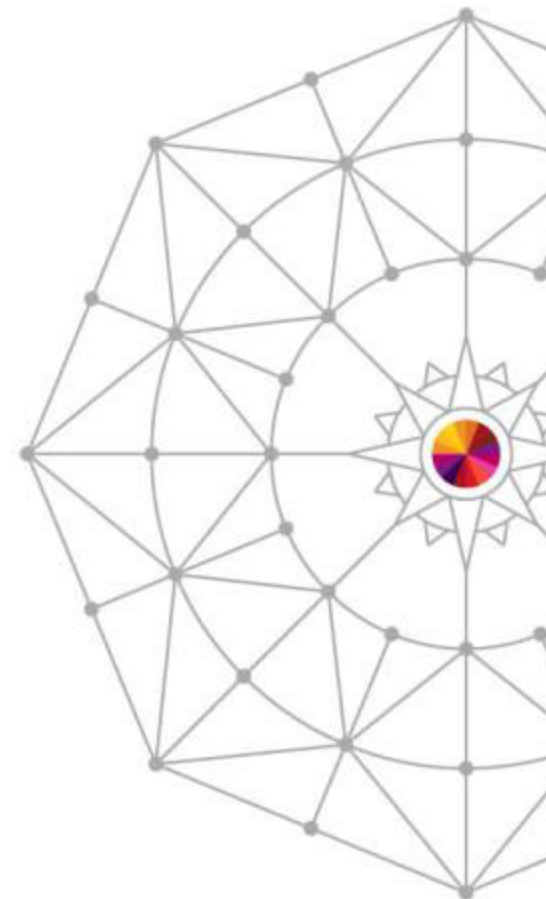


14950: z/OS Communications Server Usage of HiperSockets

Linda Harrison
lharriso@us.ibm.com
IBM ATS

March 13, 2014 8am
Session 14950



Abstract and Trademarks

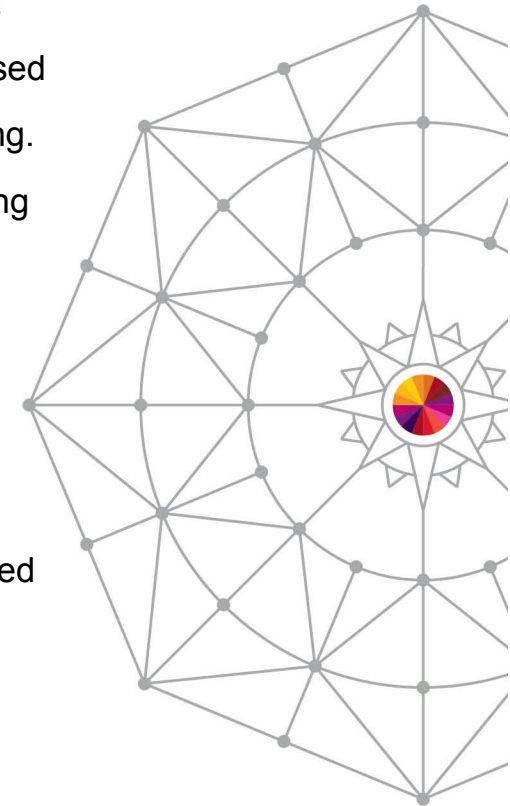


Abstract

- zEnterprise HiperSockets are virtual LANs provided by the zEnterprise platform without any additional fee. Because they are internal and virtual LANs on the zEnterprise system there is no exposed cable or wire and therefore provide a secure connection between LPARs in the same zEnterprise machine. This session will detail what zEnterprise HiperSockets are and how they can be used between LPARs. HiperSockets implementation on z/OS will be explained. Some Linux on System z implementations are resistant to use dynamic routing. Without dynamic routing HiperSockets and OSA connections between z/OS and Linux on System z will need static routing definitions. HiperSockets routing options between z/OS and Linux on System z will be discussed.

Trademarks

- The following are Registered Trademarks of the International Business Machines Corporation in the United States and/or other countries.
 - IBM
 - z/OS
- The following are trademarks or registered trademarks of other companies.
 - Microsoft is a registered trademark of Microsoft Corporation in the United States and other countries.
- All other products may be trademarks or registered trademarks of their respective companies.
- Refer to www.ibm.com/legal/us for further legal information.



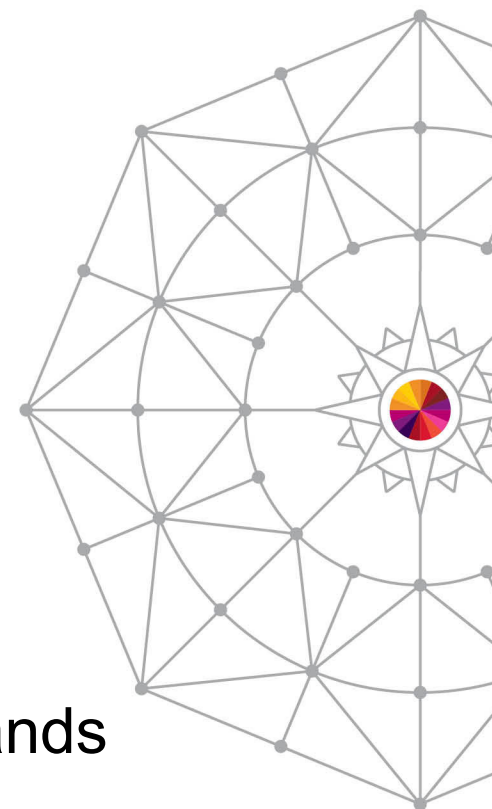
Note: z/VM has virtual HiperSockets support for guests, which is not discussed in this presentation.



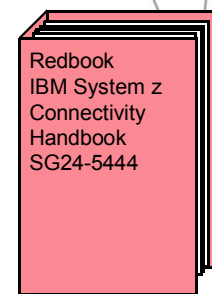
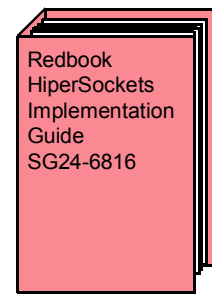
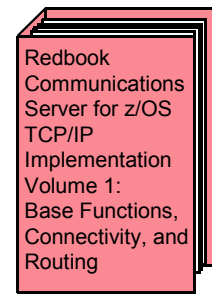
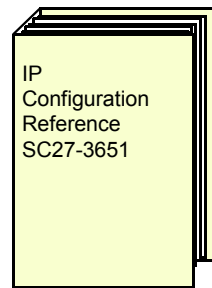
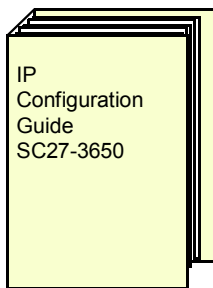
Agenda



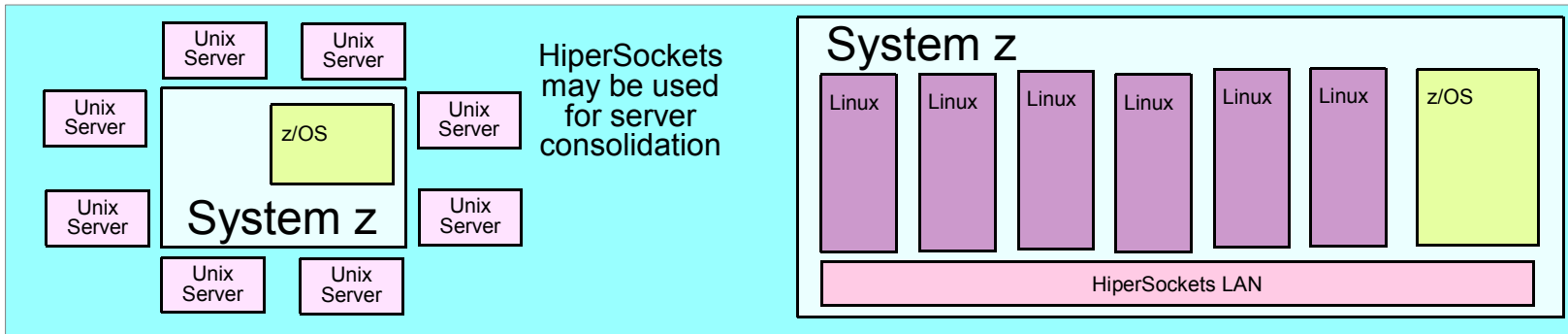
- HiperSockets Overview
 - Overview
 - Coding
 - Capabilities
- Traffic Over HiperSockets or OSA
- OSA versus HiperSockets versus RoCE
- More Information
- Appendicies
 - More Configuration Details and Commands
 - TCP/IP Routing
 - DynamicXCF Details



HiperSockets Overview



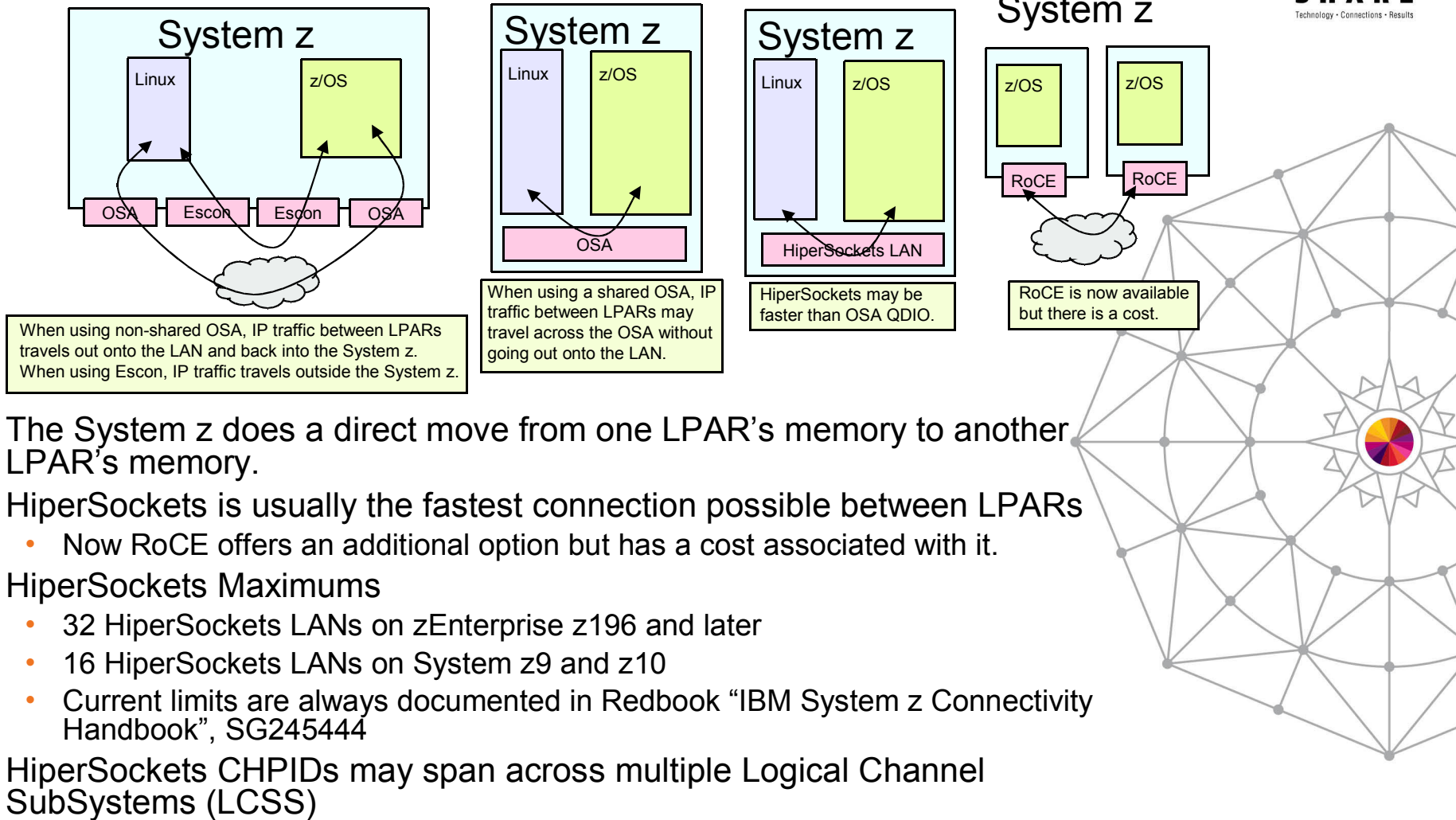
HiperSockets (iQDIO)



- HiperSockets = Internal QDIO (iQDIO (Queued Direct I/O))
 - Was developed from the QDIO architecture
 - Is limited to TCP/IP protocol only to/from z/OS (z/VM and Linux support Layer 2)
 - Also known as HiperSockets device or System z internal virtual LAN or HiperSockets LAN
 - LPAR to LPAR communication via shared memory
 - High speed, low latency, similar to cross-address-space memory move using memory bus
 - Provides better performance than channel protocols for network access.
 - Multiple HiperSockets may be configured as internal LANs on the System z box.
 - A HiperSockets LAN may be configured to be part of TCP/IP DynamicXCF.
 - A TCP/IP stack may only define a single HiperSockets LAN for DynamicXCF.
 - *Some TCP/IP stacks may use one HiperSockets LAN for DynamicXCF connectivity while other TCP/IP stacks use a different HiperSockets LAN for DynamicXCF connectivity.*
 - Not recommended for some LPARs to define a HiperSockets LAN for DynamicXCF and other LPARs to manually define the same HiperSockets LAN.
 - Common Lookup Table is stored in the Hardware System Area, the same as QDIO.
 - HiperSockets LAN, IP address, TCP/IP stack

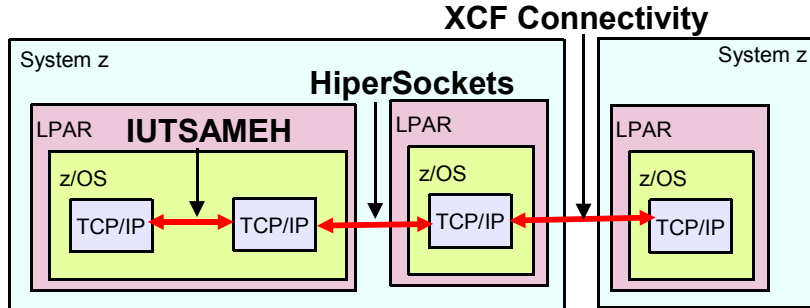


TCP/IP LPAR to LPAR Communication Path

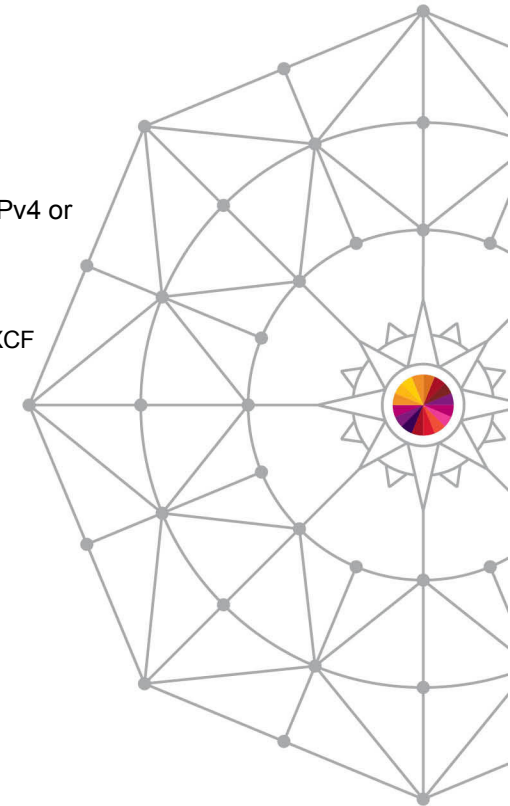


- The System z does a direct move from one LPAR's memory to another LPAR's memory.
- HiperSockets is usually the fastest connection possible between LPARs
 - Now RoCE offers an additional option but has a cost associated with it.
- HiperSockets Maximums
 - 32 HiperSockets LANs on zEnterprise z196 and later
 - 16 HiperSockets LANs on System z9 and z10
 - Current limits are always documented in Redbook "IBM System z Connectivity Handbook", SG245444
- HiperSockets CHPIDs may span across multiple Logical Channel SubSystems (LCSS)

