



z/OS Communications Server Performance Functions Update

*David Herr – dherr@us.ibm.com
IBM Raleigh, NC*

*Wednesday, August 12, 2015: 11:15 AM - 12:15 PM,
Dolphin, Asia 2*



#SHAREorg



**SHARE is an independent volunteer-run information technology association
that provides education, professional networking and industry influence.**

Copyright (c) 2015 by SHARE Inc. Except where otherwise noted, this work is licensed under
<http://creativecommons.org/licenses/by-nc-sa/3.0/>





Trademarks

The following are trademarks of the International Business Machines Corporation in the United States and/or other countries.

AIX*	DB2*	HiperSockets*	MQSeries*	PowerHA*	RMF	System z*	zEnterprise*	z/VM*
BladeCenter*	DFSMS	HyperSwap	NetView*	PR/SM	Smarter Planet*	System z10*	z10	z/VSE*
CICS*	EASY Tier	IMS	OMEGAMON*	PureSystems	Storwize*	Tivoli*	z10 EC	
Cognos*	FICON*	InfiniBand*	Parallel Sysplex*	Rational*	System Storage*	WebSphere*	z/OS*	
DataPower*	GDPS*	Lotus*	POWER7*	RACF*	System x*	XIV*		

* Registered trademarks of IBM Corporation

The following are trademarks or registered trademarks of other companies.

Adobe, the Adobe logo, PostScript, and the PostScript logo are either registered trademarks or trademarks of Adobe Systems Incorporated in the United States, and/or other countries.

Cell Broadband Engine is a trademark of Sony Computer Entertainment, Inc. in the United States, other countries, or both and is used under license therefrom.

Intel, Intel logo, Intel Inside, Intel Inside logo, Intel Centrino, Intel Centrino logo, Celeron, Intel Xeon, Intel SpeedStep, Itanium, and Pentium are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States and other countries.

IT Infrastructure Library is a registered trademark of the Central Computer and Telecommunications Agency which is now part of the Office of Government Commerce.

ITIL is a registered trademark, and a registered community trademark of the Office of Government Commerce, and is registered in the U.S. Patent and Trademark Office.

Java and all Java based trademarks and logos are trademarks or registered trademarks of Oracle and/or its affiliates.

Linear Tape-Open, LTO, the LTO Logo, Ultrium, and the Ultrium logo are trademarks of HP, IBM Corp. and Quantum in the U.S. and

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

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

OpenStack is a trademark of OpenStack LLC. The OpenStack trademark policy is available on the [OpenStack website](#).

TEALEAF is a registered trademark of Tealeaf, an IBM Company.

Windows Server and the Windows logo are trademarks of the Microsoft group of countries.

Worklight is a trademark or registered trademark of Worklight, an IBM Company.

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

* Other product and service names might be trademarks of IBM or other companies.

Notes:

Performance is in Internal Throughput Rate (ITR) ratio based on measurements and projections using standard IBM benchmarks in a controlled environment. The actual throughput that any user will experience will vary depending upon considerations such as the amount of multiprogramming in the user's job stream, the I/O configuration, the storage configuration, and the workload processed. Therefore, no assurance can be given that an individual user will achieve throughput improvements equivalent to the performance ratios stated here.

IBM hardware products are manufactured from new parts, or new and serviceable used parts. Regardless, our warranty terms apply.

All customer examples cited or described in this presentation are presented as illustrations of the manner in which some customers have used IBM products and the results they may have achieved. Actual environmental costs and performance characteristics will vary depending on individual customer configurations and conditions.

This publication was produced in the United States. IBM may not offer the products, services or features discussed in this document in other countries, and the information may be subject to change without notice. Consult your local IBM business contact for information on the product or services available in your area.

All statements regarding IBM's future direction and intent are subject to change or withdrawal without notice, and represent goals and objectives only.

Information about non-IBM products is obtained from the manufacturers of those products or their published announcements. IBM has not tested those products and cannot confirm the 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.

Prices subject to change without notice. Contact your IBM representative or Business Partner for the most current pricing in your geography.

This information provides only general descriptions of the types and portions of workloads that are eligible for execution on Specialty Engines (e.g. zIIPs, zAAPs, and IFLs) ("SEs"). IBM authorizes customers to use IBM SE only to execute the processing of Eligible Workloads of specific Programs expressly authorized by IBM as specified in the "Authorized Use Table for IBM Machines" provided at www.ibm.com/systems/support/machine_warranties/machine_code/aut.html ("AUT"). No other workload processing is authorized for execution on an SE. IBM offers SE at a lower price than General Processors/Central Processors because customers are authorized to use SEs only to process certain types and/or amounts of workloads as specified by IBM in the AUT.

Agenda

- ❑ V2R2 Performance Enhancements
- ❑ V2R1 Performance Enhancements
- ❑ Optimizing inbound communications using OSA-Express
- ❑ Optimizing outbound communications using OSA-Express
- ❑ z/OS Communications Server Performance Summaries
- ❑ Appendixes



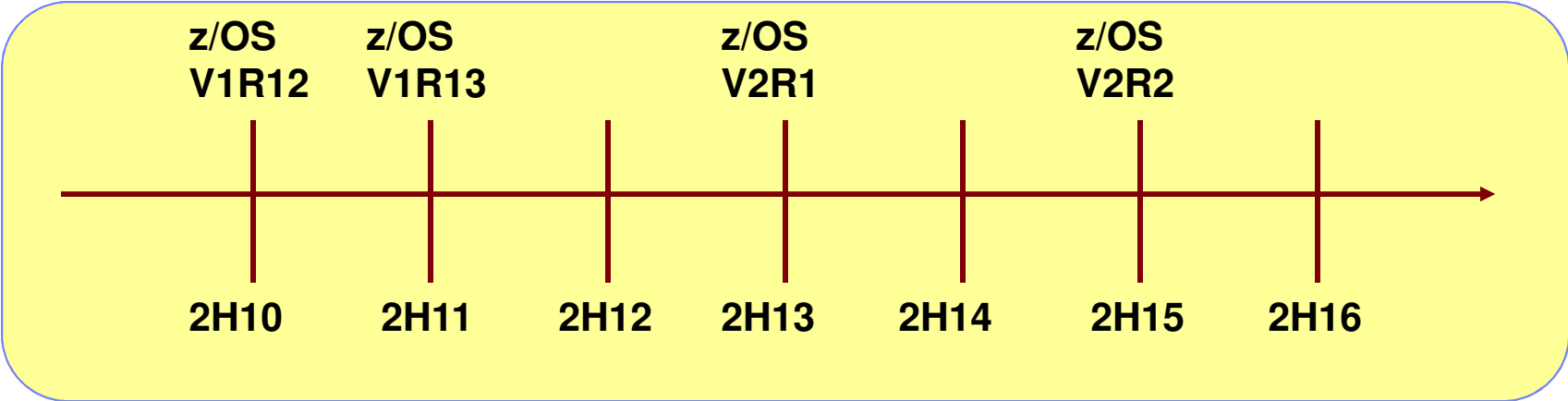
Disclaimer: All statements regarding IBM future direction or intent, including current product plans, are subject to change or withdrawal without notice and represent goals and objectives only. All information is provided for informational purposes only, on an “as is” basis, without warranty of any kind.

V2R2 Performance Enhancements – Tomorrowland

z/OS V2R2 Communications Server disclaimer

- ⑩ Plans for the z/OS Communications Server are subject to change prior to general availability.
- ⑩ Information provided in this presentation may not reflect what is actually shipped by z/OS Communications Server.
- ⑩ This presentation includes an early overview of selected future z/OS Communications Server performance enhancements.
- ⑩ Full set of results available:
 - ⑩ Future Share sessions
 - ⑩ z/OS Communications Server release summary report:
 - ⑩ <http://www-01.ibm.com/support/docview.wss?rs=852&uid=swg27005524>

Plans may change before GA of z/OS V2R2



Statements regarding IBM future direction and intent are subject to change or withdrawal, and represent goals and objectives only.

New performance and scalability functions

- New/improved TCP/IP autonomies
 - Avoid DELAYACK timer processing
 - AUTODELAYACK option
 - Eliminate occasional delays
 - Dynamic Right Sizing (DRS) improvement
 - Remain enabled
 - Unlimited runway to enablement
 - Outbound “Dynamic Right Sizing (ORS)”
 - Allow send buffer to grow when sending streaming (bulk) data
 - Improved response to lost vs. out-of-order packets
 - Don’t reduce slow start threshold for out-of-order
 - VIPAROUTE fragmentation avoidance
 - GLOBALCONFIG parameter – AUTOADJUSTMSS



New performance and scalability functions

- Improved SMC-R sending queued data algorithm
 - Send up to three queued writes in one interrupt
 - Benefit streaming/bulk type workloads
- Default to use 1MB (largest) RMBE for streaming/bulk data receive side
 - FTP uses 180K receive buffer
- SMC-R virtualization – LPARs sharing a RoCE Express feature
- Enhanced Enterprise Extender (EE) scalability
 - Large number EE connections (thousands)
 - Improved caching – improved latency and cpu
- For more details on these functions please refer to Share session:
 - [z/OS V2R2 Communications Server Technical Update \(Parts 1 & 2\)](#)

TCP Delayack Processing and Nagle's Algorithm "Catch-22"

```

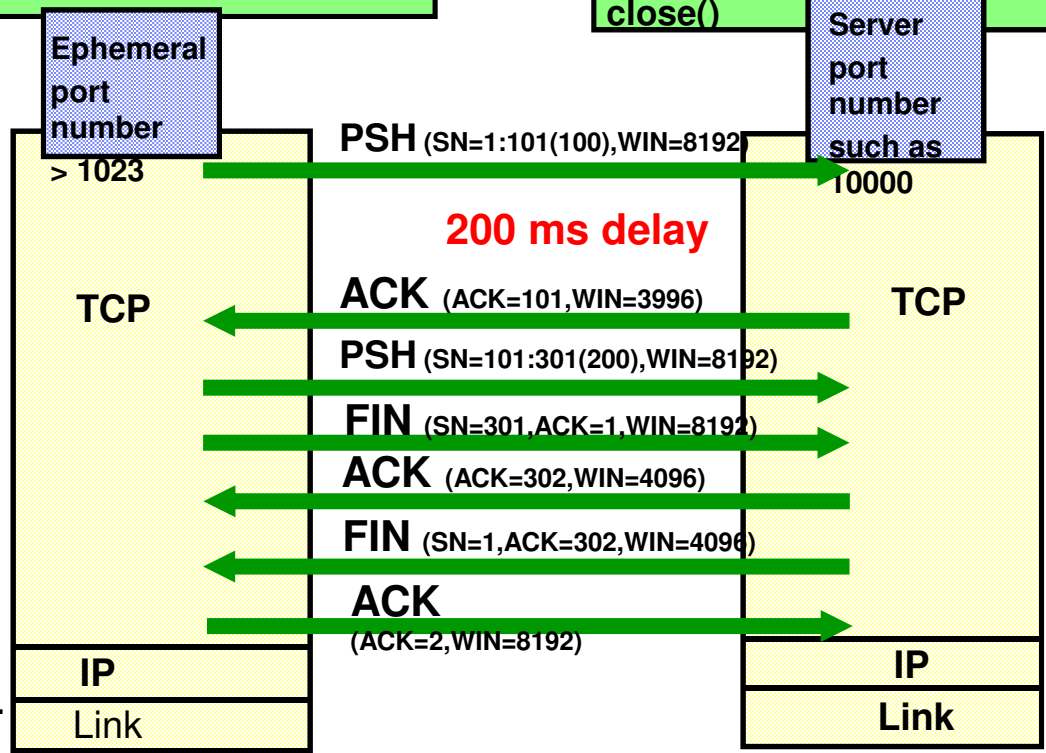
    socket()
    connect()
    send(100 bytes of data)
    send(100 bytes of data)
    send(100 bytes of data)
    close()

    socket(), bind(),
    listen(), accept()
    Loop while data to be received
    recv(4K of data)
    process data
    end Loop
    close()
  
```

Nagle's Algorithm (RFC896)

As long as there is a sent packet for which the sender has received no acknowledgment, the sender should keep buffering its output until it has a full packet's worth of output, so that output can be sent all at once.

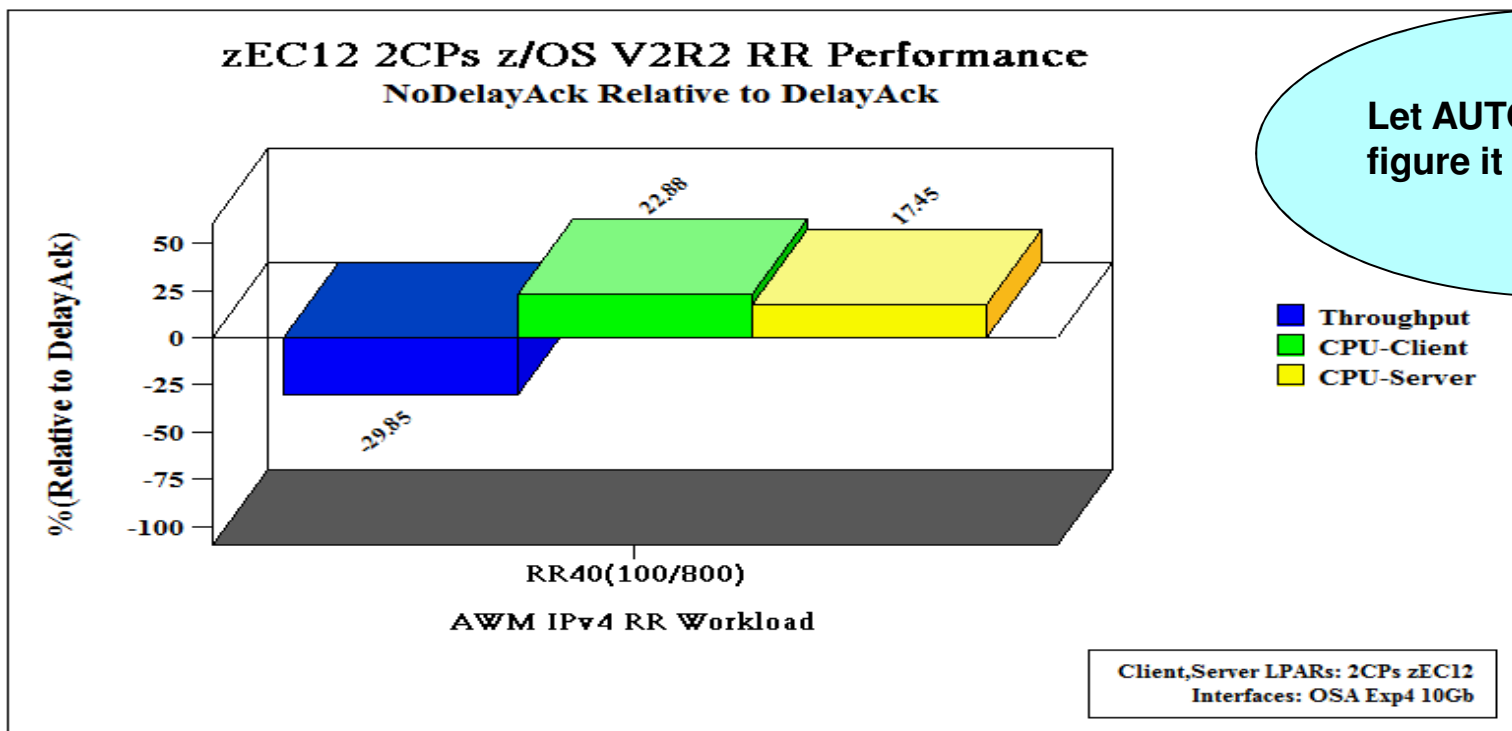
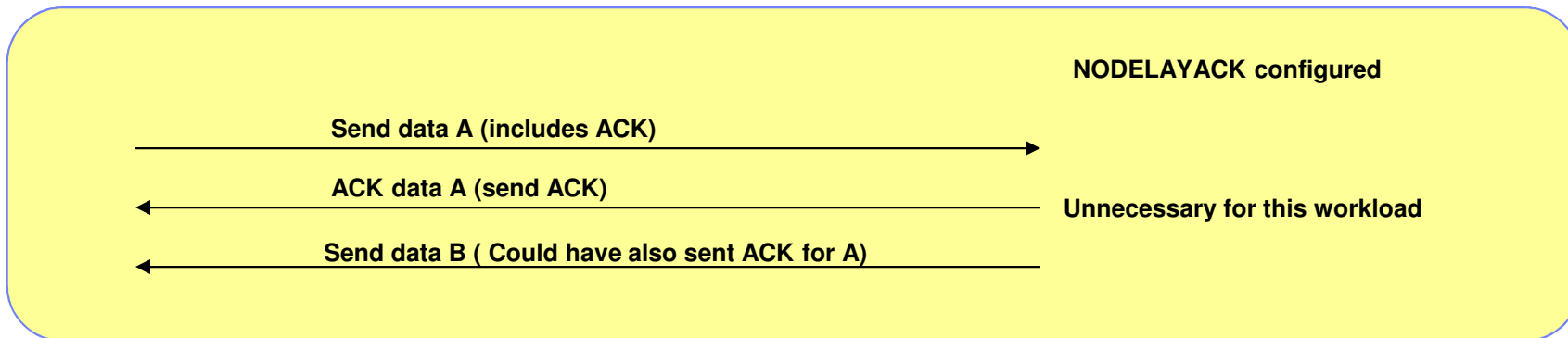
Goal: Minimize congestion



Delayed Acks (RFC1122)

A host that is receiving a stream of TCP data segments can increase efficiency in both the Internet and the hosts by sending fewer than one ACK segment per data segment received; this is known as a "delayed ACK"

NODELAYACK issue



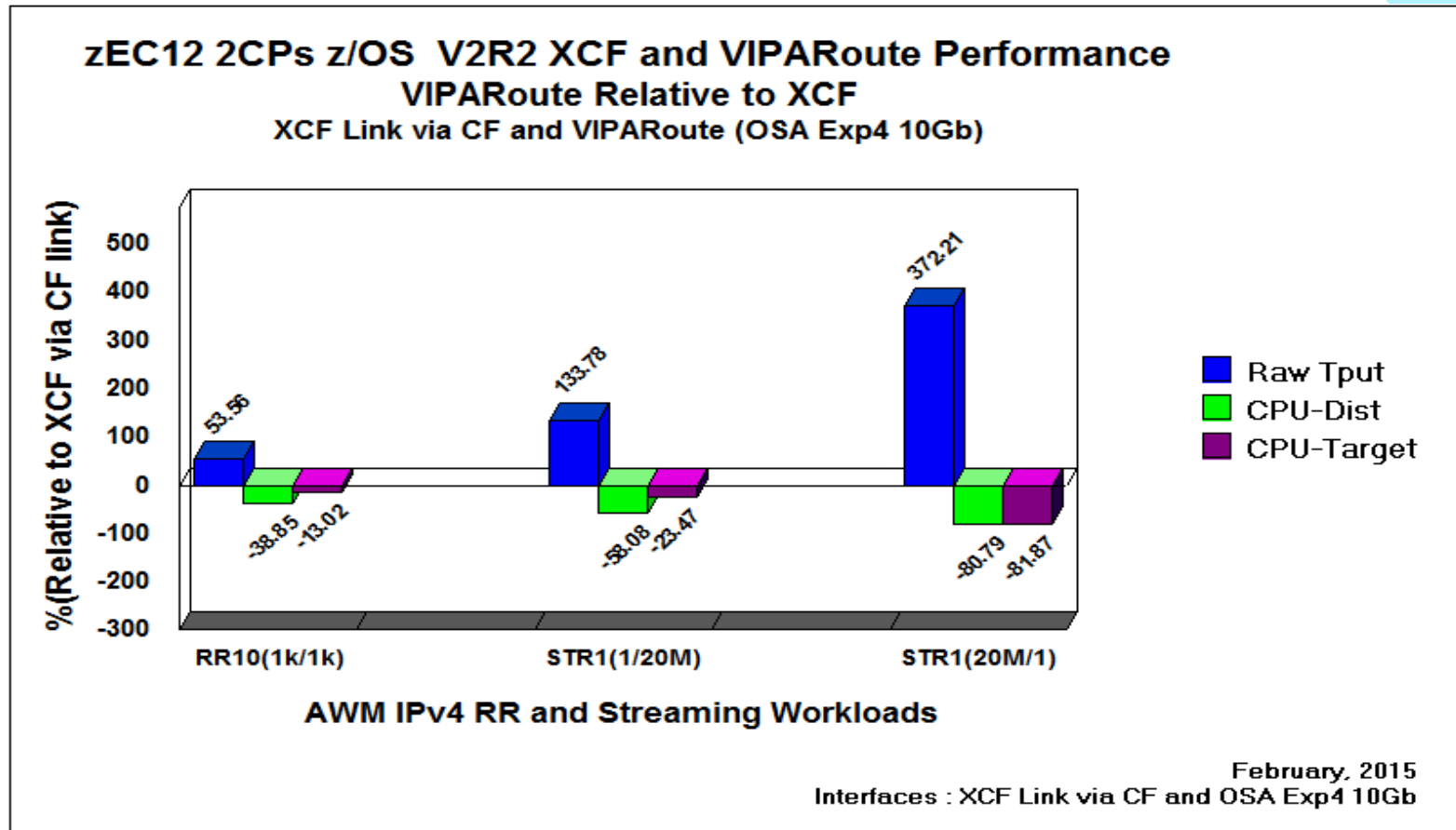
Let AUTODELAYACK figure it out for you!

AUTOADJUSTMSS for VIPARROUTE: New GLOBALCONFIG parameter



- Automatically reduce the MSS (Maximum Segment Size) of a distributed connection by the length of the GRE header
 - Eliminate fragmentation
- Why you should be using VIPARROUTE for Sysplex Distributor Workloads:

Improved in V2R2!



