



IBM z/VM Development Lab – Endicott, NY

## z/VM Performance Update

Bill Bitner, [bitnerb@us.ibm.com](mailto:bitnerb@us.ibm.com)  
Brian Wade, [bkw@us.ibm.com](mailto:bkw@us.ibm.com)

# Trademarks

## Trademarks

The following are trademarks of the International Business Machines Corporation in the United States and/or other countries. For a complete list of IBM Trademarks, see [www.ibm.com/legal/copytrade.shtml](http://www.ibm.com/legal/copytrade.shtml): AS/400, DBE, e-business logo, ESCO, eServer, FICON, IBM, IBM Logo, iSeries, MVS, OS/390, pSeries, RS/6000, S/30, VM/ESA, VSE/ESA, Websphere, xSeries, z/OS, zSeries, z/VM

The following are trademarks or registered trademarks of other companies

Lotus, Notes, and Domino are trademarks or registered trademarks of Lotus Development Corporation  
Java and all Java-related trademarks and logos are trademarks of Sun Microsystems, Inc., in the United States and other countries  
Linux is a registered trademark of Linus Torvalds  
UNIX is a registered trademark of The Open Group in the United States and other countries.  
Microsoft, Windows and Windows NT are registered trademarks of Microsoft Corporation.  
SET and Secure Electronic Transaction are trademarks owned by SET Secure Electronic Transaction LLC.  
Intel is a registered trademark of Intel Corporation  
\* All other products may be trademarks or registered trademarks of their respective 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.

References in this document to IBM products or services do not imply that IBM intends to make them available in every country.

Any proposed use of claims in this presentation outside of the United States must be reviewed by local IBM country counsel prior to such use.

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

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

Permission is hereby granted to SHARE to publish an exact copy of this paper in the SHARE proceedings. IBM retains the title to the copyright in this paper, as well as the copyright in all underlying works. IBM retains the right to make derivative works and to republish and distribute this paper to whomever it chooses in any way it chooses.

## Notice Regarding Specialty Engines (e.g., zIIPs, zAAPs and IFLs):

Any information contained in this document regarding Specialty Engines ("SEs") and SE eligible workloads 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). 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](http://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 SEs 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.

## Acknowledgements – Your z/VM Performance Team

- **Dean DiTommaso**
- **Bill Guzior**
- **Steve Jones**
- **Virg Meredith**
- **Patty Rando**
- **Dave Spencer**
- **Susan Timashenka – Dept Manager**
- **Xenia Tkatschow**
- **Brian Wade**

# Agenda

- **z/VM 6.2 thoughts**
  - LGR and SSI
    - Performance notes
    - Management and monitoring thoughts
  - Various other line items
  - Monitor record changes
  - Performance-related service
- **z/VM 6.3 Preview**

## z/VM 6.2 Highlights – A Performance View

- **Regression performance**
- **SSI and LGR considerations**
- **Memory management improvements**
- **MONDCSS and SAMPLE CONFIG increases**
- **STORBUF changes**
- **z/CMS and implications**
- **CPU Measurement Facility exploitation**
- **Monitor records**
- **z/VM Performance Toolkit changes**

## z/VM 6.2 Regression Performance

- **Ran our standard library of workloads**
  - CMS interactive, various Apache configurations
- **Results are within usual 5% regression criteria**
- **Some workloads will see improvements:**
  - Overprovisioned for logical PUs compared to utilization
  - Storage-constrained with heavy contention for <2 GB real storage
  - High virtual CPU to logical CPU overcommit with virtual CPUs often in a ready-to-run state



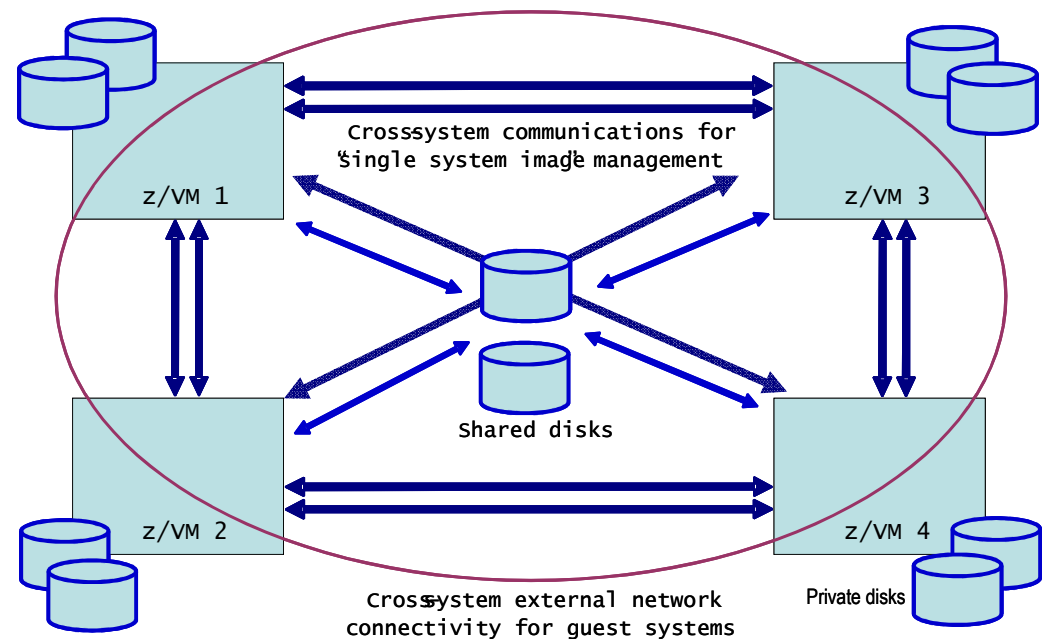
IBM z/VM Development Lab – Endicott, NY

## SSI and LGR Thoughts

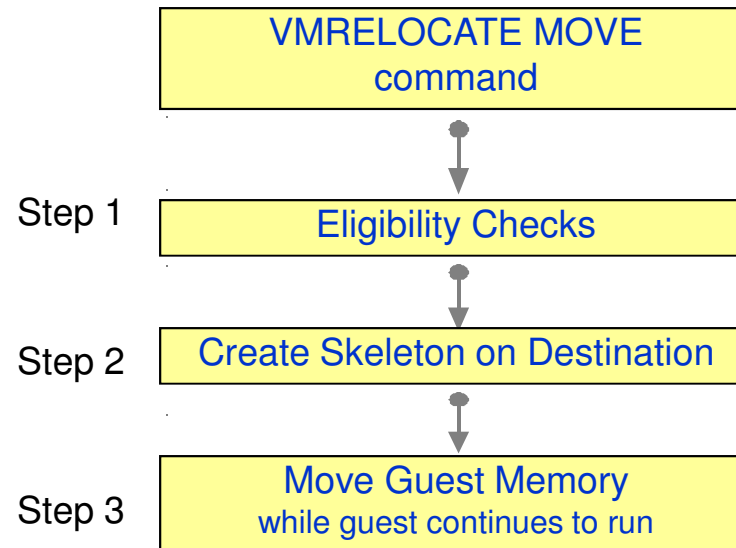


# Single System Image Feature Clustered Hypervisor with Live Guest Relocation

- Provided as an optional priced feature.
- Connect up to four z/VM systems as members of a Single System Image (SSI) cluster
- Provides a set of shared resources for member systems and their hosted virtual machines
- Cluster members can be run on the same or different System z servers
- Simplifies systems management of a multi-z/VM environment
  - Single user directory
  - Cluster management from any member
    - Apply maintenance to all members in the cluster from one location
    - Issue commands from one member to operate on another
  - Built-in cross-member capabilities
  - Resource coordination and protection of network and disks

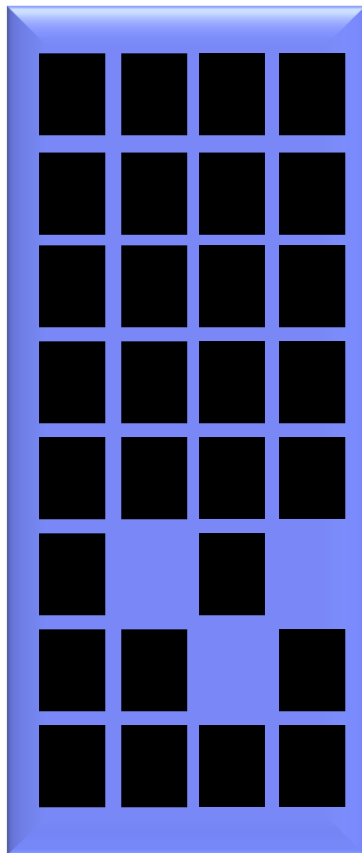


# Stages of a Live Guest Relocation



# LGR, High-Level View of Memory Move

Source



Guest Address Space

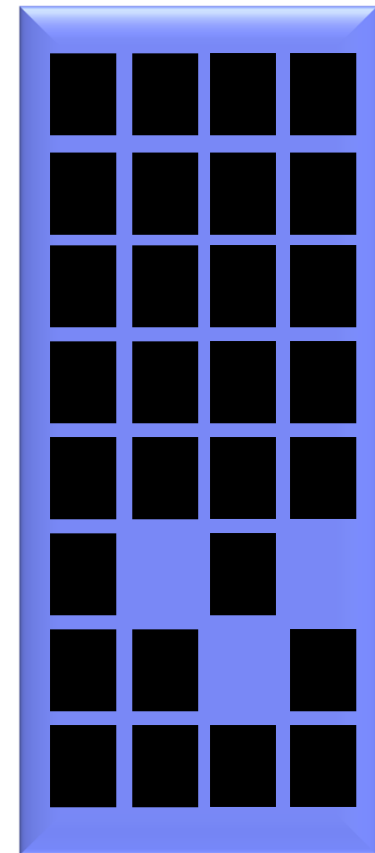
## PUSH with resend

Pass 1



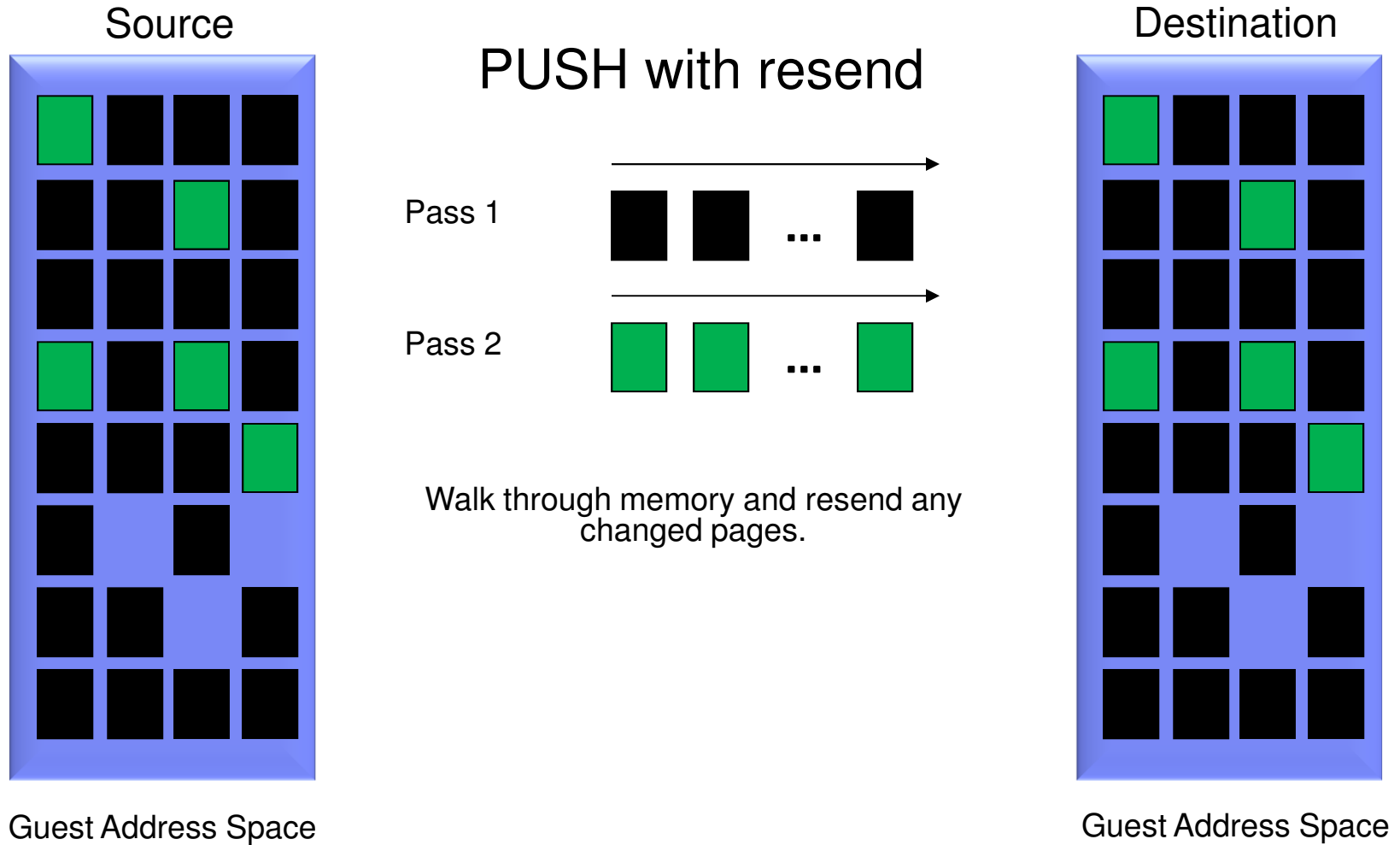
Walk through guest memory moving all non-zero pages