TCP/IP DynamicXCF Transport Choices



- TCP/IP DynamicXCF is capable of dynamically creating multiple device, link, and interfaces all with the same IPv4 or IPv6 address:
 - Same host: device IUSAMEH, link EZASAMEMVS (IPv4), interface EZ6SAMEMVS (IPv6)
 - Hipersockets: device IUTIQDIO, link IQDLNKnnnnnnnn (IPv4), interface IQDIOINTF6 (IPv6)
 - where nnnnnnnn is the hexadecimal representation of the IP address specified on the IPCONFIG DYNAMICXCF statement
 - XCF connectivity: device CPName, link EZAXCFnn (IPv4), interface EZ6XCFnn (IPv6)
 - where nn is the 2-character &SYSC clone value
- TCP/IP DynamicXCF automatically chooses the fastest path:
 - Same host is used between TCP/IP stacks inside the same LPAR
 - HiperSockets is used between TCP/IP stacks inside same System z CEC (when configured)
 - XCF connectivity is used between TCP/IP stacks outside the CEC
- VTAM start-options required for TCP/IP DynamicXCF use of HiperSockets:
 - IQDCHPID=nn (where nn is the HiperSockets LAN CHPID, ie. FA)
 - XCFINIT=YES (required for TCP/IP DynamicXCF with or without HiperSockets)
 - Generates XCF device ISTLSXCF in VTAM
 - Requires prerequisite Start Options like HPR=RTP, which requires minimal APPN enablement.
- TCP/IP Profile options required to define DynamicXCF:
 - IPCONFIG DYNAMICXCF 10.1.2.101 ...
 - IPCONFIG6 DYNAMICXCF 2001:0DB8:1:0:50C9:C2D4:0:1 ...
- When HiperSockets DynamicXCF is configured:
 - HiperSockets DEVICE/LINK/INTERFACE/HOME and VTAM TRLE are all dynamically built



Defining Manual HiperSockets LANs



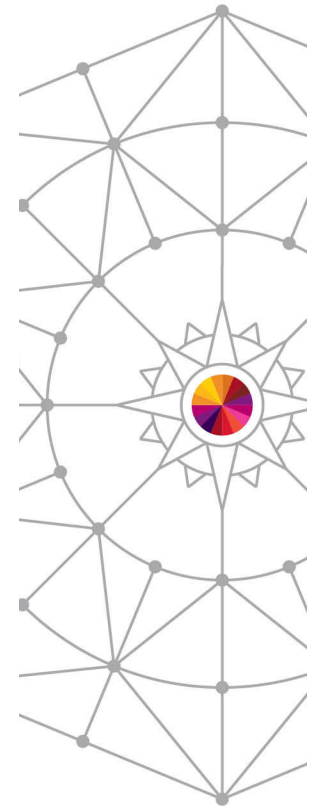
- TCP/IP Profile options required to define manual (user defined) (non-DynamicXCF) HiperSockets LAN using CHPID FB:
 - Example DEVICE, LINK, HOME, START:
 - `DEVICE IUTIQDFB MPCIPA`
 - *Device name must be IUTIQDxx where xx = CHPID*
 - *LINK HIPERLFB IPAQIDIO IUTIQDFB*
 - *HOME 10.1.1.22 HIPERLFB*
 - *START IUTIQDFB*
 - Example INTERFACE, START:
 - `INTERFACE HIPER6FB DEFINE IPAQIDIO6 CHPID FB`
 - `IPADDR 2001:0DB8:1:0:50C9:C2D4:220:21`
 - `START HIPER6FB`
- Routing customization required for static or dynamic routing



Defining HiperSockets (iQDIO)



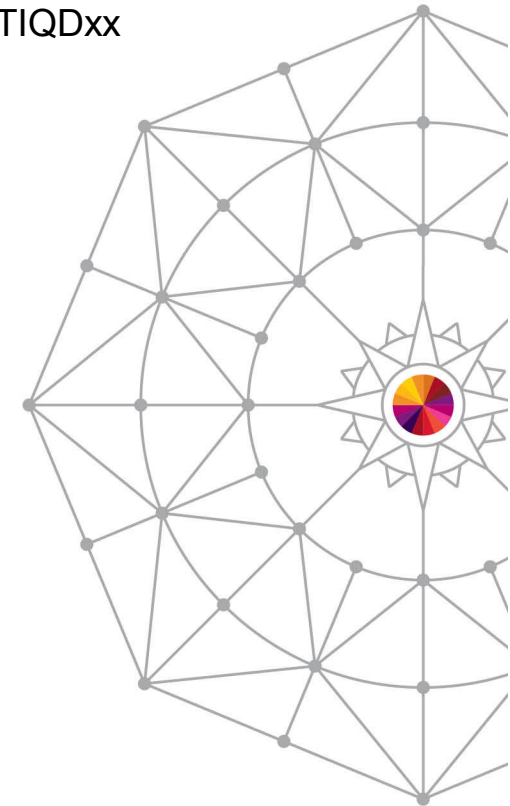
- Define HiperSockets LAN in HCD/IOCP
 - CHPID TYPE=IQD
 - CHPARM determines MTU
- VTAM Start Options
 - IQDIOSTG – optionally modify the amount of Read Storage
 - IQDCHIP – Do not take the default!
 - IQDCHIP=NONE – if you do not want a HiperSockets LAN used for DynamicXCF traffic.
 - IQDCHIP=chpid – if you do want a HiperSockets LAN used for DynamicXCF traffic.
- PROFILE.TCPIP
 - INTERFACE IPAQIDIO (or IPAQIDIO6) and START
 - or DEVICE MPCIPA, LINK IPAQIDIO, HOME, and START
 - BEGINROUTES ROUTE – define HiperSockets LAN to static routing
 - GLOBALCONFIG
 - AUTOIQDX – enables HiperSockets Integration with IEDN
 - IQDMULTIWRITE – enables HiperSockets Multi-Write
 - IQDVLANID – to define a DynamicXCF HiperSockets VLAN ID
 - ZIIP IQDIOMULTIWRITE – offloads HiperSockets processing to zIIP engine(s) (requires IQDMULTIWRITE)
 - IPCONFIG
 - QDIOACCELERATOR – enables QDIO/iQDIO Accelerator
 - IQDIOROUTING – outdated parameter for original HiperSockets Accelerator support
 - Use QDIOACCELERATOR instead
- OSPF configuration file
 - Define HiperSockets LAN to dynamic routing



Dynamically Created



- TRLEs in ISTTRL Major Node (where xx = CHPID)
 - Manual HiperSockets
 - Device/Link MPCIPA/IPAQIDIO or Interface IPAQIDIO6 TRLE = IUTIQDxx
 - Interface IPAQIDIO TRLE = **IUTIQ4xx**
 - IPv4 DynamicXCF HiperSockets TRLE = **IUTIQDIO**
 - IPv6 DynamicXCF HiperSockets TRLE = **IQDIOINTF6**
 - IPv4 HiperSockets Integration with IEDN (IQDX) TRLE = **IUTIQXxx**
 - IPv6 HiperSockets Integration with IEDN (IQDX) TRLE = **IUTIQ6xx**
- Device, Link, and Home (where nnnnnnnn is hexadecimal representation of the IP address)
 - IPv4 DynamicXCF
 - DEVICE = **IUTIQDIO**
 - LINK = **IQDIOLNKnnnnnnnn**
 - HOME for IQDIOLNKnnnnnnnn
- Interface (where xx = CHPID)
 - IPv6 DynamicXCF INTERFACE = **IQDIOINTF6**
 - HiperSockets Integration with IEDN (IQDX)
 - IPv4 Interface IPAQIQDX = **EZAIQXxx**
 - IPv6 Interface IPAQIQDX6 = **EZ6IQXxx**

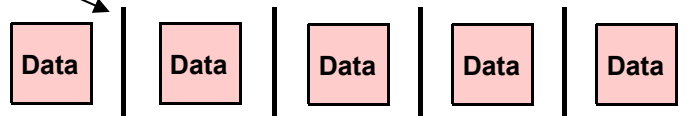


Multiple Write Facility and zIIP Offload

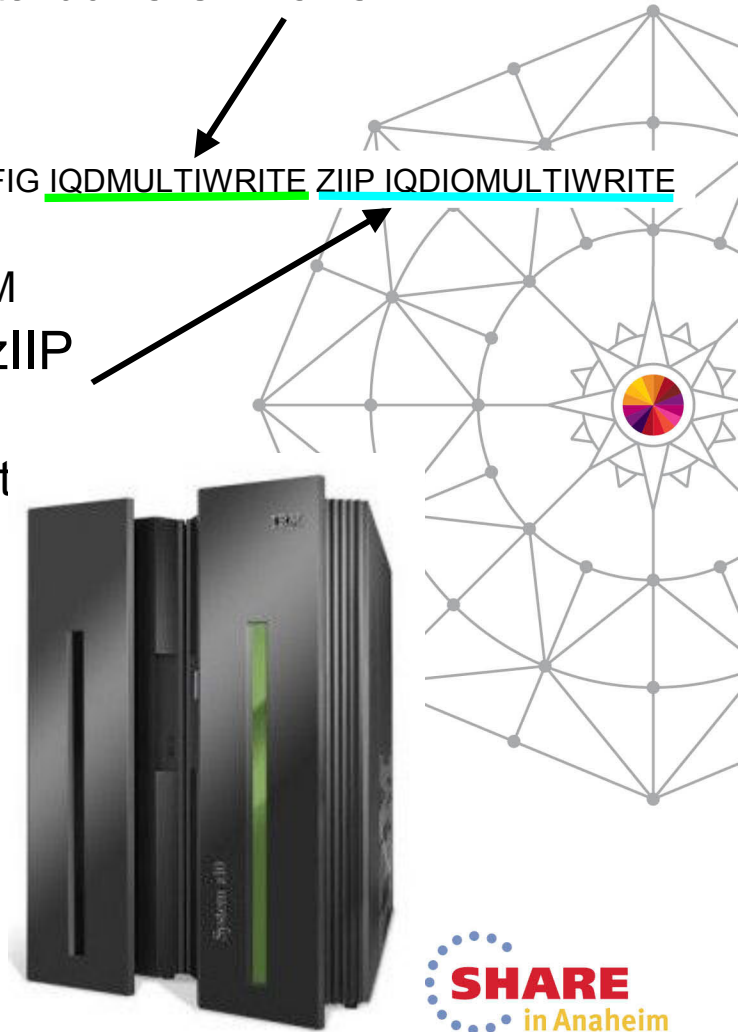
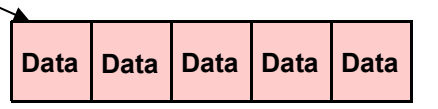
- HiperSockets can move multiple output data buffers in one write operation
 - Reduces CPU utilization
 - z/OS V1.10+
 - Requires System z10 or later
 - Not supported when z/OS is a guest under z/VM
- Multiwrite operation can be offloaded to a zIIP
 - Requires System z10 and z/OS V1.10+
 - Only for TCP traffic that originates in this host
 - Only large TCP outbound messages
 - (32KB and larger)

GLOBALCONFIG IQDMULTIWRITE ZIIP IQDIOMULTIWRITE

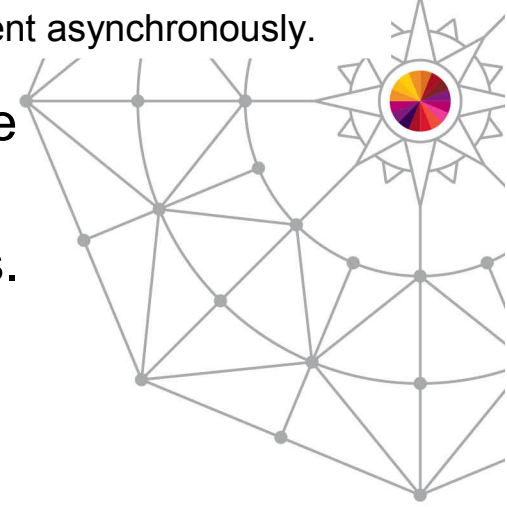
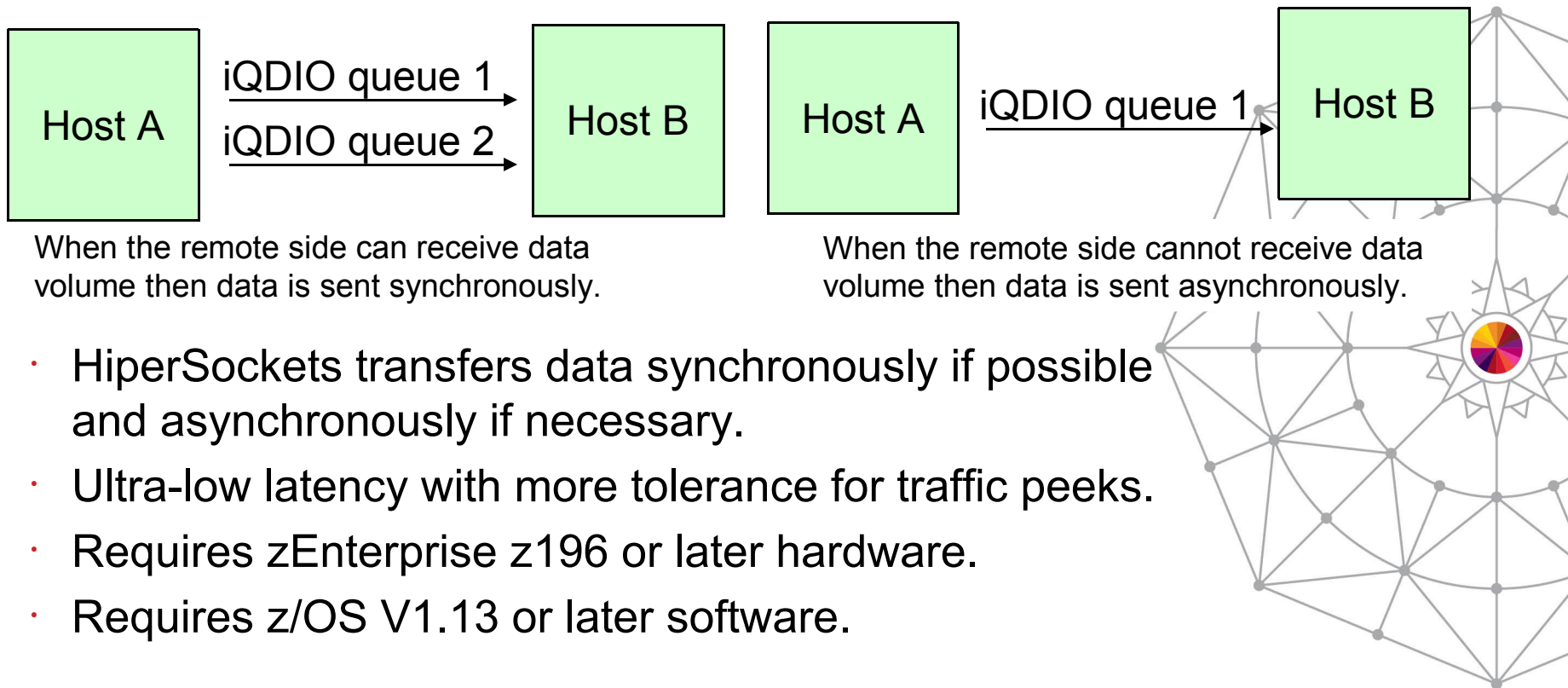
Write operation (System z9)



