  LENGTH: FULL
IPADDR/PREFIXLEN: 10.1.9.11/32
PORTNUM: *
JOBNAME: *          TRRECCNT: 00000000  LENGTH: FULL
IPADDR/PREFIXLEN: 10.1.9.11/32
PORTNUM: 21

```

Start and end records in data flow

z/OS Communications Server provides socket data trace indicating data flow start and end. The state field is used for writing these start and end records. A start record shows first socket read/write, and an end record shows that socket closed. See Example 8-19 for data flow start and end records in a formatted data trace

Example 8-19 Sample formatted trace for start and end records

```

COMPONENT TRACE SHORT FORMAT
  SYSNAME(SC30)
  COMP(SYSTCPDA)SUBNAME((TCPIPA))
  z/OS TCP/IP Packet Trace Formatter, Copyright IBM Corp. 2000, 2010; 2010.067
  DSNNAME('SYS1.SC30.TCPIPA.CTRACE')
  **** 2010/09/28
  RcdNr Sysname Mnemonic Entry Id   Time Stamp   Description
  -----
  129 SC30      DATA      00000005 13:19:18.657635 Data Trace
  To Jobname    : FTPDA                                     Full=0
  Tod Clock     : 2010/09/28 13:19:18.657635                Cid: 00000205
  Domain        : AF_Inet6      Type: Stream                Protocol: TCP
  State        : API Data Flow Starts
  Segment #     : 0                               Flags: Out
  Source        : ::ffff:10.1.1.10
  Destination   : ::ffff:10.1.100.223
  Source Port   : 20                               Dest Port: 1141  Asid: 005B TCB: 007FF1D
  -----
  131 SC30      DATA      00000005 13:19:18.661580 Data Trace
  From Jobname   : FTPDA                                     Full=0
  Tod Clock      : 2010/09/28 13:19:18.661580                Cid: 00000205
  Domain         : AF_Inet6      Type: Stream                Protocol: TCP
  State        : API Data Flow Ends
  Segment #      : 0                               Flags: None
  Source         : ::ffff:10.1.100.223
  Destination    : ::ffff:10.1.1.10
  Source Port    : 1141                               Dest Port: 20   Asid: 005B TCB: 007FF1D
  -----

```

Data trace records for the socket data flow start and end are only supported on *TCP* and *UDP* sockets, they are not supported on *RAW* sockets.

8.4.4 OMPROUTE trace (SYSTCPRT)

To diagnose OMPROUTE problems, z/OS Communications Server provides the debug and trace parameter that can be defined during OMPROUTE initialization. The resulting output is

written to the OMPROUTE log and can cause increased overhead. This performance issue can be solved by using the CTRACE facility. To do so, we highly recommend that you use the OMPROUTE option (DEBUGTRC) in the startup procedure, which changes the output destination of the OMPROUTE trace. In this section we briefly describe how to define and use CTRACE to debug OMPROUTE problems.

The OMPROUTE CTRACE can be started anytime by using the command TRACE CT, or it can be activated during OMPROUTE initialization. If not defined, OMPROUTE component trace is started with a buffer size of 1 MB and the MINIMUM tracing option.

A parmlib member can be used to customize the parameters and to initialize the trace. The default OMPROUTE Component Trace parmlib member is the SYS1.PARMLIB member CTIORA00. The parmlib member name can be changed by using the OMPROUTE_CTRACE_MEMBER environment variable.

In addition to specifying the trace options, you can also change the OMPROUTE trace buffer size. The buffer size can be changed only at OMPROUTE initialization. The maximum OMPROUTE trace buffer size is 100 MB. The OMPROUTE REGION size in the OMPROUTE catalog procedure must be large enough to accommodate a large buffer size.

Here we shown the necessary steps to start the CTRACE for OMPROUTE during OMPROUTE initialization using the parmlib member CTIORA00 and directing the trace output to an external writer.

1. Prepare the SYS1.PARMLIB member CTIORA00 to get the desired output data. Example 8-20 shows a sample of CTIORA00 contents.

Example 8-20 CTIORA00 sample

```

TRACEOPTS
/* ----- */
/*  Optionally start external writer in this file (use both      */
/*  WTRSTART and WTR with same wtr_procedure)                  */
/* ----- */
      WTRSTART(CTWTR)
/* ----- */
/*  ON OR OFF: PICK 1                                           */
/* ----- */
      ON
/*  OFF                                                         */
/* ----- */
/*  BUFSIZE: A VALUE IN RANGE 128K TO 100M                      */
/*          CTRACE buffers reside in OMPROUTE Private storage  */
/*          which is in the regions address space.              */
/* ----- */
      BUFSIZE(50M)
      WTR(CTWTR)
/* ----- */
/*  OPTIONS: NAMES OF FUNCTIONS TO BE TRACED, OR "ALL"         */
/* ----- */
/*  OPTIONS(                                                    */
/*          'ALL'                                                */
/*          , 'MINIMUM'                                          */
/*          , 'ROUTE'                                            */
/*          , 'PACKET'                                           */
/*          , 'OPACKET'                                          */
/*          , 'RPACKET'                                          */

```

```

/*          , 'IPACKET '          */
/*          , 'SPACKET '          */
/*          , 'DEBUGTRC'          */
/*                                     )          */

```

2. Start the OMROUTE procedure using the desired Debug and Trace options, as shown in Example 8-21.

Example 8-21 OMROUTE procedure

```

//OMPC32 PROC STDENV=OMPENC&SYSCONE
//OMPC32 EXEC PGM=OMROUTE,REGION=OM,TIME=NOLIMIT,
//          PARM=(' POSIX(ON) ALL31(ON) ',
//          'ENVAR("_BPXK_SETIBMOPT_TRANSPORT=TCPIPC"',
//          '" _CEE_ENVFILE=DD:STDENV")/-d1 -t2')
//STDENV DD DISP=SHR,DSN=TCPIPC.TCPPARMS(&STDENV)
//SYSPRINT DD SYSOUT=*
//SYSOUT DD SYSOUT=*
//CEEDUMP DD SYSOUT=*,DCB=(RECFM=FB,LRECL=132,BLKSIZE=132)

```

The description for the tag shown in Example 8-21 is as follows:

1The parameters **-t** (trace) and **-d** (debug) define how detailed we want the output data to be. We recommend using **-t2** and **-d1**.

3. Verify that CTRACE has been started as expected, issuing the console command as shown:

```

D TRACE,COMP=SYSTCPRT,SUB=(OMPC)
IEE843I 16.31.37 TRACE DISPLAY 058
          SYSTEM STATUS INFORMATION
ST=(ON,0256K,00512K) AS=ON BR=OFF EX=ON MO=OFF MT=(ON,024K)
TRACENAME
=====
SYSTCPRT
                                MODE BUFFER HEAD SUBS
                                =====
                                OFF          HEAD    1
          NO HEAD OPTIONS
SUBTRACE                                MODE BUFFER HEAD SUBS
-----
OMPC          ON    0010M
  ASIDS      *NONE*
  JOB NAMES  *NONE*
  OPTIONS    MINIMUM ,DEBUGTRC
  WRITER     CTWTR

```

4. You can also use TRACE CT command to define the options we want after OMROUTE has been initialized, as shown:

```

TRACE CT,ON,COMP=SYSTCPRT,SUB=(OMPC)
R 66,OPTIONS=(ALL),END
IEE600I REPLY TO 066 IS;OPTIONS=(ALL),END
ITTO38I ALL OF THE TRANSACTIONS REQUESTED VIA THE TRACE CT COMMAND
WERE SUCCESSFULLY EXECUTED.

```

5. Reproduce the problem.

6. Disconnect the external writer.

```
TRACE CT,ON,COMP=SYSTCPRT,SUB=(OMPC)
067 ITT006A SPECIFY OPERAND(S) FOR TRACE CT COMMAND.
R 67,WTR=DISCONNECT,END
IEE600I REPLY TO 067 IS;WTR=DISCONNECT,END
ITT038I ALL OF THE TRANSACTIONS REQUESTED VIA THE TRACE CT COMMAND
WERE SUCCESSFULLY EXECUTED.
```

7. Stop the component trace.

```
TRACE CT,OFF,COMP=SYSTCPRT,SUB=(OMPC)
ITT038I ALL OF THE TRANSACTIONS REQUESTED VIA THE TRACE CT COMMAND
WERE SUCCESSFULLY EXECUTED.
```

8. Stop the external writer.

```
TRACE CT,WTRSTOP=CTWTR
ITT038I ALL OF THE TRANSACTIONS REQUESTED VIA THE TRACE CT COMMAND
WERE SUCCESSFULLY EXECUTED.
ITT111I CTRACE WRITER CTWTR TERMINATED BECAUSE OF A WTRSTOP REQUEST.
IEF196I IEF142I CTWTR CTWTR - STEP WAS EXECUTED - COND CODE 0000
IEF196I IEF285I SYS1.SC32.CTRACE CATALOGED
```

9. Change the OMPROUTE Debug and Trace level to avoid performance problems using the MODIFY command, as shown in:

```
F OMPC,TRACE=0
EZZ7866I OMPROUTE MODIFY COMMAND ACCEPTED
F OMPC,DEBUG=0
EZZ7866I OMPROUTE MODIFY COMMAND ACCEPTED
```

After these steps, the generated trace file must be formatted using the IPCS. For further information about OMPROUTE diagnosis, refer to *z/OS Communications Server: IP Diagnosis Guide*, GC31-8782.

8.4.5 Resolver trace (SYSTCPRE)

z/OS Communications Server provides component trace support for the resolver. A default minimum component trace is always started during resolver initialization. To customize the parameters used to initialize the trace, update SYS1.PARMLIB member CTIRES00. In addition to specifying the trace options, you can change the resolver trace buffer size. Note that the buffer size can be changed *only* at resolver initialization.

After resolver initialization, you must use the TRACE CT command to change component trace options.

To gather the component trace for the resolver, use the commands listed in 8.4.1, “Taking a component trace” on page 316 and, in step 2 on page 317, specify the **comp=** parameter with the resolver component name, SYSTCPRE and the **sub=** parameter with the resolver proc_name.

The generated trace file created after the problem is reproduced must be formatted using the IPCS. For further information about resolver diagnosis, refer to *z/OS Communications Server: IP Diagnosis Guide*, GC31-8782.

8.4.6 IKE daemon trace (SYSTCPIK)

z/OS Communications Server provides component trace support for the IKE daemon and, as other components, a default minimum component trace is always started during IKE daemon initialization. Use a parmlib member to customize the parameters that are used to initialize the trace. The default IKE daemon component trace parmlib member is the SYS1.PARMLIB member CTIIKE00. The parmlib member name can be changed using the IKED_CTRACE_MEMBER environment variable.

Tip: The IKE daemon reads the IKED_CTRACE_MEMBER environment variable only during initialization. Changes to IKED_CTRACE_MEMBER after daemon initialization have no affect. After IKE daemon initialization, you must use the TRACE CT command to change component trace options.

After IKE daemon is initialized you can start CTRACE to modify trace options or send data to an external writer, using the commands listed in 8.4.1, “Taking a component trace” on page 316 and, in step 2 on page 317, specify the **comp=** parameter with the IKE daemon component name, SYSTCPIK and the **sub=** parameter with the iked proc_name.

The generated trace file created after the problem is reproduced must be formatted using the IPCS. For further information about IKE daemon diagnosis, refer to *z/OS Communications Server: IP Diagnosis Guide*, GC31-8782.

8.4.7 Intrusion detection services trace (SYSTCPIS)

When the TCP/IP stack starts, it reads SYS1.PARMLIB member CTIIDS00, which contains trace options for the SYSTCPIS trace. Packets are traced based on the IDS policy defined in LDAP. Refer to “Intrusion Detection Services” in *z/OS Communications Server: IP Configuration Guide*, SC31-8775, for information about defining policy.

If the EZZ4210I message indicates the parmlib member name CTIIDS00, then the IDS CTRACE space is set up using the default BUFSIZE of 32 M.

The CTIIDS00 member is used to specify the IDS CTRACE parameters. To eliminate this message, ensure that a CTIIDS00 member exists within Parmlib and that the options are correctly specified. A sample CTIIDS00 member is shipped with z/OS Communications Server.

See *z/OS Communications Server: IP Diagnosis Guide*, GC31-8782 for details about the intrusion detection services trace and *z/OS Communications Server: IP Configuration Guide*, SC31-8775 for information about defining policy.

8.4.8 OSAENTA trace (SYSTCPOT)

TCP/IP Services component trace is also available for use with the OSA-Express Network Traffic Analyzer (OSAENTA) trace facility. The OSAENTA trace is a diagnostic method for obtaining frames flowing to and from an OSA adapter. You can use the OSAENTA statement to copy frames as they enter or leave an OSA adapter for an attached host. The host can be an LPAR with z/OS, VM, or Linux. For more information about OSAENTA, refer to 8.5, “OSA-Express3 Network Traffic Analyzer” on page 333.

8.4.9 Queued Direct I/O Diagnostic Synchronization

Communications Server provides support for a new Queued Direct I/O Diagnostic Synchronization (QDIOSYNC) facility. It provides the ability to synchronize OSA-Express3 diagnostic data with host diagnostic data. The QDIOSYNC facility also provides for optional filtering of the OSA-Express3 diagnostic data. If a filter is specified, the OSA-Express3 adapter honors the filter by limiting the types of diagnostic data collected. Although the QDIOSYNC trace differs from a traditional VTAM TRACE command, you use the VTAM MODIFY TRACE and NOTRACE commands to control it. The DISPLAY TRACES command is modified to show the state of the QDIOSYNC trace.

The QDIOSYNC trace is not a traditional trace in which output is generated based on specific events. Instead, the QDIOSYNC trace freezes and captures (logs) OSA-Express3 diagnostic data in a timely manner. In addition to (or instead of) using the hardware management console (HMC) to manually capture the diagnostic data, you can arm the OSA-Express3 adapter to automatically capture diagnostic data when one of the following occurs:

- ▶ The OSA-Express3 adapter detects an unexpected loss of host connectivity.
Unexpected loss of host connectivity occurs when the OSA-Express3 adapter receives an unexpected halt signal from the host or when the host is unresponsive to OSA requests.
- ▶ The OSA-Express3 adapter receives a CAPTURE signal from the host.
A CAPTURE signal is sent by the host when one of the following occurs:
 - The VTAM-supplied message processing facility (MPF) exit (IUTLLCMP) is driven.
 - Either the VTAM or TCP/IP functional recovery routine (FRR) is driven with ABEND06F. (ABEND06F is the result of a SLIP PER trap that specifies ACTION=RECOVERY).

When arming an OSA-Express3 adapter for QDIOSYNC, you can specify an optional filter that alters what type of diagnostic data is collected by the OSA-Express3 adapter. This filtering reduces the overall amount of diagnostic data collected, and therefore decreases the likelihood that pertinent data is lost.

Note: Do not use QDIOSYNC to unconditionally arm an OSA-Express3 adapter when it is shared by other operating systems and those operating systems might use this function. In this case, the function should be coordinated between all sharing operating systems.

If you have several OSAs to arm, but you do not want to arm all of them, consider first arming all OSAs and then individually disarm those you do not want armed.

For more information about how to set up the trace, refer to:

- ▶ *z/OS Communications Server: SNA Diagnosis Vol. 1, Techniques and Procedures*, GC31-6850
- ▶ *z/OS Communications Server: SNA Operation*, SC31-8779
- ▶ *MVS Installation Exits*, SA22-7593

8.4.10 Network security services server trace (SYSTCPNS)

z/OS Communications Server provides component trace support for the Network Security Services (NSS) and, as with other components, a default minimum component trace is always started during NSS server initialization. Use a parmlib member to customize the parameters that are used to initialize the trace. The default NSS server component trace parmlib member is the SYS1.PARMLIB member CTINSS00. In addition to specifying the

trace options, you can also change the NSS trace buffer size. The buffer size can be changed only at NSS initialization and has a maximum of 256 MB.

You can change the parmlib member name using the NSSD_CTRACE_MEMBER environment variable.

Tip: The NSS server reads the NSSD_CTRACE_MEMBER environment variable *only* during initialization. Changes to NSSD_CTRACE_MEMBER after server initialization have no effect.

After the NSS server is initialized, you can start CTRACE to modify trace options or send data to a external writer by using the commands listed in 8.4.1, “Taking a component trace” on page 316 and, in step 2 on page 317, specify the **comp=** parameter with the NSS server component name, SYSTCPNS and the **sub=** parameter with the nss_proc_name.

The generated trace file created after the problem iss reproduced must be formatted using the IPCS. For more information about NSS server diagnosis, refer to *z/OS Communications Server: IP Diagnosis Guide*, GC31-8782.

8.4.11 Obtaining component trace data with a dump

If the TCP/IP or user's address space abends, TCP/IP recovery dumps the home ASID, the primary ASID, the secondary ASID, and the TCPIPDS1 data space to the data sets defined within your MVS environment. The TCPIPDS1 data space contains the trace data for SYSTCPIP, SYSTCPDA, and SYSTCPIS components.

To obtain a dump of the TCP/IP stack when no abend has occurred, use the DUMP command. Remember to specify the data space name, which is always TCPIPDS1, because it contains the trace data for the SYSTCPIP, SYSTCPDA, and SYSTCPIS components. Be sure to include “region” (RGN) in the SDATA dump options, as shown here:

```
DUMP COMM=(enter_dump_title_here)
Rxx,JOBNAME=tcpproc,DSPNAME=('tcpproc'.TCPIPDS1),CONT
Rxx,SDATA=(CSA,LSQA,NUC,PSA,RGN,SQA,SUM,SQA,TRT),END
```

To obtain a dump of the OMPROUTE, RESOLVER, or TELNET address space (which contains the trace table), use the DUMP command as shown here:

```
DUMP COMM=(enter_dump_title_here)
Rxx,JOBNAME=proc_started_task_name,SDATA=(RGN,CSA,ALLPSA,SQA,SUM,TRT,ALLNUC),END
```

For more information about how to get a dump, refer to *z/OS Communications Server: IP Diagnosis Guide*, GC31-8782.

8.4.12 Analyzing a trace

You can format component trace records using IPCS panels or a combination of IPCS panels and the CTRACE command, either from a dump or from external writer files. You can also use IPCS in batch to print a component trace.

The primary purpose of the component trace is to capture data that the IBM Support Center can use in diagnosing problems. There is little information in the documentation on interpreting trace data. If you want to analyze the packet trace or data trace, you can do so by formatting the trace data using a z/OS tool in TSO called IPCS. For more information about trace and dump analysis using IPCS, refer to *z/OS Communications Server: IP Diagnosis Guide*, GC31-8782.

8.4.13 Configuration profile trace

You can use the ITRACE statement in the PROFILE.TCPIP data set to activate TCP/IP runtime tracing for configuration, the TCP/IP SNMP subagent, commands, and the autolog subtask. ITRACE should only be set at the direction of an IBM Service representative. For more information, refer to *z/OS Communications Server: IP Diagnosis Guide*, GC31-8782.

8.5 OSA-Express3 Network Traffic Analyzer

When data problems occur in a LAN environment, multiple traces are usually required. A sniffer trace might be required to see the data as it was received from or sent to the network. An OSA hardware trace might be required if the problem is suspected in the OSA, and z/OS Communications Server traces are required to diagnose VTAM or TCP/IP problems.

To assist in problem diagnosis, the OSA-Express network traffic analyzer (OSAENTA) function provides a way to trace inbound and outbound frames for an OSA-Express3 feature. The OSAENTA trace function is controlled and formatted by z/OS Communications Server, but is collected in the OSA at the network port.

Note: To enable the OSA-Express network traffic analyzer, you must be running at least an IBM System z9 EC or z9 BC and OSA-Express3 feature in QDIO mode (CHPID type OSD). See the 2094DEVICE Preventive Service Planning (PSP) and the 2096DEVICE Preventive Service Planning (PSP) buckets for more information about these topics.

This section discusses the steps that are necessary for setting up and using OSAENTA:

- ▶ Determining the microcode level for OSA-Express3
- ▶ Defining TRLE definitions
- ▶ Checking TCPIP definitions
- ▶ Customizing OSA-Express Network Traffic Analyzer
- ▶ Defining a resource profile in RACF
- ▶ Allocating a VSAM linear data set
- ▶ Starting the OSAENTA trace

8.5.1 Determining the microcode level for OSA-Express3

There are two ways to determine the OSA-Express3 microcode level: from the Hardware Management Console (HMC), or by issuing the D NET,TRL,TRLE=OSA2080P command. Each method is discussed in more detail in this section. From the HMC:

1. Select your system.
2. Double-click **OSA Advanced Facilities**.
3. Select appropriate PCHID.
4. Select **View code level**.

Figure 8-7 shows the microcode level installed in one of our OSA-Express3 features.

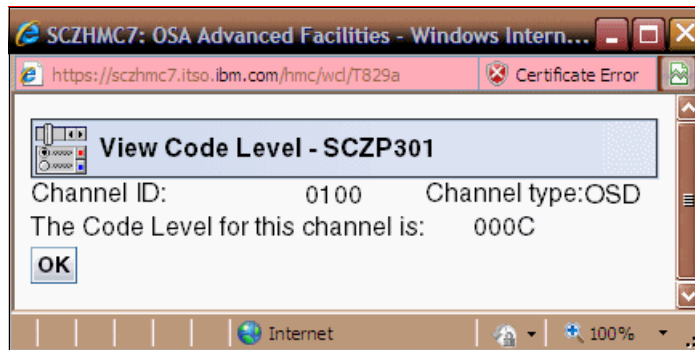


Figure 8-7 View code level

Alternatively, you can issue the D NET,TRL,TRLE=OSA2080T command. Example 8-22 shows the output.

Example 8-22 Output Display TRL

```

NAME = OSA2080T, TYPE = TRLE 068
  TRL MAJOR NODE = OSA2080
STATUS= ACTIV, DESIRED STATE= ACTIV
TYPE = LEASED           , CONTROL = MPC , HPDT = YES
MPCLEVEL = QDIO        MPCUSAGE = SHARE
PORTNAME = OSA2080     PORTNUM = 0   OSA CODE LEVEL = 000C
CHPID TYPE = OSD       CHPID = 02
HEADER SIZE = 4096 DATA SIZE = 0 STORAGE = ***NA***
WRITE DEV = 2081 STATUS = ACTIVE      STATE = ONLINE
HEADER SIZE = 4092 DATA SIZE = 0 STORAGE = ***NA***
READ  DEV = 2080 STATUS = ACTIVE      STATE = ONLINE
-----
DATA  DEV = 2082 STATUS = ACTIVE      STATE = N/A
I/O TRACE = OFF TRACE LENGTH = *NA*
```

8.5.2 Defining TRLE definitions

Use the D U,,,2080,16 command to ensure that you have defined enough devices, as shown in Example 8-23.

Example 8-23 Verifying the number of OSA devices

```

D U,,,2080,16
IEE457I 16.50.55 UNIT STATUS 833
UNIT TYPE STATUS      VOLSER    VOLSTATE
2080 OSA  A-BSY
2081 OSA  A
2082 OSA  A-BSY
2083 OSA  0
2084 OSA  0
2085 OSA  0
2086 OSA  0
2087 OSA  0
2088 OSA  0
2089 OSA  0
208A OSA  0
208B OSA  0
```

```

208C OSA 0
208D OSA 0
208E OSA 0
208F OSAD 0-RAL

```

The OSA-Express3 needs an additional “DATAPATH” statement on the TRL (see Example 8-24).

Example 8-24 TRL definition

```

OSA2080 VBUILD TYPE=TRL
OSA2080T TRLE LNCTL=MPC,
              READ=2080,
              WRITE=2081,
              DATAPATH=(2082-208E),
              PORTNAME=OSA2080,
              MPCLEVEL=QDIO

```

8.5.3 Checking TCPIP definitions

An excerpt of TCP/IP profile, displayed in Example 8-25, shows the information that you need when starting the OSAENTA trace in a later step. Keep this information available.

Example 8-25 TCP/IP definitions

```

;OSA DEFINITION
DEVICE OSA2080 MPCIPA
LINK OSA2080L IPAQENET OSA2080 VLANID 10
HOME
    10.1.2.11 OSA2080L
START OSA2080

```

After TCP/IP is started, you can also see the OAT entries using OSA/SF (see Example 8-26).

Example 8-26 OAT entries

				Image 2.3 (A16)		CULA 0				
00	(2080)	*	MPC	N/A		OSA2080P	(QDIO control)	SIU	ALL	
02	(2082)		MPC	00	No4	No6	OSA2080P	(QDIO data)	SIU	ALL
				VLAN		10	(IPv4)			
				Group Address		Multicast Address				
				01005E000001		224.000.000.001				
				VMAC		IP address				
HOME		00096B1A7490		010.001.000.010						
HOME		00096B1A7490		010.001.001.010						
HOME		00096B1A7490		010.001.002.010						
HOME		00096B1A7490		010.001.002.011						
REG		00096B1A7490		010.001.002.012						
REG		00096B1A7490		010.001.003.011						
REG		00096B1A7490		010.001.003.012						
REG		00096B1A7490		010.001.004.011						
REG		00096B1A7490		010.001.005.011						
REG		00096B1A7490		010.001.006.011						
REG		00096B1A7490		010.001.007.011						

REG	00096B1A7490	010.001.008.010
REG	00096B1A7490	010.001.008.020

03 (2083) N/A

N/A CSS

8.5.4 Customizing OSA-Express Network Traffic Analyzer

Use this task to select an OSA-Express Network Traffic Analyzer (NTA) support element control, to customize the OSA-Express NTA settings in Advanced Facilities, or to check the current OSA-Express NTA authorization.

Customizing OENTA allows the following activities for the support element:

- ▶ Set up the OSA LAN Analyzer traces and capture data to the support element hard disk
- ▶ Change authorization to allow host operating systems to enable the NTA traces outside their own partition

Note: The OSA-Express Network Traffic Analyzer is mutually exclusive with the OSA LAN Analyzer for tracing on a specified CHPID. Only one or the other can be enabled for a specified CHPID at any one time.

To accomplish this customization, change the current OSA-Express NTA control. Using the HMC, follow these steps:

1. Log on to the Support Element (SE) on the Hardware Management Console (HMC) through Single Object Operations (SOO).

Important: Enabling the OENTA support can allow tracing of sensitive information. Therefore, the user ID used to do the following steps must have the “Access Administrator Tasks” role assigned.

2. Select the CPC you want to work with, as shown in Figure 8-8 on page 336.

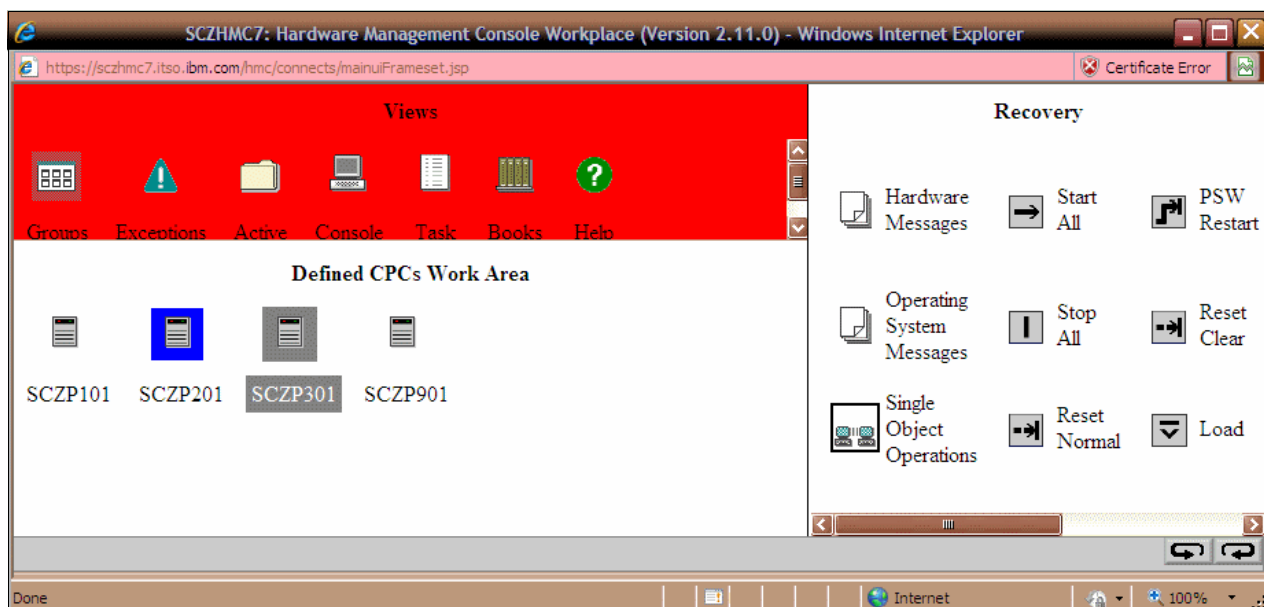


Figure 8-8 From the HMC, log on to SE

3. Select and open the Service task list, as shown in Figure 8-9.

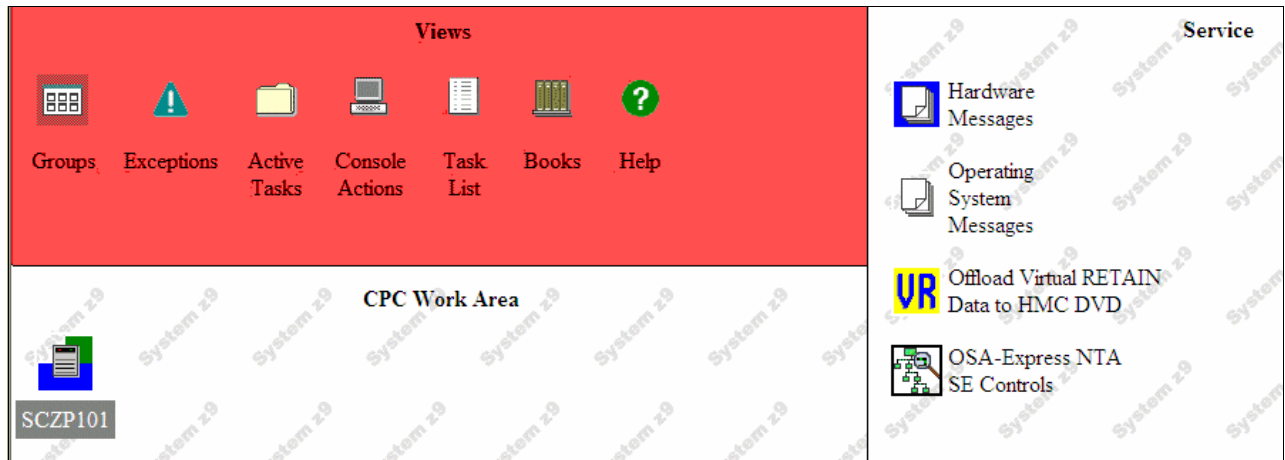


Figure 8-9 OSA-Express NTA

4. Double-click the OSA-Express NTA SE Controls task; see Figure 8-10.

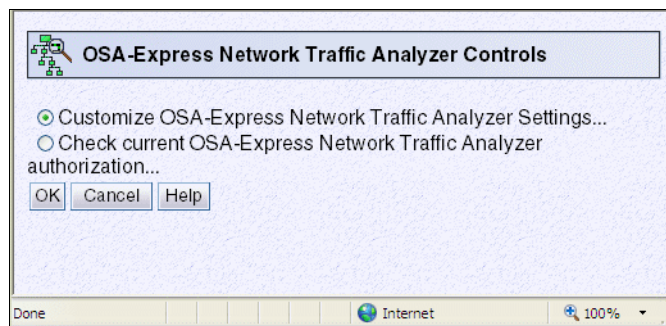


Figure 8-10 OSA NTA Controls

5. Select the control to work with:

- Customize OSA-Express Network Traffic Analyzer Settings... provides the capability to allow or disallow the support element to change authorization to allow host operating systems to enable the Network Traffic Analyzer to trace outside their own partition.
- Check current OSA-Express Network Traffic Analyzer authorization... allows the support element to scan all the OSAs and reports back which OSAs are authorized for NTA to trace outside its own partition.

6. Click **OK** to change the current OSA-Express NTA control; see Figure 8-11.

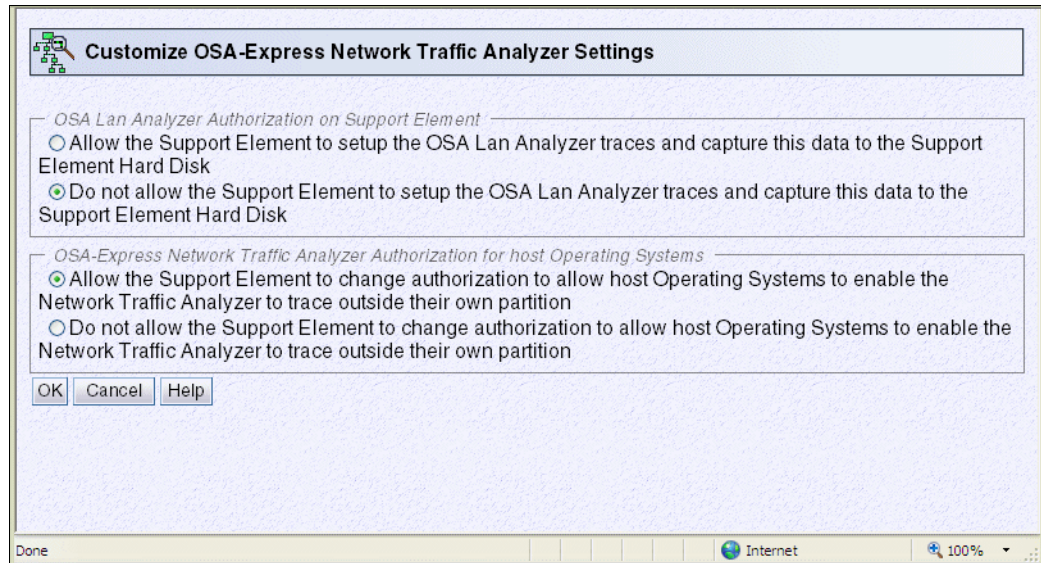


Figure 8-11 Change the current OSA-Express NTA control

7. Click **Allow the Support Element to allow Host Operating System to enable NTA**.
8. Click **OK**; see Figure 8-12.

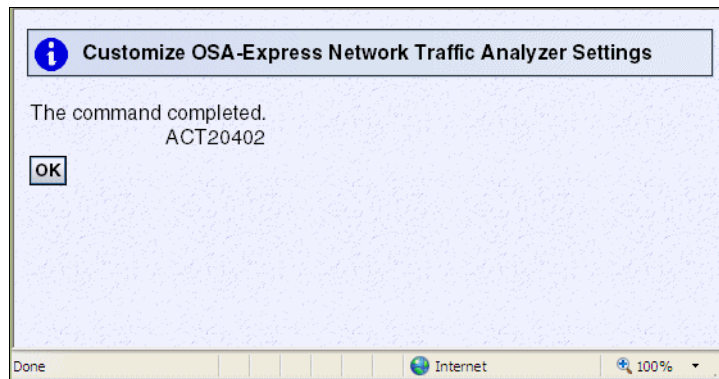


Figure 8-12 Command completed

9. Log off from the SE and from the HMC.
10. Log on to the SE on the HMC through SOO (Single Object Operations) using the SYSPROG user ID.
11. Select **Channels work area** (on the left side of the window) **Channel Operation** (on the right side of the window).

12. Select the channel that you want to manage (see Figure 8-13).

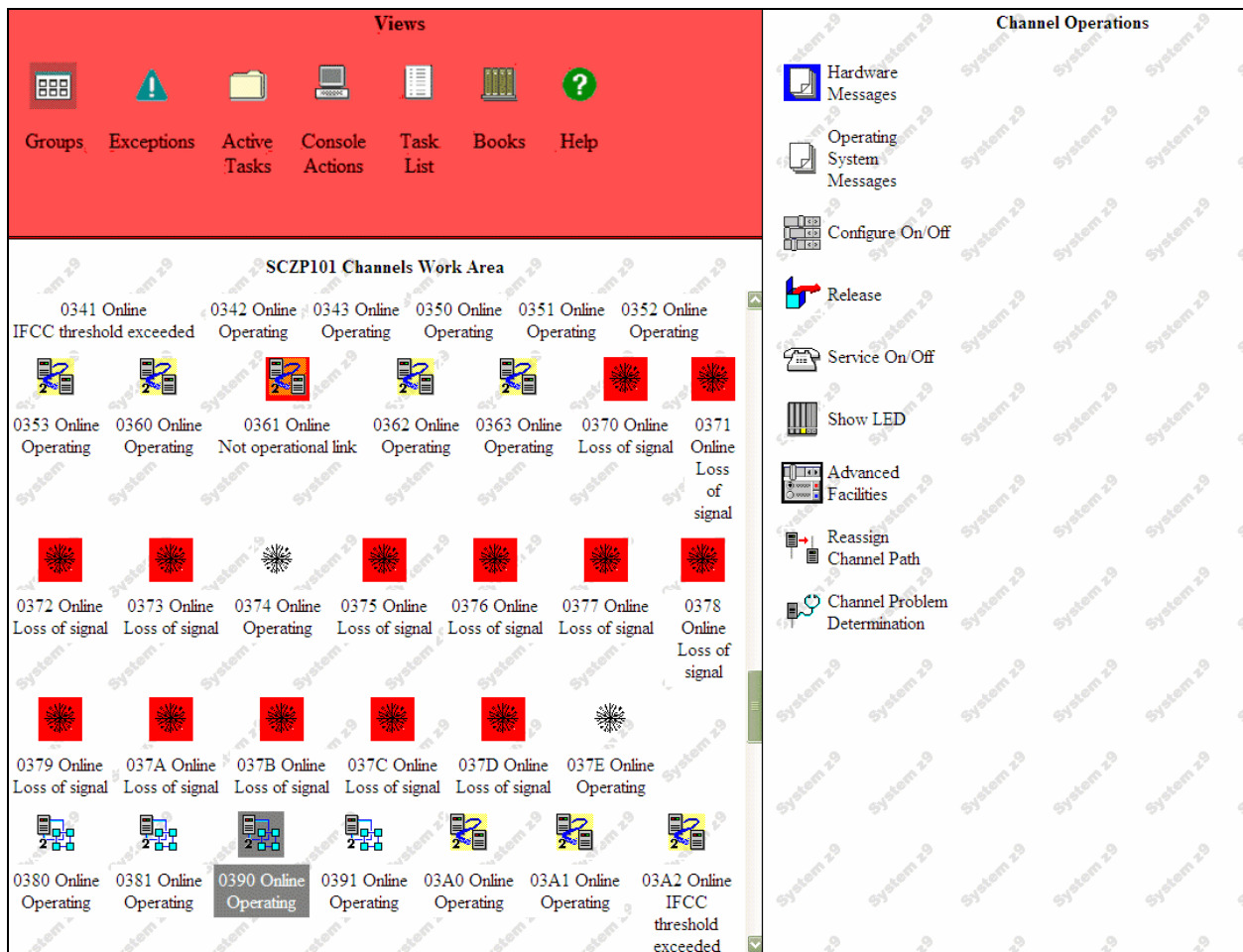


Figure 8-13 Channel Operations menu

13. In our case we selected PCHID 0390 (CHPID 02). We double-clicked **Advanced Facilities**; see Figure 8-14.

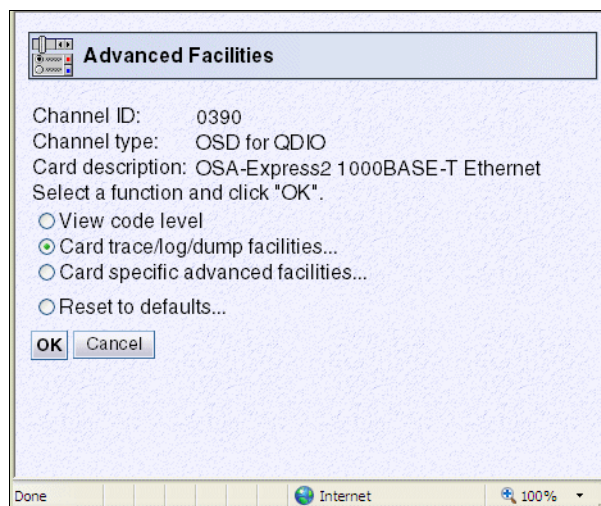


Figure 8-14 Advanced Facilities options

14. Select **Card trace/log/dump facilities**, then click **OK**; see Figure 8-15.

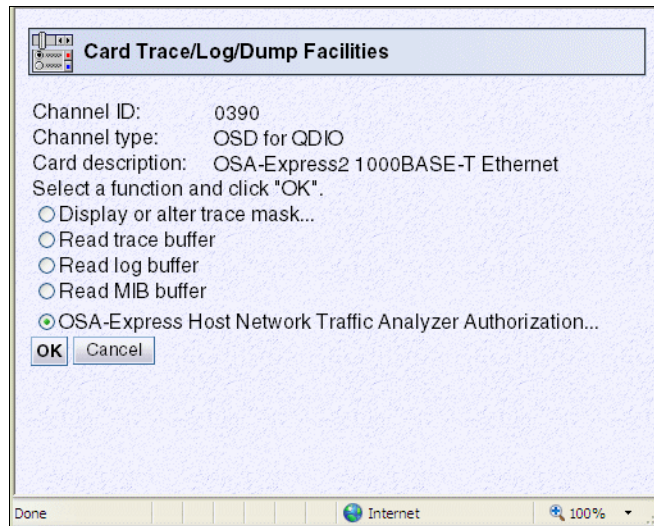


Figure 8-15 Card Trace/Log/Dump Facilities

15. Select **OSA-Express Host Network Traffic Analyzer Authorization**, then click **OK**; see Figure 8-16.

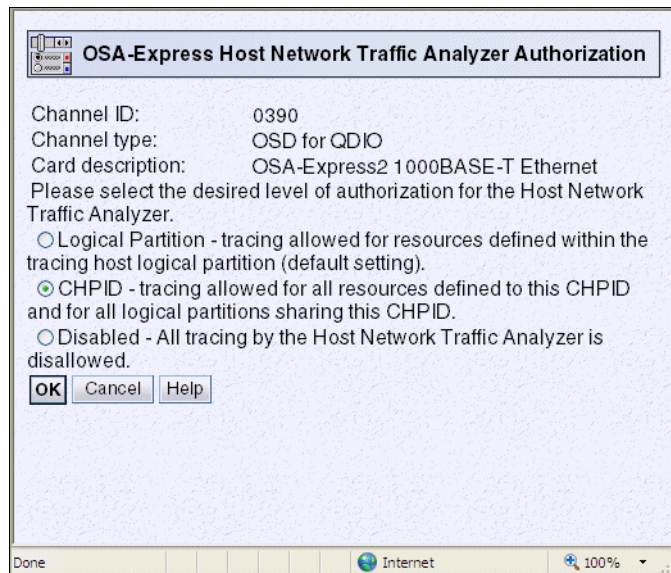


Figure 8-16 NTA Authorization

16. If your CHPID is shared between several LPARs, we suggest you take the second option shown in Figure 8-16, then click **OK**. Figure 8-17 shows the results.

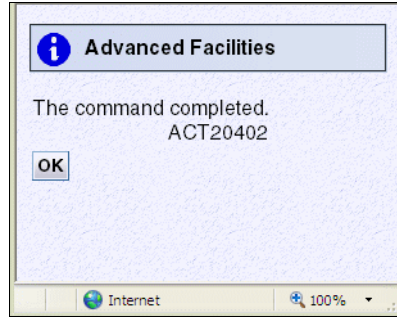


Figure 8-17 Command completed

17. To verify if the command has been set as required:

- Log off the SYSPROG user ID.
- Log on to the SE on the HMC through SOO; see Figure 8-8 on page 336.

Important: For checking the authorization of OENTA support, the user ID must have the Access Administrator Tasks role assigned.

- Select **Check current OSA-Express Network Traffic Analyzer Authorization**, as shown in Figure 8-10 on page 337.
- Click **OK**; see Figure 8-18.

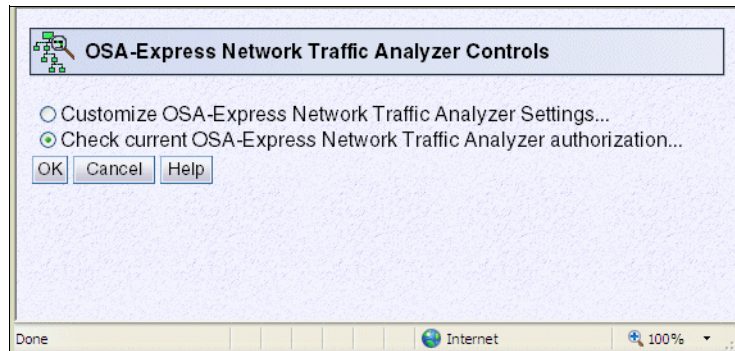


Figure 8-18 OSA-Express NTA controls

18. Figure 8-19 shows that PCHID 0390 is allowed to be traced.

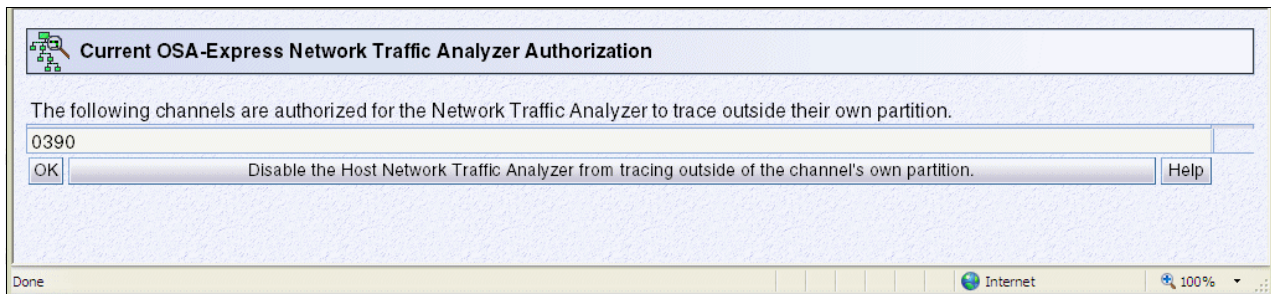


Figure 8-19 PCHID NTA Authorization

8.5.5 Defining a resource profile in RACF

See Example 8-27 for the RACF commands needed to allow users to issue the VARY TCPIP command.

Example 8-27 RACF commands

