

IBM z Systems z13 Simultaneous Multi-Threading (R)Evolution

Daniel Rosa IBM Poughkeepsie dvrosa@us.ibm.com







#SHAREorg

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

© Copyright IBM Corp. 2015

Permission is granted to SHARE Inc. to publish this presentation paper in the SHARE Inc. proceedings; IBM retains the right to distribute copies of this presentation to whomever it chooses.

Trademarks

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

IBM* IBM Logo*

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

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.

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.

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

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

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.

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 other countries.

* 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.







Agenda

- SMT Overview
- General Industry SMT Exploitation
- z Systems SMT Exploitation
- z/OS SMT Exploitation with Runtime Capacity and Accounting
- z/OS SMT Control Implications
- References



z Systems pre-z13 Background



- A "ready" CPU executes a "ready" instruction (not resolving a cache miss) immediately
- A "not ready" CPU generally originates from a "not ready" instruction (resolving a cache miss)
 - CPU unproductive until instruction "ready"
 - z Systems workloads generally experience many cache misses

Legend: M=Resolve Cache Miss, I=Execute Ready Instruction



z Systems New z13 Machine





- A core contains multiple threads resembling CPUs
 - Each thread has its own state (PSW, regs, CPU Timer, etc)
 - "Ready" threads compete for use of shared core execution resources
 - Every cycle core arbitrates "ready" threads fairly, winner executes
 - "Not ready" threads do not compete, still unproductive
- z13 supports SMT-2 (2 threads per core)
 - Core has "ready" instructions more frequently, leads to higher core throughput





Effects of SMT



- Faster access, slower execution, higher core throughput
- Thread Density=2 (2 non-waiting threads, abbreviated TD=2), Workload dependent:
 - Capacity gain (benchmarks observe +10% to +40%)
 - Speed loss (benchmarks observe -30% to -40% <like z196>)
 - Thread Density=1 (1 waiting & 1 non-waiting thread, abbr TD=1)
 - Performs like SMT-1: no capacity gain, no speed loss
 - Workload dependent capacity in use, free
 - Core Waiting (all threads waiting)







- Generally focuses on maximizing core throughput
- Hypervisor provides SMT transparently to SMT-unaware OS
 - Dispatches OS logical CPU to physical thread, OS receives:
 - Less predictable capacity (core floats b/w TD=1, TD=2)
 - Less meaningful accounting (disconnected from capacity)
 - More latency, response time variability





- Increase core throughput predictably and repeatably
 - PR/SM supports SMT for SMT-aware OS via core dispatching
 - Limits SMT variability to a single z/OS workload
 - z/OS controls and manages whole core (all threads) to:
 - Maximize core throughput (fill running, spill to waiting cores)
 - Maximize core availability (meet goals using fewest cores)
- OS receives more predictable and repeatable capacity, accounting, latency, response time



Complete your session evaluations online at www.SHARE.org/Seattle-Eval

z/OS MT Exploitation with z13

MT available via PTFs on z/OS V2R1

- MT ZIIP MODE=1|2 active threads per online zIIP core Enables MT=2 always, or MT=2 (OLTP) and MT=1 (batch)

Define LOADxx PROCessor VIEW for the life of the IPL PROCVIEW CORE – MT environment, controls

Architecture and z/OS refer to MT, synonymous with SMT

- Pre-z13, must use MT controls (common look and feel)
- z13 becomes MT-2 config (defines 2 threads per core)
- PROCVIEW <u>CPU</u> Pre-z13 environment, controls
- Without IPL, IEAOPTxx supports operator changing MT Mode=x (active threads per online core, abbreviated MT=x) by core type:





Capacity and Accounting Overview



03/04/1

- Apply existing concepts to MT environment by:
 - Reinterpreting existing fields, using new MT metrics
- Generally requires the following over some time interval:
 - Transactions complete (External Throughput Rate)
 - Capacity max, free, in use (utilization) by job (accounting)
- Capacity in use per transaction, different for every workload
 - MT Mode=1, capacity max known, free & in use calculable
 - MT Mode=2, capacity max, free, in use workload dependent
 - Measures workload behavior at TD=1, TD=2, delivers meaningful capacity and accounting <u>at runtime</u>