Destination

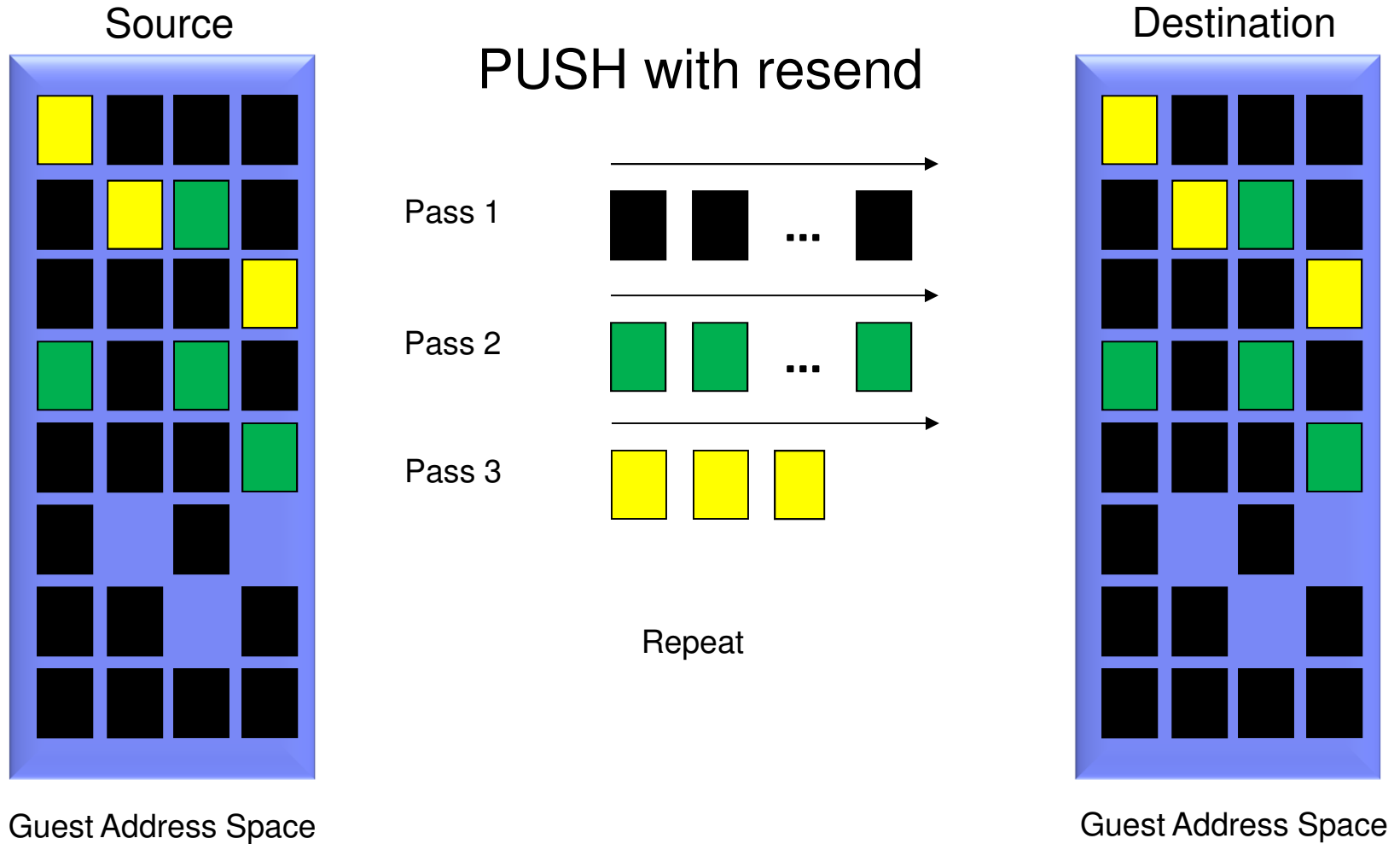


Guest Address Space

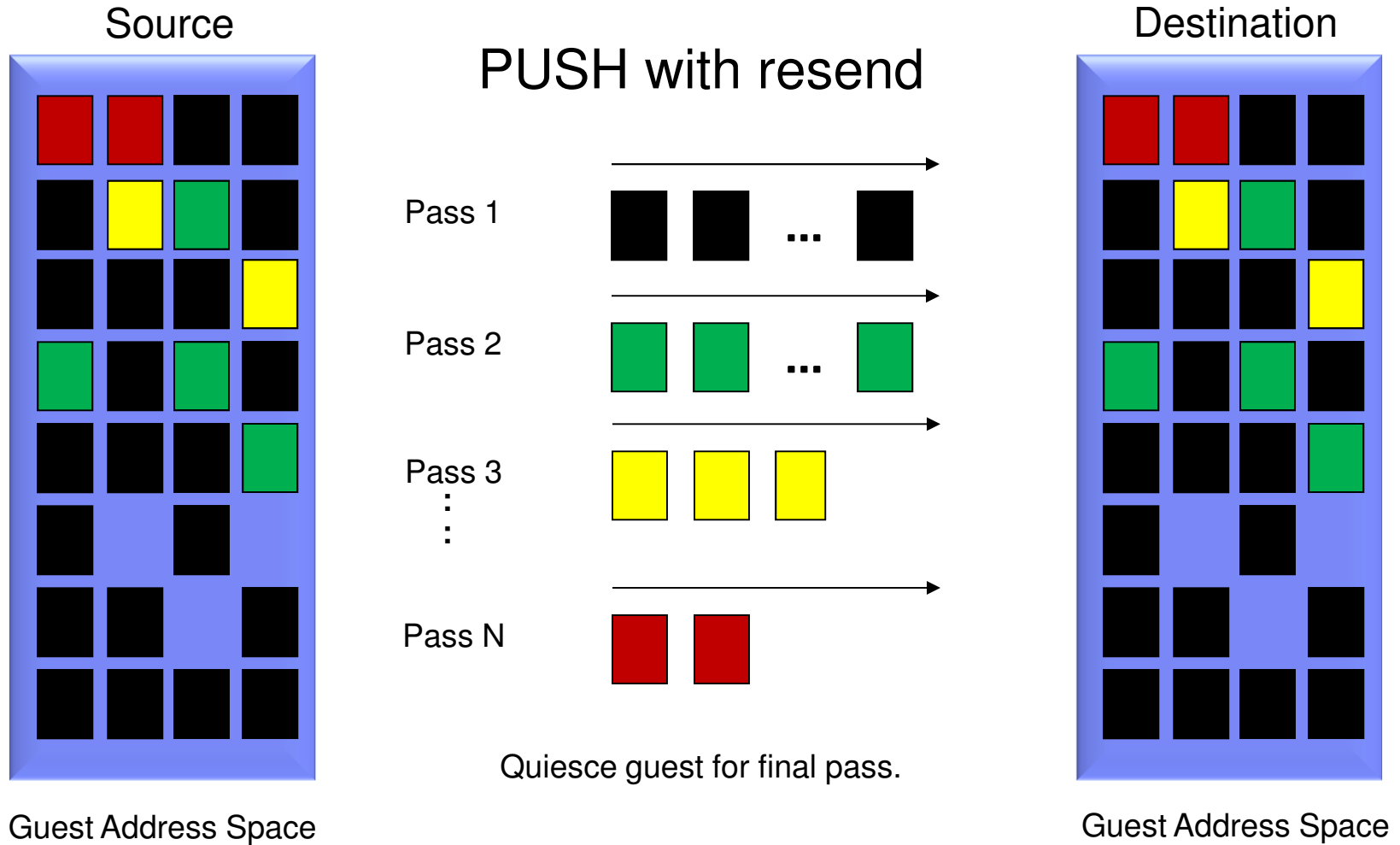
# LGR, High-Level View of Memory Move



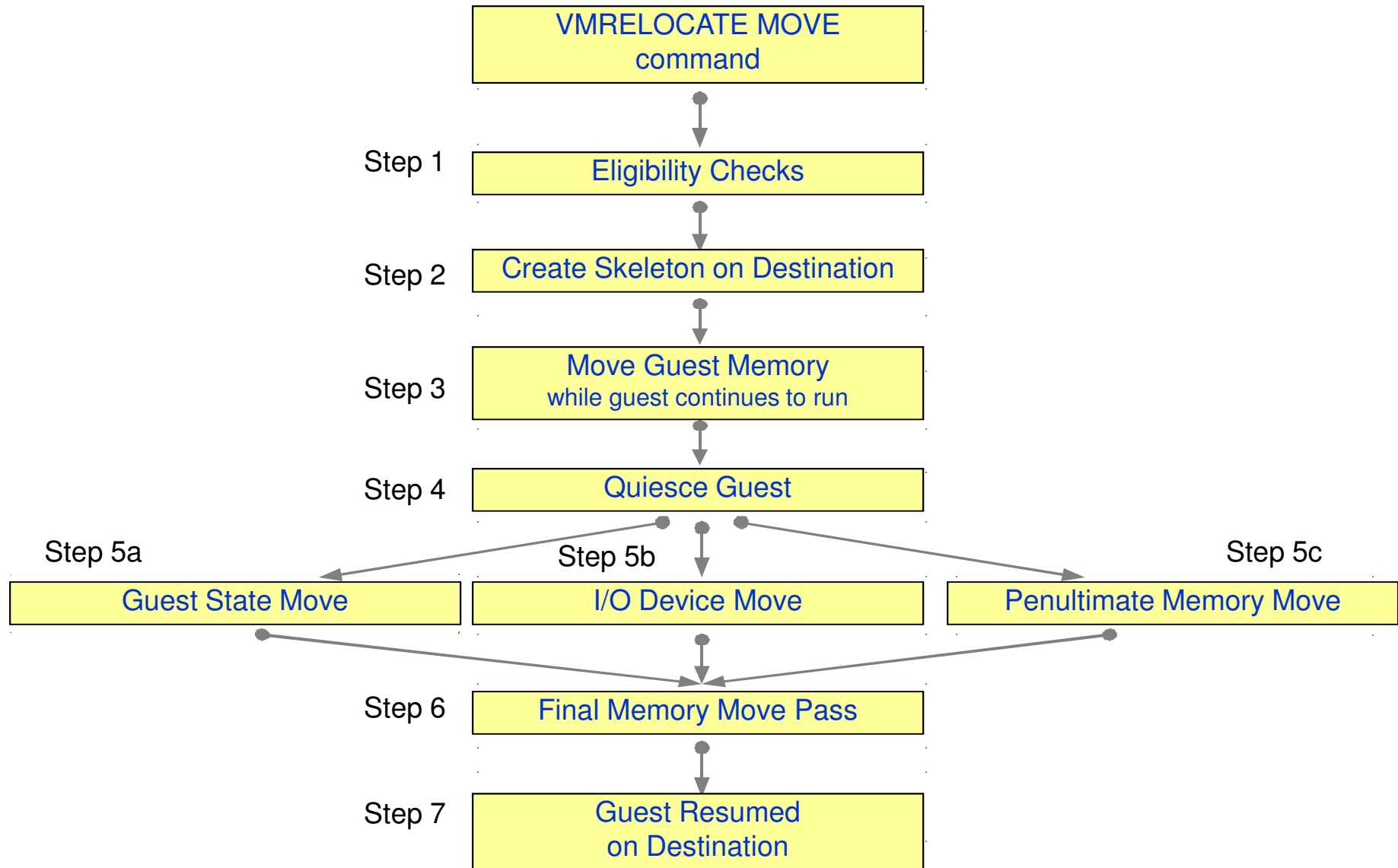
# LGR, High-Level View of Memory Move



# LGR, High-Level View of Memory Move



# Stages of a Live Guest Relocation



## Live Guest Relocation – Key Performance Metrics

### ■ Quiesce Time (QT)

- Elapsed time that the guest is stopped (stunned) so z/VM can move the guest's last set of storage pages – probably the frequently-changed ones
- To tolerate relocation, the guest and its applications must tolerate the quiesce time
- VMRELOCATE can be invoked with a specified maximum quiesce time
  - If the quiesce would run past the maximum, z/VM cancels the relocation

### ■ Relocation Time (RT)

- Elapsed time from when the VMRELOCATE command is issued to when the guest is successfully restarted on the destination system.
- Elapsed time must fit within the customer's window of time for planned outages for system maintenance, etc.

Bottom line: there are some scenarios where LGR is not feasible as a result of the requirements for relocation time and quiesce time



## LGR: Factors Affecting QT and RT

- **Size of the guest**
  - Amount of memory to move, time required to walk its DAT tables
- **How broadly or frequently the guest changes its pages**
  - It's an iterative memory push from source to destination
- **Time needed to relocate the guest's I/O configuration**
  - I/O device count, I/Os to quiesce, OSA recovery on target side
- **Capacity of the ISFC logical link**
  - Number of chpids, their speeds, number of RDEVs
- **Storage constraints on source and target systems**
- **Performance of paging subsystem**
- **Other work the systems are doing**
- **Other relocations happening concurrently with the one of interest**
- **Delays injected when LGR throttles itself back to prevent abends and other problems.**
  - End-to-end LGR throttling – triggered by paging intensities
  - Memory-move endpoint throttling – triggered by memory consumption
  - ISFC logical link throttling – triggered by ISFC running out of queued traffic buffers