```
RDEFINE OPERCMDS MVS.VARY.TCPIP.OSAENTA UACC(NONE) PERMIT MVS.VARY.TCPIP.OSAENTA  
ACCESS(CONTROL) CLASS(OPERCMDS) ID(CS03) SETR GENERIC(OPERCMDS) REFRESH SETR  
RACLIST(OPERCMDS) REFRESH
```

8.5.6 Allocating a VSAM linear data set

Example 8-28 shows how to create the VSAM linear data set. This VSAM linear data set is optional; however, we recommend its use.

Example 8-28 Allocate VSAM linear data set

```
//DEFINE EXEC PGM=IDCAMS
//SYSPRINT DD SYSOUT=*
//SYSIN DD *
DELETE +
(CS03.CTRACE.LINEAR) +
CLUSTER
DEFINE CLUSTER( +
NAME(CS03.CTRACE.LINEAR) +
LINEAR +
MEGABYTES(10) +
VOLUME(CPDLB0) +
CONTROLINTERVALSIZE(32768) +
) +
DATA( +
NAME(CS03.CTRACE.DATA) +
)
LISTCAT ENT(USER41.CTRACE.LINEAR)
ALL
```

8.5.7 Starting the OSAENTA trace

The OSAENTA statement dynamically defines a QDIO interface to the OSA-Express being traced, called an OSAENTA interface. That interface is used exclusively for capturing OSA-Express Network Traffic Analyzer traces.

The OSAENTA statement enables an installation to trace data from other hosts connected to OSA-Express.

Important: The trace data collected should be considered confidential and TCP/IP system dumps and external trace files containing this trace data should be protected.

To see the complete syntax of the OSAENTA command, refer to *z/OS Communications Server: IP Configuration Reference*, SC31-8776.

Components involved for z/OS CTRACE

The CTRACE component for collecting NTA trace data is called SYSTCPOT. The member in SYS1.PARMLIB is named CTINTA00. This member is used to define the size of the buffer space in the TCPIPDS1 data space reserved for OSAENTA CTRACE. The size can range from 1 M to 624 M, with a default of 64 M.

Note: Update CTINTA00 to set the CTRACE buffer size. Keep in mind that this will use up auxiliary page space storage.

Using the OSAENTA command

An internal interface is created when PORTNAME is defined on the OSAENTA statement. The dynamically-defined interface name is EZANTA concatenated with the port name. These EZANTA interfaces are displayed at the end of the NETSTAT DEV output.

When the ON keyword of the OSAENTA parameter is used, VTAM allocates the next available TRLE data path associated with the port. This data path is used only for inbound trace data.

When the OFF keyword of the OSAENTA parameter is used (or the trace limits of the TIME, DATA, or FRAMES keyword are reached), the data path is released.

Setting the OSAENTA traces

You can set the OSAENTA trace in two ways: by coding the OSAENTA statement in the profile TCP/IP, or by issuing a command in z/OS. These methods are explained in this section.

- To code the OSAENTA statement in the profile TCP/IP, see Example 8-29.

Example 8-29 TCP/IP profile

```
; set up the filters to trace for TCP packets on PORT 2323 with a source
;or destination
; IP address of 10.1.2.11 over MAC address 00096B1A7490
OSAENTA PORTNAME=OSA2080 PROT=TCP IP=10.1.2.11 PORTNUM=2323
OSAENTA PORTNAME=OSA2080 MAC=00096B1A7490
; activate the tracing (the trace will self-deactivate after 20,000 frames)
OSAENTA PORTNAME=OSA2080 ON FRAMES=20000
; deactivate the tracing
OSAENTA OFF PORTNAME=OSA2080
```

In this case, OSAENTA traces the portname OSA2080 only for traffic matching the following filters:

- Protocol = UDP
- IP address = 10.1.2.11
- Port number = 2323

There are seven filters available to define the packets to be captured:

- MAC address
- VLAN ID
- Ethernet frame type
- IP address (or range)
- IP protocol
- Device ID
- TCP/UDP

Note: Use filters to limit the trace records to prevent over consumption of the OSA CPU resources, the LPAR CPU resources, the TCPIPDS1 trace data space, memory, auxiliary page space and the IO subsystem writing trace data to disk.

- To issue the following command in z/OS:
V TCPIP,TCPIPA,OSAENTA,ON,PORTNAME=OSA2080,IP=10.1.2.11,PORTNUM=2323

The messages you receive in response to this command are shown in Figure 8-20.

```
RESPONSE=SC30      EZZ0060I PROCESSING COMMAND: VARY TCPIP,TCPIPA,OSAENTA,ON,  
RESPONSE=PORTNAME=OSA2080,IP=10.1.2.11,PORTNUM=2323  
RESPONSE=SC30      EZZ0053I COMMAND VARY OSAENTA COMPLETED SUCCESSFULLY
```

Figure 8-20 OSAENTA results

Important: If you receive ERROR CODE 0003 it means that an attempt was made to enable OSA-Express Network Traffic Analyzer (OSAENTA) tracing for a specified OSA, but the current authorization level does not permit it.

Refer to 8.5.4, “Customizing OSA-Express Network Traffic Analyzer” on page 336 for directions about how to change the authorization to allow OSAENTA to be used on this specified OSA. Also read *Support Element Operations Guide*, SC28-6860, for complete information about this topic.

The command NETSTAT DEVLINKS has been enhanced to show the OSAENTA definition; see Example 8-30.

Example 8-30 NETSTAT DEVLINKS command output

```
OSA-EXPRESS NETWORK TRAFFIC ANALYZER INFORMATION:  
OSA PORTNAME: OSA2080      OSA DEVSTATUS: READY  
OSA INTFNAME: EZANTAOSA2080  OSA INTFSTATUS: READY  
OSA SPEED: 1000      OSA AUTHORIZATION: CHPID  
OSAENTA CUMULATIVE TRACE STATISTICS:  
  DATAMEGS: 0      FRAMES: 0  
  DATABYTES: 0      FRAMESDISCARDED: 0  
  FRAMESLOST: 0  
OSAENTA ACTIVE TRACE STATISTICS:  
  DATAMEGS: 0      FRAMES: 0  
  DATABYTES: 0      FRAMESDISCARDED: 0  
  FRAMESLOST: 0      TIMEACTIVE: 0  
OSAENTA TRACE SETTINGS:      STATUS: ON  
  DATAMEGSLIMIT: 1024      FRAMESLIMIT: 2147483647  
  ABBREV: 224      TIMELIMIT: 10080  
  DISCARD: EXCEPTION  
OSAENTA TRACE FILTERS:      NOFILTER: NONE  
  DEVICEID: *  
  MAC: *  
  VLANID: *  
  ETHTYPE: *  
  IPADDR: 10.1.2.11/32  
  PROTOCOL: * TCP  
  PORTNUM: * 2323
```

The NETSTAT display for devices shows the Network Traffic Analyzer interfaces. The interface name has prefixed the OSA port name with EZANTA.

To display a specific NTA interface, use the INTFName=EZANTAosaportname keyword.

Traces are placed in an internal buffer, which can then be written out using a CTRACE external writer. The MVS TRACE command must also be issued for component SYSTCPOT to activate the OSAENTA trace.

Attention: If you receive ERROR CODE 0005 it means that an attempt was made to enable OSA-Express Network Traffic Analyzer tracing for a specified OSA that already has either OSAENTA or OSA LAN Analyzer tracing enabled elsewhere on the system for this OSA.

Only one instance of active tracing (either OSAENTA or LAN Analyzer) for a specified OSA is permitted on the system at any one time.

When the trace is started from OSA/SF, you can see that another device has been allocated for trace; see Example 8-31.

Example 8-31 OAT with OSAENTA started

Image 2.3 (A16)	CULA 0						
00(2080)* MPC	N/A	OSA2080P	(QDIO control)	SIU	ALL		
02(2082) MPC	00 No4 No6	OSA2080P	(QDIO data)	SIU	ALL		
	VLAN 10	(IPv4)					
		Group Address	Multicast Address				
		01005E000001	224.000.000.001				
		VMAC	IP address				
HOME	00096B1A7490	010.001.000.010					
HOME	00096B1A7490	010.001.001.010					
HOME	00096B1A7490	010.001.002.010					
HOME	00096B1A7490	010.001.002.011					
REG	00096B1A7490	010.001.002.012					
REG	00096B1A7490	010.001.003.011					
REG	00096B1A7490	010.001.003.012					
REG	00096B1A7490	010.001.004.011					
REG	00096B1A7490	010.001.005.011					
REG	00096B1A7490	010.001.006.011					
REG	00096B1A7490	010.001.007.011					
REG	00096B1A7490	010.001.008.010					
REG	00096B1A7490	010.001.008.020					
03(2083) MPC	00 No4 No6	OSA2080P	(QDIO data)	SIU	ALL		

You can also use the D NET,TRL,TRLE=OSA2080 command, as shown in Example 8-32.

Example 8-32 Output Display TRLE

```

TRL MAJOR NODE = OSA2080
STATUS= ACTIV, DESIRED STATE= ACTIV
TYPE = LEASED           , CONTROL = MPC , HPDT = YES
MPCLEVEL = QDIO         MPCUSAGE = SHARE
PORTNAME = OSA2080     LINKNUM = 0   OSA CODE LEVEL = 087A
HEADER SIZE = 4096 DATA SIZE = 0 STORAGE = ***NA***
WRITE DEV = 2081 STATUS = ACTIVE      STATE = ONLINE
HEADER SIZE = 4092 DATA SIZE = 0 STORAGE = ***NA***
READ  DEV = 2080 STATUS = ACTIVE      STATE = ONLINE
DATA  DEV = 2082 STATUS = ACTIVE      STATE = N/A
I/O TRACE = OFF TRACE LENGTH = *NA*
ULPID = TCPIPA
IQDIO ROUTING DISABLED

```

```

READ STORAGE = 4.0M(64 SBALS)
PRIORITY1: UNCONGESTED PRIORITY2: UNCONGESTED
PRIORITY3: UNCONGESTED PRIORITY4: UNCONGESTED
DEVICEID PARAMETER FOR OSAENTA TRACE COMMAND = 02-03-00-02
UNITS OF WORK FOR NCB AT ADDRESS X'0F4E7010'
P1 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
P2 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
P3 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
P4 CURRENT = 0 AVERAGE = 2 MAXIMUM = 3
TRACE DEV = 2083 STATUS = ACTIVE STATE = N/A

```

Starting the CTRACE

To print out the internal trace data, start the CTRACE using these steps:

1. Start the external writer (CTRACE writer).

```
TRACE CT,WTRSTART=CTWTR
```

2. Start the CTRACE and connect to the external writer.

```
TRACE CT,ON,COMP=SYSTCPOT,SUB=(TCPIPA)
R xx,WTR=CTWTR,END
```

3. Display the active component trace options with this command:

```
DISPLAY TRACE,COMP=SYSTCPOT,SUB=(TCPIPA)
```

Example 8-33 shows the output of this command.

Example 8-33 Display Trace output

```

RESPONSE=SC30
IEE843I 16.45.15 TRACE DISPLAY 165
      SYSTEM STATUS INFORMATION
ST=(ON,0256K,00512K) AS=ON BR=OFF EX=ON MO=OFF MT=(ON,024K)
TRACENAME
=====
SYSTCPOT
                                MODE BUFFER HEAD SUBS
                                =====
                                OFF          HEAD    2
      NO HEAD OPTIONS
SUBTRACE                        MODE BUFFER HEAD SUBS
-----
TCPIPA                          ON    0128M
ASIDS                          *NONE*
JOBNAMES                       *NONE*
OPTIONS                       MINIMUM
WRITER                        CTWTR

```

To display information about the status of the component trace for all active procedures, issue the following command:

```
DISPLAY TRACE,COMP=SYSTCPOT,SUBLEVEL,N=8
```

Example 8-34 displays the output.

Example 8-34 Status of Component Trace

```

IEE843I 10.35.04 TRACE DISPLAY 821
      SYSTEM STATUS INFORMATION
ST=(ON,0256K,00512K) AS=ON BR=OFF EX=ON MO=OFF MT=(ON,024K)
TRACENAME
=====
SYSTCPOT

                                MODE BUFFER HEAD SUBS
                                =====
                                OFF          HEAD    2

      NO HEAD OPTIONS
SUBTRACE                                MODE BUFFER HEAD SUBS
-----
TCPIPA                                ON    0128M
ASIDS                                *NONE*
JOBNAMES                            *NONE*
OPTIONS                            MINIMUM
WRITER                              CTWTR
-----
TCPIP                                MIN    0016M
ASIDS                                *NONE*
JOBNAMES                            *NONE*
OPTIONS                            MINIMUM
WRITER                              *NONE*

```

4. Reproduce the problem.

5. Disconnect the external writer.

```

TRACE CT,ON,COMP=SYSTCPOT,SUB=(TCPIPA)
R xx,WTR=DISCONNECT,END

```

6. Stop the component trace.

```

TRACE CT,OFF,COMP=SYSTCPOT,SUB=(TCPIPA)

```

7. Stop the external writer.

```

TRACE CT,WTRSTOP=CTWTR

```

Analyzing the trace

You can format the CTRACE using two methods, which we describe in this section.

Use IPCS to format CTRACE

You can format component trace records using IPCS panels or a combination of IPCS panels and the CTRACE command, either from a dump or from external writer files.

From the IPCS PRIMARY OPTION MENU, select: **0 DEFAULTS - Specify default dump and options**; see Example 8-35 on page 349 for details.

Example 8-35 IPCS default value

```
----- IPCS Default Values -----
Command ==>

You may change any of the defaults listed below. The defaults shown before
any changes are LOCAL. Change scope to GLOBAL to display global defaults.

Scope ==> LOCAL    (LOCAL, GLOBAL, or BOTH)

If you change the Source default, IPCS will display the current default
Address Space for the new source and will ignore any data entered in
the Address Space field.

Source ==> DSNAME('SYS1.SC30.CTRACE')
Address Space ==>
Message Routing ==> NOPRINT TERMINAL
Message Control ==> CONFIRM VERIFY FLAG(WARNING)
Display Content ==> NOMACHINE REMARK REQUEST NOSTORAGE SYMBOL
```

Modify the DSNAME and OPTIONS to match your environment, then select the following options:

- ▶ 2 ANALYSIS - Analyze dump contents
- ▶ 7 TRACES - Trace formatting
- ▶ 1 CTRACE - Component trace
- ▶ D DISPLAY - Specify parameters to display CTRACE entries

Fill in the parameters necessary to format the OSAENTA trace; see Example 8-36.

Example 8-36 CTRACE parameters

```
----- CTRACE DISPLAY PARAMETERS -----
COMMAND ==>

System      ==>          (System name or blank)
Component   ==> SYSTCPOT (Component name (required))
Subnames    ==> TCPIPA

GMT/LOCAL   ==> G          (G or L, GMT is default)
Start time  ==>          (mm/dd/yy,hh:mm:ss.dddddd or
Stop time   ==>          mm/dd/yy,hh.mm.ss.dddddd)
Limit       ==> 0          Exception ==>
Report type ==> SHORT      (SHort, SUMmary, Full, Tally)
User exit   ==>          (Exit program name)
Override source ==>
Options      ==>

To enter/verify required values, type any character
Entry IDs ==>  Jobnames ==>  ASIDs ==>  OPTIONS ==>  SUBS ==>

CTRACE COMP(SYSTCPOT) SUB((TCPIPA)) SHORT
```

ENTER = update CTRACE definition. END/PF3 = return to previous panel.
S = start CTRACE. R = reset all fields.

On the command line, enter the S command. Example 8-37 shows the trace formatted by IPCS.

Example 8-37 TRACE format

COMPONENT TRACE SHORT FORMAT					
SYSNAME(SC30)					
COMP(SYSTCPOT)SUBNAME((TCPIPA))					
z/OS TCP/IP Packet Trace Formatter, (C) IBM 2000-2006, 2007.052					
DSNAME('SYS1.SC30.CTRACE')					
**** 2007/09/11					
RcdNr	Sysname	Mnemonic	Entry Id	Time Stamp	Description

365	SC30	OSAENTA	00000007	15:01:23.356987	OSA-Express NTA
To Interface		:	EZANTAOSA2080		Full=86
Tod Clock		:	2010/09/24 14:20:25.931533		
Sequence #		:	0	Flags: Pkt Out Nta Vlan Lpar L3	
Source		:	10.1.2.11		
Destination		:	224.0.0.5		
Source Port		:	0	Dest Port: 0	Asid: 0000 TCB: 00000000
Frame: Device ID		:	02030002	Sequence Nr: 372	Discard: 0 (OK)
Ethernet II		:	8100	IEEE 802.1 Vlan	Len: 0x0044 (68
Destination Mac		:	01005E-000005	()	
Source Mac		:	00096B-1A7490	(IBM)	
Vlan_id		:	10	Priority: 0	Type: 0800 (Int
IpHeader: Version		:	4	Header Length: 20	
Tos		:	00	QOS: Routine Normal Service	
Packet Length		:	68	ID Number: 0AFD	
Fragment		:		Offset: 0	
TTL		:	1	Protocol: OSPFIGP	Checksum: C253
Source		:	10.1.2.11		
Destination		:	224.0.0.5		

366	SC30	OSAENTA	00000007	15:01:33.360143	OSA-Express NTA
To Interface		:	EZANTAOSA2080		Full=86
Tod Clock		:	2010/09/24 14:20:35.933003		
Sequence #		:	0	Flags: Pkt Out Nta Vlan Lpar L3	
Source		:	10.1.2.11		
Destination		:	224.0.0.5		
Source Port		:	0	Dest Port: 0	Asid: 0000 TCB: 00000000
Frame: Device ID		:	02030002	Sequence Nr: 373	Discard: 0 (OK)
Ethernet II		:	8100	IEEE 802.1 Vlan	Len: 0x0044 (68
Destination Mac		:	01005E-000005	()	
Source Mac		:	00096B-1A7490	(IBM)	
Vlan_id		:	10	Priority: 0	Type: 0800 (Int
IpHeader: Version		:	4	Header Length: 20	
Tos		:	00	QOS: Routine Normal Service	
Packet Length		:	68	ID Number: 0B07	
Fragment		:		Offset: 0	
TTL		:	1	Protocol: OSPFIGP	Checksum: C249
Source		:	10.1.2.11		
Destination		:	224.0.0.5		

Use a batch job to format CTRACE

We used a batch job to generate the TRACE file, as shown in Example 8-38.

Example 8-38 CTRACE batch job format

```
//PKT2SNIF JOB (999,P0K),'CS03',NOTIFY=&SYSUID,
//  CLASS=A,MSGCLASS=T,TIME=1439,
//  REGION=OM,MSGLEVEL=(1,1)
//  SET INDUMP='SYS1.SC30.CTRACE'
//IPCSBTCH EXEC PGM=IKJEFT01,DYNAMNBR=30
//IPCSDDIR DD DISP=SHR,DSN=SYS1.DDIR
//IPCSDUMP DD *
//SYSTSPRT DD SYSOUT=*
//SYSPRINT DD SYSOUT=*
//INDMP DD DISP=SHR,DSN=&INDUMP.
//IPCSPRNT DD DSN=ENTA.CTRACE.SHORT,UNIT=SYSDA,
//  DISP=(NEW,CATLG),LRECL=133,SPACE=(CYL,(10,1)),RECFM=VBS,DSORG=PS
//IPCSTOC DD SYSOUT=*
//SYSUDUMP DD SYSOUT=*
//SYSTSIN DD *
PROFILE MSGID
IPCS NOPARM
SETD PRINT NOTERM LENGTH(160000) NOCONFIRM FILE(INDMP)
DROPD
CTRACE COMP(SYSTCPOT) SUB((TCPIPA)) SHORT
END
```

We received the output shown in Example 8-39 from the batch job.

Example 8-39 Output from the batch job

```
COMPONENT TRACE SHORT FORMAT
SYSNAME(SC30)
COMP(SYSTCPOT)SUBNAME((TCPIPA))
z/OS TCP/IP Packet Trace Formatter, (C) IBM 2000-2006, 2007.052
FILE(INDMP)
**** 2007/09/11
RcdNr Sysname Mnemonic Entry Id Time Stamp Description
-----
365 SC30 OSAENTA 00000007 15:01:23.356987 OSA-Express NTA
To Interface : EZANTAOSA2080 Full=86
Tod Clock : 2010/09/24 14:20:25.931533
Sequence # : 0 Flags: Pkt Out Nta Vlan Lpar L3
Source : 10.1.2.11
Destination : 224.0.0.5
Source Port : 0 Dest Port: 0 Asid: 0000 TCB: 00000000
Frame: Device ID : 02030002 Sequence Nr: 372 Discard: 0 (OK)
EtherNet II : 8100 IEEE 802.1 Vlan Len: 0x0044 (68)
Destination Mac : 01005E-000005 ( )
Source Mac : 00096B-1A7490 (IBM)
Vlan_id : 10 Priority: 0 Type: 0800 (Inter
IpHeader: Version : 4 Header Length: 20
Tos : 00 QOS: Routine Normal Service
Packet Length : 68 ID Number: 0AFD
Fragment : Offset: 0
```

```

TTL           : 1                      Protocol: OSPFIGP      CheckSum: C253 FF
Source        : 10.1.2.11
Destination   : 224.0.0.5
-----
366 SC30      OSAENTA 00000007 15:01:33.360143 OSA-Express NTA
To Interface   : EZANTA0SA2080          Full=86
Tod Clock      : 2010/09/24 14:20:35.933003
Sequence #     : 0                      Flags: Pkt Out Nta Vlan Lpar L3
Source         : 10.1.2.11
Destination    : 224.0.0.5
Source Port    : 0                      Dest Port: 0      Asid: 0000 TCB: 00000000
Frame: Device ID : 02030002          Sequence Nr: 373    Discard: 0 (OK)
Ethernet II     : 8100                IEEE 802.1 Vlan    Len: 0x0044 (68)
Destination Mac : 01005E-000005      ()
Source Mac      : 00096B-1A7490      (IBM)
Vlan_id         : 10                  Priority: 0        Type: 0800 (Inter
IpHeader: Version : 4                  Header Length: 20
Tos             : 00                  QOS: Routine Normal Service
Packet Length   : 68                  ID Number: 0B07
Fragment        :                      Offset: 0
TTL             : 1                      Protocol: OSPFIGP      CheckSum: C249 FF
Source          : 10.1.2.11
Destination     : 224.0.0.5
-----

```

8.5.8 Operator command to query and display OSA information

Communications Server provides a new DISPLAY TCPIP,OSAINFO command that you can use to retrieve information about an interface from an OSA-Express feature that is in QDIO mode. This command is an alternative to using OSA/SF, which lacks information about many of the latest enhancements to the OSA-Express feature and to z/OS Communications Server.

Figure 8-21 on page 355, shows the scope of the **Display OSAINFO,INTFN**. In our example it is Interface3.

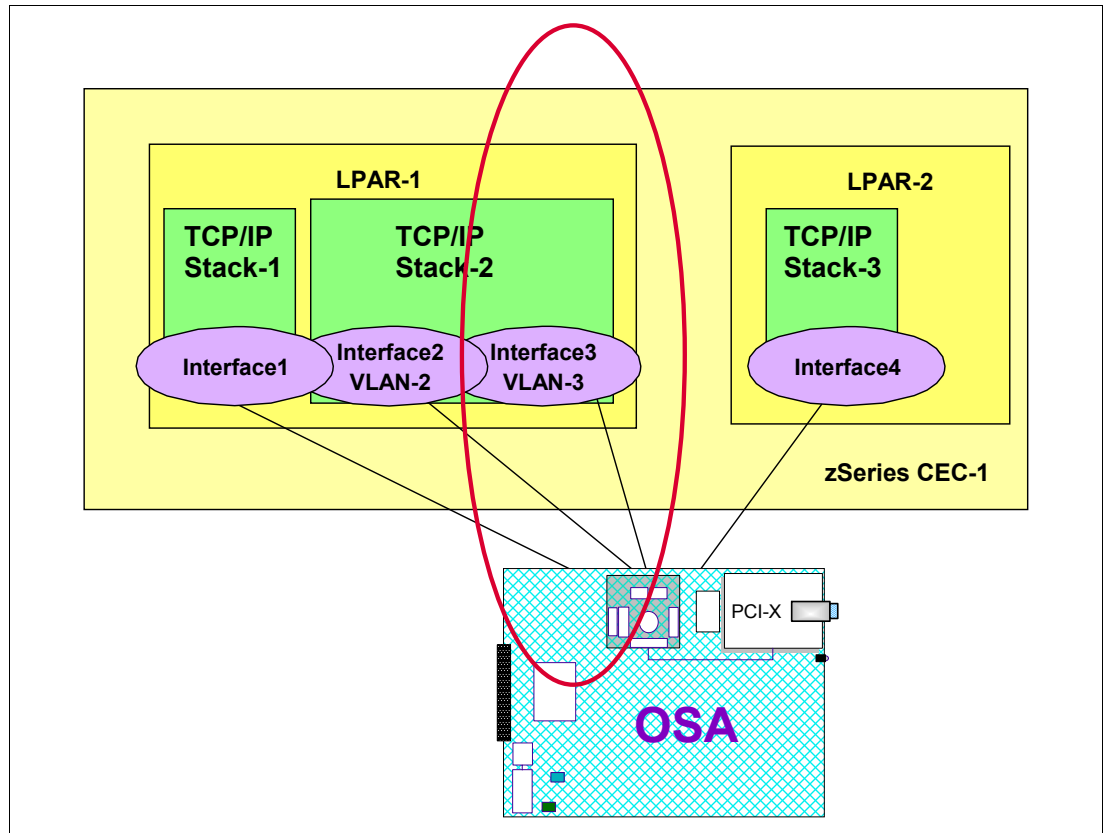


Figure 8-21 Scope of OSAINFO command

Example 8-40 shows the output of the following command:

D TCPIP,TCPIPA,OSAINFO,INTFNAME=OSA2080I,MAX=200

Example 8-40 Output of OSAINFO command

```

EZZ0053I COMMAND DISPLAY TCPIP,,OSAINFO COMPLETED SUCCESSFULLY
EZD0031I TCP/IP CS V1R12 TCPIP Name: TCPIPA 14:04:54 684
Display OSAINFO results for IntfName: OSA2080I
PortName: OSA2080 PortNum: 00 Datapath: 2082 RealAddr: 0002
PCHID: 0530 CHPID: 02 CHPID Type: OSD OSA code level: 000C
Gen: OSA-E3 Active speed/mode: 100 mb/sec full duplex
Media: Copper Jumbo frames: No Isolate: No
PhysicalMACAddr: 00145E776872 LocallyCfgMACAddr: 000000000000
Queues defined Out: 4 In: 1 Ancillary queues in use: 0
Connection Mode: Layer 3 IPv4: Yes IPv6: No
SAPSup: 000FF603 SAPEna: 00082603
IPv4 attributes:
  VLAN ID: 10 VMAC Active: Yes
  VMAC Addr: 020005776873 VMAC Origin: OSA VMAC Router: Local
  AsstParmsEna: 00300C57 OutCkSumEna: 0000001A InCkSumEna: 0000001A
Registered Addresses:
  IPv4 Unicast Addresses:
    ARP: No Addr: 10.1.1.10
    ARP: Yes Addr: 10.1.2.10
    ARP: Yes Addr: 10.1.2.11
    ARP: No Addr: 10.1.2.12
  
```

```

ARP: No   Addr: 10.1.2.14
ARP: No   Addr: 10.1.3.11
ARP: No   Addr: 10.1.3.12
ARP: No   Addr: 10.1.4.11
ARP: No   Addr: 10.1.5.11
ARP: No   Addr: 10.1.6.11
ARP: No   Addr: 10.1.7.11
ARP: No   Addr: 10.1.8.23
Total number of IPv4 addresses:      12
IPv4 Multicast Addresses:
MAC: 01005E000001  Addr: 224.0.0.1
MAC: 01005E000005  Addr: 224.0.0.5
Total number of IPv4 addresses:      2

```

If you have multiple interfaces in the same stack, you have to issue this command for each interface.

8.5.9 OSM and OSX information

OSA/SF cannot manage OSM and OSX CHPIDs. The **D TCPIP,,OSAINFO** command can be used to get information about those CHPID types.

Displaying the CHPID type

We use the z/OS command **D M=CHP** to see how the system shows these CHPIDs. Example 8-41 shows the output.

In the IOCDs, CHPIDs 0A and 0B are OSM; 18 and 19 are OSX.

Example 8-41 Output of z/OS command “D M=CHP”

```

CHANNEL PATH STATUS
      0 1 2 3 4 5 6 7 8 9 A B C D E F
0  + + + + + + + + + + + + + + + +
1  + + + + . . . . + + - - . . . .

CHANNEL PATH TYPE STATUS
      0 1 2 3 4 5 6 7 8 9 A B C D E F
0  11 14 11 11 11 11 11 11 11 11 31 31 11 14 11 11
1  11 11 11 11 11 00 00 00 00 30 30 11 11 00 00 00

30 OSA ZBX DATA          OSX
31 OSA ZBX MANAGEMENT    OSM

```

In our case the CHPIDs are online and operating.

Example 8-42 shows the output of **D U,,2340,15** (OSM CHPID).

Example 8-42 Output of Display Unit command for OSM CHPID

```

D U,,2340,15
IEE457I 16.22.42
UNIT TYPE STATUS
2340 OSAM 0-RAL
2341 OSAM 0-RAL
2342 OSAM 0-RAL

```

```
2343 OSAM 0-RAL
2344 OSAM 0-RAL
2345 OSAM 0-RAL
2346 OSAM 0-RAL
2347 OSAM 0-RAL
2348 OSAM 0-RAL
2349 OSAM 0-RAL
234A OSAM 0-RAL
234B OSAM 0-RAL
234C OSAM 0-RAL
234D OSAM 0-RAL
234E OSAM 0-RAL
```

Example 8-43 shows the output of “D U,,,2300,15” (OSX CHPID).

Example 8-43 Output of Display Unit command for OSX CHPID

```
D U,,,2300,15
IEE457I 16.23.10
UNIT TYPE STATUS
2300 OSAX 0
2301 OSAX 0
2302 OSAX 0
2303 OSAX 0
2304 OSAX 0
2305 OSAX 0
2306 OSAX 0
2307 OSAX 0
2308 OSAX 0
2309 OSAX 0
230A OSAX 0
230B OSAX 0
230C OSAX 0
230D OSAX 0
230E OSAX 0
```

The devices are online and allocated.

Note: You can use the OSA-Express Network Traffic Analyzer (OSAENTA) trace facility to debug any problem with OSX and OSM CHPID types.

When you use NTA you must have a “Datapath” available for each NTA you start.

For an OSM CHPID type you have nine Datapaths, while an OSX CHPID type supports 17 Datapaths. For more information refer to Chapter 9, “z/OS in an ensemble” on page 371.

8.6 Additional tools for diagnosing CS for z/OS IP problems

IBM and other vendors have developed tools to assist in diagnosing problems in the network from the perspective of z/OS. The tools often run as GUIs on a workstation, but retrieve their problem diagnosis information using data from SNMP, SMF, and from MVS control blocks. Some of these tools also interface with the Network Management Interface API, provided by IBM.

8.6.1 Network Management Interface API

Figure 8-22 depicts a high level view of the Network Management Interface (NMI) and its interfaces to network management products.

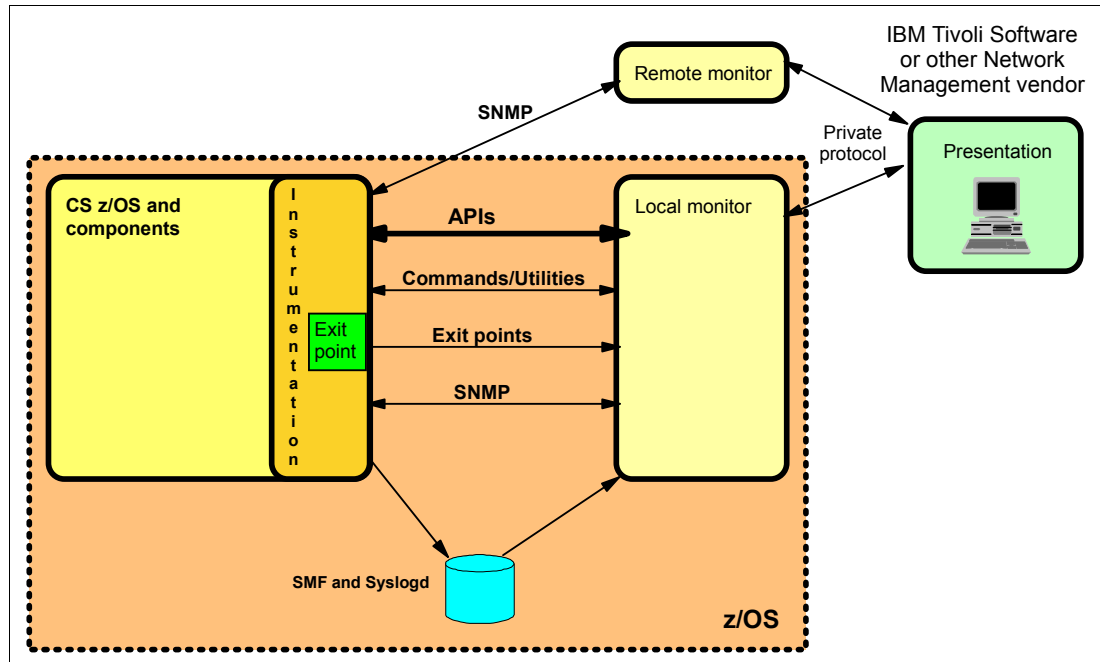


Figure 8-22 Network Management Interface Architecture

The NMI API can interface with Tivoli® OMEGAMON® XE for Mainframe Networks (or other products) to provide the following types of functions:

- ▶ Trace assistance:
 - Real-time tracing and formatting for packet and data traces (including OSA trace)
- ▶ Information gathering:
 - TCP connection initiation and termination notifications
 - API for real-time access to TN3270 server and FTP event data and to IPSec
 - APIs to poll information about currently active connections
 - TCP listeners (server processes)
 - TCP connections (detailed information about individual connections and UDP endpoints)
 - Communications Server storage usage
 - API to receive and poll for Enterprise Extender management data
 - Information and statistics for IP filtering and IPSec security associations on the local TCP/IP stacks.
 - Information and statistics for IP filtering and IPSec security associations on remote Network Security Services (NSS) clients when using the NSS server.
- ▶ Control activities
- ▶ Control the activation and inactivation of IPSec tunnels

- ▶ Loading policies for IP filtering and IPsec security associations on the local TCP/IP stacks
 - Loading policies for IP filtering and IPsec security associations on remote Network Security Services (NSS) clients when using the NSS server
 - Drop one or multiple TCP connections or UDP endpoints

8.6.2 Systems Management Facilities accounting records

Another technique that is often used to verify the state of the z/OS Communications Server - TCP/IP component in a stack or even in a Sysplex environment is to list and analyze the Systems Management Facilities (SMF) records.

In general, SMF records are created for deferred processing and analysis. SMF recording is generally not used for real-time monitoring purposes. In a TCP/IP environment, real-time monitoring is implemented using the NMI API and SNMP protocol, but on z/OS a lot of the information that is written in SMF records is useful from a real-time monitoring perspective.

The objective of the TCP/IP product is to define and generate the lowest level of detail that is needed by all disciplines. A customer has to use other products such as RMF™, Performance Reporter for z/OS (PR), MVS Information Control System (MICS), or SAS-based tools. In many cases, there are customer-written programs to generate the reports to collect and analyze the SMF Records created by TCP/IP.

Note: TCP/IP- produced SMF records should not be viewed in isolation. Other components in MVS produce SMF records for the same purposes as those produced by TCP/IP. An installation is likely to combine information from a series of subsystems when performing detailed performance or capacity planning.

The contents of SMF records can be used to generate reports in customized formats that help customers to perform tasks such as:

- ▶ Performance management:

Customized reports can be generated to verify if the defined service levels are met and, if not, to identify possible causes. These reports are usually a set of time intervals, ranging from weeks through days matching the SMF interval. Some examples of potential reports related to performance management are:

- TCP connection elapsed time per server port number per time of day (potentially broken down by source IP address or netmask)
- Number of TCP connections per server port number per time of day (potentially broken down by source IP address or netmask)
- Number of inbound/outbound bytes transferred in TCP connections per time of day (potentially broken down in various ways: by destination or source port, by destination IP address, netmask, or in total)
- Events related to dynamic VIPA environment such as:
 - Status changes
 - DVIPA removed or added
 - Changes on the target server (stop/start)

- ▶ Capacity planning:

Capacity planning can be done using the SMF records to generate reports showing trends for resource utilization of central processing power, memory, channel-based I/O subsystem, network attachments, and network bandwidth, over a period of time. These trends can help with planned launches of new applications or use of existing applications,

in order to predict capacity needs in the future. Some examples of potential reports related to capacity planning are:

- Total number of TCP connections per reserved server port number per day including analysis of average and variations around average during daily peak periods
- Total number of UDP inbound/outbound UDP datagrams per reserved server port number per day including average and variations around average during daily peak periods
- Number of bytes or packets transferred inbound and outbound per interface (LINK) per time of day (potentially broken down into unicasts, broadcasts, and multicasts)

► Auditing

Auditing involves tasks that are related to identifying and proving that individual events have taken place. Some examples of potential reports related to auditing are:

- Detailed information about specific TCP connections or UDP sockets, IP addresses, server/client identification, duration, number of bytes, and so on
- Details about activity that involves a specific client or server.
- Details about a given application session based on server-specific SMF recording, such as individual Telnet sessions or FTP sessions
- An SMF 119 record for recording TCP/IP configuration updates
- Changes on the dynamic VIPA environment

► Accounting

Accounting involves tasks that are related to calculating how much each individual user or organizational unit should be charged for use of the shared central IS resources. Input for these reports can be based on CPU cycle use, data quantities, bandwidth usage, and memory use. For TCP/IP, additional metrics can be defined, such as type of service (FTP, Web server, TN3270 and so on), and TCP connection-related information (number of connections, duration, byte transfer counts, and so on).

Some examples of potential reports related to accounting are:

- Aggregated number of connections to a given server from a given source in terms of a specific client IP address, or netmask
- Accumulated connect time to a given server from a given source in terms of a specific client IP address, or netmask
- Number of bytes transferred to or from a given source in terms of a specific client IP address, or netmask
- Amount of data protected by specific manual or dynamic tunnels
- For IKED: Information about IKE tunnels
- For TN3270: Number of sessions and session type (TN3270/TN3270E/LINEMODE)

Depending on the configuration for the z/OS Communications Server - TCP/IP component, SMF records can be cut at multiple levels in the TCP/IP protocol stack, and the type of information that can be included depends on where the SMF record is created:

► At the IP and interface layer

Information about ICMP activity, IP packet fragmentation and reassembly activity, IP checksum errors, IGMP activity, and ARP activity. This information is important to generate reports related either to performance or capacity management.

► At the transport protocol layer

Information about IP addresses, port numbers, and host names. It has also information about TCP connections, such as byte counts, connection times, reliability metrics, and performance metrics. For UDP-related workload, each UDP datagram is a separate entity; the only way to aggregate information for UDP is on a UDP socket level, where SMF records could be created every time a UDP socket is closed.

► At the application layer

Currently, application-layer SMF recording is done for the TN3270 Telnet server (Telnet), the FTP server, and the IKE daemon, but not for any other servers.

SMF record types used by Communications Server for z/OS IP

Communications Server for z/OS IP generates SMF records using two types of records: SMF record type 118 and SMF record type 119. TCP/IP SMF records written using record type 118 are created to reflect information related to the events shown in Table 8-2.

Table 8-2 Events logged using SMF record type 118

Events	Subtype records
TCP API connection initiation	1
TCP API connection termination	2
FTP client requests	3
Telnet client connection initiation and termination	4
TCP/IP statistics	5
TN3270 server session initiation and termination	20-21
FTP Server related information	70-75

SMF record type 118 provides basic information and does not have information related to the TCP/IP stack. In a multiple stack environment it is not easy to determine which SMF records relate to which TCP/IP stack.

SMF record type 119 contain additional values that identify the TCP/IP stack, which solves the record 118 problem. It also provides other advantages such as uniformity of date and time (UTC), common record format (self-defining section and TCP/IP identification section), and support for IPv6 addresses and expanded field sizes (64 bit versus 32 bit) for some counters. The SMF record type 119 subtype records available are shown in Table 8-3:

Table 8-3 Events logged using SMF record type 119

Events	Subtype records
TCP connection initiation	1
TCP connection termination	2
FTP client transfer completion	3
TCPIP PROFILE information	4
TCP/IP, interface, and server port statistics	5-7
TCP/IP stack start/stop	8
UDP socket close	10
TN3270 server session initiation and termination	20-21

Events	Subtype records
TSO Telnet client session initiation and termination	22-23
DVIPA status changes	32
DVIPA removed	33
DVIPA target added	34
DVIPA target removed	35
DVIPA target server started	36
DVIPA target server ended	37
FTP server transfer completion	70
FTP server logon failure	72
IKE tunnel activation, refresh, deactivation, and expire	73,74
Dynamic tunnel activation, refresh, installation, and removal	75-78
Manual tunnel activation and deactivation	79,80

Customizing the SMF records data collection

Depending on the type of information needs to be gathered, You can control the collection of these records using the SMFCONFIG statements in PROFILE.TCPIP, the SMF statements in the FTP.DATA for the FTP server configuration, and the SMF119 statement in IKE daemon configuration file. For more information about configuring those statements, see:

- ▶ “SMFCONFIG” on page 428
- ▶ *IBM z/OS Communications Server TCP/IP Implementation Volume 2: Standard Applications*, SG24-7897
- ▶ *Communications Server for z/OS TCP/IP Implementation Volume 4: Security and Policy-Based Networking*, SG24-7899

8.7 MVS console support for selected TCP/IP commands

In this section, we discuss MVS console support for selected TCP/IP commands.

8.7.1 Concept

A new command, EZACMD, has been implemented, and its function allows you to run selected z/OS Communications Server UNIX commands from other command environments, such MVS console, IBM Tivoli NetView for z/OS, and TSO.

The new command is used as a common interface for the execution of specific z/OS Communications Server TCP/IP infrastructure policy-related commands (**pasearch**, **trmdstat**, **nssctl**, **ipsec**, and **ping**) from other environments beyond the z/OS CS UNIX.

When operators of z/OS or NetView wanted to query or change policy aspects of z/OS Communications Server's infrastructure, they had to log in to the z/OS CS UNIX shell. Now, commands related to policies can be used on NetView, MVS console, or TSO.

However, if you want to use this new feature, you must enable EZACMD. The following sections describe how to enable the EZACMD functions.

Note: You must configure and enable the System REXX component to use EZACMD.

8.7.2 Commands and environments supported by EZACMD

The commands shown in Table 8-4 and their respective environments are the only ones supported by EZACMD.

Table 8-4 Policy-related z/OS UNIX infrastructure commands and environments

Environment	Commands
z/OS TSO	pasearch, trmdsat, nssctl, and ipsec
z/OS console	pasearch, trmdsat, nssctl, ipsec, and ping
z/OS NetView	pasearch, trmdsat, nssctl, ipsec, and ping

8.7.3 When to use EZACMD

Whenever you need information about the infrastructure's policy and cannot log in to z/OS CS UNIX, use the commands shown in Table 8-5.

Table 8-5 The commands and their functions

Command	What it does
pasearch	Queries information from the z/OS Communications Server policy agent.
nssctl	Displays information about network security services (NSS) clients that are currently connected to the local NSS server.
trmdstat	Displays a consolidated view of log messages that have been written out by the Traffic Regulation Management daemon (TRMD).
ipsec	Displays and modifies IP security information for a local TCP/IP stack and the IKE daemon. It is also used for the network security services (NSS) IPSec client that uses the IPSec network management service, of the local NSS server.
ping	Tests the connectivity between devices and the z/OS system.

8.7.4 How to use the EZACMD command

The EZACMD is in the System REXX as a group of REXX libraries:

- ▶ SYS1.SAXREXEC contains the REXX system's library. It is used by MVS console. This is a VB,LRECL=255 library.
- ▶ tcpip.SEZAEXEC contains the z/OS Communications Server REXX's library. It is used by TSO and NetView. This is a FB,LRECL=80 library.

Table 8-6 shows the syntax for the EZACMD command.

Table 8-6 Syntax of EZACMD

Command name	Command options	MAX= *
<ul style="list-style-type: none"> ▶ This is one of the supported z/OS UNIX commands: pasearch, trmdstat, nssctl, ipsec, and ping. ▶ The command name is not case sensitive. 	Refer to the options for z/OS UNIX commands in <i>z/OS Communications Server: IP System Administrator's Commands</i> , SC31-8781, <i>z/OS Communications Server: Quick Reference</i> , SX75-0124, and <i>z/OS Communications Server: IP Configuration Guide</i> , SC31-8775. Options are case sensitive and must be entered in the required case.	<ul style="list-style-type: none"> ▶ This is the optional maximum number of output lines. ▶ The default is 100, and the maximum is 64000.

Note: Each environment has specific requirements and characteristics for using EZACMD, which are discussed in 8.7.5, “Configuring z/OS for using the EZACMD” on page 362.

8.7.5 Configuring z/OS for using the EZACMD

Here we discuss configuring z/OS to use EZACMD. Perform the following steps:

1. Configure SYS1.PARMLIB member AXRnn with a System REXX command recognition string, as shown in Figure 8-23. We use ‘REXX&SYSCclone’ and the AXRnn definition CPF(‘REXX&SYSCclone.’,SYSPLEX).