Write operation (System z10)



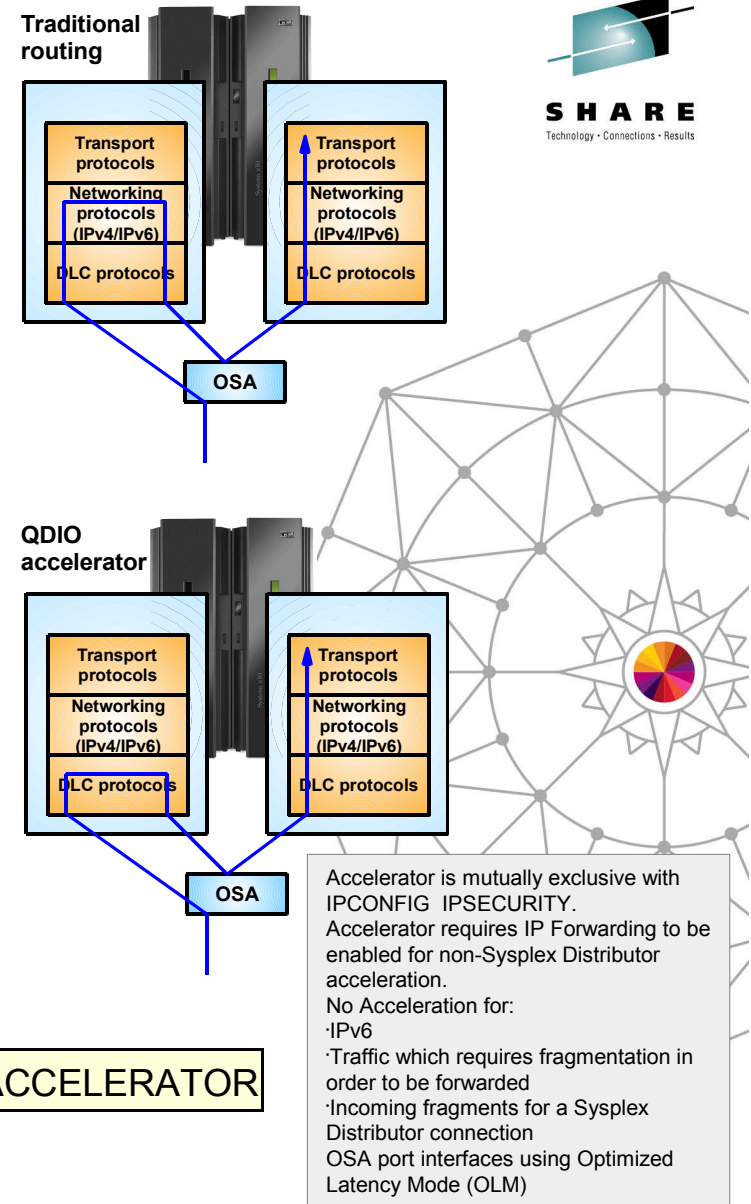
HiperSockets Completion Queue



QDIO Accelerator



- Accelerator support includes all combinations of QDIO and iQDIO traffic
 - When traffic is routed through z/OS.
 - Inbound over OSA or HiperSockets and Outbound over OSA or HiperSockets
- The first packet will travel up thru QDIO to the Accelerator stack and down thru iQDIO device drivers to reach the backend LPAR IP address. After that first packet, all the rest of the packets flow via the accelerated path through the DLC layer, thus bypassing the IP layer in z/OS and reducing path length and improving performance.



	Outbound QDIO	Outbound iQDIO
Inbound QDIO	Yes	Yes
Inbound iQDIO	Yes	Yes

IPCONFIG QDIOACCELERATOR

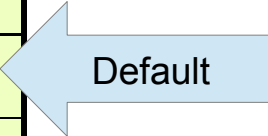
- Supports Sysplex Distributor (SD)
 - When traffic to target stack is sent over HiperSockets Dynamic XCF or QDIO as a result of VIPAROUTE definition.



MTU is Configured in HCD/IOCP



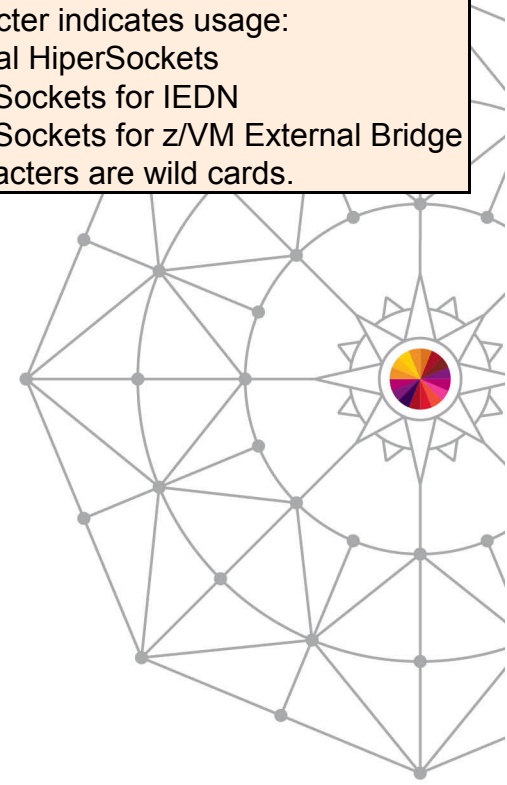
CHPID Parameter	MFS	MTU
CHPARAM=00	16K	8K
CHPARAM=40	24K	16K
CHPARAM=80	40K	32K
CHPARAM=C0	64K	56K



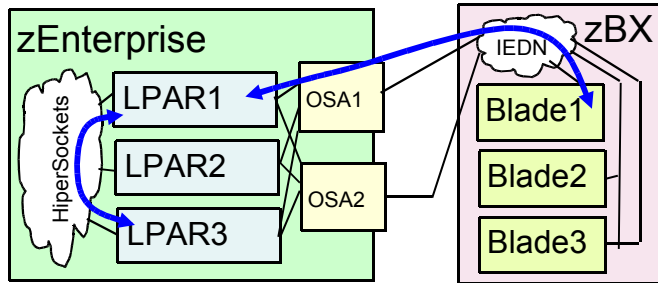
"CHPARAM" parameter was originally "OS" parameter.

On z196 and later processor the CHPID is also used to identify usage. First character still indicates frame size: 0x, 4x, 8x, and Cx (as documented on the left) The second character indicates usage: x0 indicates Normal HiperSockets x2 indicates HiperSockets for IEDN x4 indicates HiperSockets for z/VM External Bridge Where all "x" characters are wild cards.

- Each CHPID has configurable frame size (16K, 24K, 40K, 64K)
 - Allows optimization per HiperSockets LAN for small packets versus large streams
 - Affects MTU size of 8K, 16K, 32K, 56K
- HiperSockets LANs
 - Each HiperSockets LAN has its own CHPID (type IQD)
 - IBM recommends starting from x"FF" and working your way backwards through the CHPID numbers, picking addresses from the high range to avoid addressing conflicts
 - May be shared by all defined LPARs
 - Delivered as object code only (OCO)
 - *HiperSockets CHPIDs do not reside physically in the hardware but these CHPIDs cannot be used by other devices.*
 - *No physical media constraint, so no priority queuing or cabling required .*
 - Each Operating System image configures its own usage of available HiperSockets *CHPIDs*.



HiperSockets Integration with OSA for IEDN



CHPARM=x2 See previous discussion in this presentation.

GlobalConfig AUTOIQDX ALLTRAFFIC
 or GlobalConfig AUTOIQDX NOLARGEDATA
 or GlobalConfig NOAUTOIQDX

Dynamically created TRLE is IUTIQXxx or IUTIQ6xx and dynamically created interface is EZAIQXxx or EZ6IQXxx where xx is the OSX CHPID.

- A single HiperSockets LAN may be defined such that it is automatically used when the destination is an LPAR on the same CEC belonging to the same OSX/HiperSockets LAN.

- The OSA OSX devices are assigned IP Addresses.
- The HiperSockets LAN (IQDX) is not assigned an IP Address.
- Requires zEnterprise z196 or later processor
- Requires z/OS V1.13

- **Background**

- **With VIPA and Dynamic Routing**

- The application may bind to the VIPA.
- Dynamic routing causes traffic between LPARs to be routed over HiperSockets.
- Dynamic routing causes traffic to remote partners (outside the CEC) to be routed over OSA.

- **Without VIPA and Dynamic Routing**

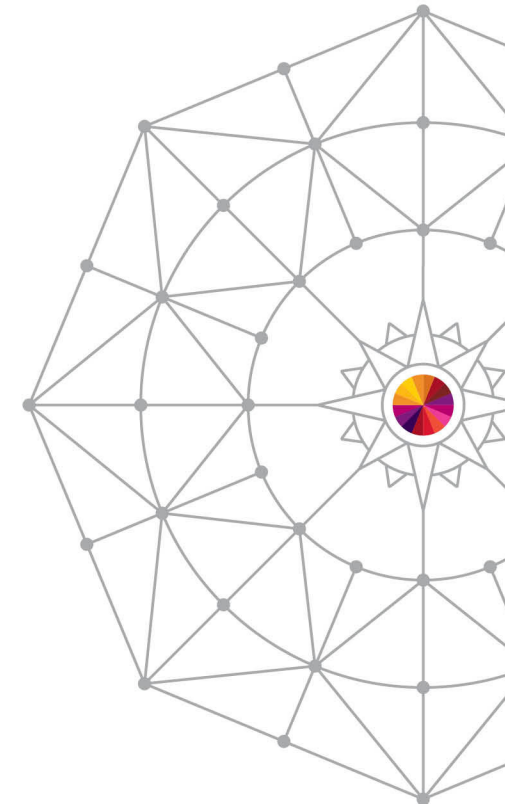
- It is a challenge to cause same CEC traffic to flow over HiperSockets at the same time that remote traffic flows over OSA.
- Static Host routes may be used.
 - *The application binds to the OSA IP address.*
 - *A static route is used on each LPAR such that when the other LPAR OSA address is the destination then the HiperSockets LAN is used to route the traffic.*
 - *With a large number of LPARs the administration of these static host routes is onerous.*

HCD (IOCP)

Define 10 subchannel addresses for each IQDX CHPID that is in use for IPv4.
 Define 10 subchannel addresses for each IQDX CHPID that is in use for IPv6.
 Multiple VLAN does not affect the required number of subchannel addresses.



HiperSockets Supported Features	z/OS	z/VM	Linux on System z	z/VSE
IPv4 Support	Yes	Yes	Yes	Yes
IPv6 Support	Yes	Yes	Yes	Yes
VLAN Support	Yes	Yes	Yes	Yes
Network Concentrator	No	No	Yes	No
Layer 2 Support	No	Yes	Yes	No
Multiple Write Facility	Yes	No	No	No
zLIP Assisted Multiple Write Facility	Yes	No	No	No
HiperSockets NTA (Network Traffic Analyzer)	No	No	Yes	No
Integration with IEDN (IQDX)	Yes	No*	Yes	No
Virtual Switch Bridge Support	No	Yes	No	No
Fast Path to Linux (LFP) Support / IUCV over HiperSockets	No	No	Yes	Yes
Completion Queue	Yes	No	Yes	Yes
* Depends upon the z/VM release				



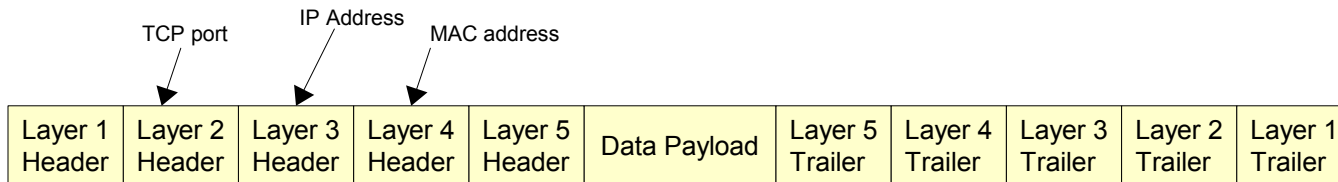
Non-z/OS Support



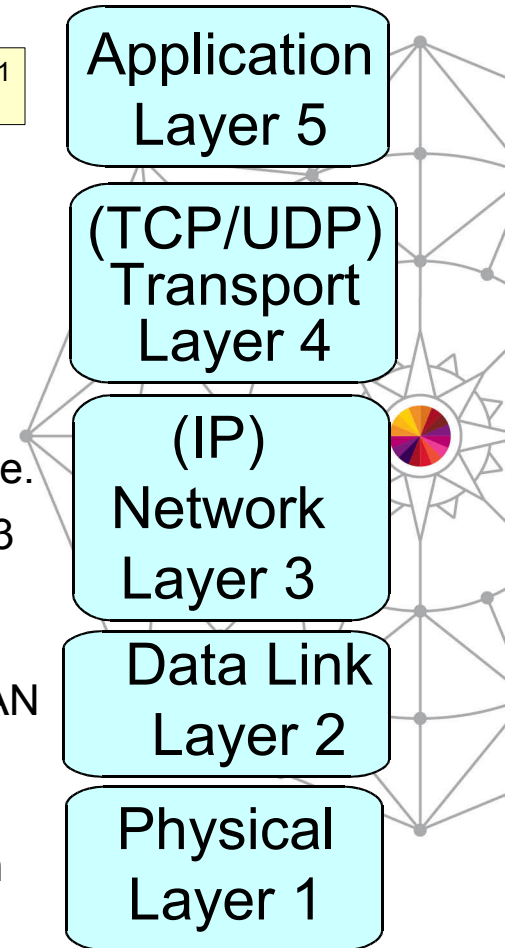
- z/VM Virtual HiperSockets Support
 - In addition to CEC HiperSockets that are available to all LPARs, z/VM is able to support virtual HiperSockets available to all guests running on that z/VM image.
- z/VM HiperSockets Virtual Switch Bridge Support (also referred to as External Bridge)
 - A single HiperSockets LAN may be defined such that it is automatically used when the destination is an LPAR on the same CEC belonging to the same OSD/HiperSockets or OSX/HiperSockets Virtual Switch.
 - The OSA OSD or OSX devices are assigned IP Addresses.
 - The HiperSockets LAN is not assigned an IP Address.
- zLinux HiperSockets Network Concentrator
 - zLinux with HiperSockets and OSD devices is able to bridge traffic without routing overhead providing increased performance.
- zLinux HiperSockets Network Traffic Analyzer
 - Allows Linux on System z to control tracing of the internal virtual LAN.
 - Requires System z10 or later hardware.
 - Requires Linux for System z software.
- z/VM and zLinux Layer 2 Support
 - Both z/VM and zLinux support HiperSockets Layer 2 as well as Layer 3
 - z/OS only supports Layer 3
- z/VSE Fast Path to Linux (LFP) Support
 - Allows communications of z/VSE TCP/IP applications to Linux without a TCP/IP stack on z/VSE.
 - Requires:
 - z196 or later
 - z/VSE V5.1.1
 - LFP in an LPAR
 - HiperSockets Completion Queue



What is Layer 2?