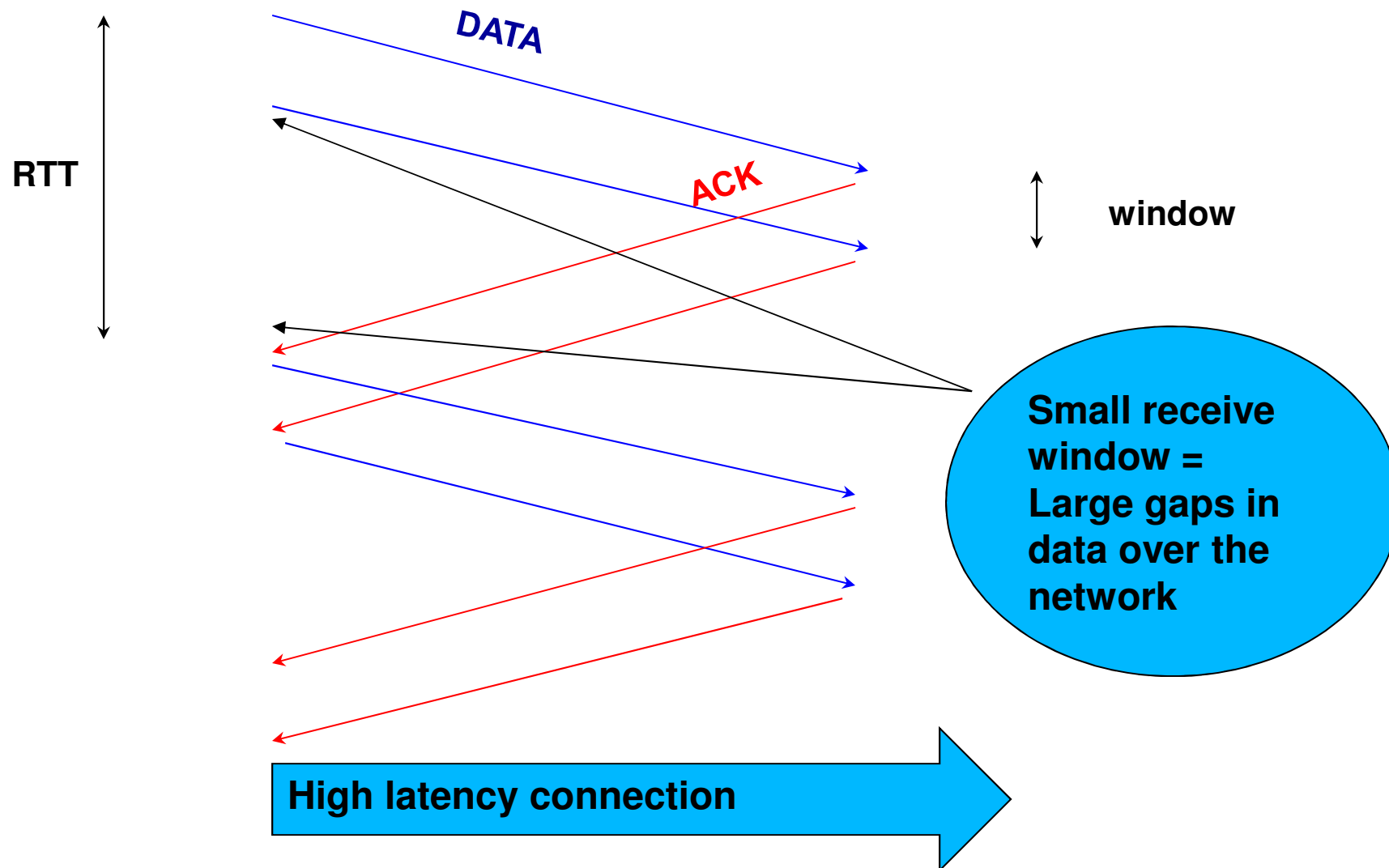
Dynamic Right Sizing

- Keep the pipe full and prevent sender from being constrained by the advertised window
- Improves performance for inbound streaming workloads over high latency networks with large bandwidth-delay product
- Function enabled automatically (no configuration)
- Has proven very helpful in several installations
- Stack 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
 - Allows window size to grow as high as 2M
- Allow enablement for connections that don't start as streamers
- Don't disable if not storage-constrained

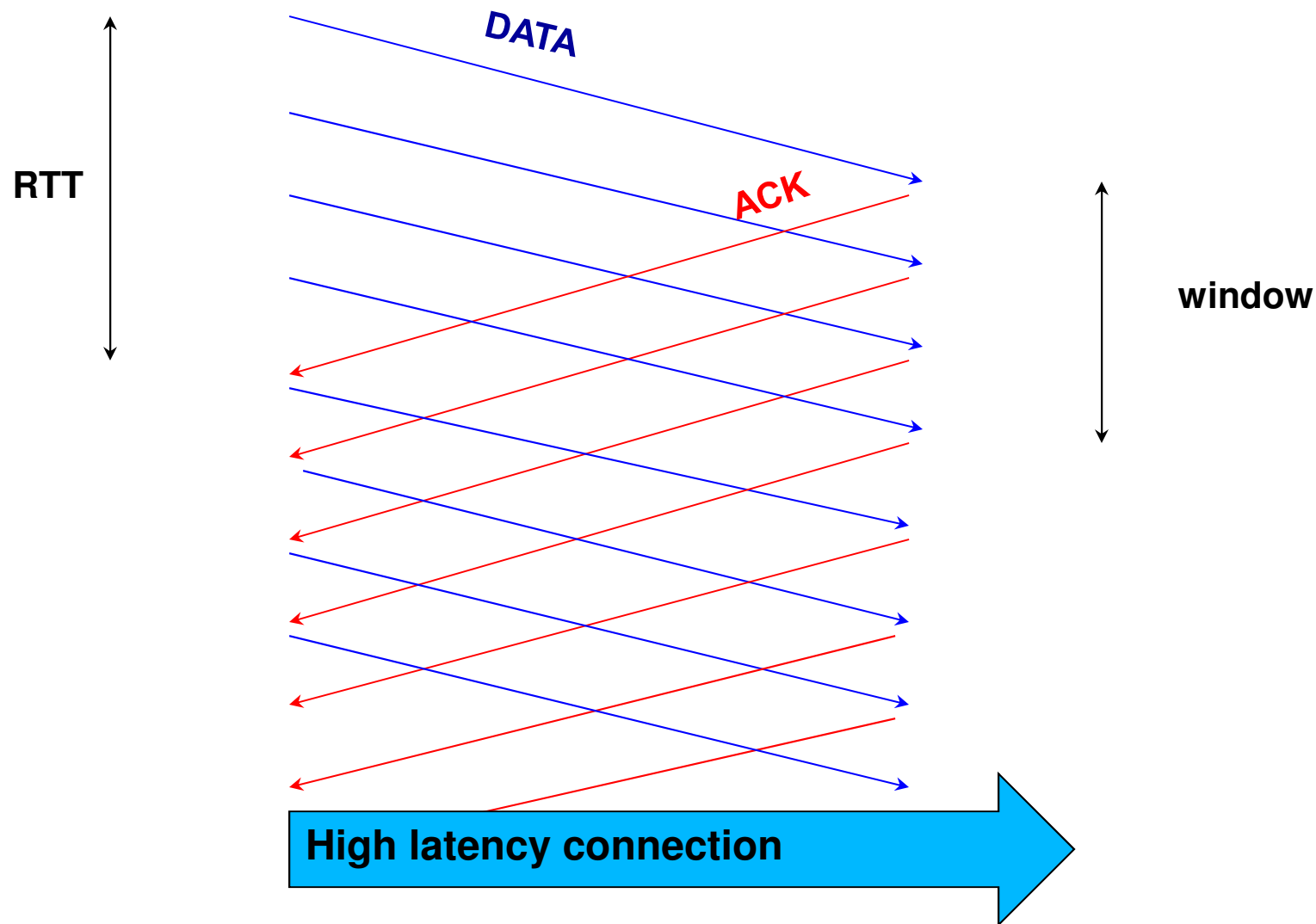


**Improved
in V2R2!**

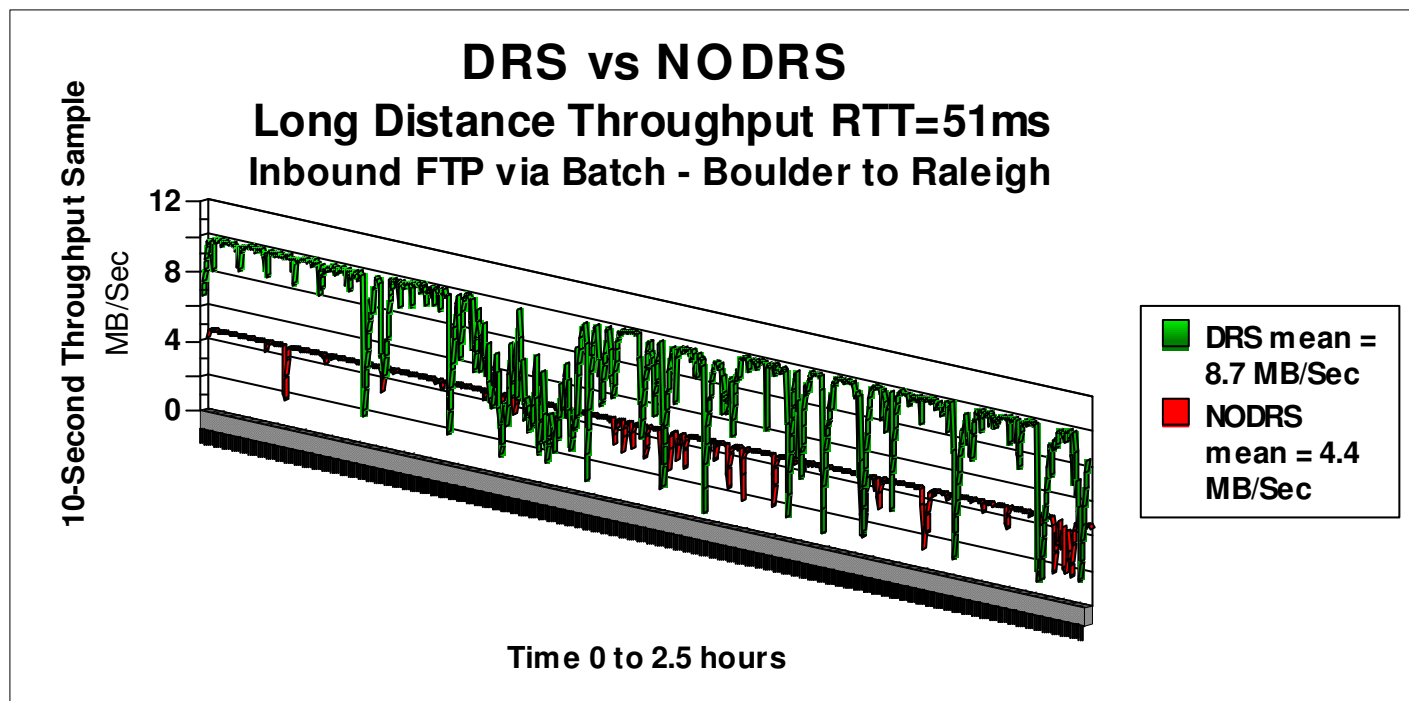
V2R2 Dynamic Right Sizing improvement



V2R2 Dynamic Right Sizing improvement



Dynamic Right Sizing



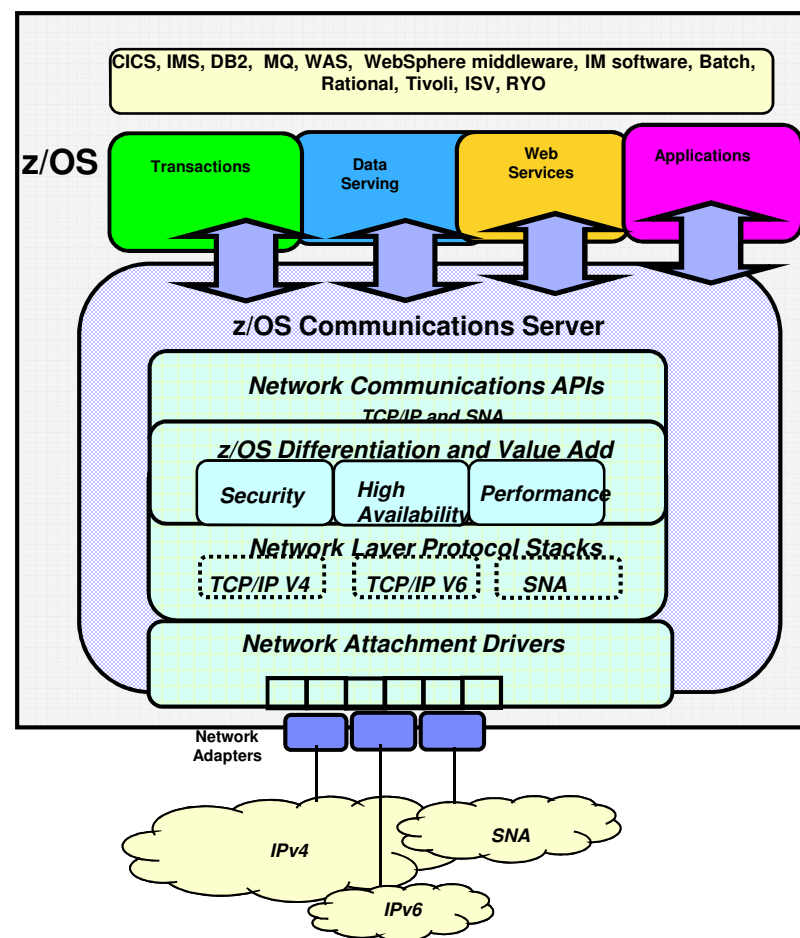
Over an extended 2.5 hour experiment, the DRS enabled receiver averaged double the throughput compared to no DRS.

This experiment repeatedly transferred a 2.8 GB file, and DRS never disabled over the 2.5 hour period.

If DRS had disabled at any point, the numbers would have been roughly identical during the period of DRS disablement.

64-Bit enablement of the TCP/IP stack and strategic DLCs

- TCP/IP has supported 64-bit applications since 64-bit support was introduced on the platform
 - But the mainline path has been 31-bit with extensive use of AR mode
- As systems become more powerful, customers have increased the workloads on the systems which in turn increases the storage demands placed on the systems.
- The storage in 31-bit addressing mode (below the bar) has been of special concern. Over the past several releases work has started to move storage that used to be obtained below the bar to 64-bit addressing mode (above the bar).
- Some of these changes, such as V1R13's move of the CTRACE and VIT above the bar, were visible to customers, while others were just changes in internal "plumbing".
- The next step is a large one: To move most of the remaining storage above the bar without incurring an unacceptable overhead in switching between AMODE(31) and AMODE(64) requires the complete 64-bit enablement of the TCP/IP stack and strategic device drivers (DLCs).



64 bit storage savings – Telnet run, large number TCP connections

Address Space and Storage Type	31 - Bit V2R1 (KB)	64 - Bit V2R2 (KB)	% change from V2R1
Telnet ECSA	1,575	145	-91
TCP/IP ECSA	9,188	6,593	-28
TCP/IP Private	275,338	43,332	-84

Removing workload growth constraints from below-the-bar Private and Common storage:

ECSA savings:

- **Dynamic storage/data buffers moved**

Private savings:

- **Key connection control block moved**

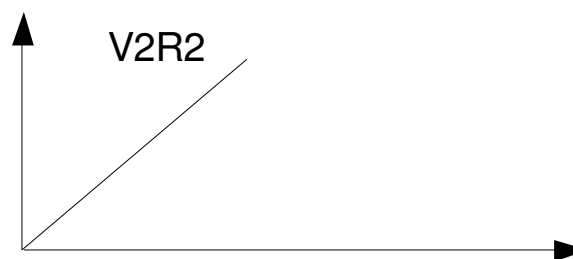
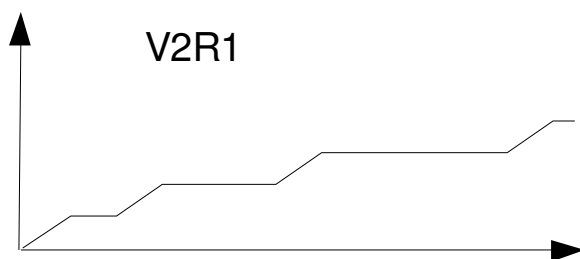
* Based on benchmarks of modeled z/OS TCP sockets based workloads with telnet traffic patterns. The benefits any user will experience will vary.

Enhanced IKED scalability (IPSEC)

- Current IKED design relies heavily on a single thread to perform the bulk of the work
- When a very large number (multiple thousands) of remote IKE peers simultaneously initiate negotiations with a single z/OS IKED, the z/OS daemon struggles to keep up with the load
- In V2R2, z/OS IKED is modified to handle heavy bursts of negotiations from very large numbers (multiple thousands) of IKE peers
 - A new thread pool is added to parallelize handling of IKE messages from different peers
 - Logic is added to minimize the amount of effort IKED spends processing retransmitted messages from peers
 - Transparent to the vast majority of current IKED users
 - Improvement will be most noticeable to users with very large numbers (multiple thousands) of IKE peers

Enhanced IKED scalability – Early results

- Preliminary z/OS V2R2 performance results show significant performance improvements in establishing SAs when a large number of concurrent client requests arrive in a small interval of time
 - IKE V1 (more messages exchanged):
 - Up to 6.8X improvement in throughput (rate of SA activations) and up to 75% reduction in CPU cost *
 - IKE V2:
 - Up to 3.8X improvement in throughput (rate of SA activations) and up to 57% reduction in CPU cost *



*Note: The performance measurements were collected in IBM internal tests using a dedicated system environment. 4,200 clients simulated using 4 Linux for System z images running under zVM. IKE v1 benchmarks performed with PSK. IKE v2 benchmarks performed with RSA. The results obtained in other configurations or operating system environments may vary significantly depending upon environments used. Therefore, no assurance can be given, and there is no guarantee that an individual user will achieve performance or throughput improvements equivalent to the results stated here.

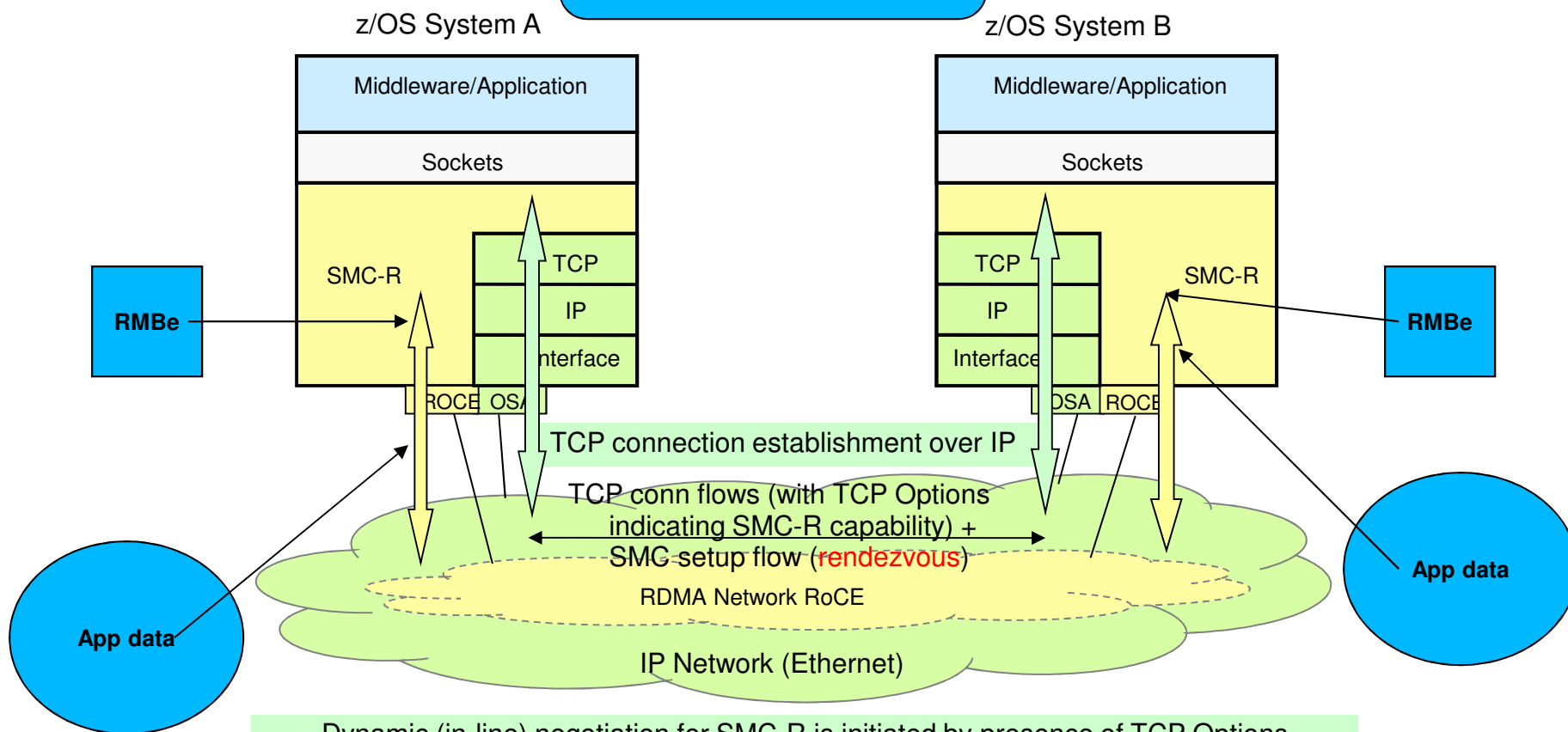
V2R1 Performance Enhancements

Shared Memory Communications – Remote (SMC-R)

SMC-R Background

Both TCP and SMC-R “connections” remain active

V2R1



SMC-R - RDMA

V2R1

- Key attributes of RDMA
 - Enables a host to read or write directly from/to a remote host's memory ***without*** involving the remote host's CPU
 - By registering specific memory for RDMA partner use
 - Avoids TCP ACK processing
 - Avoids TCP “packetizing” the data
 - **Interrupts still required for notification (i.e. CPU cycles are not completely eliminated)**
 - Reduced networking stack overhead by using streamlined, low level, RDMA interfaces
 - Key requirements:
 - A reliable “lossless” network fabric (LAN for layer 2 data center network distance)
 - An RDMA capable NIC (RNIC) and RDMA capable switch

SMC-R - Solution

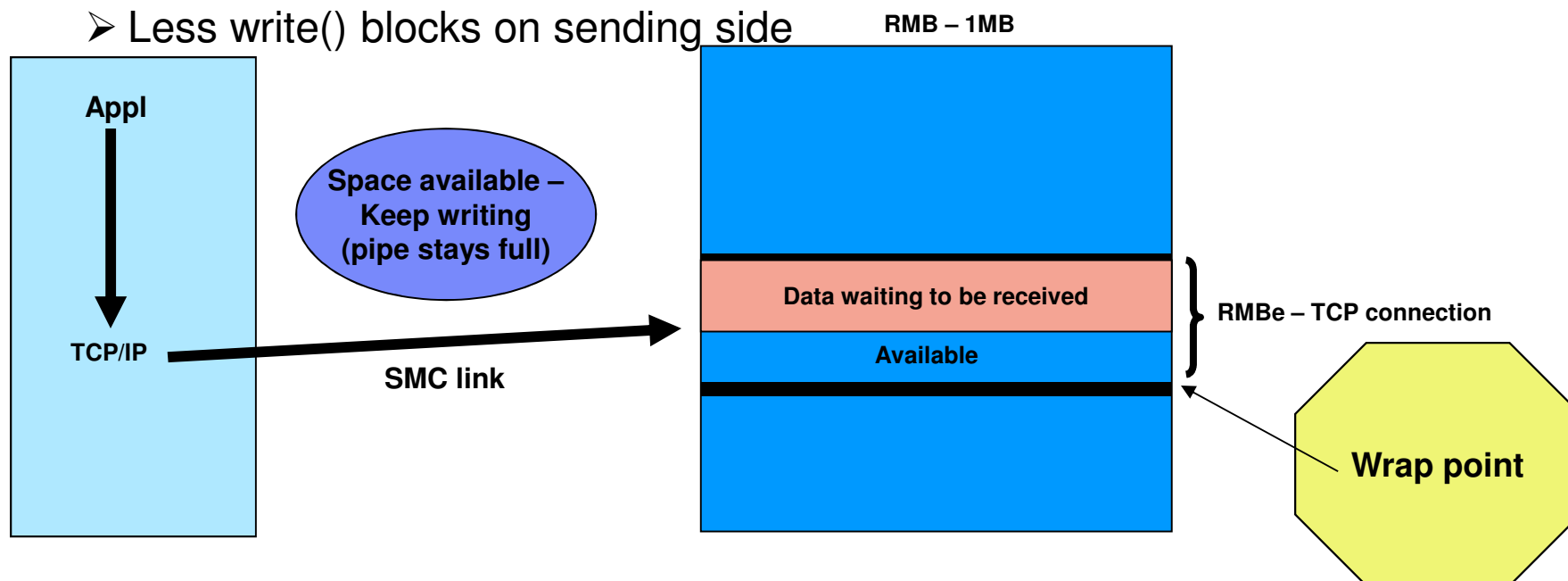
V2R1

- Shared Memory Communications over RDMA (SMC-R) is a protocol that allows *TCP sockets* applications to transparently exploit RDMA (RoCE)
- SMC-R is a “hybrid” solution that:
 - Uses TCP connection (3-way handshake) to establish SMC-R connection
 - Each TCP end point exchanges TCP options that indicate whether it supports the SMC-R protocol
 - SMC-R “rendezvous” (RDMA attributes) information is then exchanged within the TCP data stream (similar to SSL handshake)
 - Socket application data is exchanged via RDMA (write operations)
 - TCP connection remains active (controls SMC-R connection)
 - This model preserves many critical existing operational and network management features of TCP/IP

SMC-R – Role of the RMBe (buffer size)

- The RMBe is a slot in the RMB buffer for a specific TCP connection
 - Based on TCPRCVBufsize – NOT equal to
 - Can be controlled by application using setsockopt() SO_RCVBUF
 - 5 sizes – 32K, 64K, 128K, 256K and 1024K (1MB)
 - Will use 1MB for TCPRCVBufsize > 128K
 - Depending on the workload, a larger RMBe can improve performance
 - Streaming (bulk) workloads
 - Less wrapping of the RMBe = less RDMA writes
 - Less frequent “acknowledgement” interrupts to sending side
 - Less write() blocks on sending side

New for V2R2!