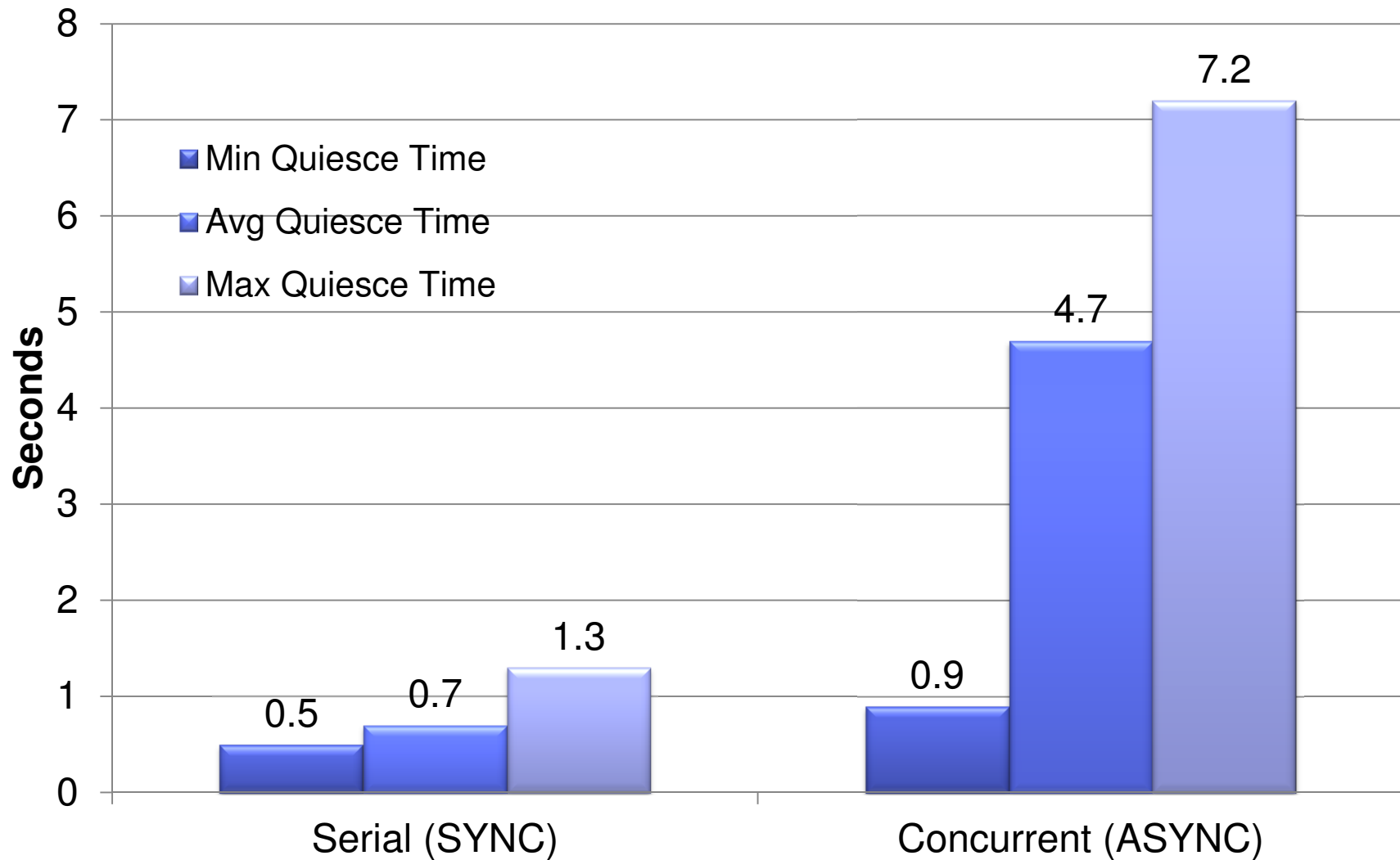
## LGR: Serial vs. Concurrent Relocations

- **By default, the VMRELOCATE command operates synchronously.**
- **There is a command option (ASYNCH) to run it asynchronously (a la SPXTAPE)**
- **You could also achieve concurrent relocations by:**
  - Use the asynchronous version of VMRELOCATE multiple times.
  - Run VMRELOCATE commands in multiple users concurrently.

The best practice, though, is to run only one relocation at a time.

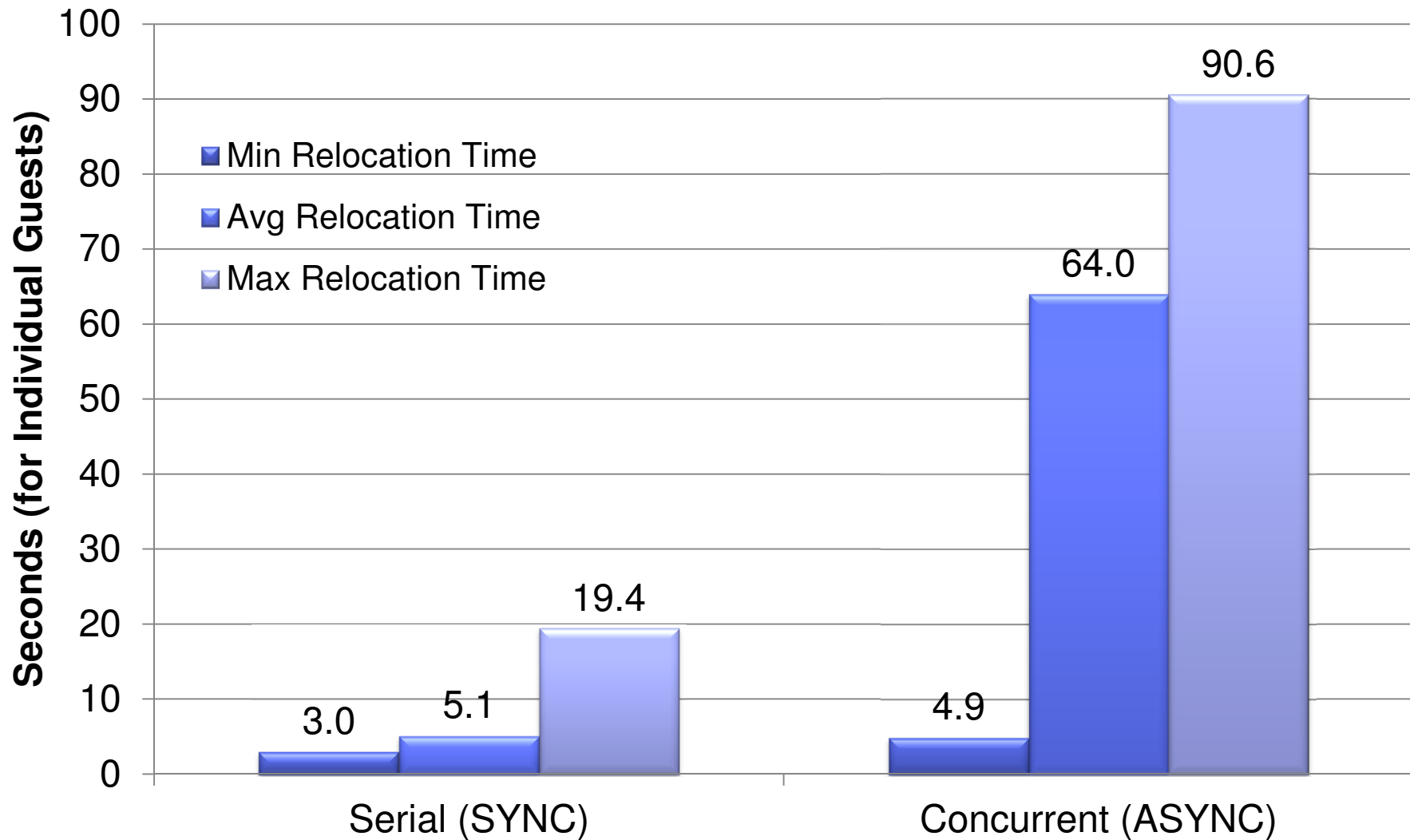
- **QT and individual RT improves substantially when relocations are done serially**
  - ... and total RT elongates only slightly

## Effect of Serial vs. Concurrent on Quiesce Time



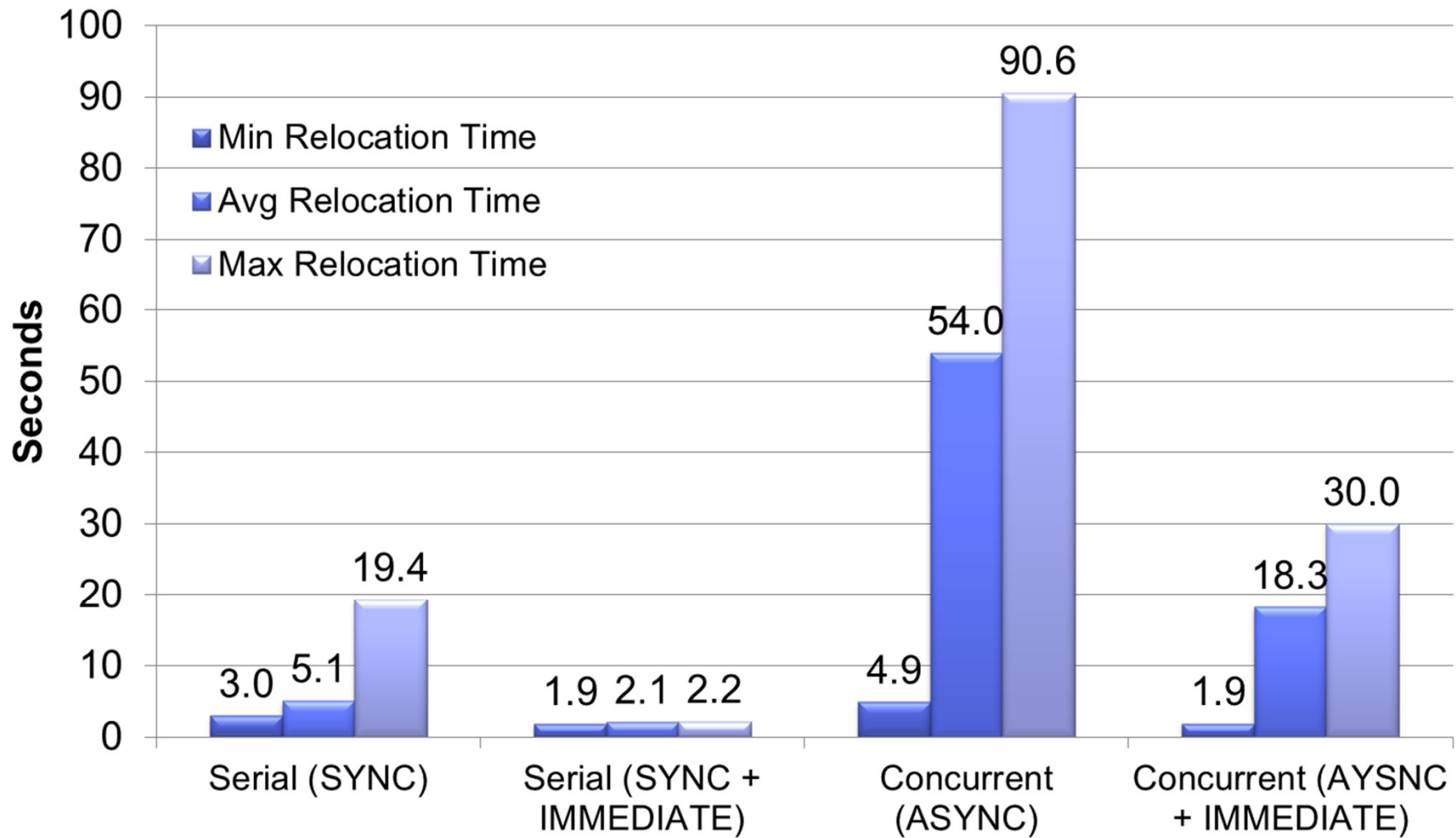
**Relocation Parameters – 25 4GB Linux Guests**

# Effect of Serial vs. Concurrent on Relocation Time



**Relocation Parameters – 25 4GB Linux Guests**

# Effect of IMMEDIATE option on Relocation Time



**Relocation Parameters – 25 4GB Linux Guests**

## VMRELOCATE Options Summary

- **Best total relocation time for all virtual machines**
  - Concurrent (ASYNCH) + IMMEDIATE
- **Best individual relocation time**
  - Serial (SYNCH) + IMMEDIATE
- **Best quiesce times**
  - Serial (SYNCH)
- **Worst quiesce times**
  - Concurrent (ASYNCH) + IMMEDIATE

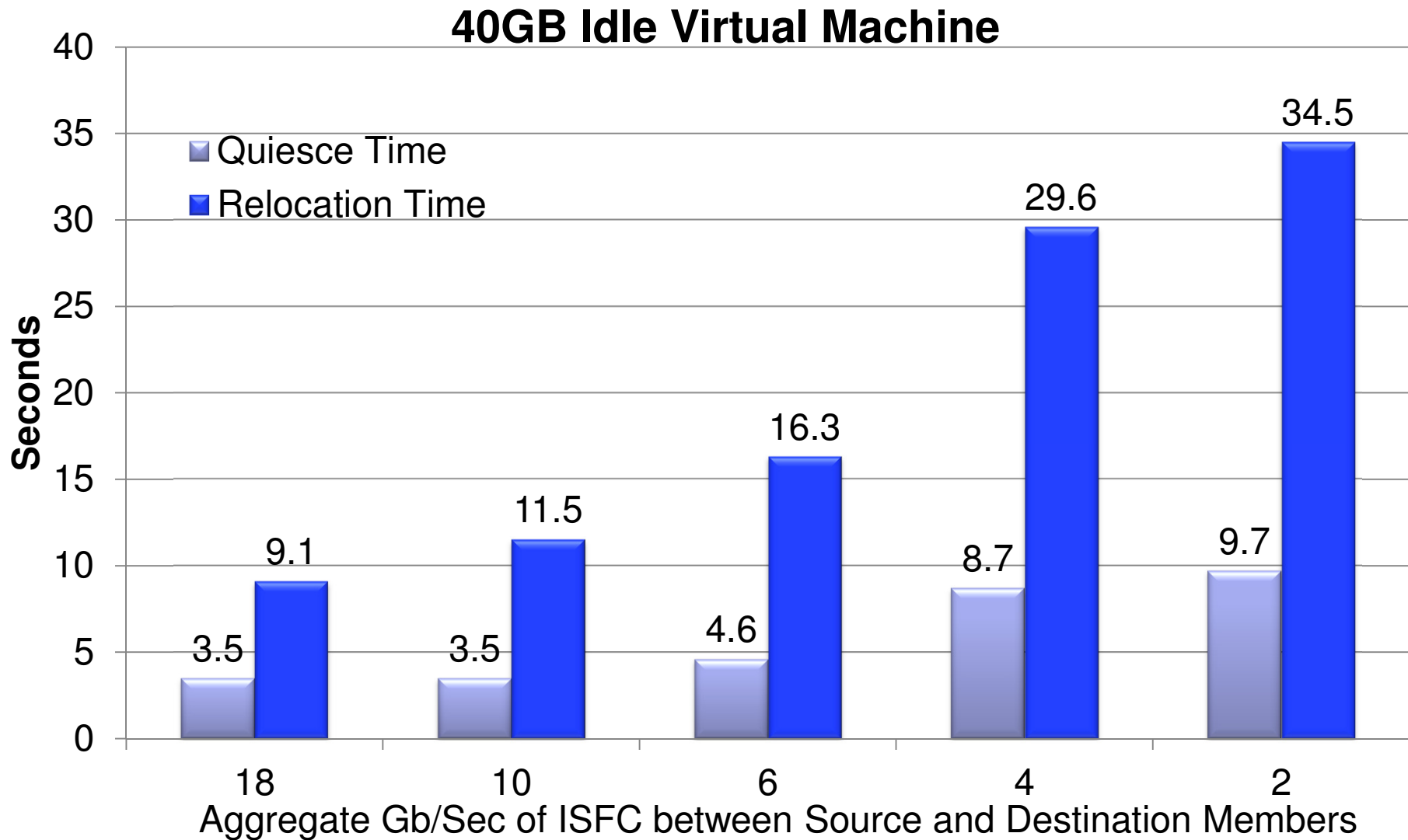
## Background on ISFC Capacity Test

**Table 3. Evaluated ISFC Logical Link Configurations.**

ISFC Logical Link CHPIDs	ISFC Capacity Factor *	CTCs/FICON CHPID	Total CTCs
1-2Gb, 2-4Gb, 1-8Gb	18	4	16
1-2Gb, 2-4Gb	10	4	12
1-2Gb, 1-4Gb	6	4	8
1-4Gb	4	4	4
1-2Gb	2	4	4