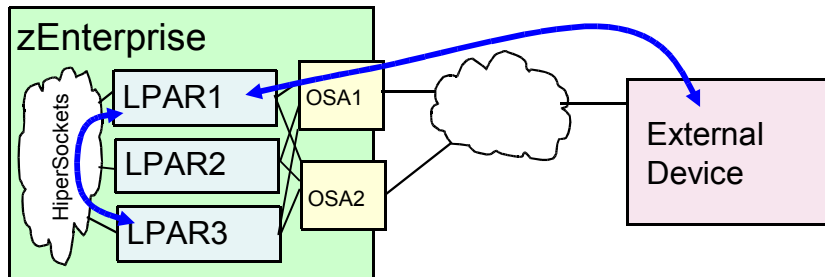
- z/OS Operating System only supports Layer 3 protocol transmission. The data sent/received is either TCP/IP data or SNA (OSA OSE). The Layer 1 and 2 Header and Trailer have been stripped off the packet before it is passed to z/OS.
- Linux on System z supports Layer 2 protocol transmission. Data received is passed to Linux with the Layer 2 Header and Trailer in place.
 - Layer 2 is also referred to as protocol agnostic since the Layer 3 protocol type does not matter (it could be IP, SNA, Apple Talk, anything).
- Can z/OS communicate to Linux even though z/OS is using Layer 3 LAN attachment and Linux is using Layer 2 attachment? YES
- A Layer 2 LAN can refer to a LAN between devices that does not have an IP router (also referred to as a Layer 3 router) in the communication path, all devices are in the same IP subnet.
 - z/OS can attach to a Layer 2 LAN.



Traffic Over HiperSockets or OSA



z/OS to z/OS

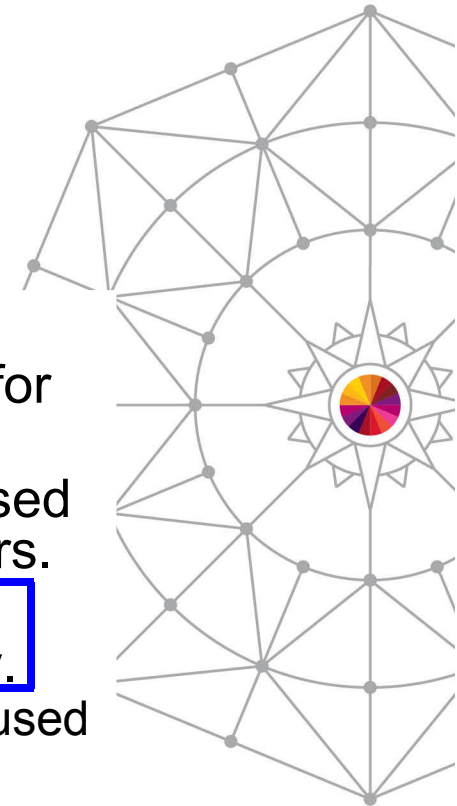


- VIPA and Dynamic routing

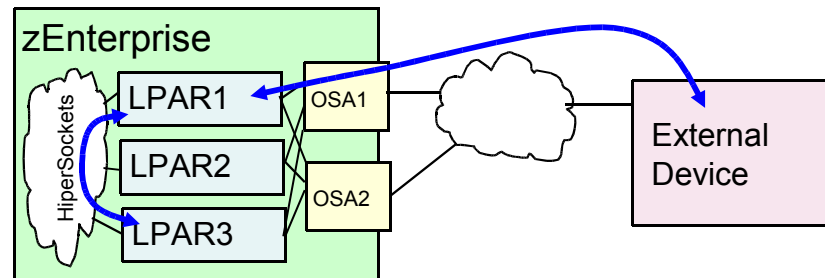
- A single z/OS source IP Address, VIPA, may be used for both same CEC partners as well as external partners.
- A single z/OS destination IP Address, VIPA, may be used by both same CEC partners as well as external partners.

- **It is possible to use HiperSockets between same CEC partners and OSA with non-CEC partners concurrently.**

- OSA between same CEC partners may hypothetically be used as backup in the event that the HiperSockets fails. “Hypothetically” since a HiperSockets failure may indicate serious issues on the CEC.

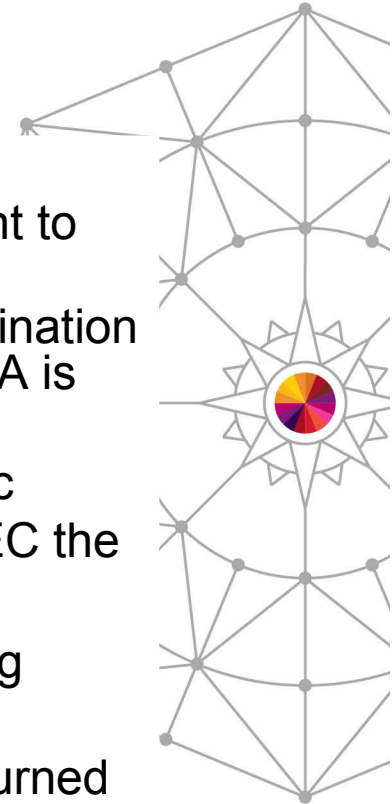


Between z/OS and Linux on System z

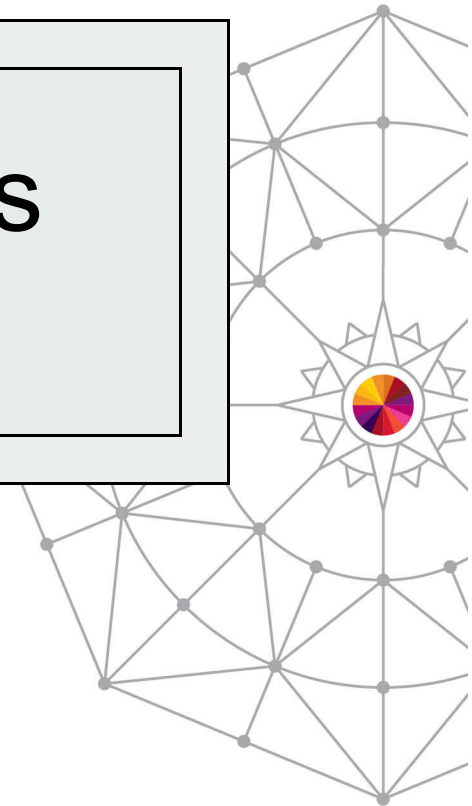


Static routing

- Linux supports dynamic routing but customers are often reluctant to implement it because of the small overhead required.
- Static Host Routes may be used to route traffic with source/destination IP address of the OSA address over the HiperSockets LAN. OSA is then still used for non-CEC partners.
 - Static Routes are manual routes that do not offer dynamic backup, however since HiperSockets is internal to the CEC the lack of backup may be acceptable.
 - Static Routes may be administratively onerous, depending upon the number of host partners.
- Different DNS may be used so that HiperSockets address is returned for LPARs and OSA address is returned for non-CEC partners.
 - As long as Linux application is not bound to OSA Address.



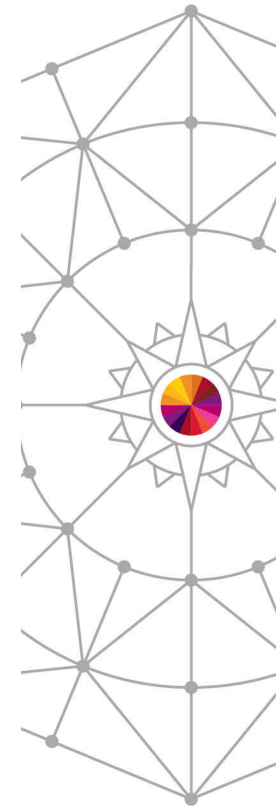
OSA vs. HiperSockets vs. RoCE



OSA, HiperSockets, RoCE



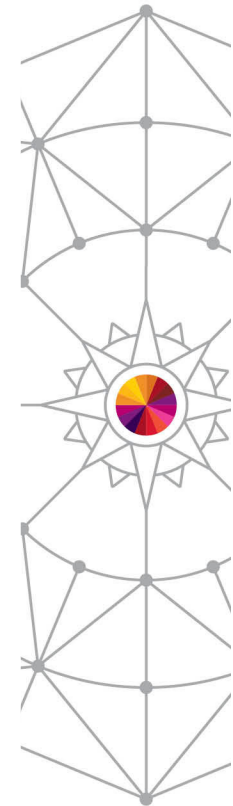
- OSA, HiperSockets, and RoCE are all used for data transfer between hosts.
- Cross CEC zEnterprise traffic
 - OSA and RoCE can be used to send traffic between CECs.
 - HiperSockets only supports traffic between LPARs on a single CEC.
- Different zEnterprise Operating Systems
 - OSA and HiperSockets are supported by multiple Operating Systems (ie. z/OS, z/VM, Linux on System z, etc.)
 - RoCE is only supported by z/OS.
- Traffic outside of a single zEnterprise CEC (and might therefore require additional security measures)
 - OSA traffic can go over the card if shared between LPARs on a single CEC.
 - OSA traffic goes over a network if used to a non-shared-OSA partner.
 - RoCE traffic goes over a 10GbE Layer 2 LAN.
 - Security exposure is limited if contained in a secure location.
 - HiperSockets traffic never goes outside a single CEC.
- Firewall with stateful packet inspection (a PCI (Payment Card Industry) requirement for some traffic)
 - OSA traffic can be sent over a LAN to a Firewall with stateful packet inspection support.
 - HiperSockets traffic cannot be sent over a Firewall with stateful packet inspection support (unless routed traffic).
 - RoCE traffic cannot be sent over a Firewall with stateful packet inspection support (does not support routed traffic).



OSA, HiperSockets, RoCE (cont.)



- Protocol Support
 - OSA supports all TCP/IP protocols and even supports SNA protocol (natively in OSE mode, or when sent with UDP (Enterprise Extender (EE))).
 - HiperSockets only supports IP traffic, it does not support native SNA (so EE must be used to send SNA).
 - RoCE only supports TCP traffic (except IPsec), it does not support native SNA or UDP (EE).
- Required hardware feature
 - OSA feature is required.
 - HiperSockets is part of System z Firmware so it does not require any additional hardware / adapter card purchase.
 - RoCE feature is required. A minimum of 2 per LPAR is recommended.
- CP Overhead
 - OSA provides many different types of offload to the adapter that reduces CP overhead (ARP, Check Sum, Segmentation, etc.)
 - HiperSockets supports zIIP offload to reduce associated cost.
 - RoCE reduces TCP/IP overhead by using RDMA protocol.
- Storage Usage
 - OSA and HiperSockets use CSM fixed storage backed by 64-bit real for data buffers and use Hardware System Area (HAS) memory for routing table.
 - RoCE uses pinned Fixed Memory (mostly 64-bit Common), not CSM managed memory. The maximum amount of memory available for SMC-R with RoCE is definable. When the maximum is reached, new connections will not be SMC-R RoCE eligible.



OSA, HiperSockets, RoCE (cont.)



IP Routing

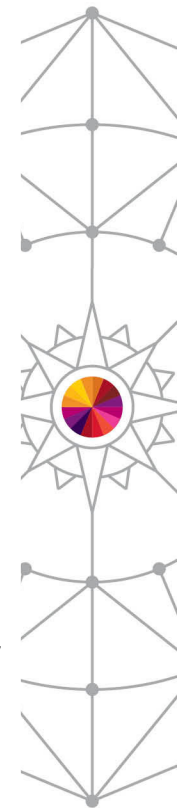
- OSA requires IP routing. OSA does provide dynamic backup with multiple OSAs to the same subnet and Multipath. When multiple OSAs are attached to different subnets or OSA and other attachments, like HiperSockets, are used there is no dynamic backup without a Dynamic Routing protocol (ie. OSPF).
- HiperSockets requires IP routing so there is no dynamic backup without a Dynamic Routing protocol.
 - When HiperSockets is defined as part of a DynamicXCF network routing is handled automatically.
 - When HiperSockets is defined with Integration with IEDN (IQDX) routing is handled automatically.
 - Backup might not be a requirement because if HiperSockets fails there is probably major CEC problems occurring.
- Traffic is automatically / transparently switched from OSA to RoCE. If RoCE connection is not available traffic is sent over OSA. RoCE has automatic / transparent backup to a different RoCE path if one exists.
 - Active RoCE sessions are dropped if the RoCE connection fails and there is no alternate RoCE path, but new sessions will flow using OSA.

IP Routed Traffic

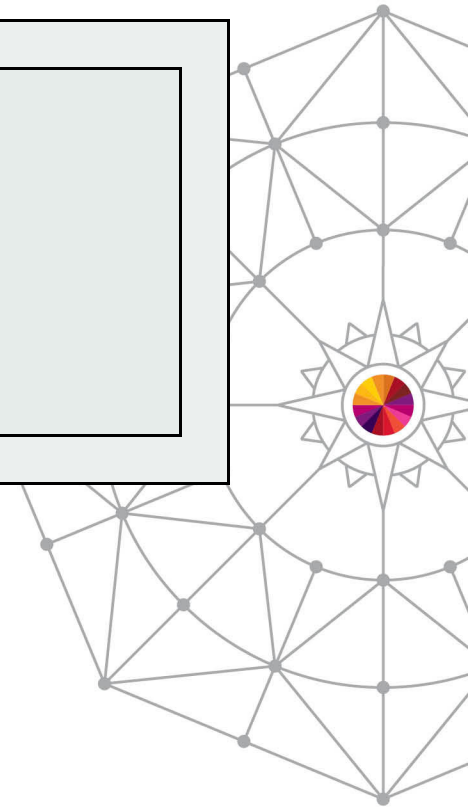
- OSA supports routed traffic. Routed traffic is optimized on z/OS with QDIO Accelerator.
- HiperSockets supports routed traffic (ie. traffic may come in over OSA and then be routed over HiperSockets). Routed traffic is optimized on z/OS with QDIO Accelerator.
- RoCE does not support routed traffic. All connections are over a Layer 2 network.

VLAN Support

- OSA supports VLAN tagging.
- HiperSockets supports VLAN tagging.
- RoCE inherits VLAN IDs from all associated OSAs and supports VLAN tagging.



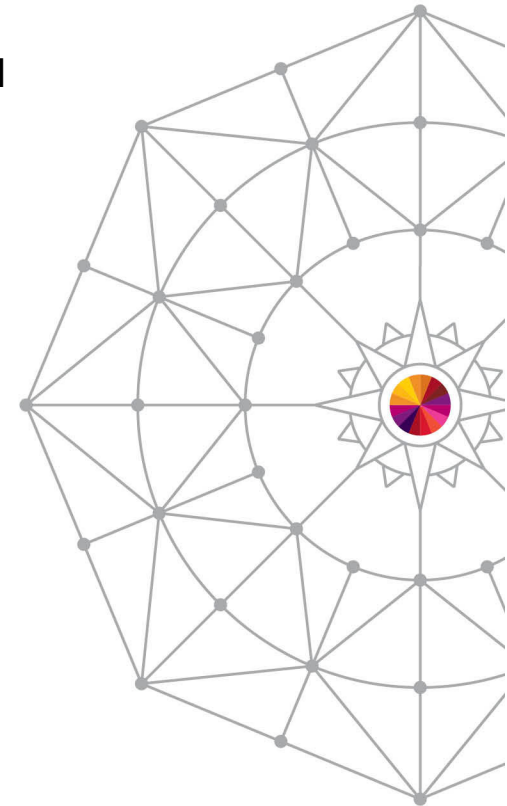
More Information



Web Information



- System z HiperSockets web page:
 - <http://www.ibm.com/systems/z/hardware/networking/products.html>
- IBM ATS Technical Documents:
 - <http://www.ibm.com/support/techdocs>
- z/OS Communications Server
 - <http://www.ibm.com/software/network/commserver/zos>
- IBM Information Center
 - <http://www.ibm.com/support/documentation/us/en>
- IBM Education Assistant
 - <http://www.ibm.com/software/info/education/assistant>
- z/OS Communications Server Publications
 - <http://www.ibm.com/systems/z/os/zos/bkserv>
- IBM Redbooks
 - <http://www.redbooks.ibm.com>
- System z main web site:
 - <http://www.ibm.com/systems/z/hardware/>



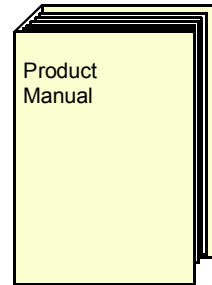
IBM Product Manuals and Redbooks



Some manuals located on the web sites listed on the previous page...