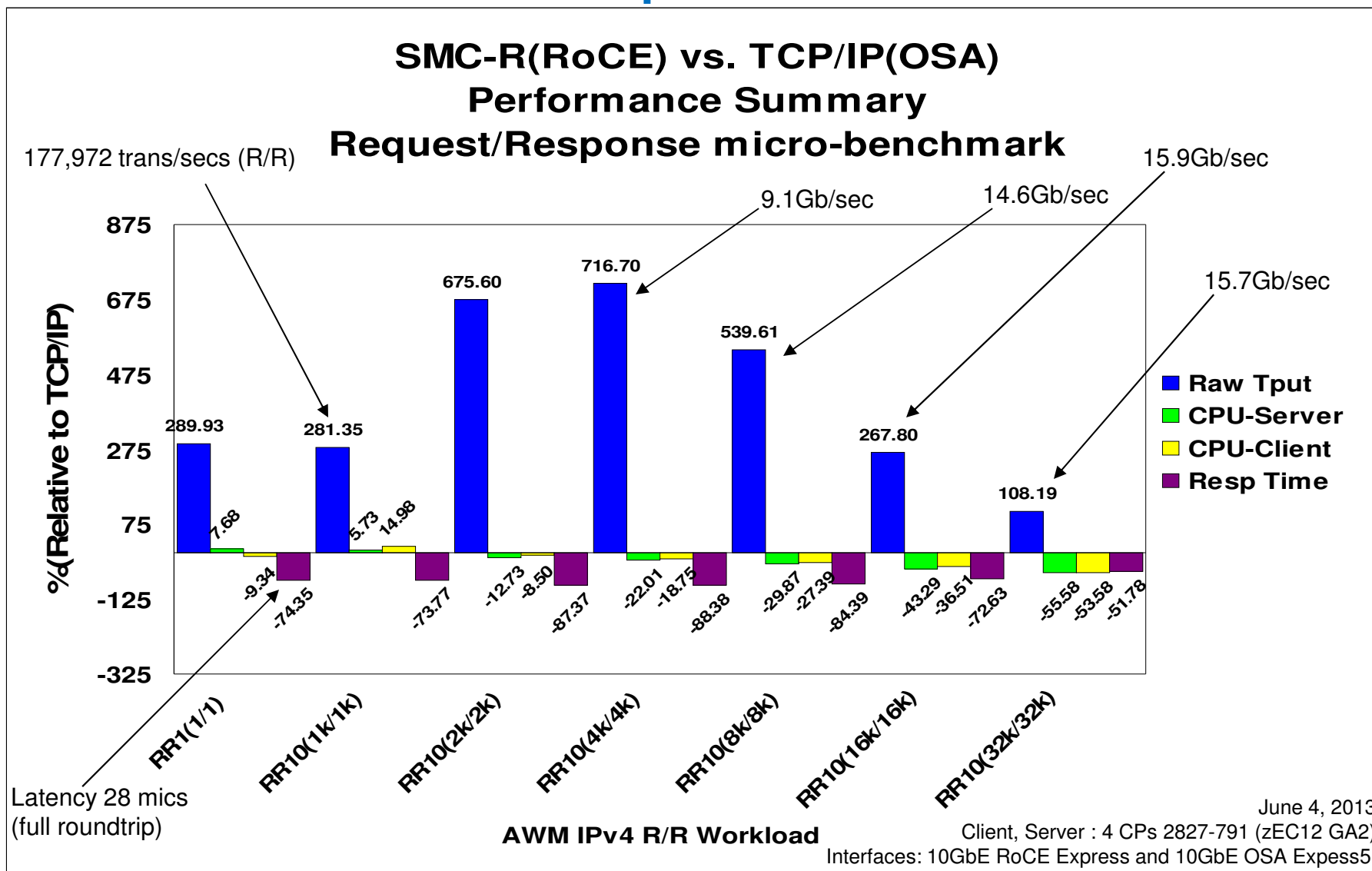
SMC-R – Micro benchmark performance results

V2R1

- Response time/Throughput and CPU improvements
- Workload:
 - Using AWM (Application Workload Modeler) to model “socket to socket” performance using SMC-R
 - AWM very lightweight - contains no application/business logic
 - Stresses and measures the networking infrastructure
 - Real workload benefits **will be smaller** than the improvements seen in AWM benchmarks!
 - MTU: RoCE (1K and 2K) OSA (1500 and 8000)
 - SEGMENTATIONOFFLOAD enabled for some of the TCP/IP streaming runs
 - RR1(1/1): Single interactive session with 1 byte request and 1 byte reply
 - RR10: 10 concurrent connections with various message sizes
 - STR1(1/20M): Single Streaming session with 1 byte request (Client) and 20,000,000 bytes reply (Server)
 - Used large RMBs – 1MB

SMC-R – Micro benchmark performance results

V2R1



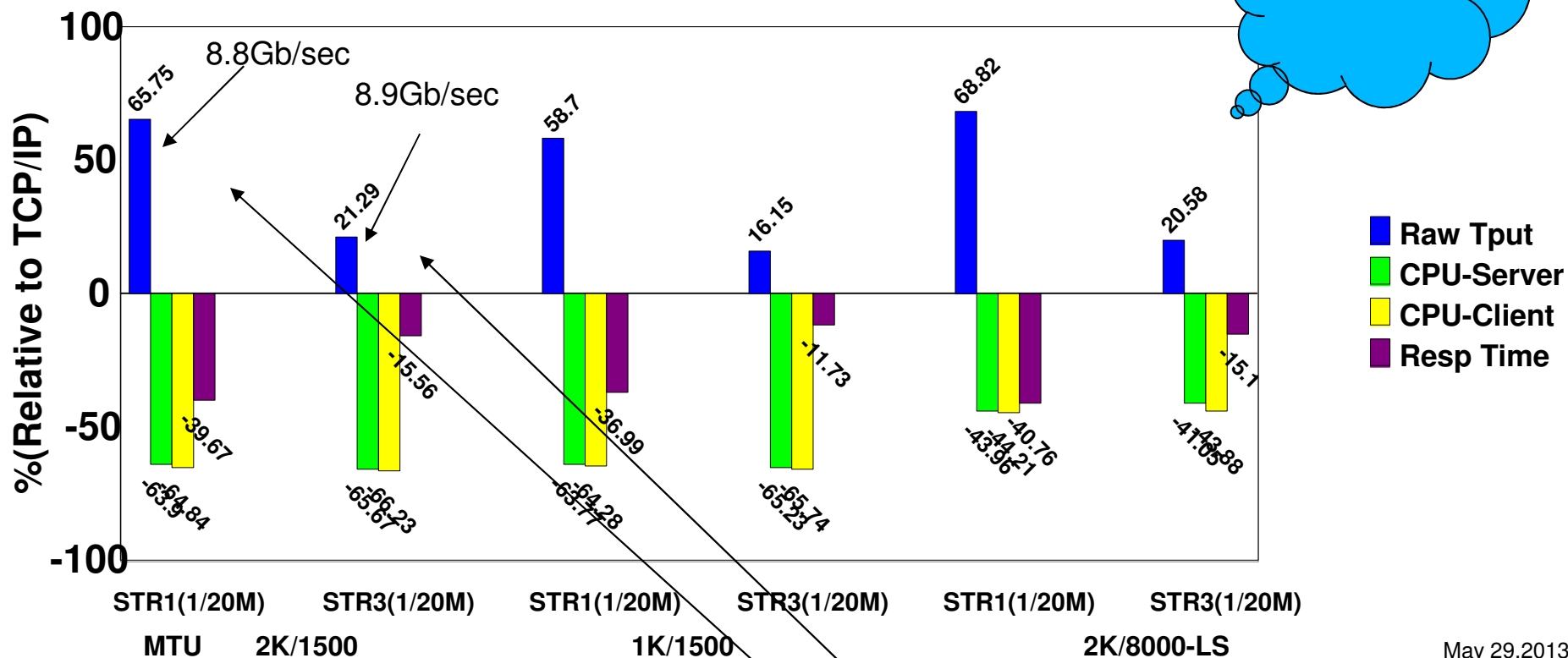
Significant Latency reduction across all data sizes (52-88%)
Reduced CPU cost as payload increases (up to 56% CPU savings)
Impressive throughput gains across all data sizes (Up to +717%)

Note: vs typical OSA customer configuration
 MTU (1500), Large Send disabled
 RoCE MTU: 1K

SMC-R – Micro benchmark performance results

V2R1

z/OS V2R1 SMC-R vs TCP/IP
Streaming Data Performance Summary (AWM)



1MB RMBs

Saturation reached

May 29, 2013
Client, Server: 2827-791 2CPs
Interfaces: 10GbE RoCE Express and 10GbE

- Notes:
- Significant throughput benefits and CPU reduction benefits
 - Up to 69% throughput improvement
 - Up to 66% reduction in CPU costs
 - 2K RoCE MTU does yield throughput advantages
 - LS – Large Send enabled (Segmentation offload)

SMC-R – Micro benchmark performance results

V2R1

- Summary –
 - Network latency for z/OS TCP/IP based OLTP (request/response) workloads reduced by up to 80%*
 - Networking related CPU consumption reduction for z/OS TCP/IP based OLTP (request/response) workloads increases as payload size increases
 - Networking related CPU consumption for z/OS TCP/IP based workloads with streaming data patterns reduced by up to 60% with a network throughput increase of up to 60%**
 - CPU consumption can be further optimized by using larger RMBe sizes
 - Less data consumed processing
 - Less data wrapping
 - Less data queuing

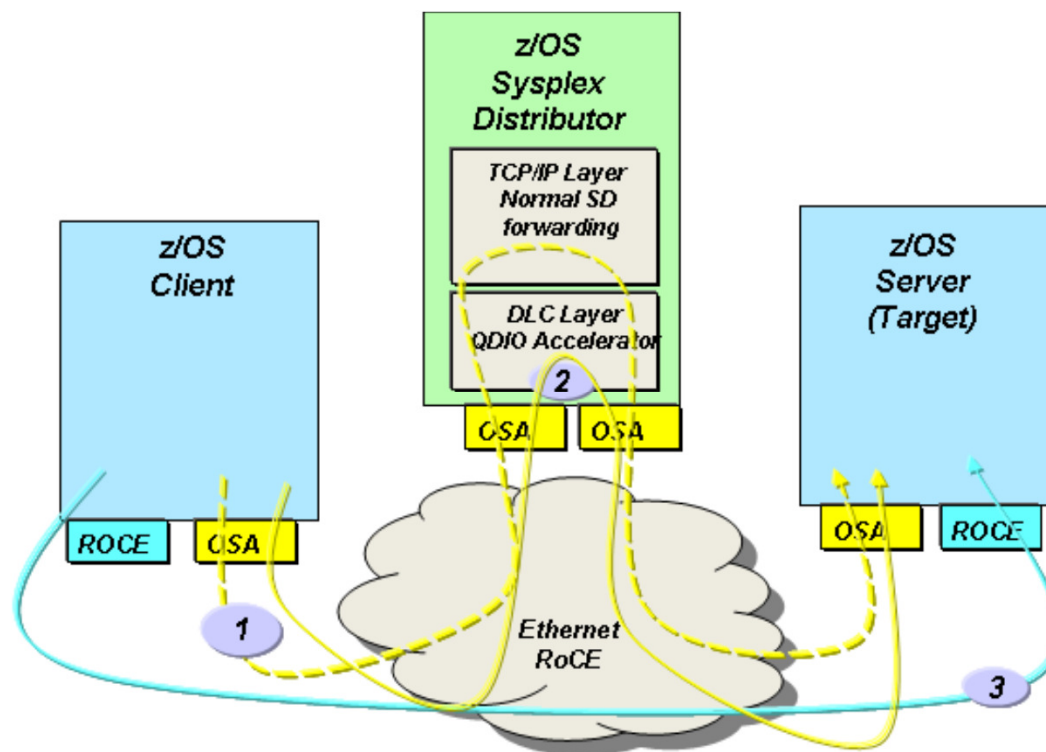
* Based on benchmarks of modeled z/OS TCP sockets based workloads with request/response traffic patterns using SMC-R vs. TCP/IP. The actual response times and CPU savings any user will experience will vary.

** Based on benchmarks of modeled z/OS TCP sockets based workloads with streaming data patterns using SMC-R vs. TCP/IP. The benefits any user will experience will vary

SMC-R – Sysplex Distributor performance results

V2R1

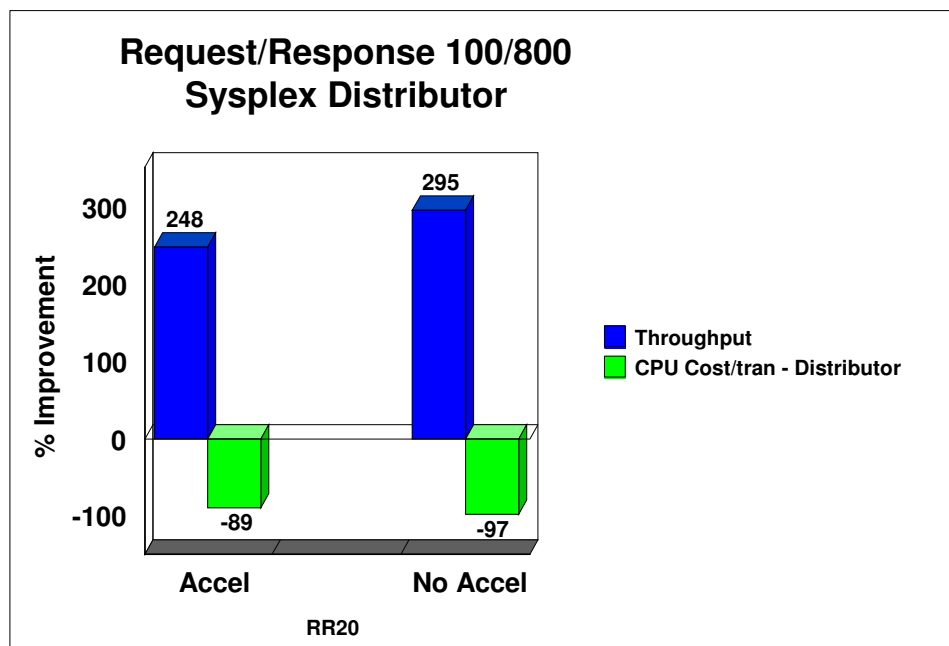
With SMC-R the distributing stack is bypassed for inbound data. Connection setup and SMC-R rendezvous packets will be the only inbound traffic going through the distributing stack. Remember that all outbound traffic bypasses the distributing stack for all scenarios.



- Line 1 - TCP/IP distributed connections without QDIO Accelerator
- Line 2 - TCP/IP distributed connections utilizing QDIO Accelerator
- Line 3 - SMC-R distributed connections

SMC-R – Sysplex Distributor performance results

V2R1



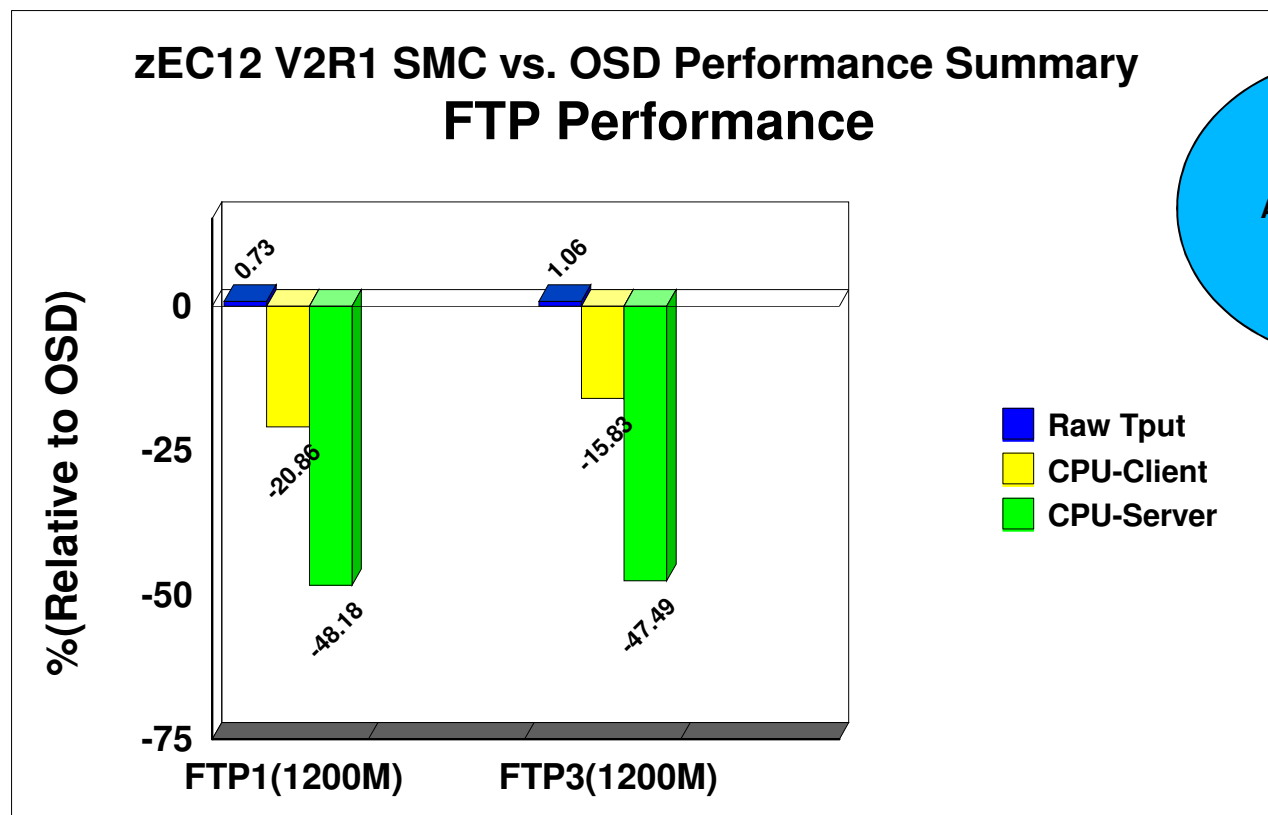
Results from Sysplex distributing Stack perspective

SMC-R removes virtually all CP processing on Distributing stack

250%+ throughput improvement

Workload – 20 simultaneous request/response connections sending 100 and receiving 800 bytes. Large data workloads would yield even bigger performance improvements.

SMC-R – FTP performance summary



- FTP binary PUTs to z/OS FTP server, 1 and 3 sessions, transferring 1200 MB data
- OSD – OSA Express4 10Gb interface
- Reading from and writing to DASD datasets – Limits throughput

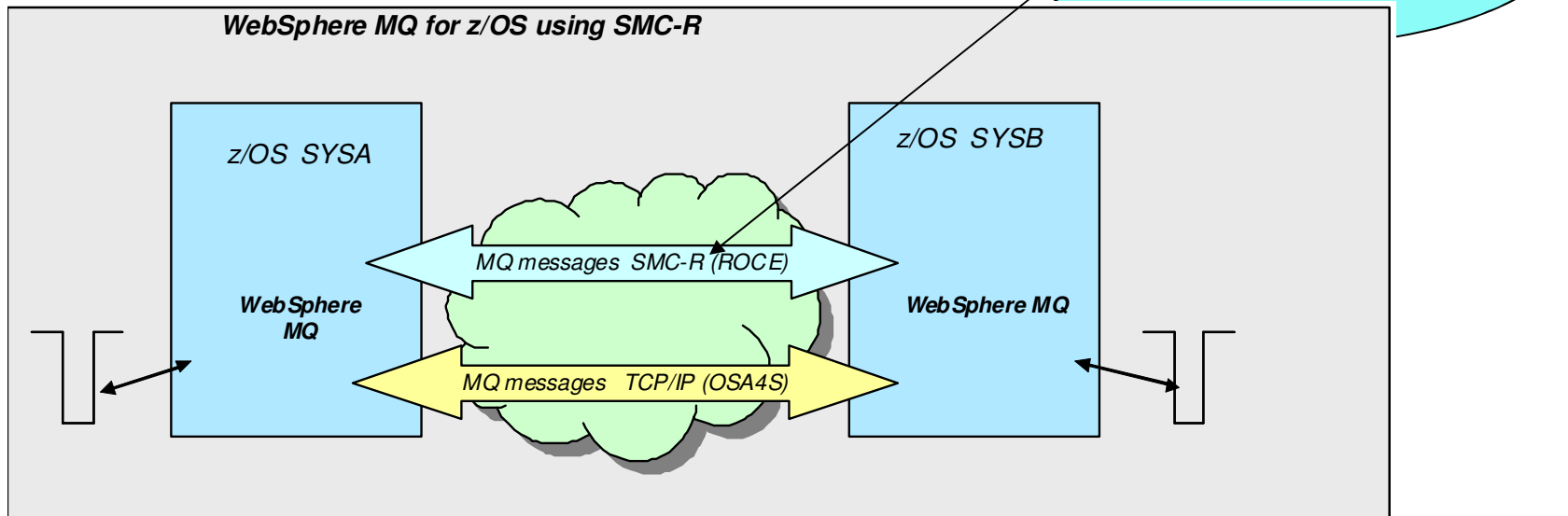
The performance measurements discussed in this document were collected using a dedicated system environment. The results obtained in other configurations or operating system environments may vary significantly depending upon environments used.

SMC-R - WebSphere MQ for z/OS performance improvement



V2R1

- Latency improvements



▪ WebSphere MQ for z/OS realizes **up to a 3x increase** in messages per second it can deliver across z/OS systems when using SMC-R vs standard TCP/IP for 64K messages over 1 channel *

*Based on internal IBM benchmarks using a modeled WebSphere MQ for z/OS workload driving non-persistent messages across z/OS systems in a request/response pattern. The benchmarks included various data sizes and number of channel pairs. The actual throughput and CPU savings users will experience may vary based on the user workload and configuration.

