

SmarterComputing



The World of z/OS Dispatching on the zEC12 Processor

Glenn Anderson, IBM Lab Services and Training



Summer SHARE 2013
Session 14040



© 2013 IBM Corporation

What I hope to cover.....

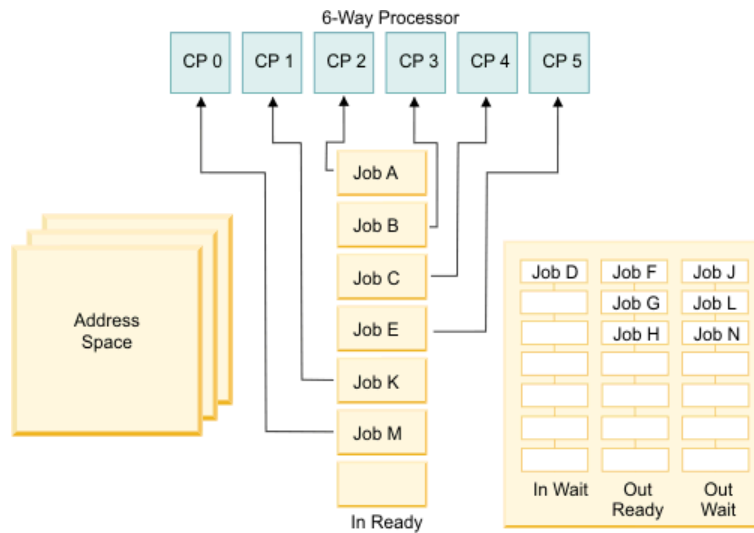
- What are dispatchable units of work on z/OS
- How WLM manages dispatchable units of work
- The role of HiperDispatch
- Dispatching work to zIIP and zAAP engines
- RMF Work Unit Analysis Report

z/OS Dispatchable Units

- **There are different types of Dispatchable Units (DU's) in z/OS**

- Preemptible Task (TCB)
- Non Preemptible Service Request (SRB)
- Preemptible Enclave Service Request (enclave SRB)
 - Independent - a new transaction
 - Dependent – extend existing address space
 - Work-dependent – extend existing independent enclave

z/OS Dispatching Work

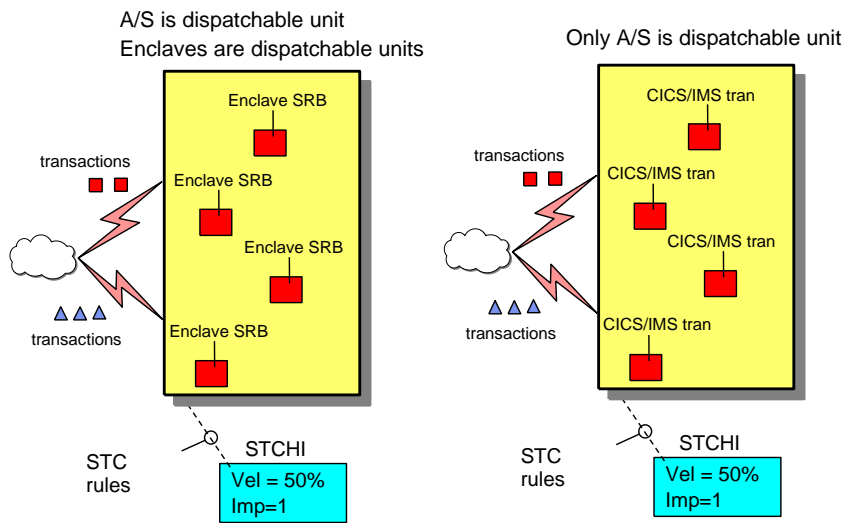


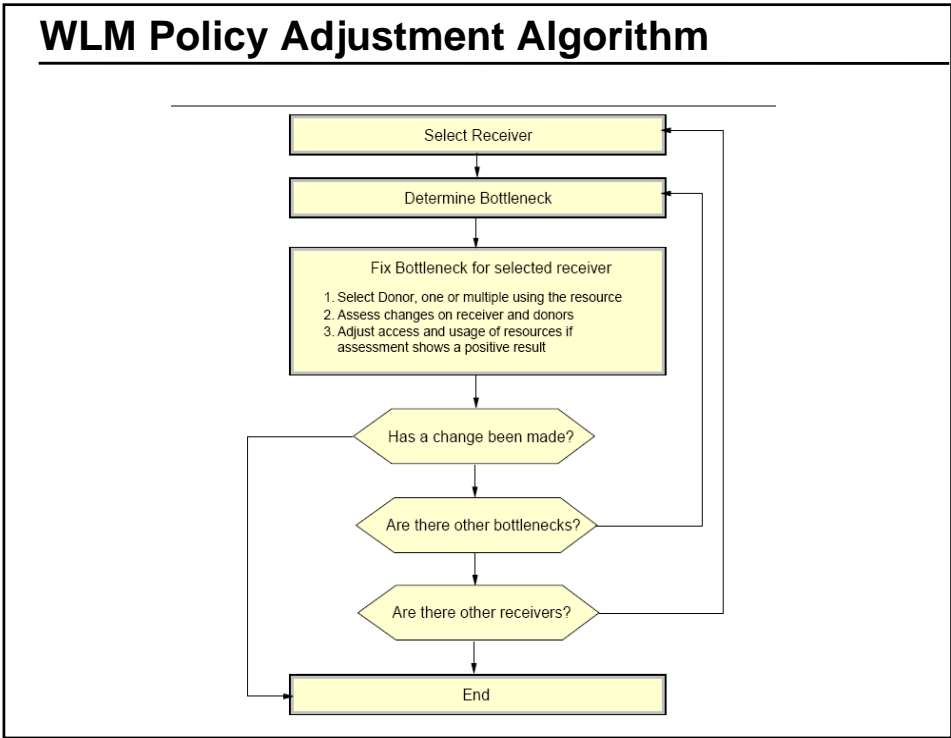
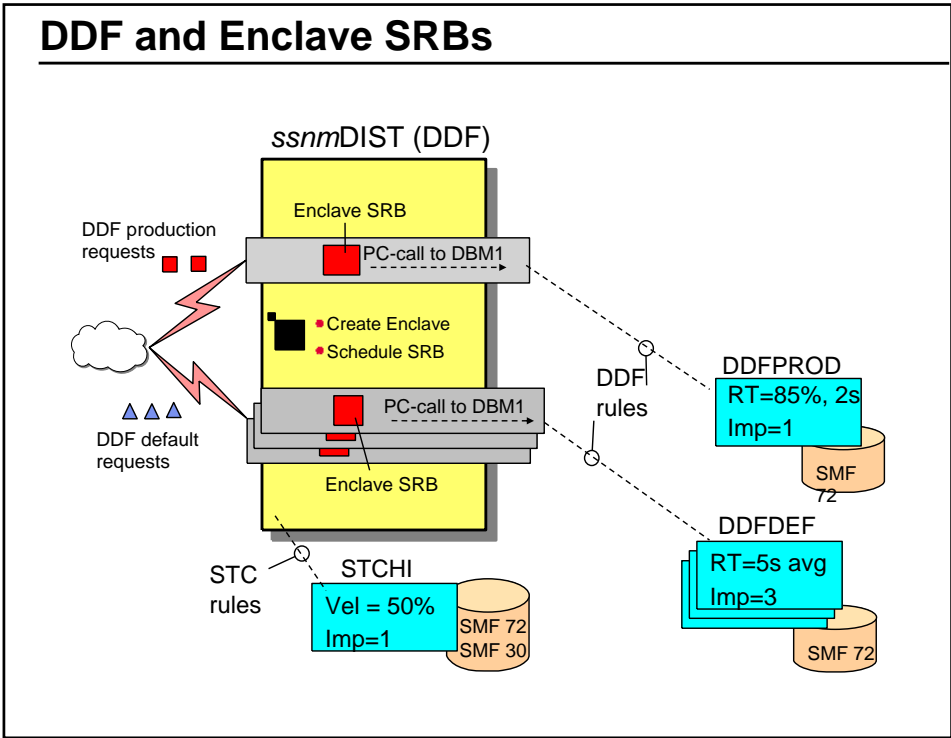
What is a WLM Transaction?