**Note:** \* ISFC capacity factor is the sum of speeds of the FICON CTCs between the SSI member systems.

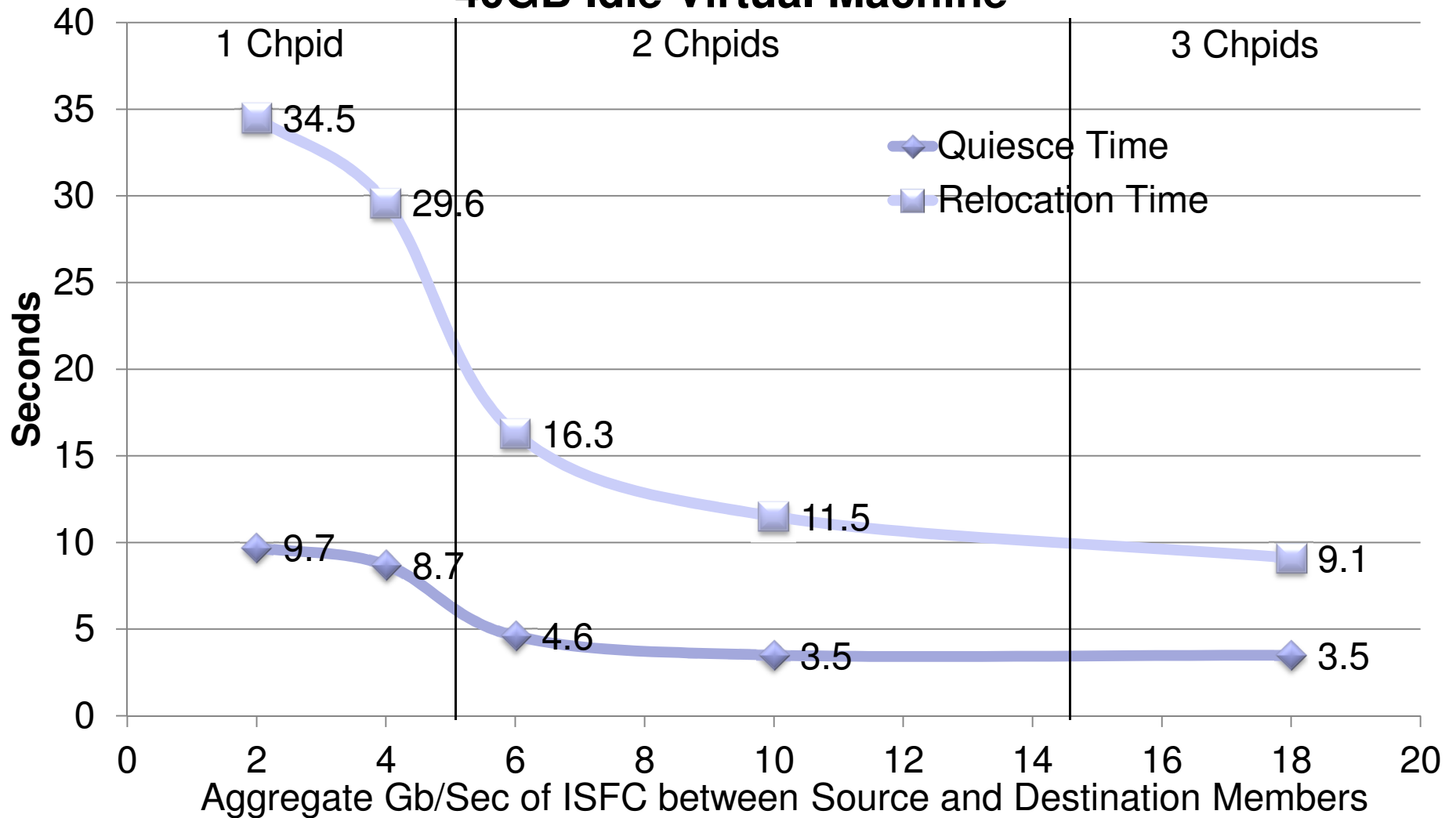
# Effect of CTC Bandwidth on LGR



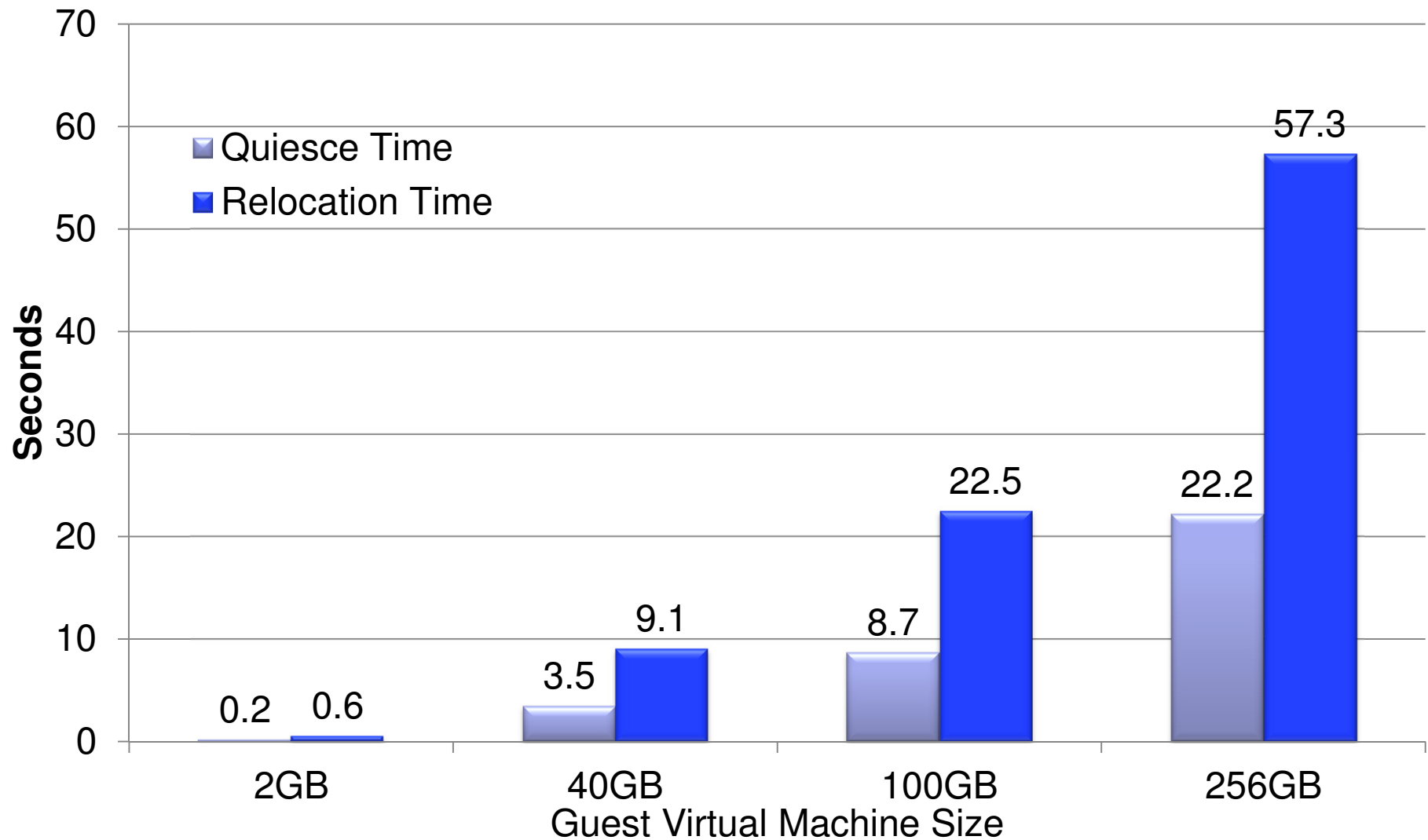


# Effect of CTC Bandwidth on LGR

## 40GB Idle Virtual Machine

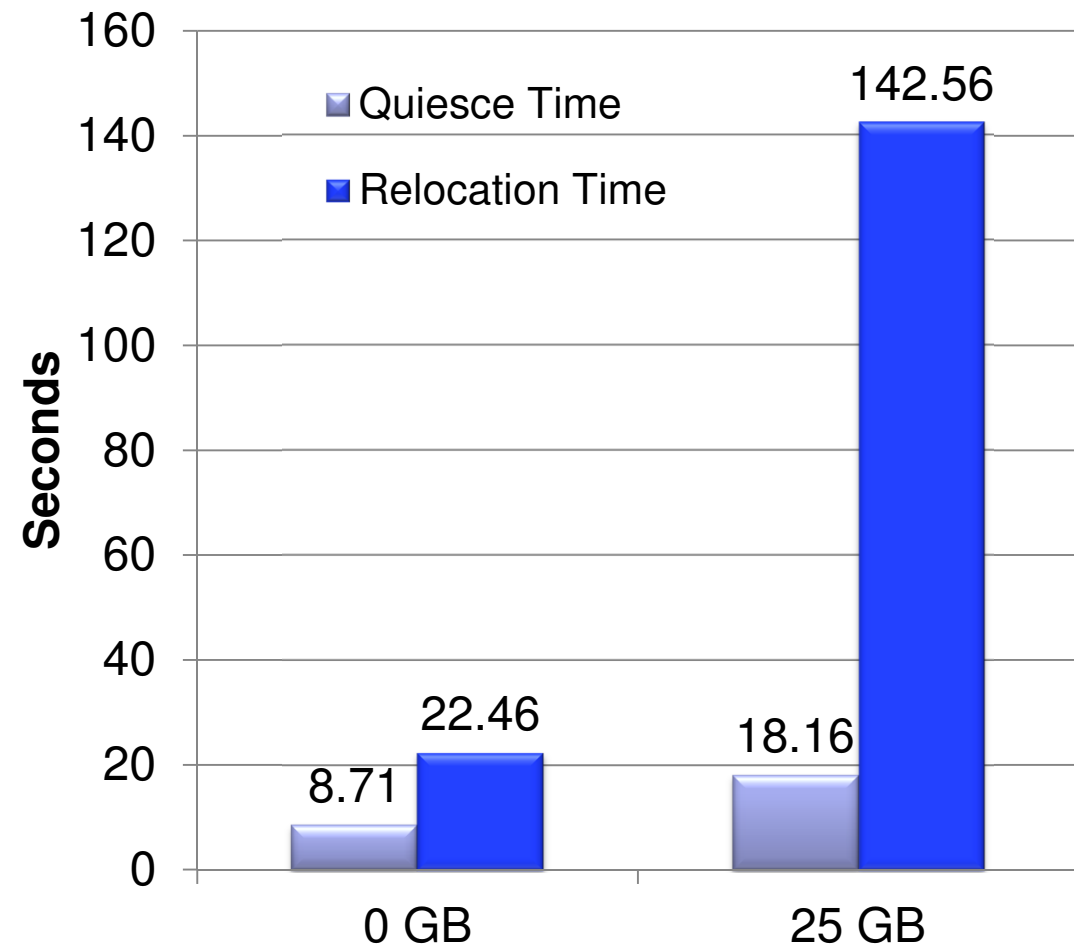


## Effect of Virtual Machine Size on LGR



# Impact of Virtual Machine Changing Memory on LGR

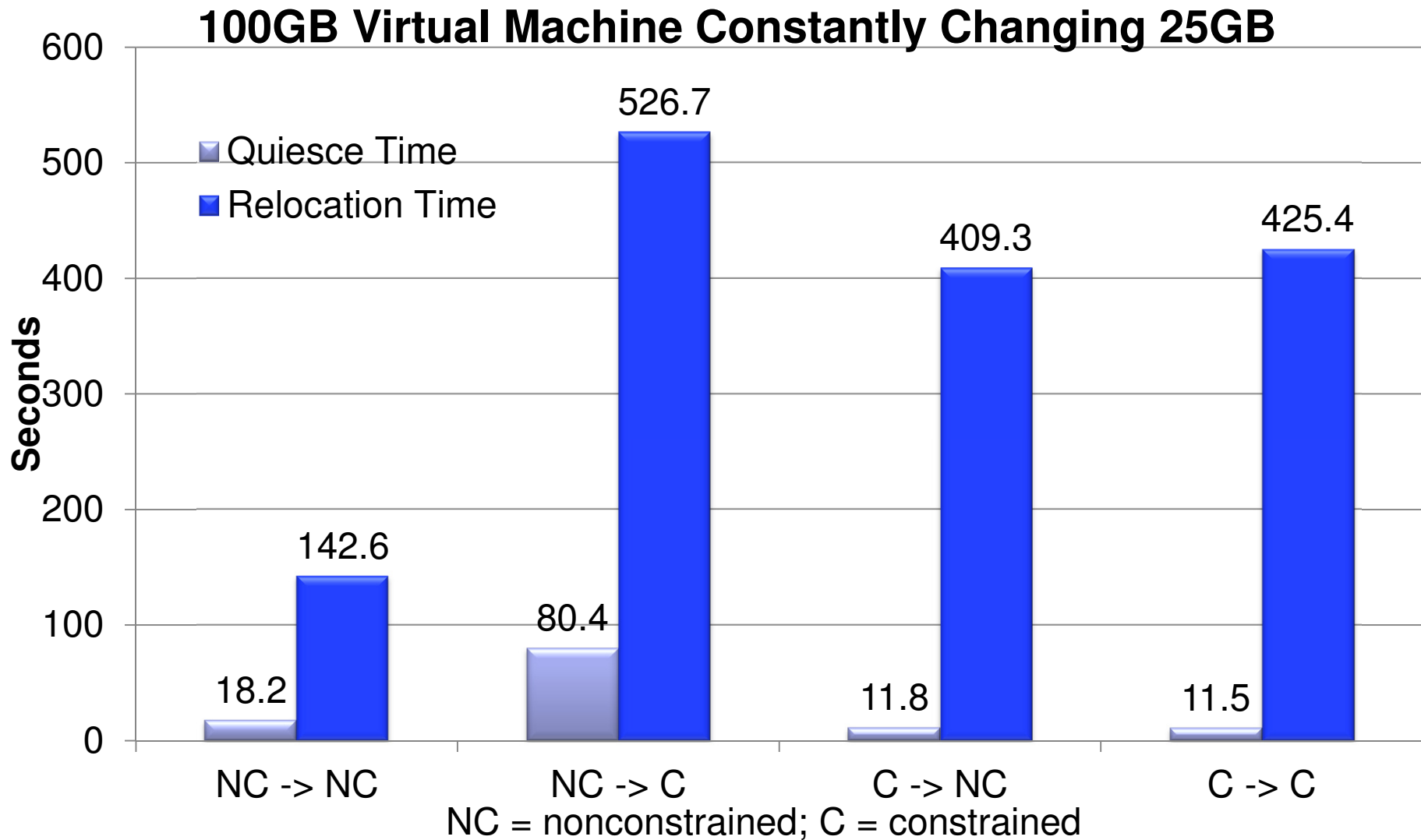
- **Idle case (0GB changing) there is less memory to move and fewer Memory Move Passes**
- **Number of Passes**
  - 0GB: 4
  - 25GB: 8
- **Total Memory Moved**
  - 0GB: 4.9GB
  - 25GB: 160GB