- Max Capacity Factor: How much work a core can complete
- MT=1 Max Capacity Factor is 1.0 (100%)
- MT=2 Max Capacity Factor represents how much work a MT=2 core at TD=2 completes relative to a MT=1 core
- Expect Max Cap Factor between 1.1-1.4 (110%-140%)

Theoretical highest Max Capacity Factor=2.0, can get < 1.0

03/04/



- Capacity Factor: How much work core actually completes
- MT=1 Capacity Factor is 1.0 (100%)
- MT=2 Capacity Factor represents how much work a MT=2 core is <u>actually completing</u> relative to a MT=1 core. Considers:
 - MT=1, MT=2 Max Cap Factors. Time at TD=1, TD=2
- Expect Cap Factor 1.0-1.4 (100%-140%), can be < 1.0 (100%)
 - More TD=1 nears 1.0,

more TD=2 nears max Cap Factor



- Core Busy Time: Core execution time, includes:
 - Time any thread on the core is executing
 - Time core executes at TD=1, TD=2
- Average Thread Density: Average executing threads during Core Busy Time
 - MT=1 Average Thread Density is 1.0 (100%)
 - MT=2 Average Thread Density is 1.0-2.0 (100%-200%)

03/04/







z/OS Accounting With MT Mode=2

Similar capacity use should yield similar accounting



03/04/15

- MT Mode=1, ready instructions execute immediately
 - Account for capacity use with CPU Timer delta
- MT Mode=2, ready instructions compete to execute
 - CPU Timer decrements on both threads each cycle win or lose
 - Convert CPU Timer delta to MT=1 Equivalent Time
 - (MT=1 Core Capacity use) = (Capacity Factor)*(Core Busy Time)
 - Use Avg Thread Density to distribute over threads
 - If Cap Factor=1.25, Core Busy Time=4s, Avg TD=2.0, core used 1.25 * 4s = 5s of MT=1 Equivalent Time and each thread used 5s / 2.0=2.5s of MT=1 Equivalent Time
 - All accounting fields (SMF 30, 72, etc), services (TimeUsed) transparently become MT=1 Equivalent Times

HISMT Service



- Provides MT metrics between current, previous call. Requires:
 - Authorization: Supervisor State key 0
 - Dispatchable Unit: Task or SRB
 - Cross Memory Mode: Any HASN, any PASN, any SASN
 - Amode: 31-bit
 - ASC Mode: Primary
 - Interrupts: Enabled/Disabled for I/O & external interrupts
 - Control Parameters, save area must be addressable from primary and be in Disabled REFerence or Fixed storage
- Callable for any PROCVIEW, any MT Mode, HIS address space can be active/inactive



HISMT Syntax



- HISMT INTVAREA=xintvarea
 - , INTVAREALEN=xintvarealen
 - [, PRODCLASS=NO|YES]
 - [, PRODCORE=NO|YES]
 - [,CAPCLASS=NO|YES]
 - [,MAXCAPCLASS=NO|YES]

[,COREBUSYTIME=NO|YES] ← requires z13 PROCVIEW CORE (see CVT macro field

CvtMultiCpusPerCore)

- [,AVGTDCLASS=NO|YES]
- [,RETCODE=xretcode]
- [,RSNCODE=xrsncode]



MT=1	RMF	Сри	Activity F	Repor	t		
CPU		[FIME %		MT	%	LOG PROC
NUM TYPI	E ONLINE	LPAR BU	JSY MVS BUSY	PARKED	PROD	UTIL	
• • •							
4 II!	P 100.00	70.47	70.32	0.00	100.00	70.47	100.0
5 II	P 100.00	55.40	55.32	0.00	100.00	55.40	100.0
6 II	P 100.00	71.62	71.49	0.00	100.00	71.62	100.0
7 II	P 100.00	57.77	57.71	0.00	100.00	57.77	100.0
TOTAL/AV	VERAGE	63.81	63.71		100.00	63.81	400.0
	MUI	LTI-THRE	EADING ANALYS	SIS			
CPU TYI	PE MOI	DE I	MAX CF	CF	A۱	/G TD	
CI	2	1	1.000	1.000	.000 1.000		
I	IP_	1	1.000	1.000	1	L.000	

LPAR Busy = PR/SM dispatching logical CPU to physical CPU

03/04/15

- MVS Busy = Unparked logical CPU not waiting
- Parked = Logical CPU parked
- CPU utilization = (sum of CPUs' LPAR Busy)