- **A WLM transaction represents a WLM "unit of work"**
 - basic workload entity for which WLM collects a resource usage value
 - foundation for statistics presented in workload activity report
 - represents a single subsystem "work request"
- **Subsystems can implement one of three transaction types**
 - **Address Space:**
 - ▶ WLM transaction measures all resource used by a subsystem request in a **single address space**
 - ▶ Used by JES (a batch job), TSO (a TSO command), OMVS (a process), STC (a started task) and ASCH (single APPC program)
 - **Enclave:**
 - ▶ Enclave created and destroyed by subsystem for each work request
 - ▶ WLM transaction measures resources used by a single subsystem request across **multiple address spaces and systems**
 - ▶ Exploited by subsystems - Component Broker(WebSphere), DB2, DDF, IWEB, MQSeries Workflow, LDAP, NETV, TCP
 - **CICS/IMS Transactions**
 - ▶ Neither address space or enclave oriented - special type
 - ▶ WLM transaction measures resource used by a single CICS/IMS transaction program request

The WLM View

Address Spaces, and the transactions inside

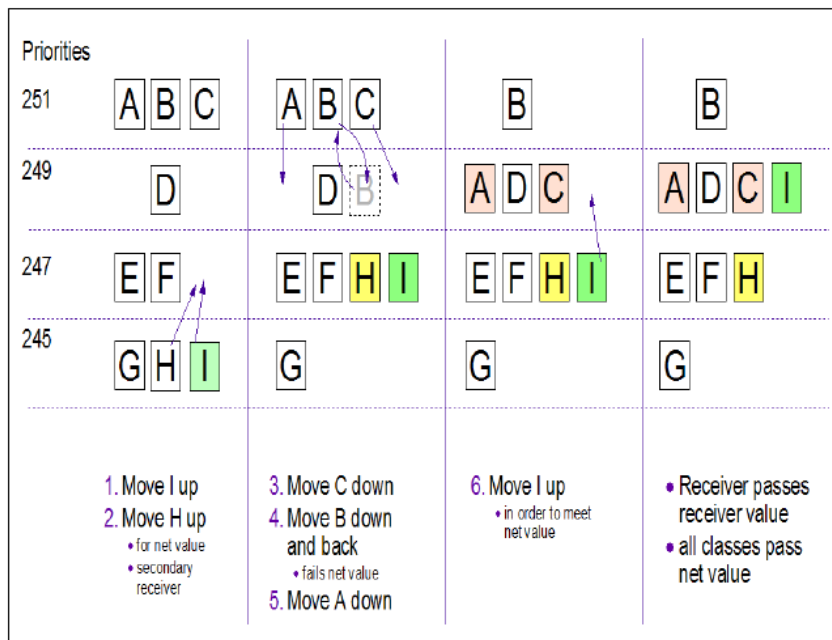




WLM Dispatching Priority Usage

255	SYSTEM
254	SYSSTC
253	<i>Small Consumer</i>
252	Priorities for dynamic policy adjustment
208	
207	
202	Not used
201	Discretionary work Mean Time to wit algorithm
192	

WLM CPU Management Example



WLM Blocked Workload Support: IEAOPTxx

BLWLTRPCT	<p>Percentage of the CPU capacity of the LPAR to be used for promotion</p> <ul style="list-style-type: none"> Specified in units of 0.1% Default is 5 (=0.5%) Maximum is 200 (=20%) Would only be spent when enough units of work exist which need promotion
BLWLINTHD	<p>Specifies threshold time interval for which a blocked address space or enclave must wait before being considered for promotion.</p> <ul style="list-style-type: none"> Minimum is 5 seconds. Maximum is 65535 seconds. Default is 20 seconds.