## LGR: CPU and Memory Use Habits

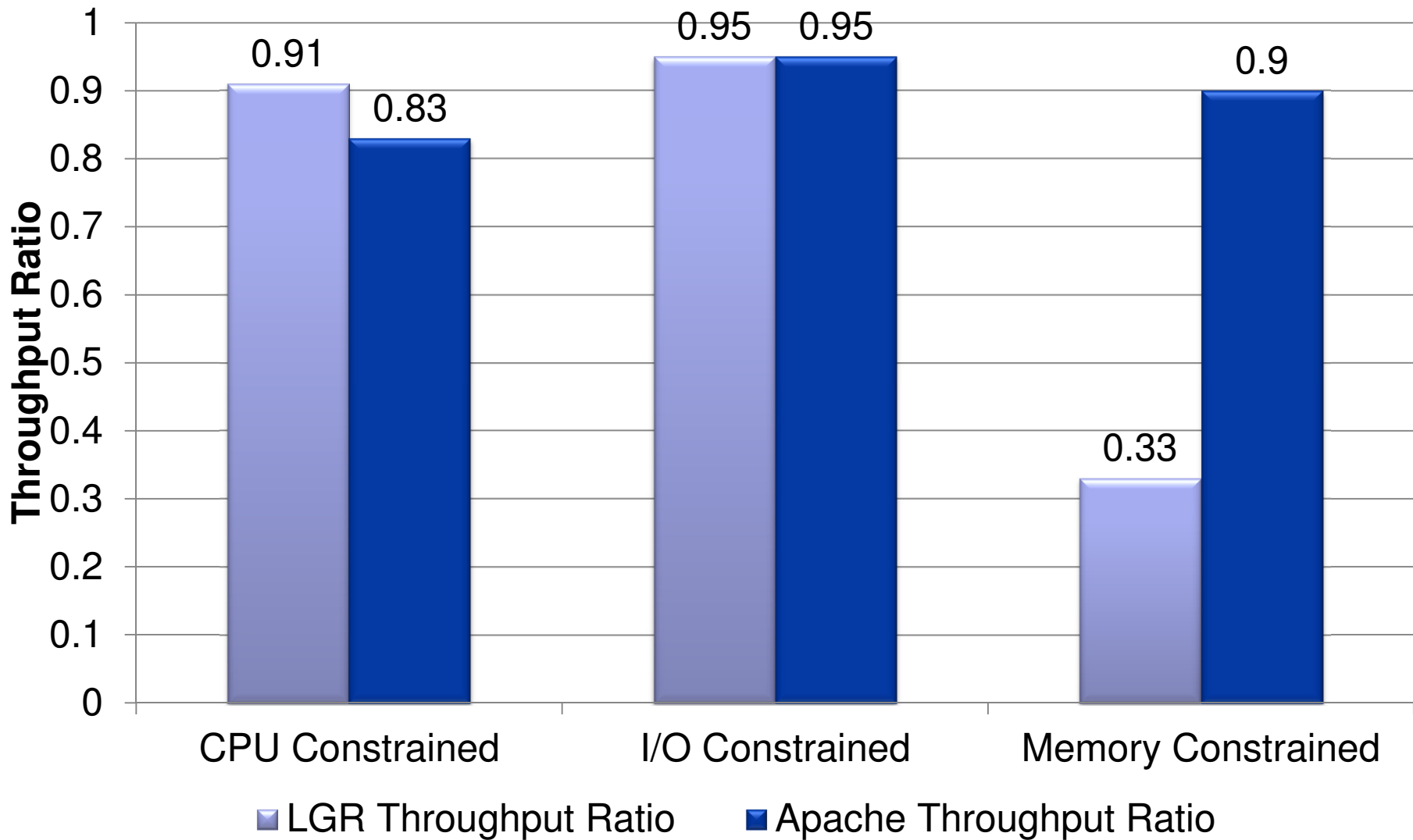
- **CPU: generally LGR gets what it needs**
  - Taken “off the top” compared to your workload
- **Memory: CP tries really hard not to interfere**
  - End-to-end throttling, ISFC buffer limits, ...
  - Socket memory-move throttling – triggered by memory consumption
  - ISFC logical link throttling – triggered by ISFC running out of queued traffic buffers
  - Considers effect on paging, memory use for specific relocations, ...

# Effect of System Memory Constraint on LGR



# Effect of LGR on Existing Workloads

LGR Bounce and Apache Web Serving Workloads



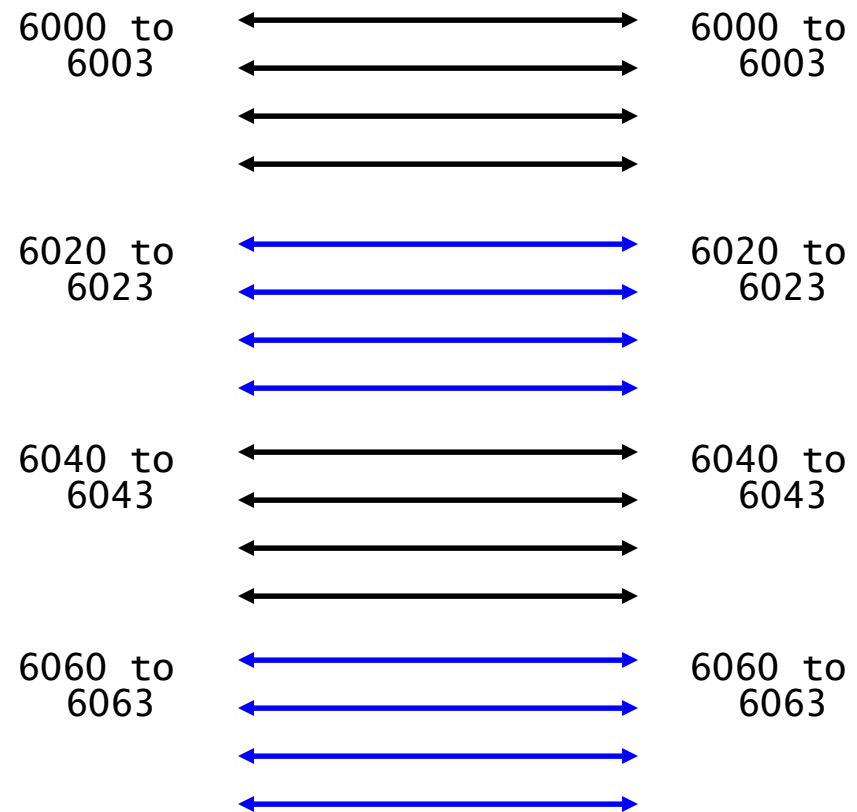
## LGR: Keep These in Mind...

- **Charge back:** can your procedures handle guests that suddenly disappear and then reappear somewhere else?
- **Second-level schedulers:** do you have them? Can they handle guest motion?
- **VMRM:** if VMRM-A tweaks the guest and then the guest moves to system B, what happens? And then what happens when the guest comes back?

Best practice is not to include relocating guests in VMRM-managed groups.

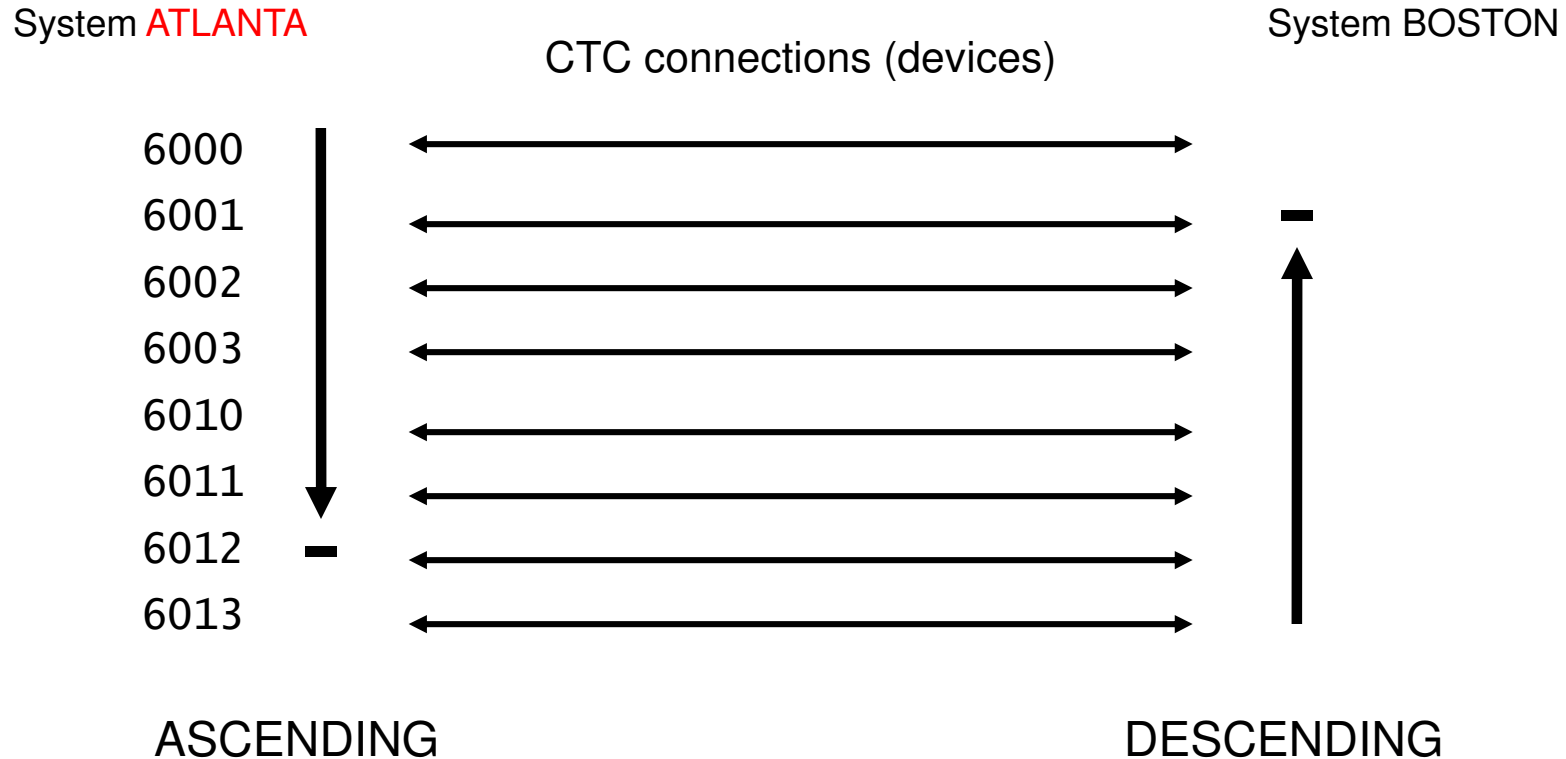
## SSI: ISFC Logical Link Configuration Best Practices

- **Use multiple FICON chpids of all the same speed. Up to 4 chpids.**
- **Use four CTC devices per chpid**
- **Use same RDEV numbers on both ends**
- **More esoteric configurations are certainly possible**
- **Can share the chpids but requires capacity planning**





## SSI: ISFC Logical Link Write Scheduling, under the covers



Moral: put the fast chpids in the middle of ATLANTA's RDEV range.  
 Selection of where to start in selecting write path is alphabetical.