SMC-R - WebSphere MQ for z/OS performance improvement

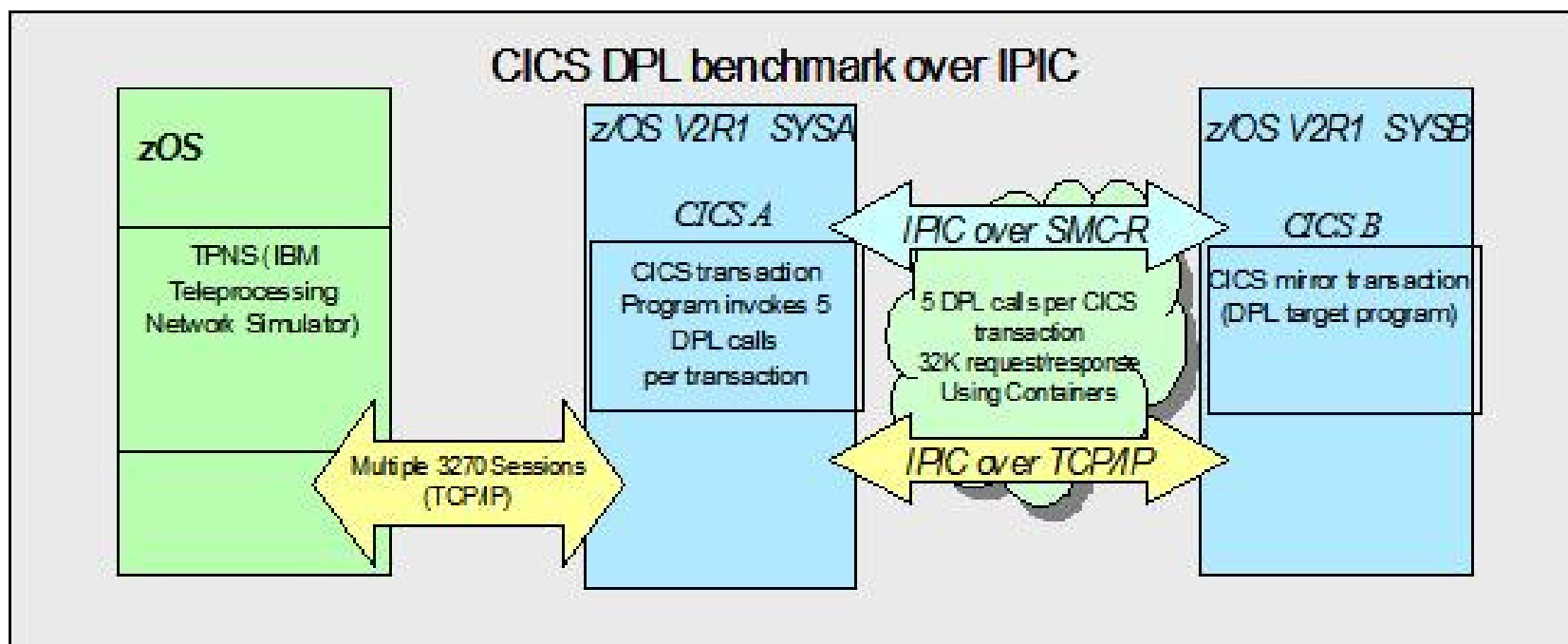


V2R1

- Latency improvements
- Workload
 - Measurements using WebSphere MQ V7.1.0
 - MQ between 2 LPARs on zEC12 machine (10 processors each)
 - Request/Response workload
 - On each LPAR, a queue manager was started and configured with 50 outbound sender channels and 50 inbound receiver channels, with default options for the channel definitions (100 TCP connections)
 - Each configuration was run with message sizes of 2KB, 32KB and 64KB where all messages were non-persistent
 - Results were consistent across all three message sizes

SMC-R – CICS performance improvement

V2R1



- Benchmarks run on z/OS V2R1 with zEC12 and 10GbE RoCE Express feature
 - Compared use of SMC-R (10GbE RoCE Express) vs standard TCP/IP (10GbE OSA Express4S) with CICS IPIC communications for DPL (Distributed Program Link) processing
 - **Up to 48% improvement in CICS transaction response time** as measured on CICS system issuing the DPL calls (CICS A)
 - **Up to 10% decrease in overall z/OS CPU consumption** on the CICS systems

SMC-R – CICS performance improvement

- Response time and CPU utilization improvements
- Workload - Each transaction
 - Makes 5 DPL (Distributed Program Link) requests over an IPIC connection
 - Sends 32K container on each request
 - Server program Receives the data and Send back 32K
 - Receives back a 32K container for each request

IPIC - IP Interconnectivity

- Introduced in CICS TS 3.2/TG 7.1
- TCP/IP based communications
- Alternative to LU6.2/SNA for Distributed program calls

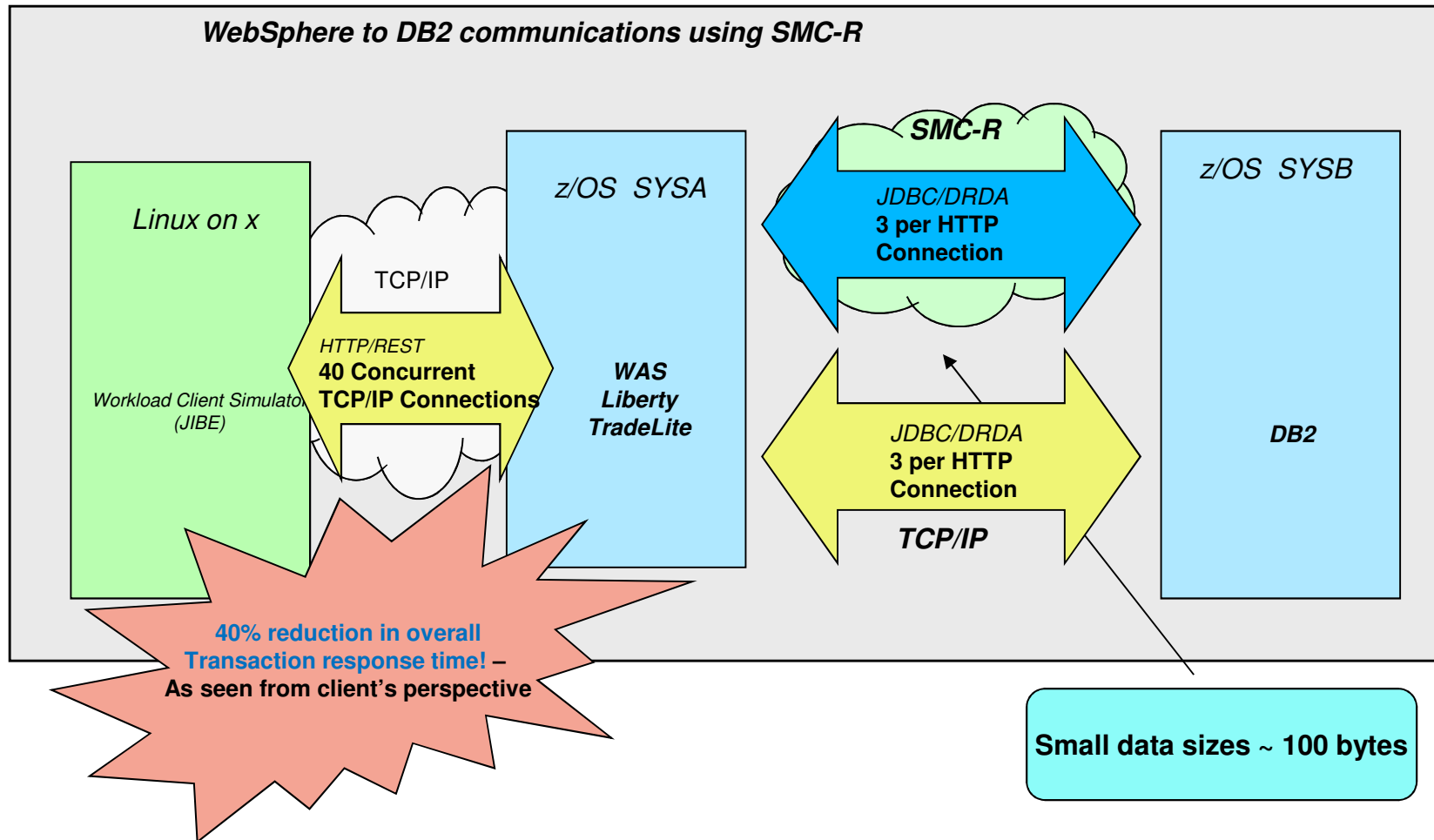
Note: Results based on internal IBM benchmarks using a modeled CICS workload driving a CICS transaction that performs 5 DPL calls to a CICS region on a remote z/OS system, using 32K input/output containers. Response times and CPU savings measured on z/OS system initiating the DPL calls. The actual response times and CPU savings any user will experience will vary.

SMC-R – Websphere to DB2 communications performance improvement



- Response time improvements

V2R1



Based on projections and measurements completed in a controlled environment. Results may vary by customer based on individual workload, configuration and software levels.

SMC-R and RoCE performance benchmarks at distance

- Initial statement of support for SMC-R and RoCE Express
 - 300 meters maximum distance from RoCE Express port to 10GbE switch port using OM3 fiber cable
 - 600 meters maximum when sharing the same switch across 2 RoCE Express features
 - Distance can be extended across multiple cascaded switches
 - All initial performance benchmarks focused on short distances (i.e. same site)
- Updated testing for RoCE and SMC-R over long distances
 - IBM System z™ Qualified Wavelength Division Multiplexer (WDM) products for Multi-site Sysplex and GDPS® solutions qualification testing updated to include RoCE and SMC-R. Several vendors already certified their DWDM solution for SMC-R and RoCE Express:

To monitor the latest products qualified refer to:

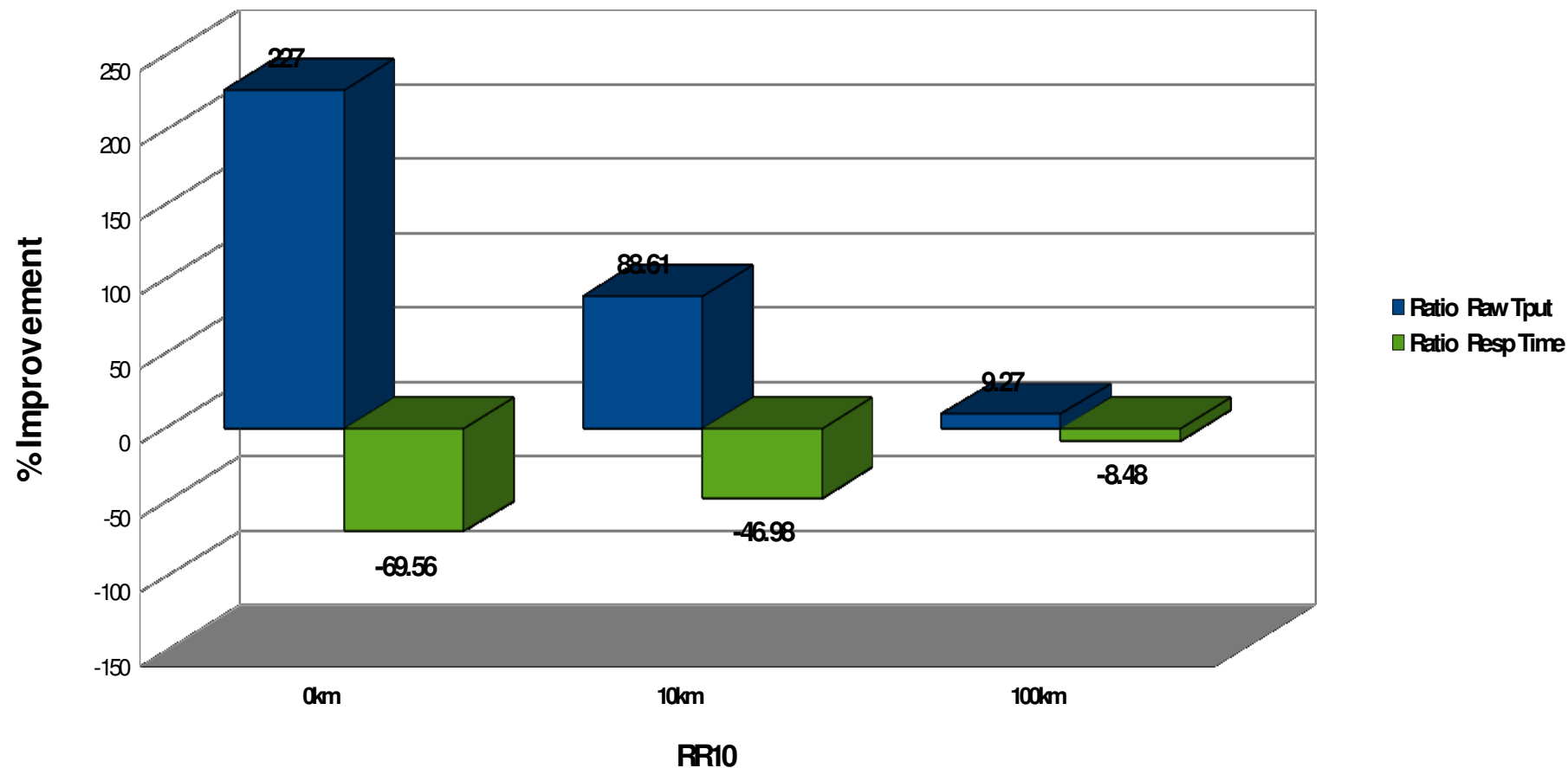
<https://www-304.ibm.com/servers/resourcelink/lib03020.nsf/pages/systemzQualifiedWdmProductsForGdpSolutions?OpenDocument>

- *So, how does SMC-R and RoCE perform at distance?*

SMC-R RoCE performance at distance - Request/Response Pattern (small data)

Request/Response-1KB/1KB

SMC-R vs. TCP/IP



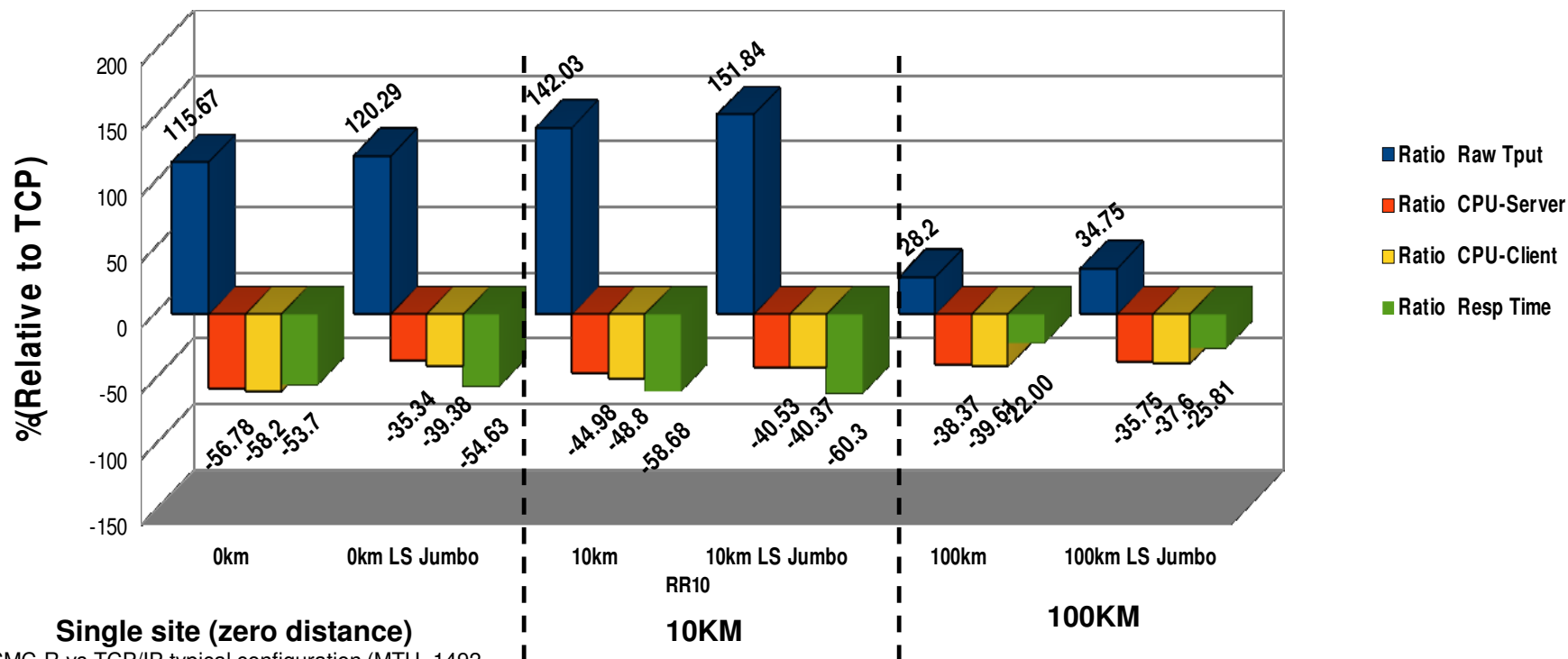
Notes:

- RR10(1K/1K): 10 persistent TCP connections simulating request/response data pattern, client sends 1KB request, server responds with 1KB
- **Substantial response time (i.e. latency) improvements at 10KM, benefits drop off at 100km**

SMC-R RoCE performance at distance – Request/Response Pattern

Request/Response -32KB/32KB

SMC-R vs. TCP/IP



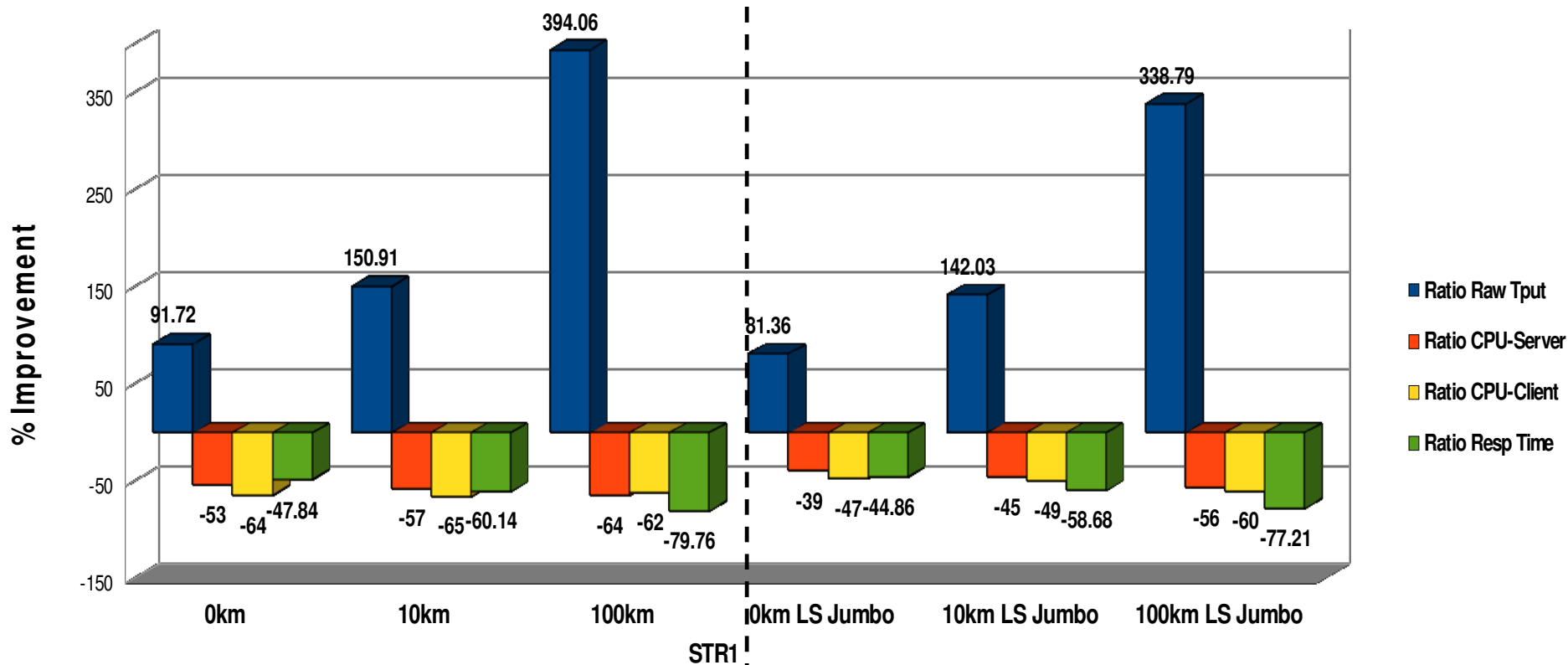
- Single site (zero distance)**
- 1) SMC-R vs TCP/IP typical configuration (MTU=1492, Large Send Disabled)
 - 2) SMC-R vs TCP/IP optimal configuration optimal TCP/IP configuration (MTU=8000, Large Send Enabled)

Typical TCP/IP configuration
MTU=1492, Large Send disabled

- Notes:
 - RR10(32K/32K): 10 persistent TCP connections simulating request/response data pattern, client sends 32KB request, server responds with 32KB .
 - **CPU benefits of SMC-R for streaming connections unaffected by distance (and in several cases better at longer distances)**
 - **Significant response time improvement**

SMC-R RoCE performance at distance – Streaming/Bulk Data (1 session)

Streaming (Bulk Data) 1/20M
SMC-R vs. TCP/IP



Typical TCP/IP configuration
MTU=1492, Large Send disabled

Optimal TCP/IP configuration
MTU=8000, Large Send Enabled

- Notes:
 - STR1: Single TCP connection simulating streaming data pattern, client sends 1 byte, server responds with 20MB of data.
 - **CPU benefits of SMC-R for streaming connections unaffected by distance (and in several cases better at longer distances)**
 - **Significant throughput improvements at distance (improving overall response time significantly)**

Evaluating SMC-R applicability and benefits – SMC Applicability Tool (SMCAT)

As customers express interest in SMC-R and RoCE Express one of the initial questions asked is:

- “What benefit will SMC-R provide in my environment?”
 - Some users are well aware of significant traffic patterns that can benefit from SMC-R
 - But others are unsure of how much of their TCP traffic (in their environment) is:
 - z/OS to z/OS
 - IPSEC?
 - Traffic well suited to SMC-R?
- Reviewing SMF records, using Netstat displays, Ctrace analysis and reports from various Network Management products can provide these insights...

This approach can be a time consuming activity that requires significant expertise.

SMC Applicability Tool Introduction

A new tool called SMC Applicability Tool (SMCAT) has been created that will help customers determine the ***potential*** value of SMC-R in their environment with minimal effort and minimal impact

- SMCAT is integrated within the TCP/IP stack:
Gather new statistics that are used to project SMC-R applicability and benefits for the current system
 - Minimal system overhead, no changes in TCP/IP network flows
 - Produces reports on potential benefits of enabling SMC-R

- Available via the service stream on existing z/OS releases as well
 - z/OS V2R1 - APAR PI39612, PTF UI28867
 - z/OS V1R13 - APAR PI41713, PTF UI29684
 - Does not require:
 - SMC-R code or RoCE hardware to use
 - Any changes in IP configuration (i.e. captures your normal TCP/IP workloads)

SMCAT Usage Overview

Activated by Operator command

(*Vary TCPIP,,SMCAT,dsn(smcatconfig)*) – Input dataset contains:

- Interval Duration, list of IP addresses or IP subnets of peer z/OS systems ((i.e. systems that we can use SMC-R for)
 - If subnets are used, the entire subnet must be comprised of z/OS systems that are SMC-R eligible
 - It is important that all the IP addresses used for establishing TCP connections are specified (including DVIPAs, etc.)

- At the end of the interval a summary report is generated that includes:
 1. **Percent of traffic eligible for SMC-R** (*% of TCP traffic that is eligible for SMC-R*)
 - *All traffic that matches configured IP addresses (not using IPsec or FRCA)*
 2. **Percent of traffic well suited for SMC-R** (*your eligible traffic that is also “well suited” to SMC-R, excludes workloads with very short lived TCP connections that have trivial payloads*)
 - *Includes break out of TCP traffic send sizes (i.e. how large is the payload of each send request)*
 - *Helps users quantify SMC-R benefit (reduced latency vs reduced CPU cost)*

SMCAT Usage Overview (continued)

The Summary Report includes 2 sections based on the specified IP addresses/subnets defined in SMCAT configuration file:

1. Potential benefit:

All TCP traffic that matches the configuration - Includes TCP traffic that could not use SMC-R without changes (TCP traffic that does not meet the direct IP route connectivity requirement)

This represents the opportunity of re-configuring routing topology to enable SMC-R

2. Immediate benefit:

The TCP traffic that can use SMC-R immediately / as is (meets SMC-R direct route connectivity requirements). Subset of section 1.