```
CPF('REXX&SYSCclone.',SYSPLEX) /* Defines REXXnn as a sysplex      */ 1
                               /* wide cpf value                  */
AXRUSER(AXRUSER)              /* ?AXREXX security=axruser results in the
                               exec running in a security environment
                               defined by the userid AXRUSER      */
REXXLIB ADD DSN(SYS1.SAXRExec) VOL(&SYSR1.)
```

Figure 8-23 A SYS1.PARMLIB(AXR00)

Where 1 REXX is the constant and SYSCclone is the variable referencing the sysplex.

2. Follow the System REXX documentation for defining JCL procedures and RACF definitions. Table 8-7 displays the required configurations steps.

Table 8-7 Steps to configure the MVS support for selected TCP/IP commands

Task	How to do it	Reference
Enable the use of EZACMD from the MVS console.	Configure and enable the System REXX component on z/OS.	Chapter 8, “AXR00 (default System REXX data set concatenation)” in <i>MVS Programming: Authorized Assembler Services Guide</i> , SA22-7608, chapter 31 “System REXX”. <i>MVS Initialization and Tuning Reference</i> , SA22-7592

Task	How to do it	Reference
Call z/OS Communications Server UNIX policy-related commands from the MVS console, TSO, or NetView environments.	Use the new EZACMD command, followed by a specific policy-related command, such as pasearch , trmdstat , nssctl , ipsec , or ping , as input.	<i>z/OS Communications Server: IP System Administrator's Commands</i> , SC31-8781.

8.7.6 Using the EZACMD command in the z/OS console

To use the EZACMD command in the z/OS console, perform the following steps:

1. Go to **ISPF** → **SDSF** → **LOG (System log)**, or directly to DA (display active users), and input the variable for the environment (in this case, we use REXX30) followed by the EZACMD command.
2. Next, enter the specific TCP/IP policy command and its options within quotes to avoid having the input translated to upper-case, as shown in Figure 8-24.

```

Display Filter View Print Options Help
-----
SDSF DA SC30      (ALL)    PAG 0 CPU/L/Z  2/  2/  0  LINE 1-10 (10)
COMMAND INPUT ==> /+                                SCROLL ==> CSR
NP

      System Command Extension

Type or complete typing a system command, then press Enter.

==> REXX30EZACMD 'ipsec -f display -r short -p tcpipa [MAX=20]'
==>

Place the cursor on a command and press Enter to retrieve it.
More:      +

=> REXX30EZACMD 'ipsec -f display -r short -p tcpipa'
=> REXX30EZACMD 'IPSEC -F DISPLAY -R DETAIL -P TCPIPA'
=> REXX30EZACMD 'PING -v 9.12.6.50'
=> D PROG,APF
=> D TCPIP,TCPIPA,N,VCRT
=> S TN3270A
=> S TCPIPA

- Wait 1 second to display responses (specify with SET DELAY)
- Do not save commands for the next SDSF session

F1=Help  F5=FullScr  F7=Backward F8=Forward  F11=ClearLst  F12=Cancel

```

Figure 8-24 EZACMD syntax

Example 8-44 shows the response to the EZACND command in Figure 8-24 on page 363.

Example 8-44 Response of Ipsec command through EZACMD via MVS console

```

REXX30EZACMD 'ipsec -f display -r short -p tcpipa MAX=10'
System REXX EZACMD: ipsec command - start - userID=CS02
System REXX EZACMD: ipsec -f display -r short -p tcpipa

```

```

CS V1R12 ipsec Stack Name: TCPIPA Wed Sep 2 16:23:25 2009
Primary: Filter          Function: Display          Format: Short

```

```
Source: Stack Profile Scope: Current TotAvail: 4
Logging: On Predecap: Off DVIPSec: No
NatKeepAlive: 20
Defensive Mode: Inactive
```

```
FilterName |FilterNameExtension
            |GroupName
System REXX EZACMD: Maximum number of output lines (10) has been
reached.
System REXX EZACMD: ipsec command - end - RC=4
```

Note: If you help while in the console, type **(pref)EZACMD ? /-? /help**, where (pref) is your current sysplex/partition.

8.7.7 Preparing the EZACMD command in z/OS TSO and z/OS NetView

To prepare the EZACMD command in z/OS TSO and z/OS NetView, follow the steps shown in Table 8-8.

Table 8-8 Steps to configure the use of EZACMD by TSO and NetView

Task	How to do it	Reference
Enable the use of EZACMD from the z/OS TSO.	<ol style="list-style-type: none"> 1. Copy EZACMD to a REXX library that is used by TSO. 2. Concatenate tcpip.SEZAEXEC to the SYSEXEC or SYSPROC DD name. 	<i>z/OS Communications Server: IP System Administrator's Commands, SC31-8781</i>
Enable the use of EZACMD from the z/OS NetView.	<p>Ensure that you:</p> <ol style="list-style-type: none"> 1. Have copied EZACMD to a REXX library that is used by NetView. 2. Concatenate tcpip.SEZAEXEC to the DSICLD DD name. 	<i>z/OS Communications Server: IP System Administrator's Commands, SC31-8781</i>

Note: To preserve the case of the entered arguments, prefix the EZACMD command with the NetView NETVASIS command:

```
netvasis ezacmd ping -v w3.ibm.com max=20
```

8.7.8 Using EZACMD command from z/OS TSO

Go to ISPF menu 6 (TSO command) or to native line-mode TSO and type EZACMD, as shown in Figure 8-25 on page 365.

```

Menu List Mode Functions Utilities Help
ISPf Command Shell
Enter TSO or Workstation commands below:

==> ex 'tcpip.sezaexec(ezacmd)' 'ipsec -p tcpipa -f display max=20'

Place cursor on choice and press enter to Retrieve command

```

Figure 8-25 EZACMD syntax for TSO

Example 8-45 shows the response to the EZACMD command in Figure 8-25.

Example 8-45 Response to Ipsec command through EZACMD in TSO

```

TSO REXX EZACMD: ipsec command - start - userID=CS02
TSO REXX EZACMD: ipsec -p tcpipa -f display

CS V1R12 ipsec Stack Name: TCPIPA Tue Sep 24 17:42:18 2010
Primary: Filter Function: Display Format: Detail
Source: Stack Profile Scope: Current TotAvail: 4
Logging: On Predecap: Off DVIPSec: No
NatKeepAlive: 0
Defensive Mode: Inactive

***
FilterName: SYSDEFAULTRULE.1
FilterNameExtension: 1
GroupName: n/a
LocalStartActionName: n/a
VpnActionName: n/a
TunnelID: 0x00
Type: Generic
DefensiveType: n/a
State: Active
Action: Permit
Scope: Local
Direction: Outbound
TSO REXX EZACMD: Maximum number of output lines (20) has been reached.
TSO REXX EZACMD: ipsec command - end - RC=4
***

```

8.7.9 Integrating EZACMD into REXX programs in TSO and NetView

The EZACMD command can easily be integrated with other automation-based REXX logic in NetView by using the PIPE command. To preserve the case when used in REXX, use the address NETVASIS prefix for the PIPE command, as shown in Example 8-46.

Example 8-46 EZACMD integrated in REXX via the NetView PIPE command

```
/* NetView REXX */
cmd = 'EZACMD ping -v 127.0.0.1'
address NETVASIS 'PIPE NETV 'cmd' | Corwait 10 | Stem cmdout.'
if cmdout.0 > 0 then do nvix=1 to cmdout.0
say '**' || cmdout.nvix
end
exit(0)
```

There are no specific requirements for using EZACMD in a TSO REXX program. It can be invoked like any other TSO command by using an address command, as shown in Example 8-47.

Example 8-47 EZACMD integrated in REXX through TSO address command

```
/* TSO REXX */
x = outtrap('cmdout.')
address TSO 'ezacmd ipsec -f display max=10'
x = outtrap('OFF')
if cmdout.0 > 0 then do xi=1 to cmdout.0
say '**' || cmdout.xi
end
```

8.7.10 Protecting the EZACMD command

This section discusses protecting the EZACMD command from being issued by unauthorized users.

Console command security

You can protect the EZACMD command in the z/OS console by using a normal RACF OPERCMDS class, as shown in Example 8-48.

Example 8-48 EZACMD protected by RACF

```
CLASS NAME
-----
OPERCMDS MVS.SYSREXX.EXECUTE.EZACMD
LEVEL OWNER UNIVERSAL ACCESS YOUR ACCESS WARNING
-----
00 USER1 NONE READ NO
USER ACCESS ACCESS COUNT
-----
USER1 READ 000000
```

Create an OPERCMDS resource profile with the name of MVS.SYSREXX.EXECUTE.EZACMD to protect the EZACMD command. Only logged in console users who are authorized with READ access to that profile can use the EZACMD command from the z/OS console. This level of security applies to the z/OS console environment only. The z/OS UNIX command security described in Table 8-9 applies to all environments in which the EZACMD command is used.

Table 8-9 SERVAUTH profile applicable to EZACMD

The SERVAUTH profiles	Function
EZB.IPSECCMD.sysname.stackname.command_type	Provides the ability to control ipsec command usage in general.
EZB.IPSECCMD.sysname.DMD_GLOBAL.command_type	Controls whether a user can display (command_type=DISPLAY) or update (command_type=CONTROL) the defensive filters on a system.
EZB.NETMGMT.sysname.clientname.IPSEC.CONTROL	Controls whether a user can issue the ipsec command with the -z option to perform a management action on an NSS IPsec client (for example, to activate and deactivate options).
EZB.NETMGMT.sysname.clientname.IPSEC.DISPLAY	Controls whether a user can issue the ipsec command with the -z option to display options for an NSS IPsec client.
EZB.NETMGMT.sysname.sysname.NSS.DISPLAY	Controls whether a user can issue the ipsec command with the -x option to display NSS IPsec client connections to the NSS server. It also controls whether a user can issue the nssct1 command to display NSS client connections to the NSS server.
EZB.NETMGMT.sysname.sysname.IKED.DISPLAY	Controls whether a user can issue the ipsec command with the -w option to display IKE daemon NSS IPsec client information.
EZB.PAGENT.sysname.image.ptype	Provides the ability to restrict pasearch command, IKE daemon, policy clients, and nslapm2 usage by policy type.

SERVAUTH profiles are especially useful when using the **ipsec** command, so consider using them for that command.

For more information about this topic, see Appendix E, “Steps for preparing to run IP security”, in *z/OS V1R12 Communications Server: IP Configuration Guide*, SC31-8775.

Security when using EZACMD via NetView

Review your NetView security setup and learn which z/OS UNIX security credentials under which the UNIX commands will run. The z/OS NetView supports five different types of operator security, as shown in Table 8-10.

Table 8-10 The five different types of operator security supports by z/OS NetView

NetView SECOPTS.OPERSEC setting	_BPX_USER ID passed to z/OS UNIX by EZACMD	z/OS UNIX command	What it does it
SAFDEF	NetView operator SAF user ID	NetView operator SAF user ID, UID, and GID	SAF checking for both logon passwords and attributes (NETVIEW segment).
SAFCHECK	NetView operator SAF user ID	NetView operator SAF user ID, UID, and GID	Logon passwords checked by SAF. Attributes specified in the DSIOPF/DSIPRF. DATASET, and OPERCMDS classes are checked at the task level.
SAFPW	NetView operator SAF user ID	NetView operator SAF user ID,UID, and GID	Logon passwords checked by SAF. Attributes specified in the DSIOPF/DSIPRF. DATASET, and OPERCMDS classes are checked at the NetView level.
NETVPW	NetView started task user ID	NetView started task SAF user ID, UID, and GID	Logon passwords defined in DSIOPF or DSIEX12. Attributes are specified in DSIOPF/DSIPRF.
MINIMAL	NetView started task user ID	NetView started task SAF user ID, UID, and GID	Ignore logon passwords and attributes.

8.7.11 Diagnosis: diagnosing the EZACMD command

If the EZACMD encounters problems, it will display various error messages. Commands that do not complete within a defined period (for example, 30 seconds) might time out.

In Example 8-49, the IP address cannot be reached from the location where the **ping** command is issued, and System REXX times out the command after 30 seconds.

Example 8-49 EZACMD timeout by REXX

```
REXX30EZACMD 'ping -v 1.12.6.51 MAX=100'
```

```
System REXX EZACMD: ping command - start - userID=CS02
System REXX EZACMD: ping -v 1.12.6.51
System REXX EZACMD: Halt trap entered (likely timeout)
System REXX EZACMD: ping command - end - RC=8
```

```
AXR0203I AXREXX INVOCATION OF EZACMD FAILED.
RETCODE=0000000C RSNCODE=042A0C0A
REQTOKEN=0000400000000000C4BAC3DFB9B67580
DIAG1=00000000 DIAG2=00000000 DIAG3=00000000 DIAG4=00000000
```

8.8 Additional information

Refer to the following content for more information regarding the use of logs, standard commands, tools, and utilities:

- ▶ *z/OS Communications Server: IP System Administrator's Commands*, SC31-8781
- ▶ *z/OS Communications Server: IP Diagnosis Guide*, GC31-8782
- ▶ *z/OS Communications Server: IP Configuration Reference*, SC31-8776
- ▶ *z/OS MVS Diagnosis: Tools and Service Aids*, GA22-7589
- ▶ *z/OS Communications Server: SNA Diagnosis Vol. 1, Techniques and Procedures*, GC31-6850
- ▶ *z/OS Communications Server: SNA Operation*, SC31-8779
- ▶ *MVS Installation Exits*, SA22-7593
- ▶ *Support Element Operations Guide*, SC28-6860

You can find information about z/OS Communications Server product support at the following address:

<http://www.ibm.com/software/network/commserver/zos/support/>

For information about IBM Tivoli OMEGAMON XE for Mainframe Networks, go to the following address:

<http://publib.boulder.ibm.com/tividd/td/IBMTivoliOMEGAMONXEforMainframeNetworks1.0.html>



z/OS in an ensemble

The zEnterprise System brings about a revolution in the end-to-end management of diverse systems, while offering expanded and evolved traditional System z capabilities.

With the zEnterprise System, virtualized resources of both the zEnterprise 196 (z196) and selected IBM blades, which are housed in the zEnterprise BladeCenter Extension (zBX), are pooled together and jointly managed through the zEnterprise Unified Resource Manager (zManager).

End-to-end solutions based on multi-platform workloads can be deployed across the zEnterprise System infrastructure and benefit from the System z traditional qualities of service, including high availability, and simplified and improved management of the virtualized resources.

This chapter discusses the following topics.

Section	Topic
9.1, “Basic concepts” on page 372	Basic concepts of zEnterprise
9.2, “Connectivity” on page 372	The connections between the z196 and the zBX
9.3, “Enabling z/OS as a member of the ensemble” on page 373	Requirements for z/OS to become a member of the ensemble
9.4, “Defining and activating the z/OS ensemble interfaces” on page 380	How to define and verify the z/OS ensemble interfaces

9.1 Basic concepts

Each z196, with its optional zBX, makes up a *node* of a zEnterprise ensemble. A zEnterprise ensemble is composed of up to eight members, with up to eight z196 servers and up to eight zBXs, dedicated integrated networks for management and data, and the zManager function. With the zManager, the z196 provides advanced end-to-end management capabilities for the diverse systems housed in the zBX.

The zBX components are configured, managed, and serviced the same way as the other components of the z196. Despite the fact that the zBX processors are not System z PUs and run specific software, including hypervisors, the software intrinsic to the zBX components does not require any additional administration effort or tuning by the user. In fact, it is handled as System z Licensed Internal Code. The zBX hardware features are part of the mainframe, not add-ons.

9.2 Connectivity

Figure 9-1 shows a zEnterprise node, consisting of a z196 and a zBX. The first rack (Rack B) in the zBX has four top of rack (TOR) switches for network connectivity: two TOR switches for the intranode management network (INMN) and two TOR switches for the intraensemble data network (IEDN).

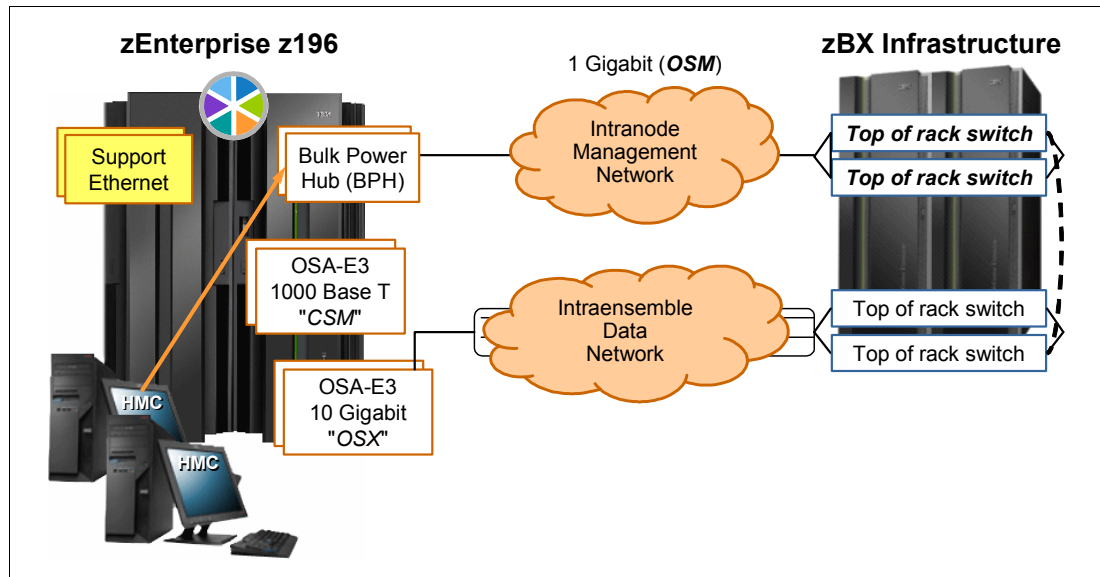


Figure 9-1 zEnterprise Node

9.2.1 Intranode management network (INMN)

The INMN is one of the ensemble's two private and secure internal networks. INMN is used by the zManager functions in the HMC. The z196 introduces the OSA-Express for Unified Resource Manager (OSM) CHPID type. The OSM connections are from OSA-Express3 ports to the Bulk Power Hubs (BPHs) in the z196. The BPHs are connected to the INMN TOR switches in zBX. The INMN requires two OSA-Express3 1000BASE-T ports from separate features.

Note: Access to INMN is restricted to authorized management applications only, and is available thru Port 0 of the OSA-Express3 1000BASE-T feature. To use the INMN, stacks must be IPv6-enabled.

9.2.2 Intraensemble data network (IEDN)

The IEDN is the ensemble's other private and secure internal network. IEDN is used for communications across the virtualized images (LPARs and virtual machines on z/VM and the IBM blades). The z196 introduces the OSA-Express for zBX (OSX) CHPID type. The OSX connection is from the z196 to the IEDN TOR switches in zBX. The IEDN requires two OSA-Express3 10 GbE ports from separate features.

9.3 Enabling z/OS as a member of the ensemble

To become a member of the ensemble, you must enable z/OS by completing the following three tasks:

1. Enable the z/OS TCP/IP stack for IPv6 so that it can participate in the INMN.
2. Specify in VTAM that this z/OS image is to participate in the ensemble. To allow z/OS Communications Server to have connectivity to the IEDN and the INMN, the parameter (ENSEMBLE=YES) must be added to the VTAM start options.
3. Define the necessary interfaces in the TCP/IP stack for connecting into the ensemble. INMN connections (via CHPID type OSM) are dynamically created when an ensemble is defined. For IEDN connections (CHPID type OSX), a VTAM TRLEs must be created in one of two ways:
 - a. Dynamically by VTAM, using a prefix and the CHPID
 - b. Manually in VTAM as a TRLE majornode

9.3.1 Enabling z/OS for IPv6

To enable the z/OS image for IPv6, add the IPv6 addressing family to the BPXPRMnn member in hlq.PARMLIB (see Example 9-1).

Example 9-1 BPXPRMnn changes to add IPv6 address family to the z/OS image

NETWORK	DOMAINNAME(AF_INET)	A
	DOMAINNUMBER(2)	
	MAXSOCKETS(10000)	
	TYPE(CINET)	
	INADDRANYPORT(10000)	
	INADDRANYCOUNT(2000)	
NETWORK	DOMAINNAME(AF_INET6)	B
	DOMAINNUMBER(19)	
	MAXSOCKETS(10000)	
	TYPE(CINET)	

With these definitions in BPXPRMnn *dual-mode* TCP/IP stacks are supported, IPv4 with AF_INET (**A**) and IPv6 with AF_INET6 (**B**). If your z/OS image contains only one TCP/IP stack, your definition is simpler, indicating a TYPE(INET) and omitting the INADDRANYPORT and INADDRANYCOUNT parameters.

Note: MAXSOCKETS is enforced independently for AF_INET and AF_INET6 sockets.

The INADDRANYPORT and INADDRANYCOUNT values for NETWORK AF_INET6 are taken from the NETWORK AF_INET statement. These values are ignored if they are specified on the NETWORK statement for AF_INET6.

- You can add the AF_INET6 NETWORK statement dynamically with a **SETOMVS RESET** command against the BPXPRMnn member to which the new statement has been added, or you can re-IPL z/OS to pick up the new statement. Once the statement is installed you must recycle TCP/IP to pick up the AF_INET6 physical file system. To verify that you have a dual-mode z/OS, issue the **D OMVS,PFS** command and examine the output, as shown in Example 9-2.

Example 9-2 Verifying that IPv6 is available in the z/OS image

```
D OMVS,PFS
BPX0068I 16.49.39 DISPLAY OMVS 052
OMVS      000F ACTIVE              OMVS=(2A)
PFS CONFIGURATION INFORMATION
PFS TYPE  ENTRY      ASNAME  DESC      ST      START/EXIT TIME
NFS       GFSCINIT   NFSCCNT  REMOTE    A       2010/11/09 14.47.46
CINET     BPXTCINT                SOCKETS   A       2010/11/09 14.47.46
UDS       BPXTUINT                SOCKETS   A       2010/11/09 14.47.46
ZFS       IOEFSCM     ZFS      LOCAL     A       2010/11/09 14.47.42
AUTOMNT   BPXTAMD                LOCAL     A       2010/11/09 14.47.42
TFS       BPXTFS                LOCAL     A       2010/11/09 14.47.42
HFS       GFUAINIT                LOCAL     A       2010/11/09 14.47.42

PFS TYPE  DOMAIN      MAXSOCK  OPNSOCK  HIGHUSED
PFS TYPE  DOMAIN      MAXSOCK  OPNSOCK  HIGHUSED
CINET     AF_INET6 1 10000    41       41
          AF_INET 10000    22       24
UDS       AF_UNIX 10000    7        7

SUBTYPES OF COMMON INET
PFS NAME  ENTRY      START/EXIT TIME      STATUS  FLAGS
TCP/IP    EZBPFINI   2010/11/09 14.47.56  ACT     SC
TCP/IPA   EZBPFINI                INACT
TCP/IPB   EZBPFINI                INACT
TCP/IPC   EZBPFINI                INACT
TCP/IPD   EZBPFINI                INACT
TCP/IPE   EZBPFINI                INACT
TCP/IPF   EZBPFINI                INACT

PFS TYPE  FILESYSTYPE PARAMETER INFORMATION
ZFS       PRM=(30,00)
HFS

          CURRENT VALUES: FIXED(0) VIRTUAL(2009)

PFS TYPE  STATUS INFORMATION
AUTOMNT   TIME=2010/11/09 14:55:23 SYSTEM=SC33      USER=OMVSKERN
          POLICY=/etc/auto.master
```

Observe at **A** that we are now running Common INET (CINET) with the IPv6 physical file system and the address family for IPv6 (AF_INET6).

Next display the TCP/IP stack's home list to verify that a LOOPBACK6 device appears there, indicating that the stack itself is enabled for IPv6 (Example 9-3).

Example 9-3 A z/OS display of the dual-mode TCP/IP stack and its home list with IPv6 enabled

```
D TCPIP,,N,HOME
EZD0101I NETSTAT CS V1R12 TCPIP 482
HOME ADDRESS LIST:
LINKNAME:  OSA2100LNK
  ADDRESS:  9.12.4.211
  FLAGS:    PRIMARY
LINKNAME:  OSA2120LNK
  ADDRESS:  9.12.4.212
  FLAGS:
LINKNAME:  TOSAME1
  ADDRESS:  192.1.1.1
  FLAGS:
LINKNAME:  VIPL090C0525
  ADDRESS:  9.12.5.37
  FLAGS:
LINKNAME:  LOOPBACK
  ADDRESS:  127.0.0.1
  FLAGS:
INTFNAME:  LOOPBACK6      A
  ADDRESS:  ::1
  TYPE:    LOOPBACK
  FLAGS:
6 OF 6 RECORDS DISPLAYED
END OF THE REPORT
```

You see five interfaces with their IPv4 addresses that existed in the TCP/IP stack prior to IPv6 enablement. You also see a sixth interface at **A**, LOOPBACK6, that has been generated because of the IPv6 enablement in the BPXPRMnn member. The format of the IPv6 loopback address includes colons (:) to indicate that the address is 128 bits long, with leading zeroes followed by 1.

9.3.2 Enabling VTAM for the ensemble

Although a z/OS LPAR can be automatically detected by the zManager firmware, resulting in the creation of a Virtual Server container for the z/OS operating system, z/OS itself cannot participate in an ensemble until it has been enabled. You must make a change to VTAM to allow this z/OS image to participate in the ensemble.

The ENSEMBLE start option must be changed to indicate *ENSEMBLE=YES*. You can enable the option in either of two ways:

- ▶ Change the VTAM Start Options to include the option setting.
- ▶ Issue a MODIFY command.

Before we enable the ENSEMBLE option, we should examine the VTAM Transport Resource List Entries (TRLEs) to determine if any have been built for the INMN network. Therefore, we display the VTAM member ISTTRL (Example 9-4 on page 376).

Example 9-4 The TRLEs prior to Ensemble enablement

```
D NET,E,ID=ISTTRL
IST097I DISPLAY ACCEPTED
IST075I NAME = ISTTRL, TYPE = TRL MAJOR NODE 811
IST1314I TRLE = IUTIQDF6 STATUS = ACTIV CONTROL = MPC
IST1314I TRLE = IUTIQDF5 STATUS = ACTIV CONTROL = MPC
IST1314I TRLE = IUTIQDF4 STATUS = ACTIV CONTROL = MPC
IST1314I TRLE = ISTT3033 STATUS = ACTIV CONTROL = XCF
IST1314I TRLE = ISTT3032 STATUS = ACTIV CONTROL = XCF
IST1314I TRLE = ISTT3031 STATUS = ACTIV CONTROL = XCF
IST1314I TRLE = IUTIQDIO STATUS = ACTIV CONTROL = MPC
IST1314I TRLE = IUTSAMEH STATUS = ACTIV CONTROL = MPC
IST314I END
```

The INMN TRLEs that are built for VTAM should have a prefix of *IUTM*; there are none in the list displayed in Example 9-5. Therefore, we must enable the VTAM ENSEMBLE start option as a first step towards obtaining the INMN TRLEs.

In our LPAR, the VTAM Start Options are SYS1.VTAMLST(ATCSTR30). We have added a new parameter to this Start Option list so that the next recycle of VTAM can make z/OS ready for the ensemble.

Example 9-5 Ensemble Start Option in ATCSTR30

```
SSCPID=30,NOPROMPT,NETID=USIBMSC,SSCPNAME=SC30M, X
CONFIG=30,SUPP=NOSUP, X
HOSTPU=SC30PU, X
NETID=USIBMSC, X
IQDCHPID=F3, X

.....

ENSEMBLE=YES, A X
.....
```

Fortunately, the start option ENSEMBLE (A) is dynamically modifiable; therefore, prior to a new IPL of VTAM, we are able to change the default of ENSEMBLE=NO to ENSEMBLE=YES. We use this command from the z/OS console to change the setting of the parameter:

```
F NET,VTAMOPTS,ENSEMBLE=YES
```

If you display the TRLEs at VTAM, you see that the dynamic TRLEs for the INMN have still not been created and look the same as they did in Example 9-4. They are initially created only with an initialization (or recycle) of the TCP/IP stack.

9.3.3 Validating the ensemble interfaces in z/OS

Our next display of the TRLEs shows that both the INMN and the IEDN TRLEs have been created. See Example 9-6.

Example 9-6 Displaying the TRLEs for the INMN and the IEDN Connections in VTAM

```
D NET,E,ID=ISTTRL
IST097I DISPLAY ACCEPTED
IST075I NAME = ISTTRL, TYPE = TRL MAJOR NODE 248
```


IST1314I	TRLE = IUTXT019	STATUS = ACTIV	CONTROL = MPC	A
IST1314I	TRLE = IUTXT018	STATUS = ACTIV	CONTROL = MPC	B
IST1314I	TRLE = IUTIQDF6	STATUS = INACT	CONTROL = MPC	
IST1314I	TRLE = IUTIQDF5	STATUS = INACT	CONTROL = MPC	
IST1314I	TRLE = IUTIQDF4	STATUS = INACT	CONTROL = MPC	
IST1314I	TRLE = ISTT3033	STATUS = ACTIV	CONTROL = XCF	
IST1314I	TRLE = ISTT3032	STATUS = ACTIV	CONTROL = XCF	
IST1314I	TRLE = ISTT3031	STATUS = ACTIV	CONTROL = XCF	
IST1314I	TRLE = IUTMT00B	STATUS = ACTIV	CONTROL = MPC	C
IST1314I	TRLE = IUTMT00A	STATUS = ACTIV	CONTROL = MPC	D
IST1314I	TRLE = IUTIQDIO	STATUS = INACT	CONTROL = MPC	
IST1314I	TRLE = IUTSAMEH	STATUS = INACT	CONTROL = MPC	
IST314I	END			

In Example 9-6 at points **A** and **B** we observe that we now have two active IEDN TRLEs by the names of *IUTXT019* and *IUTXT018*, where the last two digits of each name represent the CHPID number of the TRLEs. At points **C** and **D** we observe that we now have two active INMN TRLEs by the names of *IUTMT00B* and *IUTMT00A*. 00A and 00B are the suffixes of the TRLE names and are also our CHPIDs for the two OSM OSA ports.

A simple display of one of the ensemble TRLEs in VTAM shows us which device addresses from the IOCDS have been used to build the TRLE. As an example, we display the TRLE for an OSM CHPID (see Example 9-7).

Example 9-7 TRLE display of the devices to be used for an OSM interface

```

D NET,E,ID=IUTMT00A
IST097I DISPLAY ACCEPTED
IST075I NAME = IUTMT00A, TYPE = TRLE 281
IST486I STATUS= ACTIV, DESIRED STATE= ACTIV
IST087I TYPE = LEASED, CONTROL = MPC, HPDT = YES
IST1954I TRL MAJOR NODE = ISTTRL
IST1715I MPCLEVEL = QDIO MPCUSAGE = SHARE
IST2263I PORTNAME = IUTMP00A PORTNUM = 0 OSA CODE LEVEL = 0906
IST2337I CHPID TYPE = OSM CHPID = 0A
IST1577I HEADER SIZE = 4096 DATA SIZE = 0 STORAGE = ***NA***
IST1221I WRITE DEV = 2341 STATUS = ACTIVE STATE = ONLINE
IST1577I HEADER SIZE = 4092 DATA SIZE = 0 STORAGE = ***NA***
IST1221I READ DEV = 2340 STATUS = ACTIVE STATE = ONLINE
IST924I -----
IST1221I DATA DEV = 2342 STATUS = ACTIVE STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST1717I ULPID = TCPIP
IST2310I ACCELERATED ROUTING DISABLED
IST2331I QUEUE QUEUE READ
IST2332I ID TYPE STORAGE
IST2205I -----
IST2333I RD/1 PRIMARY 4.0M(64 SBALS)
IST2305I NUMBER OF DISCARDED INBOUND READ BUFFERS = 0
IST1757I PRIORITY1: UNCONGESTED PRIORITY2: *****NA*****
IST1757I PRIORITY3: *****NA***** PRIORITY4: *****NA*****
IST2190I DEVICEID PARAMETER FOR OSAENTA TRACE COMMAND = 01-01-00-02
IST1801I UNITS OF WORK FOR NCB AT ADDRESS X'2807E010'
IST1802I P1 CURRENT = 1 AVERAGE = 2 MAXIMUM = 4
IST924I -----
IST1221I DATA DEV = 2343 STATUS = RESET STATE = N/A

```

```

IST1724I I/O TRACE = OFF  TRACE LENGTH = *NA*
IST924I -----
.....
IST1500I STATE TRACE = OFF
IST314I END

```

The next display shows the TCP/IP home list in the netstat. Look for the INMN interfaces at points **A** and **B** in Example 9-8.

Example 9-8 Displaying the INMN and IEDN addresses and interface names in TCP/IP

```

D TCPIP,,N,HOME
EZD0101I NETSTAT CS V1R12 TCPIP 277
HOME ADDRESS LIST:
LINKNAME:  OSA2100LNK
  ADDRESS:  9.12.4.211
  FLAGS:    PRIMARY
LINKNAME:  OSA2120LNK
  ADDRESS:  9.12.4.212
  FLAGS:
LINKNAME:  TOSAME1
  ADDRESS:  192.1.1.1
  FLAGS:
LINKNAME:  LOOPBACK
  ADDRESS:  127.0.0.1
  FLAGS:
INTFNAME:  LOOPBACK6
  ADDRESS:  ::1
  TYPE:    LOOPBACK
  FLAGS:
INTFNAME:  EZ60SM01
  ADDRESS:  FE80::76FF:FE9E:C008      A
  TYPE:    LINK_LOCAL
  FLAGS:  AUTOCONFIGURED
INTFNAME:  EZ60SM02
  ADDRESS:  FE80::76FF:FE87:8009      B
  TYPE:    LINK_LOCAL
  FLAGS:  AUTOCONFIGURED
11 OF 11 RECORDS DISPLAYED
END OF THE REPORT

```

Note how the addresses for the INMN interfaces (**A** and **B**) are IPv6 LINK_LOCAL addresses beginning with the prefix of FE80. The dynamically assigned names for the *autoconfigured* addresses are *EZ60SM01* and *EZ60SM02*.

9.3.4 Displaying information about the OSM interfaces

We display one of the two OSM interfaces, EZ60SM01, in Example 9-9.

Example 9-9 Displaying the OSA Information for an OSM OSA interface

```

D TCPIP,,OSAINFO,INTFNAME=EZ60SM01      1
EZZ0053I COMMAND DISPLAY TCPIP,,OSAINFO COMPLETED SUCCESSFULLY
EZD0031I TCP/IP CS V1R12  TCPIP Name: TCPIP      15:12:35 212
Display OSAINFO results for IntfName: EZ60SM01

```