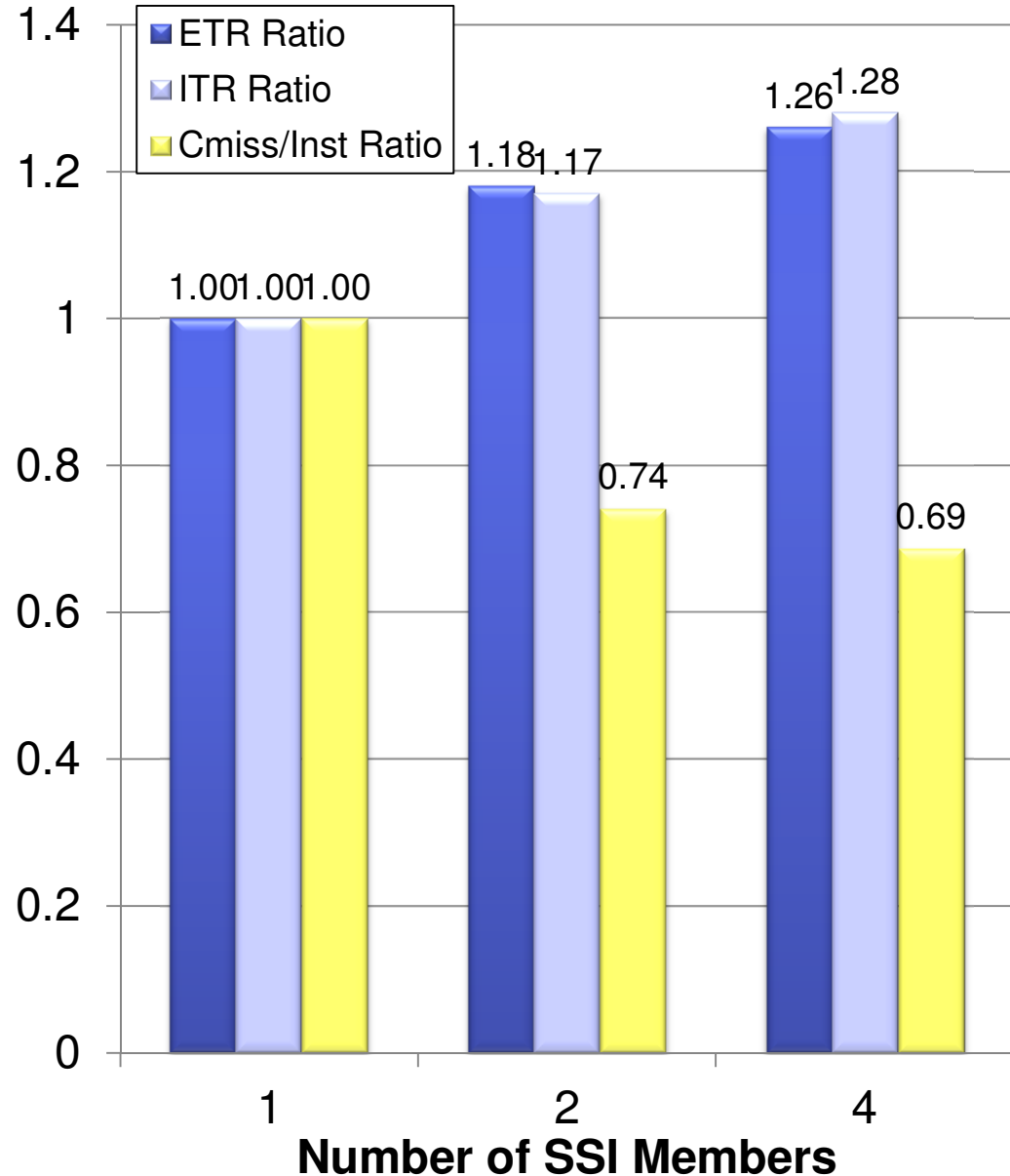
## SSI Workload Distribution Measurements

Parameters	1 Member	2 Member	4 Member
Central Storage	43 GB	22 GB	11 GB
Expanded Storage	8 GB	4 GB	2 GB
Processors	12	6	3

- Series of measurements to see how a workload spread across a number of members would run compared to one larger systems of just one member.
- Resources kept the same, as shown above.
- Apache workload where clients and servers were all virtual machines was used.
  - Varied number of client and servers and use of MDC to create different stress points.

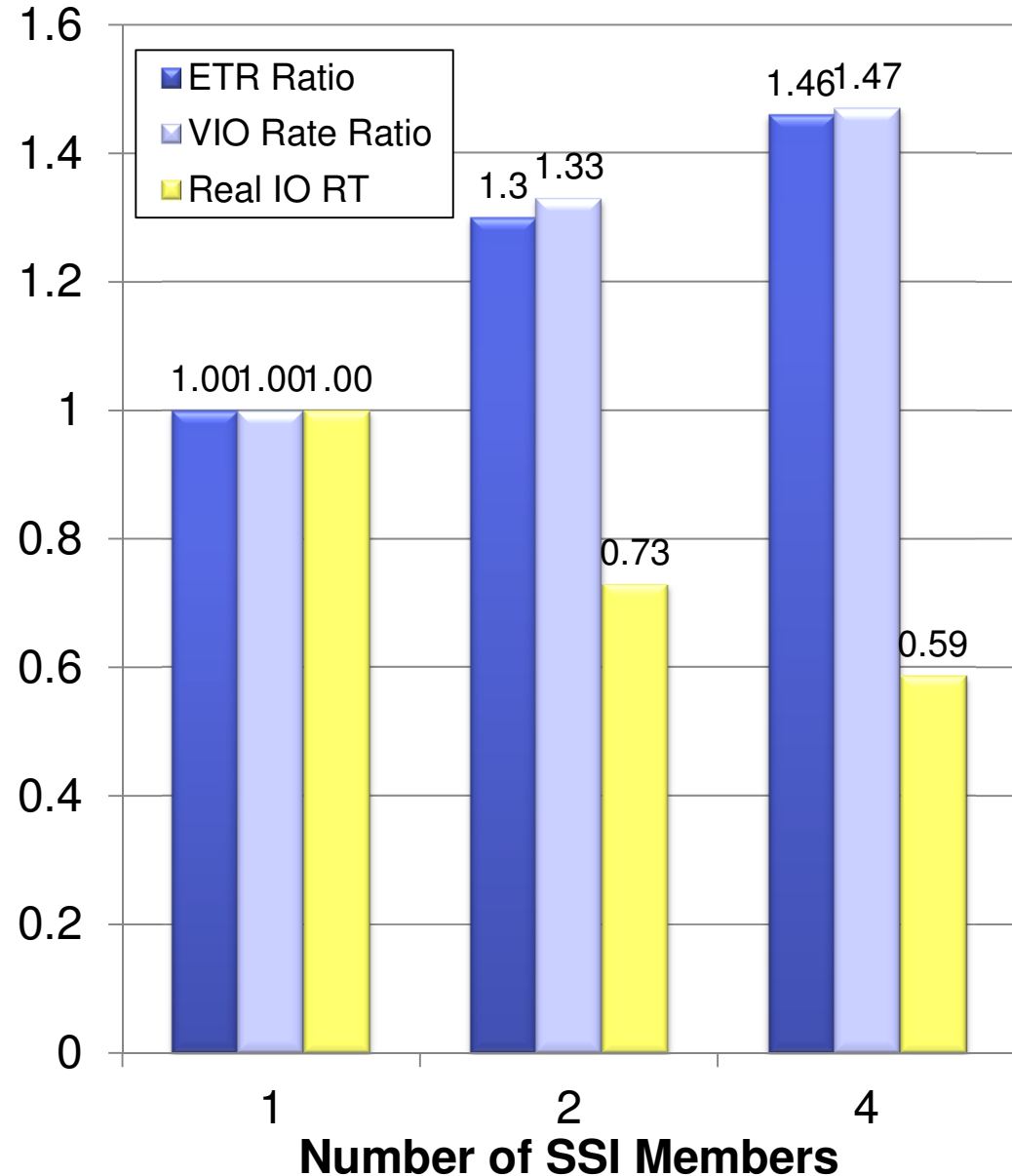
## SSI Distribution: CPU Constrained Measurement

- **Keep the physical resources the same, but distribute over 1, 2, or 4 members.**
- **Apache Web Serving with the configuration being CPU bound.**
- **Benefits from running smaller n-way partitions**



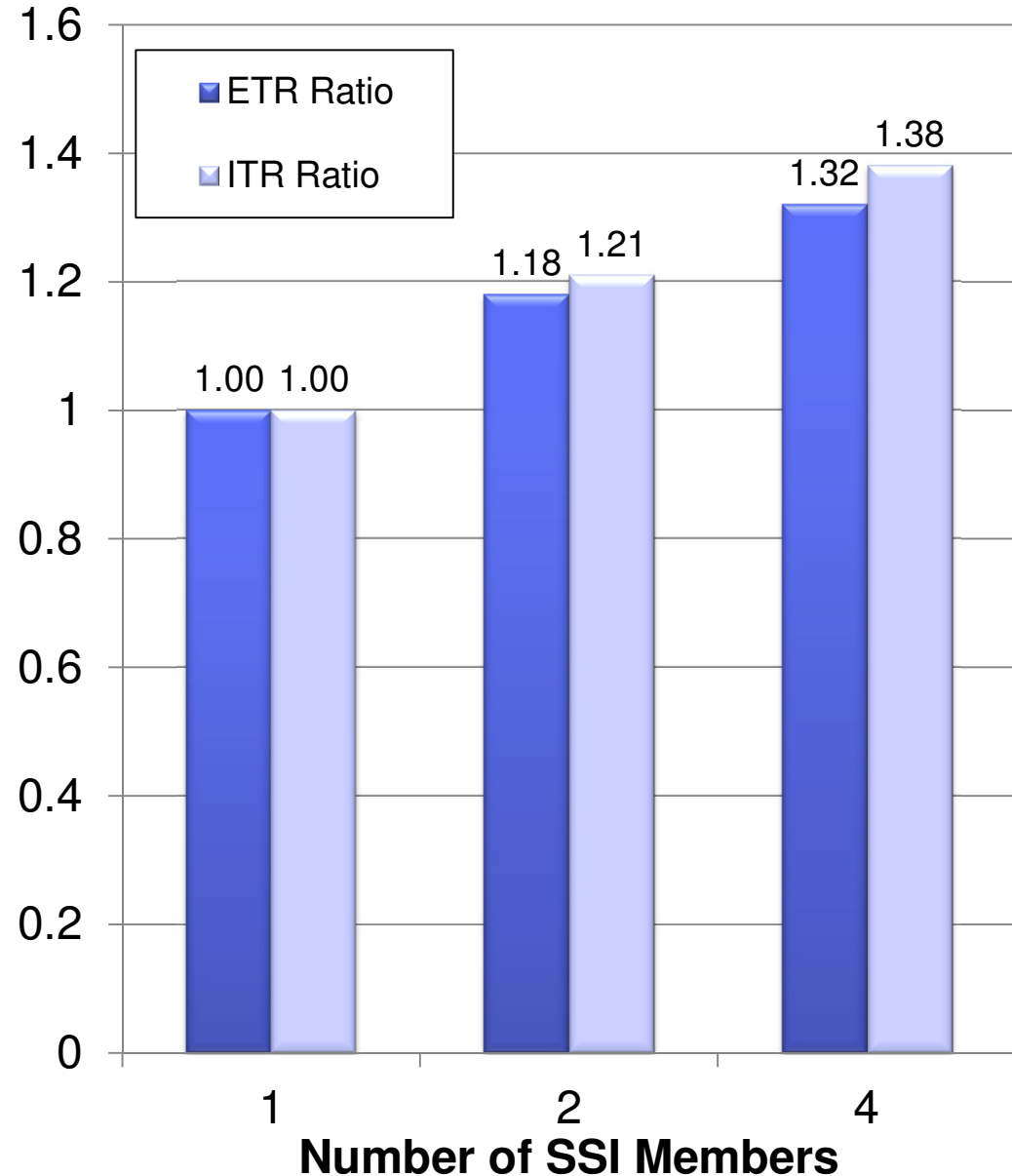
## SSI Distribution: Virtual I/O Constrained Measurement

- Keep the physical resources the same, but distribute over 1, 2, or 4 members.
- Apache Web Serving with the configuration being I/O bound due to virtual read I/O.
- PAV not used in base case, so SSI essentially gives PAV like benefits.
- Real I/O RT shown is for one of the shared Linux volumes containing files being served.



## SSI Distribution: Memory Constrained Measurement

- Keep the physical resources the same, but distribute over 1, 2, or 4 members.
- Apache Web Serving with the configuration with there being memory constraint.
- Similar savings as in CPU bound measurement.
- Additional efficiencies in memory management.



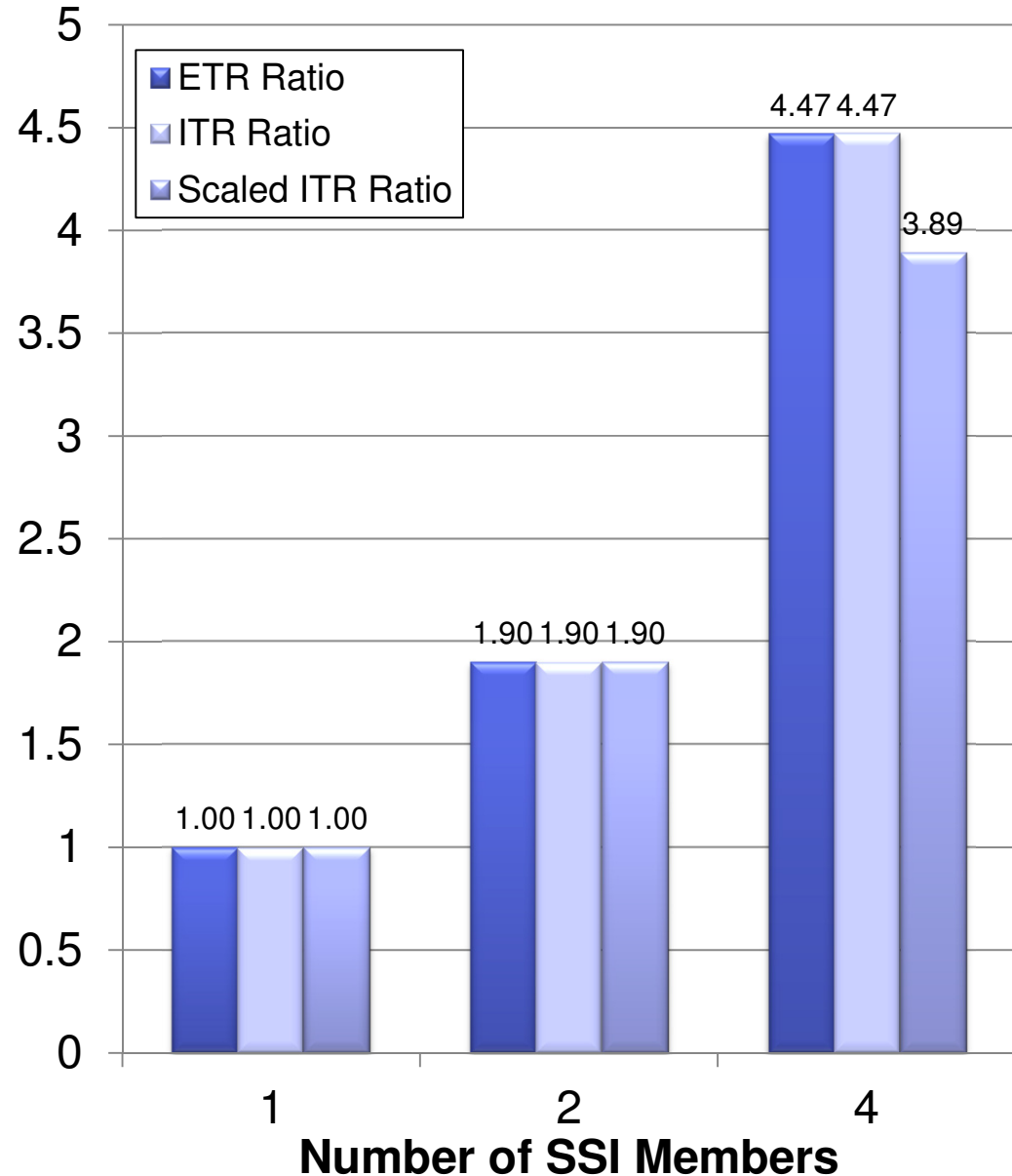
## SSI Workload Scaling Measurements

z/VM Limits	1 Member	2 Member	4 Member
Central Storage	256 GB	512 GB	1 TB
IFLs	32	64	128

- Measurements were made to see how well z/VM scales within an SSI cluster.
- Resources increased with each new member added to configuration.
- Apache workload where clients and servers were all virtual machines was used.
  - Apache clients and servers scaled accordingly.
- Needed to mix processor types to get 128 IFLs, so 1 & 2 Member runs are z10, 4 member adds in z196.
- Scaled down memory to make runs more feasible.