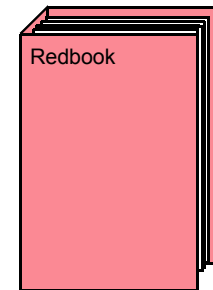
z/OS Communications Server

- IP Configuration Guide, SC27-3650
- IP Configuration Reference, SC27-3651
- IP System Administrator's Commands, SC27-3661
- SNA Network Implementation Guide, SC27-3672
- SNA Resource Definition Reference, SC27-3675
- SNA Operation, SC27-3673



z/OS Hardware Configuration Definition

- HCD Reference Summary, SX33-9032
- HCD Planning, GA32-0907
- HCD User's Guide, SC34-2669
- IOCP User's Guide, SB10-7037

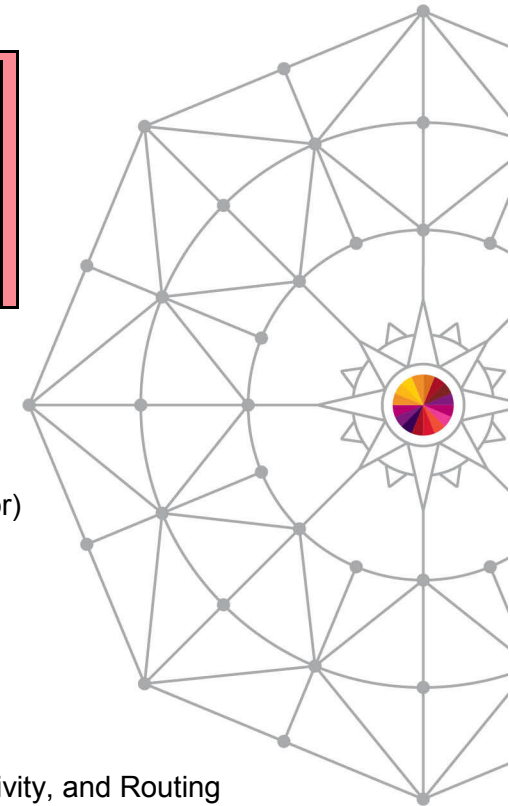


zEnterprise

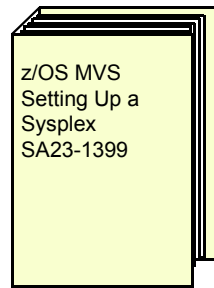
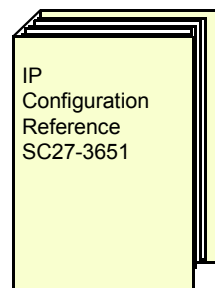
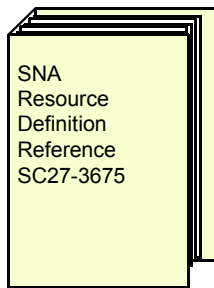
- Install Manual (EC12 model 2827 is GC28-6913, different one for each processor)
- Install for Physical Planning (EC12 model 2827 is GC28-6914, different one for each processor)
- System Overview (EC12 model 2827 is SA22-1088, different one for each processor)
- Ensemble Planning and Configuring Guide, GC27-2608
- Intro to Ensembles, GC27-2609
- HMC Operations Guide for Ensembles, SC27-2615
- HMC Ensembles, SC27-2622

Redbook

- Communications Server for z/OS TCP/IP Implementation Volume 1: Base Functions, Connectivity, and Routing (z/OS V1.13 is SG24-7996, different one for each release)
- HiperSockets Implementation Guide, SG24-6816
- IBM System z Connectivity Handbook, SG24-5444
- I/O Configuration Using z/OS HCD and HCM, SG24-7804
- Building an Ensemble Using Unified Resource Manager, SG24-7921



Appendix: More Configuration Details and Commands

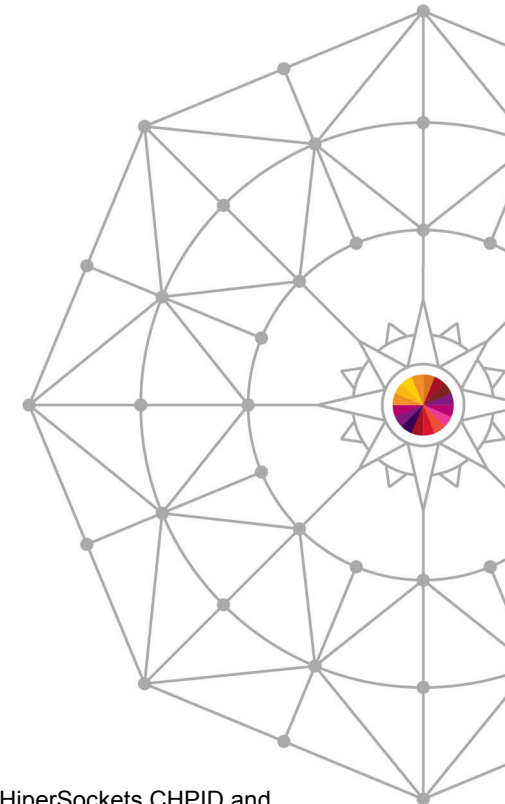


Sample IOCP for HiperSockets



```
*****
* CHPARM values are '00'=16K, '40'=24K, '80'=40K and 'C0'=64K. *
*
* Need at least 3 addresses per z/OS, maximum of 10:
*   - 2 addresses for control
*   - 1 address for data for each TCP stack (between 1 and 8) *
*****
CHPID PATH=(FA), SHARED,
      PARTITION=(LPAR1, LPAR2, LPAR3), (LPAR1, LPAR2, LPAR3)),
      TYPE=IQD, CHPARM=00
CNTLUNIT CUNUMBR=FD00, PATH=(FA), UNIT=IQD
IODEVICE ADDRESS=(FD00, 010), CUNUMBR=(FD00), UNIT=IQD
CHPID PATH=(FB), SHARED,
      PARTITION=(LPAR1, LPAR2, LPAR3), (LPAR1, LPAR2, LPAR3)),
      TYPE=IQD, CHPARM=40
CNTLUNIT CUNUMBR=FD10, PATH=(FB), UNIT=IQD
IODEVICE ADDRESS=(FD10, 010), CUNUMBR=(FD10), UNIT=IQD
CHPID PATH=(FC), SHARED,
      PARTITION=(LPAR1, LPAR2, LPAR3), (LPAR1, LPAR2, LPAR3)),
      TYPE=IQD, CHPARM=80
CNTLUNIT CUNUMBR=FD20, PATH=(FC), UNIT=IQD
IODEVICE ADDRESS=(FD20, 010), CUNUMBR=(FD20), UNIT=IQD
CHPID PATH=(FD), SHARED,
      PARTITION=(LPAR1, LPAR2, LPAR3), (LPAR1, LPAR2, LPAR3)),
      TYPE=IQD, CHPARM=C0
CNTLUNIT CUNUMBR=FD30, PATH=(FD), UNIT=IQD
IODEVICE ADDRESS=(FD30, 010), CUNUMBR=(FD30), UNIT=IQD
```

CHPID	MFS	MTU
CHPARM=00	16K	8K
CHPARM=40	24K	16K
CHPARM=80	40K	32K
CHPARM=C0	64K	56K



- Hardware Configuration Definition (HCD) or I/O Configuration Program (IOCP) must be used to create an IOCDs with the HiperSockets CHPID and subchannel (I/O device) definitions. Because HiperSockets are shared among LPARs, the CHPIDs must be defined as shared in the hardware definitions. A minimum of three subchannel addresses must be configured. One read control device, one write control device, and one data device. Each TCP/IP stack (max of 8) on a single LPAR requires a single data device. The read device must be an even number. The write device must be the read device number plus 1. The data devices can be any device numbers; it does not need to be the next sequential number after the read and write device numbers.
- The Maximum Frame Size (MFS) to be used on a HiperSockets network is set by specifying the "CHPARM=" parameter.
- TCP/IP coding for Maximum Transmission Unit (MTU), will need to correspond to this setting.



QDIO/iQDIO Read Storage

Amount of storage for read processing:

HiperSockets MFS=16K	2 Meg
HiperSockets MFS=24K	3 Meg
HiperSockets MFS=40K	5 Meg
HiperSockets MFS=64K	8 Meg

Configurable value via VTAM Start Option or Link/Interface keyword

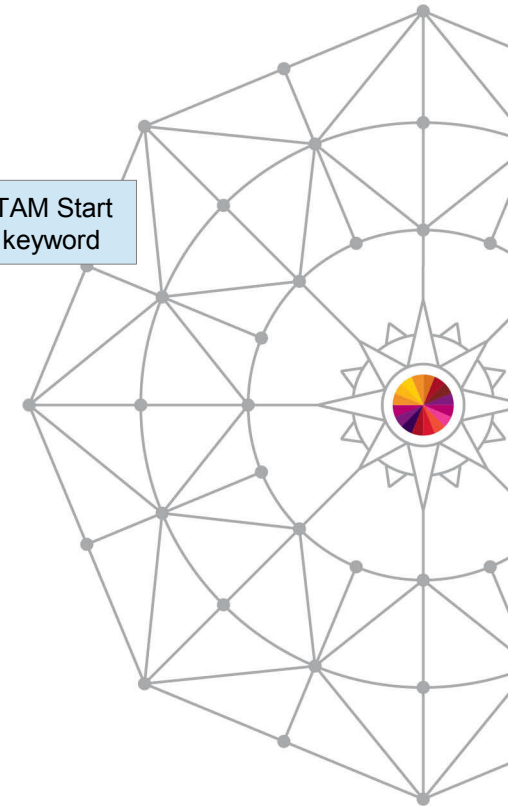
The storage used for read processing is allocated from the CSM data space 4K pool, and is fixed storage backed by 64-bit real. (CSM fixed storage defined in PARMLIB member IVTPRMxx)

OSA QDIO

- 64 SBALs (storage block address lists) x 64K = 4M

HiperSockets

- 126 SBALs x 16K = 2M
- 126 SBALs x 24K = 3M
- 126 SBALs x 40K = 5M
- 126 SBALs x 64K = 8M



VTAM Start Options to Define Storage

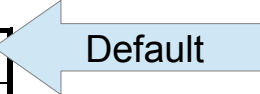
OSA QDIO Read Storage VTAM Start Option QDIOSTG

- Defines how much storage VTAM keeps available for read processing for all OSA QDIO devices

```

+---QDIOSTG==MAX-----+
>>-----+-----+-----+
+---QDIOSTG==+---MAX---+
          +---AVG---+
          +---MIN---+
          +---nnn---+
    
```

MAX	64 SBALs x 64K = 4M
AVG	32 SBALs x 64K = 2M
MIN	16 SBALs x 64K = 1M



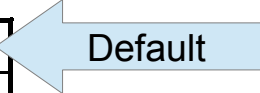
HiperSockets Read Storage VTAM Start Option IQDIOSTG

- Defines how much storage VTAM keeps available for read processing for all HiperSockets devices that use an MFS of 64K

```

+---IQDIOSTG==MAX-----+
>>-----+-----+-----+
+---IQDIOSTG==+---MAX---+
          +---AVG---+
          +---MIN---+
          +---nnn---+
    
```

MAX	126 SBALs x 64K = 8M
AVG	96 SBALs x 64K = 6M
MIN	64 SBALs x 64K = 4M

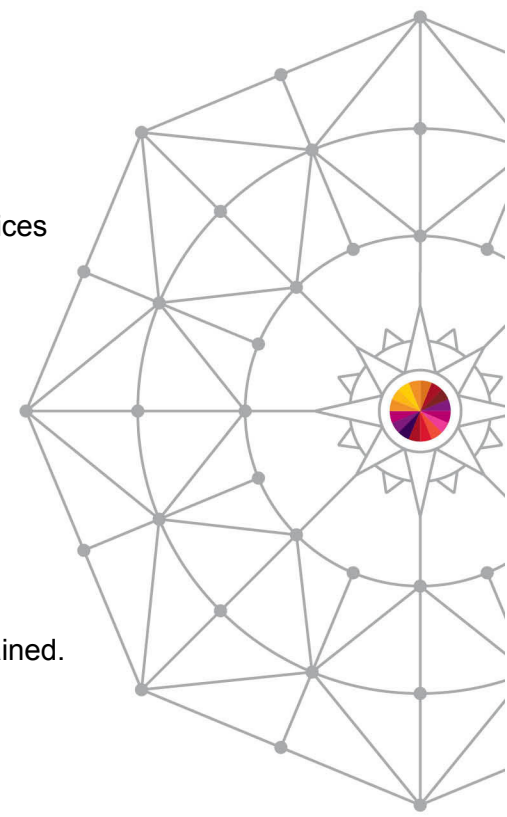


Storage units are defined in terms of QDIO SBALs (QDIO read buffers)

- nnn is the exact number of SBALs in the range 8-126
- MAX allows for the best performance (for example, throughput), but requires more storage.
- MIN may be used for devices with lighter workloads or where system storage might be constrained.
- The amount of storage used is times the number of active QDIO data devices.

Start Option defaults are appropriate for most environments

- Review CSM specifications in PARMLIB member IVTPRMxx and increase, if appropriate
- Use the D NET,CSM to display CSM usage
- Modify storage settings using Start Options, as appropriate
- Use VTAM tuning stats to evaluate needs and usage. Under a typical workload, the NOREADS counter should remain low (close to 0). If this count does not remain low you may need to consider a higher setting for QDIOSTG/IQDIOSTG.
- RMF records send failures can be an indication that the HiperSockets target LP (logical partition) does not have enough storage.



Interface/Link Keyword to Define Storage

```
>>--INTERFACE-intf_name-DEFINE IPAQIDIO-CHPID-chpid---+-----+> ...
+---READSTORAGE GLOBAL-----+
+---READSTORAGE-----+ +---VLANID id---+
+---MAX---+
+---AVG---+
+---MIN---+

>>--LINK--linkname--IPAQIDIO-dev_name---+-----+> ...
+---READSTORAGE GLOBAL-----+
+---READSTORAGE-----+ +---VLANID id---+
+---MAX---+
+---AVG---+
+---MIN---+
```



- Keyword READSTORAGE on LINK and INTERFACE
 - Link statement for IPAQIDIO (IPv4 HiperSockets)
 - Interface statement for IPAQIDIO (IPv4 HiperSockets) or IPAQIDIO6 (IPv6 HiperSockets)
- Override VTAM Start option IQDIOSTG for a specific iQDIO device.
- Global causes the IQDIOSTG VTAM start option values to be used.
 - This is the default.
- MAX, AVG, and MIN
 - Causes the MAX, AVG, or MIN VTAM Start option MAX, AVG, or MIN values to be used.

HiperSockets VLAN ID Summary



- HiperSockets LAN may be divided into separate Virtual LANs (VLANs)
 - DynamicXCF HiperSockets VLANs are required for different TCP/IP stacks in an LPAR to belong to different TCP/IP Subplexes.
 - TCP/IP Profile GLOBALCONFIG IQDVLANid nnnn
 - Manual HiperSockets VLANs may be defined with VLANID on LINK and INTERFACE to divide a HiperSockets LAN into multiple subsets just as OSA QDIO VLAN support.