```

PortName: IUTMP00A 2 PortNum: 00 3 Datapath: 2342 4 RealAddr: 0002
PCHID: 0531 5 CHPID: 0A 6 CHPID Type: OSM 7 OSA code level: 0906 8
Gen: OSA-E3 Active speed/mode: 1000 mb/sec full duplex
Media: Copper 9 Jumbo frames: Yes 10 Isolate: Yes 11
PhysicalMACAddr: 00145E7769EC 12 LocallyCfgMACAddr: 000000000000
Queues defined 13 Out: 1 In: 1 Ancillary queues in use: 0
Connection Mode: Layer 2 14
SAPSup: 0009F603 SAPEna: 00082603
Layer 2 attributes:
  VLAN ID: N/A 15 VMAC Active: Yes 16
  VMAC Addr: 0200769EC008 17 VMAC Origin: OSA 18
15 of 15 lines displayed
End of report

```

Example 9-9 provides valuable information with the OSAINFO display:

- 1 Syntax of command to display a single OSM, dynamically generated interface.
- 2 Dynamically assigned portname for an OSM TRLE and interface.
- 3 OSM must be on Port number 0 of the OSM adapter.
- 4 The datapath assignment correlates with the IOCDS for the generated TRLE.
- 5, 6, 7 The PCHID, CHPID number, and CHPID Type correlate with the IOCDS.
- 8 MCL level of the OSA port.
- 9, 10, 11 This is a Copper OSA, capable of Jumbo frames, operating in ISOLATE mode.
- 12 The physical MAC address of the OSA port.
- 13 The management OSA does not perform priority queuing in either direction.
- 14 The management OSA operates only in Layer 2 mode with no Layer 3 routing.
- 15 The management OSA port is operating in ACCESS mode; the TOR switch assigns the VLAN ID; the stack is unaware of any VLAN ID.
- 16, 17 A Virtual MAC is active and its address is displayed with the Ensemble prefix.
- 18 The Virtual MAC was fully generated by the OSA itself using the Ensemble prefix.

The next display is a typical NETSTAT output display. We use it for more information about the OSM OSA port.

Example 9-10 Device display for an OSM interface

```

D TCPIP,,N,DEV,INTFNAME=EZ60SM01
EZD0101I NETSTAT CS V1R12 TCPIP 214
INTFNAME: EZ60SM01 INTFTYPE: IPAQENET6 INTFSTATUS: READY
PORTNAME: IUTMP00A DATAPATH: 2342 DATAPATHSTATUS: READY
CHPIDTYPE: OSM
QUESIZE: 0 SPEED: 0000001000
VMACADDR: 0200769EC008 VMACORIGIN: OSA VMACROUTER: ALL
DUPADDRDET: 1
CFGMTU: NONE ACTMTU: 1500
VLANID: NONE VLANPRIORITY: DISABLED
READSTORAGE: GLOBAL (4096K)
INBPERF: BALANCED
SECCLASS: 255 MONSYSPLEX: NO
ISOLATE: YES OPTLATENCYMODE: NO
TEMPPREFIX: NONE
MULTICAST SPECIFIC:
MULTICAST CAPABILITY: YES
GROUP: FF02::1:FF9E:C008
GROUP: FF02::1:FF9E:C008
REFCNT: 000000001 SRCFLTMD: EXCLUDE

```

```

SRCADDR: NONE
GROUP:      FF01::1
REFCNT: 0000000001 SRCFLTMD: EXCLUDE
SRCADDR: NONE
GROUP:      FF02::1
REFCNT: 0000000001 SRCFLTMD: EXCLUDE
SRCADDR: NONE
INTERFACE STATISTICS:
BYTESIN                      = 1497186636
INBOUND PACKETS              = 59183
INBOUND PACKETS IN ERROR     = 39
INBOUND PACKETS DISCARDED    = 0
INBOUND PACKETS WITH NO PROTOCOL = 0
BYTESOUT                     = 42964196
OUTBOUND PACKETS             = 59294
OUTBOUND PACKETS IN ERROR    = 0
OUTBOUND PACKETS DISCARDED   = 0
IPV6 LAN GROUP SUMMARY
LANGROUP: 00005
NAME          STATUS      NDOWNER      VIPAOWNER
-----
EZ60SM01      ACTIVE      EZ60SM01     YES
EZ60SM02      ACTIVE      EZ60SM02     NO
1 OF 1 RECORDS DISPLAYED
END OF THE REPORT

```

9.4 Defining and activating the z/OS ensemble interfaces

There are two ways to define OSX INTERFACES for z/OS:

1. You can allow VTAM to build the TRLEs for the IP interfaces dynamically by referring to the CHPID number.
2. You can predefine the VTAM TRLEs with a PORTNAME, and then code the IP interface definitions using the PORTNAME.

Important: The PORTNAMEs assigned to IEDN CHPIDs must be consistent across all sharing LPARs for z/OS. Therefore, you should standardize the TRL definition type for each OSA CHPID. That is, if you have four z/OS LPARs sharing an OSA port, all four must define the TRLEs in the same way: either through dynamic definition of the PORTNAME (choice 1 in the previous list), or through a predefinition (choice 2). Do not attempt to mix the definition types on a single OSA port; this will result in a failure of the IEDN interfaces definition and activation.

Example 9-11 shows how to define the INTERFACES using the CHPID number. We inserted the statements for the interfaces on IEDN VLANs 99 and 1034 into the TCP/IP profile.

Example 9-11 IEDN interface statements added to SYS1.TCPPARMS(PROF30)

```

; ----- IEDN INTERFACES FOR ENSEMBLE ATTACHMENTS -----
; ---VLAN 99---for DB2---
INTERFACE OSX2300
DEFINE IPAQENET CHPIDTYPE OSX

```

```

CHPID 18    VLANID 99
MTU 8992      IPADDR 172.30.99.1/24
;;
INTERFACE OSX2320
  DEFINE IPAQENET CHPIDTYPE OSX
  CHPID 19    VLANID 99
MTU 8992      IPADDR 172.30.99.2/24
;;
;; ---VLAN 1034 ---for Administration---
INTERFACE OSX2300A
  DEFINE IPAQENET CHPIDTYPE OSX
  CHPID 18    VLANID 1034
MTU 8992      IPADDR 172.30.10.1/24
;;
INTERFACE OSX2320A
  DEFINE IPAQENET CHPIDTYPE OSX
  CHPID 19    VLANID 1034
MTU 8992      IPADDR 172.30.10.2/24
;
; -----END:    IEDN INTERFACES FOR ENSEMBLE ATTACHMENTS -----

```

Note: The IEDN interfaces can also be dynamically added with an OBEYFILE, but the INMN connections require a stack initiation for the initial dynamic creation.

Example 9-12 shows the TCP/IP home list in the netstat.

Example 9-12 Displaying the INMN and IEDN addresses and interface names in TCP/IP

```

D TCPIP,,N,HOME
EZD0101I NETSTAT CS V1R12 TCPIP 277
HOME ADDRESS LIST:
LINKNAME:  OSA2100LNK
  ADDRESS:  9.12.4.211
  FLAGS:    PRIMARY
LINKNAME:  OSA2120LNK
  ADDRESS:  9.12.4.212
  FLAGS:
LINKNAME:  TOSAME1
  ADDRESS:  192.1.1.1
  FLAGS:
LINKNAME:  LOOPBACK
  ADDRESS:  127.0.0.1
  FLAGS:
INTFNAME:  OSX2300           1
  ADDRESS:  172.30.99.1
  FLAGS:
INTFNAME:  OSX2320           2
  ADDRESS:  172.30.99.2
  FLAGS:
INTFNAME:  OSX2300A          3
  ADDRESS:  172.30.10.1
  FLAGS:
INTFNAME:  OSX2320A          4
  ADDRESS:  172.30.10.2
  FLAGS:

```

```

INTFNAME: LOOPBACK6
ADDRESS: ::1
TYPE: LOOPBACK
FLAGS:
INTFNAME: EZ60SM01
ADDRESS: FE80::76FF:FE9E:C008
TYPE: LINK_LOCAL
FLAGS: AUTOCONFIGURED
INTFNAME: EZ60SM02
ADDRESS: FE80::76FF:FE87:8009
TYPE: LINK_LOCAL
FLAGS: AUTOCONFIGURED
11 OF 11 RECORDS DISPLAYED
END OF THE REPORT

```

Observe the names of the OSX interfaces at **1** through **4**. These are the names that we preassigned in our INTERFACE definitions in the TCP/IP profile. The interfaces have been assigned the IPv4 addresses (C) that we had planned for.

9.4.1 Displaying information about the OSX interfaces

The **D TCPIP,,OSAINFO,INTFNAME** and the **D TCPIP,,N,DEV,INTFNAME** commands can be used to provide information about the OSA interface (Example 9-13 and Example 9-14).

Example 9-13 Displaying the OSA Information for an OSX OSA interface

```

D TCPIP,,OSAINFO,INTFNAME=OSX2300
EZZ0053I COMMAND DISPLAY TCPIP,,OSAINFO COMPLETED SUCCESSFULLY
EZD0031I TCP/IP CS V1R12 TCPIP Name: TCPIP 15:06:03 203
Display OSAINFO results for IntfName: OSX2300
PortName: IUTXP018 PortNum: 00 Datapath: 2302 RealAddr: 0002
PCHID: 0590 CHPID: 18 CHPID Type: OSD OSA code level: 0D0A
Gen: OSA-E3 Active speed/mode: 10 gigabit full duplex
Media: Multimode Fiber Jumbo frames: Yes Isolate: No
PhysicalMACAddr: 001A643B2135 LocallyCfgMACAddr: 000000000000
Queues defined Out: 4 In: 1 Ancillary queues in use: 0
Connection Mode: Layer 3 IPv4: Yes IPv6: No
SAPSup: 000FF603 SAPEna: 0008A603
IPv4 attributes:
VLAN ID: 99 VMAC Active: Yes
VMAC Addr: 02BECB000002 VMAC Origin: Cfg VMAC Router: All
AsstParmsEna: 00200C57 OutCkSumEna: 0000001A InCkSumEna: 0000001A
Registered Addresses:
IPv4 Unicast Addresses:
IPv4 Unicast Addresses:
ARP: Yes Addr: 172.30.99.1
Total number of IPv4 addresses: 1
IPv4 Multicast Addresses:
MAC: 01005E000001 Addr: 224.0.0.1
Total number of IPv4 addresses: 1
23 of 23 lines displayed
End of report

```

Example 9-14 Device display for an OSX interface

```

D TCPIP,,N,DEV,INTFNAME=OSX2300
EZD0101I NETSTAT CS V1R12 TCPIP 216
INTFNAME: OSX2300          INTFTYPE: IPAQENET    INTFSTATUS: READY
    PORTNAME: IUTXP018      DATAPATH: 2302        DATAPATHSTATUS: READY
    CHPIDTYPE: OSX          CHPID: 18
    SPEED: 0000010000
    IPBROADCASTCAPABILITY: NO
    VMACADDR: 02BECB000002  VMACORIGIN: OSA      VMACROUTER: ALL
    ARPOFFLOAD: YES         ARPOFFLOADINFO: YES
    CFGMTU: 8992            ACTMTU: 8992
    IPADDR: 172.30.99.1/24
    VLANID: 99              VLANPRIORITY: DISABLED
    DYNVLANREGCFG: NO       DYNVLANREGCAP: YES
    READSTORAGE: GLOBAL (4096K)
    INBPERF: DYNAMIC
        WORKLOADQUEUEING: NO
    CHECKSUMOFFLOAD: YES
    SECCLASS: 255           MONSYSPLEX: NO
    SECCLASS: 255           MONSYSPLEX: NO
    ISOLATE: NO             OPTLATENCYMODE: NO
MULTICAST SPECIFIC:
MULTICAST CAPABILITY: YES
GROUP          REFCNT          SRCFLTMD
-----
224.0.0.1      0000000001      EXCLUDE
    SRCADDR: NONE
INTERFACE STATISTICS:
BYTESIN                = 168
INBOUND PACKETS        = 2
INBOUND PACKETS IN ERROR = 0
INBOUND PACKETS DISCARDED = 0
INBOUND PACKETS WITH NO PROTOCOL = 0
BYTESOUT               = 0
OUTBOUND PACKETS       = 0
OUTBOUND PACKETS IN ERROR = 0
OUTBOUND PACKETS DISCARDED = 0
OUTBOUND PACKETS DISCARDED = 0
IPV4 LAN GROUP SUMMARY
LANGROUP: 00008
NAME          STATUS      ARPOWNER      VIPAOWNER
-----
OSX2300      ACTIVE      OSX2300      YES
1 OF 1 RECORDS DISPLAYED
END OF THE REPORT

```

9.5 References

For more information regarding the zEnterprise System and z/OS Communications Server ensemble setup, refer to the following documentation:

- ▶ z/OS Communications Server ensemble implementation:
 - *z/OS Communications Server V1R12 SNA Network Implementation Guide*, SC31-8777
 - *z/OS Communications Server V1R12 SNA Network Definition Reference*, SC31-8778
 - *z/OS Communications Server V1R12 IP Configuration Guide*, SC31-8775
- ▶ IPv6 information:
 - *z/OS Communications Server V1R12 IPv6 Network and Application Design Guide*, SC31-8885
- ▶ IBM Redbooks publications:
 - *IBM zEnterprise Technical Introduction*, SG24-7832
 - *IBM zEnterprise Technical Guide*, SG24-7833
 - *IBM zEnterprise Configuration Setup*, SG24-7834
 - *IBM System p Advanced POWER Virtualization Best Practices*, REDP-4194
 - *IBM BladeCenter JS12 and JS22 Implementation Guide*, SG24-7655



IPv6 support

IPv6 has gained greater acceptance in the industry in recent years. IPv6 is attractive because it resolves some deficiencies of IPv4 in the following areas:

- ▶ IP address space
- ▶ Auto-configuration
- ▶ Security
- ▶ Quality of service
- ▶ Anycast and multicast

This appendix discusses the following topics.

Section	Topic
"Overview of IPv6" on page 386	Basic concepts of IPv6
"Importance of IPv6" on page 386	Key characteristics of IPv6 and why it might be important in your environment
"Common design scenarios for IPv6" on page 387	Commonly implemented IPv6 design scenarios, their dependencies, advantages, considerations, and our recommendations
"How IPv6 is implemented in z/OS Communications Server" on page 389	Selected implementation scenarios, tasks, configuration examples, and problem determination suggestions

Overview of IPv6

Internet Protocol Version 6 (IPv6) is the next generation of the Internet protocol designed to replace the current Internet Protocol Version 4 (IPv4).

IPv6 was developed to resolve impending problems related to the limitations of IPv4 and the rapidly-growing demand for IP resources and functionality; the most significant issue is the diminishing supply and expected shortages of IPv4 addresses.

Using IPv4 32-bit addressing allows for over 4 million nodes, each with a globally unique address. This current IPv4 space will be unable to satisfy the huge expected increase in the number of users on the Internet. The expected shortage will be exacerbated by the requirements of emerging technologies such as PDAs, HomeArea Networks, and Internet-connected commodities such as automotive and integrated telephone services. IPv6 uses 128-bit addressing and will generate a space large enough to last for the foreseeable future.

Importance of IPv6

IPv6 is important because it addresses the limitations of IPv4, such as:

- ▶ 128-bit addressing
This quadruples the network address bits from 32 to 128, thereby significantly increasing the number of possible unique IP addresses to comfortably accommodate on the Internet. This huge address space obviates the need for private addresses and Network Address Translators (NATs).
- ▶ Simplified header formats
This allows for more efficient packet handling and reduced bandwidth cost.
- ▶ Hierarchical addressing and routing
This keeps routing tables small and backbone routing efficient by using address prefixes rather than address classes.
- ▶ Improved support for options
This changes the way IP header options are encoded, allowing more efficient forwarding and greater flexibility.
- ▶ Address auto-configuration
This allows stateless IP address configuration without a configuration server.
- ▶ Security
IPv6 brings greater authentication and privacy capabilities through the definition of extensions.
- ▶ Quality of Service (QoS)
QoS is provided through a traffic class byte in the header.

Common design scenarios for IPv6

Although as explained, there are predictable improvements over IPv4, the success of any IPv6 implementation depends on the ability to have IPv6 *coexist* with IPv4. Because of the pervasiveness of IPv4, this coexistence will be around for some time. Therefore, the development of technologies and mechanisms to facilitate coexistence is as important as the deployment strategy for IPv6. In the following sections we will discuss some of the *coexistence* technologies available today.

Tunneling

Tunneling is the transmission of IPv6 traffic encapsulated within IPv4 packets over an IPv4 connection. Tunnels are used primarily to connect remote IPv6 networks, or to simply connect an IPv6 network over an IPv4 network infrastructure.

Dependencies

All tunnel mechanisms require that the endpoints of the tunnel run in dual-stack mode. A dual-stack router is a router running *both* versions of IP. There are other dependencies based on the tunneling mechanism used.

For example, an IPv6 *Manually* configured tunnel requires an ISP-registered IP address. The *Automatic* tunnel mechanism requires IPv6 prefixes. Intra-Site Automatic Tunnel Addressing Protocol (ISATAP) tunnels only require a dual-stack router, but they are not yet commercially available and *6over4* tunnels are not supported by vendor router software.

Advantages

Tunneling allows the implementation of IPv6 without any significant upgrades to the existing infrastructure, and therefore does not risk interrupting the existing services provided by the IPv4 network.

Considerations

There are various tunneling mechanisms designed to do primarily different things, so you must give careful consideration to the mechanism you choose. Some are primarily used for stable and secure links for regular communications. Others are primarily used for single hosts or small sites, with low data traffic volumes.

Dedicated data links

Network architects can choose a separate ATM, or frame relay Permanent Virtual Circuits (PVCs), or separate optical media, to run IPv6 traffic across. The only requirement is the reconfiguration of the routers (with IPv6 support). Note that these links can *only* be used for IPv6 traffic.

Dependencies

Dual-stack routers with IPv6 and IPv4 addresses are required to provide access to the WAN. Access to a Domain Name System (DNS) is needed to resolve IPv6 names and addresses.

Advantages

Use of the existing Layer 2 infrastructure makes this implementation less complex and immediate. This implementation is not disruptive, apart from a schedule change for router configuration, and therefore there is little impact to the status quo.

Considerations

All routers on the WAN need to support IPv6 over dedicated data links. Additional costs for the those links will be incurred until the environment is completely migrated over to IPv6.

MPLS backbones

An MPLS IPv4 core network can enable IPv6 domains to communicate over Multi-Protocol Label Switching (MPLS) backbones. It is therefore primarily used by enterprises and service providers. There are various implementations of this strategy, ranging from no changes and no impact to changes and risks. This means that the closer the IPv6 implementation is to the client edge, the less expensive it becomes. In contrast, the closer the IPv6 implementation is to the service provider edge, the more expensive it becomes.

Dependencies

Dependencies vary from router configuration to specific hardware requirements to software upgrades, depending on the service provider solution.

Advantages

Use of this strategy requires minor modifications to the infrastructure and minor reconfigurations of the core routers. It is therefore a strategy that could have little or no impact to your environment, involving low costs and low risks.

Considerations

Considerations also vary, depending on the strategy chosen. For example, using the Circuit Transport over MPLS strategy does not support a mix of IPv4 and IPv6 traffic. IPv6 on service provider edge routers do not support Virtual Private Networks (VPNs) or Virtual Routing and Forwarding (VRF), currently.

Dual-stack backbones

A dual-stack backbone is a core network with all routers configured to support dual-stacks. It essentially consists of two network types existing side by side. The IPv4 stack routes the IPv4 traffic through the IPv4 network. The IPv6 stack routes the IPv6 traffic through the IPv6 network. This is a very basic approach to routing both IPv4 and IPv6 traffic through a network.

Dependencies

Each site has the appropriate entries in a DNS to resolve both IPv4 and IPv6 names and IP addresses.

Advantages

This is a basic and simple strategy for routing IPv4 and IPv6 traffic in a network.

Considerations

All routers in the network require a software upgrade to support dual-stack. Having dual-stack requires additional router management of a dual addressing scheme and additional router memory.

Dual-mode stack

A dual-mode stack is a stack configured to support both IPv4 and IPv6 protocols. It is a single stack (not two stacks) configured to support IPv4 and IPv6 simultaneously. Both IPv4 and IPv6 interfaces are capable of receiving and sending IPv4 and IPv6 packets over corresponding interfaces.

Dependencies

A z/OS Communications Server that is configured to support IPv6 requires OSA-Express ports to be running in QDIO mode.

Advantages

There are no additional software or hardware requirements for users in a z/OS environment configured with OSA-Express features. Dual-mode allows IPv4 and IPv6 applications to coexist indefinitely. However, any application can be migrated one at a time or at the user's convenience from IPv4 to IPv6. This is therefore an inexpensive, low risk, low impact deployment strategy.

Considerations

The only link layer protocol that supports IPv6 is MPC+. The devices that use the MPC+ protocol are XCF, MPCPTP, and MPCIPA (for example, OSA-Express3 in QDIO mode and HiperSockets on the System z196).

Recommendation

Using dual-mode stacks is the recommended strategy for application migration from IPv4 to IPv6.

How IPv6 is implemented in z/OS Communications Server

IPv6 is implemented in the z/OS Communications Server through a series of configuration tasks. We configure the stack to support IPv6 in a similar fashion to the steps performed for IPv4 configuration. However, before you start to configure the stack to support IPv6 traffic, you need to understand a few things about IPv6.

IPv6 addressing

An IPv6 address is a 128-bit number written in colon hexadecimal notation. This scheme is hexadecimal and consists of eight 16-bit pieces of the address.

Alternate notations described in RFC 2373 are acceptable; for example:

FEDC:BA98:7654:3210:FEDC:BA98:7654:3210

The following conventional forms represent IPv6 addresses as text strings:

- The preferred form is xxxxxxxx, where the x's indicate the hexadecimal value of the eight 16-bit pieces of the address, for example:

FEDC:BA98:7654:3210:FEDC:BA98:7654:3210
1080:0:0:0:8:800:200C:417A

Note: It is not necessary to write the leading zeros in an individual field, but there must be at least one numeral in every field, *except* for the case described in the next list item.

- ▶ Due to some methods of allocating certain styles of IPv6 addresses, it is common for addresses to contain long strings of zero bits. To simplify the writing of addresses that contain zero bits, a special syntax is available to compress the zeros. The use of :: indicates multiple groups of 16 bits of zeros. The :: can appear only *once* in an address. The :: can also be used to compress the leading or trailing zeros in an address.

Consider the following addresses:

1080:0:0:0:8:800:200C:417A (unicast address)
FF01:0:0:0:0:0:101 (multicast address)
0:0:0:0:0:0:0:1 (loopback address)
0:0:0:0:0:0:0:0 (unspecified addresses)

They can be represented as:

1080::8:800:200C:417A (unicast address)
FF01::101 (multicast address)
::1 (loopback address)
:: (unspecified addresses)

- ▶ An alternative form that is sometimes more convenient to use when dealing with a mixed environment of IPv4 and IPv6 nodes is to use x:x:x:x:x:d.d.d.d. Here, the x's are the hexadecimal values of the six high-order 16-bit pieces of the address. The d's are the decimal values of the four low-order 8-bit pieces of the address (standard IPv4 representation).

Consider this example:

0:0:0:0:0:0:13.1.68.3
0:0:0:0:0:0:FFFF:129.144.52.38

In compressed form, it is written as:

::13.1.68.3
::FFFF:129.144.52.38

Stateless address autoconfiguration

IPv6 addresses can be manually defined or *autoconfigured*.

With minimal router configuration and no manual configuration of local addresses, a host can generate its own IPv6 addresses. An IPv6 public autoconfigured address is the combination of a router advertised prefix and the interface ID provided by the OSA-Express QDIO adapter or manually configured using the INTFID parameter on the INTERFACE statement.

Routers advertise prefixes that identify the subnets associated with a LAN. In the absence of routers or manual configuration, a host can only generate link local addresses. However, link local addresses are sufficient for allowing communication among nodes attached to the same LAN.

Defining or adding an address on the INTERFACE statement indicates that stateless autoconfiguration is not wanted. Only a link local address is generated. See Figure A-1 on page 391 to examine the layout of a link local address.

"Prefix" or "Full Routing Prefix"			
Format Prefix Scope Prefix			
Unicast	0	9	63
	10 bits	54 bits	64 bits Interface ID
Link-local Scope	1111 1110 10 FE80	0...0	MAC, Other Interface ID
	0	9	63
	10 bits	54 bits	64 bits Interface ID
Site-local Scope (deprecated)	1111 1110 11 FEC0	0...0	MAC, Other Interface ID
	0	4	63
Global Scope RFC 2373 vs. RFC 3513	3 bits	61 bits	64 bits Interface ID
	001 (anything else)	variable "subnet"	MAC, Other Interface ID
	0	4	63
	3 bits	61 bits	64 bits Interface ID
	001 (anything else)	variable "subnet"	MAC, Other Interface ID
	0	4	63

Figure A-1 Unicast IPv6 Addressing Formats

Guidelines for IP filter rules and security associations with IPv6

If you use autoconfiguration, your IPv6 addresses might not be predictable. To configure IP filter rules for dynamic security associations with autoconfigured IPv6 addresses, you need to specify the IP addresses using wild cards.

Manual security associations typically use specific IP addresses for the endpoints. You can use wild cards for the security endpoint addresses so that the data endpoints and security endpoints are considered identical. Alternatively, you can use predictable IPv6 addresses for the security endpoints. You can obtain predictable IPv6 addresses by configuring full 128-bit IPv6 addresses on your INTERFACE statements by specifying the INTFID keyword on your INTERFACE statements or by using VIPAs.

Security with IPv6 autoconfiguration

RFC 4941, Privacy Extensions for Stateless Address Autoconfiguration in IPv6, addresses a potential security concern that can arise with the use of stateless address autoconfiguration. The static interface ID in an autoconfigured address makes it possible to correlate independent transactions to and from the system using the adapter even if the overall IPv6 address changes.

RFC 4941 defines a mechanism to generate a random interface ID that changes over time, thus eliminating the potential security exposure caused by the usual predictability of the interface ID. This random interface ID can be used in place of the static interface ID in generating temporary autoconfigured addresses. Based on new configuration parameters, the random interface ID and temporary addresses are regenerated periodically.

Temporary autoconfigured addresses have the same characteristics as public autoconfigured addresses. The address is generated as the result of a received router advertisement. The

address is deprecated at the end of its preferred lifetime. The address is deleted at the end of its valid lifetime.

Temporary autoconfigured addresses are designed to be used by short-lived client applications to make it more difficult to correlate activity. Temporary addresses should not be used for a server because the server needs to have a known IP address (or DNS name) so that clients can reach it. Temporary addresses should not be used with a long-lived client connection because the connection can become unusable if its source IP address is deleted while the connection is active.

To enable temporary address support for a TCP/IP stack, specify TEMPADDRS on the IPCONFIG6 statement in the TCP/IP profile. TEMPPREFIX on an interface definition specifies the set of prefixes for which temporary IPv6 addresses can be generated.

IPv6 TCP/IP Network part (prefix)

Designers have defined some address types, known as “address scopes”, and have left room for future definitions, because unknown requirements might arise. RFC 2373, IP version 6 Addressing Architecture (July 1998), defines the current addressing scheme. Figure A-1 on page 391 shows the layout of three types or “scopes” of addresses:

- ▶ The Link-local scope
- ▶ The Site-local scope (now deprecated)
- ▶ The Global scope

Each address begins with a format or scope prefix of 10 bits, followed by a second field and then an interface identifier field. Each of these addresses serves a unique purpose:

- ▶ Link local scope

These are special addresses that are only valid on a *link* of an interface. Using this address as the destination, the packet never passes through a router. A packet with a link-local source or destination address will not leave its originating LAN. A router receiving the packet will not forward it onto another physical LAN. An address of this type bears the prefix of fe80.

A link-local address is assigned to each IPv6-enabled interface after stateless auto-configuration, commonly used in IPv6 implementations. The link-local address is used for link communications such as:

- Neighbor discovery, that is, discovering whether there is anyone else on this link
- Communication with a neighbor when a router is unnecessary

Consider the example shown in Figure A-2 on page 393.

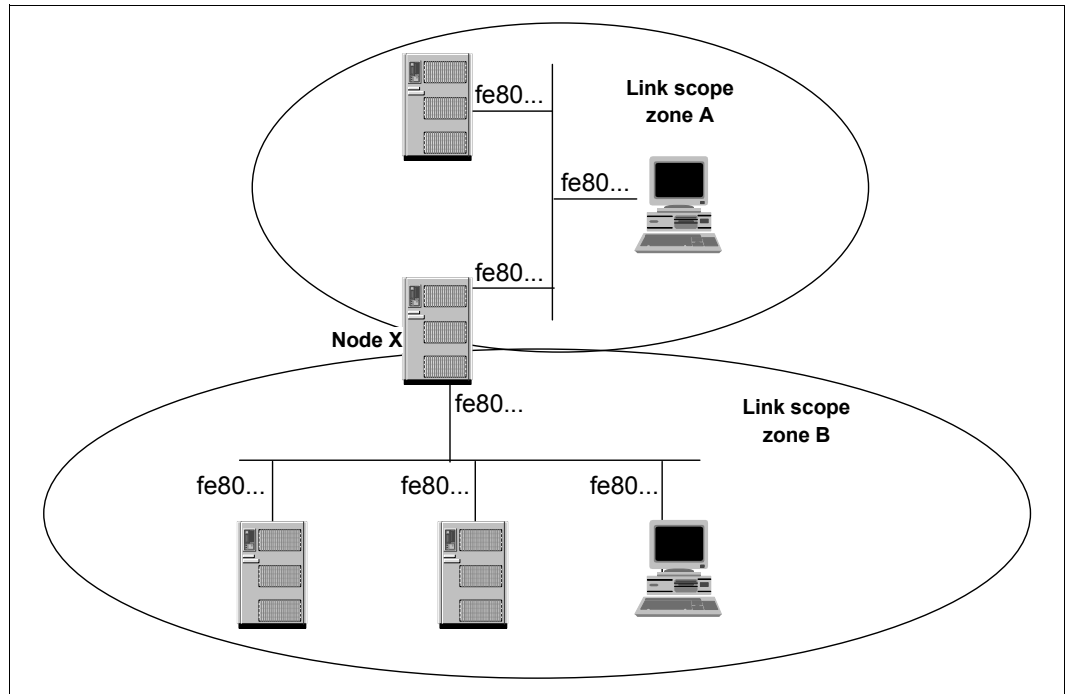


Figure A-2 Link-local addresses and link-local scope zones

Figure A-2 shows a LAN environment separated into two LAN segments, which are represented by Link scope zone A with three nodes and Link scope zone B with four nodes. The link local addresses in each zone begin with the prefix “fe80”.

Within a zone, nodes communicate with each other using link-local addresses. Across zones, nodes must communicate with each other using global scope addresses, which are discussed later.

Note that Node X has link-local addresses in two zones: in zone A and in zone B. Because link-local addresses use the same prefix value, it is necessary to understand which zone a packet should be sent to, particularly when a default route is to be used. So if a route exists on Node X for any destination address with a prefix of “fe80”, then the routing table needs to be able to distinguish between “fe80” in zone A and “fe80” in zone B.

Therefore, both the address and the zone index value need to be specified in the routing table. The *zone index* is a value assigned by the stack to represent the correct entry (or interface) in the routing table. If the zone index is not present, then the stack uses the “default route” for this configuration.

If the default route uses the interface that matches the IPv6 link-local address that was specified, everything works just fine. If, however, the default route does *not* use the correct interface for the specified IPv6 link-local address, then a routing error is encountered and the application request fails or times out. So the zone index helps the stack to distinguish whether the routing path should flow into zone A or into zone B.

z/OS Communications Server supports scope zone information about Getaddrinfo and Getnameinfo invocations, and also on the z/OS Socket APIs that support IPv6, thus satisfying requirements for IPv6 compliance. In addition, scope zone information can be included on command-line operations and in configuration files for **ftp**, **ping**, **traceroute**, **rexec**, **orexec**, **rsh**, and **orsh**.

- Site local address type

These addresses are now deprecated, that is, no longer recommended by the IETF. Use and deployment difficulties caused by the use of such addresses has led the IETF to discourage their use.

Originally, site-local addresses were used to communicate across routers or zones within the same intranet. They were similar to the RFC 1918/Address Allocation for Private Internets in IPv4 today, such as the address ranges represented by 10.0.0.0/8, 172.16.0.0/16 - 172.31.0.0/16, and 192.168.0.0/24. Since their deprecation, they are treated as global unicast addresses.

This deprecated address scope begins with the following prefixes:

fec0: (most commonly used)
fed0:
fee0:

- 6bone test addresses

These addresses were the first global addresses defined and used for testing purposes. They all begin with the following prefix:

3ffe:

- 6to4 addresses

These addresses were designed for a special tunneling mechanism (RFC 3056/Connection of IPv6 Domains using IPv4 Clouds and RFC 2893/Transition Mechanisms for IPv6 Hosts and Routers). They encode a given IPv4 address and a possible subnet. They begin with the following prefix:

2002:

Consider this example, representing 192.168.1.1/5:

2002:c0a8:0101:5::1

- Assigned by a provider for hierarchical routing

These addresses are delegated to Internet service providers (ISP) and begin with the following prefix:

2001:

- Multicast addresses

Multicast addresses are used for related services and always begin with the prefix *ffxx::*. Here, *xx* is the scope value.

- Anycast addresses

Anycast addresses are special addresses used to cover such items as the nearest Domain Name System (DNS) server, nearest Dynamic Host Configuration Protocol (DHCP) server, or similar dynamic groups. Addresses are taken out of the unicast address space, and can be aggregated globally or site-local at the moment. The anycast mechanism (client view) is handled by dynamic routing protocols.

Note: Anycast addresses cannot be used as source addresses. They are used only as destination addresses.

A simple example of an anycast address is the *subnet-router anycast address*. Assuming that a node has the following global assigned IPv6 address:

3ffe:ffff:100:f101:210:a4ff:fee3:9566/64

The subnet-router anycast address is created by blanking the suffix (least significant 64 bits) completely:

3ffe:ffff:100:f101::/64

IPv6 implementation in z/OS

The z/OS Communications Server provides support for both IPv4 and IPv6 protocols to coexist. In this book we review how this strategy can be implemented in a z/OS networking environment.

Further details about configuration options not referenced here are available in *z/OS Communications Server: IP Configuration Reference*, SC31-8776 and *z/OS Communications Server: IPv6 Network and Application Design Guide*, SC31-8885.

Table A-1 summarizes the z/OS TCP/IP stack-related functions and the level of support, based on the current release of the z/OS Communications Server. You can use this table to determine whether a given function is applicable and supported.

Table A-1 z/OS TCP/IP stack function support

z/OS TCP/IP stack function	IPv4 support	IPv6 support	Comments
Link-layer device support			
OSA-Express in QDIO mode	Y	Y	Related configuration statements: <ul style="list-style-type: none"> ▶ DEVICE MPCIPA and LINK IPAQENET ▶ INTERFACE IPAQENET ▶ INTERFACE IPAQENET6
CTC	Y	N	
LCS	Y	N	
CLAW	Y	N	
CDLC (3745/3746)	Y	N	
SNALINK LU0 and LU6.2	Y	N	
X.25 NPSI	Y	N	
NSC HyperChannel	Y	N	
MPC Point-Point	Y	Y	Related configuration statements: <ul style="list-style-type: none"> ▶ INTERFACE MPCPTP6
ATM	Y	N	
HiperSockets	Y	Y	Related configuration statements: <ul style="list-style-type: none"> ▶ INTERFACE IPAQIDIO6 ▶ IPCONFIG6 DYNAMICXCF
XCF	Y	Y	Related configuration statements: <ul style="list-style-type: none"> ▶ INTERFACE MPCPTP6 ▶ IPCONFIG6 DYNAMICXCF
Virtual IP addressing support			
Virtual Device/Interface Configuration for static VIPA	Y	Y	Related configuration statements: <ul style="list-style-type: none"> ▶ DEVICE and LINK VIRTUAL ▶ INTERFACE VIRTUAL6

z/OS TCP/IP stack function	IPv4 support	IPv6 support	Comments
Sysplex support			
Sysplex distributor integration with Cisco MNLB	Y	N	
Sysplex Wide Security Associations (SWSA)	Y	N	
IP routing functions			
Dynamic routing - OSPF and RIP	Y	Y	OMPROUTE supports OSPFv3 and RIPng.
Multipath Routing Groups	Y	Y	Related configuration statements: ► IPCONFIG ► IPCONFIG6
Policy-based Routing	Y	N	
Static Route Configuration by way of BEGINROUTES statement	Y	Y	Related configuration statement: ► BEGINROUTES
Static Route Configuration by way of GATEWAY statement	Y	N	Related configuration statement: ► GATEWAY
Miscellaneous IP/IF-layer functions			
Path MTU discovery	Y	Y	Path MTU discovery is mandatory in IPv6. Related configuration statements: ► IPCONFIG ► IPCONFIG6
Configurable Device or Interface Recovery Interval	Y	Y	
Link-Layer Address Resolution	Y	Y	In IPv4, this is performed using Address Resolution Protocol (ARP). In IPv6, this is performed using the neighbor discovery protocol. Related configuration statements: DEVICE and LINK (LAN Channel Station and OSA devices) INTERFACE (IPAQENET6 interfaces)
ARP/Neighbor Cache PURGE Capability	Y	Y	Use the V TCPIP,PURGECACHE command. For information see <i>z/OS Communications Server: IP System Administrator's Commands</i> , SC31-8781.
Datagram Forwarding Enable/Disable	Y	Y	Related configuration statements: ► IPCONFIG ► IPCONFIG6
Hipersockets Accelerator	Y	N	Related configuration statements: ► IPCONFIG QDIOROUTING
QDIO Accelerator	Y	N	Related configuration statements: ► IPCONFIG QDIOACCELERATOR
Checksum offload	Y	N	
Segmentation offload	Y	N	

z/OS TCP/IP stack function	IPv4 support	IPv6 support	Comments
QDIO inbound workload queueing	Y	N	
Transport-layer functions			
Fast Response Cache Accelerator	Y	N	
Enterprise Extender	Y	Y	IPv6 Enterprise Extender support requires a virtual IP address and IUTSAMEH. Related configuration statements: <ul style="list-style-type: none"> ▶ INTERFACE VIRTUAL6 and MPCPTP6 ▶ IPCONFIG6 DYNAMICXCF
Server-BIND control	Y	Y	Related configuration statement: <ul style="list-style-type: none"> ▶ PORT
UDP Checksum Disablement Option	Y	N	UDP checksum is required when operating over IPv6. Related configuration statement: <ul style="list-style-type: none"> ▶ UDPCONFIG
Network management and accounting functions			
SNMP	Y	Y	SNMP applications can communicate over an IPv6 connection. IPv6 management data includes added support for the version-neutral (both IPv4 and IPv6) MIB data in the following new, IETF Internet drafts: <ul style="list-style-type: none"> ▶ IP-MIB: draft-ietf-ipv6-rfc2011-update-01.txt ▶ IP-FORWARD-MIB: draft-ietf-ipv6-rfc2096-update-02.txt ▶ TCP-MIB: draft-ietf-ipv6-rfc2012-update-01.txt
SNMP agent	Y	Y	
TCP/IP subagent	Y	Y	No IPv6 UDP support
Network SLAPM2 subagent	Y	Y	
Distributed Protocol Interface	Y	Y	
OMPROUTE subagent	Y	N	
Trap forwarder daemon	Y	Y	
Policy-Based Networking	Y	Y	IPv6 support in Policy Agent: <ul style="list-style-type: none"> ▶ IPv6 source and destination IP addresses are allowed to be specified in policy rules (LDAP and configuration files). ▶ Interfaces in policy rules and subnet priority TOS masks are allowed to be specified by name. <ul style="list-style-type: none"> – Allowed for both IPv4 and IPv6 interfaces – IPv6 interfaces <i>must</i> be specified by name ▶ TOS in policy definitions means IPv4 Type of Service or IPv6 Traffic Class.
SMF	Y	Y	Related configuration statement: <ul style="list-style-type: none"> ▶ SMFCONFIG
TN3270 subagent	Y	Y	

z/OS TCP/IP stack function	IPv4 support	IPv6 support	Comments
Security function			
IPSec	Y	Y	Related configuration statements: <ul style="list-style-type: none"> ▸ IPCONFIG ▸ IPCONFIG6
IP filtering	Y	Y	
IKE daemon	Y	Y	
NAT traversal	Y	N	
Network Access Control	Y	Y	Related configuration statement: <ul style="list-style-type: none"> ▸ NETACCESS
Stack and Port Access Control	Y	Y	Related configuration statements: <ul style="list-style-type: none"> ▸ PORT ▸ DELETE
Application Transparent TLS	Y	Y	
Intrusion Detection Services	Y	N	
Server applications			
Rpcbind server	Y	Y	IPv6-enabled RPC applications require a Rpcbind server. The following RPC facilities are not IPv6 enabled, and they do not support RPC binding protocols Version 3 and Version 4: <ul style="list-style-type: none"> ▸ rpcgen and orpcgen ▸ rpcinfo and orpcinfo ▸ RPC library for the z/OS Communications Server environment ▸ RPC library for the z/OS UNIX System Services environment For more information see <i>z/OS Communications Server: IP Configuration Guide</i> , SC31-8775

Based on the discussion and recommendation in “Common design scenarios for IPv6” on page 387, here we concentrate on a single stack environment running in dual-mode. A single stack environment is one TCP/IP stack running in an LPAR.

Dual-mode stack

As previously discussed, a TCP/IP stack that supports both IPv4 and IPv6 interfaces and is capable of receiving and sending IPv4 and IPv6 packets over the corresponding interfaces is referred to as a dual-mode stack. A dual-mode stack is a single stack supporting IPv4 and IPv6 protocols, which is different from dual-stack mode that uses two TCP/IP stacks running side by side, each supporting only one of the protocols (either IPv4 or IPv6).

The z/OS Communications Server can be configured to support an IPv4-only stack or a dual-mode stack (IPv4 and IPv6). There is no support for an IPv6-only stack. By default, IPv6-enabled applications can communicate with both IPv4 and IPv6 peers in a dual-mode environment.

A z/OS dual-mode stack is enabled when both AF_INET and AF_INET6 are coded in SYS1.PARMLIB(BPXPRMxx). You cannot code AF_INET6 without specifying AF_INET, and doing so will cause the TCP/IP stack initialization to fail.

Note that AF_INET6 support can be dynamically enabled by configuring AF_INET6 in BPXPRMxx and then issuing the SETOMVS RESE= command to activate the new configuration.

IPv6 application on a dual-mode stack

An IPv6 application on a dual-mode stack can communicate with IPv4 and IPv6 partners as long as it does *not* bind to a native IPv6 address. If it binds to a native IPv6 address, then the native IPv6 address cannot be converted to an IPv4 address.

If a partner is IPv6, then all communication will use IPv6 packets.

If a partner is IPv4, then the following will occur:

- ▶ Both source and destination will be IPv4-mapped IPv6 addresses.
- ▶ On inbound, the transport protocol layer will map the IPv4 address to its corresponding IPv4-mapped IPv6 address before returning to the application with AF_INET6 addresses.
- ▶ On outbound, the transport protocol layer will convert the IPv4-mapped address to the native IPv4 addresses and send IPv4 packets.

IPv4 application on a dual-mode stack

An IPv4 application running on a dual-mode stack can communicate with an IPv4 partner. The source and destination addresses will be native IPv4 addresses and the packet will be an IPv4 packet.

If a partner is IPv6 enabled and running on an IPv6-only stack, then communication will fail. The partner only has a native IPv6 address (not an IPv4-mapped IPv6 address). The native IPv6 address for the partner cannot be converted into a form that the AF_INET application will understand.

Older AF_INET applications are only able to communicate using IPv4 addresses. IPv6-enabled applications that use AF_INET6 sockets can communicate using both IPv4 and IPv6 addresses (on a dual-mode stack). AF_INET and AF_INET6 applications can thus communicate with one another, but only using IPv4 addresses.

If the socket libraries on the IPv6-enabled host are updated to support IPv6 sockets (AF_INET6), applications can be IPv6 enabled. When an application on a dual mode stack is IPv6 enabled, the application is able to communicate with both IPv4 and IPv6 partners. This is true for both clients and server on a dual-mode stack.

IPv6-enabling both sockets libraries and applications on dual-mode stack therefore becomes a migration concern. As soon as IPv6-only hosts are being deployed in a network, applications on those IPv6-only partners cannot communicate with the IPv4-only applications on the dual mode hosts, unless one of the multiple migration technologies is implemented either on intermediate nodes in the network or directly on the dual mode hosts.

Table A-2 summarizes the application communication rules when running in dual-mode.

Table A-2 Dual-mode communication

Partner	Application communication on a dual-mode TCP/IP stack	
	IPv4 only	IPv6 enabled
IPv4-only	Yes	Yes
IPv6-only	No	Yes

Figure A-3 depicts a dual-mode stack, which is the IPv6 configuration we implemented in our networking environment. The following sections walk you through the setup.

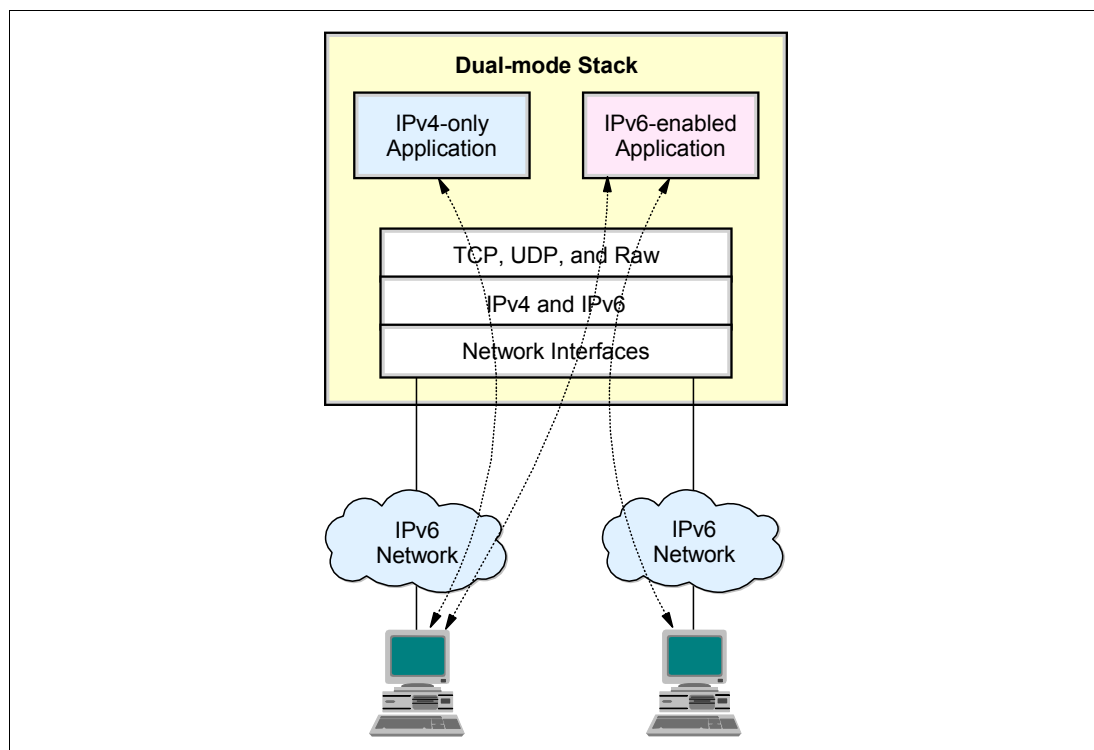


Figure A-3 Dual-mode TCP/IP stack

Implementation tasks for a dual-mode stack

To implement a dual-mode stack in our networking environment, we modified the following:

- ▶ BPXPRMxx definitions
- ▶ VTAM definitions
- ▶ TCP/IP definitions

BPXPRMxx definitions

IPv6 is not enabled, by default. You must specify a NETWORK statement with AF_INET6 in your BPXPRMxx member.

To support our dual-mode stack (IPv4 and IPv6), we added the NETWORK statement, as shown in Example A-1, to our BPXPRMxx member.

Example A-1 BPXPRMxx NETWORK statement