11

Blocked Workload Support: User Interface: RMF

```

...
CPU ACTIVITY
...
BLOCKED WORKLOAD ANALYSIS
OPT PARAMETERS: BLWLTRPCT (%) 0.5 PROMOTE RATE: DEFINED 50000 WAITERS FOR PROMOTE: AVG 0.001
                  BLWLINTHD 60                USED (%) 95                PEAK 15

```

- ❑ Extensions of RMF Postprocessor CPU Activity and WLMGL reports with information about blocked workloads and the temporary promotion of their dispatching priority
- ❑ SMF record 70-1 (CPU activity) and SMF 72-3 (Workload activity)

IBM

WORKLOAD ACTIVITY

z/OS VIRI3 SYSPLX SYPLEX3 DATE 09/28/2011 INTERVAL 15.00.003 MODE = GOAL PAGE 1
 RPT VERSION VIRI3 RMF TIME 17.00.00

POLICY ACTIVATION DATE/TIME 09/14/2011 11.08.09

----- SERVICE CLASS(ES)

REPORT BY: POLICY=BASEPOL WORKLOAD=STC_WLD SERVICE CLASS=STCLOW RESOURCE GROUP=NONE
 CRITICAL =NONE
 DESCRIPTION =Low priority for STC workloads

-TRANSACTIONS-	TRANS-TIME	HHH.MM.SS.TTT	--DASD	I/O--	---SERVICE---	SERVICE TIME	---APPL %---	--PROMOTED--	----STORAGE----	
AVG 153.37	ACTUAL	3.02.885	SSCHRT	56.9	IOC	3964	CPU 805.697	CP 92.24	BLK 1.489	AVG 1195.43
MPL 152.35	EXECUTION	3.02.391	RESP	15.1	CPU	15184K	SRB 13.850	AAPCP 0.00	ENQ 0.046	TOTAL 182122.4
ENDED 599	QUEUED	494	CONN	1.3	MSD	0	RCT 9.995	IIPCP 0.00	CRM 5.593	SHARED 230.59
END/S 0.67	R/S AFFIN	0	DISC	0.3	SRB	261005	IIT 0.576		LCK 0.000	
#SHAPS 3391	INELIGIBLE	0	Q-PEND	4.5	TOT	15449K	HST 0.000	AAP 0.000	SUP 0.000	-PAGE-IN RATES-
EXCTD 0	CONVERSION	5.188	IOSQ	9.0	/SEC	17202	AAP 0.000	IIP 0.000		SINGLE 0.0
AVG ENC 0.00	STD DEV	3.27.429								BLOCK 0.0
REM ENC 0.00					ABSRPTN	113				SHARED 0.0
MS ENC 0.00					TRX SERV	112				HSP 0.0

-----SERVICE CLASSES BEING SERVED-----

DB2LOW

IBM

© 2013 IBM Corporation

13

IBM

SERVICE TIME	---APPL %---	--PROMOTED--	----STORAGE----
CPU 805.697	CP 92.24	BLK 1.489	AVG 1195.43
SRB 13.850	AAPCP 0.00	ENQ 0.046	TOTAL 182122.4
RCT 9.995	IIPCP 0.00	CRM 5.593	SHARED 230.59
IIT 0.576		LCK 0.000	
HST 0.000	AAP 0.00	SUP 0.000	-PAGE-IN RATES-
AAP 0.000	IIP 0.00		SINGLE 0.0
IIP 0.000			BLOCK 0.0
			SHARED 0.0
			HSP 0.0

RVED-----

IBM

© 2013 IBM Corporation

14

IBM

Promoted

CPU time in seconds that transactions in this group were running at a promoted dispatching priority, separated by the reason for the promotion:

BLK CPU time in seconds consumed while the dispatching priority of work with low importance was temporarily raised to help blocked workloads

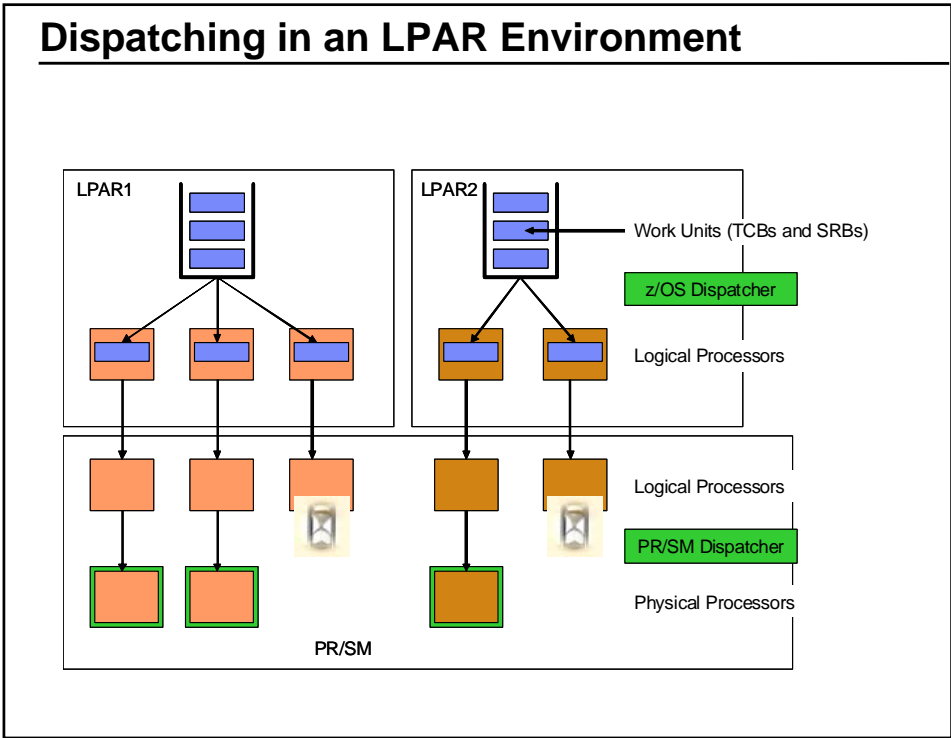
ENQ CPU time in seconds consumed while the dispatching priority was temporarily raised by enqueue management because the work held a resource that other work needed.