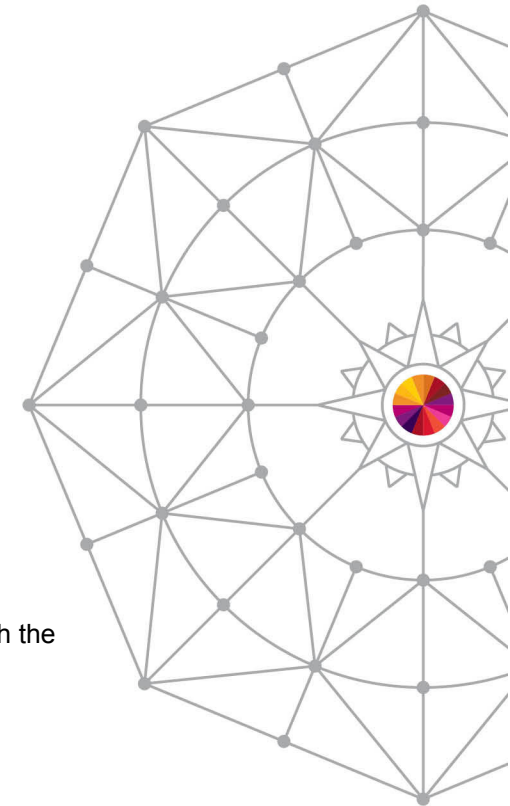
```
+++READSTORAGE GLOBAL-----+
>>--LINK--linkname--IPAQIDIO--dev_name-----+-----+> ...
+++READSTORAGE-----+-----+
                                     +---VLANID id---+
                                     +---MAX---+
                                     +---AVG---+
                                     +---MIN---+
```



HiperSockets Defined in OSPF



- Subnet mask and metric on DynamicXCF PROFILE statement is for ORouted use only.
 - OSPF_Interface in OMPROUTE config required for OSPF advertisements to be sent over the XCF network.
- Define DynamicXCF links in OMPROUTE with wildcard IP addresses
 - Avoids having to determine what the names of the various flavors of dynamic XCF links will be.
 - The syntax checker does not require that the "Name" be specified.
- Point-to-Multipoint Networks
 - MPC, XCF, IUTSAMEH
 - Unicast to Each Interface: Hello (Type 1)
 - Does not require DR election
- Broadcast Multiaccess Network
 - Token Ring, Ethernet, FDDI, LANE, **HiperSockets**
 - Multicast to 224.0.0.5: Hello (Type 1)
 - Requires DR election
 - OSPF_INTERFACE NON_BROADCAST=YES should not be defined.
- The HELLO protocol determines who the Designated Router (DR) will be.
- Role of the DR:
 - It is adjacent to all other routers on the network.
 - It generates and floods the network link advertisements on behalf of the network.
 - Reduces amount of router protocol traffic, as only the DR is responsible for flooding the network with the information.
 - It is responsible for maintaining the network topology database.
- Router with highest Router_Priority becomes DR on a broadcast multiaccess network.
 - If there is a tie, the router with the higher Router_ID becomes the DR.
 - If the Router ID is not specified, the IP address of one of the OSPF interfaces will be used as Router ID.
 - Define IP address of static VIPA or physical interface for RouterID to avoid selection of a Dynamic VIPA which could move.
- If your z/OS system is not to be used primarily for routing, consider setting Router_Priority to 0 for all non-HiperSockets interfaces so that the system is ineligible to become the DR.



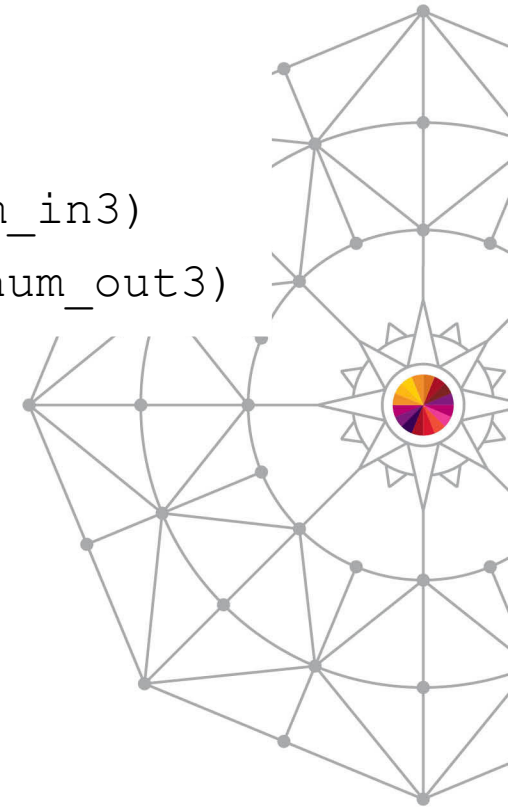
ParmLib



- **SYS1.PARMLIB(COUPLExx)**

```
PATHIN  DEVICE(dev_num_in1,dev_num_in2,dev_num_in3)
```

```
PATHOUT DEVICE(dev_num_out1,dev_num_out2,dev_num_out3)
```



VTAM Start Options

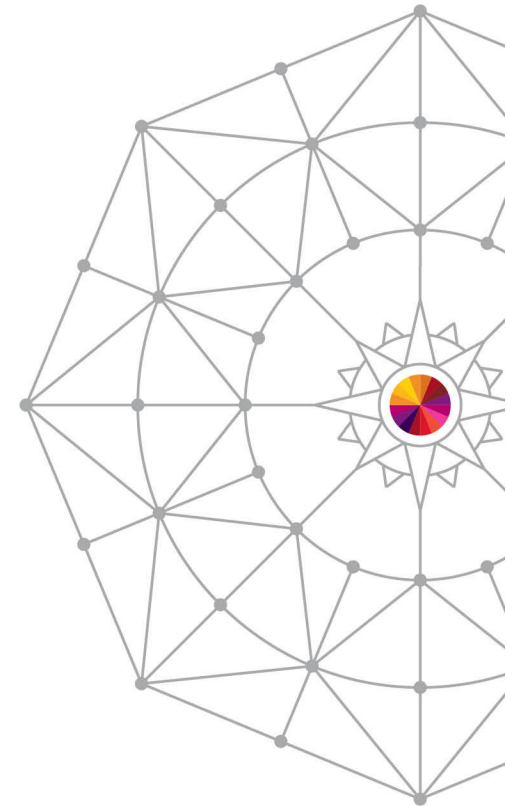


· SYS1.VTAMLST(ATCSTRxx)

```
+---IQDCHPID-----+
>>-----+-----+-----><
+---IQDCHPID=---+---ANY-----+
                +---NONE-----+
                +---chpid---+

>>-----+-----+-----><
+---XCFGRPID=group_id---+

+---XCFINIT=YES-----+
>>-----+-----+-----><
+---XCFINIT=---+---NO-----+
                +---YES-----+
                +---DEFINE---+
```

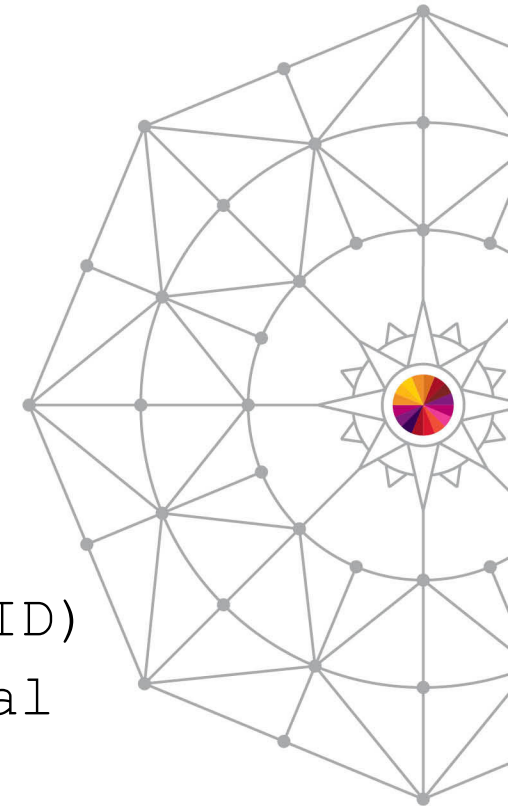


VTAM Display



- Display Commands

```
DISPLAY NET,TRL
DISPLAY NET,TRL,TRLMN=trl_maj_node
DISPLAY NET,TRL,TRLMN=ISTTRL
DISPLAY NET,TRL,TRLE=trl_entry_name
DISPLAY NET,TRL,TRLE=IUTSAMEH
DISPLAY NET,TRL,TRLE=IUTIQDIO
DISPLAY NET,TRL,TRLE=IUTIQDxx (xx=CHPID)
DISPLAY NET,TRL,TRLE=ISTTlsrs (ls=local
&SYSCLONE,rs=remote &SYSCLONE)
DISPLAY NET,TRL,XCFCP=cp_name
DISPLAY NET,VTAMOPTS
```



Device/Link



- PROFILE.TCPIP statement

```

+---NOAUTORestart---+
>>---DEVIce---device_name---MPCIPA-----+-----><
+---AUTORestart-----+

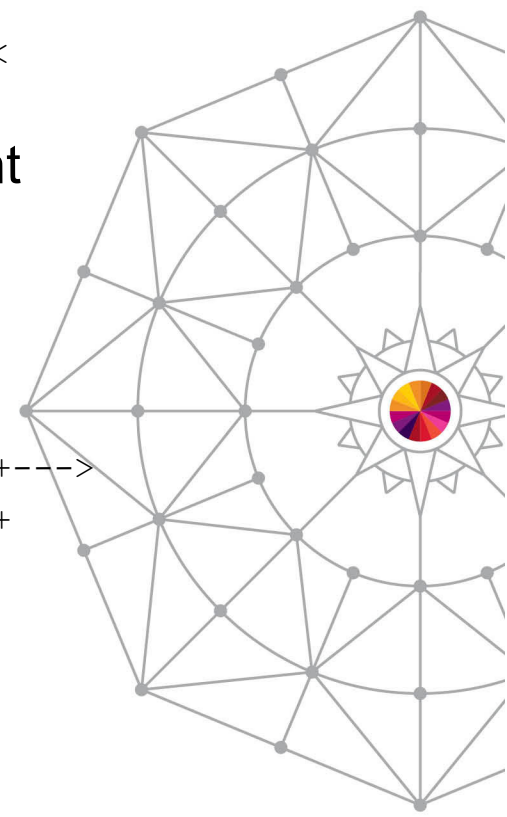
```

- The device name coded on the DEVICE statement must be IUTIQDxx (where xx = CHPID)
- VTAM TRLE is dynamically built
 - TRLE name = device_name

```

>>---LINK---link_name---IPAQIDIO---device_name---+-----+--->
+---IPBCAST---+
+---READSTORAGE GLOBAL-----+
>---+-----+-----+-----+-----+----->
+---REASTORAGE---+---MAX---+---+ +---VLANid---+
+---AVG---+
+---MIN---+
+---SECCLASS---255-----+ +---NOMONSYSPLEX---+
>---+-----+-----+-----+-----+-----><
+---SECCLASS---security_class---+ +---MONSYSPLEX-----+

```

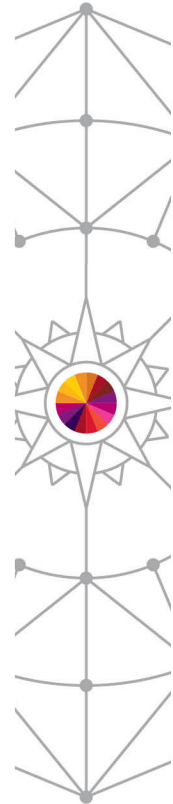


IPAQIDIO6



• PROFILE.TCPIP

```
>>---INTERFace---intf_name---+---DEFINE---IPAQIDIO6---| Interface Definition |--+---><
+---+---DELEte-----+-----+
|                               +-----+
|                               v         |
+---ADDADDR-----| IPaddr Spec |--+---+
|                               +-----+
|                               v         |
+--- ADDADDR-----| IPaddr Spec |--+---+
|                               +-----+
|                               v         |
+---ADDADDR-----| IPaddr Spec |--+---+
```



Interface Definition:

```
|---CHPID---chpid---+-----+-----+-----+-----+-----+-----+-----+-----+----->
+---INTFID---interface_id---+ |                               +-----+
+---IPADDR-----| IPaddr Spec |--+-----+

+---READSTORAGE GLOBAL-----+
>+-----+-----+-----+-----+-----+-----+-----+-----+----->
+---READSTORAGE---+---MAX---+---+ +---VLANID---id---+
+---AVG---+
+---MIN---+
>+-----+-----+-----+-----+-----+-----+-----+-----+----->
+---SOURCEVIPAINterface---vipa_name---+

+---SECCLASS---255-----+ +---NOMONSYSPLEX---+
>+-----+-----+-----+-----+-----+-----+-----+-----+-----|
+---SECCLASS---security_class---+ +---MONSYSPLEX-----+
```

IPaddr Spec:

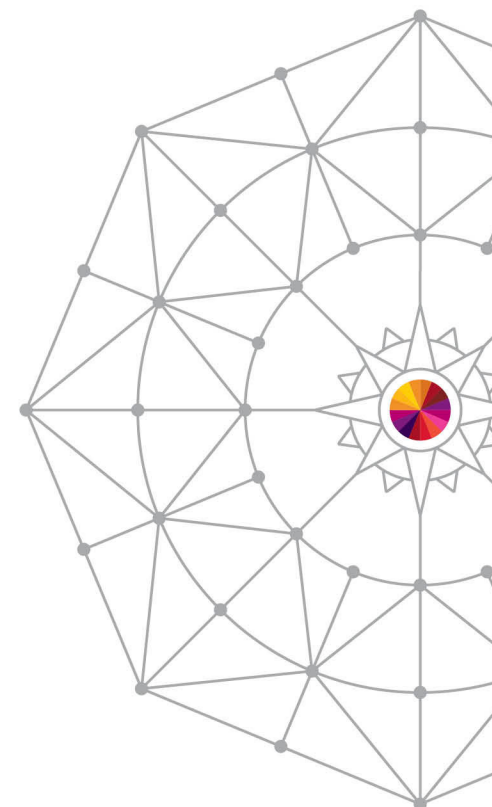
```
|---+---ipv6_address-----+---|
+---prefix/prefix_length---+
```

Start and Stop



- PROFILE.TCPIP

```
>>---START---+---device_name-----+---<<
          +---interface_name----+
>>---STOP---+---device_name-----+---<<
          +---interface_name----+
```



GlobalConfig

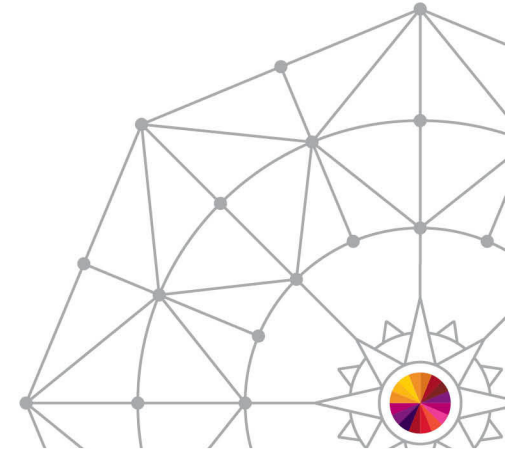
Only some GlobalConfig parameters are listed here. See the IP Configuration Reference for a full list of GlobalConfig parameters.