```

NETWORK DOMAINNAME (AF_INET6)
DOMAINNUMBER(19)
MAXSOCKETS(2000)
TYPE(INET)

```

The TYPE option in our case is INET, because we used a single stack.

Note: The BPXPRMxx member can be updated dynamically using the z/OS command SETOMVS RESET=(xx). After the reset, we received the message BPXF203I DOMAIN AF_INET6 WAS SUCCESSFULLY ACTIVATED. We then *recycled* the TCP/IP stack to pick up the change.

For more details about the definitions required in BPXPRMxx to provide a dual- stack, refer to *z/OS Communications Server: IP Configuration Guide*, SC31-8775.

VTAM definitions

As previously mentioned, one of the protocols that z/OS Communications Server TCP/IP supports is MPC+, and the MPC+ protocols are used to define the DLCs for OSA-Express devices in QDIO. OSA-Express QDIO connections are configured through a TRLE definition. Because VTAM provides the DLCs for TCP/IP, all TRLEs are defined as VTAM major nodes (see Example A-2).

Example A-2 TRLE definition

```
OSA2080  VBUILD TYPE=TRL
OSA2080T  TRLE  LNCTL=MPC,
              READ=2080,
              WRITE=2081,
              DATAPATH=(2082-2087),
              PORTNAME=OSA2080,  1
              MPCLEVEL=QDIO
```

The PORTNAME 1 is identical to the device name defined in the TCP/IP PROFILE data set on the INTERFACE statement.

TCP/IP definitions

We added one INTERFACE statement for the OSA-Express3 1000BASE-T port to support IPv6. This statement merges the DEVICE, LINK, and HOME definitions into a single statement. Several different parameters are associated with the INTERFACE statement. To determine which of them best fits your requirements, refer to *z/OS Communications Server: IP Configuration Reference*, SC31-8776.

We used the following syntax:

```
INTERFace interfname DEFINE linktype PORTNAME portname IPADDR ipaddr
```

The syntax is explained here:

- ▶ *interfname* specifies a name for the interface with no more than 16 characters in length.
- ▶ *linktype* must be IPAQENET6, which is the only DLC that currently supports IPv6.
- ▶ *portname* is specified in the VTAM TRLE definition for the QDIO interface.
- ▶ *ipaddr* is optional for link type IPAQENET6. If not specified, TCP/IP enables auto-configuration for the interface. If used, one or more prefixes or full IPv6 addresses can be specified.

Note: To configure a single physical device for both IPv4 and IPv6 traffic, you must use DEVICE/LINK/HOME for the IPv4 definition and INTERFACE for the IPv6 definition, so that the PORTNAME value on the INTERFACE statement matches the device name on the DEVICE statement.

The TCP/IP IPv6 PRFOFILE for a single stack is illustrated in Example A-3. The INTERFACE statement defines the configuration of the OSA-Express device (OSA2080) that we used for network connectivity. The PORTNAME must be identical to the PORTNAME defined in the TRLE. The TRLE is defined as a VTAM major node in the VTAM definition data set.

Example A-3 shows the TCP/IP profile for our environment, using SYSTEM SYMBOLS and INCLUDE statements. The &SYSCONE that you see throughout the example will result in a two-digit value (30 in our example, for system SC30) being inserted. By doing this we can use the same profile for each of several systems, each time translating to the appropriate system value (systems 30, 31, and 32). The &SYSCONE value is defined in SYS1.PARMLIB.

Example A-3 Profile definition with the use of SYSTEM SYMBOLS and INCLUDE

```

ARPAGE 20
;
GLOBALCONFIG NOTCPIPSTATISTICS
;
IPCONFIG NODATAGRAMFWD SOURCEVIPA      1
IPCONFIG6 NODATAGRAMFWD SOURCEVIPA     2
;
SOMAXCONN 240
;
TCPCONFIG TCPSENDBFRSIZE 64K TCPRCVBUFRSIZE 64K SENDGARBAGE FALSE
TCPCONFIG TCPMAXRCVBUFRSIZE 256K
TCPCONFIG RESTRICTLOWPORTS
;
UDPCONFIG RESTRICTLOWPORTS
;
INCLUDE TCPIPE.TCPPARMS(HOME&SYSCONE.V6)
INCLUDE TCPIPE.TCPPARMS(STAT&SYSCONE.V6)
;
AUTOLOG 5
      FTPDE&SYSCONE JOBNAME FTPDE&SYSCONE.1
ENDAUTOLOG
;
PORT
    20 TCP * NOAUTOLOG      ; FTP Server
    21 TCP FTPDE&SYSCONE.1  ; control port
    23 TCP TN3270XE NOAUTOLOG ; MVS Telnet Server
    23 TCP OMVS             ; Telnet Server
    25 TCP SMTP             ; SMTP Server
    514 UDP OMVS            ; UNIX Syslogd daemon
;
SACONFIG ENABLED COMMUNITY public AGENT 161
;
SMFCONFIG
    FTPCLIENT TN3270CLIENT
    TYPE119 FTPCLIENT TN3270CLIENT
;

```

In this example, the numbers correspond to the following information:

- 1** Defines the IPv4 environment.
- 2** Defines the IPv6 environment.

Example A-4 shows the DEVICE, LINK, HOME, INTERFACE, and IPADDR definitions we used to support IPv4 and IPv6 and their addressing schemes.

Example A-4 Interface and address definitions

```

DEVICE OSA2080      MPCIPA                                1
LINK   OSA2080LNK  IPAQENET  OSA2080 VLANID 12
INTERFACE LNK62080 DEFINE IPAQENET6 PORTNAME OSA2080      1
IPADDR FEC0:0:0:1::3302
        FEC0:0:0:1001::3302
;
DEVICE STAVIPA1     VIRTUAL 0
LINK   STAVIPA1LNK VIRTUAL 0      STAVIPA1
;
HOME
    192.168.1.10    STAVIPA1LNK
    192.168.2.10    OSA2080LNK
;
START OSA2080
START LNK62080

```

In this example, the numbers correspond to the following information:

1 Defines the same device (OSA2080) to support IPv4 and IPv6 addresses.

Example A-5 show static routes in a flat network (no dynamic routing protocol).

Example A-5 Static route definitions

```

BEGINRoutes
; Direct Routes - Routes that are directly connected to my interfaces
;   Destination   Subnet Mask   First Hop   Link Name   Packet Size
ROUTE FEC0::0/10      =                LNK62080    MTU 1492
ROUTE 192.168.2.0    255.255.255.0 =                OSA2080LNK MTU 1492
; Default Route - All packets to an unknown destination are routed
;through this route.
;   Destination           First Hop   Link Name   Packet Size
ROUTE DEFAULT            192.168.2.240 OSA2080LNK MTU 1492
ENDRoutes

```

The messages shown in Example A-6 were written to the z/OS console when the TCP/IP stack of TCPPIPE was initializing on SC30. We also manually started our external TN3270E server (TN3270XE).

Example A-6 TCP/IP stack and TN3270E server initializations

```

S TCPPIPE
$HASP100 TCPPIPE  ON STCINRDR
IEF695I START TCPPIPE  WITH JOBNAME TCPPIPE  IS ASSIGNED TO USER
TCPPIP  , GROUP TCPGRP
$HASP373 TCPPIPE  STARTED
IEE252I MEMBER CTIEZB00 FOUND IN SYS1.IBM.PARMLIB
IEE252I MEMBER CTIIDS00 FOUND IN SYS1.IBM.PARMLIB
IEE252I MEMBER CTINTA00 FOUND IN SYS1.PARMLIB
EZZ7450I FFST SUBSYSTEM IS NOT INSTALLED
EZZ0162I HOST NAME FOR TCPPIPE IS WTSC30E
EZZ0300I OPENED INCLUDE FILE 'TCPPIPE.TCPPARMS(HOME30V6) '

```

```

EZZ0300I OPENED INCLUDE FILE 'TCPIPE.TCPPARMS(STAT30V6)'
EZZ0300I OPENED PROFILE FILE DD:PROFILE
EZZ0309I PROFILE PROCESSING BEGINNING FOR DD:PROFILE
EZZ0309I PROFILE PROCESSING BEGINNING FOR TCPIPE.TCPPARMS(HOME30V6)
EZZ0316I PROFILE PROCESSING COMPLETE FOR FILE 'TCPIPE.TCPPARMS(HOME30V
6) '
EZZ0304I RESUMING PROCESSING OF FILE DD:PROFILE
EZZ0309I PROFILE PROCESSING BEGINNING FOR TCPIPE.TCPPARMS(STAT30V6)
EZZ0316I PROFILE PROCESSING COMPLETE FOR FILE 'TCPIPE.TCPPARMS(STAT30V
6) '
EZZ0304I RESUMING PROCESSING OF FILE DD:PROFILE
EZZ0316I PROFILE PROCESSING COMPLETE FOR FILE DD:PROFILE
IEF196I IEF237I 2084 ALLOCATED TO TP2084
EZZ0334I IP FORWARDING IS DISABLED
EZZ0351I SOURCEVIPA SUPPORT IS ENABLED
EZZ0699I IPV6 FORWARDING IS DISABLED
EZZ0702I IPV6 SOURCEVIPA SUPPORT IS ENABLED
EZZ4313I INITIALIZATION COMPLETE FOR DEVICE OSA2080
EZZ0338I TCP PORTS 1 THRU 1023 ARE RESERVED
EZZ0338I UDP PORTS 1 THRU 1023 ARE RESERVED
EZZ0613I TCPIPSTATISTICS IS DISABLED
EZZ4248E TCPIPE WAITING FOR PAGENT TTLS POLICY
EZZ4202I Z/OS UNIX - TCP/IP CONNECTION ESTABLISHED FOR TCPIPE
BPXF206I ROUTING INFORMATION FOR TRANSPORT DRIVER TCPIPE HAS BEEN
INITIALIZED OR UPDATED.
EVB6473I TCP/IP STACK FUNCTIONS INITIALIZATION COMPLETE.
EZAIN11I ALL TCPIP SERVICES FOR PROC TCPIPE ARE AVAILABLE.
EZZ4340I INITIALIZATION COMPLETE FOR INTERFACE LNK62080
EZD1176I TCPIPE HAS SUCCESSFULLY JOINED THE TCP/IP SYSPLEX GROUP
EZBTCPCS
S FTPDE30
$HASP100 FTPDE30 ON STCINRDR
IEF695I START FTPDE30 WITH JOBNAME FTPDE30 IS ASSIGNED TO USER
TCPIP , GROUP TCPGRP
$HASP373 FTPDE30 STARTED
$HASP395 FTPDE30 ENDED

```

1
1

2

```

-----
We manually started our external TN3270E server.
-----
S TN3270XE
$HASP100 TN3270XE ON STCINRDR
IEF695I START TN3270XE WITH JOBNAME TN3270XE IS ASSIGNED TO USER
TCPIP , GROUP TCPGRP
$HASP373 TN3270XE STARTED
IEE252I MEMBER CTIEZBTN FOUND IN SYS1.IBM.PARMLIB
EZZ6001I TN3270XE SERVER STARTED
EZZ6044I TN3270XE PROFILE PROCESSING BEGINNING FOR FILE 897
TCPIPE.TCPPARMS(TN3270XE)
EZZ6045I TN3270XE PROFILE PROCESSING COMPLETE FOR FILE
TCPIPE.TCPPARMS(TN3270XE)
EZZ6003I TN3270XE LISTENING ON PORT 23

```

In this example, 1 and 2 indicate that IPv6 support is enabled and that the interface is initialized with IPv6 addresses.

Verification

Next, we verified our environment. Because the TRLE must be active before the interface is started, we ensured that the TRLE is in an active state. The results are shown in Example A-7. You can also verify the OSA-Express code level with this command.

Example A-7 OSA-Express status and code level

```
D NET,TRL,TRLE=OSA2080T
IST097I DISPLAY ACCEPTED
IST075I NAME = OSA2080T, TYPE = TRLE 071
IST1954I TRL MAJOR NODE = OSA2080
IST486I STATUS= ACTIV, DESIRED STATE= ACTIV 1
IST087I TYPE = LEASED, CONTROL = MPC, HPDT = YES
IST1715I MPCLEVEL = QDIO MPCUSAGE = SHARE
IST2263I PORTNAME = OSA2080 PORTNUM = 0 OSA CODE LEVEL = 000C 4
IST2337I CHPID TYPE = OSD CHPID = 02
IST1577I HEADER SIZE = 4096 DATA SIZE = 0 STORAGE = ***NA***
IST1221I WRITE DEV = 2081 STATUS = ACTIVE STATE = ONLINE 2
IST1577I HEADER SIZE = 4092 DATA SIZE = 0 STORAGE = ***NA***
IST1221I READ DEV = 2080 STATUS = ACTIVE STATE = ONLINE 2
IST924I -----
...
IST924I -----
IST1221I DATA DEV = 2084 STATUS = ACTIVE STATE = N/A 3
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST1717I ULPID = TCPIPE
IST2310I ACCELERATED ROUTING DISABLED
IST2331I QUEUE QUEUE READ
IST2332I ID TYPE STORAGE
IST2205I -----
IST2333I RD/1 PRIMARY 4.0M(64 SBALS)
IST2305I NUMBER OF DISCARDED INBOUND READ BUFFERS = 0
IST1757I PRIORITY1: UNCONGESTED PRIORITY2: UNCONGESTED
IST1757I PRIORITY3: UNCONGESTED PRIORITY4: UNCONGESTED
IST2190I DEVICEID PARAMETER FOR OSAENTA TRACE COMMAND = 01-01-00-04
IST1801I UNITS OF WORK FOR NCB AT ADDRESS X'25A28010'
IST1802I P1 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P2 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P3 CURRENT = 0 AVERAGE = 0 MAXIMUM = 0
IST1802I P4 CURRENT = 0 AVERAGE = 1 MAXIMUM = 2
IST924I -----
...
IST314I END
```

In this example, the numbers correspond to the following information:

- **1** indicates the state of the TRLE major node.
- **2** and **3** are the desired and required states.
- **4** indicates the OSA code level, which is a four-digit number that relates to a specific microcode engineering change (EC) and patch level (MCL).

Example A-8 on page 406 displays the HOME addresses after initialization.

Example A-8 HOME addresses displayed

```
D TCPIP,TCPIPE,N,HOME
EZD0101I NETSTAT CS V1R12 TCPIPE 078
HOME ADDRESS LIST:
LINKNAME: STAVIPA1LNK
ADDRESS: 192.168.1.10
FLAGS: PRIMARY
LINKNAME: OSA2080LNK
ADDRESS: 192.168.2.10      1
FLAGS:
LINKNAME: LOOPBACK
ADDRESS: 127.0.0.1
FLAGS:
INTFNAME: LNK62080
ADDRESS: FEC0:0:0:1::3302  2
TYPE: GLOBAL
FLAGS:
ADDRESS: FEC0:0:0:1001::3302 2
TYPE: GLOBAL
FLAGS:
ADDRESS: FE80::14:5E00:D77:6872 3
TYPE: LINK_LOCAL
FLAGS: AUTOCONFIGURED
INTFNAME: LOOPBACK6
ADDRESS: ::1              4
TYPE: LOOPBACK
FLAGS:
7 OF 7 RECORDS DISPLAYED
END OF THE REPORT
```

In this example, the numbers correspond to the following information:

- ▶ **1** This is the IPv4 address assigned to the OSA Express device (OSA2080).
- ▶ **2** These are the IPv6 addresses assigned to the same OSA Express device (OSA2080) defined with the INTERFACE statement. However, these SITE_LOCAL addresses are not generally recommended.
- ▶ **3** This is an auto-configured LINK_LOCAL address for the same OSA Express device.
- ▶ **4** This is the IPv6 Loopback address.

Example A-9 shows the output of NETSTAT DEV, with the IPv6 Loopback and device interfaces shown as READY.

Example A-9 Device display, using NETSTAT

```
D TCPIP,TCPIPE,N,DEV
EZD0101I NETSTAT CS V1R12 TCPIPE 083
DEVNAME: LOOPBACK      DEVTYPE: LOOPBACK
DEVSTATUS: READY
LNKNAME: LOOPBACK      LNKTYPE: LOOPBACK  LNKSTATUS: READY
ACTMTU: 65535
ROUTING PARAMETERS:
MTU SIZE: N/A          METRIC: 00
DESTADDR: 0.0.0.0      SUBNETMASK: 0.0.0.0
MULTICAST SPECIFIC:
MULTICAST CAPABILITY: NO
```

```

LINK STATISTICS:
  BYTESIN                      = 2286
  INBOUND PACKETS              = 36
  INBOUND PACKETS IN ERROR     = 0
  INBOUND PACKETS DISCARDED    = 0
  INBOUND PACKETS WITH NO PROTOCOL = 0
  BYTESOUT                    = 2286
  OUTBOUND PACKETS            = 36
  OUTBOUND PACKETS IN ERROR   = 0
  OUTBOUND PACKETS DISCARDED  = 0
INTFNAME: LOOPBACK6          INTFTYPE: LOOPBACK6  INTFSTATUS: READY  1
  ACTMTU: 65535
MULTICAST SPECIFIC:
  MULTICAST CAPABILITY: NO
INTERFACE STATISTICS:
  BYTESIN                      = 0
  INBOUND PACKETS              = 0
  INBOUND PACKETS IN ERROR     = 0
  INBOUND PACKETS DISCARDED    = 0
  INBOUND PACKETS WITH NO PROTOCOL = 0
  BYTESOUT                    = 0
  OUTBOUND PACKETS            = 0
  OUTBOUND PACKETS IN ERROR   = 0
  OUTBOUND PACKETS DISCARDED  = 0
DEVNAME: OSA2080             DEVTYPE: MPCIPA
DEVSTATUS: READY
LNKNAME: OSA2080LNK          LNKTYPE: IPAQENET   LNKSTATUS: READY  2
  SPEED: 0000001000
  IPBROADCASTCAPABILITY: NO
  CFGROUTER: NON              ACTROUTER: NON
  ARPOFFLOAD: YES              ARPOFFLOADINFO: YES
  ACTMTU: 8992
  VLANID: 12                   VLANPRIORITY: DISABLED
  DYNVLANREGCFG: NO            DYNVLANREGCAP: YES
  READSTORAGE: GLOBAL (4096K)
  INBPERF: BALANCED
  CHECKSUMOFFLOAD: YES
  SECCLASS: 255                MONSYSPLEX: NO
ROUTING PARAMETERS:
  MTU SIZE: N/A                METRIC: 00
  DESTADDR: 0.0.0.0            SUBNETMASK: 255.255.255.0
MULTICAST SPECIFIC:
  MULTICAST CAPABILITY: YES
  GROUP          REFCNT        SRCFLTMD
  -----
  224.0.0.1      0000000001    EXCLUDE
  SRCADDR: NONE
LINK STATISTICS:
  BYTESIN                      = 0
  INBOUND PACKETS              = 0
  INBOUND PACKETS IN ERROR     = 0
  INBOUND PACKETS DISCARDED    = 0
  INBOUND PACKETS WITH NO PROTOCOL = 0
  BYTESOUT                    = 0
  OUTBOUND PACKETS            = 0

```

```

OUTBOUND PACKETS IN ERROR          = 0
OUTBOUND PACKETS DISCARDED         = 0
INTFNAME: LNK62080                 INTFTYPE: IPAQENET6  INTFSTATUS: READY
PORTNAME: OSA2080                  DATAPATH: 2084      DATAPATHSTATUS: READY
CHIPIDTYPE: OSD
QUESIZE: 0      SPEED: 0000001000
MACADDRESS: 00145E776872
DUPADDRDET: 1
CFGROUTER: NON                     ACTROUTER: NON
CFGMTU: NONE                       ACTMTU: 9000
VLANID: NONE                       VLANPRIORITY: DISABLED
READSTORAGE: GLOBAL (4096K)
INBPERF: BALANCED
SECCLASS: 255                      MONSYSPLEX: NO
ISOLATE: NO                        OPTLATENCYMODE: NO
TEMPPREFIX: ALL
MULTICAST SPECIFIC:
MULTICAST CAPABILITY: YES
GROUP:      FF02::1:FF00:3302
  REFCNT: 0000000002 SRCFLTMD: EXCLUDE
  SRCADDR: NONE
GROUP:      FF02::1:FF77:6872
  REFCNT: 0000000001 SRCFLTMD: EXCLUDE
  SRCADDR: NONE
GROUP:      FF01::1
  REFCNT: 0000000001 SRCFLTMD: EXCLUDE
  SRCADDR: NONE
GROUP:      FF02::1
  REFCNT: 0000000001 SRCFLTMD: EXCLUDE
  SRCADDR: NONE
INTERFACE STATISTICS:
BYTESIN          = 0
INBOUND PACKETS  = 0
INBOUND PACKETS IN ERROR = 0
INBOUND PACKETS DISCARDED = 0
INBOUND PACKETS WITH NO PROTOCOL = 0
BYTESOUT         = 888
OUTBOUND PACKETS = 9
OUTBOUND PACKETS IN ERROR = 0
OUTBOUND PACKETS DISCARDED = 0
DEVNAME: STAVIPA1      DEVTYPE: VIPA
DEVSTATUS: READY
LNKNAME: STAVIPA1LNK   LNKTYPE: VIPA      LNKSTATUS: READY
ROUTING PARAMETERS:
  MTU SIZE: N/A        METRIC: 00
  DESTADDR: 0.0.0.0    SUBNETMASK: 255.255.255.0
MULTICAST SPECIFIC:
  MULTICAST CAPABILITY: NO
IPV4 LAN GROUP SUMMARY
LANGROUP: 00002
  NAME          STATUS      ARPOWNER      VIPAOWNER
  ----          -
  OSA2080LNK    ACTIVE      OSA2080LNK    YES
IPV6 LAN GROUP SUMMARY
LANGROUP: 00001

```


NAME	STATUS	NOWNER	VIPAOWNER
----	-----	-----	-----
LNK62080	ACTIVE	LNK62080	YES

OSA-EXPRESS NETWORK TRAFFIC ANALYZER INFORMATION:
 NO OSA-EXPRESS NETWORK TRAFFIC ANALYZER INTERFACES ARE DEFINED
 5 OF 5 RECORDS DISPLAYED
 END OF THE REPORT

If the device does not have a LNKSTATUS or INTFSTATUS of READY (as with **1**, **2**, and **3**), you must resolve this before you continue. There are several factors that might cause the LNKSTATUS or INTFSTATUS to not be READY. For example, the device cannot be varied online or defined to z/OS correctly, or the device cannot be defined in the TCP/IP profile correctly, and so on.

Example A-10 shows the FTP server **1** and the TN3270E server **2** bound to the IPv6 unspecified address (in6addr_any). They can now be accessed by another IPv6-enabled client across an IPv4 network.

Example A-10 Sockets in the stack

```
D TCPIP,TCPIPE,N,CONN
EZD0101I NETSTAT CS V1R12 TCPIPE 900
USER ID  CONN      STATE
FTPDE301 0000001F LISTEN          1
  LOCAL SOCKET:  ::..21
  FOREIGN SOCKET: ::..0
JES2S001 00000013 LISTEN
  LOCAL SOCKET:  ::..175
  FOREIGN SOCKET: ::..0
TN3270XE 00000027 LISTEN          2
  LOCAL SOCKET:  ::..23
  FOREIGN SOCKET: ::..0
TCPIPE   0000002C UDP
  LOCAL SOCKET:  ::..1030
  FOREIGN SOCKET: *.*
4 OF 4 RECORDS DISPLAYED
END OF THE REPORT
```

We used the TSO **ping** command to verify locally IPv4 and IPv6 interfaces (see Example A-11).

Example A-11 Results of the tso ping command

```
ping fe80::14:5e00:d77:6872
CS V1R12: Pinging host FE80::14:5E00:D77:6872
Ping #1 response took 0.000 seconds.

ping FEC0:0:0:1001::3302
CS V1R12: Pinging host FEC0:0:0:1001::3302
Ping #1 response took 0.000 seconds.

ping 192.168.2.10
CS V1R12: Pinging host 192.168.2.10
Ping #1 response took 0.000 seconds.
```




Additional parameters and functions

This appendix contains examples and a discussion of the following topics:

- ▶ MVS System symbols
- ▶ Reusable Address Space ID (REUSASID) function
- ▶ PROFILE.TCPIP parameters
- ▶ TCP/IP built-in security functions

MVS System symbols

One of the many strengths of the z/OS technology is that it allows multiple TCP/IP stacks (instances) to be configured in the same MVS system or across multiple MVS systems. If you need to run many stacks, you need to ensure that each profile configuration data set is unique. For example, if you are running your TCP/IP stacks in a sysplex, you would need to maintain one configuration for each stack on each of the systems. As more systems are added to the sysplex, more TCP/IP configuration files need to be maintained and synchronized.

In our case, we used MVS System symbols to enable us to share the definitions for our TCP/IP stacks across LPAR SC30, SC31, SC32, and SC33. MVS System symbols are used in creating shared definitions for systems that are in a sysplex. With this facility, you use the symbols defined during system startup as variables in configuring your TCP/IP stack. This means that you only need to create and maintain a template file for all the systems in the sysplex.

MVS System symbols processing

Use of MVS system symbols in the following files or environment variables is automatically supported:

- ▶ PROFILE.TCPIP file
- ▶ Resolver SETUP file
- ▶ TCPIP.DATA file
- ▶ OMPROUTE configuration file
- ▶ Resolver environment variables, such as RESOLVER_CONFIG and RESOLVER_TRACE

For the use of MVS system symbols in other configuration files, use the symbol translator utility, EZACFSM1. EZACFSM1 reads an input file which includes the system symbols, and creates an output file with the symbols translated to the system-specific values. This process is done before the files are read by TCP/IP.

The sample JCL for EZACFSM1 utility is included in hlq.SEZAINST(CONVSYM), as shown in Example B-1.

Example B-1 JCL for EZACFSM1

```
//CONVSYM JOB (accounting,information),programmer.name,
//          MSGLEVEL=(1,1),MSGCLASS=A,CLASS=A
//*
//STEP1 EXEC PGM=EZACFSM1,REGION=OK
//SYSIN DD DSN=TCP.DATA.INPUT,DISP=SHR
//*SYSIN DD PATH='/tmp/tcp.data.input'
//*          The input file can be either an MVS data set or an z/OS
//*          UNIX file.
//*
//*
//*
//SYSOUT DD DSN=TCP.DATA.OUTPUT,DISP=SHR
//*SYSOUT DD PATH='/tmp/tcp.data.output',PATHOPTS=(OWRONLY,OCREAT),
//*          PATHMODE=(SIRUSR,SIWUSR,SIRGRP,SIWGRP)
```

The input to EZACFSM1 is your template data set that contains the system symbols and the definitions that you need. The output data set will be the parameter files, such as TCPIP.DATA, that the TCP/IP stack or CS for z/OS IP application will use during its startup and operation. You need to run the utility on each of the systems where you need to have the symbols translated.

Symbols definitions

The variable &SYSCZONE is defined in the IEASYMxx member of SYS1.PARMLIB. As shown in Example B-2, the value for &SYSCZONE is derived from &SYSNAME. The variable &SYSNAME could be defined either in the IEASYSxx member or in the LOADxx member used during IPL. In our case, &SYSNAME was defined in IEASYSxx, which we used to IPL our MVS images. Refer to Example B-3 for a sample of the IEASYSxx that we used for the startup of SC30. You can find further information about system symbols in *z/OS V1R1.0-V1R2.0 MVS Initialization and Tuning Guide*, SA22-7591.

Example B-2 &SYSCZONE definition in SYS1.PARMLIB

```
SYSDEF          SYSCZONE(&SYSNAME(3:2)) 1
               SYMDEF(&SYSR2='037RZ1')
               SYMDEF(&SYSR3='&SYSR2(1:5).2')
               SYMDEF(&SYSR4='&SYSR2(1:5).3')
```

In this example, the numbers correspond to the following information:

- 1 The value of SYSCZONE is defined as two characters starting from the third character of SYSNAME. Our SYSNAME is SC30, so SYSCZONE resolves to 30.

Example B-3 IEASYSxx definition

```
COUPLE=00,
OMVS=7A,
PROD=01,
PROG=(A0,S0,D0,C1,L0),  Authorization list
SMF=00,
SMS=00,
SYSNAME=SC30, 1
SSN=00,
VAL=00,
```

In this example, the numbers correspond to the following information:

- 1 The SYSNAME is defined as SC30 in this LPAR (LPAR A11).

You can also define and use your own variable in configuring CS for z/OS IP, aside from &SYSNAME or &SYSCZONE. Refer to *z/OS Communications Server: IP Configuration Guide*, SC31-8775, for information about creating symbols output data set.

Include files

Together with the MVS System symbols support, we also used a facility (INCLUDE) to help us organize and share our stack configuration. By using the include configuration statement, we were able to structure our configuration better by putting different sections of PROFILE.TCPIP in separate files. During the stack's initialization, the contents of the file pointed to by the include statement are read and processed. These include statements are treated as though they were coded in PROFILE.

Sample PROFILE.TCPIP definition using MVS System symbols

Example B-4 shows the use of MVS system symbols in TCP/IP profile. Because &SYSCZONE is unique in each system, it ensures that the files and IDs that will be generated when the stacks initialize are also unique.

Important: The system symbols are stored in upper case by MVS. Because you can code the TCP/IP configuration statements in either upper case or lower case, you must ensure that you code the system symbol name in upper case.

In our environment, all stacks across LPARs shared the same OSAs and used the same HyperSockets interfaces. We could share the device-related definitions: DEVICE, LINK, BEGINROUTES, and START. We could not share the definitions for HOME and VIPADynamic statements because they are unique in each TCP/IP stack, so we made them into separate members and used the INCLUDE statement. We used the SYSCZONE value to point to those members (the members name must include SYSCZONE).

Example B-4 Use of system symbols in our TCP/IP profile

```
.....
;*****
; Include the stack-specific Dynamic VIPA definitions
;*****
INCLUDE TCPIPA.TCPPARMS(DVIPA&SYSCZONE) 1
;
;*****
; Include the stack-specific HOME definitions
;*****
INCLUDE TCPIPA.TCPPARMS(HOME&SYSCZONE) 1
;
..... Lines
deleted
;*****
; start the ftp daemon in each of the A stack
;*****
AUTOLOG 5
    FTPD&SYSCZONE JOBNAME FTPD&SYSCZONE.1 2
; OMP&SYSCZONE          ; OSPF daemon
; SMTP                  ; SMTP Server
ENDAUTOLOG
;*****
; Include the stack-specific PORT definitions
;*****
INCLUDE TCPIPA.TCPPARMS(PORT&SYSCZONE) 1
.....
START OSA2080I
START OSA20C0I
START OSA20E0I
START OSA20A0I
START IUTIQDF4
START IUTIQDF5
START IUTIQDF6
```

In this example, the numbers correspond to the following information:

- 1** Include file for system-specific device definitions.
- 2** Defines the AUTOLOG with unique FTP server daemon Jobname.

Note: A dot (.) is needed at the end of &SYSCONE because the next character is not a space.

Example B-5 shows the sample definition of a separate member for a stack-specific statement. It contains only the HOME statement for system SC30, called HOME30. This member is included in the PROFILE.TCPIP file in SC30 system. Likewise, define separate members for other LPARs.

Example B-5 Included device file HOME30 for SC30

```
;*****  
; TCIPA.TCPPARMS(HOME30)  
; HOME definitions for stack on SC30 image  
;*****  
HOME  
    10.1.1.10      VIPA1L  
    10.1.2.11      OSA2080I  
    10.1.3.11      OSA20C0I  
    10.1.3.12      OSA20E0I  
    10.1.2.12      OSA20A0I  
    10.1.4.11      IUTIQDF4L  
    10.1.5.11      IUTIQDF5L  
    10.1.6.11      IUTIQDF6L  
    10.1.2.10      VIPA2L  
    10.1.&SYSCONE..10 VIPA3L
```

Reusable Address Space ID (REUSASID) function examples

In this section, we detail sample definitions of the REUSASID function and results of its usage. Example B-6 shows how to enable this function in PARMLIB.

Example B-6 Sample DIAGXX member in PARMLIB

```
002800 VSM TRACK CSA(ON) SQA(ON)  
002900 VSM TRACE GETFREE(OFF)  
003000 REUSASID(YES) 1
```

In Example B-6, the number corresponds to the following information:

- 1** Parameter to code in member DIAGxx of PARMLIB (Example B-7) to enable REUSASID

Example B-7 Enabling the new DIAGXX definition

```
T DIAG=88  
IEE252I MEMBER DIAG88   FOUND IN SYS1.IBM.PARMLIB  
IEE536I DIAG      VALUE 88 NOW IN EFFECT
```

Without REUSASID, the old ASID is unavailable and a new ASID is assigned, as shown in Example B-8.

Example B-8 Without REUSASID TCPIP the old ASID is unavailable and a new ASID is assigned

```

D A,TCIPA
IEE115I 14.42.09 2010.298 ACTIVITY 318
  JOBS      M/S      TS USERS      SYSAS      INITS      ACTIVE/MAX VTAM      OAS
00004      00024      00003      00034      00019      00003/00030      00022
  TCIPA      TCIPA      TCIPA      NSW      SO A=0082 PER=NO SMC=000 1
                                PGN=N/A DMN=N/A AFF=NONE
                                CT=000.223S ET=333.691S
                                WUID=STC09685 USERID=TCPIP
                                WKL=SYSTEM SCL=SYSSTC P=1
                                RGP=N/A SRVR=NO QSC=NO
                                ADDR SPACE ASTE=062B7080
                                DSPNAME=00000EDC ASTE=093D7500
                                DSPNAME=TCPIPDS1 ASTE=7EE44C00

P TCIPA
EZZ4201I TCP/IP TERMINATION COMPLETE FOR TCIPA
IEF352I ADDRESS SPACE UNAVAILABLE 2
$HASP395 TCIPA ENDED

S TCIPA 3
$HASP100 TCIPA ON STCINRDR
IEF695I START TCIPA WITH JOBNAME TCIPA IS ASSIGNED TO USER
TCPIP , GROUP TCPGRP
$HASP373 TCIPA STARTED

D A,TCIPA
IEE115I 14.43.10 2010.298 ACTIVITY 446
  JOBS      M/S      TS USERS      SYSAS      INITS      ACTIVE/MAX VTAM      OAS
00004      00023      00002      00034      00019      00002/00030      00021
  TCIPA      TCIPA      TCIPA      NSW      SO A=0085 PER=NO SMC=000 4
                                PGN=N/A DMN=N/A AFF=NONE
                                CT=000.107S ET=030.909S
                                WUID=STC09689 USERID=TCPIP
                                WKL=SYSTEM SCL=SYSSTC P=1
                                RGP=N/A SRVR=NO QSC=NO
                                ADDR SPACE ASTE=062B7140
                                DSPNAME=00000EDC ASTE=093D7500
                                DSPNAME=TCPIPDS1 ASTE=7EE44C00

```

In Example B-8, the numbers correspond to the following information:

- 1** TCIPA ASID=0082.
- 2** Without REUSASID, the address space is unavailable.
- 3** TCIPA restarted without REUSASID parameter.
- 4** TCIPA new ASID=0085.

When REUSASID is enabled, the old ASID is available and reused, as shown in Example B-9 on page 417.

Example B-9 With REUSASID enabled the old ASID is available and reused

```
S TCIPA,REUSASID=YES 1
$HASP100 TCIPA  ON STCINRDR
IEF695I START TCIPA  WITH JOBNAME TCIPA  IS ASSIGNED TO USER
TCPIP  , GROUP TCPGRP
$HASP373 TCIPA  STARTED

D A,TCIPA
IEE115I 14.49.38 2010.298 ACTIVITY 711
  JOBS      M/S      TS USERS      SYSAS      INITS      ACTIVE/MAX VTAM      OAS
00004      00023      00002      00034      00019      00002/00030      00021
  TCIPA      TCIPA      TCIPA      NSW SO  A=0085  PER=NO  SMC=000 2
                                     PGN=N/A  DMN=N/A  AFF=NONE
                                     CT=000.121S  ET=069.808S
                                     WUID=STC09694  USERID=TCPIP
                                     WKL=SYSTEM  SCL=SYSSTC  P=1
                                     RGP=N/A      SRVR=NO  QSC=NO
                                     ADDR SPACE  ASTE=062B7140
                                     DSPNAME=00000EDC  ASTE=093D7500
                                     DSPNAME=TCPIPDS1  ASTE=7EE44C00

P TCIPA
EZZ4201I TCP/IP TERMINATION COMPLETE FOR TCIPA 3
$HASP395 TCIPA  ENDED

S TCIPA,REUSASID=YES 4
$HASP100 TCIPA  ON STCINRDR
IEF695I START TCIPA  WITH JOBNAME TCIPA  IS ASSIGNED TO USER
TCPIP  , GROUP TCPGRP
$HASP373 TCIPA  STARTED