Available CPU capacity = (Total Log Proc Share %) -(Sum of CPUs' LPAR Busy)

MT=2 RMF Cpu Activity Report											
CE NUM	PU TYPE	ONLINE	LPAR	TIME BUSY	8 MVS	BUSY	PARKED	M' PROD	F % UTIL	LOG PROC SHARE %	SHARE Educate · Network · Influence
4	IIP	100.00	78.23		67. 58.	24 40	0.00	87.30	68.29	100.0	
5	IIP	100.00	59.46		50. 41.	57 88	0.00	85.64	50.92	100.0	
6	IIP	100.00	80.77		70.	34 20	0.00	88.38	71.38	100.0	
7	IIP	100.00	63.67		55. 45.	08 52	0.00	86.43	55.03	100.0	
TOTA	AL/AVE	ERAGE MU	70.53 LTI-TH	READI	56. NG	41 ANALY	SIS	86.94	61.41	400.0	
CPU	J TYPI	e mo	DE	MAX	CF		CF	Ž	AVG TD		
	CP		1	1.0	00		1.000		1.000		
	III	2	2	1.4	73		1.283		1.600		
LPAR Busy = PR/SM dispatching logical core to physical core										l core	

- MVS Busy = Unparked logical CPU not waiting
- Parked = Logical CPU parked



M	Γ=2	RMF	Cpu Ac	ctivity F	Repoi	rt			
CI NUM	PU TYPE	ONLINE	TIM LPAR BUSY	E % MVS BUSY	PARKED	PROD	C % UTIL	LOG PROC SHARE %	SHARE Educate · Network · Influence
•••• 4	IIP	100.00	78.23	67.24	0.00	87.30	68.29	100.0	
5	IIP	100.00	59.46	58.40 50.57 41.88	0.00	85.64	50.92	100.0	
6	IIP	100.00	80.77	70.34	0.00	88.38	71.38	100.0	
7	IIP	100.00	63.67	55.08 45.52	0.00	86.43	55.03	100.0	
TOT <i>I</i>	AL/AVE	ERAGE MUI	70.53 LTI-THREAD	56.41 ING ANALYS	SIS	86.94	61.41	400.0	
CPU	J TYPI CP	e Moi	DE MAX 1 1.	CF 000	CF 1.000	1	AVG TD 1.000		
· C	III Ore I	<u>.</u> Itilizatio	2 1. on (% M	473 FUJTII) =	1.283 = [(P	AR Bi	1.600	Productiv	vitv)]
- Total zIIP (MT=2) core utilization (% MT UTIL): 245.62%									
Available core capacity = (Total Log Proc Share %) - (Sum of									
cores' MT % UTIL)									
-10 otal ZIIP (IVI I = 2) available: 400%-245.62%=154.38%									
Complete your session evaluations online at www.SHARE.org/Seattle-Eval								03/04/15	Seattle 2015
© Copyright IBM Corp. 2015								21	

MT=2 Workload Activity Report



- #SWAPS HST 0.000 AAP N/A (N/A EXCTD AAP TTP 220.69 0 AVG ENC 12.69 IIP 1995.640
- Service Times are in MT=1 Equivalent Time units
- APPL % is % of core relative to its max Capacity Factor
 - IIP APPL %=(1995.64s)*100/ [(614s)*(1.473)] = 220.65 %
 - CP APPL % = (3545.743s-1995.640s) * 100 / [(614s) * (1.0)] = 252.46 %



Complete your session evaluations online at www.SHARE.org/Seattle-Eval

AVG

MPT.