PROFILE.TCPIP

```

>>---GLOBALCONFig-----+----->>
      |
      |      +---ALLTRAFFIC-----+
      |      | +---AUTOIQDX---+-----+
      |      | | +---NOLARGEDATA---+ |
      |      +-----+
      |      | +---NOAUTOIQDX-----+
      |      |
      |      | +---NOIQDMULTIWRITE---+
      |      | +---IQDMULTIWRITE-----+
      |      +-----+
      |      | +---IQDVLANid---vlan_id---+
      |      |
      |      | +---NOSEGMENTATIONOFFLoad---+
      |      | +---SEGMENTATIONOFFLoad-----+
      |      +-----+
      |      | +---SYSPLEXMONitor---|Sysplex Options|-----+
      |      |
      |      | +---NOWLMPRIORITYQ-----+
      |      |
      |      | +---default_control_values-----+
      |      | | +---WLMRIORITYQ---+-----+
      |      | | +---IOPRIn---control_values---+
      |      |
      |      |
      |      | +-----+
      |      | | +---NOIPSECURITY---+
      |      | | +---NOIQDIOMULTIWRITE---+
      |      | | +---IQDIOMULTIWRITE-----+
      |      +-----+
      +---ZIIP-----+-----+

```



Sysplex Options:

```

+---NOAUTOREJOIN---+
| +---AUTOREJOIN-----+
| +---NODELAYJOIN---+
+-----+
| +---DELAYJOIN---+
+-----+
| +---NOJOIN---+
| +---NOMONINTERFACE---NODYNROUTE-----+
+-----+
| +---NODYNROUTE---+
| +---NOMONINTERFACE---+-----+
| +---DYNROUTE-----+
| +---MONINTERFACE---+-----+
| +---NODYNROUTE---+
+-----+
| +---NORECOVERY---+
| +---RECOVERY-----+
| +---TIMERSECS---60-----+
+-----+
| +---TIMERSECS---seconds---+
+-----+

```

IPConfig

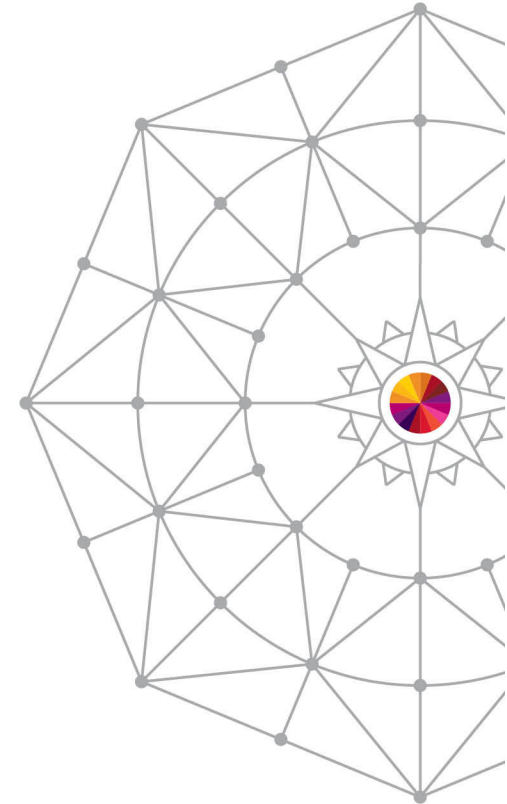
PROFILE.TCPIP

Only some IPConfig parameters are listed here. See the IP Configuration Reference for a full list of IPConfig parameters.



```

>>---IPCONFig-----+-----+----->>
| +---ARPTO---1200-----+
| +---ARPTO---ARP_cache_timeout---+
| +---CHECKSUMOFFLoad-----+
| +---NOCHECKSUMOFFLoad---+
|
| +---NOFWMULTipath-----+
| +---DATAGramfwd---+ +---FWMULTipath---PERPacket---+
|
| +---NODATAGramfwd-----+
| +---DEVRETRYDURation---90-----+
| +---DEVRETRYDURation---dev_retry_duration---+
|
| +---NOIQDIORouting-----+
| +---IQDIORouting---+ +---QDIOPriority---1-----+
| +---QDIOPriority---priority---+
| +---NOMULTIPATH-----+
| +---MULTIPATH---+ +---PERConnection---+
| +---PERPacket-----+
| +---NOPATHMTUDIScovery---+
| +---PATHMTUDIScovery-----+
| +---NOQDIOACCElerator-----+
| +---QDIOACCElerator---+ +---QDIOPriority---1-----+
| +---QDIOPriority---priority---+
| +---REASSEMBLytimeout---60-----+
| +---REASSEMBLytimeout---reassembly_timeout---+
| +---NOSEGMENTATIONOFFLoad---+
| +---SEGMENTATIONOFFLoad-----+
| +---NOSOURCEVIPA---+
| +---SOURCEVIPA---+
|
| +---NOSYSPLERouting---+
| +---SYSPLERouting-----+
| +---NOTCPSTACKSOURCEVipa-----+
| +---TCPSTACKSOURCEVipa---vipa_addr---+
| +---TTL---64-----+
| +---TTL---time_to_live---+
  
```



TCP/IP NetStat

- Display Commands

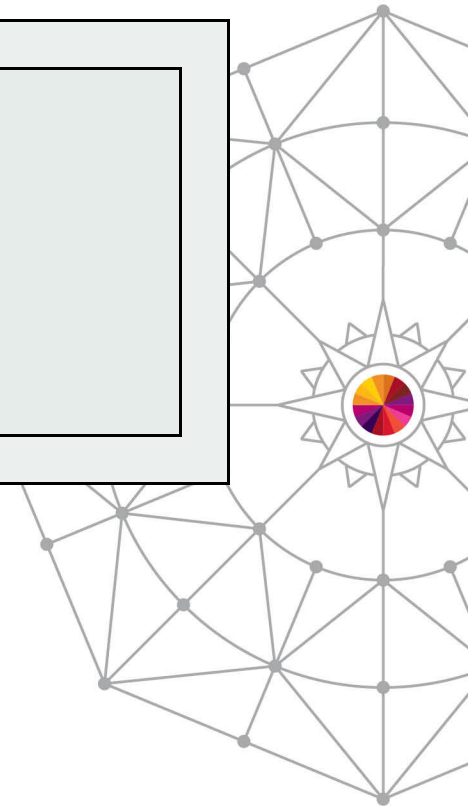
```

>>---Display---TCPiP,tcp_proc,Netstat,-----DEVlinks-----+-----
|               |               |               |               |               |               |
+---Home-----+ +---,INTFName=intfname---+ +---,FORMat=---LONG---+ +---,MAX=*-----+
|               |               |               |               |               |               |
+---SHORT---+ +---,MAX=recs---+
|
+-----+
|   v   |
>>---Display---TCPiP,tcp_proc,ROUTE-----+-----+----->>
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---,ADDRTYpe=---IPV4---+ +---IPAddr=---ipaddr-----+ +---,FORMat=---LONG---+ +---,MAX=*-----+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---DETAIL-----+---IPV6---+ +---ipaddr/prefixLen---+ +---SHORT---+ +---,MAX=recs---+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---IQDIO-----+ +---ipaddr/subnetmask---+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---PR---+---ALL-----+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---QDIOACCEL---+---prname---+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---RADV-----+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---RSTAT-----+
|
TSO
>>---NETSTAT---+---DEVlinks-----+---TCp---tcpname---+---FORMat---+---LONG---+---(INTFName---intfname---+>>
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---Home-----+ +---TCp---tcpname---+ +---FORMat---+---LONG---+ +---(INTFName---intfname---+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---SHORT---+ +---REPort-----+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---STACK---+---DSN---dsname---+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---TITLes---+ +---HLQ---hlqname---+
|
+-----+
|   v   |
>>---NETSTAT---ROUTE-----+-----+----->>
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---ADDRTYpe---+---IPV4---+ +---TCp---tcpname---+ +---FORMat---+---LONG---+ +---(IPAddr---+---ipaddr-----+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---DETAIL-----+---IPV6---+ +---REPort---+---SHORT---+ +---IPAddr---+---ipaddr/prefixLen---+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---IQDIO-----+ +---ipaddr/subnetmask---+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---PR---+---ALL-----+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---QDIOACCEL---+---prname---+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---RADV-----+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---RSTAT-----+
|
z/OS Unix
>>---netstat---+--- -d ---+ +--- -p ---tcpname---+ +--- -M ---+---LONG---+ +--- -K ---intfname---+>>
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+--- -h ---+ +--- -p ---tcpname---+ +--- -M ---+---LONG---+ +--- -K ---intfname---+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---SHORT---+
|
>>---netstat--- -r -----+-----+----->>
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+--- -I ---+ +--- -V ---+ +--- -I ---+ +--- -V ---+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---ADDRTYpe---+---IPV4---+ +--- -p ---tcpname---+ +--- -M ---+---LONG---+ +--- -I ---+ +--- -V ---+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---DETAIL-----+---IPV6---+ +--- -p ---tcpname---+ +--- -M ---+---LONG---+ +--- -I ---+ +--- -V ---+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---IQDIO-----+ +---ipaddr/prefixLen---+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---PR---+---ALL-----+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---QDIOACCEL---+---prname---+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---RADV-----+
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
+---RSTAT-----+

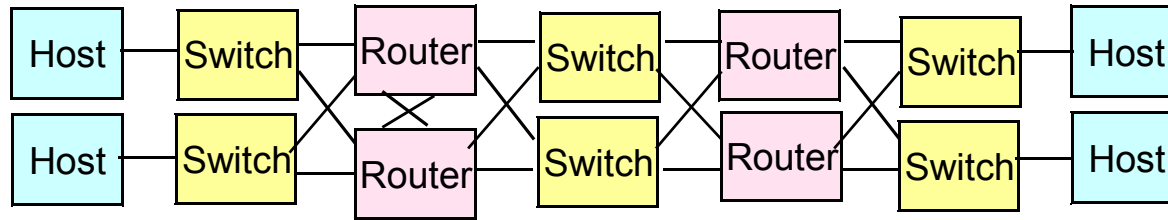
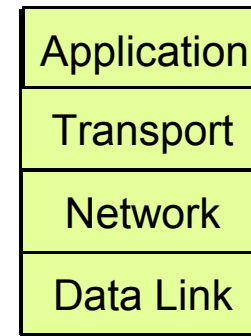
```



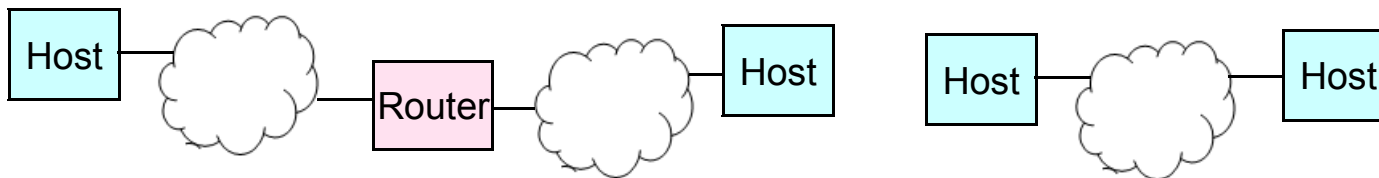
Appendix: TCP/IP Routing



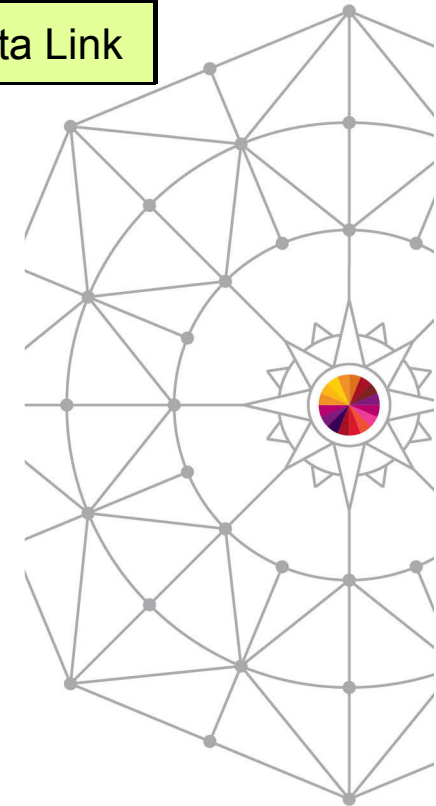
Physical and Logical Network



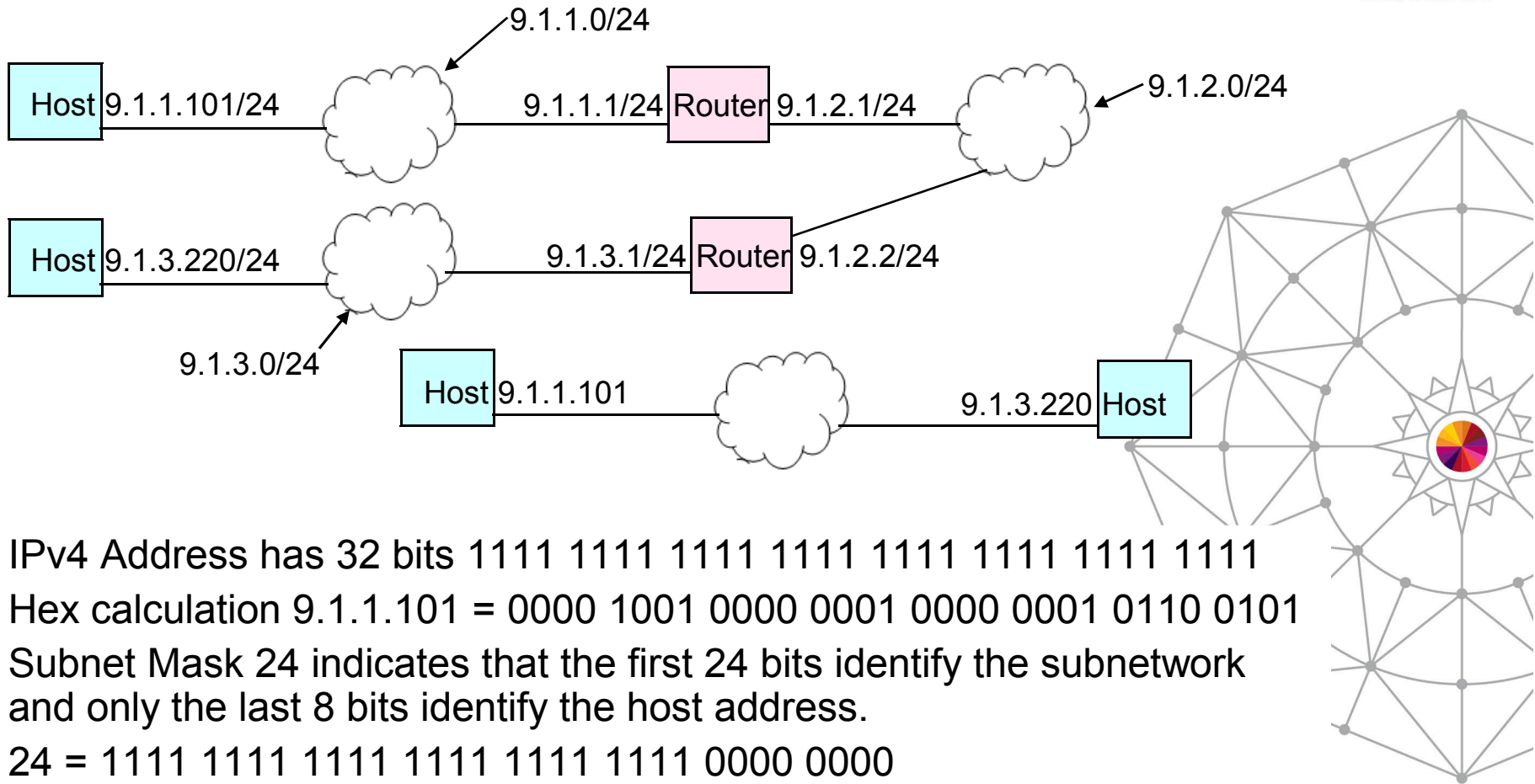
- There is some physical network.
- It can be shared between TCP/IP and SNA, and even other protocols concurrently.
- Devices are connected with cables.
- Each device network connection has a Media Access Control (MAC) address that is used for sending and receiving messages.



- There are logical network views that use the underlying physical network.
 - The level of detail depends upon the discussion.
- Each device network connection has a Logical address that is used for sending and receiving messages.