D A,TCIPA
IEE115I 14.56.01 2010.298 ACTIVITY 868
  JOBS      M/S      TS USERS      SYSAS      INITS      ACTIVE/MAX VTAM      OAS
00004      00023      00001      00034      00019      00001/00030      00021
  TCIPA      TCIPA      TCIPA      NSW SO  A=0085  PER=NO  SMC=000 5
                                     PGN=N/A  DMN=N/A  AFF=NONE
                                     CT=000.111S  ET=028.495S
                                     WUID=STC09698  USERID=TCPIP
                                     WKL=SYSTEM  SCL=SYSSTC  P=1
                                     RGP=N/A      SRVR=NO  QSC=NO
                                     ADDR SPACE  ASTE=062B7140
                                     DSPNAME=00000EDC  ASTE=093D7500
                                     DSPNAME=TCPIPDS1  ASTE=7EE44C00
```

In Example B-9, the numbers correspond to the following information:

- 1 TCIPA started with REUSASID parameter.
- 2 TCIPA using ASID=0085.
- 3 When TCIPA is terminated the message ADDRESS SPACE UNAVAILABLE is no longer issued.
- 4 TCIPA restarted with REUSASID parameter.
- 5 The old ASID=0085 is being reused.

PROFILE.TCPIP statements

In this section we show PROFILE.TCPIP statements that are not always necessary, but are important. For detailed descriptions of statement, refer to *z/OS Communications Server: IP Configuration Guide*, SC31-8775. The syntax for the statement in the PROFILE can be found in *z/OS Communications Server: IP Configuration Reference*, SC31-8776.

IPCONFIG statements

This section provides information about IPCONFIG statements.

SOURCEVIPA

When the packet is sent to the destination host, the source IP address is included in the packet. In most cases the source IP address of the packet is used as the destination IP address of the returning packet from the other host. For the inbound traffic, z/OS Communications Server sets the destination IP address of the incoming packet to the source IP address of the return packet. However, for outbound traffic, the source IP address is determined by several parameters.

By default (IPCONFIG NOSOURCEVIPA), z/OS Communications Server sets the IP address of the interface which is used to send out a packet to a specific destination as the source IP address. The sending interface is selected depending on the routing table of the TCP/IP stack.

When IPCONFIG SOURCEVIPA is set, outbound datagrams use the virtual IP address (VIPA) for the source IP address of the packet instead of the physical interface IP address. By using VIPA as the source IP address, and therefore the destination IP address of the return packets from other hosts, SOURCEVIPA provides the tolerance of device and adapter failures.

The order of the HOME list is important if SOURCEVIPA is specified. The source IP address is the first static VIPA listed above the interface chosen for sending the packet. In Example B-10, if OSA20C0 2 is chosen as the actual physical interface for sending the outbound packet, then the IP address of the first VIPA above the HOME list, 10.1.2.10, is the source IP address.

Example B-10 Source IP Address selection with IPCONFIG SOURCEVIPA

```
....
IPCONFIG SOURCEVIPA

HOME
  10.1.1.10      VIPA1L
  10.1.2.10      VIPA2L 1
  10.1.2.11      OSA2080I
  10.1.3.11      OSA20C0I 2
....
```

Note: The source IP address selection can be overridden with SRCIP statement. Refer to “SRCIP” on page 437 for details.

SOURCEVIPA has no effect on OSPF or RIP route information exchange packets generated by the OMPROUTE routing daemon, which means that it is only applicable for data diagrams.

MULTIPATH

With the IPCONFIG MULTIPATH statement, packets can be load balanced on routes that have been defined to be of equal cost. These routes could either be learned dynamically or defined statically in your routing program (OMPROUTE). With multipath enabled, TCP/IP will select a route to that destination network or host on a round-robin basis. TCP/IP can select a route on a per-connection or per-packet basis, but we recommend that you do not use the per-packet basis because it requires high CPU processing for reassembly of out-of-order packets at the receiver. See Chapter 5, “Routing” on page 205 for details about this topic.

By default (IPCONFIG NOMULTIPATH), there is no multipath support and all connections use the first active route to the destination network or host even if there are other, equal-cost routes available.

PATHMTUDISCOVERY

Coding IPCONFIG PATHMTUDISCOVERY prevents the fragmentation of datagrams. It tells TCP/IP to discover dynamically the Path Maximum Transfer Unit (PMTU), which is the smallest of the MTU sizes of each hop in the path between two hosts.

When a connection is established, TCP/IP uses the minimum MTU of the sending host as the starting segment size and sets the Don't Fragment (DF) bit in the IP header. Any router along the route that cannot process the MTU will return an ICMP message requesting fragmentation and will inform the sending host that the destination is unreachable. The sending host can then reduce the size of its assumed PMTU. You can find more information about PMTU discovery in RFC 1191, Path MTU Discovery.

The default is IPCONFIG NOPATHMTUDISCOVERY. Aside from enabling PMTU during stack initialization, you could also enable or disable PMTU discovery by using VARY OBEYFILE.

IQDIOROUTING

When IPCONFIG IQDIOROUTING is configured, the inbound packets that are to be forwarded by this TCP/IP stack use HiperSockets (also known as Internal Queued Direct I/O or IQDIO) and Queued Direct I/O (QDIO) directly and bypass the TCP/IP stack. This type of routing is called *HiperSockets Accelerator* because it allows you to concentrate external network traffic over a single OSA-Express QDIO connection and then accelerates the routing over a HiperSockets link, bypassing the TCP/IP stack. The default is NOIQDIOROUTING. For further information about HiperSockets, see Chapter 4, “Connectivity” on page 117.

ARPTO

IPCONFIG ARPTO and ARPAGE statements have the same function: they specify the time interval between the creation or revalidation and deletion of an entry in the ARP table. The value of IPCONFIG ARPTO is specified in seconds, and the value of ARPAGE is specified in minutes. ARP cache entries for MPCIPA and MPCOSA are not affected by ARPTO or ARPAGE because they use the ARP offload function. The ARP cache timer for ARP offload is set to 20 minutes. It is hard-coded and not configurable. For more information about devices that are affected by ARPTO, refer to *z/OS Communications Server: IP Configuration Guide*, SC31-8775.

The UNIX shell command **onetstat -R** displays the current ARP cache entries. The upper case R in the option is required for this display. A third parameter can be coded that would specify the IP address of the entry you want to display, as the **NETSTAT ARP ip_addr** command does from TSO. If you want to display the entire ARP cache, you can specify the third parameter with the reserved word ALL (again, all in upper case letters). If you do not specify in upper case letters, the reserved word is not recognized (see Example B-11).

Example B-11 ARP display

```

D TCPIP,TCPIPA,N,ARP
QUERYING ARP CACHE FOR ADDRESS 10.1.2.11
INTERFACE: OSA2080I          ETHERNET: 020002749925
QUERYING ARP CACHE FOR ADDRESS 10.1.2.32
INTERFACE: OSA2080I          ETHERNET: 00145E749924
QUERYING ARP CACHE FOR ADDRESS 10.1.2.31
INTERFACE: OSA2080I          ETHERNET: 00145E749924
QUERYING ARP CACHE FOR ADDRESS 10.1.2.41
INTERFACE: OSA2080I          ETHERNET: 00145E749924
QUERYING ARP CACHE FOR ADDRESS 10.1.2.42
INTERFACE: OSA2080I          ETHERNET: 00145E749924
..... Lines deleted
QUERYING ARP CACHE FOR ADDRESS 10.1.3.13
INTERFACE: OSA20E0I          ETHERNET: 020003749A7F
QUERYING ARP CACHE FOR ADDRESS 10.1.4.12
INTERFACE: IUTIQDF4L
QUERYING ARP CACHE FOR ADDRESS 10.1.4.11
INTERFACE: IUTIQDF4L
31 OF 31 RECORDS DISPLAYED
END OF THE REPORT

```

GLOBALCONFIG statements

This section provides information about GLOBALCONFIG statements.

TCPIPSTATISTICS

This statement prints the values of several TCP/IP counters to the output data set designated by the CFGPRINT JCL statement. These counters include the number of TCP retransmissions and the total number of TCP segments sent from the MVS TCP/IP system. These TCP/IP statistics are written to the designated output data only during termination of the TCP/IP address space.

The TCPIPSTATISTICS parameter is confirmed by the message:

```
EZZ0613I TCPIPSTATISTICS IS ENABLED
```

This parameter should be specified in the GLOBALCONFIG section. Note that the SMFCONFIG TCPIPSTATISTICS parameter serves a different purpose; it requests that SMF records of subtype 5 containing TCP/IP statistics be created.

SEGMENTATIONOFFLOAD

When sending or receiving packets over OSA-Express in QDIO mode with checksum offload support, TCP/IP offloads most IPv4 (outbound and inbound) checksum processing (IP header, TCP, and UDP checksums) to the OSA. The TCP/IP stack still performs checksum processing in the cases where checksum cannot be offloaded.

When sending packets over OSA-Express in QDIO mode with TCP segmentation offload support, TCP/IP offloads most IPv4 outbound TCP segmentation processing to the OSA. The TCP/IP stack still performs TCP segmentation processing in the cases where segmentation cannot be offloaded.

Tip: Applications that use large TCP send buffers will obtain the most benefit from TCP segmentation offload. The size of the TCP receive buffer on the other side of the TCP connection also affects the negotiated buffer size.

You can control the size of these buffers using the TCPSENBFRSIZE and TCPRCVBFRSIZE parameters on the TCPCONFIG statement to set the default TCP send/receive buffer size for all applications. However, an application can use the SO_SNDBUF socket option to override the default TCP send buffer sizes (example FTP).

The segmentation offload feature decreases host CPU utilization and increases data transfer efficiency for IPv4 packets. The z/OS Communications Server provides the Offloads feature for IPv4 segmentation processing to OSA-Express2 in QDIO mode. This enhances the data transfer efficiency of IPv4 packets while decreasing host CPU utilization.

The OFFLOAD feature is supported by OSA-Express in QDIO mode.

Example B-12 displays the NETSTAT DEVLINKS of an OSA-Express that has SegmentationOffload enabled.

Example B-12 Segmentation Offload enabled

```

D TCPIP,TCPIPA,N,DE
..... Lines deleted
EZD0101I NETSTAT CS V1R12 TCPIPA 896
INTFNAME: OSA2080I          INTFTYPE: IPAQENET  INTFSTATUS: READY
PORTNAME: OSA2080  DATAPATH: 2082  DATAPATHSTATUS: READY
CHPIDTYPE: OSD
SPEED: 0000001000
IPBROADCASTCAPABILITY: NO
VMACADDR: 020010776873  VMACORIGIN: OSA  VMACROUTER: LOCAL
ARPOFFLOAD: YES 1          ARPOFFLOADINFO: YES 1
CFGMTU: 1492          ACTMTU: 1492
IPADDR: 10.1.2.11/24
VLANID: 10          VLANPRIORITY: DISABLED
DYNVLANREGCFG: NO  DYNVLANREGCAP: YES
READSTORAGE: GLOBAL (4096K)
INBPERF: BALANCED
CHECKSUMOFFLOAD: YES 1      SEGMENTATIONOFFLOAD: YES 1
SECCLASS: 255          MONSYSPLEX: NO
ISOLATE: NO          OPTLATENCYMODE: NO
MULTICAST SPECIFIC:
MULTICAST CAPABILITY: YES
..... Lines deleted

```

In this example, the numbers correspond to the following information:

1 Indicates the enabled features of ARP, Segmentation and Checksum Offload.

IQDMULTIWRITE | NOIQDMULTIWRITE

This statement allows the HiperSockets to move multiple buffers of data with a single write operation. HiperSockets multiple write can reduce CPU use and increase throughput for outbound streaming-type workloads, such as FTP transfers.

This parameter applies to all HiperSockets interfaces, including IUTIQDIO and IQDIOINTF6 interfaces created for Dynamic XCF.

NOIQDIOMULTIWRITE | IQDIOMULTIWRITE

This statement tells to TCP/IP should displace CPU cycles for HiperSockets multiple write workload to a zIIP. Example B-13 shows the output of z/OS command NETSTAT CONFIG.

Example B-13 Output of DISPLAY NETSTAT CONFIG

```
EZD0101I NETSTAT CS V1R12 TCPIPA 932
TCP CONFIGURATION TABLE:
DEFAULTRCVBUFSIZE: 00131072  DEFAULTSNDBUFSIZE: 00131072
DEFLTMAXRCVBUFSIZE: 00262144  SOMAXCONN: 0000000010
MAXRETRANSMITTIME: 120.000  MINRETRANSMITTIME: 0.500
ROUNDTRIPGAIN: 0.125  VARIANCEGAIN: 0.250
VARIANCEMULTIPLIER: 2.000  MAXSEGLIFETIME: 30.000
DEFAULTKEEPAIVE: 00000120  DELAYACK: YES
RESTRICTLOWPORT: NO  SENDGARBAGE: NO
TCPTIMESTAMP: YES  FINWAIT2TIME: 600
TTLS: NO
UDP CONFIGURATION TABLE:
DEFAULTRCVBUFSIZE: 00065535  DEFAULTSNDBUFSIZE: 00065535
CHECKSUM: YES
RESTRICTLOWPORT: NO  UDPQUEUELIMIT: NO
IP CONFIGURATION TABLE:
FORWARDING: YES  TIMETOLIVE: 00064  RSMTIMEOUT: 00060
IPSECURITY: NO
ARPTIMEOUT: 01200  MAXRSMSSIZE: 65535  FORMAT: LONG
IGREDIRECT: YES  SYSPLXROUT: YES  DOUBLENOP: NO
STOPCLAWER: NO  SOURCEVIPA: YES
MULTIPATH: CONN  PATHMTUDSC: YES  DEVRTRYDUR: 0000000090
DYNAMICXCF: YES
  IPADDR: 10.1.7.11  SUBNET: 255.255.255.0  METRIC: 08
  SECCCLASS: 255
QDIOACCEL: YES  QDIOACCELPRIORITY: 1
IQDIOROUTE: N/A
TCPSTACKSRCVIPA: NO
IPV6 CONFIGURATION TABLE:
FORWARDING: YES  HOPLIMIT: 00255  IGREDIRECT: NO
SOURCEVIPA: NO  MULTIPATH: NO  ICMPERRRLIM: 00003
IGRTRHOPLIMIT: NO
IPSECURITY: NO
DYNAMICXCF: NO
TCPSTACKSRCVIPA: NO
TEMPADDRESSES: NO
SMF PARAMETERS:
TYPE 118:
  TCPINIT: 00  TCPTERM: 00  FTPCLIENT: 00
  TN3270CLIENT: 00  TCPIPSTATS: 00
TYPE 119:
```

```

TCPINIT:      YES  TCPTERM:    YES  FTPCLIENT:    YES
TCPIPSTATS:   YES  IFSTATS:   NO   PORTSTATS:    NO
STACK:        NO   UDPTERM:    NO   TN3270CLIENT: YES
IPSECURITY:   NO   PROFILE:    NO   DVIPA:          NO
GLOBAL CONFIGURATION INFORMATION:
TCPIPSTATS: NO   ECSALIMIT: 0000000K  POOLLIMIT: 0000000K
MLSCHKTERM: NO   XCFGRPID: 21          IQDVLANID: 21
SEGOFFLOAD: YES  SYSPLEXWLMPOLL: 060  MAXRECS: 100
EXPLICITBINDPORTRANGE: 00000-00000  IQDMULTIWRITE: YES  1
WLMPPRIORITYQ: NO
SYSPLEX MONITOR:
TIMERSECS: 0060  RECOVERY: YES  DELAYJOIN: YES  AUTOREJOIN: NO
MONINTF:  NO     DYNROUTE: NO   JOIN:        YES
ZIIP:
IPSECURITY: NO   IQDIOMULTIWRITE: YES  2
NETWORK MONITOR CONFIGURATION INFORMATION:
PKTTRCSRV: NO   TCPCNNSRV: NO   NTASRV: NO
SMFSRV:  YES
IPSECURITY: YES  PROFILE: YES  CSSMTP: YES  CSMAIL: NO   DVIPA: YES
AUTOLOG CONFIGURATION INFORMATION: WAIT TIME: 0300
PROCNAME: FTPDA  JOBNAME: FTPDA1
  PARMSTRING:
  DELAYSTART: NO
PROCNAME: OMPA   JOBNAME: OMPA
  PARMSTRING:
  DELAYSTART: NO
PROCNAME: IOASRV JOBNAME: IOASRV
  PARMSTRING:
  DELAYSTART: NO
END OF THE REPORT

```

In this example, the numbers correspond to the following information:

- 1** Indicates the enabled features of Hipersocket Multi Write is “on.”
- 2** Indicates the enabled features of zIIP Assisted Hipersocket Multiple Write is “on.”

PORT statement

This section discusses uses of the PORT statement.

Port sharing (TCP only)

If you want to run multiple instances of a listener for performance reasons, you can share the same port between them. TCP/IP will select the listener with the fewest connections (both active and in the backlog) at the time when a client request comes in. A typical application using this feature is the Internet Connection Secure Server. If the load gets high, additional servers are started by the Workload Manager.

An example of a shared port is:

```
PORT
      80    TCP    WEBSRV1 SHAREPORT
      80    TCP    WEBSRV2
      80    TCP    WEBSRV3
```

BIND control for INADDR_ANY

The BIND option associates the server job name with a specific IP address when the server binds to INADDR_ANY. This new function can be used to change the BIND for INADDR_ANY to a BIND for a specific IP address.

Telnet, for example, is a server that binds to INADDR_ANY. Previously, an client that wants to access both Telnet servers, TN3270 and UNIX Telnet, would connect to different ports or different TCP/IP stacks, depending on which Telnet server it wanted to connect to. This led to cases where either one server used a different, nonstandard port, or multiple TCP/IP stacks had to be used. With this function you do not need to have two different ports or TCP/IP stacks. You use the same port 23 for both TN3270 and UNIX Telnet. All that is needed is to code the BIND keyword in the PORT statement for each server:

```
PORT
      23 TCP    TN3270A BIND 10.1.1.10
      23 TCP    OMVS      BIND 10.1.1.20
```

In this case, the TN3270A is a jobname for TN3270 server. When it BINDs to port 23 and INADDR_ANY, it is associated with IP address 10.1.1.10. The OMVS job name identifies any UNIX server, including the UNIX Telnet server. When UNIX Telnet Server BINDs to port 23 and INADDR_ANY, it is associated with IP address 10.1.1.20.

Both IP addresses can be dynamic VIPA addresses, static VIPA addresses, or real interface addresses. You also can code a wild card for the job name. Note that this function will work only for servers that bind to INADDR_ANY, and it is not valid with the PORTRANGE statement.

TCPCONFIG/UDPCONFIG RESTRICTLOWPORTS

Port numbers that are not specified on a PORT profile statement are considered unreserved ports. You can restrict the use of unreserved ports below 1024 to programs that are APF-authorized or have OMVS superuser authority. You might decide not to explicitly reserve all well-known ports by defining the UNRESTRICTLOWPORTS option on the TCPCONFIG and UDPCONFIG statements. This would allow any socket application to acquire a well-known port. See Example B-14 on page 425.

TCPCONFIG UNRESTRICTLOWPORTS

UDPCONFIG UNRESTRICTLOWPORTS

If you want the well-known ports to be used only by predefined application processes or superuser-authorized application processes, then you can define the RESTRICTLOWPORTS option on the TCPCONFIG and UDPCONFIG statements. This prevents any non-authorized socket application from acquiring a well-known port.

PORT

The PORT reservations that are defined in the PROFILE data set are the ports that are used by specific applications. You control access to particular ports by port number, by reserving the port using the PORT or PORTRANGE profile statements. You can also use the optional SAF parameter to provide additional access control.

You then need to explicitly define PORT statements to reserve each port or define the process with superuser authority in RACF. The reserved ports indicate that the port is not available for use by any user. However, the unreserved port numbers from 1024 through 65535 are available for use by any application that issues an explicit bind to a specific unreserved port. These port numbers are also used by the stack to provide stack-selected ephemeral ports.

Controlling access to unreserved ports

You can also use the PORT statement to control application access to unreserved ports by configuring one or more PORT statements in which the port number is replaced by the keyword UNRSV. The UNRSV keyword refers to any unreserved port (any port number that has not been reserved by a PORT or PORTRANGE statement). If you configure the RESTRICTLOWPORTS parameter on the TCPCONFIG or UDPCONFIG profile statement, PORT UNRSV statements for the corresponding protocol control access only to unreserved ports above port 1023. If you do not configure the RESTRICTLOWPORTS parameter, PORT UNRSV statements control access to all unreserved ports in the range 1 to 65535.

This new type of entry is identified by the keyword 'UNRSV' and it too is used to specify the jobname or user IDs that are allowed to run applications that use an application-specified unreserved port.

Note: this new control (UNRSV) does not affect the use of ports that are selected by the stack either as a local ephemeral port or as a sysplex-wide port for a distributed DVIPA.

You reserve the ports with PORT, PORTRANGE, and UNRSV commands using the keyword OMVS with a job name of the process or a wild card job name such as *. UNIX applications. The job can fork() another address space with a different name (for example, inetd or FTP server). Example B-15 shows the access control to the ports.

Example B-15 PROFILE.TCPIP: PORT, PORTRANGE and UNRSV

```

TCPCONFIG
  RESTRICTLOWPORTS
UDPCONFIG
  RESTRICTLOWPORTS
PORT
  20 TCP OMVS      NOAUTOLOG    ; FTP Server 1
  21 TCP FTPDA1   BIND 10.1.1.10 ; FTP Server 2
  23 TCP TN3270A                ; Telnet Server
  25 TCP SMTP                ; SMTP Server
  514 UDP OMVS                ; UNIX SyslogD Server 3
  520 UDP OMPA NOAUTOLOG      ; OMPROUTE RIP IPv4
  521 UDP OMPA NOAUTOLOG      ; OMPROUTE RIP IPv6
PORTRANGE 10000 2000 TCP OMVS ; TCP 10000 - 11999 4
PORTRANGE 10000 2000 UDP OMVS ; UDP 10000 - 11999 4
PORT  UNRSV  UDP *  DENY      5

```

Normally you can specify either OMVS or the job name in the PORT statement. However, certain daemons have special considerations on this matter.

When the FTP server starts, it forks the listener process to run in the background, requiring that the name of the forked address space (FTPDA1, in this example), not the original procedure name, be used on the PORT statement of the control connection 2. You must specify OMVS as the name on the PORT for FTP's PORT 20 1, which is used for the data connection managed by the child process. If you specify the forked name on the data connection (Port 20), the data connections will fail.

Note that you can also reserve UDP port 514 3 to OMVS. This port is used by the SyslogD server in OMVS to receive log messages from other SyslogD servers in the TCP/IP network. The PORTRANGE statements 4 reserve a range of ephemeral TCP and UDP ports for UNIX System Services and the PORT UNRSV statement 5 denies UDP explicit bind access to application-specified unreserved ports by any job.

In Example B-15, ports 10000 to 11999 are reserved. The range must match the INADDRANYPORT and INADDRANYCOUNT in your BPXPRMxx member 6 (see Example B-16).

Example B-16 INADDRANYPORT and INADDRANYCOUNT in BPXPRMxx member

```

NETWORK DOMAINNAME(AF_INET)
  DOMAINNUMBER(2)
  MAXSOCKETS(10000)
  TYPE(INET)
  INADDRANYPORT(10000) 6
  INADDRANYCOUNT(2000) 6

NETWORK DOMAINNAME(AF_INET6)
  DOMAINNUMBER(19)
  MAXSOCKETS(10000)
  TYPE(INET)

```

To display the PORT reservation list, use the TSO/E command NETSTAT PORTL, the MVS command D TCPIP,procname,NETSTAT PORTL command, or the UNIX shell command **onetstat -p procname -o**. Example B-17 shows these MVS commands.

Example B-17 Viewing port reservation list

```

D TCPIP,TCPIPA,N,PORTL
PORT#  PROT  USER      FLAGS    RANGE      SAF  NAME
UNRSV  UDP    *          XL
20      TCP    OMVS       D
21      TCP    FTPDA1     DABU
        BINDSPECIFIC: 10.1.1.10
23      TCP    TN3270A    DA 500    UDP  IKED      DA
520     UDP    OMPROUTE   D
521     UDP    OMPROUTE   D
4500    UDP    IKED       DA
7 OF 7 RECORDS DISPLAYED
END OF THE REPORT

```

TCPCONFIG TCPSEENDBFRSIZE

TCPCONFIG TCPSEENDBFRSIZE specifies the TCP send buffer size. This value is used as the default send buffer size for those applications that do not explicitly set the buffer size using SETSOCKOPT(). The default is 16384 (16 000).

TCPCONFIG TCPRCVBUFRSIZE

TCPCONFIG TCPRCVBUFRSIZE specifies the TCP receive buffer size. This value is used as the default receive buffer size for those applications that do not explicitly set the buffer size using SETSOCKOPT(). You can specify value from 256 and TCPMAXRCVBUFRSIZE. The default is 16384 (16 000).

TCPCONFIG TCPMAXRCVBUFRSIZE

TCPCONFIG TCPMAXRCVBUFRSIZE specifies the TCP maximum receive buffer size an application can set as its receive buffer size using SETSOCKOPT(). You can use this parameter to limit the receive buffer size that an application can set. The minimum value you can specify is TCPRCVBUFRSIZE, and the maximum is 512 KB. The default is 256 KB.

Note: The FTP server and client applications override the default settings and use 64 KB as the TCP window size and 180 KB for send/recv buffers. No changes are required in the TCPCONFIG statement for the FTP server and client.

TCPCONFIG FINWAIT2TIME

TCPCONFIG FINWAIT2TIME parameter allows you to specify the number of seconds a TCP connection should remain in the FINWAIT2 state. When this time limit is reached, the system waits a further 75 seconds before dropping the connection. The default is 600 seconds, but you can specify a value as low as 60 seconds, which will reduce the time a connection remains in the FINWAIT2 status, and thereby free up resources for future connections.

TCPCONFIG TCPTIMESTAMP

The TCP time stamp option is exchanged during connection setup. This option is enabled (by default) using the TCPCONFIG TCPTIMESTAMP parameter. Enabling the TCP time stamp allows TCP/IP to better estimate the Route Trip Response Time (RTT), which helps avoid unnecessary retransmissions and helps protect against the wrapping of sequence numbers.

IDYNAMICXCF

You have the option of either defining the DEVICE, LINK, HOME, and START statements for MPC XCF connections to another z/OS, or letting TCP/IP dynamically define them for you. Dynamic XCF devices and links, when activated, appear to the stack as though they had been defined in the TCP/IP profile. They can be displayed using standard commands, and they can be stopped and started. For multiple stack environments, IUTSAMEH links are dynamically created for same-LPAR links. Refer to *IBM z/OS Communications Server TCP/IP Implementation Volume 3: High Availability, Scalability, and Performance*, for further details.

SACONFIG (SNMP subagent)

The SACONFIG statement provides subagent support for SNMP. Through the subagent support you can manage an ATM OSA network interface. Refer to the SNMP subagent chapter of *IBM z/OS Communications Server TCP/IP Implementation Volume 2: Standard Applications*, SG24-7897, for further information. Example B-18 shows an example of this statement.

Example B-18 The SACONFIG statement

```
SACONFIG
  COMMUNITY public    ; Community string
  OSASF 760           ; OSASF port number
; AGENT 161           ; Agent port number
  ENABLED
  SETSENAbleD
  ATMENableD
```

SMFCONFIG

The SMFCONFIG statement is used to turn on SMF logging. It defines the type 118 and type 119 records to be collected (the default format is type 118). The Example below shows the SMFCONFIG statement to provide SMF logging for TCP stack activity, TCP connection initialization, TCP connection termination TCP/IP statistics, when a IPSEC dynamic tunnel is added and removed and when a IPSEC manual tunnel is activated or deactivated:

```
SMFCONFIG TYPE119 DVIPA TCPSTACK TCPINIT TCPTerm TCPIPSTATISTICS IPSECURITY
```

The SMFPARMS statement can also be used to turn on SMF logging. However, you are encouraged to migrate to SMFCONFIG, which has the following advantages over the SMFPARMS statement:

- ▶ Using SMFCONFIG means that SMF records are written using standard subtypes. With SMFPARMS, you have to specify the subtypes to be used.
- ▶ SMFCONFIG allows you to record both type 118 and type 119 records. With SMFPARMS, only type 118 records can be collected.
- ▶ SMFCONFIG enables you to record a wider variety of information.
- ▶ By using SMFCONFIG, you gain support for dynamic reconfiguration, for all environments under which CS for z/OS IP is executing (SRB mode, reentrant, XMEM mode, and so on), and you can avoid duplicate SMF exit processes.

In the following example, type 118 FTP client records, type 119 TN3270 client records, and type 119 IPSEC records are collected:

```
SMFCONFIG TYPE118 FTPCLIENT
          TYPE119 TN3270CLIENT IPSECURITY
```

The preceding example can also be coded this as shown here, because type 118 records are collected by default:

```
SMFCONFIG FTPCLIENT
          TYPE119 TN3270CLIENT IPSECURITY
```

SMFCONFIG is coded in the PROFILE.TCPIP, but it has related entries in both Telnet and in FTP. (See *IBM z/OS Communications Server TCP/IP Implementation Volume 2: Standard Applications*, SG24-7897.)

We use the **NETSTAT CONFIG** command to check if the SMFCONFIG setup above is right, the Example B-19 shows it.

Example B-19 Output NETSTAT CONFIG

```
D TCPIP,TCPIPA,N,CONFIG
SMF PARAMETERS:
TYPE 118:
  TCPINIT:      00  TCPTERM:    00  FTPCLIENT:    00
  TN3270CLIENT: 00  TCPIPSTATS: 00
TYPE 119:
  TCPINIT:      YES TCPTERM:    YES FTPCLIENT:    YES
  TCPIPSTATS:   YES IFSTATS:    NO PORTSTATS:    NO
  STACK:        NO  UDPTERM:    NO  TN3270CLIENT: YES
  IPSECURITY:   NO  PROFILE:    NO  DVIPA:          YES
```

The only SMF exit supported in CS for z/OS IP is the FTP server SMF exit, FTPSMFEX. This exit is only called for type 118 records. If you need to access type 119 FTP SMF records, use the standard SMF exit facilities, IEFU83, IEFU84, and IEFU85.

For further information about TCP/IP SMF record layouts and standardized subtype numbers, refer to *z/OS Communications Server: IP Configuration Reference*, SC31-8776.

Netmonitor

Use the NETMONITOR statement to activate or deactivate selected real-time TCP/IP network management interfaces (NMI).

The NETMONITOR parameters, TCPCONNSERVICE and SMFSERVICE, provide two functions:

- They control the availability of the real-time SMF services that are associated with each parameter. T
- They control the creation of the SMF 119 records that are supported by each service.

If you want your application to process only SMF 119 records by using these real-time SMF services, you need to configure only the NETMONITOR profile statement. You do not need to request support for these SMF 119 records on the SMFCONFIG profile statement.

The SMFSERVICE parameter can be used to configure the real-time TCP/IP network monitoring NMI to support the new SMF 119 event records, subtypes 32– 37, which provide sysplex event information, specifying the subparameter DVIPA, as shown in Example B-20:

Example B-20 TCPIPA profile contents, NETMONITOR option

```

;
NETMONITOR SMFSERVICE DVIPA
;

```

To verify the configuration is implemented as expected, we used the comand **NETSTAT,CONFIG**, as shown in Example B-21:

Example B-21 Using NETSTAT CONFIG command to verify the network monitor statements

```

D TCPIP,TCPIPA,N,CONFIG
NETWORK MONITOR CONFIGURATION INFORMATION:
PKTTRCSRV: NO   TCPCNSRV: NO   NTASRV: NO
SMFSRV:    YES
  IPSECURITY: YES  PROFILE: YES  CSSMTP: YES  CSMAIL: NO   DVIPA: YES

```

For further information about the NETMONITOR usage, refer to *z/OS Communications Server: IP Programmer's Guide and Reference*, SC31-8787

INTERFACE statement

You can use the INTERFACE statement to define either IPv4 or IPv6. If used for IPv4, the statement combines the definitions of the DEVICE/LINK/HOME statements. See Example B-22.

Example B-22 INTERFACE statement to define IPv4 interfaces

```

INTERFACE OSA20A0I  1
  DEFINE IPAQENET    2
  PORTNAME OSA20A0   3
  IPADDR 10.1.2.12/24 4
  MTU 1492           5
  VLANID 20          6
  VMAC               7
  SOURCEVIPAINTE VIPA2L 8
;
INTERFACE OSA20A0X  1
  DEFINE IPAQENET    3
  PORTNAME OSA20A0   4
  IPADDR 10.1.10.16/24 6
  MTU 1492           7
  VLANID 21          7
  VMAC               7

```

In Example B-22, the numbers correspond to the following information:

- 1** The INTERFACE statement replaces the DEVICE and LINK statements. The INTERFACE statement label must be unique.
- 2** In IPv4 the INTERFACE statement can be used for IPAQENET devices only.

- 3** The PORTNAME operand as defined in TRL node. For multiple VLAN configurations the same PORTNAME can be defined several times.
- 4** The IPADDR operand replaces the HOME statement. The optional subnetmask definition replaces similar definition coded in BEGINROUTES.
- 5** The optional MTU operand replaces similar definition coded in BSDROUTINGPARMS.
- 6** The optional VLANID operand is required when defining multiple VLANs.
- 7** The optional VMAC operand, with or without set values, is required when defining VLANs.
- 8** SOURCEVIPAINTE defines VIPA associated with this INTERFACE.

Note: If SOURCEVIPAINTE is coded, the whole INTERFACE definition block must be defined in PROFILE *after* the VIPA DEVICE and LINK statements are defined.

Example B-23 shows the output of the **netstat dev (-d)** command.

Example B-23 Display netstat dev (-d)

```

D TCP/IP,TCPIP,A,N,DE
..... Lines deleted
DEVNAME: OSA2080 1 DEVTYPE: MPCIPA
DEVSTATUS: READY
LNKNAME: OSA2080I 2 LNKTYPE: IPAQENET LNKSTATUS: READY
SPEED: 0000001000
IPBROADCASTCAPABILITY: NO
CFGROUTER: NON ACTROUTER: NON
ARPOFFLOAD: YES ARPOFFLOADINFO: YES
ACTMTU: 8992
VLANID: 10 VLANPRIORITY: DISABLED
DYNVLANREGCFG: NO DYNVLANREGCAP: YES
READSTORAGE: GLOBAL (4096K) INBPERF: BALANCED
CHECKSUMOFFLOAD: YES
SECCLASS: 255 MONSYSPLEX: NO
BSD ROUTING PARAMETERS:
MTU SIZE: 1492 METRIC: 90
DESTADDR: 0.0.0.0 SUBNETMASK: 255.255.255.0
MULTICAST SPECIFIC:
MULTICAST CAPABILITY: YES
..... Lines deleted
INTFNAME: OSA20A0I 3 INTFTYPE: IPAQENET INTFSTATUS: READY
PORTNAME: OSA20A0 4 DATAPATH: 20A2 5 DATAPATHSTATUS: READY
SPEED: 0000001000
IPBROADCASTCAPABILITY: NO
VMACADDR: 020012749661 VMACORIGIN: OSA 6 VMACROUTER: ALL
ARPOFFLOAD: YES ARPOFFLOADINFO: YES
CFGMTU: 1492 ACTMTU: 1492
IPADDR: 10.1.2.12/24 7
VLANID: 20 VLANPRIORITY: DISABLED
DYNVLANREGCFG: NO DYNVLANREGCAP: YES
READSTORAGE: GLOBAL (4096K) INBPERF: BALANCED
CHECKSUMOFFLOAD: YES
SECCLASS: 255 MONSYSPLEX: NO
MULTICAST SPECIFIC:
MULTICAST CAPABILITY: YES
..... Lines deleted

```

```

INTFNAME: OSA20A0X 3 INTFTYPE: IPAQENET INTFSTATUS: READY
PORTNAME: OSA20A0 4 DATAPATH: 20A3 5 DATAPATHSTATUS: READY
SPEED: 0000001000
IPBROADCASTCAPABILITY: NO
VMACADDR: 020013749661 VMACORIGIN: OSA 6 VMACROUTER: ALL
ARPOFFLOAD: YES ARPOFFLOADINFO: YES
CFGMTU: 1492 ACTMTU: 1492
IPADDR: 10.1.10.16/24 7
VLANID: 21 VLANPRIORITY: DISABLED
DYNVLANREGCFG: NO DYNVLANREGCAP: YES
READSTORAGE: GLOBAL (4096K) INBPERF: BALANCED
CHECKSUMOFFLOAD: YES
SECCCLASS: 255 MONSYSPLEX: NO
MULTICAST SPECIFIC:
MULTICAST CAPABILITY: YES

```

Compare resulting displays of resources defined with the DEVICE/LINK/HOME statement to resources defined with the INTERFACE statement. In Example B-23 on page 431, the numbers correspond to the following information:

- 1 Device and link names.
- 2 Device and link names.
- 3 Interface name.
- 4 PORTNAME in use matching the TRLE PORTNAME definition. The same PORTNAME is defined several times for multiple VLANs (that is not possible with DEVICE/LINK).
- 5 DATAPATH device address in use. One DATAPATH device is needed per each INTERFACE defined on the same physical OSA port.
- 6 Virtual MAC address (VMAC) dynamically assigned by OSA.
- 7 IP Address and subnet mask.

Note: The `netstat dev (-d)` command always return the resources defined with the DEVICE/LINK statements first and the resources defined with the INTERFACE statement later.

Example B-24 shows the OBEYFILE definition to delete an INTERFACE.

Example B-24 OBEYFILE definition to delete an INTERFACE

```

INTERFACE OSA20A0I
DELETE 1

```

In Example B-24, the numbers correspond to the following information:

- 1 Parameter to code to delete an INTERFACE. Note syntax differences from DEVICE/LINK deletion coding.

Example B-25 shows the process to delete an INTERFACE.

Example B-25 Process to delete an INTERFACE

```

V TCPIP,TCPIPA,STOP,OSA20A0I 1
EZZ0060I PROCESSING COMMAND: VARY TCPIP,TCPIPA,STOP,OSA20A0I
EZZ0053I COMMAND VARY STOP COMPLETED SUCCESSFULLY

D TCPIP,TCPIPA,N,DE,INTFN=OSA20A0I 2
INTFNAME: OSA20A0I          INTFTYPE: IPAQENET  INTFSTATUS: NOT ACTIVE 2
PORTNAME: OSA20A0  DATAPATH: UNKNOWN  DATAPATHSTATUS: NOT ACTIVE 2
..... Lines deleted
IPV4 LAN GROUP SUMMARY
LANGROUP: 00002
NAME          STATUS      ARPOWNER      VIPAOWNER
----          -
OSA2081L      ACTIVE      OSA2081L      YES
OSA20A0I 2    NOT ACTIVE  OSA2081L      NO
1 OF 1 RECORDS DISPLAYED
END OF THE REPORT

V TCPIP,TCPIPA,0,TCPIPA.TCPPARMS(OBDELINT) 3
EZZ0060I PROCESSING COMMAND: VARY TCPIP,TCPIPA,0,TCPIPA.TCPPARMS(OBDELINT)
EZZ0300I OPENED OBEYFILE FILE 'TCPIPA.TCPPARMS(OBDELINT)'
EZZ0309I PROFILE PROCESSING BEGINNING FOR 'TCPIPA.TCPPARMS(OBDELINT)'
EZZ0316I PROFILE PROCESSING COMPLETE FOR FILE 'TCPIPA.TCPPARMS(OBDELINT)'
EZZ0053I COMMAND VARY OBEY COMPLETED SUCCESSFULLY

D TCPIP,TCPIPA,N,DE,INTFN=OSA20A0I 4
0 OF 0 RECORDS DISPLAYED 4
END OF THE REPORT

```

In Example B-25 on page 433, the numbers correspond to the following information:

- 1** Stop the interface.
- 2** Check that the interface is not active.
- 3** Enter the OBEYFILE command.
- 4** Check that the interface has been deleted.

Example B-26 shows the TRL nodes definition in VTAMLST for OSA-Express 3.

Example B-26 TRL nodes definition in VTAMLST for OSA-Express 3

```

OSA20A0  VBUILD TYPE=TRL
OSA20A0P TRLE  LNCTL=MPC,          *
                READ=20A0,          1 *
                WRITE=20A1,         1 *
                DATAPATH=(20A2-20A7), 1 2 *
                PORTNAME=OSA20A0,    3 *
                PORTNUM=0,           4 *
                MPCLEVEL=QDIO
OSA20A1  VBUILD TYPE=TRL
OSA20A1P TRLE  LNCTL=MPC,          *
                READ=20A8,          1 *
                WRITE=20A9,         1 *
                DATAPATH=(20AA-20AE), 1 2 *

```

```

PORTNAME=OSA20A1,      3      *
PORTNUM=1,             4      *
MPCLEVEL=QDIO

```

In Example B-26, the numbers correspond to the following information:

- 1 OSA-Express 3 devices defined on the same CHPID (see Example B-29 on page 436).
- 2 Multiple DATAPATH device addresses to allow for multiple INTERFACE statements.
- 3 PORTNAME to be referenced in VTAM TRL node.
- 4 Port to be used on OSA-Express 3.

Note: The two TRLE resources associated with the two ports defined on the same CHPID can be defined on either the same or different TRL major nodes.

Example B-27 shows the TRL nodes.

Example B-27 Display TRL nodes

```

D NET,TRL,TRLE=OSA20A0P
IST097I DISPLAY ACCEPTED
IST075I NAME = OSA20A0P, TYPE = TRLE 003
IST1954I TRL MAJOR NODE = OSA20A0
IST486I STATUS= ACTIV, DESIRED STATE= ACTIV
IST087I TYPE = LEASED           , CONTROL = MPC , HPDT = YES
IST1715I MPCLEVEL = QDIO        MPCUSAGE = SHARE
IST2263I PORTNAME = OSA20A0     PORTNUM = 0 1 OSA CODE LEVEL = 000C
IST2337I CHPID TYPE = OSD       CHPID = 03
IST1577I HEADER SIZE = 4096 DATA SIZE = 0 STORAGE = ***NA***
IST1221I WRITE DEV = 20A1 STATUS = ACTIVE 2 STATE = ONLINE
IST1577I HEADER SIZE = 4092 DATA SIZE = 0 STORAGE = ***NA***
IST1221I READ  DEV = 20A0 STATUS = ACTIVE 2 STATE = ONLINE
IST924I -----
IST1221I DATA DEV = 20A2 STATUS = ACTIVE 2 STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST1717I ULPID = TCPIPA
IST2309I ACCELERATED ROUTING ENABLED
IST2331I QUEUE  QUEUE      READ
IST2332I ID     TYPE       STORAGE
IST2205I -----
IST2333I RD/1    PRIMARY   4.0M(64 SBALS)
IST2305I NUMBER OF DISCARDED INBOUND READ BUFFERS = 0
IST1757I PRIORITY1: UNCONGESTED PRIORITY2: UNCONGESTED
IST1757I PRIORITY3: UNCONGESTED PRIORITY4: UNCONGESTED
IST2190I DEVICEID PARAMETER FOR OSAENTA TRACE COMMAND = 01-01-00-02
IST1801I UNITS OF WORK FOR NCB AT ADDRESS X'0F3A4010'