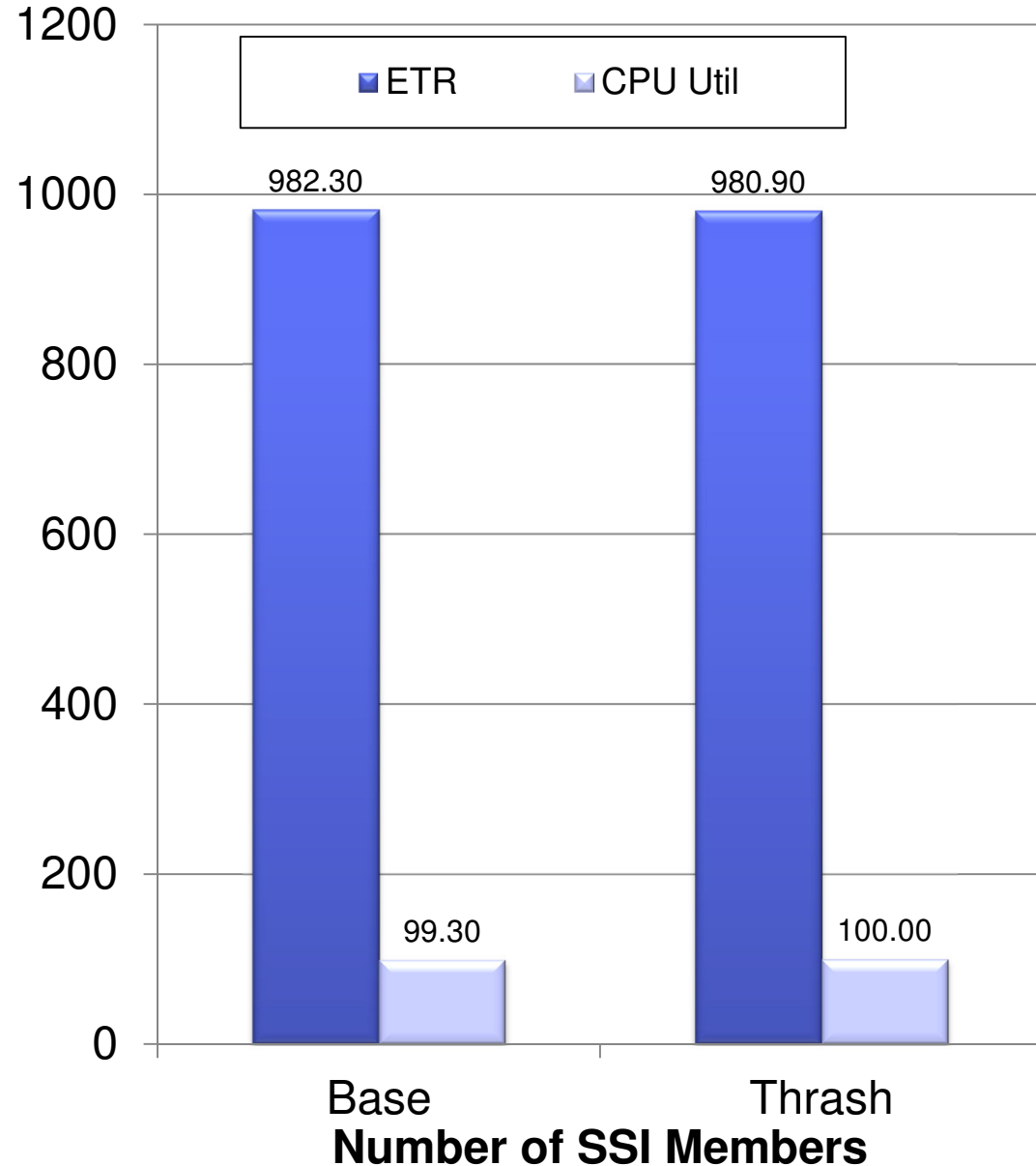
## SSI Scaling Measurements

- The SSI Cluster overhead for a running environment is very low.
- Note: z196s were added to get the 3<sup>rd</sup> and 4<sup>th</sup> Member.
- “Scaled ITR Ratio is an estimate of the Ratio if the entire cluster were on z10 processors.



## SSI Transition Measurement

- **Measurement to determine if activity or Cluster management would influence performance.**
- **Four Member environment where 3 of the members are constantly transitioning through states:**
  - Joined
  - Leaving
  - Down
  - Joining
  - *repeat*





## SSI: Performance Toolkit, Considerations

- **Performance Toolkit continues to run separately on each member of the cluster**
  - There continues to be a unique z/VM monitor data stream for each member.
  - There will be a PERFSVM virtual machine on each member
- **Configuration and usage**
  - Configure so that you will log onto or connect to a different PERFSVM on each system.
  - Configure Performance Toolkit to use the Remote Performance Monitoring Facility, which allows local and remote performance monitoring from a single screen.
- **In general, Performance Toolkit does not produce “cluster view” reports**
  - DASD device-busy view, for example

## SSI: Performance Toolkit, New Reports

- **New Reports for SSI**

- SSICONF: SSI configuration
- SSISCHLG: SSI state change synchronization activity log
- SSISMILG: SSI state/mode information log

- **New ISFC reports related to SSI**

- ISFECONF: ISFC end point configuration
- ISFEACT: ISFC end point activity
- ISFLCONF: ISFC logical link configuration
- ISFLACT: ISFC logical link activity
- ISFLALOG: ISFC logical link activity log

## SSI: MONWRITE Considerations

- **IBM often asks you to run MONWRITE**
  - PMR diagnosis, for example
- **You should be running MONWRITE anyway**
- **You should now be running MONWRITE on every member of the cluster**
- **Make sure it's easy to go find the MONWRITE data for all members for a specified time interval**

## SSI: Dump and PMR Considerations

- **To solve your PMR,**
- **... IBM might need concurrently-taken dumps.**
  
- **Just be prepared:**
  - Know how to take a SNAPDUMP. Practice.
  - Know the effect of SNAPDUMP on your workload.
  - Know how to take a restart dump.



IBM z/VM Development Lab – Endicott, NY

## z/VM 6.2 – More Than Just SSI and LGR

## Memory Management: Needle-in-Haystack Searches

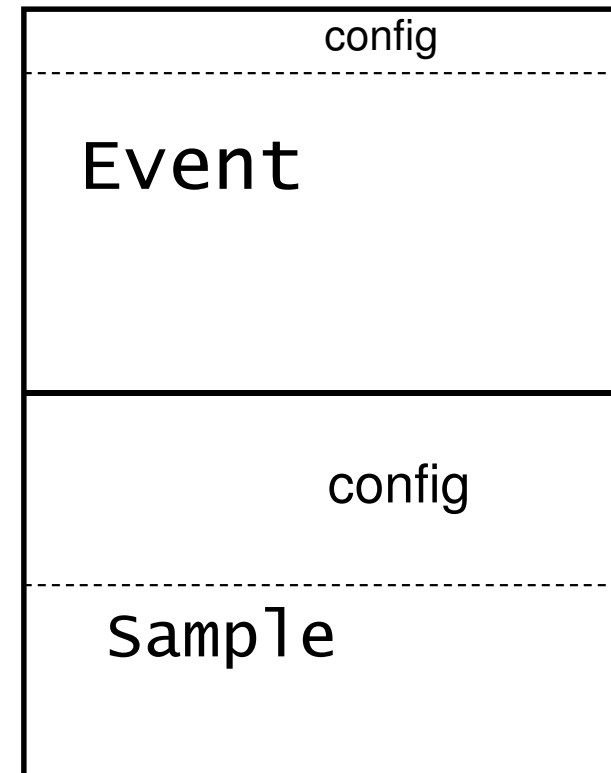
- **Searching for a below-2-GB frame in lists dominated by above-2-GB frames**
  - In months of study we identified about 10 of these searches
  - Development prototype that shut off all unnecessary use of <2GB storage gave us tremendous results
- **z/VM now does not allocate pageable buffers <2GB if:**
  - Dynamically, usable >2GB to usable <2GB is beyond a certain threshold
  - Statically, if the partition is beyond a certain size, for the life of the IPL
- **Result: no more needle searches**
- **Practically speaking, systems with 128 GB or more of real memory use below-2-GB memory only when it is architecturally required.**

## MONDCSS and SAMPLE CONFIG Changes

- The old defaults are too small for most systems nowadays
- So we have changed the default layout
- **MONDCSS is 64 MB now (16384 pages)**
  - Half (32 MB) for EVENT
  - Half (32 MB) for SAMPLE
    - Half (16 MB) for SAMPLE CONFIG
- As before, empty pages are not instantiated
- Remember, config pages evaporate after a short time
- **MONWRITE 191 disk also increased to 300 cylinders.**

If you use your own MONDCSS, the new default SAMPLE CONFIG size may be too large, requiring you to set it manually or to change your MONDCSS.

MONDCSS – 16384 pages



## Default STORBUF Changes

- **Many parties were noticing that the old defaults of 125 105 95 were not appropriate for Linux workloads**
- **We considered several different proposals**
  - From IBM ATS
  - From vendors
  - From Redbooks
  - From customer data
- **After careful consideration by “top people” we came to 300 250 200 as new defaults**

- If you already override defaults, the only impact would be if you also use SET SRM STORBUF INITIAL at some point.
- For CMS-intensive workloads, the old defaults might be more appropriate, and you should validate the settings for these workloads when you migrate to z/VM 6.2



## z/CMS

- **Prior to z/VM 6.2, z/CMS was supplied as a sample.**
- **z/VM 6.2 supports z/CMS as an optional alternative to the standard CMS that runs in ESA and XC mode virtual machines and 31-bit addressing.**
- **z/CMS can run in a z/Architecture guest**
  - Allows programs to use z/Architecture instructions, including 64-bit addressing
- **Standard CMS function does not exploit memory above 2GB**
- **Remember that z/Architecture is not XC**
  - No VM Data Spaces
  - No SFS DIRCONTROL-in-data-space
  - No DB/2-for-VM data space use
- **The standard, usual, XC-mode CMS is still there**

## CPU Measurement Facility Counters

- **CPU MF counters are a System z hardware facility that characterizes the performance of the CPU and nest**
  - Instructions, cycles, cache misses, and other processor related information
- **Available on z10 EC/BC, z196, and z114**
- **The CPU MF counter values:**
  - Help IBM to understand how your workload stresses a CEC for future design
  - Help IBM to map your workload into the LSPR curves for better sizing results
  - Help IBM better understand your system when there is a processor performance related problem.
- **z/VM 6.2, 6.1, and 5.4 can all collect the CPU MF counters from the hardware**
  - z/VM 5.4 and 6.1: VM64961, UM33440 (5.4), UM33442 (6.1)
  - Counters are put in new z/VM monitor record
- **We want volunteers to send us MONWRITE data!**
  - Your contributions will help us to understand customer workloads!

## IBM Wants Your CPU MF Counter Data

- **Your data will help IBM to build a library of customer workloads**
- **Collect an hour's worth of MONWRITE data...**
  - From a peak period,
  - With CPU MF counters enabled,
  - With one-minute sample intervals
- **Contact Richard Lewis at [rflewis at us.ibm.com](mailto:rflewis@us.ibm.com)**
- **Richard will send you instructions on how to transmit the data to IBM**
- **No deliverable will be returned to you**
- **We will be ever grateful for your contribution**

## z196 and z114 Support for Energy Savings

- **Processor performance (capability) can change due to over heating condition or static energy savings mode.**
- **Reflected in monitor data and QUERY CAPABILITY command.**

*Response (may only get first line on system with no changes):*

```
CAPABILITY: PRIMARY 696          SECONDARY 696          NOMINAL 696
CAPACITY-ADJUSTMENT INDICATION 100    CAPACITY-CHANGE REASON 0
RUNNING AT NOMINAL CAPACITY.
```

*Response for static power savings mode:*

```
RUNNING WITH REDUCED CAPACITY DUE TO A MANUAL CONTROL SETTING.
```

*Response possible for ambient temperature exceeded specified maximum:*

```
RUNNING WITH REDUCED CAPACITY DUE TO AN EXTERNAL EXCEPTION CONDITION.
```

## z/VM 6.2: Service Integrated in Base of z/VM 6.2

- **VM64774 SET/QUERY REORDER command**
- **All of the SSL scaling fixes**
- **VM64721 LIMITHARD now works**
  - SET SRM LIMITHARD CONSUMPTION is default now
- **VM64767/64876 VARY PROCESSOR causes hangs**
- **VM64850 VSWITCH failover buffer mixup**
- **VM64795 Enhanced Contiguous Frame Handling**
- **VM64927 Spin Lock Manager Improvement**
- **VM64887 Erratic System Performance (PLDV overflow)**
- **VM64756 Long CPEBK Chains, Master-only work, and SYSTEMMP**



IBM z/VM Development Lab – Endicott, NY

## z/VM Performance: April 2012 SPEs

## High Performance FICON: Outline

- **VM65041 lets z/VM guests use transport-mode I/O if channels and control units are so capable**
- **Transport-mode I/O uses a simpler command word structure that is easier for the channel subsystem and FICON adapter to handle, compared to conventional command-mode I/O**
- **On comparably configured workloads, transport-mode I/O gave us:**
  - About 35% increase in I/O rate
  - About 18% decrease in I/O service time
  - About 45% to 75% in CP CPU time per I/O
- **Workloads doing large I/Os tended to benefit most**
- **[www.vm.ibm.com/perf/reports/zvm/html/620jb.html](http://www.vm.ibm.com/perf/reports/zvm/html/620jb.html)**

## High Performance FICON: Some Numbers

<b>Transport-mode I/O vs. command-mode I/O, 67% reads, 1 record per I/O</b>					
<b>Guests/vol</b>	<b>Run Name</b>	<b>I/Os/vol/sec</b>	<b>Serv/I/O (msec)</b>	<b>%Busy/vol</b>	<b>%CP-CPU/I/O</b>
1	JB001238 (c)	3605.1	0.2326	83.8508	0.00156
	JB001239 (t)	4914.6	0.1855	91.1570	0.00090
	Delta	1309.5	-0.0471	7.3062	-0.00066
	%Delta	36.32	-20.25	8.71	-42.33
<b>Notes:</b> (c) denotes command-mode I/O. (t) denotes transport-mode I/O.					