CRM CPU time in seconds consumed while the dispatching priority was temporarily raised by chronic resource contention management because the work held a resource that other work needed

LCK In HiperDispatch mode, the CPU time in seconds consumed while the dispatching priority was temporarily raised to shorten the lock hold time of a local suspend lock held by the work unit.

SUP CPU time in seconds consumed while the dispatching priority for a work unit was temporarily raised by the z/OS supervisor to a higher dispatching priority than assigned by WLM.

© 2013 IBM Corporation
15



HiperDispatch Introduction

- Design Objective
 - Keep work as much as possible local to a physical processor to optimize the usage of the processor caches
 - As a result systems with high number of physical processors provide a much better scalability
- Function: HiperDispatch
 - Interaction between z/OS and the PR/SM Hypervisor to optimize work unit and logical processor placement to physical processors
 - Consists of 2 parts
 - In z/OS (sometimes referred as Dispatcher Affinity)
 - In PR/SM (sometimes referred as Vertical CPU Management)

17

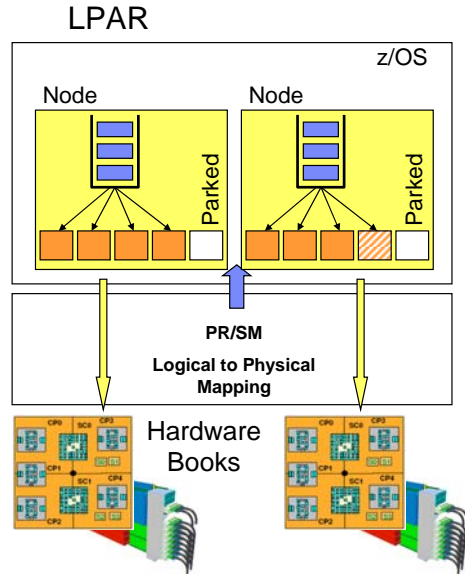


HiperDispatch Mode

- PR/SM
 - Supplies topology information/updates to z/OS
 - Ties *high priority* logicals to physicals (gives 100% share)
 - Distributes remaining share to *medium priority* logicals
 - Distributes any additional service to unparked *low priority* logicals
- z/OS
 - Ties tasks to small subsets of logical processors
 - Dispatches work to *high priority* subset of logicals
 - Parks *low priority* processors that are not need or will not get service

HiperDispatch: z/OS Part

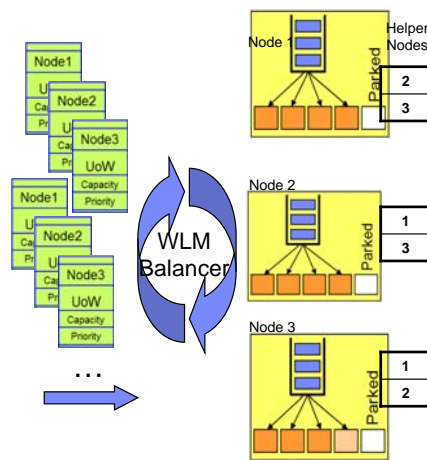
- z/OS obtains the logical to physical processor mapping in Hiperdispatch mode
 - Whether a logical processor has high, medium or low share
 - On which book the logical processor is located
- z/OS creates dispatch nodes
 - The idea is to have 4 high share CPUs in one node
 - Each node has TCBs and SRBs assigned to the node
 - Optimizes the execution of work units on z/OS



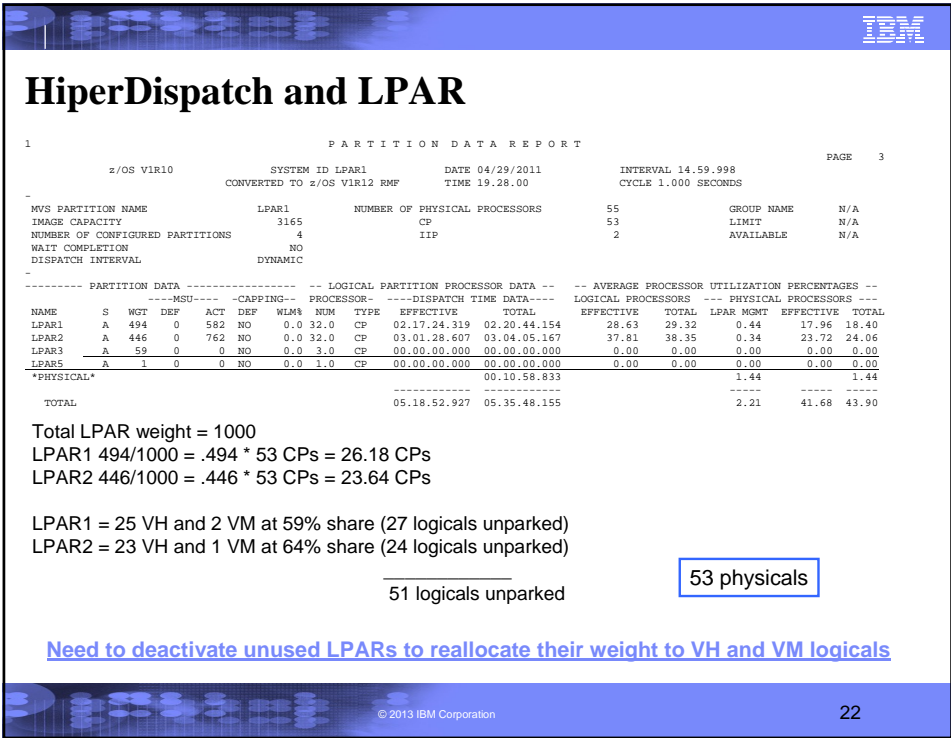
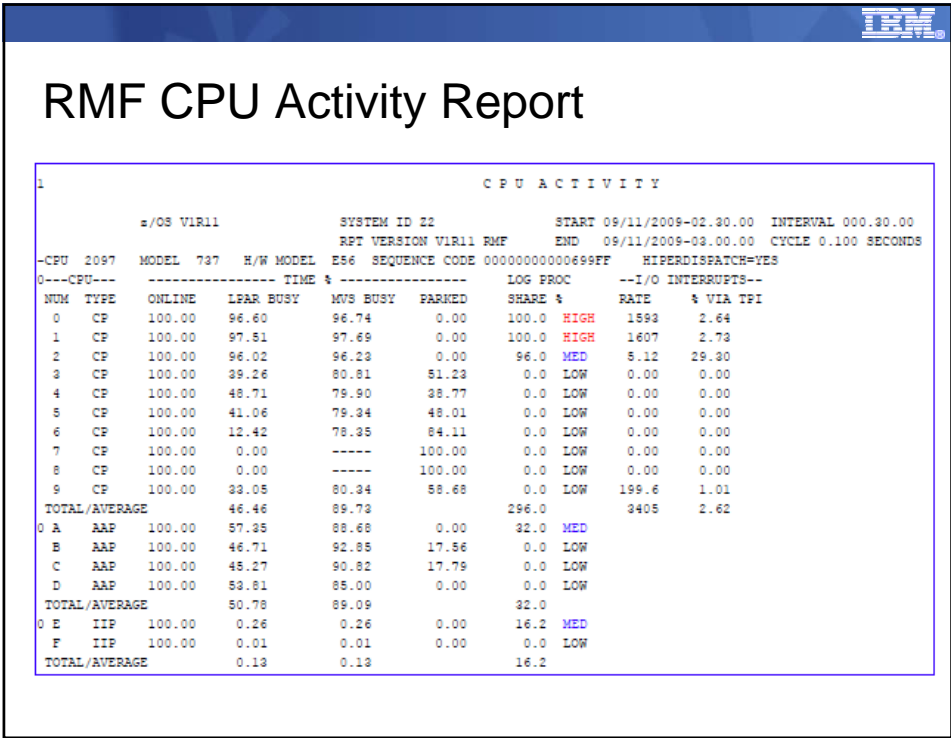
19