Detected by the tool automatically (non-routed traffic)

SMC Applicability Tool Sample Report (Part 1. Direct Connections)

TCP SMC traffic analysis for matching direct connections

Connections meeting direct connectivity requirements

- 20% of all TCP connections can use SMC (eligible)
- 90% of eligible connections are well-suited for SMC
- 18% of all TCP traffic (segments) is well-suited for SMC
- 20% of outbound traffic (segments) is well-suited for SMC
- 16% of inbound traffic (segments) is well-suited for SMC

Interval Details:

Total TCP Connections:	120
Total SMC eligible connections:	24
Total SMC well-suited connections:	22
Total outbound traffic (in segments)	110000
SMC well-suited outbound traffic (in segments)	22000
Total inbound traffic (in segments)	100000
SMC well-suited inbound traffic (in segments)	16000

Application send sizes used for well-suited connections:

Size	# sends	Percentage
----	-----	-----
1500 (<=1500):	200	39%
4K (>1500 and <=4k):	150	29%
8K (>4k and <= 8k):	0	0%
16K (>8k and <= 16k):	0	0%
32K (>16k and <= 32k):	0	0%
64K (>32k and <= 64k):	50	10%
256K (>64K and <= 256K):	109	22%
>256K:	0	0%

End of report

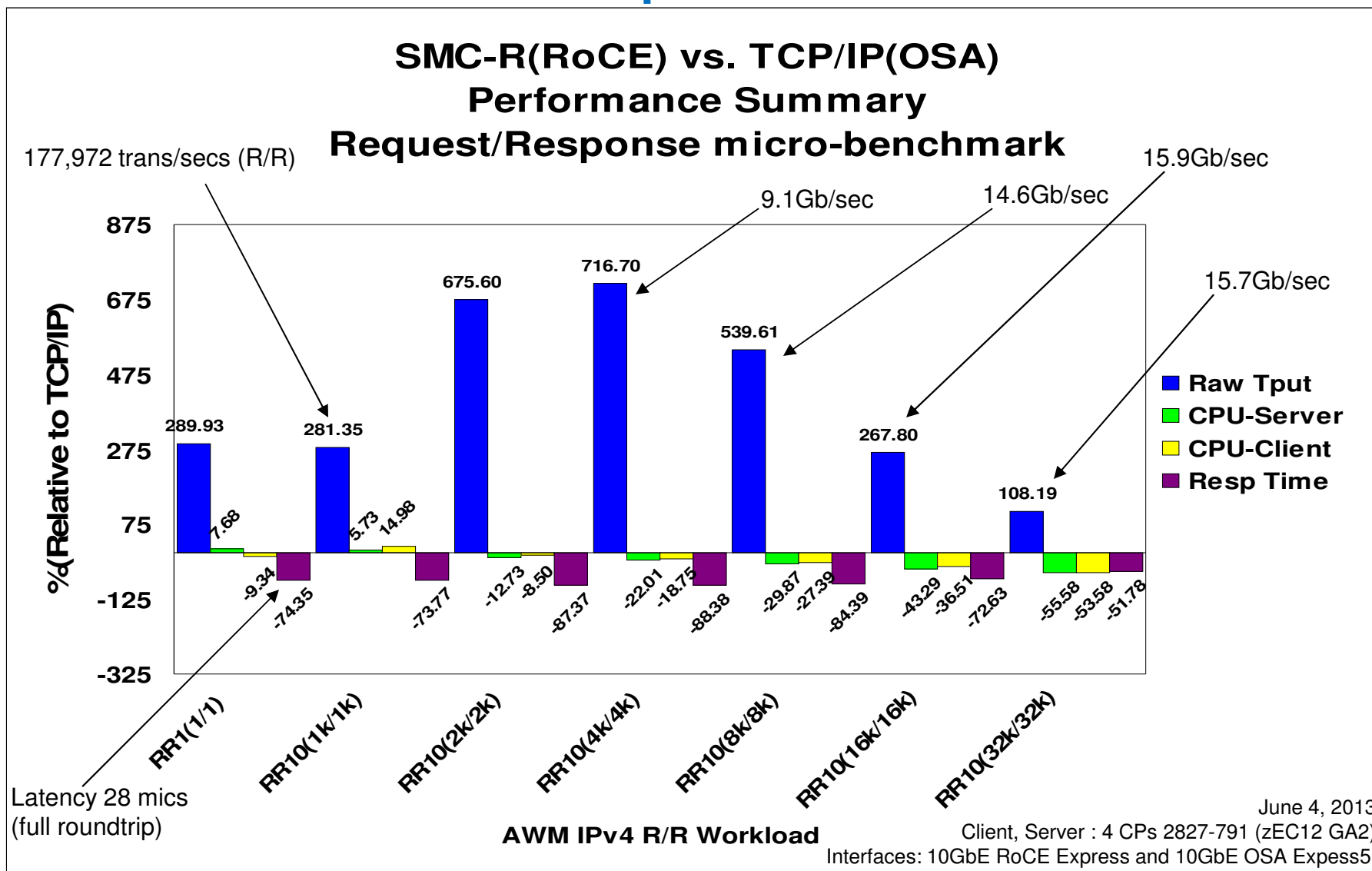
How much of my TCP workload can benefit from SMC?

What kind of CPU savings can I expect from SMC?

This report is repeated for indirect IP connections (different subnets)

SMC-R – Micro benchmark performance results

V2R1



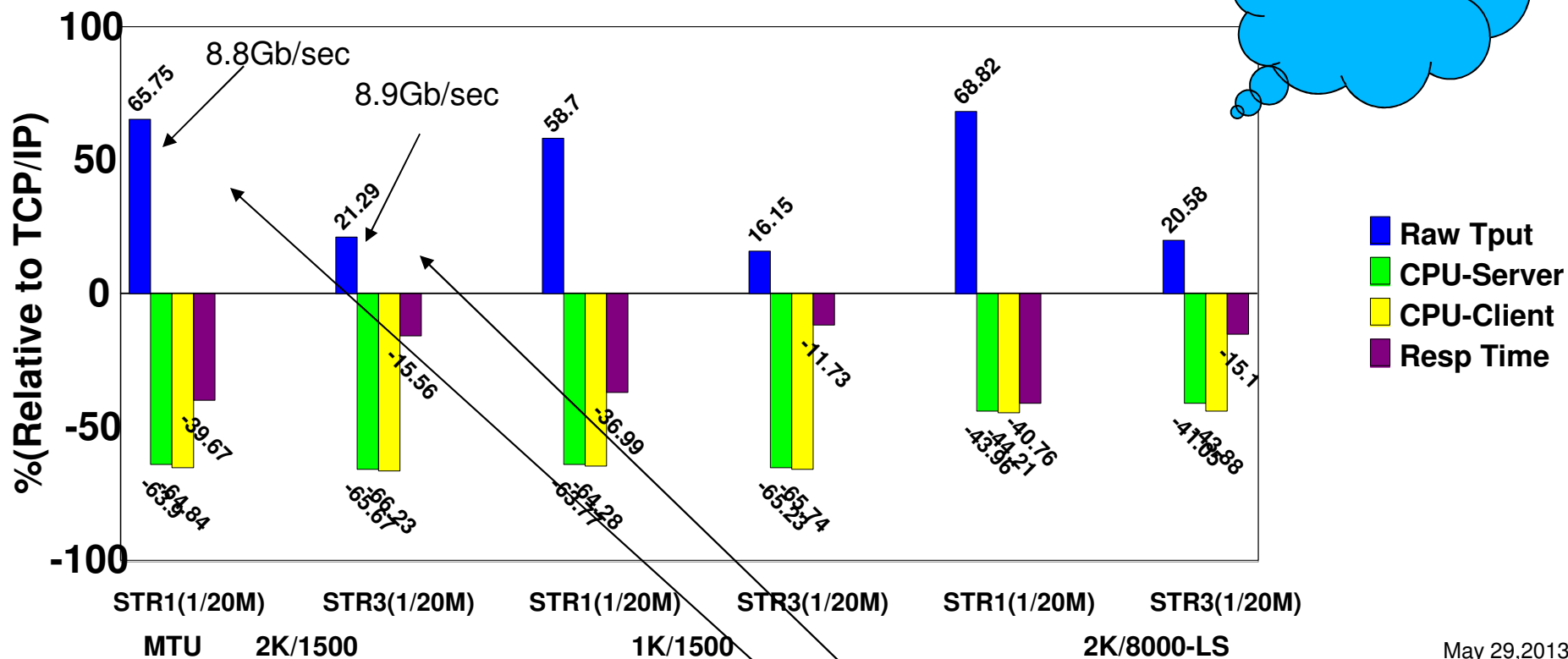
Significant Latency reduction across all data sizes (52-88%)
Reduced CPU cost as payload increases (up to 56% CPU savings)
Impressive throughput gains across all data sizes (Up to +717%)

Note: vs typical OSA customer configuration
 MTU (1500), Large Send disabled
 RoCE MTU: 1K

SMC-R – Micro benchmark performance results

V2R1

z/OS V2R1 SMC-R vs TCP/IP
Streaming Data Performance Summary (AWM)



Saturation reached

May 29, 2013
Client, Server: 2827-791 2CPs
Interfaces: 10GbE RoCE Express and 10GbE

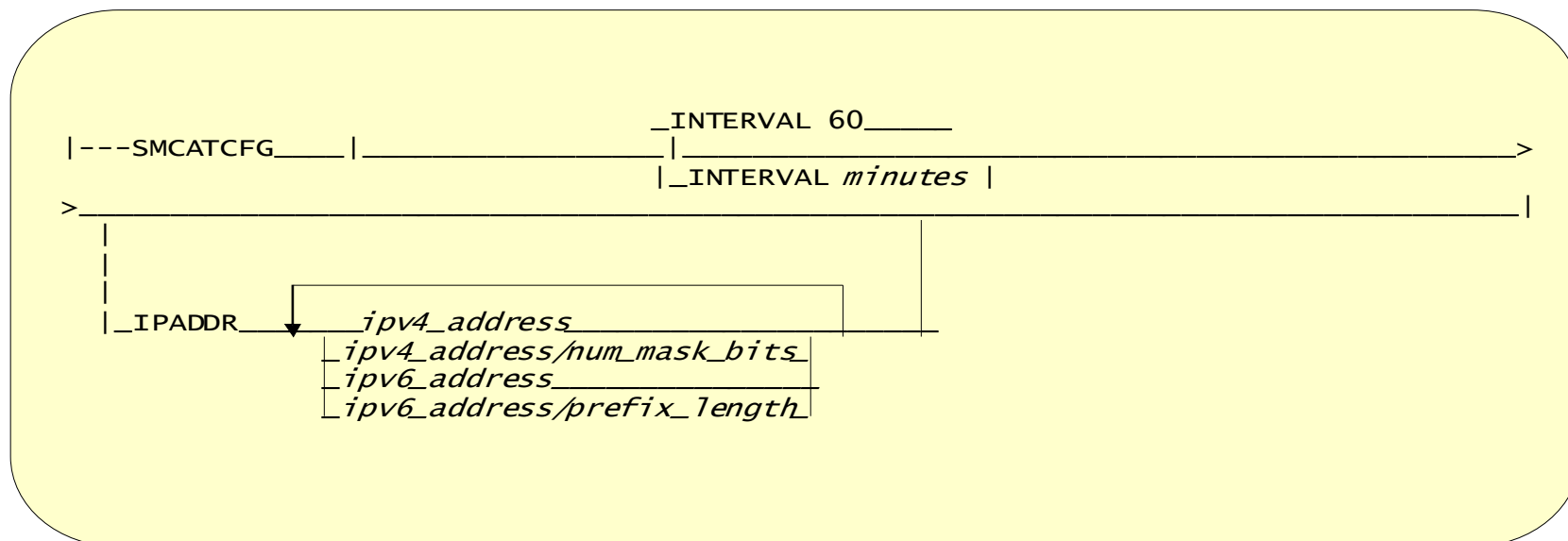
Notes:

- Significant throughput benefits and CPU reduction benefits
 - Up to 69% throughput improvement
 - Up to 66% reduction in CPU costs
- 2K RoCE MTU does yield throughput advantages
- LS – Large Send enabled (Segmentation offload)

Configuring the SMCAT Dataset

SMCAT data set configuration

- Interval defaults to 60 minutes
Max interval is 1440 minutes (24 hours)
- IPADDR is a list of IPv4 and Ipv6 addresses and subnets
256 max combination of addresses and subnets



SMCAT Dataset Example

```
SMCATCFG INTERVAL 120  
IPADDR  
C5::1:2:3:4/126  
9.67.113.61
```



Simple!

When SMCAT is started using this SMCAT configuration data set it will:

- **Monitor TCP traffic for 2 hours for:**
 - IPv6 prefix C5::1:2:3:4/126 and
 - IPv4 address 9.67.113.61

Starting and Stopping SMCAT

Vary TCPIP,,SMCAT command starts and stops the monitoring tool:

- **datasetname** value indicates that SMCAT is being turned on
- **datasetname** contains the SMCATCFG statement that specifies monitoring interval and IP addresses or subnets to be monitored
- **OFF** will stop SMCAT monitoring and generate report

```
>> __Vary__TCPIP,_____,SMCAT,____datasetname____><  
                [__procname__]      [__,OFF__]
```

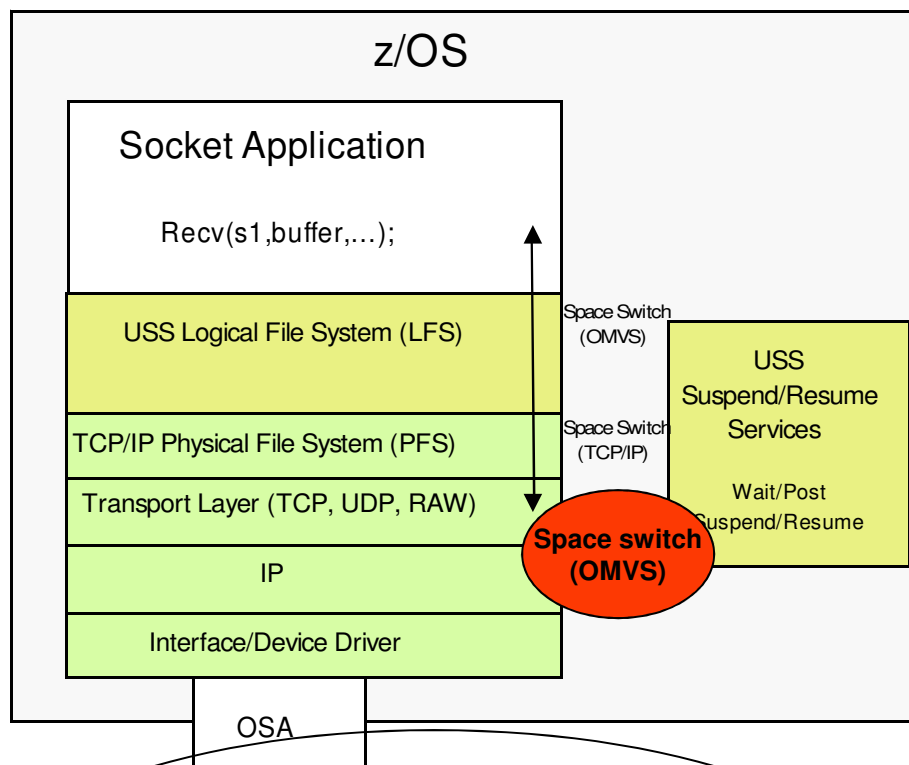
```
VARY TCPIP,TCPPROC,SMCAT,USER99.TCPIP.SMCAT1
```

SMCAT Usage Notes:

- When you have many instances of hosts that provide similar workloads (similar application servers) consider measuring a subset of the hosts and then extrapolating the SMCAT results of your sample across your enterprise data center
- Run the SMCAT tool at different intervals to measure changing workloads

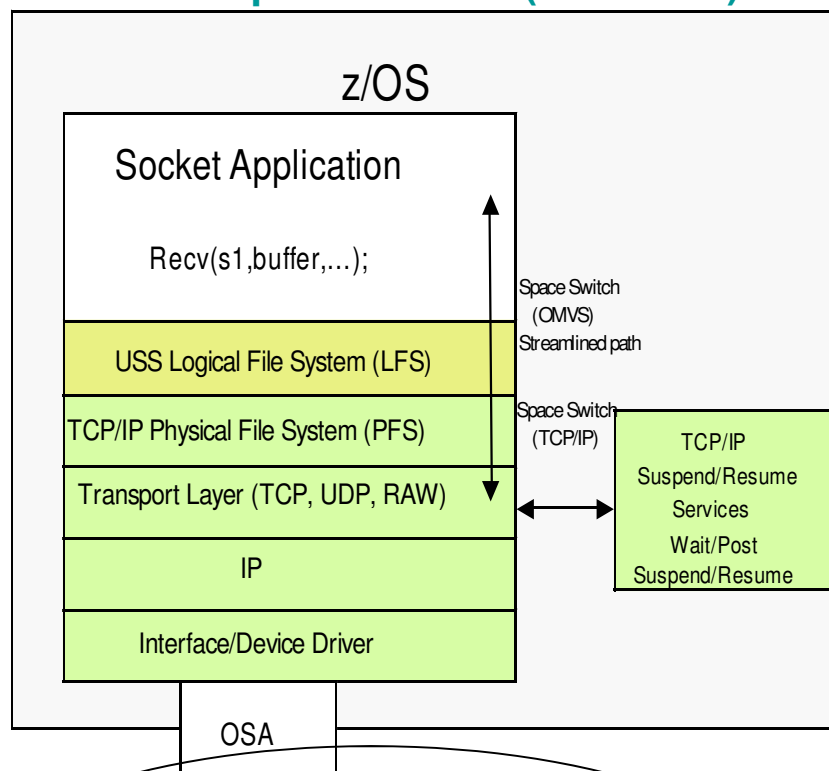
TCP/IP Enhanced Fast Path Sockets

TCP/IP sockets (normal path)



- Full function support for sockets, including support for Unix signals, POSIX compliance
- When TCP/IP needs to suspend a thread waiting for network flows, USS suspend/resume services are invoked

TCP/IP fast path sockets (Pre-V2R1)



- Streamlined path through USS LFS for selected socket APIs
- TCP/IP performs the wait/post or suspend/resume inline using its own services
- Significant reduction in path length

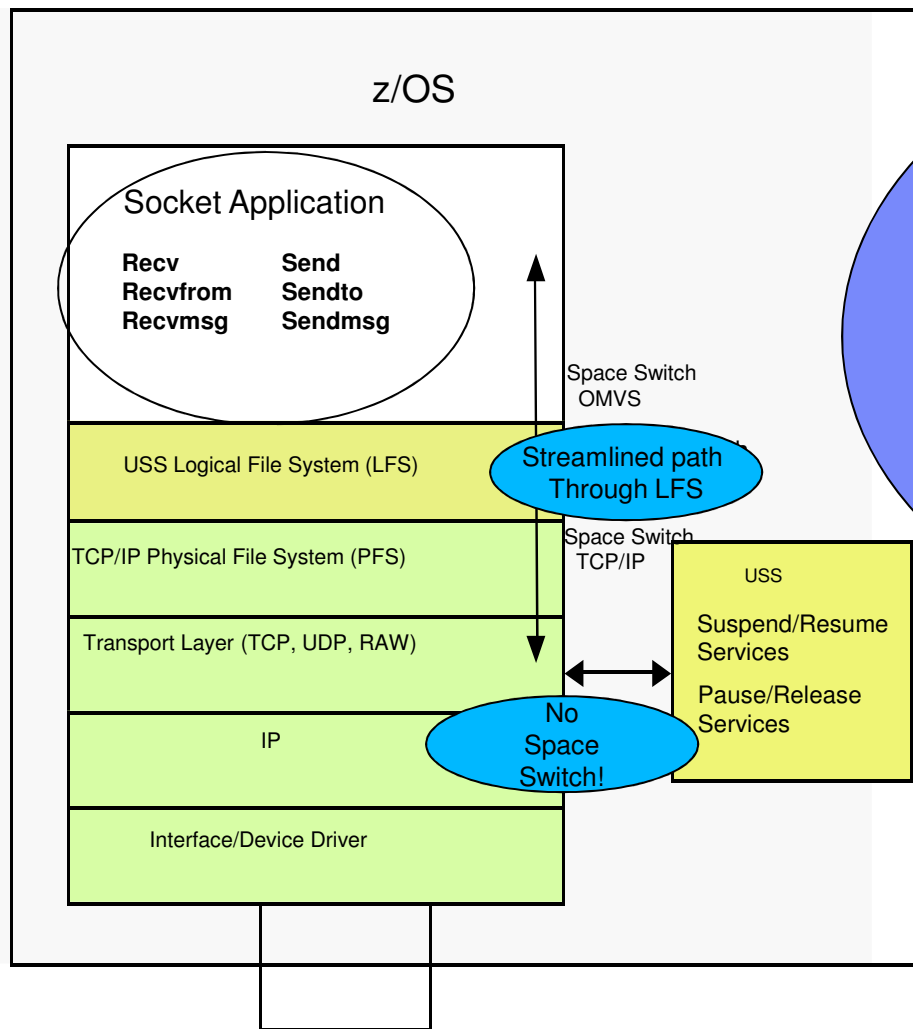
TCP/IP Enhanced Fast Path Sockets

Pre-V2R1 fast path provided CPU savings but not widely adopted:

- No support for Unix signals (other than SIGTERM)
 - Only useful to applications that have no requirement for signal support
- No DBX support (debugger)
- Must be explicitly enabled!
 - BPXK_INET_FASTPATH environment variable
 - `lcc#FastPath` IOCTL
- Only supported for UNIX System Services socket API or the z/OS XL C/C++ Run-time Library functions

TCP/IP Enhanced Fast Path Sockets

V2R1



Fast path sockets performance without all the conditions!