MT=2 Workload Activity Report

	OPKI				7		_
2R1 SYSPLE	EX PATPLX	29 DATE 0	1/21/20)15 INTE	ERVAL 10.	.14.053	SHARE Educate · Network · Influence
RPT VERSI	ION V2R1	RMF	TIME 2	21.46.55	5		
BY: POLIC	CY=PATPLE	X WORKI	OAD=WAS	SWKLD	SERVICE	CLASS=WA	STRANS
RESOURCE	GROUP=*N	ONE	CRITIC	CAL=CPU			
ACTIONS-	SERV	ICE TIME	API	PL %			
12.69	CPU	3545.743	CP	252.44			
12.69	SRB	0.000	AAPCP	0.00			
5502452	RCT	0.000	IIPCP	0.43			
8960.87	IIT	0.000					
0	HST	0.000	AAP	N/A			
0	AAP	N/A	IIP	220.69			
C 12.69	IIP	1995.640					
	W 2R1 SYSPLE RPT VERSI BY: POLIC RESOURCE ACTIONS- 12.69 12.69 5502452 8960.87 0 0 12.69	W O R K L 2R1 SYSPLEX PATPLX RPT VERSION V2R1 BY: POLICY=PATPLE RESOURCE GROUP=*N ACTIONS SERV 12.69 CPU 12.69 SRB 5502452 RCT 8960.87 IIT 0 HST 0 AAP C 12.69 IIP	WORKLOAD A 2R1 SYSPLEX PATPLX29 DATE 0 RPT VERSION V2R1 RMF BY: POLICY=PATPLEX WORKL RESOURCE GROUP=*NONE ACTIONS- SERVICE TIME 12.69 CPU 3545.743 12.69 SRB 0.000 5502452 RCT 0.000 0 HST 0.000 0 AAP N/A 12.69 IIP 1995.640	WORKLOAD ACTI 2R1 SYSPLEX PATPLX29 DATE 01/21/20 RPT VERSION V2R1 RMF TIME 2 BY: POLICY=PATPLEX WORKLOAD=WAS RESOURCE GROUP=*NONE CRITIC ACTIONS- SERVICE TIME APH 12.69 CPU 3545.743 CP 12.69 SRB 0.000 AAPCP 5502452 RCT 0.000 IIPCP 8960.87 IIT 0.000 AAP 0 AAP N/A IIP 2 12.69 IIP 1995.640	W O R K L O A D A C T I V I T Y 2R1 SYSPLEX PATPLX29 DATE 01/21/2015 INTER RPT VERSION V2R1 RMF TIME 21.46.55 BY: POLICY=PATPLEX WORKLOAD=WASWKLD RESOURCE GROUP=*NONE CRITICAL=CPU ACTIONS- SERVICE TIME APPL % 12.69 CPU 3545.743 CP 252.44 12.69 SRB 0.000 AAPCP 0.00 5502452 RCT 0.000 IIPCP 0.43 8960.87 IIT 0.000 AAP N/A 0 AAP N/A IIP 220.69 C 12.69 IIP 1995.640 12.69 12.69	W O R K L O A D A C T I V I T Y 2R1 SYSPLEX PATPLX29 DATE 01/21/2015 INTERVAL 10 RPT VERSION V2R1 RMF TIME 21.46.55 BY: POLICY=PATPLEX WORKLOAD=WASWKLD SERVICE RESOURCE GROUP=*NONE CRITICAL=CPU ACTIONS- SERVICE TIME APPL % 12.69 CPU 3545.743 CP 252.44 12.69 SRB 0.000 AAPCP 0.00 5502452 RCT 0.000 IIPCP 0.43 8960.87 IIT 0.000 AAP N/A 0 AAP N/A IIP 220.69 C 12.69 IIP 1995.640 12.69 IIP	WORKLOAD ACTIVITY 2R1 SYSPLEX PATPLX29 DATE 01/21/2015 INTERVAL 10.14.053 RPT VERSION V2R1 RMF TIME 21.46.55 BY: POLICY=PATPLEX WORKLOAD=WASWKLD SERVICE CLASS=WA RESOURCE GROUP=*NONE CRITICAL=CPU ACTIONS- SERVICE TIME APPL % 12.69 CPU 3545.743 CP 252.44 12.69 SRB 0.000 AAPCP 0.00 5502452 RCT 0.000 IIPCP 0.43 8960.87 IIT 0.000 AAP N/A 0 AAP N/A IIP 220.69 C 12.69 IIP 1995.640 12.69

- zIIP Capture Ratio = Sum(service classes' APPL %) * 100 / Sum(cores' MT % UTIL). If WASTRANS is the only service class using zIIPs:
 - zIIP Capture Ratio: 220.69 * 100 / 245.62 = 89.9%



Capacity and Accounting Observations Production zIIP transactional workloads 20-90% Busy (Washington Systems Center recommends 1-2 zIIPs run <=50% Busy, >=3 zIIPs run <=70% busy) generally receive optimal metric results