HiperDispatch: z/OS Affinity Dispatching

- Affinity dispatching
 - WLM balances the units of work (TCBs/SRBs) across the nodes to equalize the utilization of nodes
 - And to assure that each node has work of different priorities
 - Balancing takes place every 2 seconds
 - For unbalanced situations in between
 - WLM creates lists of helper nodes for each node
 - These helper nodes can be asked to select work from a given node in cases the node is overloaded
 - Helper nodes are sorted to avoid book crossing if possible
- Low share processors
 - WLM parks and un-parks these processors based on demand and utilization of the CEC



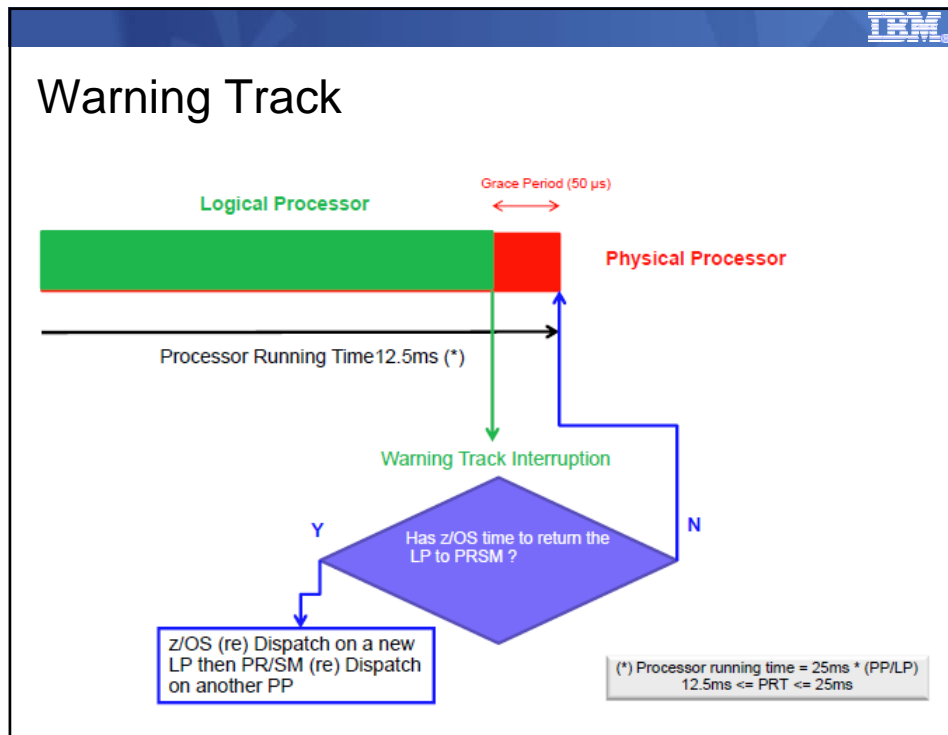
20



Warning Track

- ▶ In a PR/SM™ environment the LPAR hypervisor assigns physical engines to logical engines accordingly to the weighting factors of the partitions.
- ▶ Once the time slice for a logical engine is expired the currently executing work is suspended until a physical engine is assigned to the logical engine again.
- ▶ The Warning Track Interruption Facility notifies the operating system that PR/SM™ will undispach a certain logical processor within the next 50 microseconds (grace period).
- ▶ z/OS is now able to save status for the running unit of work and re-dispatch the work unit on a different logical processor within the grace period.
- ▶ z/OS now signals to PR/SM via Diagnose x'9C' that the logical processor can be un-dispatched.
- ▶ Warning Track processing is only supported in HyperDispatch=YES environments.
- ▶ A high benefit can be achieved for Low Share processors which might be parked by WLM.


© 2013 IBM Corporation 23



Warning Track Statistics

- ▶ RMF keeps track of the number of times PR/SM issued a warning-track interruption to a logical processor and z/OS was able/unable to return the logical processor within the grace period.
- ▶ RMF measures the amount of time in microseconds that a processor was yielded to PR/SM due to Warning-track processing.

SMF record type 70 subtype 1 (CPU Activity) – CPU data section				
Offset	Name	Length	Format	Description
80 x50	SMF70WTS	4	Binary	The number of times PR/SM issued a warning-track interruption to a logical processor and z/OS was able to return the logical processor within the grace period.
84 x54	SMF70WTU	4	Binary	The number of times PR/SM issued a warning-track interruption to a logical processor and z/OS was unable to return the logical processor within the grace period.
88 x58	SMF70WTI	4	Binary	Amount of time in microseconds that a logical processor was yielded to PR/SM due to Warning Track processing.