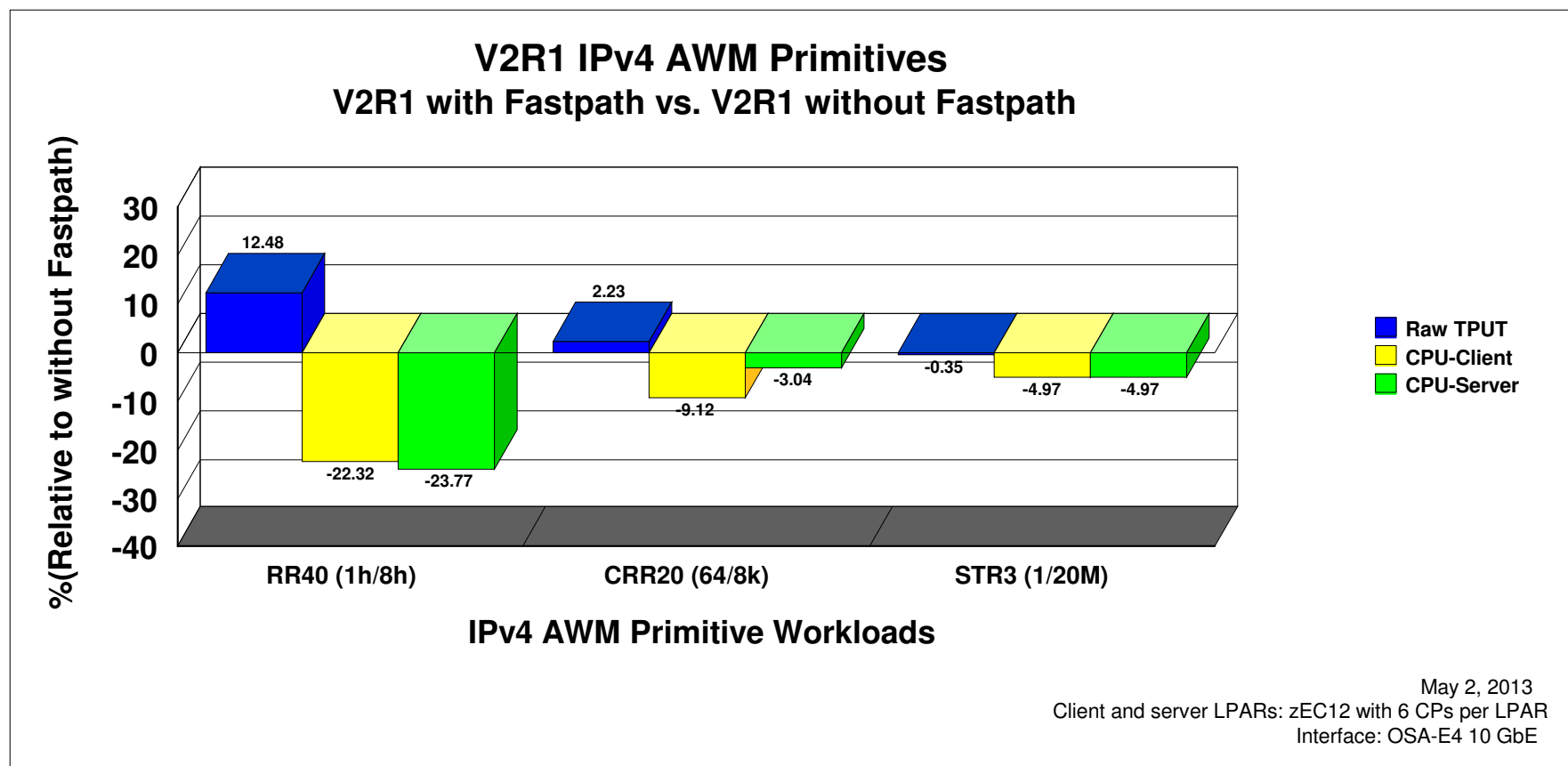
- Enabled by default
- Full POSIX compliance, signals support and DBX support
- Valid for **ALL** socket APIs (with the exception of the Pascal API)

TCP/IP Enhanced Fast Path Sockets

- No new externals
- Still supports “activating Fast path explicitly” to avoid migration issues
 - Provides performance benefits of enhanced Fast Path sockets
 - Keeps the following restrictions:
 - Does not support POSIX signals (blocked by z/OS UNIX)
 - Cannot use dbx debugger

TCP/IP Enhanced Fast Path Sockets

V2R1



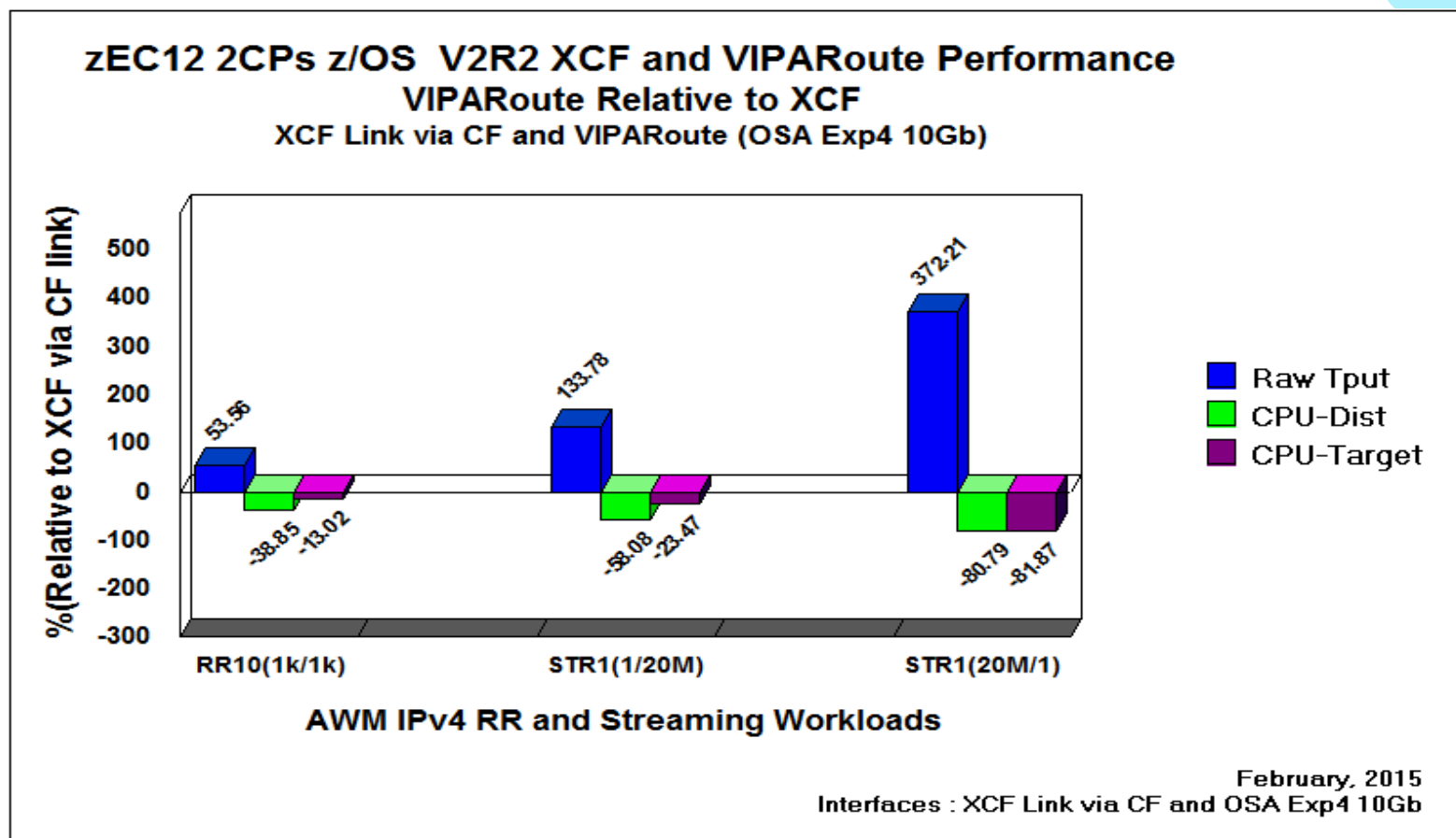
Note: The performance measurements discussed in this presentation are z/OS V2R1 Communications Server numbers and were collected using a dedicated system environment. The results obtained in other configurations or operating system environments may vary.

Communications Server performance: Best practices and other things I find interesting

VIPARROUTE: New GLOBALCONFIG parameter – AUTOADJUSTMSS

➤ Why you should be using VIPARROUTE for Sysplex Distributor Workloads:

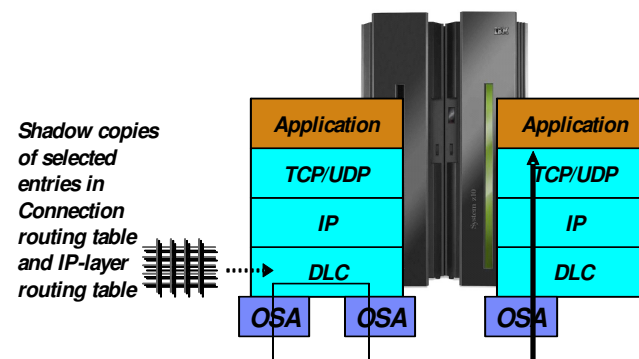
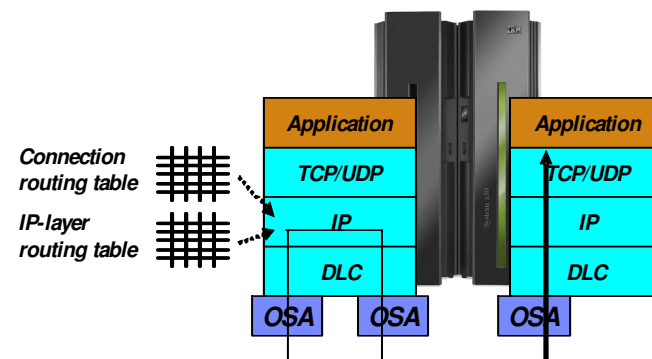
Improved in
V2R2!



Background information: QDIO Accelerator

- Provides fast path IP forwarding for these DLC combinations
 - Inbound QDIO, outbound QDIO or HiperSockets
 - Inbound HiperSockets, outbound QDIO or HiperSockets
- Sysplex Distributor (SD) acceleration
 - Inbound packets over HiperSockets or OSA-E QDIO
 - When SD gets to the target stack using either
 - Dynamic XCF connectivity over HiperSockets
 - VIPAROUTE over OSA-E QDIO
- Improves performance and reduces processor usage for such workloads

V1R11



V2R1

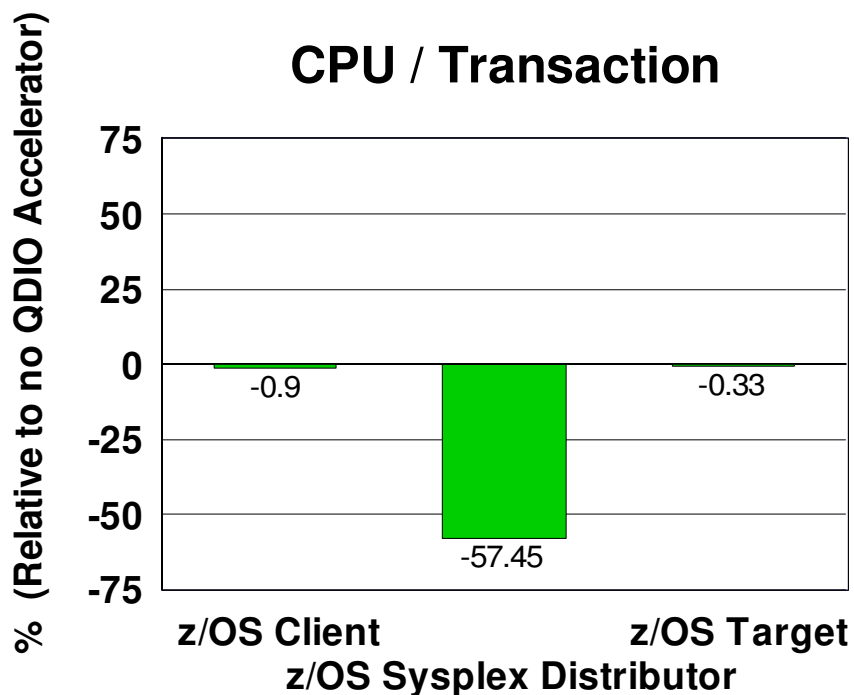
Always allow acceleration for Sysplex Distributed traffic when IP Security enabled

QDIO Accelerator: IPCONFIG syntax and performance results

```

...
|  _NOQDIOACCErator_____ |
|  _____ |
|  _QDIOPriority 1_____ |
|  | _QDIOACCErator_ | _____ |
|  | _____ | _QDIOPriority priority_ |
|  | _____ |
...

```



**Request-Response workload
RR20: 20 sessions, 100 / 800**

FTP using zHPF – Improving throughput

- There are many factors that influence the transfer rates for z/OS FTP connections. Some of the more significant ones are (in order of impact):
 - **DASD read/write access**
 - Data transfer type (Binary, ASCII..)
 - Dataset characteristics (e.g., fixed block or variable)

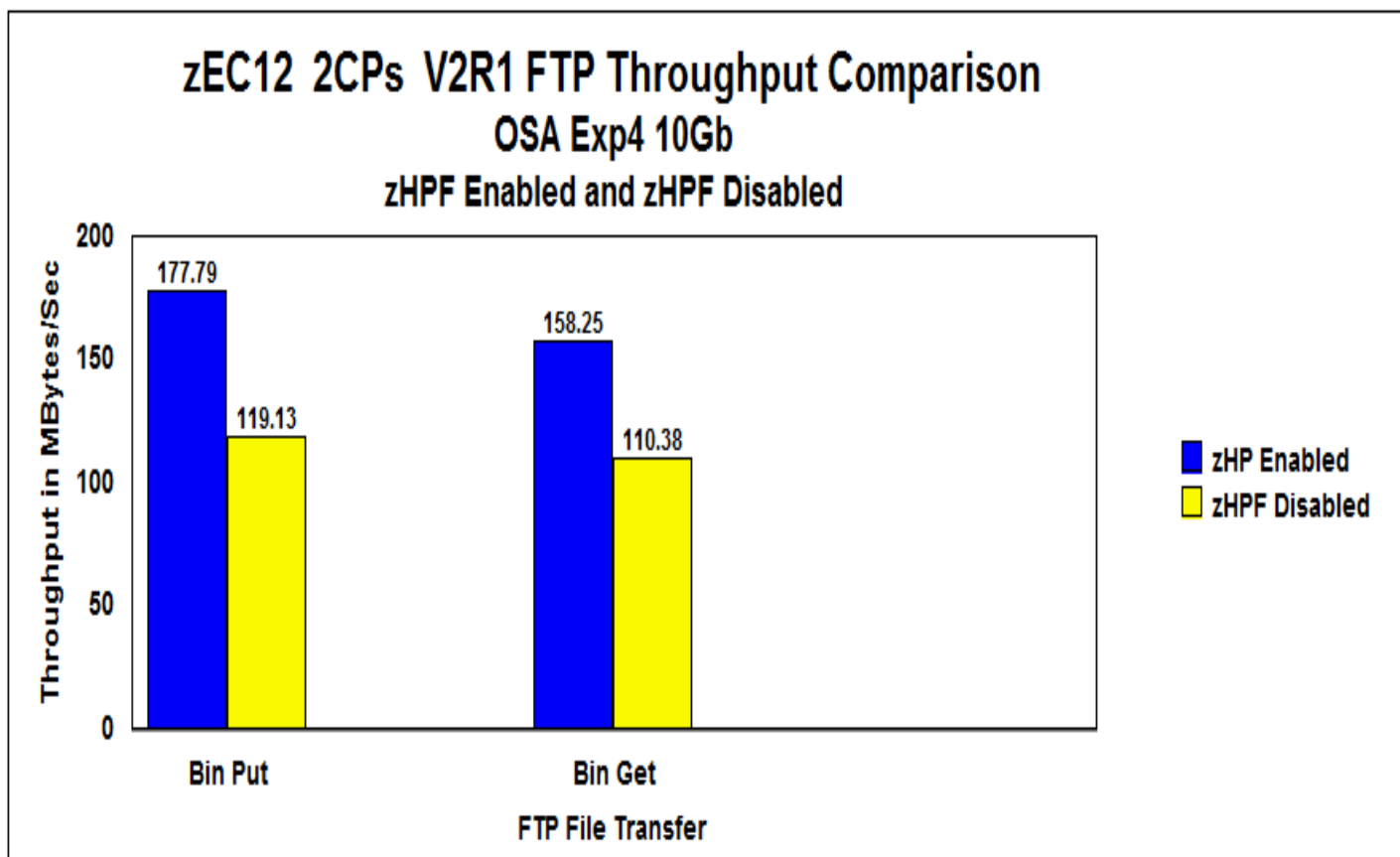
*Note the network (Hipersockets, OSA, 10Gb, SMC-R) characteristics have very little impact when reading from, and writing to, DASD as you will see in our results section.
- zHPF FAQ link
 - <http://www-03.ibm.com/support/techdocs/atsmastr.nsf/WebIndex/FQ127122>
 - Works with DS8000 storage systems

FTP using zHPF – Improving throughput

- FTP Workload
 - z/OS FTP client GET or PUT 1200 MB data set from or to z/OS FTP server
 - DASD to DASD (read from or write to)
 - zHPF enabled/disabled
 - Single file transfer
 - Used Variable block data set for the test
 - Organization PS
 - Record Format ...VB
 - Record Length ...6140
 - Block size23424
 - For Hipersocket
 - Configure GLOBALCONFIG IQDMULTIWRITE

FTP using zHPF – Improving throughput

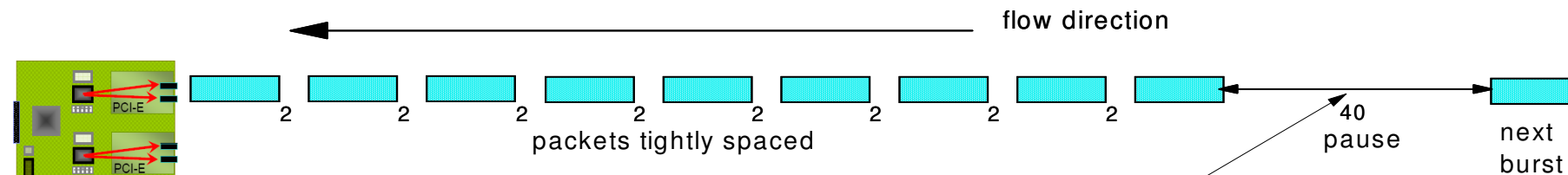
Throughput is improved by 43-49% with Enabling zHPF



Optimizing inbound communications using OSA-Express

Timing Considerations for Various Inbound workloads...

Inbound Streaming Traffic Pattern



receiving OSA-Express3

For inbound streaming traffic, it's most efficient to have OSA defer interrupting z/OS until it sees a pause in the stream.... To accomplish this, we'd want the OSA **LAN-Idle timer** set fairly high (e.g., don't interrupt unless there's a traffic pause of at least 20 microseconds)

Interactive Traffic Pattern



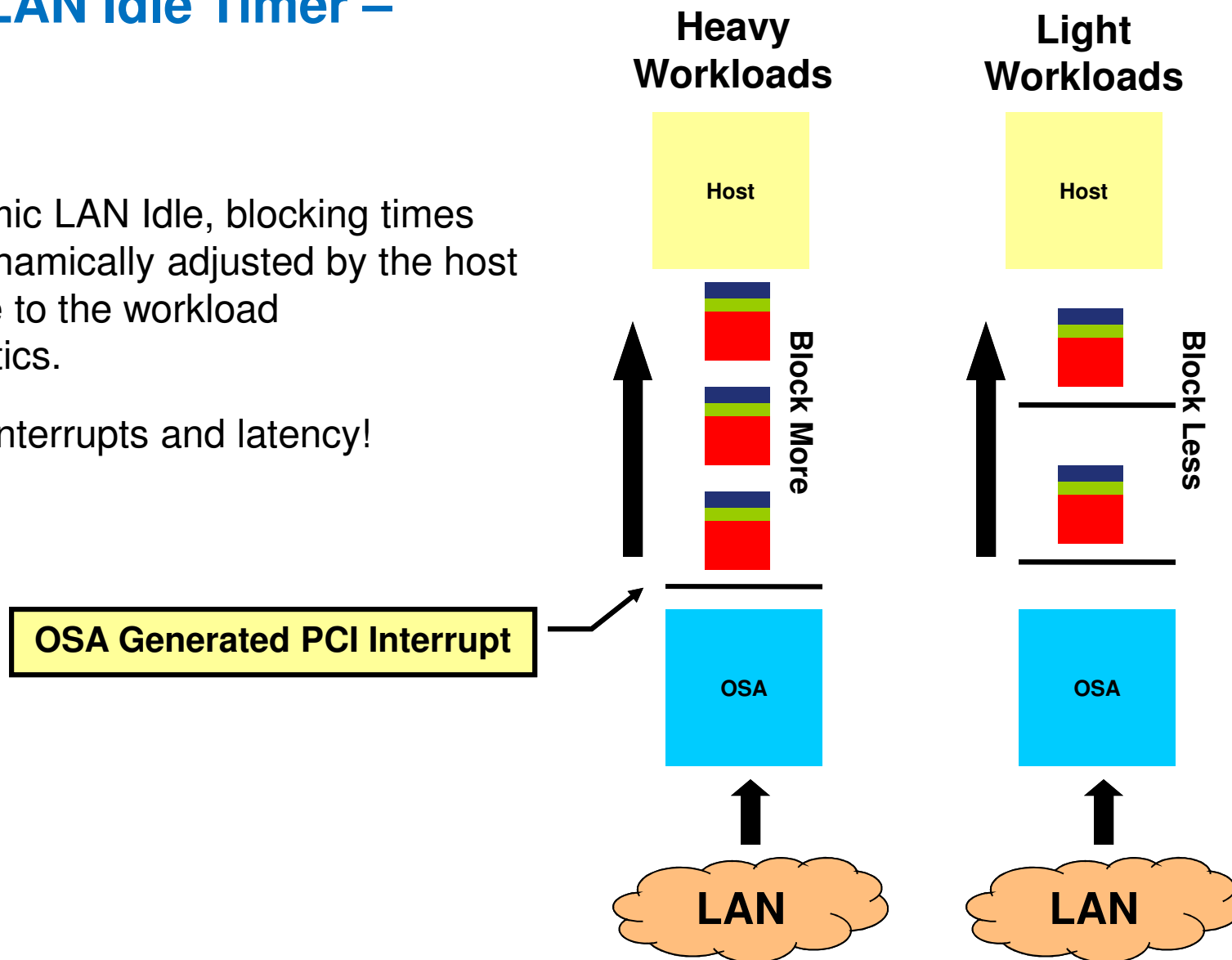
But for interactive traffic, response time would be best if OSA would interrupt z/OS immediately.... To accomplish this, we'd want the OSA **LAN-Idle timer** set as low as it can go (e.g., 1 microsecond)

Read-Side interrupt frequency is all about the LAN-Idle timer!

For detailed discussion on inbound interrupt timing, please see Part 1 of "z/OS Communications Server V1R12 Performance Study: OSA-Express3 Inbound Workload Queueing". <http://www-01.ibm.com/support/docview.wss?uid=swg27005524>

Dynamic LAN Idle Timer –

- With Dynamic LAN Idle, blocking times are now dynamically adjusted by the host in response to the workload characteristics.
- Optimizes interrupts and latency!



Dynamic LAN Idle Timer: Configuration

- Configure INBPERF DYNAMIC on the INTERFACE statement

```

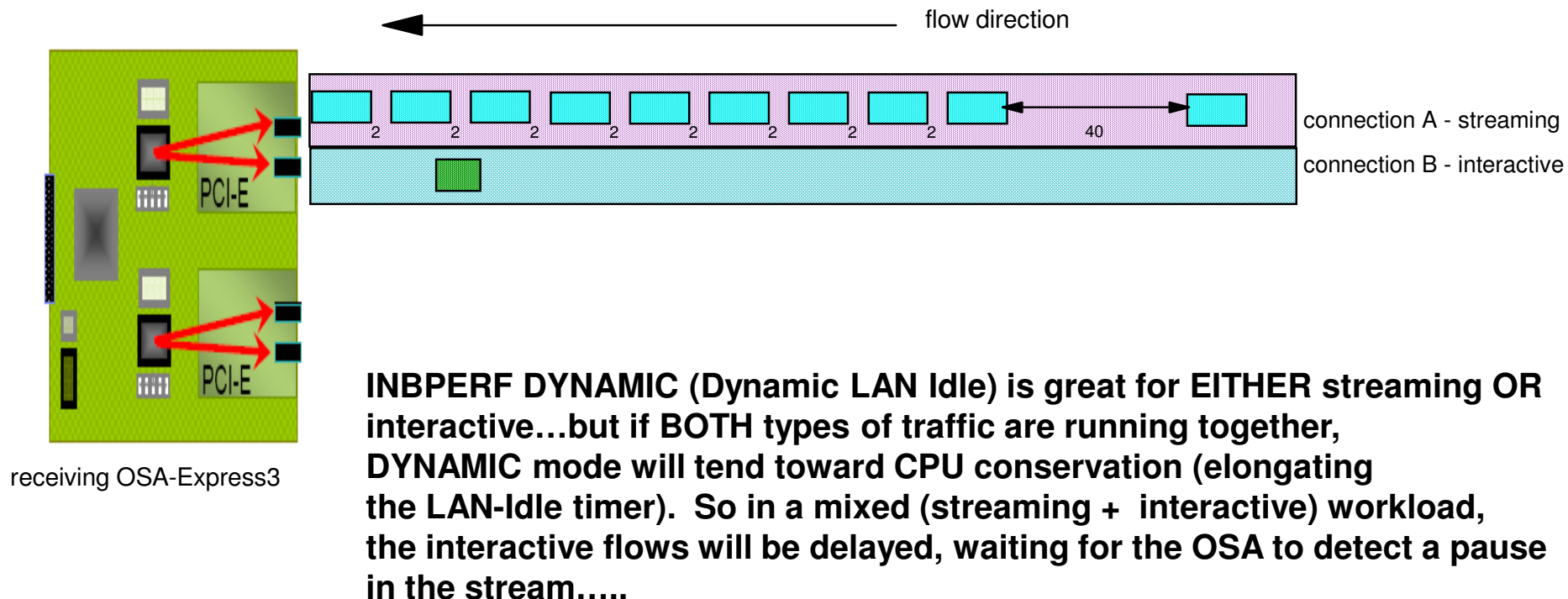
>>-INTERFace--intf_name----->
.
.-INBPERF BALANCED-----.
>+-----+----->
  '-INBPERF--+-DYNAMIC-----+'
      +-MINCPU-----+
      '-MINLATENCY-'
.
  