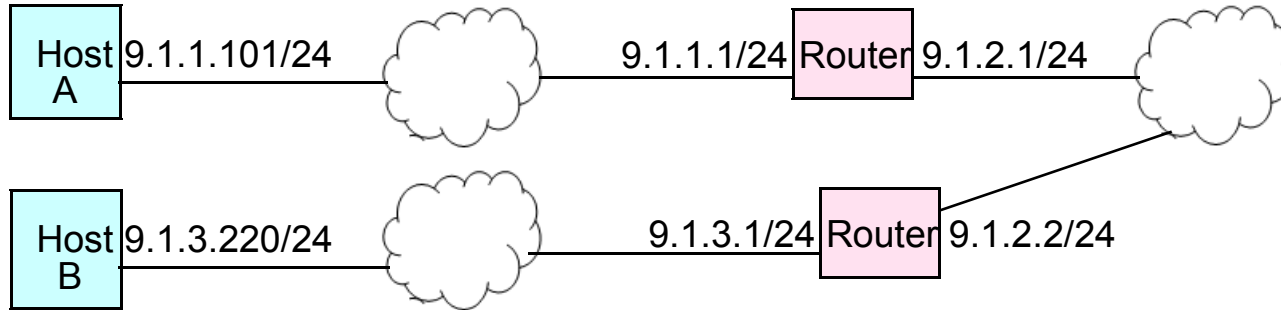


TCP/IP Network



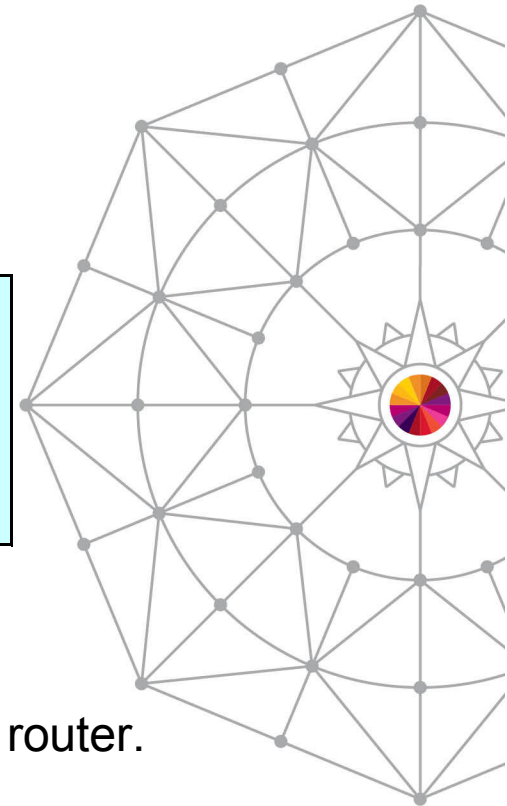
- IPv4 Address has 32 bits 1111 1111 1111 1111 1111 1111 1111 1111
- Hex calculation 9.1.1.101 = 0000 1001 0000 0001 0000 0001 0110 0101
- Subnet Mask 24 indicates that the first 24 bits identify the subnetwork and only the last 8 bits identify the host address.
- 24 = 1111 1111 1111 1111 1111 1111 0000 0000
 - Also referred to as 255.255.255.0

TCP/IP Routing



Host A:
Direct Route to 9.1.1.0/24
Indirect Route to 9.1.3.0/24 via 9.1.1.1
Default Route to 9.1.1.1

- Direct Routes
 - Subnetworks that this host connects to.
- Indirect Routes
 - Subnetworks that this host can reach by routing through a router.
- Default Routes
 - Where to send messages when there is no explicit route.



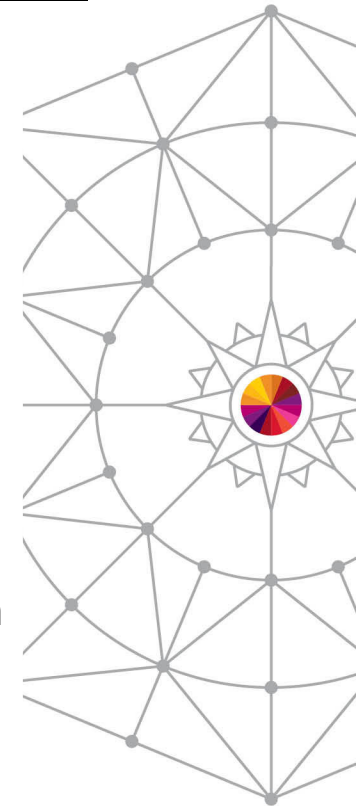
Routing Table

Subnet mask of 32 indicates Host Routes, the most specific type of route possible. More specific routes always take precedence over less specific routes. Subnet or Network routes, in this case subnet route with mask of 24, are less specific than host routes but more specific than default routes. Default routes are the least specific routes possible and are therefore only used when no other route matches the destination.

z/OS B Routing Table

Destination	Gateway	Interface
10.3.1.202/32	10.20.1.202	HiperB
10.10.4.0/24	0	OsaB
10.20.1.0/24	0	HiperB
DEFAULT	10.10.4.1	OsaB

- The destination IP Address, IP subnet, or IP network is in the first column above.
 - DEFAULT is a special destination keyword indicating where to send packets when the destination does not match any other route in the table.
- Gateway IP Address is in the second column above. A zero, “0”, indicates there is no next hop gateway.
 - A route without a next hop gateway is known as a direct route.
 - A route with a next hop gateway is known as an indirect route.
- The interface name or link name is in the third column above to indicate which interface the packet is to be sent over.
- Each route in the table may be read as “If this destination, then send packet to this gateway over this interface”.



Source and Destination IP Addresses

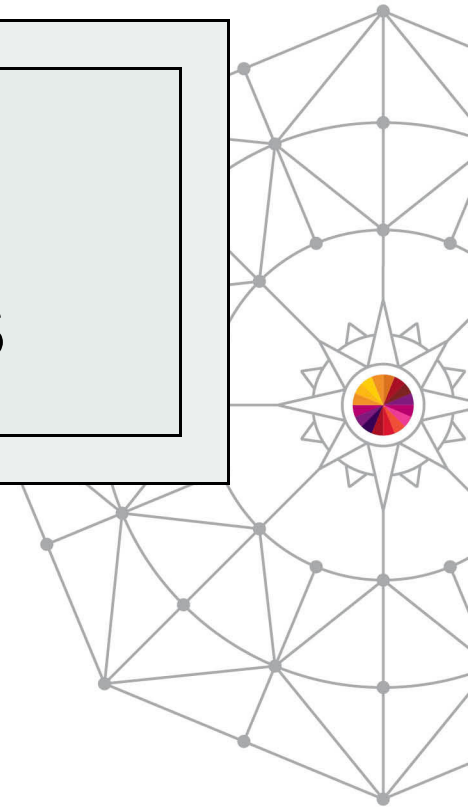


TCP/IP Application Session

- One partner initiates the connection and sends a connection request packet. For TCP connections the host that sends the connection request is the client.
 - What destination IP address to connect to?
 - *If hostname is given then IP address is resolved by DNS or Local Host Name file (IPNodes).*
 - What source IP address to use in the initiate packet?
 - *Hierarchy of source IP address is detailed in z/OS Communication Server (CS) IP Configuration Guide, SC31-8775.*
 1. *Sendmsg() using the IPV6_PKTINFO ancillary option specifying a nonzero source address (RAW and UDP sockets only)*
 2. *Setsockopt() IPV6_PKTINFO option specifying a nonzero source address (RAW and UDP sockets only)*
 3. *Explicit bind to a specific local IP address*
 4. *bind2addrsel socket function (AF_INET6 sockets only)*
 5. *PORT profile statement with the BIND parameter*
 6. *SRCIP profile statement (TCP connections only)*
 7. *TCPSTACKSOURCEVIPA parameter on the IPCONFIG or IPCONFIG6 profile statement (TCP connections only)*
 8. *SOURCEVIPA: Static VIPA address from the HOME list or from the SOURCEVIPAINTERFACE parameter*
 9. *HOME IP address of the link over which the packet is sent*
- The other partner receives the connection request packet and responds. For TCP connections the host that receives the connection request is the server.
 - The host that receives the connection request takes the source IP address from the connection packet and uses that for the destination IP address in the packet that it sends back.
 - The host that receives the connection request takes the destination IP address from the connection packet and uses that for the source IP address in the packet that it sends back.
 - These source and destination assignments are usually used for the life of the connection.



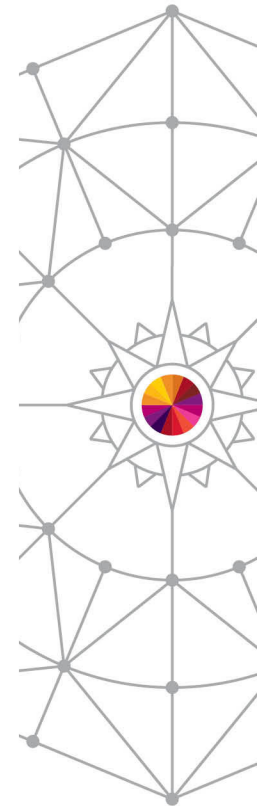
Appendix: DynamicXCF Details



VTAM & TCP/IP Sysplex Support Background



- Each VTAM in the Sysplex:
 - Joins ISTXCF and ISTCFS01 Sysplex groups
 - Establishes DynamicXCF Connectivity (XCFINIT start option)
 - Can access Generic Resource and Multinode Persistent Session (MNPS) structures (STRGR & STRMNPS start options)
- Each TCP/IP stack in the Sysplex:
 - Joins EZBTCPCS Sysplex group
 - Exchanges IP address information
 - Coordinates DVIPA movement
 - Can be a Sysplex Distributor target
 - Can access Sysplex-Wide Security Associations (SWSA) and Sysplexports structures
 - Sets up TCP/IP DynamicXCF connectivity
 - Dynamic Same Host - if stacks are in the same LPAR
 - Dynamic HyperSockets - if stacks are on the same CEC
 - VTAM's DynamicXCF Connectivity
- DynamicXCF requires cross-system coupling facility (XCF) messaging support. The system parameter PLEXCFG (in the IEASYSxx Member of Parmlib) defines XCF support (XCF messaging and signaling).
- Some APPN enablement is required for TCP/IP DynamicXCF. If APPN searching and routing are not desired, the following VTAM Start Options should be researched:
 - NODETYPE=EN
 - CONNTYPE=LEN
 - CPCP=NO
 - HPR=RTP (default)
 - XCFINIT=YES (default)
 - The above Start Options should allow enough APPN function to be present in VTAM to support DYNAMICXCF but allow searching, session initiation, and routing to remain subarea only).



Sysplex and Subplex Background



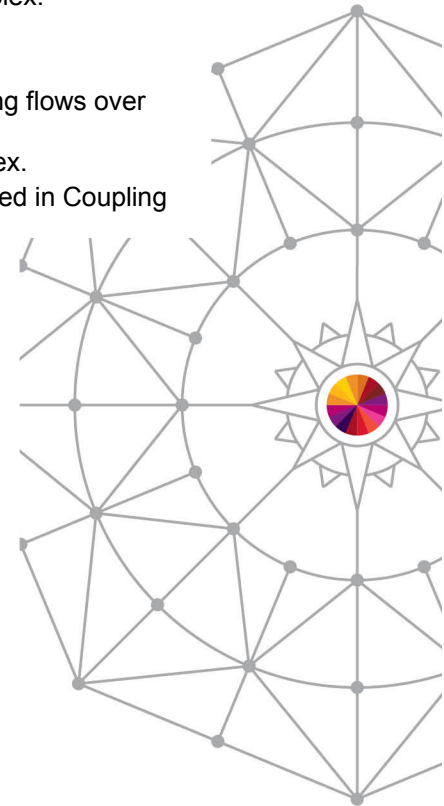
	Base Sysplex	Parallel Sysplex
XCF signaling or messaging	Yes	Yes
Coupling Facility installed	No	Yes
Structures defined	No	Yes
DVIPA	Yes	Yes

Base Sysplex

- In Base Sysplex XCF signaling or messaging flows over ESCON CTCs, or ICP (Internal Coupling Peer Channel) in CEC (COUPLExx).
- Coupling Facility is not required for Base Sysplex.
- Structures are not defined in a Base Sysplex.

Parallel Sysplex

- In Parallel Sysplex XCF signaling or messaging flows over Coupling Facility (COUPLExx).
- Coupling Facility is required for Parallel Sysplex.
- Structures are defined in COUPLExx and stored in Coupling Facility in a Parallel Sysplex.



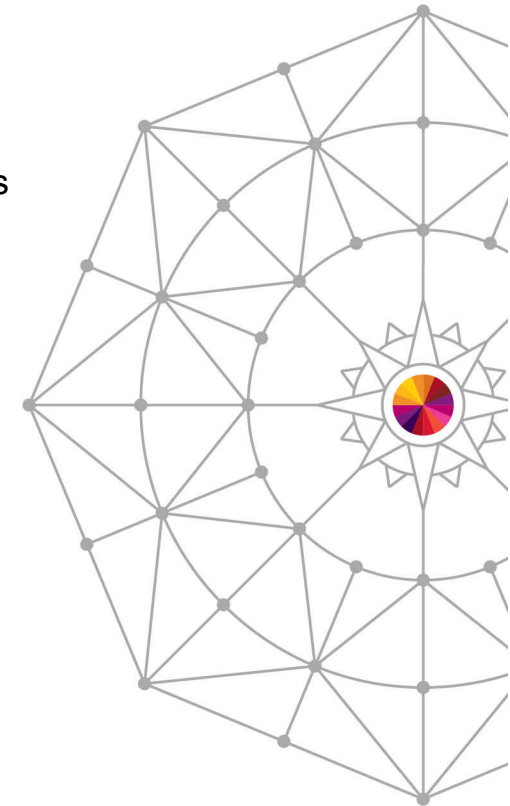
- Partitioning the Sysplex into Subplexes restricts use of Sysplex functions to multiple security areas.
- A Subplex is a subset of VTAM nodes or TCP stacks in the Sysplex.
 - All the members of a subplex can establish dynamic connectivity through the Sysplex with other members of the subplex.
 - Members of the subplex cannot establish dynamic connectivity through the Sysplex with non-members of the subplex.
- Functions restricted to within a Subplex (Subplex Scope):
 - VTAM Generic Resources (GR) & Multi-Node Persistent Session (MNPS) resources
 - Automatic connectivity - IP connectivity (DynamicXCF) and VTAM connectivity over XCF
 - TCP/IP stack IP address awareness and visibility (including DVIPA)
 - DVIPA movement candidates
 - Sysplex Distributor target candidates



Subplex Support



- Each VTAM can belong to only one subplex
 - The subplex is defined by the name of the Sysplex groups that VTAM joins
 - One VTAM Subplex per LPAR
- Each TCP/IP stack can belong to only one subplex
 - The subplex is defined by the name of the Sysplex group that TCP/IP joins
 - A TCP/IP Subplex cannot span multiple VTAM Subplexes
 - Different TCP/IP stacks in an LPAR may belong to different TCP/IP Subplexes
- VTAM Subplex Support
 - Start Option XCFGRPID vv
 - where vv is a number between 2 and 31
 - VTAM joins ISTXCFvv and ISTCFSvv Subplex groups
 - STRGR and STRMNPS CF structure names are suffixed with vv
- TCP/IP Subplex Support
 - Profile GLOBALCONFIG statement
 - XCFGRPID tt to subplex TCP/IP
 - where tt is a numeric value between 2 and 31
 - **IQDVLANID nn to subplex HiperSockets for DynamicXCF connectivity**
 - where nn is a numeric value between 1 and 4094
 - TCP/IP will join Sysplex group EZBTvvtt
 - where vv is the VTAM subplex number
 - SWSA and Sysplexports structure names will be suffixed by vvtt
 - EZBDVIPAvvtt and EZBEPORtvvtt
 - **Profile keyword on LINK and INTERFACE statements for HiperSockets**
 - VLANID nn
 - where nn is numeric value between 1 and 4094



The End