RMF Postprocessor Overview Conditions		
Name	Qualifier	Description
WTRKCP (WTRKAAP) (WTRKIIP)	cpu-id	The percentage of times PR/SM issued a warning-track interruption to a processor and z/OS was able to return it to PR/SM within the grace period.
WTRKTCP (WTRKTAAP) (WTRKTIIP)	cpu-id	Time in microseconds that a purpose processor was yielded to PR/SM due to Warning Track processing.

IBM System z Application Assist Processor (zAAP)

Specialty assist processor dedicated exclusively to execution of Java workloads under z/OS® – e.g. WebSphere®, CICS, IMS, DB2

- Available on IBM zEnterprise 196 and 114, IBM System z10, IBM System z9 and zSeries z990 and z890
- Used by any workload with Java cycles, e.g. WebSphere, CICS, Batch, DB2®.
- Executes Java code with no changes to applications
- **Objective: Enable integration of new Java based Web applications with core z/OS backend database environment for high performance, reliability, availability, security, and lower total cost of ownership**

Standard Processor

WebSphere

Execute JAVA Code

JVM

Switch to zAAP

z/OS Dispatcher

Suspend JVM task on z/OS standard logical processor

z/OS Dispatcher

Dispatch JVM task on z/OS standard logical processor

JVM

Switch to zAAP

WebSphere

zAAP

z/OS Dispatcher

Dispatch JVM task on z/OS zAAP logical processor

JVM

Switch to zAAP

Java Application Code

Executing on a zAAP logical processor

JVM

Switch to standard processor

z/OS Dispatcher

Suspend JVM task on z/OS zAAP logical processor


IBM System z Integrated Information Processor (zIIP)

Specialty assist processor dedicated exclusively to execution of any workloads under z/OS®

- Available on IBM zEnterprise 196 and 114, IBM System z10, IBM System z9
- Requires work to be run as an enclave SRB
- Exploiters work with z/OS to determine how much capacity should be redirected to a zIIP
- Exploiters must use licensed interface to enable exploitation


Current IBM Exploitation of zAAPs and zIIPs

Specialty CP	Eligible	Major Users
zAAP	Any Java Execution	Websphere CICS Native apps XMLSS
zIIP	Enclave SRBs	DRDA over TCPIP DB2 Parallel Query DB2 Utilities Load, Reorg, Rebuild DB2 V9 z/OS remote native SQL procedures TCPIP - IPSEC XMLSS zIIP Assisted HiperSockets Multiple Write Virtual Tape Facility Mainframe (VTFM) Software z/OS Global Mirror (XRC), System Data Mover (SDM) z/OS CIM Server



What is "Needs Help"

- Determination zIIP or zAAP work is being delayed and additional resources should help process the work
 - ▶ Requires xxPHONORPRIORITY=YES to be set
- If help is required:
 - ▶ The zxxP CP signals waiting zxxP to help
 - ▶ When all zxxP CPs are busy the zxxP asks for help from the GCP
 - All available speciality engines (of the same type) must be busy before help is asked of the GCPs
 - IF the zxxPs needs help and all zxxPs are busy help is obtained from 1 GCP
 - IF zxxPs continue to need help additional CPs may be asked to help
 - ▶ Help is always provided in dispatch priority order



Specialty CP work running in a WLM Service Class

```

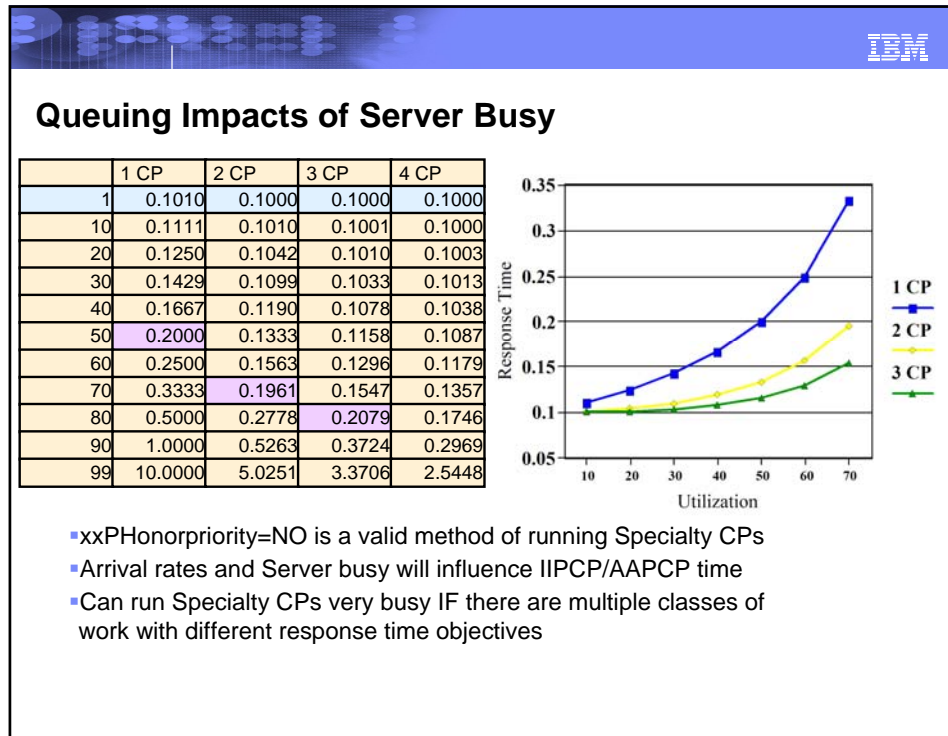
REPORT BY: POLICY=WLMPOL      WORKLOAD=BAT_WKL      SERVICE CLASS=BATSPEC      RESOURCE GROUP=BATMAXRG
TRANSACTIONS  TRANS-TIME HHH.MM.SS.TTT  --DASD I/O--  ---SERVICE---  SERVICE TIMES  ---APPL %---
AVG           0.98  ACTUAL           6.520  SSCHRT  11.5  IOC       8326  CPU      24.7  CP       0.97
MPL           0.98  EXECUTION          6.128  RESP    7.0  CPU      662386  SRB      0.0  AAPCP   0.01
ENDED         10  QUEUED             391    CONN    6.9  MSO       0    RCT      0.0  IIPCP   0.00
END/S         0.17  R/S AFFIN           0    DISC    0.0  SRB      965    IIT      0.0
#SWAPS        0    INELIGIBLE          0    Q+PEND  0.1  TOT     671677  HST      0.0  AAP     40.27
EXCTD         0    CONVERSION          0    IOSQ    0.0  /SEC    11195  AAP      24.2  IIP      0.00
AVG ENC       0.00  STD DEV             0
    
```