```

-
-
-

DYNAMIC - a dynamic interrupt-timing value that changes based on current inbound workload conditions ← **Generally Recommended!**

Note: These values cannot be changed without stopping and restarting the interface

Dynamic LAN Idle Timer: But what about mixed workloads?

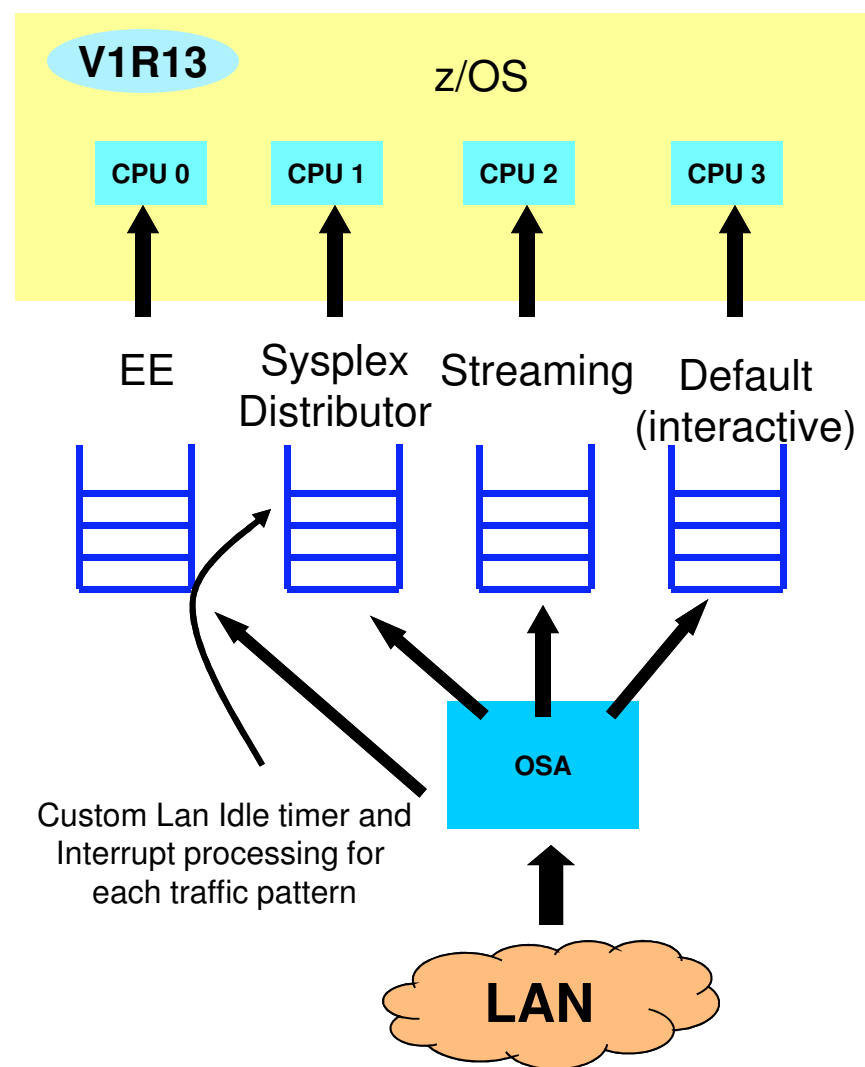


Inbound Workload Queuing

Starting with OSA-Express3S IWQ and z/OS V1R12, OSA now directs streaming traffic onto its own input queue – transparently separating the streaming traffic away from the more latency-sensitive interactive flows...

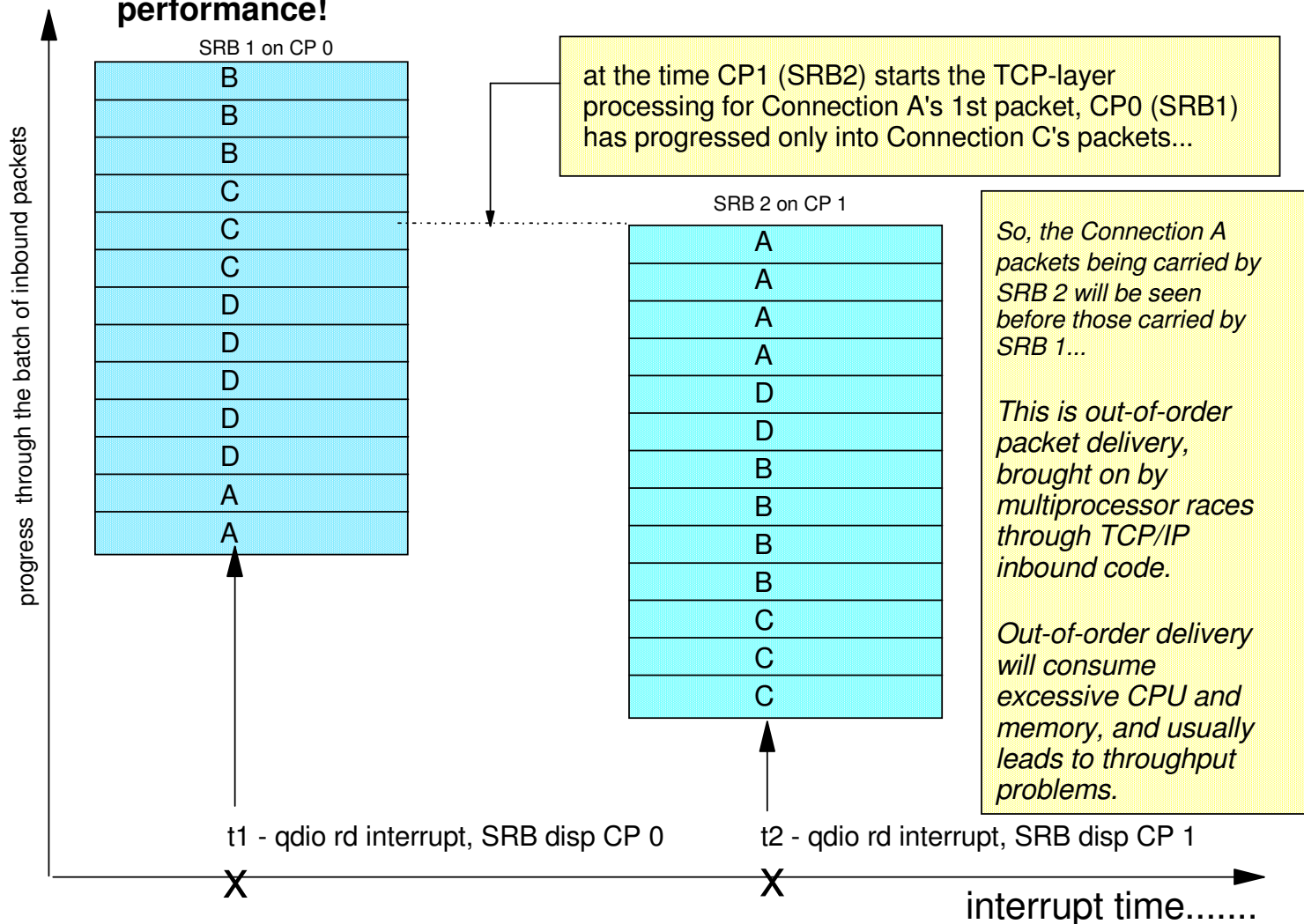
And each input queue has its own LAN-Idle timer, so the Dynamic LAN Idle function can now tune the streaming (bulk) queue to conserve CPU (high LAN-idle timer setting), while generally allowing the primary queue to operate with very low latency (minimizing its LAN-idle timer setting). So interactive traffic (on the primary input queue) may see significantly improved response time.

The separation of streaming traffic away from interactive also enables new streaming traffic efficiencies in Communications Server. This results in improved in-order delivery (better throughput and CPU consumption).



Improved Streaming Traffic Efficiency With IWQ

Before we had IWQ, Multiprocessor races would degrade streaming performance!



IWQ does away with MP-race-induced ordering problems!

With streaming traffic sorted onto its own queue, it is now convenient to service streaming traffic from a single CP (i.e., using a single SRB).

So with IWQ, we no longer have inbound SRB races for streaming data.

QDIO Inbound Workload Queuing – Configuration

- INBPERF DYNAMIC WORKLOADQ enables QDIO Inbound Workload Queuing (IWQ)

```

>>-INTERFace--intf_name----->
.
.-INBPERF BALANCED-----
>--+-----+-->
  |                    .-NOWORKLOADQ-. |
  |   '-INBPERF+--DYNAMIC+-----+--+'
  |           |           '-WORKLOADQ----' |
  |   +-MINCPU-----+
  |   '-MINLATENCY-----'

```

- INTERFACE statements only - no support for DEVICE/LINK definitions
- QDIO Inbound Workload Queuing requires VMAC

QDIO Inbound Workload Queuing - Monitoring

- Display OSAINFO command (V1R12) shows you what's registered in OSA

```

D TCPIP,,OSAINFO,INTFN=V6O3ETHG0
.
Ancillary Input Queue Routing Variables:
Queue Type: BULKDATA Queue ID: 2 Protocol: TCP
Src: 2000:197:11:201:0:1:0:1..221
Dst: 100::101..257
Src: 2000:197:11:201:0:2:0:1..290
Dst: 200::202..514
Total number of IPv6 connections: 2
Queue Type: SYSDIST Queue ID: 3 Protocol: TCP
Addr: 2000:197:11:201:0:1:0:1
Addr: 2000:197:11:201:0:2:0:1
Total number of IPv6 addresses: 2
36 of 36 Lines Displayed
End of report
  
```

- BULKDATA queue registers 5-tuples with OSA (streaming connections)
- SYSDIST queue registers Distributable DVIPAs with OSA

QDIO Inbound Workload Queuing: Netstat DEvlinks/-d

- Display TCPIP,,NETSTAT,DEVlinks to see whether QDIO inbound workload queuing is enabled for a QDIO interface

```
D TCPIP,,NETSTAT,DEVlinks,INTFNAME=QDIO4101L
EZD0101I NETSTAT CS V1R12 TCPCS1
INTFNAME: QDIO4101L          INTFTYPE: IPAQENET   INTFSTATUS: READY
PORTNAME: QDIO4101  DATAPATH: 0E2A      DATAPATHSTATUS: READY
CHPIDTYPE: OSD
SPEED: 0000001000
...
READSTORAGE: GLOBAL (4096K)
INBPERF: DYNAMIC
WORKLOADQUEUEING: YES
CHECKSUMOFFLOAD: YES
SECCLASS: 255                MONSYSPLEX: NO
ISOLATE: NO                  OPTLATENCYMODE: NO
...
1 OF 1 RECORDS DISPLAYED
END OF THE REPORT
```


QDIO Inbound Workload Queuing: Display TRLE

- Display NET,TRL,TRLE=trlename to see whether QDIO inbound workload queuing is in use for a QDIO interface

```

D NET,TRL,TRLE=QDIO101
IST097I DISPLAY ACCEPTED
...
IST2263I PORTNAME = QDIO4101    PORTNUM =    0    OSA CODE LEVEL = ABCD
...
IST1221I DATA  DEV = 0E2A STATUS = ACTIVE        STATE = N/A
IST1724I I/O TRACE = OFF  TRACE LENGTH = *NA*
IST1717I ULPID = TCPCS1
IST2310I ACCELERATED ROUTING DISABLED
IST2331I QUEUE    QUEUE    READ
IST2332I ID      TYPE     STORAGE
IST2205I -----  -----  -----
IST2333I RD/1    PRIMARY   4.0M(64 SBALS)
IST2333I RD/2    BULKDATA  4.0M(64 SBALS)
IST2333I RD/3    SYSDIST   4.0M(64 SBALS)
IST2333I RD/4    EE        4.0M(64 SBALS)
...
IST924I -----
IST314I END

```

QDIO Inbound Workload Queuing: Netstat ALL/-A

- Display TCPIP,,NETSTAT,ALL to see whether QDIO inbound workload BULKDATA queueing is in use for a given connection

```
D TCPIP, , NETSTAT, ALL, CLIENT=USER1
EZD0101I NETSTAT CS V1R12 TCPCS1
CLIENT NAME: USER1                CLIENT ID: 00000046
LOCAL SOCKET:  ::FFFF:172.16.1.1..20
FOREIGN SOCKET:  ::FFFF:172.16.1.5..1030
  BYTESIN:           00000000000023316386
  BYTESOUT:          00000000000000000000
  SEGMENTSIN:       00000000000000016246
  SEGMENTSOUT:      00000000000000000922
  LAST TOUCHED:    21:38:53          STATE:          ESTABLISH
...
Ancillary Input Queue: Yes
BulkDataIntfName: QDIO4101L
...
APPLICATION DATA:  EZAFTPOS D USER1      C      PSSS
-----
1 OF 1 RECORDS DISPLAYED
END OF THE REPORT
```

QDIO Inbound Workload Queuing: Netstat STATS/-S

- Display TCPIP,,NETSTAT,STATS to see the total number of TCP segments received on BULKDATA queues

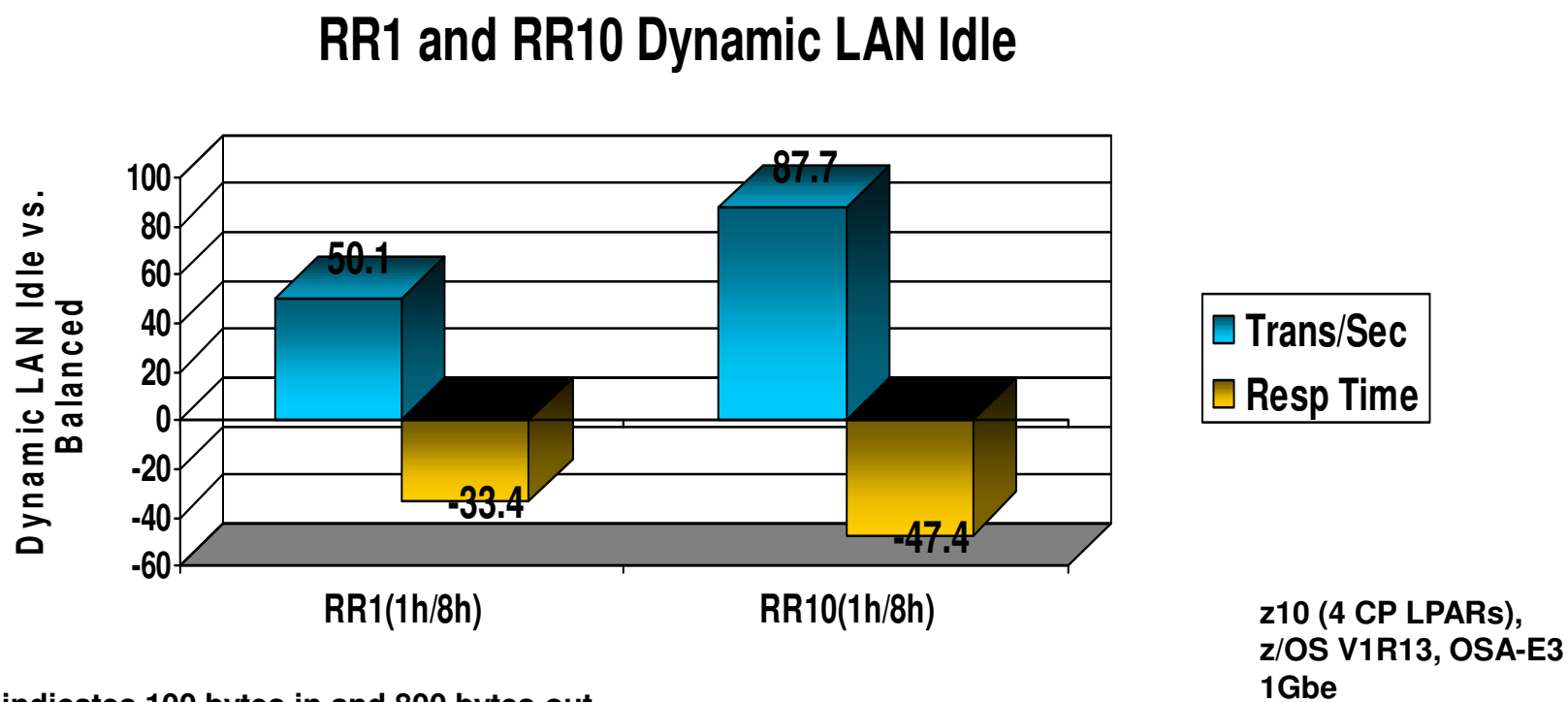
```
D TCPIP, , NETSTAT, STATS, PROTOCOL=TCP
EZD0101I NETSTAT CS V1R12 TCPCS1
TCP STATISTICS
CURRENT ESTABLISHED CONNECTIONS      = 6
ACTIVE CONNECTIONS OPENED             = 1
PASSIVE CONNECTIONS OPENED            = 5
CONNECTIONS CLOSED                     = 5
ESTABLISHED CONNECTIONS DROPPED        = 0
CONNECTION ATTEMPTS DROPPED            = 0
CONNECTION ATTEMPTS DISCARDED          = 0
TIMEWAIT CONNECTIONS REUSED           = 0
SEGMENTS RECEIVED                      = 38611
...
SEGMENTS RECEIVED ON OSA BULK QUEUES= 2169
SEGMENTS SENT                          = 2254
...
END OF THE REPORT
```

Quick INBPERF Review Before We Push On....

- The original static INBPERF settings (MINCPU, MINLATENCY, BALANCED) provide sub-optimal performance for workloads that tend to shift between request/response and streaming modes.
- We therefore **recommend customers specify INBPERF DYNAMIC**, since it self-tunes, to provide excellent performance even when inbound traffic patterns shift.
- Inbound Workload Queueing (IWQ) mode is an extension to the Dynamic LAN Idle function. IWQ improves upon the DYNAMIC setting, in part because it provides finer interrupt-timing control for mixed (interactive + streaming) workloads.

Dynamic LAN Idle Timer: Performance Data

Dynamic LAN Idle improved RR1 TPS 50% and RR10 TPS by 88%. Response Time for these workloads is improved 33% and 47%, respectively.



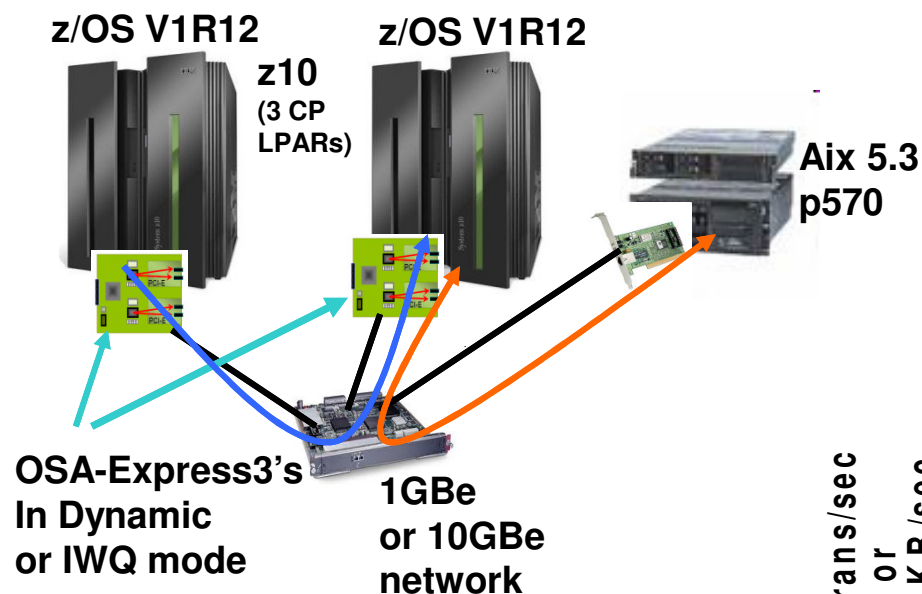
1h/8h indicates 100 bytes in and 800 bytes out

Note: The performance measurements discussed in this presentation are z/OS V1R13 Communications Server numbers and were collected using a dedicated system environment. The results obtained in other configurations or operating system environments may vary.

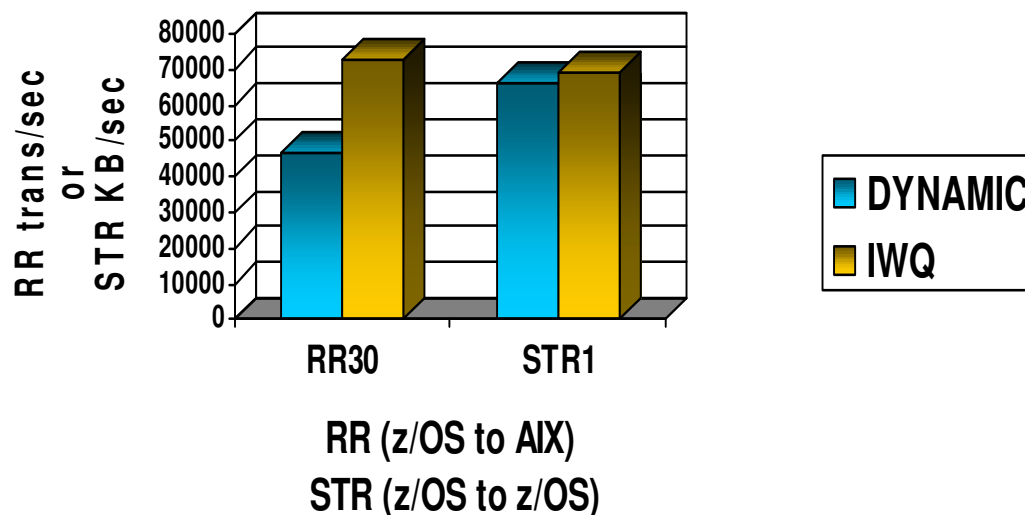
Inbound Workload Queuing: Performance Data

IWQ: Mixed Workload Results vs DYNAMIC:

- z/OS<->AIX R/R Throughput improved 55% (Response Time improved 36%)
- Streaming Throughput also improved in this test: +5%

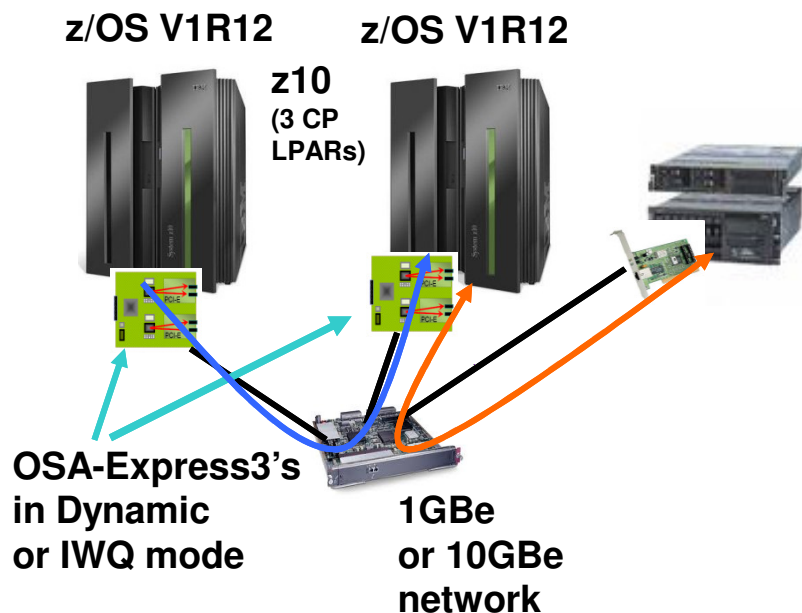


Mixed Workload (IWQ vs Dynamic)



For z/OS outbound streaming to another platform, the degree of performance boost (due to IWQ) is relative to receiving platform's sensitivity to out-of-order packet delivery. For streaming INTO z/OS, IWQ will be especially beneficial for multi-CP configurations.