- Yields workload representative sample at TD=1, TD=2
- Cases where MT metrics may yield suboptimal results:
 - High/low utilization can yield insufficient samples at TD=1/2
 - Non-production, non-transactional, or unfriendly MT workloads like compute intensive, small memory footprint
- Indicators of suboptimal metric results may include:
 - Frozen metrics (typically insufficient samples at TD=1/2)
 - Max Cap Factor, Cap Factor at min (0.5) or max (2.0)

HIS Support Of (Core) MT-Diagnostic Counter Set



03/04/

- Hardware actives, enables on z13 w/ PROCVIEW CORE SHARE
- Modify HIS command updates for MT-Diagnostic Counter Set:
 - CTRSET=(MTD or MTDIAG or MTDIAGNOSTIC) collects it
 - CTRSET=HDWR or CTRSET=COMPLETE includes it
- SMF Type 113 Record Changes
 - Subtype 1 counter data sections are CPU or core specific
 - MTDiag deltas in subtype 1 for 1st online CPU per core
 - MTDiag counters not in subtype 2
 - Counter data sections for existing counter sets for Subtype 1 or 2 remain CPU specific
- USS output file
 - New counter data section for MT-diagnostic counter set for first online CPU per core

PROCVIEW CORE Control Implications



- System programmers manage cores (including MT Mode)^{S H A R E}
 - z/OS manages threads accordingly
- Parmlib changes:
 - CONFIGxx must specify CORE keyword e.g.:
 - CORE(0),ONLINE,STANDARD
 - LOADxx DYNCPADD nnnn cores z/OS prepares to dynamically add
- Stand Alone Dump:
 - Continue to follow the same procedure
 - Includes appropriate CPU/thread status for any PROCVIEW



PROCVIEW CORE Control Implications



- MODE command controls machine check type recovery action via:
 - RECORD=ALL (just record error to LOGREC)
 - RECORD=n,INTERVAL=s, If n machine checks within s seconds:
 - Perform ACR or Timer recovery depending on machine check type
- With PROCVIEW CORE on any hardware:
 - ACR eligible machine check types now RECORD=ALL, CPU=ALL for:
 - PD (Instruction processing Damage)
 - SD (System Damage)
 - IV (InValid PSW or registers)
 - TC (Tod Clock damage)
 - PT (Processor Timer damage)
 - CC (Clock Comparator damage)
 - PS (Primary clock Synchronization)
 - AD (ETR attachment)
 - SL (Switch to Local synchronization)
 - Specifying RECORD=n or INTERVAL=s results in:
 - IGF960I MODE COMMAND REJECTED. PROCVIEW
 CORE IS IN EFFECT





- Config Core(x), Online Configs core online for MT Mode
 - MT Mode=1, 1st thread online, 2nd thread offline
 - MT Mode=2, both threads online
- Config Core(x),Offline Configs all threads on core offline
- Config Member=xx Configs cores according to CONFIGxx
- Config Online or Config Offline Lists eligible cores to config
 - Reply to IEE522D accepts CORE(x) to configure
- Display Matrix=Core Displays core status (new message)
- Display Matrix=Config(xx) CONFIGxx vs system differences



Display Matrix=Core Sample Output



CP cores 0-4 online, MT=1, thread 0 online, thread 1 offline (N means offline and cannot be brought online)

ZIIP cores 5-8 online, MT=2, both threads online



z/OS MT Support Summary



- z/OS V2R1 MT support requires applying the following PTFs:
 - UA90753 (OA43366) Supervisor, SADump, HIS, Reconfig, IPL/NIP, LoadWait
 - UA90762 (OA43622) SRM / WLM
 - UA76154 (OA44439) XCF
- z/OS V2R1 MT optional PTFs include:
 - UA76026 (OA44101) RMF which will require:
 - UA90772 (OA44624) USS



References



- z/OS MVS Authorized Assembler Services Reference EDT-IXG
- z/OS MVS Initialization and Tuning Reference
- z/OS MVS System Codes
- z/OS MVS System Commands
- z/OS MVS System Management Facilities (SMF)
- z/OS MVS System Messages, Vol 7 (IEB-IEE)
- z/OS RMF Report Analysis

IBM Redbooks Point-Of-View Publication: z Systems Simultaneous Multithreading Revolution



Questions?





Complete your session evaluations online at www.SHARE.org/Seattle-Eval

© Copyright IBM Corp. 2015