GOAL: EXECUTION VELOCITY 35.0% VELOCITY MIGRATION: I/O MGMT 99.2% INIT MGMT 92.2%

```

RESPONSE TIME EX  PERF  AVG  ----- USING% -----  EXECUTION DELAYS %  -----
SYSTEM           VEL%  INDX  ADRSP  CPU  AAP  IIP  I/O  TOT  CPU
SYSDB           --N/A--  99.2  0.4  1.0  0.8  45.9  0.0  3.9  0.4  0.4
    
```

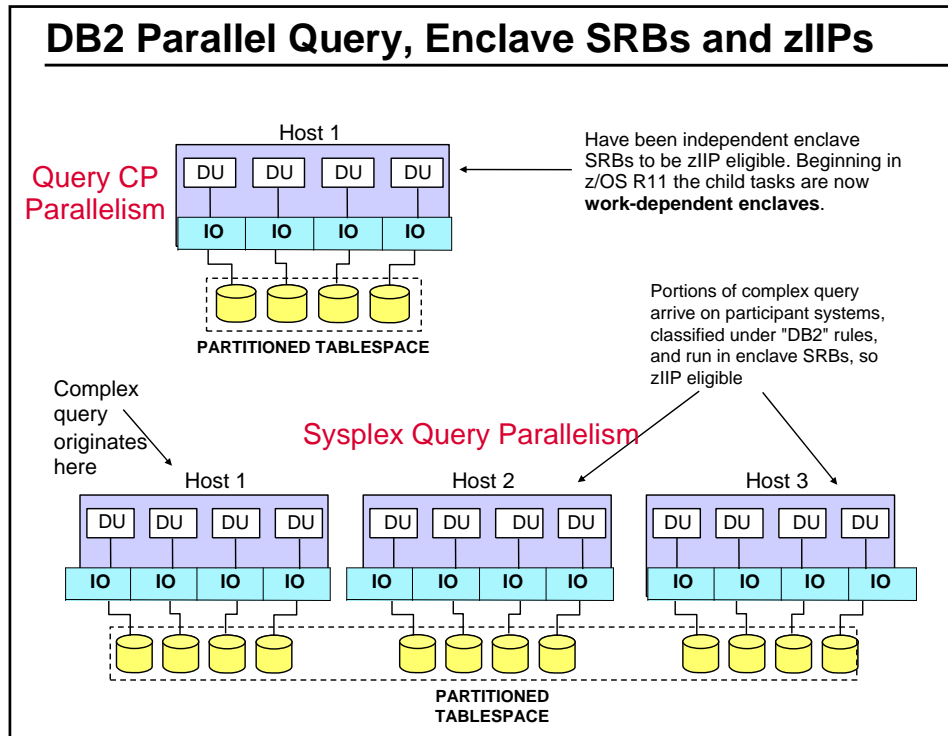
RMF report is at 1 minute interval



IBM

zAAPs on zIIPs

- When a [processor](#) has no zAAPs installed but the LPAR has zIIP(s), treat zAAP-eligible work instead as zIIP-eligible work
 - ▶ Exception: Under z/VM can define a guest without zAAPs, on a processor with zAAPs, to allow testing
- `IEASYSxx` system parm indicates whether zAAP on zIIP is enabled
 - ▶ `ZZ=NO` | YES for z/OS 1.9 and z/OS 1.10
 - ▶ `ZAAPZIIP=NO` | [YES](#) for z/OS 1.11 (also supports alias of `ZZ`)
 - ▶ The enablement state of zAAP on zIIP cannot be changed after IPL
- Timing fields within SMF will show offload time under zIIP if a zIIP is installed
 - ▶ No method to determine what portion of zIIP time originated due to a zAAP request
 - ▶ Accounting and Capacity planning may need updating to use zIIP field
- New operator command in z/OS R12 to display ZIIPZAAP eligibility
- New planning White Paper [WP101642](#)