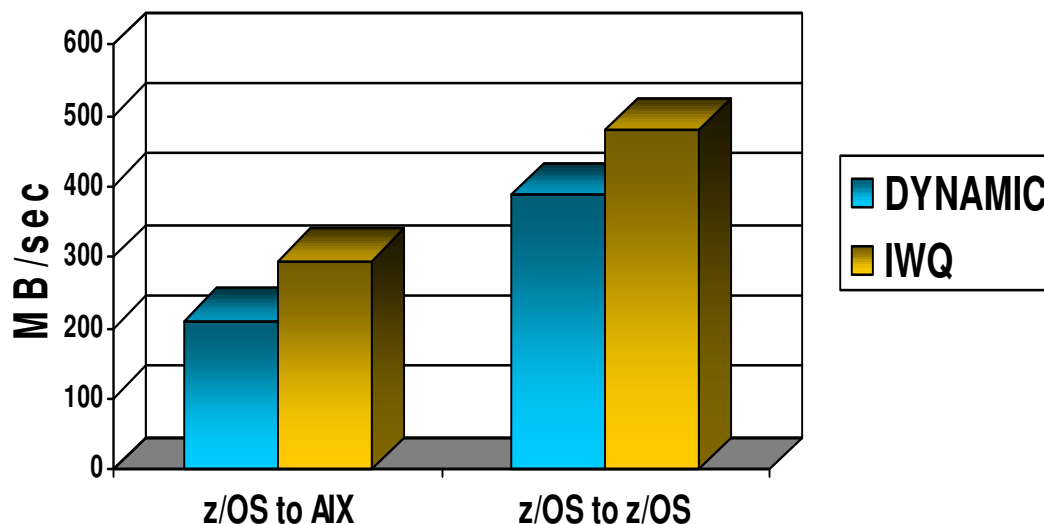
Inbound Workload Queuing: Performance Data



IWQ: Pure Streaming Results vs DYNAMIC:

- z/OS<->AIX Streaming Throughput improved 40%
- z/OS<->z/OS Streaming Throughput improved 24%

Pure Streaming (IWQ vs Dynamic)



For z/OS outbound streaming to another platform, the degree of performance boost (due to IWQ) is relative to receiving platform's sensitivity to out-of-order packet delivery. For streaming INTO z/OS, IWQ will be especially beneficial for multi-CP configurations.

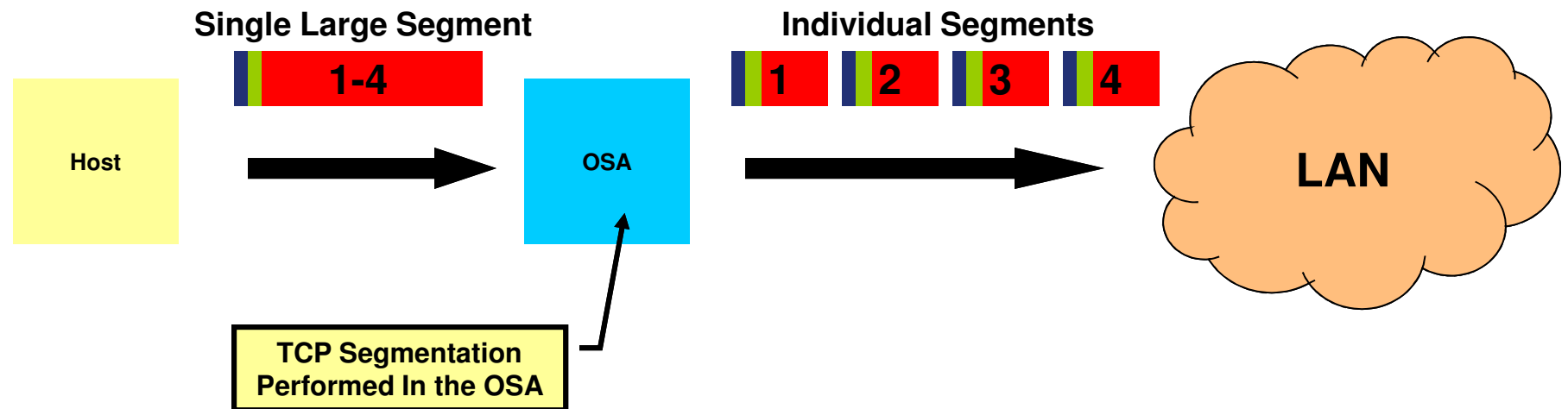
IWQ Usage Considerations:

- Minor ECSA Usage increase: IWQ will grow ECSA usage by 72KBytes (per OSA interface) if Sysplex Distributor (SD) or EE is in use; 36KBytes if SD and EE are not in use
- IWQ requires OSA-Express3/OSA-Express4/OSA-Express5 in QDIO mode running on zEnterprise 196/ zEC12(for OSAE5).
- IWQ must be configured using the INTERFACE statement (not DEVICE/LINK)
- IWQ is not supported when z/OS is running as a z/VM guest with simulated devices (VSWITCH or guest LAN)

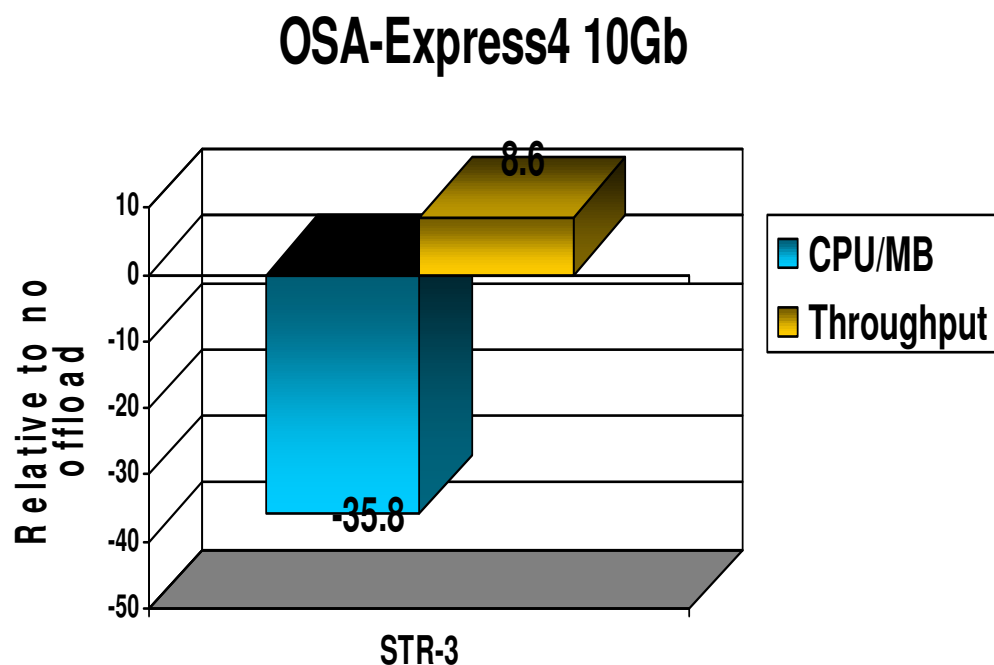
Optimizing outbound communications using OSA- Express

TCP Segmentation Offload

- Segmentation consumes (high cost) host CPU cycles in the TCP stack
- Segmentation Offload (also referred to as “Large Send”)
 - Offload most IPv4 and/or IPv6 TCP segmentation processing to OSA
 - Decrease host CPU utilization
 - Increase data transfer efficiency
 - Checksum offload also added for IPv6



z/OS Segmentation Offload performance measurements



Send buffer size: 180K for streaming workloads

Segmentation offload may significantly reduce CPU cycles when sending bulk data from z/OS!

Note: The performance measurements discussed in this presentation are z/OS V1R13 Communications Server numbers and were collected using a dedicated system environment. The results obtained in other configurations or operating system environments may vary.

TCP Segmentation Offload: Configuration

- Enabled with IPCONFIG/IPCONFIG6 SEGMENTATIONOFFLOAD

```

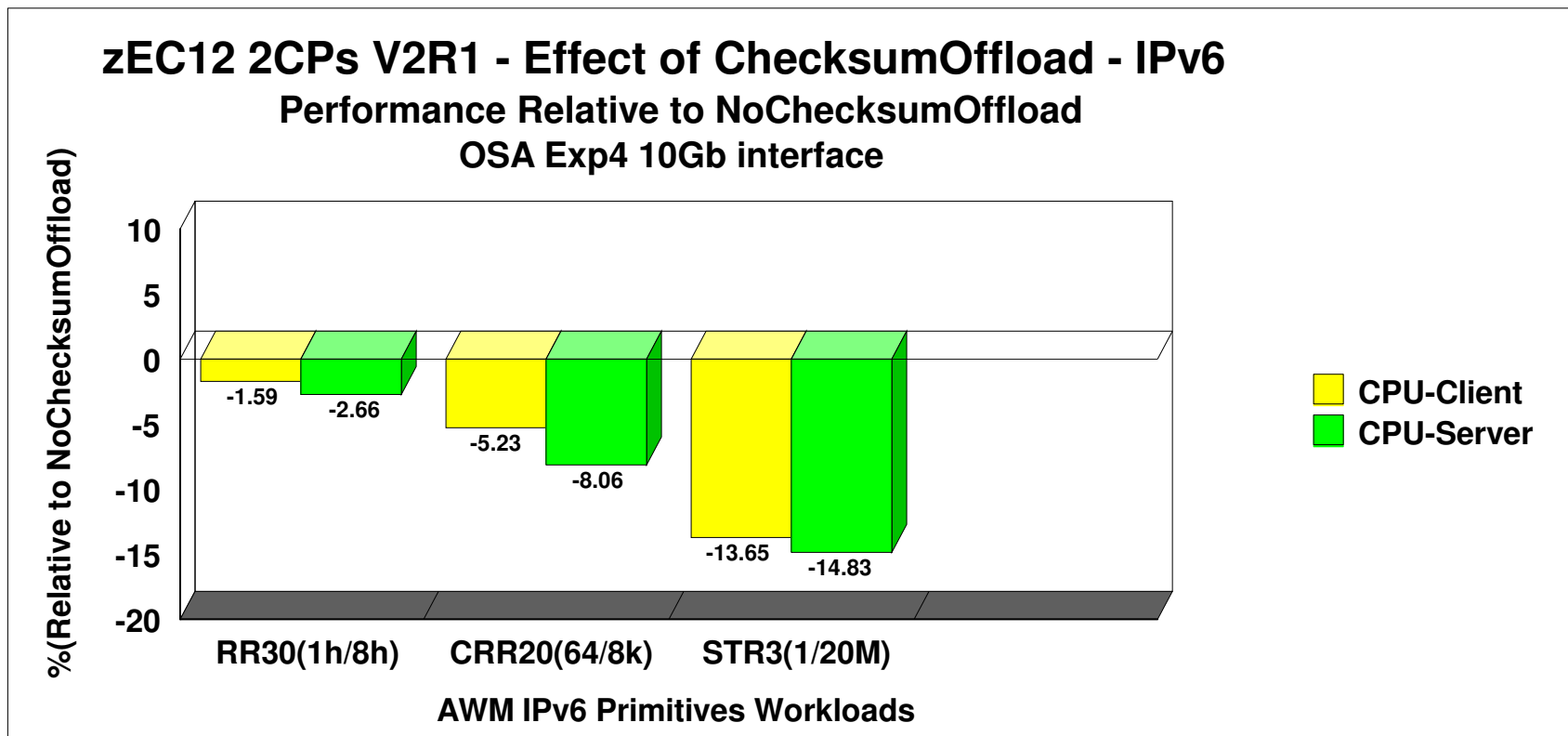
>>-IPCONFIG----->
.
.
>-----+-----+----->
       | .-NOSEGMENTATIONOFFLoad-. |
       +-----+-----+-----+
       | '-SEGMENTATIONOFFLoad---' |
  
```

- Disabled by default
- Previously enabled via GLOBALCONFIG
- Segmentation cannot be offloaded for
 - Packets to another stack sharing OSA port
 - IPSec encapsulated packets
 - When multipath is in effect (unless all interfaces in the multipath group support segmentation offload)

Reminder!
Checksum Offload
enabled by default

z/OS Checksum Offload performance measurements

V1R13



Note: The performance measurements discussed in this presentation are z/OS V2R1 Communications Server numbers and were collected using a dedicated system environment. The results obtained in other configurations or operating system environments may vary.

z/OS Communications Server Performance Summaries

z/OS Communications Server Performance Summaries

- Performance of each z/OS Communications Server release is studied by an internal performance team
- Summaries are created and published online
 - <http://www-01.ibm.com/support/docview.wss?rs=852&uid=swg27005524>
- Recently added:
 - The z/OS V2R1 Communications Server Performance Summary
 - Release to release comparisons
 - Capacity planning information
 - IBM z/OS Shared Memory Communications over RDMA: Performance Considerations – Whitepaper
- Coming soon:
 - The z/OS V2R2 Communications Server Performance Summary

z/OS Communications Server Performance Website

www-01.ibm.com/support/docview.wss?uid=swg27005524

IBM z/OS Communication x

www-01.ibm.com/support/docview.wss?rs=852&uid=swg27005524

United States

IBM Industries & solutions Services Products Support & downloads My IBM Search

← Go to IBM Support Portal

z/OS Communications Server performance index

Tags
Add a tag | Search all tags
Add a tag >
My tags | All tags
View as cloud | list

White paper

Abstract
z/OS Communications Server performance summary reports

Content

- [z/OS V2R1 Communications Server Performance Summary](#)
- [IBM z/OS Shared Memory Communications over RDMA: Performance Considerations](#)
- [z/OS V1R13 Communications Server Performance Summary](#)
- [z/OS V1R12 Communications Server Performance Summary](#)
- [z/OS V1R12 Communications Server Performance Study: OSA Express3 Inbound Workload Queueing](#)
- [z/OS V1R11 Communications Server Performance Summary](#)
- [z/OS V1R10 Communications Server Large Send Performance Summary](#)
- [z/OS V1R10 Communications Server Performance Summary](#)
- [z/OS V1R9 Communications Server TN3270 Capacity Planning 2008](#)
- [z/OS V1R9 Communications Server Performance Summary](#)
- [z/OS IP Network Security: Capacity Planning for zLIP Assisted IPSec](#)
- [z/OS V1R8 Communications Server Performance Summary](#)

Rate this page:
Average rating ★★★★★ (2 users)

Add comments +

Document information

More support for:
[z/OS Communications Server](#)
All

Software version:
1.8, 1.9, 1.10, 1.11, 1.12, 1.13, 2.1

Operating system(s):
z/OS

Reference #:
7005524

Modified date:
2013-01-28

Translate my page
Select Language ▾

Rate this page:
Average rating ★★★★★ (2 users)

Please fill out your session evaluation

- z/OS CS Performance Improvements
- QR Code:



Find us on Facebook at
<http://www.facebook.com/IBMCommserver>

Follow us on Twitter at
http://www.twitter.com/IBM_Commserver

Read the z/OS Communications Server blog at
<http://tinyurl.com/zoscsblog>

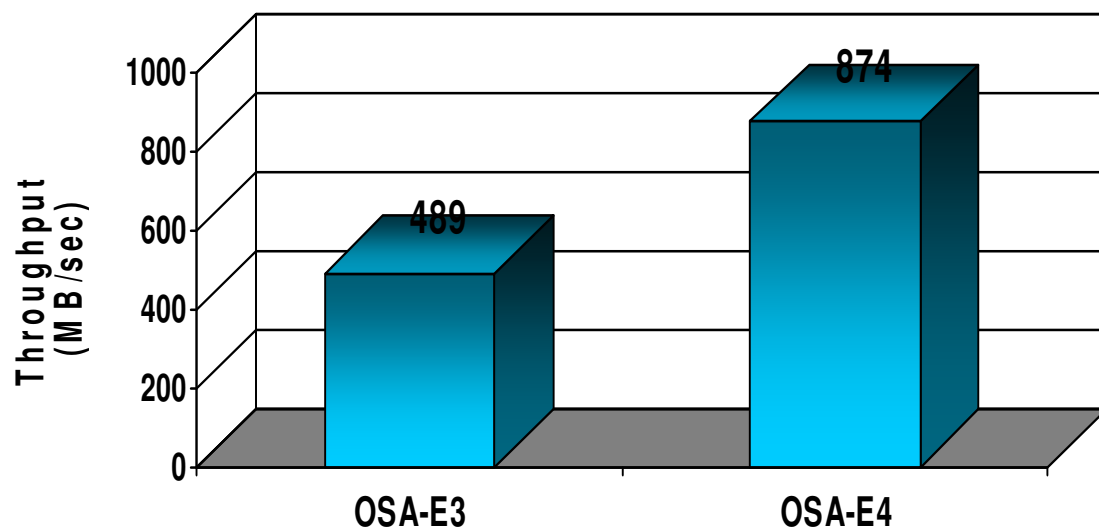
Visit the z/OS CS YouTube channel at
<http://www.youtube.com/user/zOSCommServer>

Appendix OSA-Express4/5

OSA-Express4/5 Enhancements – 10GB improvements

- Improved on-card processor speed and memory bus provides better utilization of 10GB network

OSA 10GBe - Inbound Bulk traffic



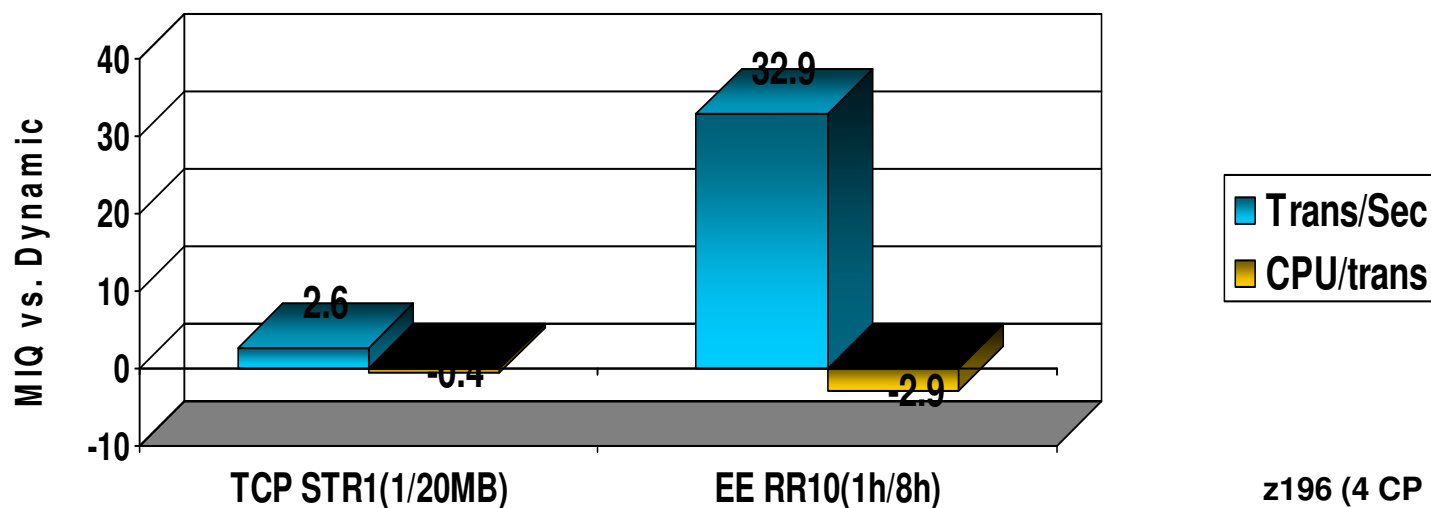
z196 (4 CP LPARs),
z/OS V1R13, OSA-
E3/OSA-E4 10Gbe

Note: The performance measurements discussed in this presentation are z/OS V1R13 Communications Server numbers and were collected using a dedicated system environment. The results obtained in other configurations or operating system environments may vary.

OSA-Express4 Enhancements – EE Inbound Queue

- Enterprise Extender queue provides internal optimizations
 - EE traffic processed quicker
 - Avoids memory copy of data

OSA 1Gbe - mixed TCP and EE workloads



z196 (4 CP LPARs),
z/OS V1R13, OSA-
E3/OSA-E4 1Gbe

Note: The performance measurements discussed in this presentation are z/OS V1R13 Communications Server numbers and were collected using a dedicated system environment. The results obtained in other configurations or operating system environments may vary.

OSA-Express4 Enhancements – Other improvements

- Checksum Offload support for IPv6 traffic
- Segmentation Offload support for IPv6 traffic

Appendix

Detailed Usage Considerations for IWQ and OLM

IWQ Usage Considerations:

- Minor ECSA Usage increase: IWQ will grow ECSA usage by 72KBytes (per OSA interface) if Sysplex Distributor (SD) is in use; 36KBytes if SD is not in use
- IWQ requires OSA-Express3 in QDIO mode running on IBM System z10 or OSA-Express3/OSA-Express4 in QDIO mode running on zEnterprise 196.
 - For z10: the minimum field level recommended for OSA-Express3 is microcode level- Driver 79, EC N24398, MCL006
 - For z196 GA1: the minimum field level recommended for OSA-Express3 is microcode level- Driver 86, EC N28792, MCL009
 - For z196 GA2: the minimum field level recommended for OSA-Express3 is microcode level- Driver 93, EC N48158, MCL009
 - For z196 GA2: the minimum field level recommended for OSA-Express4 is microcode level- Driver 93, EC N48121, MCL010
- IWQ must be configured using the INTERFACE statement (not DEVICE/LINK)
- IWQ is not supported when z/OS is running as a z/VM guest with simulated devices (VSWITCH or guest LAN)
- Make sure to apply z/OS V1R12 PTF UK61028 (APAR PM20056) for added streaming throughput boost with IWQ

OLM Usage Considerations(1): OSA Sharing

- Concurrent interfaces to an OSA-Express port using OLM is limited.
 - If one or more interfaces operate OLM on a given port,
 - Only four total interfaces allowed to that single port
 - Only eight total interfaces allowed to that CHPID
 - All four interfaces can operate in OLM
 - An interface can be:
 - Another interface (e.g. IPv6) defined for this OSA-Express port
 - Another stack on the same LPAR using the OSA-Express port
 - Another LPAR using the OSA-Express port
 - Another VLAN defined for this OSA-Express port
 - Any stack activating the OSA-Express Network Traffic Analyzer (OSAENTA)

OLM Usage Considerations (2):

- QDIO Accelerator or HiperSockets Accelerator will not accelerate traffic to or from an OSA-Express operating in OLM
- OLM usage may increase z/OS CPU consumption (due to “early interrupt”)
 - Usage of OLM is therefore not recommended on z/OS images expected to normally be running at extremely high utilization levels
 - OLM does not apply to the bulk-data input queue of an IWQ-mode OSA. From a CPU-consumption perspective, OLM is therefore a more attractive option when combined with IWQ than without IWQ
- Only supported on OSA-Express3 and above with the INTERFACE statement
- Enabled via PTFs for z/OS V1R11
 - PK90205 (PTF UK49041) and OA29634 (UA49172).