..... Lines deleted
IST1221I DATA DEV = 20A3 STATUS = ACTIVE 2 STATE = N/A
..... Lines deleted
IST1221I DATA DEV = 20A4 STATUS = RESET 3 STATE = N/A
..... Lines deleted
IST314I END

D NET,TRL,TRLE=OSA20A1P

```

```

IST097I DISPLAY ACCEPTED
IST075I NAME = OSA20A1P, TYPE = TRLE
IST1954I TRL MAJOR NODE = OSA20A1
IST486I STATUS= ACTIV, DESIRED STATE= ACTIV
IST087I TYPE = LEASED , CONTROL = MPC , HPDT = YES
IST1715I MPCLEVEL = QDIO MPCUSAGE = SHARE
IST2263I PORTNAME = OSA20A1 PORTNUM = 1 1 OSA CODE LEVEL = 000C
IST1577I HEADER SIZE = 4096 DATA SIZE = 0 STORAGE = ***NA***
IST1221I WRITE DEV = 20A9 STATUS = ACTIVE 2 STATE = ONLINE
IST1577I HEADER SIZE = 4092 DATA SIZE = 0 STORAGE = ***NA***
IST1221I READ DEV = 20A8 STATUS = ACTIVE 2 STATE = ONLINE
IST1221I DATA DEV = 20AA STATUS = ACTIVE 2 STATE = N/A
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST1724I I/O TRACE = OFF TRACE LENGTH = *NA*
IST1717I ULPID = TCPIPA
..... Lines removed
IST1221I DATA DEV = 20AB STATUS = ACTIVE 2 STATE = N/A
..... Lines removed
IST1221I DATA DEV = 20AC STATUS = ACTIVE 2 STATE = N/A
..... Lines removed
IST1221I DATA DEV = 20AD STATUS = RESET 3 STATE = N/A
..... Lines removed
IST314I END

```

In Example B-27 on page 434, the numbers correspond to the following information:

- 1** OSA-Express 3 Port number.
- 2** Read, Write and Datapath device addresses in use. Multiple DATAPATH devices are needed if multiple INTERFACES and multiple VLANs are defined on the same OSA port.
- 3** Datapath device not in use.

Example B-28 shows the OSA-Express 3 devices online and allocated by NET.

Example B-28 Display OSA-Express 3 devices online and allocated by NET

```

D U,,,20A0,16
IEE457I 09.34.57 UNIT STATUS 249
UNIT TYPE STATUS      VOLSER  VOLSTATE
20A0 OSA  A-BSY
20A1 OSA  A
20A2 OSA  A-BSY
20A3 OSA  A-BSY
20A4 OSA  0
20A5 OSA  0
20A6 OSA  0
20A7 OSA  0
20A8 OSA  A-BSY
20A9 OSA  A
20AA OSA  A-BSY
20AB OSA  A-BSY
20AC OSA  A-BSY
20AD OSA  0
20AE OSA  0
20AF OSAD 0-RAL

```

Example B-29 shows the OSA-Express 3 CHPID.

Example B-29 Display OSA-Express 3 CHPID

```
D M=CHP(03)
IEE174I 14.59.34 DISPLAY M 586
CHPID 03: TYPE=11, DESC=OSA DIRECT EXPRESS, ONLINE
DEVICE STATUS FOR CHANNEL PATH 03
      0 1 2 3 4 5 6 7 8 9 A B C D E F
020A + + + + + + + + + + + + + + +
SWITCH DEVICE NUMBER = NONE
PHYSICAL CHANNEL ID = 0581
```

Note: All devices OSA-Express 3 of either port 0 and port 1 are defined under the same CHPID. Additional device addresses can be defined through HCD if required (see *OSA-Express Customer's Guide and Reference*, SA22-7935).

DEVICE and LINK statements

DEVICE and LINK statements now support features to support VLAN IDs and OFFLOAD processing to the OSA-Express adapter.

VLAN ID

Support is provided for virtual local area network standard IEEE 802.1q (VLAN). Implementing VLAN allows a physical LAN to be logically subdivided into separate logical LANs. With VLANID specified, the TCP/IP stacks that share an OSA can have an IP address assigned from different IP subnets.

It is configured and implemented in the z/OS environment through the LINK definitions in the PROFILE.TCPIP for OSA-Express in QDIO mode. VLANs support ARP takeover in a *flat network* (no routing protocol) when connected appropriately. Refer to Chapter 4, "Connectivity" on page 117 for more information about this implementation.

Example B-30 shows a link definition example of OSA2080I attached to virtual LAN 10.

Example B-30 Link definition example

```
INTERFACE OSA2080I
  DEFINE IPAQENET
  PORTNAME OSA2080
  IPADDR 10.1.2.11/24
  MTU 1492
  VLANID 10
  VMAC
```

Example B-31 displays the NETSTAT DEVLINKS display of an OSA-Express that has VLAN ID enabled.

Example B-31 VLAN ID enabled

```
D TCPIP,TCPIPA,N,DE
..... Lines deleted
INTFNAME: OSA2080I          INTFTYPE: IPAQENET  INTFSTATUS: READY
PORTNAME: OSA2080  DATAPATH: 2082  DATAPATHSTATUS: READY
SPEED: 0000001000
IPBROADCASTCAPABILITY: NO
VMACADDR: 020005749925  VMACORIGIN: OSA  VMACROUTER: ALL
ARPOFFLOAD: YES          ARPOFFLOADINFO: NO
CFGMTU: 1492             ACTMTU: 1492
IPADDR: 10.1.2.11/24
VLANID: 10 1             VLANPRIORITY: DISABLED
DYNVLANREGCFG: NO        DYNVLANREGCAP: YES
READSTORAGE: GLOBAL (4096K)  INBPERF: BALANCED
CHECKSUMOFFLOAD: YES       SEGMENTATIONOFFLOAD: YES
SECCLASS: 255            MONSYSPLEX: NO
..... Lines deleted
```

In this example, the numbers correspond to the following information:

- 1 VLAN tagging is enabled on this device (VLAN 10).

SRCIP

For inbound packets, the source IP address of a returning packet is always the destination IP address of a receiving packet. For outbound packets, the default source IP address is the HOME IP address of the interface chosen for sending the packet according to the routing table. If you specify IPCONFIG SOURCEVIPA, the source IP address is the first static VIPA listed above the interface chosen for sending the packet.

Alternatively you can designate the source IP addresses to be used for outbound TCP connections initiated by specified jobs or destined for specified IP addresses, networks, or subnets, by using the SRCIP statement, as described here:

- *Job-specific* source IP addressing by using the JOBNAME option in the SRCIP statement
- *Destination-specific* source IP addressing by using the DESTINATION option in the SRCIP statement

These source IP address definitions override any other source IP address specification in TCP/IP profile. However, the use of SRCIP can also be overridden directly by an application through the use of specific socket API options.

A distributed DVIPA cannot be specified as the source IP address on the DESTINATION statement. The TCP/IP client application issues a `connect()` socket call to start a TCP/IP connection, and optionally issues a `bind()` socket call prior to `connect()`. A problem occurs when a client application issues an explicit `bind()` socket call with `INADDR_ANY` and port 0 to have a port assigned prior to `connect()`. Until a `connect()` socket call that includes the destination IP address is issued, z/OS Communications Server is unable to determine the source IP address and therefore fails to choose which sysplex port pool the port should be assigned from.

You can relieve this restriction by adding the option `EXPLICITBINDPORTRANGE`. Unlike the sysplexport pools for which each pool is associated with a specific distributed DVIPA, the port range specified by `EXPLICITBINDPORTRANGE` is not associated with any specific distributed DVIPA, and can be used for any distributed DVIPA.

The EZBEPOR vv tt structure in the Coupling Facility, where vv is the 2-character VTAM group ID suffix specified on the XCFGRPID start option and tt is the TCP group ID suffix specified on the GLOBALCONFIG statement in the TCP/IP profile, coordinates this port range among all members of the sysplex. The port range should be identical in all members of the sysplex.

Note: The use of EXPLICITBINDPORTRANGE has a restriction in CINET environment. It is only available when the application has an affinity to a specific TCPIP stack, or when only one stack is managed by CINET.

Example B-32 shows a sample definition of the SRCIP statement.

Example B-32 SRCIP definition

```
GLOBALCONFIG EXPLICITBINDPORTRANGE 7000 1024 3

SRCIP
  JOBNAME      *           10.1.1.10 CLIENT 1
  JOBNAME      CUST*       10.1.2.10 SERVER 1
  DESTINATION  10.1.2.240   10.1.1.10 2
  DESTINATION  10.1.2.0/24  10.1.2.10 2
  DESTINATION  10.1.100.0/24 10.1.8.10 3
ENDSRCIP
```

In this example, the numbers correspond to the following information:

- 1** This is a sample definition of a job-specific source IP address. The the SERVER option listens to server applications while the CLIENT option indicates support for client applications, and is the default.
- 2** This is a sample definition of a destination-specific source IP address feature. The most specific match is applied.
- 3** This example uses a distributed DVIPA for the source IP address on the DESTINATION option. Define GLOBALCONFIG EXPLICITBINDPORTRANGE to reserve 1024 ports starting from 7000 for any distributed DVIPAs that are to be the source IP addresses. Ensure that the ports specified for EXPLICITBINDPORTRANGE are same among all sysplex members and do not overlap with any other port reservations: PORT, PORTRANGE, SYSPLEXPORTS, or BPXPRMxx INADDRANYPORT.

We used the NETSTAT,SRCIP command to verify our configuration, as shown in Example B-33.

Example B-33 NETSTAT SRCIP display

```
D TCPIP,TCPIPA,N,SRCIP
SOURCE IP ADDRESS BASED ON JOB NAME:
JOB NAME  TYPE  FLG  SOURCE
-----  ---  ---  -----
*         IPV4  C    10.1.1.10
CUST*     IPV4  S    10.1.2.10

SOURCE IP ADDRESS BASED ON DESTINATION
DESTINATION: 10.1.100.0/24
SOURCE:      10.1.8.10
DESTINATION: 10.1.2.240
SOURCE:      10.1.1.10
```

DESTINATION: 10.1.2.0/24
SOURCE: 10.1.2.10
5 OF 5 RECORDS DISPLAYED
END OF THE REPORT

To verify the destination-specific source IP address feature functions correctly, we issued the TSO Telnet command with an IP address configured in an L3 Switch. Example B-34 on page 439 shows results of the **show tcp brief** command issued for the L3 Switch.

Example B-34 Show tcp brief commands

```
telnet 10.1.2.240
Router1#sh tcp bri
TCB      Local Address      Foreign Address      (state)
46303D58  10.1.2.240.23            10.1.1.10.1036      ESTAB

telnet 10.1.2.220
Router2#sh tcp bri
TCB      Local Address      Foreign Address      (state)
423B1414  10.1.2.220.23            10.1.2.10.1037      ESTAB
```

We see a different source IP address is used for each specific destination IP address.

TCP/IP built-in security functions

z/OS Communications Server has built-in security functions that can be activated and used to control specific areas:

- ▶ Simple Mail Transfer Protocol (SMTP) provides a secure mail gateway option that allows an installation to create a database of registered network job entry (NJE) users who are allowed to send mail through SMTP to a TCP/IP network recipient.
- ▶ The FTP server gives you the opportunity to code security exits, in which you can extend control over the functions performed by the FTP server. Using these exits you can control:
 - The use of the FTP server based on IP addresses and port numbers
 - The use of the FTP server based on user IDs
 - The use of individual FTP subcommands
 - The submission of batch jobs through the FTP server
- ▶ FTP server logins act in the following ways:
 - You can configure the FTP server to restrict the users that can log in to the FTP server to only those users who are granted READ access to a resource profile in the SERVAUTH class.
 - When logging into the FTP server using the protected port (the port defined by the TLSPORT configuration statement), the FTP server and client initiate a TLS handshake without using the AUTH command. In previous releases, the FTP server had interoperability issues with non-z/OS FTP clients that connect to the protected port. With this enhancement, you can configure the FTP server to support non-z/OS FTP clients that connect to the protected port.
- ▶ z/OS Communications Server provides an SNMP agent that supports community-based security such as SNMPv1 and SNMPv2C, and user-based security such as SNMPv3. If you are concerned about sending SNMP data in a less secure environment, consider

implementing SNMPv3, whose messages have data integrity and data origin authentication.

Both the IMS sockets and CICS sockets support provide a user exit that you can use to validate each IMS or CICS transaction received by the Listener function. How you code this exit, and what data you require to be present in the transaction initiation request, is your decision.



Examples used in our environment

This appendix provides the examples that we used in the configuration of our environment.

Resolver

This section discusses how to set up the resolver. Example C-1 through Example C-6 on page 443 show the required procedures.

Resolver cataloged procedure

Example C-1 The resolver cataloged procedure

```
/* *****
/* SYS1.PROCLIB(RESOLV30)
/* *****
//RESOLV30 PROC PARM='CTRACE(CTIRES00)'
//EZBREINI EXEC PGM=EZBREINI,REGION=OM,TIME=1440,PARM=&PARMS
//*   SETUP contains resolver setup parameters.
//*   See the section on "Understanding Resolvers" in the
//*   IP Configuration Guide for more information. A sample of
//*   resolver setup parameters is included in member RESSETUP
//*   of the SEZAINST data set.
//*
//SETUP   DD   DSN=TCPIPA.TCPPARMS(RESOLV&SYSCZONE),DISP=SHR,FREE=CLOSE
```

BPXPRMxx

Example C-2 Specifying the resolver procedure to be started

```
/* *****
/* SYS1.PARMLIB(BPXPRM00)
/* *****
/* RESOLVER_PROC is used to specify how the resolver address space */
/* is processed during Unix System Services initialization.          */
/* The resolver address space is used by Tcp/Ip applications        */
/* for name-to-address or address-to-name resolution.              */
/* In order to create a resolver address space, a system must be    */
/* configured with an AF_INET or AF_INET6 domain.                  */
/* RESOLVER_PROC(procname|DEFAULT|NONE)                             */
/*   procname - The name of the address space for the resolver.     */
/*               In this case, this is the name of the address      */
/*               space as well as the procedure member name         */
/*               in SYS1.PROCLIB. procname is 1 to 8 characters     */
/*               long.                                              */
/*   DEFAULT - An address space with the name RESOLVER will        */
/*               be started. This is the same result that will      */
/*               occur if the RESOLVER_PROC statement is not        */
/*               specified in the BPXPRMxx profile.                 */
/*   NONE     - Specifies that a RESOLVER address space is         */
/*               not to be started.                                  */
/*                                                                 */
/*                                                                 @DAA*/
/* *****
RESOLVER_PROC(RESOLV30)
```

Resolver SETUP data set

Example C-3 Resolver address space SETUP data set

```
; *****
; TCPIPA.TCPPARMS(RESOLV30)
; *****
GLOBALTCPIPDATA('TCPIPA.TCPPARMS(GLOBAL)')
DEFAULTTCPIPDATA('TCPIPA.TCPPARMS(DEFAULT)')
GLOBALIPNODES('TCPIPA.TCPPARMS(IPNODES)')
COMMONSEARCH
```

Global TCPIP.DATA file

Example C-4 Global TCPIP.DATA file

```
; *****
; TCPIPA.TCPPARMS(GLOBAL)
; *****
DOMAINORIGIN  ITS0.IBM.COM
NSINTERADDR   10.12.6.7
NSPORTADDR    53
RESOLVEVIA    UDP
RESOLVERTIMEOUT 10
RESOLVERUDPRETRIES 1
LOOKUP        LOCAL DNS
```

Default TCPIP.DATA file

Example C-5 Default TCPIP.DATA file

```
; *****
; TCPIPA.TCPPARMS(DEFAULT)
; *****
TCPIPJOBNAME TCPIP
HOSTNAME WTSC30
```

Global ETC.IPNODES data set

Example C-6 GLOBALIPNODES data set

```
; *****
; TCPIPA.TCPPARMS(IPNODES)
; *****
10.12.6.7  OURDNS
10.1.1.10  WTSC30A
10.1.1.20  WTSC31B
10.1.1.30  WTSC32C
10.1.2.240 router1
```

TCP/IP stack

This section lists some examples that define the TCP/IP stack. Example C-7 through Example C-10 on page 447 show the required procedures.

TCPIPA stack started procedure

Example C-7 TCPIPA procedure

```
/* *****
/* SYS1.PROCLIB(TCPIPA)
/* *****
//TCPIPA    PROC  PARM='CTRACE(CTIEZB00),IDS=00',
//          PROFILE=PROFA&SYSCLONE.,TCPDATA=DATAA&SYSCLONE
//TCPIPA    EXEC  PGM=EZBTCPIP,REGION=OM,TIME=1440,
//          PARM=('&PARMS',
//          'ENVAR("RESOLVER_CONFIG=/' TCPIPA.TCPPARMS(&TCPDATA) '")')
//SYSPRINT DD  SYSOUT=*,DCB=(RECFM=VB,LRECL=132,BLKSIZE=136)
//ALGPRT DD  SYSOUT=*,DCB=(RECFM=VB,LRECL=132,BLKSIZE=136)
//CFGPRNT DD  SYSOUT=*,DCB=(RECFM=VB,LRECL=132,BLKSIZE=136)
//SYSOUT DD  SYSOUT=*,DCB=(RECFM=VB,LRECL=132,BLKSIZE=136)
//CEEDUMP DD  SYSOUT=*,DCB=(RECFM=VB,LRECL=132,BLKSIZE=136)
//SYSERROR DD  SYSOUT=*
//PROFILE DD  DISP=SHR,DSN=TCPIPA.TCPPARMS(&PROFILE.)
```

TCPIP.DATA file for TCPIPA stack

Example C-8 TCPIP.DATA file DATAA30

```
; *****
; TCPIPA.TCPPARMS(DATAA30)
; *****
TCPIPJOBNAME TCPIPA
HOSTNAME WTSC30A
DATASETPREFIX TCPIPA
MESSAGECASE MIXED
```

PROFILE.TCPIP (for static routing)

Example C-9 PROFILE.TCPIP (for static routing)

```
; This profile is for static routing
ARPAGE 20
;
GLOBALCONFIG NOTCPIPSTATISTICS IQDMULTIWRITE ZIIP IQDIOMULTIWRITE
GLOBALCONFIG SEGMENTATIONOFFLOAD
GLOBALCONFIG EXPLICITBINDPORTRANGE 7000 3
GLOBALCONFIG XCFGRPID 21 IQDVLANID 21
GLOBALCONFIG SYSPLEXMONITOR DELAYJOIN RECOVERY TIMERSECS 60
;
IPCONFIG
  ARPTO 1200
  SOURCEVIPA
  IGNOREREDIRECT
  DATAGRAMFWD
  SYSPLEXROUTING
  MULTIPATH PERCONNECTION
```

```

PATHMTUDISCOVERY
DYNAMICXCF 10.1.7.11    255.255.255.0 8
;
TCPCONFIG TCPSENDBFRSIZE 256K TCPRCVBUFRSIZE 256K SENDGARBAGE FALSE
TCPCONFIG TCPMAXRCVBUFRSIZE 512K
TCPCONFIG UNRESTRICTLOWPORTS
;
UDPCONFIG UNRESTRICTLOWPORTS UDPSENDBFRSIZE 65535 UDPRCVBUFRSIZE 65535
UDPCONFIG NOUDPQUEUELIMIT
;
NETMONITOR SMFSERVICE
;
SOMAXCONN 10
;
;OSA definitions
;TRL MAJ NODE: OSA2080,OSA20A0,OSA20C0,and OSA20E0
;
INTERFACE OSA2080I
DEFINE IPAQENET
PORTNAME OSA2080
IPADDR 10.1.2.11/24
MTU 1492
VLANID 10
VMAC
;
INTERFACE OSA20A0I
DEFINE IPAQENET
PORTNAME OSA20A0
IPADDR 10.1.2.12/24
MTU 1492
VLANID 10
VMAC
;
INTERFACE OSA20C0I
DEFINE IPAQENET
PORTNAME OSA20C0
IPADDR 10.1.3.11/24
MTU 1492
VLANID 11
VMAC
;
INTERFACE OSA20E0I
DEFINE IPAQENET
PORTNAME OSA20E0
IPADDR 10.1.3.12/24
MTU 1492
VLANID 11
VMAC
;
;HiperSockets definitions
DEVICE IUTIQDF4 MPCIPA
LINK IUTIQDF4L IPAQIDIO IUTIQDF4
DEVICE IUTIQDF5 MPCIPA
LINK IUTIQDF5L IPAQIDIO IUTIQDF5
DEVICE IUTIQDF6 MPCIPA

```

```

LINK    IUTIQDF6L  IPAQIDIO    IUTIQDF6
;
;Static VIPA definitions
DEVICE VIPA1      VIRTUAL 0
LINK    VIPA1L    VIRTUAL 0    VIPA1
DEVICE VIPA2      VIRTUAL 0
LINK    VIPA2L    VIRTUAL 0    VIPA2
;
;Dynamic VIPA definitions
VIPADYNAMIC
    VIPADEFINE     MOVE IMMED 255.255.255.0 10.1.8.25    ;FTP
    VIPADISTRIBUTE DEFINE SYSPLEXPORTS DISTMETHOD BASEWLM
                        10.1.8.25    PORT 20 21
                        DESTIP 10.1.7.11
                        10.1.7.21
    VIPADEFINE     MOVE IMMED 255.255.255.0 10.1.8.21    ;General
    VIPADISTRIBUTE DEFINE SYSPLEXPORTS DISTMETHOD ROUNDROBIN
                        10.1.8.21    PORT 992 20 21 23
                        DESTIP 10.1.7.11
                        10.1.7.21
    VIPADEFINE     MOVE IMMED 255.255.255.0 10.1.8.44    ;General
    VIPADISTRIBUTE DEFINE SYSPLEXPORTS DISTMETHOD ROUNDROBIN
                        10.1.8.44    PORT 4444
                        DESTIP 10.1.7.11
                        10.1.7.21
    VIPADEFINE     MOVE IMMED 255.255.255.0 10.1.8.22    ;Admin
    VIPADISTRIBUTE DEFINE SYSPLEXPORTS DISTMETHOD BASEWLM
                        10.1.8.22    PORT 992 20 21
                        DESTIP 10.1.7.11
                        10.1.7.21
    VIPADEFINE     MOVE IMMED 255.255.255.0 10.1.8.23    ;Payro1
    VIPADISTRIBUTE DEFINE SYSPLEXPORTS DISTMETHOD SERVERWLM
                        10.1.8.23    PORT 992 20 21
                        DESTIP 10.1.7.11
                        10.1.7.21
    VIPABACKUP 200 MOVE IMMED 255.255.255.0 10.1.8.24    ; EXTRAS
    VIPADISTRIBUTE DEFINE SYSPLEXPORTS DISTMETHOD SERVERWLM
                        10.1.8.24    PORT 20 21
                        DESTIP 10.1.7.11
                        10.1.7.21
    VIPAROUTE      DEFINE 10.1.7.11    10.1.1.10    ; sc30's static vipa
ENDVIPADYNAMIC
;
HOME
    10.1.1.10      VIPA1L
    10.1.2.10      VIPA2L
    10.1.4.11      IUTIQDF4L
    10.1.5.11      IUTIQDF5L
    10.1.6.11      IUTIQDF6L
;
    PRIMARYINTERFACE VIPA1L
;
BEGINRoutes
; Direct Routes - Routes that are directly connected to the interfaces
;   Destination      Subnet Mask    First Hop Link Name    Packet Size

```

```

ROUTE 10.10.4.0      255.255.255.0 =      IUTIQDF4L      mtu defaultsize repl
ROUTE 10.10.4.0/24      =      IUTIQDF5L      mtu defaultsize repl
ROUTE 10.10.5.0/24      =      IUTIQDF6L      mtu defaultsize repl
; Default Route - All packets to an unknown destination are routed
;                      through this route.
;      Destination      First Hop      Link Name      Packet Size
ROUTE DEFAULT          10.10.2.1      OSA2080I      mtu defaultsize repl
ROUTE DEFAULT          10.10.3.1      OSA20A0I      mtu defaultsize repl
ROUTE DEFAULT          10.10.2.2      OSA20C0I      mtu defaultsize repl
ROUTE DEFAULT          10.10.3.2      OSA20E0I      mtu defaultsize repl
ENDRoutes
;
AUTOLOG 5
      FTPDA  JOBNAME FTPDA1
      OMPA
ENDAUTOLOG
;
PORT
      20 TCP OMVS      NOAUTOLOG      ; FTP Server  1
      21 TCP FTPDA1  BIND 10.1.1.10  ; control port
      23 TCP TN3270A      ; Telnet Server
      25 TCP SMTP      ; SMTP Server
      500 UDP IKED      ; @ADI
      514 UDP OMVS      ; Remote Execution Server
      520 UDP OMPA      NOAUTOLOG      ; OMPROUTE RIPV2 port
      521 UDP OMPA      NOAUTOLOG      ; OMPROUTE RIPV2 port
      4500 UDP IKED
;                      ; @ADI
PORTRANGE 10000 2000 TCP OMVS ; TCP 10000 - 11999
PORTRANGE 10000 2000 UDP OMVS ; UDP 10000 - 11999
;
PORT UNRSV TCP * DENY
;
SACONFIG ENABLED COMMUNITY j0s9m2ap AGENT 161
SMFCONFIG TYPE119 TCPINIT TCPTERM TCPIPSTATISTICS TN3270CLIENT
      FTPCLIENT
;
START OSA2080I
START OSA2081I
START OSA20A0I
START OSA20C0I
START OSA20E0I
START IUTIQDF4
START IUTIQDF5
START IUTIQDF6

```

PROFILE.TCPIP (for OMPROUTE dynamic routing)

Example C-10 PROFILE.TCPIP (for OMPROUTE dynamic routing)

```

; This profile is for static routing
ARPAGE 20
;
GLOBALCONFIG NOTCPIPSTATISTICS IQDMULTIWRITE ZIIP IQDIOMULTIWRITE
GLOBALCONFIG SEGMENTATIONOFFLOAD
GLOBALCONFIG EXPLICITBINDPORTRANGE 7000 3

```

```

GLOBALCONFIG XCFGRPID 21 IQDVLANID 21
GLOBALCONFIG SYSPLEXMONITOR DELAYJOIN    RECOVERY TIMERSECS 60
;
IPCONFIG
    ARPTO      1200
    SOURCEVIPA
    IGNOREREDIRECT
    DATAGRAMFWD
    SYSPLEXROUTING
    MULTIPATH PERCONNECTION
    PATHMTUDISCOVERY
    DYNAMICXCF 10.1.7.11    255.255.255.0 8
;
TCPCONFIG TCPSENDERBFRSIZE 256K TCPCVBUFRSIZE 256K SENDGARBAGE FALSE
TCPCONFIG TCPMAXRCVBUFRSIZE 512K
TCPCONFIG UNRESTRICTLOWPORTS
;
UDPCONFIG UNRESTRICTLOWPORTS UDPSENDERBFRSIZE 65535 UDPCVBUFRSIZE 65535
UDPCONFIG NOUDPQUEUELIMIT
;
NETMONITOR SMFSERVICE
;
SOMAXCONN 10
;
;OSA definitions
;TRL MAJ NODE: OSA2080,OSA20A0,OSA20C0,and OSA20E0
;
INTERFACE OSA2080I
DEFINE IPAQENET
PORTNAME OSA2080
IPADDR 10.1.2.11/24
MTU 1492
VLANID 10
VMAC
;
INTERFACE OSA20A0I
DEFINE IPAQENET
PORTNAME OSA20A0
IPADDR 10.1.2.12/24
MTU 1492
VLANID 10
VMAC
;
INTERFACE OSA20C0I
DEFINE IPAQENET
PORTNAME OSA20C0
IPADDR 10.1.3.11/24
MTU 1492
VLANID 11
VMAC
;
INTERFACE OSA20E0I
DEFINE IPAQENET
PORTNAME OSA20E0
IPADDR 10.1.3.12/24

```



```

MTU 1492
VLANID 11
VMAC
;
;HiperSockets definitions
DEVICE IUTIQDF4 MPCIPA
LINK IUTIQDF4L IPAQIDIO IUTIQDF4
DEVICE IUTIQDF5 MPCIPA
LINK IUTIQDF5L IPAQIDIO IUTIQDF5
DEVICE IUTIQDF6 MPCIPA
LINK IUTIQDF6L IPAQIDIO IUTIQDF6
;
;Static VIPA definitions
DEVICE VIPA1 VIRTUAL 0
LINK VIPA1L VIRTUAL 0 VIPA1
DEVICE VIPA2 VIRTUAL 0
LINK VIPA2L VIRTUAL 0 VIPA2
;
;Dynamic VIPA definitions
VIPADYNAMIC
    VIPADEFINE MOVE IMMED 255.255.255.0 10.1.8.25 ;FTP
    VIPADISTRIBUTE DEFINE SYSPLEXPORTS DISTMETHOD BASEWLM
                    10.1.8.25 PORT 20 21
                    DESTIP 10.1.7.11
                    10.1.7.21
    VIPADEFINE MOVE IMMED 255.255.255.0 10.1.8.21 ;General
    VIPADISTRIBUTE DEFINE SYSPLEXPORTS DISTMETHOD ROUNDROBIN
                    10.1.8.21 PORT 992 20 21 23
                    DESTIP 10.1.7.11
                    10.1.7.21
    VIPADEFINE MOVE IMMED 255.255.255.0 10.1.8.44 ;General
    VIPADISTRIBUTE DEFINE SYSPLEXPORTS DISTMETHOD ROUNDROBIN
                    10.1.8.44 PORT 4444
                    DESTIP 10.1.7.11
                    10.1.7.21
    VIPADEFINE MOVE IMMED 255.255.255.0 10.1.8.22 ;Admin
    VIPADISTRIBUTE DEFINE SYSPLEXPORTS DISTMETHOD BASEWLM
                    10.1.8.22 PORT 992 20 21
                    DESTIP 10.1.7.11
                    10.1.7.21
    VIPADEFINE MOVE IMMED 255.255.255.0 10.1.8.23 ;Payrol
    VIPADISTRIBUTE DEFINE SYSPLEXPORTS DISTMETHOD SERVERWLM
                    10.1.8.23 PORT 992 20 21
                    DESTIP 10.1.7.11
                    10.1.7.21
    VIPABACKUP 200 MOVE IMMED 255.255.255.0 10.1.8.24 ; EXTRAS
    VIPADISTRIBUTE DEFINE SYSPLEXPORTS DISTMETHOD SERVERWLM
                    10.1.8.24 PORT 20 21
                    DESTIP 10.1.7.11
                    10.1.7.21
    VIPAROUTE DEFINE 10.1.7.11 10.1.1.10 ; sc30's static vipa
ENDVIPADYNAMIC
;
HOME
    10.1.1.10 VIPA1L

```

```

10.1.2.10      VIPA2L
10.1.4.11      IUTIQDF4L
10.1.5.11      IUTIQDF5L
10.1.6.11      IUTIQDF6L
;
PRIMARYINTERFACE VIPA1L
;
AUTOLOG 5
  FTPDA  JOBNAME FTPDA1
  OMPA
ENDAUTOLOG
;
PORT
  20 TCP OMVS      NOAUTOLOG      ; FTP Server  1
  21 TCP FTPDA1 BIND 10.1.1.10 ; control port
  23 TCP TN3270A      ; Telnet Server
  25 TCP SMTP          ; SMTP Server
  500 UDP IKED          ; @ADI
  514 UDP OMVS          ; Remote Execution Server
  520 UDP OMPA      NOAUTOLOG      ; OMPROUTE RIPV2 port
  521 UDP OMPA      NOAUTOLOG      ; OMPROUTE RIPV2 port
  4500 UDP IKED
;
; @ADI
PORTRANGE 10000 2000 TCP OMVS ; TCP 10000 - 11999
PORTRANGE 10000 2000 UDP OMVS ; UDP 10000 - 11999
;
PORT UNRSV TCP * DENY
;
SACONFIG ENABLED COMMUNITY j0s9m2ap AGENT 161
SMFCONFIG TYPE119 TCPINIT TCPTERM TCPIPSTATISTICS TN3270CLIENT
FTPCLIENT
;
START OSA2080I
START OSA2081I
START OSA20A0I
START OSA20C0I
START OSA20E0I
START IUTIQDF4
START IUTIQDF5
START IUTIQDF6

```

OMPROUTE dynamic routing

These are the complete examples (Example C-11 through Example C-17 on page 453) that we used in our environment and discussed in Chapter 5, “Routing” on page 205.

OMPROUTE cataloged procedure

Example C-11 OMPROUTE cataloged procedure

```

//OMPA30 PROC STDENV=OMPENA&SYSCLONE
//OMPA30 EXEC PGM=OMPROUTE,REGION=OM,TIME=NOLIMIT,
//          PARM=('POSIX(ON) ALL31(ON)',
//          'ENVAR("_BPXK_SETIBMOPT_TRANSPORT=TCPIPA"',

```

```
//          '"_CEE_ENVFILE=DD:STDENV")/')
```

```
//STDENV DD DISP=SHR,DSN=TCPIP.SC&SYSCONE..STDENV(&STDENV)
//SYSOUT DD SYSOUT=*
//OMPCFG DD DSN=TCPIPA.TCPPARMS(OMPA&SYSCONE.),DISP=SHR
//CEEDUMP DD SYSOUT=*,DCB=(RECFM=FB,LRECL=132,BLKSIZE=132)
```

OMPROUTE environment variables

Example C-12 OMPROUTE environment variables

```
; *****
; TCPIP.SC30.STDENV(OMPENA30)
; *****
RESOLVER_CONFIG=/'TCPIPA.TCPPARMS(DATAA&SYSCONE.)'
;OMPROUTE_FILE=/'TCPIPA.TCPPARMS(OMPA30)'
OMPROUTE_DEBUG_FILE=/etc/omproute/debug30a
OMPROUTE_DEBUG_FILE_CONTROL=100000,5
```

TCPIP.DATA file

Example C-13 Global TCPIP.DATA file

```
; *****
; TCPIPA.TCPPARMS(GLOBAL)
; *****
DOMAINORIGIN ITS0.IBM.COM
NSINTERADDR 10.12.6.7
NSPORTADDR 53
RESOLVEVIA UDP
RESOLVERTIMEOUT 10
RESOLVERUDPRETRIES 1
LOOKUP LOCAL DNS
```

Example C-14 Local TCPIP.DATA file

```
; *****
; TCPIPA.TCPPARMS(DATAA30)
; *****
TCPIPJOBNAME TCPIPA
HOSTNAME WTSC30A
DATASETPREFIX TCPIPA
MESSAGECASE MIXED
```

Syslogd configuration file

Example C-15 Syslogd configuration file

```
##*****
#*
#* syslog.conf - Defines the actions to be taken for the specified *
#* facilities/priorities by the syslogd daemon. *
#* *
#*.OMPA*.*.err /tmp/syslog/ompa.err.log
```

OMPROUTE configuration file

Example C-16 OMPROUTE configuration file

```
Area Area_Number=0.0.0.2
  Stub_Area=YES
  Authentication_type=None
  Import_Summaries=Yes;
OSPF
  RouterID=10.1.30.10
  Comparison=Type2
  DR_Max_Adj_Attempt = 10
  Demand_Circuit=YES;
Global_Options
  Ignore_Undefined_Interfaces=YES;
;
; Static vipa
OSPF_interface ip_address=10.1.1.10
  name=VIPA1L
  subnet_mask=255.255.255.0
  Advertise_VIPA_Routes=HOST_ONLY
  attaches_to_area=0.0.0.2
  cost0=10
  mtu=65535;
OSPF_interface ip_address=10.1.2.10
  name=VIPA2L
  subnet_mask=255.255.255.0
  Advertise_VIPA_Routes=HOST_ONLY
  attaches_to_area=0.0.0.2
  cost0=10
  mtu=65535;
; OSA Qdio VLAN10
OSPF_Interface IP_address=10.1.2.*
  Subnet_mask=255.255.255.0
  Router_Priority=0
  Attaches_To_Area=0.0.0.2
  Cost0=100
  MTU=1492;
; OSA Qdio VLAN11
OSPF_Interface IP_address=10.1.3.*
  Subnet_mask=255.255.255.0
  Router_Priority=0
  Attaches_To_Area=0.0.0.2
  Cost0=100
  MTU=1492;
; Hipersockets 10.1.4.x
OSPF_Interface IP_address=10.1.4.*
  Subnet_mask=255.255.255.0
  Router_Priority=1
  Attaches_To_Area=0.0.0.2
  Cost0=90
  MTU=8192;
; Hipersockets 10.1.5.x
OSPF_Interface IP_address=10.1.5.*
  Subnet_mask=255.255.255.0
  Router_Priority=1
  Attaches_To_Area=0.0.0.2
```

```

        Cost0=90
        MTU=8192;
; Hipersockets 10.1.6.x
OSPF_Interface IP_address=10.1.6.*
        Subnet_mask=255.255.255.0
        Router_Priority=1
        Attaches_To_Area=0.0.0.2
        Cost0=90
        MTU=8192;
;Dynamic XCF
interface ip_address=10.1.7.*
        subnet_mask=255.255.255.0
        mtu=65535;
; Dynamic vipa VIPADEFINE
ospf_interface ip_address=10.1.8.*
        subnet_mask=255.255.255.0
        Advertise_VIPA_Routes=HOST_ONLY
        attaches_to_area=0.0.0.2
        cost0=10
        mtu=65535;
; Dynamic vipa VIPADEFINE
ospf_interface ip_address=10.1.9.*
        subnet_mask=255.255.255.0
        Advertise_VIPA_Routes=HOST_ONLY
        attaches_to_area=0.0.0.2
        cost0=10
        mtu=65535;

```

Router configuration

Example C-17 Router configuration