<b>Transport-mode I/O vs. command-mode I/O, 67% reads, 64 records per I/O</b>					
<b>Guests/vol</b>	<b>Run Name</b>	<b>I/Os/vol/sec</b>	<b>Serv/I/O (msec)</b>	<b>%Busy/vol</b>	<b>%CP-CPU/I/O</b>
1	JB001246 (c)	359.4	2.6454	95.0844	0.00976
	JB001247 (t)	554.4	1.7660	97.9029	0.00242
	Delta	195.0	-0.8794	2.8185	-0.00734
	%Delta	54.26	-33.24	2.96	-75.17
<b>Notes:</b> (c) denotes command-mode I/O. (t) denotes transport-mode I/O.					

z10, 4 ded, 30G/2G, 4 FICON Express8, switched, DS8800/6GB, z/VM 6.2 plus VM65041, IO3390.



## High Performance FICON: Interaction with MDC

- **Transport-mode I/O directed at a minidisk will shut off MDC for said minidisk**
- **QUERY MDC command output will reveal that it happened**

off for FRED 1234

Disabled by Transport Mode I/O

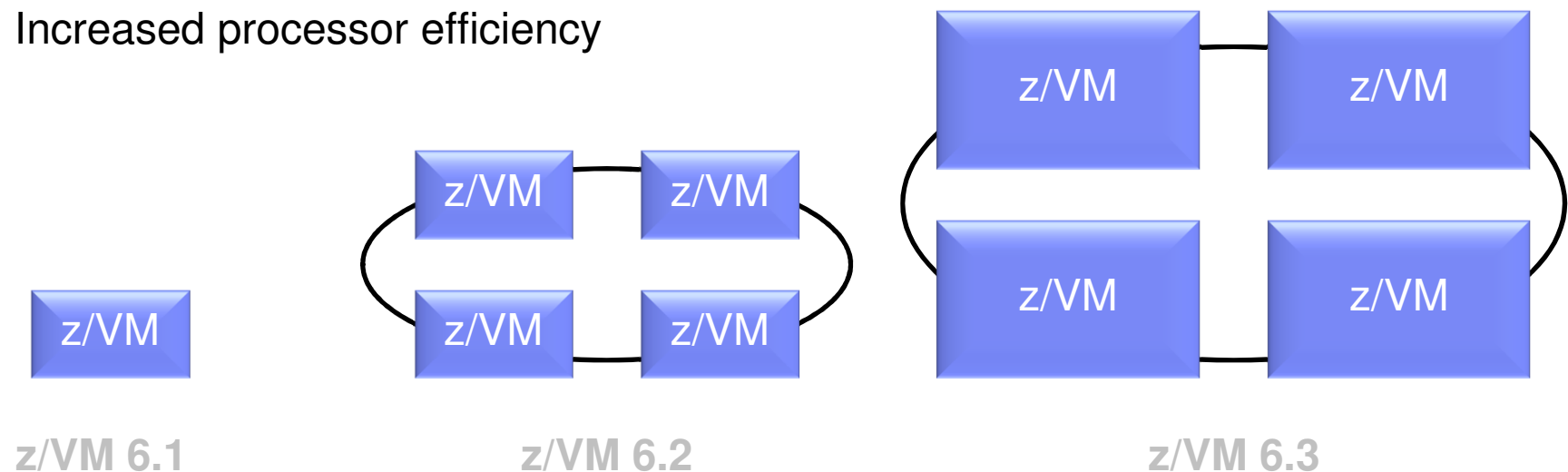


IBM z/VM Development Lab – Endicott, NY

## z/VM 6.3

## z/VM 6.3 – Making Room to Grow Your Business

- **Preview Announcement introducing z/VM 6.3 to be made February 5, 2013**
- **Planned Availability 3<sup>rd</sup> Quarter 2013**
- **Major Enhancements for Scalability and Performance**
  - Support for larger amounts of real memory
  - Increased processor efficiency



# Large Memory Support

- **Increases the real memory limit from 256GB to 1TB**
  - Proportionately increases total virtual memory based on tolerable over commitment levels and workload dependencies
  - Individual virtual machines up to 1TB
  
- **Paging DASD utilization and requirements change**
  - Proactive writing of pages to DASD increases need to have properly configured paging subsystem
  - Removed the need to double the paging space on DASD
    - Some additional space will continue to be recommended to avoid problems.
  
- **Expanded Storage continues to be supported with limit of 128GB**

# Large Memory Support

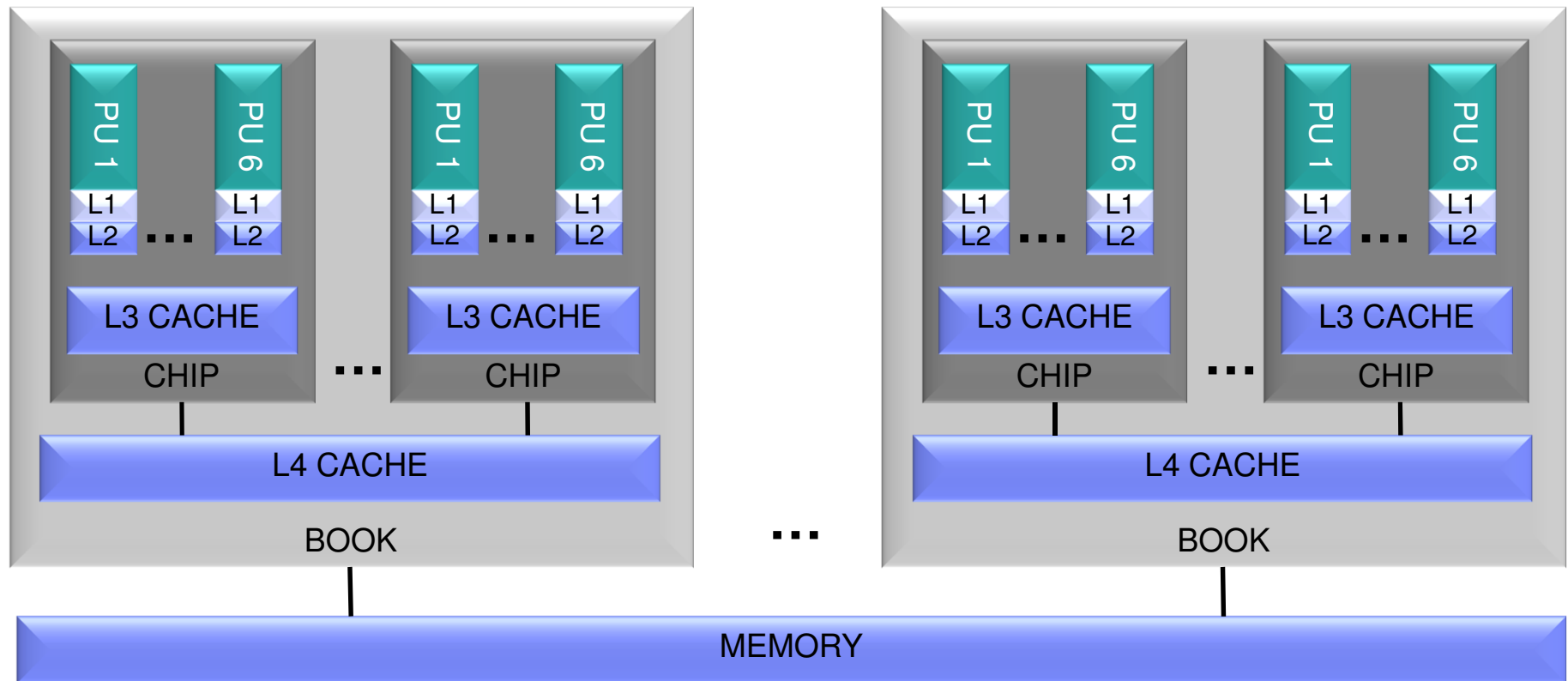
- **Page selection algorithms rewritten**
  - Reorder processing removed
  - Greater separation from the scheduler lists
    - Better handling of Linux guests that do not go truly idle
  
- **Improved effectiveness of the CP SET RESERVE command**
  - Pages protected better than previously
  - Support for reserving pages of NSS or DCSS space
    - Example: Use with the Monitor Segment (MONDCSS)
  - Ability to limit the overall number of reserved pages for a system

# HiperDispatch

- **Improves processor efficiency**
  - Better n-way curves
    - Supported processors limit remains at 32
  - Better use of processor cache to take advantage of cache-rich system design
  
- **Two components:**
  - Dispatching Affinity: dispatching cognizant of processor cache topology
  - Vertical CPU Management: cooperation with PR/SM to distribute physical processor resources to logical processors more efficiently for some configurations

## HiperDispatch: Dispatching Affinity

- Processor cache structures become increasingly complex and critical to performance
- Goal is to re-dispatch work close (in terms of topology) to where it last ran



## HiperDispatch: Dispatching Affinity

- **Dispatcher changed to be aware of the cache topology and dispatch work accordingly**
  - Dispatch virtual CPU near where its data may be in cache based on where the virtual CPU was last dispatched
  
- **Potentially increases cache efficiency, lowering processor costs by reducing CPI (Cycles Per Instruction)**
  
- **Previously, z/VM used soft affinity to processor in dispatching virtual CPUs**
  - No awareness of chip or book



## HiperDispatch: Vertical CPU Management

- **Today's “horizontal” management distributes the LPAR weight evenly distributed across the logical processors of the z/VM LPAR**
- **“Vertical” management attempts to minimize the number of logical processors, allowing LPAR to similarly manage logical CPUs**

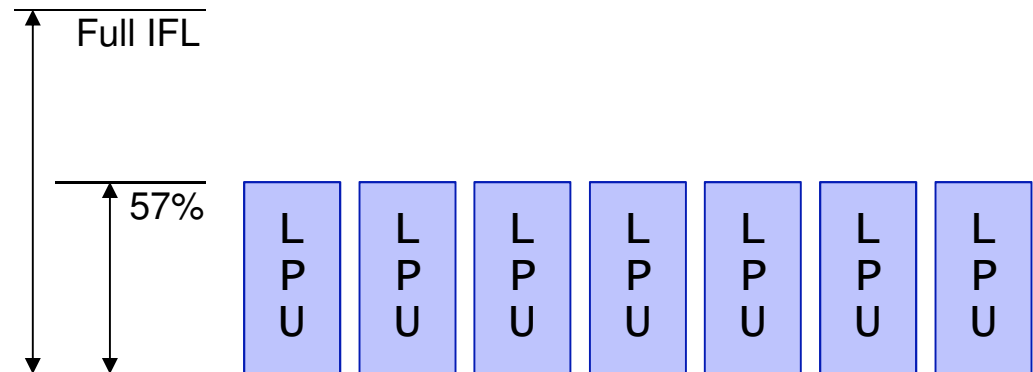
### Example:

- 10 Physical IFLs, 7 logical IFLs, weight of 400 out of 1000
  - Each logical IFL (LPU) entitled to 57% of an IFL
- When CEC is constrained, the LPAR's entitlement is reduced to 4 IFLs, so 7 is more than required
- z/VM & LPAR will cooperate
  - z/VM will concentrate the workload on a smaller number of logical processors
  - LPAR will redistribute the partition weight to give a greater portion to this smaller number of logical processors (~100% of 4 CPUs)

# Horizontal vs. Vertical CPU Management

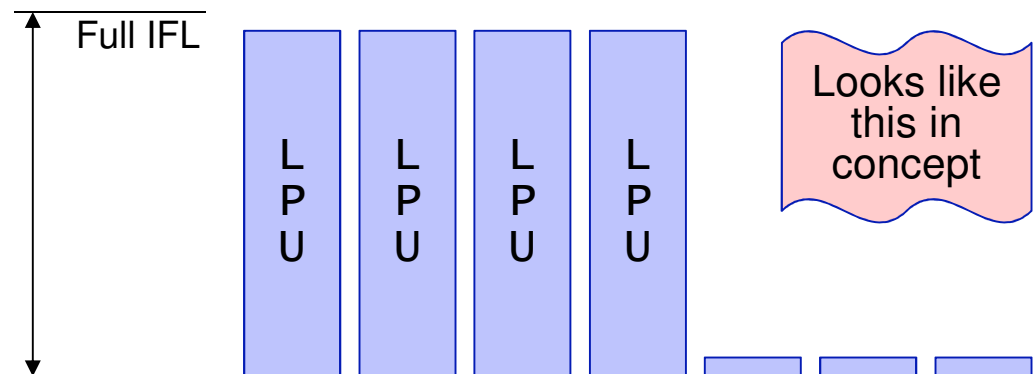
## Horizontal:

- The logical processors are all created/treated equally.
- z/VM dispatches work evenly across the 7 logical processors



## Vertical:

- The logical processors are skewed to where some get greater share of the weight.
- z/VM dispatches work accordingly to the heavier weighted workload.





IBM z/VM Development Lab – Endicott, NY

# Summary

## z/VM Performance Update: Summary

- **z/VM 6.2: SSI and LGR, plus more**
  - Loose clustering for guest mobility
  - Recognition of systems becoming larger
    - Memory management improvements
    - Better defaults: MONDCSS, SAMPLE CONFIG, STORBUF
  - CPU MF counters: help us, help you
  - Lots of good service rolled into the base
  - See <http://www.vm.ibm.com/perf/> for more details
- **The adventure continues**

**Contact Info:**

Bill Bitner

z/VM Customer Focus and Care

z/VM Development Lab – Endicott, NY

bitnerb@us.ibm.com

+1 607 -429 -3286