```

interface Loopback1
 ip address 10.1.200.1 255.255.255.0
!
interface Vlan10
 ip address 10.1.2.240 255.255.255.0
 ip ospf cost 100
 ip ospf priority 100
!
interface Vlan11
 ip address 10.1.3.240 255.255.255.0
 ip ospf cost 100
 ip ospf priority 100
!
router ospf 100
 router-id 10.1.3.240
 log-adjacency-changes
 area 2 stub no-summary
 network 10.1.2.0 0.0.0.255 area 2
 network 10.1.3.0 0.0.0.255 area 2
 network 10.1.100.0 0.0.0.255 area 2
 network 10.200.1.0 0.0.0.255 area 0
 default-information originate always metric-type 1

```



D

Our implementation environment

We wrote the four *z/OS Communications Server TCP/IP Implementation* books at the same time. Given the complexity of this project, we needed to be creative in organizing the test environment so that each team could work with minimal coordination and interference from the other teams. In this appendix, we show the complete environment that we used for the four books as well as the environment that we used for this book.

The environment used for all four books

To enable concurrent work on each of the four books, we set up and shared the test environment illustrated in Figure D-1.

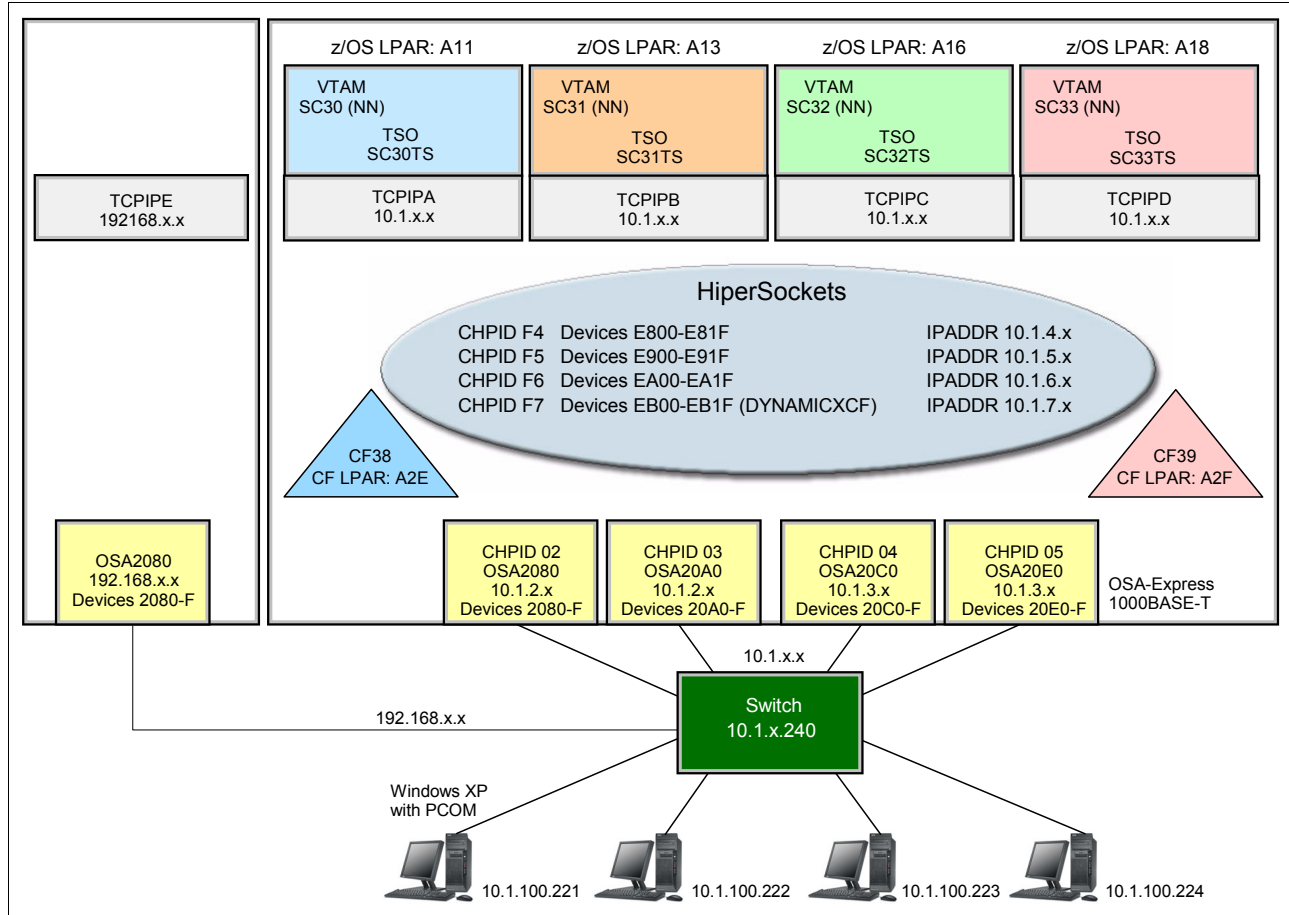


Figure D-1 Our implementation environment

We wrote our books (and ran our implementation scenarios) using four logical partitions (LPARs) on an IBM System z196-32 (referred to as LPARs A11, A13, A16, and A18). We implemented and started one TCP/IP stack on each LPAR. Each LPAR shared the following resources:

- ▶ HiperSockets inter-server connectivity
- ▶ Coupling Facility connectivity (CF38 and CF39) for Parallel Sysplex scenarios
- ▶ Eight OSA-Express3 1000BASE-T Ethernet ports connected to a switch

Finally, we shared four Windows workstations, representing corporate network access to the z/OS networking environment. The workstations are connected to the switch. For verifying our scenarios, we used applications such as TN3270 and FTP.

The IP addressing scheme that we used allowed us to build multiple subnetworks so that we would not impede ongoing activities from other team members.

VLANs were also defined to isolate the TCP/IP stacks and portions of the LAN environment (see Figure D-2).

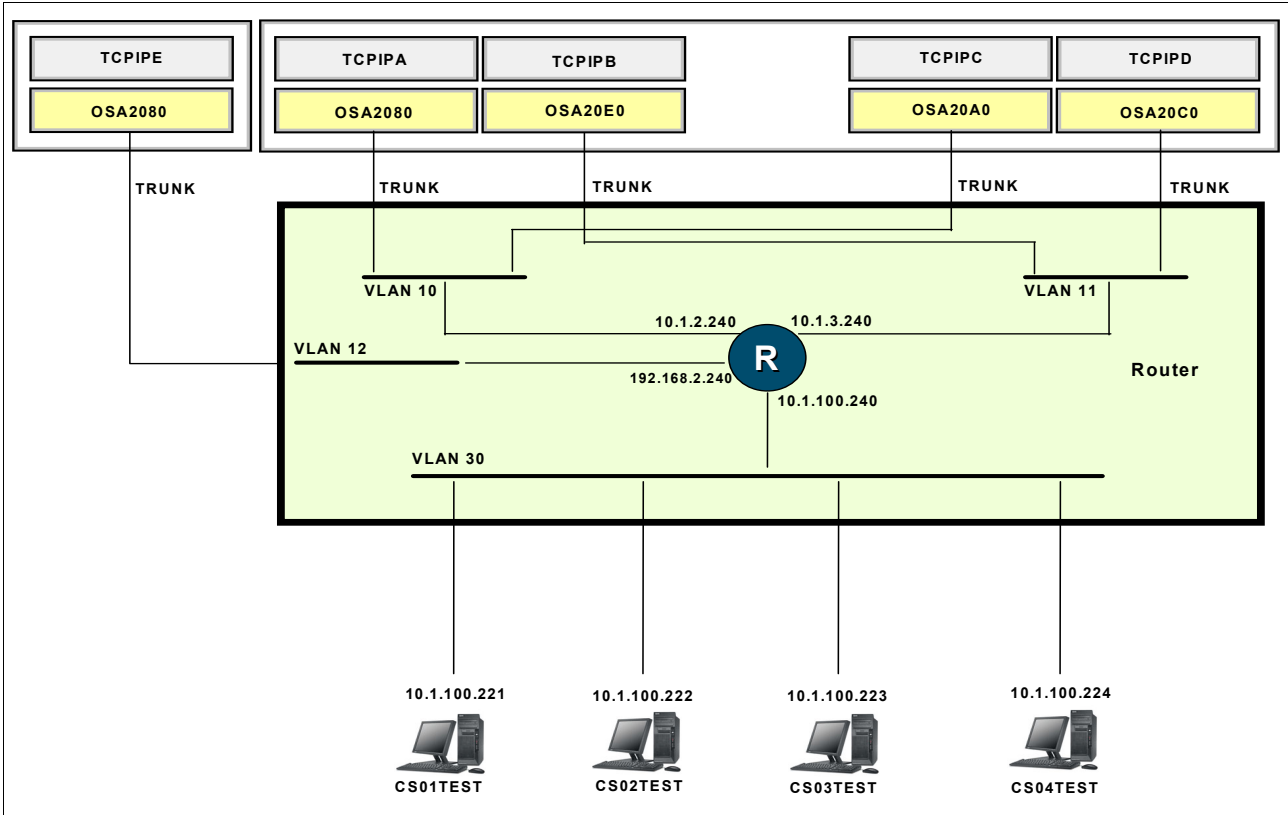


Figure D-2 LAN configuration: VLAN and IP addressing

Our focus for this book

Figure D-3 depicts the environment that we worked with, as required for our basic function implementation scenarios.

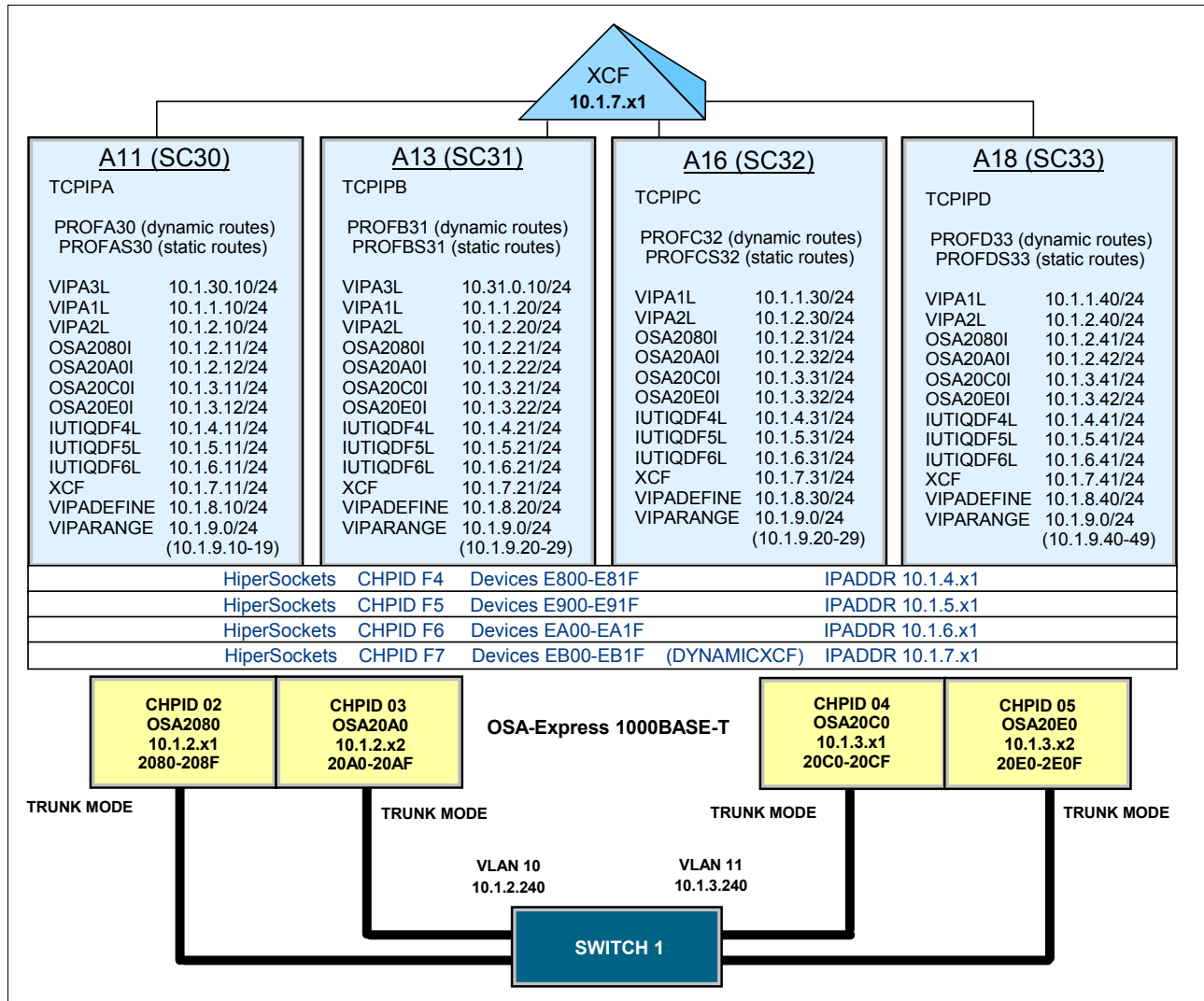


Figure D-3 Our environment for this book

Related publications

The publications listed in this section are considered particularly suitable for a more detailed discussion of the topics covered in this book.

IBM Redbooks publications

For information about ordering these publications, see “How to get IBM Redbooks publications” on page 461. Note that some of the documents referenced here might be available in softcopy only.

- ▶ *IBM z/OS Communications Server TCP/IP Implementation Volume 1: Base Functions, Connectivity, and Routing*, SG24-7896
- ▶ *IBM z/OS Communications Server TCP/IP Implementation Volume 2: Standard Applications*, SG24-7897
- ▶ *IBM z/OS Communications Server TCP/IP Implementation Volume 3: High Availability, Scalability, and Performance*, , SG24-7898
- ▶ *Communications Server for z/OS TCP/IP Implementation Volume 4: Security and Policy-Based Networking*, SG24-7899
- ▶ *IP Network Design Guide*, SG24-2580
- ▶ *HiperSockets Implementation Guide*, SG24-6816
- ▶ *Migrating Subarea Networks to an IP Infrastructure Using Enterprise Extender*, SG24-5957
- ▶ *OSA-Express Implementation Guide*, SG24-5948
- ▶ *SNA in a Parallel Sysplex Environment*, SG24-2113
- ▶ *TCP/IP Tutorial and Technical Overview*, GG24-3376
- ▶ *z/OS 1.6 Security Services Update*, SG24-6448
- ▶ *z/OS Infoprint Server Implementation*, SG24-6234

Other publications

These publications are also relevant as further information sources:

- ▶ *IBM Health Checker for z/OS: User's Guide*, SA22-7994
- ▶ *OSA-Express Customer's Guide and Reference*, SA22-7935
- ▶ *z/OS Communications Server: CSM Guide*, SC31-8808
- ▶ *z/OS Communications Server: IP Configuration Guide*, SC31-8775
- ▶ *z/OS Communications Server: IP Configuration Reference*, SC31-8776
- ▶ *z/OS Communications Server: IP Diagnosis Guide*, GC31-8782
- ▶ *z/OS Communications Server: IP Messages Volume 1 (EZA)*, SC31-8783
- ▶ *z/OS Communications Server: IP Messages Volume 2 (EZB, EZD)*, SC31-8784
- ▶ *z/OS Communications Server: IP Messages Volume 3 (EZY)*, SC31-8785

- ▶ *z/OS Communications Server: IP Messages Volume 4 (EZZ, SNM)*, SC31-8786
- ▶ *z/OS Communications Server: IP Programmer's Guide and Reference*, SC31-8787
- ▶ *z/OS Communications Server: IP and SNA Codes*, SC31-8791
- ▶ *z/OS Communications Server: IP Sockets Application Programming Interface Guide and Reference*, SC31-8788
- ▶ *z/OS Communications Server: IP System Administrator's Commands*, SC31-8781
- ▶ *z/OS Communications Server: IP User's Guide and Commands*, SC31-8780
- ▶ *z/OS Communications Server: IPv6 Network and Application Design Guide*, SC31-8885
- ▶ *z/OS Communications Server: New Function Summary*, GC31-8771
- ▶ *z/OS Communications Server: Quick Reference*, SX75-0124
- ▶ *z/OS Communications Server: SNA Network Implementation*, SC31-8777
- ▶ *z/OS Communications Server: SNA Operation*, SC31-8779
- ▶ *z/OS Communications Server: SNA Resource Definition*, SC31-8778
- ▶ *z/OS Migration*, GA22-7499
- ▶ *z/OS MVS IPCS Commands*, SA22-7594
- ▶ *z/OS MVS System Commands*, SA22-7627
- ▶ *z/OS TSO/E Command Reference*, SA22-7782
- ▶ *z/OS UNIX System Services Command Reference*, SA22-7802
- ▶ *z/OS UNIX System Services Programming: Assembler Callable Services Reference*, SA22-7803
- ▶ *z/OS UNIX System Services Command Reference*, SA22-7802
- ▶ *z/OS UNIX System Services File System Interface Reference*, SA22-7808
- ▶ *z/OS UNIX System Services Messages and Codes*, SA22-7807
- ▶ *z/OS UNIX System Services Parallel Environment: Operation and Use*, SA22-7810
- ▶ *z/OS UNIX System Services Planning*, GA22-7800
- ▶ *z/OS UNIX System Services Programming Tools*, SA22-7805
- ▶ *z/OS UNIX System Services User's Guide*, SA22-7801
- ▶ *z/OS XL C/C++ Run-Time Library Reference*, SA22-7821

Online resources

These Web sites are also relevant as further information sources:

- ▶ Mainframe networking
<http://www.ibm.com/servers/eserver/zseries/networking/>
- ▶ z/OS Communications Server product overview
<http://www.ibm.com/software/network/commserver/zos/>
- ▶ z/OS Communications Server product support
<http://www-306.ibm.com/software/network/commserver/zos/support/>
- ▶ z/OS Communications Server publications
<http://www-03.ibm.com/systems/z/os/zos/bkserv/r9pdf/#commserv>

How to get IBM Redbooks publications

You can search for, view, or download Redbooks, Redpapers, Technotes, draft publications and Additional materials, as well as order hardcopy Redbooks publications, at this Web site:

ibm.com/redbooks

Help from IBM

IBM Support and downloads

ibm.com/support

IBM Global Services

ibm.com/services

Index

Numerics

6bone test address 394
6to4 address 394

A

access mode 275
access port 274
accessing the z/OS UNIX shells 13
accounting 358
active state 405
address registration 121
Address Resolution Protocol (ARP) 279, 396
address space 8, 37, 71, 420
 abends 332
 child process 11
 name 426
adjacencies, lost 215
adjacent router 219, 225
 router advertisements 225
AF_INET 14–15
AF_INET addressing family 15
AF_INET socket 15
AF_INET transport provider 16
AF_UNIX 14
AF_UNIX addressing family 14
AF_UNIX socket 14
anchor configuration data set 86
Anycast address 394
 simple example 394
APF authorization 67
APPC/MVS 11
ARP cache 200, 420
ARP offload 122
ARPTO 428
assigning user IDs to started tasks 65
attached TCP/IP network
 other TCP/IP server host 74
auditing 358
Authorized Program Facility (APF) 67
autoconfigured addresses, public 391
autoconfigured addresses, temporary 391
AUTOLOG considerations 82
AUTOLOG statement 245
 OMPROUTE start procedure 245
 procedure name 84

B

base functions 59, 371
 common design scenarios 60, 372
batch job 2, 11, 60, 439
BEGINROUTES 214, 228
BEGINROUTES statement 396
BIND control for INADDR_ANY 424

BPXPRMxx 69, 73
BPXPRMxx definitions 400
BPXPRMxx member 11, 63, 400
 incorrect value 73
BPXPRMxx, CINET 16
BPXTIINT 73
BSDROUTINGPARMS statement 134, 183, 213
BUFFERPOOL statements 419

C

cache 28
CACHESIZE 29
canonical mode 14
capacity planning 357
CEEPRMxx 237
certain periodic updates (CPU) 61, 211
channel-to-channel (CTC) 118
CHPID 181
CICS 440
CINET 16
C-INET pre-router 37
CMD 251
command TRACE (CT) 258, 327
Common INET Physical File System (CINET) 16
Common Link Access to Workstation (CLAW) 118
Common Storage 419
Common Storage Manager (CSM) 2, 108
COMPONENT TRACE
 Status 56
 TALLY Report 57
component trace (CT) 202, 258, 315, 327, 356
 primary purpose 333
configuration data set
 selecting 63
configuration data set names 79
configuration examples 47
configuration statement 22, 213, 413
configure default TCPIP.DATA data set statements 46
configure the setup data set 44
configure the SITE table (HOSTS.LOCAL) 86
configure z/OS TCP/IP 79
connected stacks 37
Connection Isolation 158
connectivity 118
 common design scenarios 136
connectivity status
 verifying 152, 183, 190
connectivity, lost 216
controlling the device definitions 195
CPU cycle 17, 41, 61–62, 136, 211
CPU resources 64
creating multiple stacks 87
cron 11
Cross-System Coupling Facility (XCF) 8

- CSA 72, 419
- CTIEZB00 103, 315, 356
- CTIEZBxx 320
- CTRACE 315, 356
- CTRACE -- RESOLVER (SYSTCPRE) 55
- CTRACE command 332, 348
- CTWTR procedure 55, 259
- customization levels of UNIX System Services 11

D

- daemons 11
- Data Facility Storage Management Subsystem (DFSMS) 3, 60
- Data Link Control (DLC) 5, 118, 401
- data set 10, 67–68, 237, 316
 - configuration 63
 - trace data 316
- DATASETPREFIX 79, 86
- DD card 42
- DDVIPA 268
- dedicated data links 387
- default destination address selection 38
- default directory path 76
- default route 207
- default search path and symbolic links 76
- default user 13
- dependencies 61–62, 140, 158–159, 181, 188, 234, 387–389
- designated router (DR) 220
- destination address 38–39, 197, 208, 394
- device
 - adding, changing, and deleting 110
- DEVICE AND LINK 436
- DEVICE and LINK statements for each OSA-Express port 149
- DEVICE and LINK statements for HiperSockets CHPID 182
- device OSA2080 85, 92, 403
- DEVICE statement 401
- device status in TCP/IP stack
 - verifying 152, 184, 190
- device type 121
- DEVICE/LINK statement 279
- diagnosing a OMPROUTE problem 256
- diagnosing the resolver address space environment 50
- differences between RIPng and RIP-2 226
- Direct Memory Access
 - processor system memory 8
- Direct Memory Access (DMA) 7, 120
- direct route 207
- directed mode 267
- dispatch mode 267
- displaying storage usage 108
- displaying TCP/IP device resources in VTAM 152, 184
- displaying the status of devices 106
- displaying the TCP/IP configuration 104
- distance vector algorithm 223
- DNS cache 28
- Domain Name
 - Services 10, 387

- System 20, 394
- DOMAINORIGIN 22
- dual-mode stack 389, 398
 - Implementation 400
- dual-mode TCP/IP
 - application communication 399
- dual-stack backbones 388
- Dynamic Host Configuration Protocol (DHCP) 394
- dynamic route 214
- Dynamic VIPA 3
 - address 424
- DYNAMIC XCF
 - IPCONFIG definition 428
- dynamic XCF 120, 134
 - additional information 135
 - connectivity 188
 - device 134, 428
 - support 134
- DYNAMICXCF implementation 189

E

- EDNS0 39
- engineering change (EC) 405
- Enterprise Extender (EE) 397
- ENTRYPOINT 73
- ENVAR 86
- ephemeral ports 64
- ETC.IPNO Des 38
 - IPv4 addresses 38
- ETC.IPNODES data set 46
- exclusive DLCs 7
- explicit data set allocation 79
- extended common service area
 - maximum amount 108
- extended common service area (ECSA) 108
- Extension Mechanism for DNS (EDNS0) 39
- exterior gateway protocols 206
- external gateway protocol (EGP) 224
- external writer 55, 202, 258, 316
 - file 332, 348
 - procedure CTWTR 262
- EZACMD 360
- EZAZSSI 71–72
- EZAZSSI JCL procedure 74
- EZBPFINI 73
- EZZ0053I Command 196
- EZZ4203I 73
- EZZ4313I Initialization 196

F

- forked address spaces 11
- FTP
 - security 439
- FTP server
 - TCPCONFIG statement 427
- full-function mode 10, 60

G

- Generic Attribute Registration Protocol 123
- Generic Resource Encapsulation (GRE) 266
- generic server 36, 63
- getmain 72
- GID 12, 65
- Gigabit Ethernet 119
- global TCPIP.DATA
 - statement 22
- graphical mode 14
- GRE 268
- group ID 12, 65
- groups 65
- GSSAPI 4
- GVRP 123

H

- Health Checker 114
- Hierarchical File System (HFS) 12
- high level qualifier (HLQ) 79
- High Performance Data Transfer (HPDT) 7, 139
- High Performance Native Socket (HPNS) 15
- high-bandwidth and high-speed networking technologies 120
- HiperSockets 285, 419
 - microcode
 - description 120
 - multiple write facility 132
- HiperSockets (Internal Queued Direct I/O - IQDIO) 130
- HiperSockets (MPCIPA) 130
- HiperSockets Accelerator 133, 419
- HiperSockets connectivity 180
- HiperSockets DYNAMICXCF connectivity 135
- HiperSockets implementation 182
- HOME 82
- HOME address to each defined LINK 151, 183
- HOME Statement 82
- host name resolution 86
- hostname 22, 86
- hosts file
 - hosts file 86
- HOSTS.LOCAL
 - IPv4 addresses 38
- HOSTS.LOCAL 86
- HPDT
 - displaying TCP/IP device resources in VTAM 152, 184, 191
- hybrid mode 275

I

- IBM System z9 118
- ICMP redirect 161
- identity, MVS 12
- identity, UNIX 12
- IDYNAMICXCF 428
- IEASYS00 68
- IEASYSxx 72
- IEE839I St 56
- IKJTSOxx 72

- implementation tasks 42
- implicit data set allocation 79
- important and commonly used interfaces 120
- IMS 440
- INADDRANYPORT 73
- inbound packet 123
- inbound routing 122
- include files 413
- INCLUDE statement 109
- indirect route 207
- INET 16
- inetd 11
- Information Management System (IMS) 9
- Input/Output flow process 5
- installation
 - DATASETPREFIX 79
 - explicit data set allocation 79
 - high level qualifier (HLQ) 79
 - implicit data set allocation 79
 - LNKLST 71
 - LPALST 71
 - node name 71
 - PARMLIB 74
 - SCHEDxx 72
- Integrated Sockets PFS definitions 69
- Interactive Problem Control System (IPCS) 202
- INTERFACE statement 81, 213
 - PORTNAME value 401
- interior gateway protocols 206
- Internal Queued Direct I/O, (IQDIO) 8
- Internet Protocol (IP) 117, 120, 205–206, 386
 - next generation 386
- IOCP definitions 137
- IP 315
- IP address 16, 20, 64, 135, 206, 322, 385
 - configured name server 21
 - server jobname 424
- IP filter rules 391
- IP network 2, 118, 206
 - large number 206
- IP offload 121
- IP packet
 - forwarding processing 133
- IP route 210
- IP routing
 - common design scenarios 227
- IP routing algorithm 208
- IP traffic 120
- IPCONFIG 84
- IPCONFIG DYNAMICXCF 134
- IPCONFIG6 85
- IPCS command 57, 320
- ipsec 360
- IPv4 address 38, 386
- IPv4 application on a dual-mode stack 399
- IPv6 20, 391
 - address 37
 - common design scenarios 387
 - implementation 389, 395

- importance 386
- IPv6 address 38, 389, 401
 - certain styles 390
- IPv6 addressing 389
- IPv6 application on a dual-mode stack 399
- IPv6 changes to resolver processing 37
- IPv6 dynamic routing 225
- IPv6 dynamic routing using OMPROUTE 225
- IPv6 dynamic routing using router discovery 225
- IPv6 implementation in z/OS 395
- IPv6 network 225, 387
 - connectivity 14
 - IPv6 traffic 388
 - prefix 225
- IPv6 resolver statements 40
- IPv6 RIP
 - protocol 225
 - route 225
 - UDP port 239
- IPv6 support 14, 37, 69, 387, 395
- IPv6 TCP/IP Network part (prefix) 392
- IPv6 traffic 387
- IQD CHPID 131, 181
 - F7 188
 - hexadecimal value 182
- iQD chpid 131
- IQDCHPID 293
- iQDIO 130
- IQDIOROUTING 419
- IQDVLANID 290
- ISHELL 3, 10, 67
- ISOLATE 158
- IUTiQDIO device 188
- IUTSAMEH 285
- IUTSAMEH device 188
- IVTPRMxx 73

J

- job log versus syslog as diagnosis tool 114
- jobname 66, 414

K

- Kerberos 4

L

- LAN Channel Station (LCS) 7, 396
- LAN connections 118
- LCS 7
- line mode 14
- link local address type 392
- link state
 - algorithm 212
 - database 219
- link state (LS) 210
- Link State Advertisement (LSA) 219, 221
- LINK statement 85, 106, 123, 436
- link statement 149

- LINK statement using BSDROUTINGPARMS 151, 183
- link-state database 221
- link-state routing 221
- LNKLST 71
- Load Balancing Advisor (LBA) 3, 287
- Local Area Network (LAN) 118, 220
- local hosts file 86
- local settings in a multiple stack environment 34
- local settings in a single stack environment 29, 40
- locating PROFILE.TCPIP 85
- Logical File System (LFS) 8
- login name 12
- LOOPBACK 82
- loopback 78
- LPALST 71
- LPARs 61, 398
 - data traffic flow 189

M

- Management Information Base (MIB) 4
- maximum transfer unit (MTU) 419
- maximum transmission unit (MTU) 64
- MAXRECS 313
- MAXTTL 29
- message types 114
 - where to find them 114
- messages 114
- mid-level qualifier (MLQ) 80
- modifying a device 110
- modifying characteristics of a device 197
- modifying OMPROUTE 251
- MPC
 - displaying resources 152, 184, 191
- MPLS backbones 388
- MTU
 - recommendation 64
- MTU size 151, 183
- multicast address 390, 394
- MULTIPATH 419
- Multipath Channel
 - I/O process 7
- Multi-Path Channel (MPC) 4, 85, 389
- Multipath Channel+ (MPC+) 7
- multiple AF_INET transport providers 16
- multiple security zones 284
- multiple stack 16–17, 41, 60–62, 135–136, 286
 - _iptcpn() 63
 - cookbook 87
 - ephemeral ports 64
 - generic server 63
 - separate workload 62
 - setibmopt() socket call 63
- multiple TCP/IP
 - stack 16–17, 40–41, 64, 136, 213, 412
- multiple TCP/IP stack 400
- multiple write assist with IBM zIIP 132
- MVPTSSI 71
- MVS console 360
- MVS identity 12
- MVS image 14, 131

AF_INET transport providers 16

N

name and address resolution functions 38
name resolution 86
name server 37
NAT 268
NETSTAT command 66, 152, 184, 190
Network Address Translation (NAT) 266
NETWORK DOMAINNAME 69, 400, 426
network job entry (NJE) 439
network management
 programming interface 4
 tool 4
networking connectivity
 diagnose problems 254
NMI API 356
NOCACHE 28
non-canonical mode 14
NONRouter 268
NSPORTADDR 53 45, 237, 443, 451
nssctl 360

O

OAT 268
OBEYFILE and security 74
OBEYFILE command 74, 109, 197, 230, 322
OFFLOAD 420
OMEGAMON 356
OMPROUTE 235
OMPROUTE configuration file 239, 452
OMPROUTE CTRACE 258
OMPROUTE in a z/OS environment 233
OMPROUTE management 244
OMPROUTE procedure 236
OMVS segment 13, 65
 RACF user IDs 66
OMVS shell 3, 75
 interface 75
Open Shortest Path First (OSPF) 8, 210, 217
operating environment 5
operating mode 14
OPERCMD5, generic class 74
OSA 275
 Express 136
OSA Address Table (OAT) 120, 266, 279
OSA Connection Isolation 158, 228
OSA device 396
OSAENTA 333
OSA-Express 3
OSA-Express (MPCIPA) 121
OSA-Express connectivity 139, 156
OSA-Express device 85, 111, 196, 401
 DLC layer 85
OSA-EXPRESS feature 89, 120, 389
OSA-Express implementation with VLAN ID 148, 162
OSA-Express port 89, 133, 139, 254, 389
 link statements 149
OSA-Express QDIO

 connection 85, 401
 interface 120
OSA-Express router support 123
OSA-Express VLAN support 121, 132
OSA-Express3 141
OSI model 118
OSPF 213
OSPF area 219
 network topology 221
OSPF for IPv6 226
OSPF network 218
 other areas 221
 RIP routes 218
OSPF terminology 218
outbound connections 37

P

packet trace 321
PARMLIB 74
Pascal API 8
pasearch 360
Path Maximum Transfer Unit (PMTU) 419
Path Maximum Transmission Unit (PMTU)
 IPCONFIG definition 428
 RFC 1191 428
PATHMTUDISCOVERY 428
performance management 357
performance, storage 214
permission bits 13, 76
PFS 15
Physical File
 System 2, 94, 325
 System transport provider 16
Physical File System (PFS) 15, 61
Physical File System transport provider 16
physical network types 222
ping 360
PING and TRACERTE
 verifying interfaces 109
PING command 197, 409
ping command (TSO or z/OS Unix) 303
PKTTRACE 321
PMTU 305
Policy-Based Routing 162
PORT 73
port management 63
port sharing 424
port sharing (TCP only) 424
PORT statement 80
primary differences between IPv6 OSPF and IPv4
OSPFv2 226
PRIRouter 266, 268
problem determination 197, 233, 253
problems with the home directory 75
process 11
procname 252
PROFILE.TCPI P 61, 321, 411
 configuration file 80
 data 245, 333
 definition 82

- different sections 413
- file 109
- OBEY statement 109
- parameter 418
- PKTTRACE statement 321
- statement 322
- PROFILE.TCPIP 80
 - verifying 109
- PROFILE.TCPIP parameters 80
- PROGnn 71
- program properties table (PPT) 72
- protocols and devices 6

Q

- QDIO data connection isolation 158
- QDIO mode 3, 389, 420
 - OSA-Express port 140
- Quality of Service (QOS) 386
- Queued Direct I/O (QDIO) 4, 120, 419

R

- RACF 65
- RACF authorize user IDs for starting OMPROUTE 238
- RACF definition 65
- RACF environment 65
- RACF facility classes 65
- RACF implementation 65
- RACF profiles 65
- RACF resources 65
- RACF security environment 66
- RACF with z/OS Communications Server TCP/IP 67
- raw mode 14
- reconfiguring the system with z/OS commands 109
- replaceable static routes 214
- resolv.conf 86
- RESOLVE_VIA_LOOKUP compile symbol 87
- resolver address space 19
 - global definitions 34
- resolver configuration data sets 86
- Resolver CTRACE
 - analysis 57
 - initialization PARMLIB member 55
- resolver DNS cache 28
- resolver problems
 - diagnosing 51
- Resource Access Control Facility (RACF) 3, 60
- resource profiles 13
- restartable platform 71
- Reusable Address Space ID (REUSASID) 6
- REUSASID coding 6
- REXX sockets 9
- RFC 4941 391
- RIP 213
- RIP V1 211, 223
 - implementation 225
 - limitations 223
 - packet 225
 - system 225
- RIP V2 211, 224

- extension 225
- message 224
- packet 224
- protocol extension 224
- system 225
- RIPng or RIP next generation 225
- root file system 12
- ROUTE 10.10.3.0 229
- Route Trip Response Time (RTT) 427
- ROUTEALL 122
- router 206
- router configuration statements 243, 453
- routing daemons 213
- Routing Information Protocol (RIP) 210, 222
- routing protocol 8, 206, 211, 436
 - main characteristics 211
- routing table 134, 206, 208
 - network routes 224
 - static definition 210
 - static routes 228
- routing, next-hop address 228
- ROUTLCL 122

S

- S806, abend code 73
- SACONFIG (SNMP subagent) 428
- same LPAR 135
- same VLAN 89
- SAMEHOST 7
- SCHEDxx 72
- search order 85
- SECRouter 266, 268
- Security 391
- security 65, 391
 - CICS 440
 - client 74
 - FTP 439
 - IMS 440
 - server 74
 - SMTP 439
 - SNMP 439
- security associations 391
- segment
 - OMVS 65
- server
 - generic 36
- Server Application State Protocol
 - outboard load balancers 3
- Server Application State Protocol (SASP) 3
- set of messages show (SMS) 77
- setting up the resolver procedure 42
- shared DLCs 7
- sharing resolver between multiple stacks 63
- shell 10
- shell access 13
- shell interface 11
- shell, ISHELL 10
- shell, OMVS 10
- shortest path first (SPF) 226
- Simple Mail Transfer Protocol (SMTP) 439

- single AF_INET transport provider 16
- single stack 61
- site local address type 394
- SMF 357
- SMFCONFIG 428
- SMTP, security 439
- SNMP subagent 428
- SNMP, security 439
- SNMPv1 439
- SNMPv2C 439
- SNMPv2u 439
- socket address 14
- socket addressing families 14
- socket addressing families in UNIX System Services 14
- SOURCEVIPA 418, 428
- spawned address spaces 11
- SQA 72
- stack affinity 63
- stacks
 - connected 37
- START statement 134, 196
- start syslogd 238
- started task user IDs 66
- starting a device 196
- starting z/OS Communications Server IP 93
- starting z/OS TCP/IP after IPL 103
- stateless address autoconfiguration 390
- static and dynamic routing 209
- static route 206, 214
- static routing 227
- static routing scenario 228
- static routing table
 - Managing 230
- step-by-step checklist xi
- stopping a device 196
- storage shortage 214–215
- structure
 - names 292
- stub area
 - default routes 215
- subchannel device 131
 - maximum number 131
- subnet mask 151, 183
- subnet-router anycast address 394
- subnetwork 135, 209
- subplexing 283, 287
- superuser 13, 75
- superuser mode 75
- supported routing applications 8
- switch port 254
 - network traffic 254
- switch port configuration
 - verifying 148
- symbolic links 77
- SYMDEF 413
- SYS1.PARM LIB member
 - CTIEZB00 318
 - CTIEZB01 319
 - CTIORA00 258, 327
 - CTIRES00 55, 329

- SYS1.PARMLIB 68, 74
- SYSCLONE system variable
 - definition 413
- SYSDEF 413
- Sysplex Distributor 3, 134, 266
- SYSTCPD DD name 86
- SYSTEM SYMBOL 402
- system symbols
 - definition 413
- System z File System 12
- System z9 7, 117, 389
 - compute-intensive functions 7
 - memory speed 120
 - server-to-server traffic 131
 - system memory 180
- SYSEX.PARMLIB updates
 - 71

T

- TCP structures 292
- TCP/IP 1, 19, 59, 206, 315, 392, 395, 411
- TCP/IP address space 2, 66, 322, 420
- TCP/IP application
 - server 3
- TCP/IP Base Functions
 - HOSTS.LOCAL 86
 - TCPIP.DATA 86
- TCP/IP client functions 74
- TCP/IP commands 360
- TCP/IP component 79, 316
- TCP/IP configuration
 - data 61
 - data set 79
 - file 412
 - parameter 80
 - statement 414
 - verifying 104
- TCP/IP configuration data set names 86
- TCP/IP data set names 79
- TCP/IP definition 401
- TCP/IP network 8, 225, 439
 - RIP router 225
 - workstation connectivity 3
- TCP/IP profile 78, 134, 188, 401, 414
 - DYNAMICXCF definition 188
 - required connectivity definitions 195
 - START statement 196
- TCP/IP server functions 74
- TCP/IP socket
 - APIs 8
 - layer 319
- TCP/IP socket APIs 9
- TCPCONFIG 80
- TCPIP 196, 252, 322, 396, 427
- TCPIP.DATA 86
- TCPIP.DATA file 63
- TCPIP.DATA statement values in z/OS
 - verifying 108
- TCPIP.DATA statement values in z/OS UNIX
 - verifying 108

- TCPIPE.TCPP Arm 85, 414
- TCPIPJOBNAME 86
- TCPIPSTATISTICS 420
- TEMPPREFIX 392
- test environment 42, 237
- thread 11
- time-to-live (TTL) 255
- timezone 237
- TNF 71
- TRACE Ct 55, 261, 317, 328
- trace option 257, 315
- Trace Resolver 51, 53
- TRACEROUTE command
 - 256, 306
- TRANSACTION TRACE (TT) 56
- Transport Layer Security (TLS) 4
- transport providers 16
- Transport Resource List (TRL) 139
- Transport Resource List Element (TRLE) 152, 184, 191
- TRL 152, 184, 191
- TRLE 139, 152, 184, 191
 - displaying 192
- TRLE definition 85, 140, 401
 - PORTNAME value 140
- TRLE in VTAM to represent each OSA-Express port 149
- trmdstat 360
- trunk mode 275
- trunk port 274
- TSO clients 63
- TSO command 231
- TSO logon procedures
 - PROCLIB 74
- TT CMD 56
- tunneling 387
- Type of Service (TOS) 397
- Types of IP routing 207
- TZ= 237

U

- UDP datagram 37–38
- UDPCONFIG 80
- UDPCONFIG statement 424
- UID 12, 65
- UNIX client functions 75
- UNIX Hierarchical File System 12
- UNIX identity 12
- UNIX permission bits 76
- UNIX shell 3
- UNIX System Services 2, 10, 60
 - Common errors 75
 - full-function mode 11
 - minimum mode 11
 - z/OS UNIX file system interaction 11
- UNIX System Services communication 14
- UNIX System Services concepts 11
- UNIX System Services Verification 94
- update the resolver configuration file 237, 451
- user ID 12, 65, 238, 412
 - BPXROOT 102
 - OMVSKERN 66

- RACF definitions 65
- user ID defined (UID) 12, 61
- user name 12

V

- V TCPIP,PURGECACHE command 396
- VARY 323
- VARY TCPIP command 109
- verification 48, 230, 405
- verification checklist 77
- VIPA route 213
- Virtual IP
 - Address 3
- Virtual IP Address 428
- Virtual IP Addressing
 - IPCONFIG definition 428
- virtual LAN (VLAN) 274
- virtual local area network (VLAN) 117, 430, 436
- virtual MAC (VMAC) 266
- Virtual Medium Access Control (VMAC) 265
- VLAN and primary/secondary router support 124
- VLAN ID 106, 254
 - port assignment 148
 - tag 275
 - value 123
- VLAN number 149
- VLAN support 123
- VLAN support of Generic Attribute Registration Protocol - GVRP 123
- VMAC 266
- VMCF 71
- VSWITCH port isolation 158
- VTAM 152, 184, 191
- VTAM definition 401
- VTAM Resource 85

W

- Web server 11
- Workload Manager 424
- Workstation Operating Mode 14

X

- XCF links 285
- XCFGRPID 290

Z

- z/OS Communications Server
 - applications 10
 - component product 1
 - configure OMROUTE 235
 - dual-mode stack dependencies 389
 - IP routing-related terms 206
 - IPv6 router discovery 225
 - tightly coupled design 62
- z/OS Communications Server IP 79, 114
 - importance 3
- z/OS Communications Server TCP/IP
 - RACF 67

- z/OS customization for z/OS UNIX 74
- z/OS environment 1, 16, 60, 136, 205, 389, 436
 - connectivity scenario 136
 - dynamic routing 214–215
 - feature information 73
 - important data set 68
 - UNIX concepts 65
 - Using OMPROUTE 233
- z/OS image 8, 47, 188
- z/OS shell 10, 236
 - issue 251
 - superuser ID 251
- z/OS system 3, 65, 118
 - link list 71
 - space 13, 66
 - TCP/IP application availability 3
- z/OS TCP/IP 2, 60, 67, 123, 234, 395, 401
 - environment 79
 - RACF 60
- z/OS UNIX 3, 53, 60
 - address space 14
 - administrator 5
 - APIs 8
 - assembler callable service 9
 - C socket 9
 - C socket APIs 9
 - design components 69
 - element 9
 - environment 5, 60, 197
 - file system 3, 10, 60
 - file system data 12
 - file system data set 12
 - file system file 12, 86, 197, 258
 - file system file system dependancy 4
 - file system home directory 67
 - file system interaction 11
 - function 9
 - group 61
 - identity 12
 - initialization member BPXPRM7A 413
 - initialization time 68
 - interface 9
 - logon service 10
 - onetstat 109
 - operating system 12
 - resource 13
 - service 11, 60
 - shell 13
 - shell traceroute/otracert command 233, 250, 255
 - sockets programming 9
 - system 3
 - Systems Services environment 97
 - user identification 12
 - version 9
- z/OS UNIX APIs 9
- z/OS UNIX file system definitions in BPXPRMxx 12
- z/OS UNIX user identification 12
- z/OS V1R7.0 Communications Server 202
- z/OS VARY TCPIP commands 66
- zIIP 132



Redbooks

IBM z/OS V1R12 Communications Server TCP/IP Implementation: Volume 1 Base Functions, Connectivity, and Routing

(1.0" spine)

0.875" <-> 1.498"

460 <-> 788 pages



IBM z/OS V1R12 Communications Server TCP/IP Implementation: Volume 1 Base Functions, Connectivity, and Routing



Redbooks®

Discusses important z/OS Communications Server TCP/IP base function capabilities

Describes z/OS Communications Server base function implementation

Provides useful verification techniques

For more than 40 years, IBM® mainframes have supported an extraordinary portion of the world's computing work, providing centralized corporate databases and mission-critical enterprise-wide applications. The IBM System z®, the latest generation of the IBM distinguished family of mainframe systems, has come a long way from its IBM System/360 heritage. Likewise, its IBM z/OS® operating system is far superior to its predecessors in providing, among many other capabilities, world class and state-of-the-art support for the TCP/IP Internet protocol suite.

TCP/IP is a large and evolving collection of communication protocols managed by the Internet Engineering Task Force (IETF), an open, volunteer organization. Because of its openness, the TCP/IP protocol suite has become the foundation for the set of technologies that form the basis of the Internet. The convergence of IBM mainframe capabilities with Internet technology, connectivity, and standards (particularly TCP/IP) is dramatically changing the face of information technology and driving requirements for even more secure, scalable, and highly available mainframe TCP/IP implementations.

The z/OS Communications Server TCP/IP Implementation series provides understandable, step-by-step guidance about how to enable the most commonly used and important functions of z/OS Communications Server TCP/IP.

In this IBM Redbooks® publication, we provide an introduction to z/OS Communications Server TCP/IP. We then discuss the system resolver, showing the implementation of global and local settings for single and multi-stack environments. We present implementation scenarios for TCP/IP Base functions, Connectivity, Routing, Virtual MAC support, and sysplex subplexing.

INTERNATIONAL TECHNICAL SUPPORT ORGANIZATION

BUILDING TECHNICAL INFORMATION BASED ON PRACTICAL EXPERIENCE

IBM Redbooks are developed by the IBM International Technical Support Organization. Experts from IBM, Customers and Partners from around the world create timely technical information based on realistic scenarios. Specific recommendations are provided to help you implement IT solutions more effectively in your environment.

For more information:
ibm.com/redbooks

SG24-7896-00

ISBN 073843549X