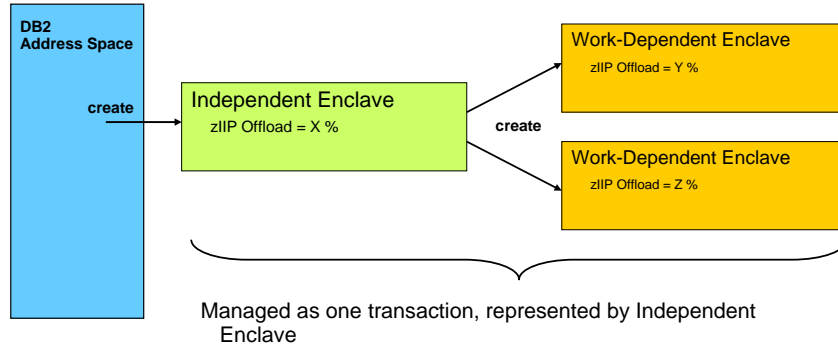


IBM

DB2 Parallelism, WLM, and zIIPs

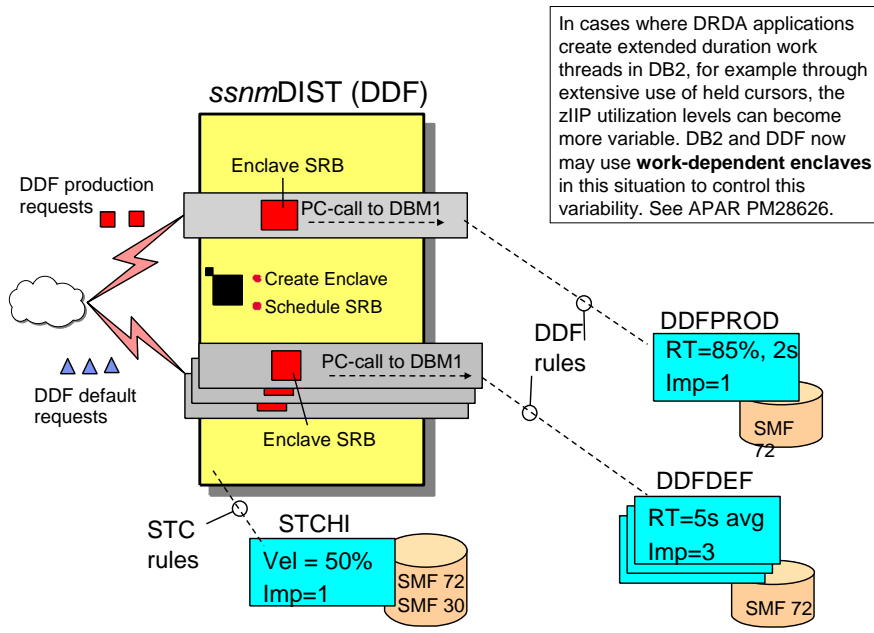
- DB2 Parallelism and zIIPs
 - ▶ Controlled by a CPU threshold. Once the threshold is met all child tasks are zIIP eligible
 - ▶ Parents are not zIIP eligible
 - ▶ Parent and child CPU time contribute to the CPU Threshold
 - ▶ Can see any kind of work, CICS, IMS, TSO, batch using zIIP resources
- DB2 will use new Work-Dependent Enclaves for Child tasks
 - ▶ APAR OA26104 for releases 1.8 and beyond
 - ▶ Without new Work Dependent Enclave support parallel enclaves must be classified using subsystem DB2
 - Unclassified work would wind up in SYSOTHER

Work-Dependent Enclaves



Implement a new type of enclave named "Work-Dependent" as an extension of an Independent Enclave. A Work-Dependent enclave becomes part of the Independent enclave's transaction but allows to have its own set of attributes (including zIIP offload percentage)

DDF and Work-dependent Enclaves



Work-dependent Enclaves in SDSF


The screenshot shows a terminal window titled 'E - TBLATT2.ws' with a menu bar (Display, Filter, View, Print, Options, Help). The main display area shows the output of the 'SDSF ENCLAVE DISPLAY' command. The output is a table with columns: NP, NAME, STATUS, TYPE, SRVCLASS, PER, RPTCLASS, CPU-TIME, OWNER, and RE. The data shows five active work-dependent enclaves (WDEP) with a velocity of 1, all running under the SYS1 owner.

NP	NAME	STATUS	TYPE	SRVCLASS	PER	RPTCLASS	CPU-TIME	OWNER	RE
2000000016		ACTIVE	IND	VEL_1	2	RC_1	0.00		32
240000001A		ACTIVE	WDEP	VEL_1	2	RC_1	0.39		32
280000001B		ACTIVE	WDEP	VEL_1	2	RC_1	0.39		32
2C00000019		ACTIVE	WDEP	VEL_1	2	RC_1	0.39		32
3000000018		ACTIVE	WDEP	VEL_1	2	RC_1	0.39		32
3400000017		ACTIVE	WDEP	VEL_1	2	RC_1	0.39		32

At the bottom of the terminal window, there is a status bar showing '04/021' and a connection message: 'Connected to remote server/host.vmttool1.pok.ibm.com using port 23'.


Statement of Direction

- zEC12 is planned to be last System z server to offer support for zAAP
- Continued support for zAAP on zIIP
- Remove restriction for zAAP on zIIP when zAAP installed on server
 - APAR OA38829 on z/OS R12 and R13



July 2013 zIIP/zAAP Announcement

- Effective July 23, 2013, IBM has modified the ratio of zIIP/zAAPs to CPs to be 2:1 for a zEC12 and/or zBC12. Customers may purchase up to two zIIP and/or up to two zAAP processors for every general purpose processor they purchase on the server. For the z196 and z114 and earlier servers, the 1:1 ratio will still be enforced; customers may purchase one zIIP and/or one zAAP for every general purpose processor they purchase on the server.



Work Unit Queue Distribution – Monitor I CPU Report

SYSTEM ADDRESS SPACE AND WORK UNIT ANALYSIS

QUEUE TYPES	MIN	MAX	AVG
IN	550	1,008	594.8
IN READY	4	438	22.7
OUT READY	0	1	0.0
OUT WAIT	0	0	0.0
LOGICAL OUT RDY	0	628	10.3
LOGICAL OUT WAIT	178	634	589.4

ADDRESS SPACE TYPES

BATCH	281	284	282.0
STC	708	763	736.4
TSO	97	98	97.9
ASCH	0	1	0.0
OMVS	43	97	68.0

-----NUMBER OF WORK UNITS-----

CPU TYPES	MIN	MAX	AVG
CP	444	888	555.5
AAP	22	33	28.8
IIP	0	0	0.0

-----DISTRIBUTION OF IN-READY WORK UNIT QUEUE-----

NUMBER OF WORK UNITS	(%)
<= N	55.5
= N + 1	4.4
= N + 2	4.0
= N + 3	3.7
<= N + 5	7.1
<= N + 10	14.7
<= N + 15	8.0
<= N + 20	1.9
<= N + 30	0.2
<= N + 40	0.0
<= N + 60	0.0
<= N + 80	0.0
<= N + 100	0.0
<= N + 120	0.0
<= N + 150	0.0
> N + 150	0.0


N = NUMBER OF PROCESSORS ONLINE UNPARKED (22.4 ON AVG)

Work unit count on CPU type level

16 Buckets representing the work unit count with regard to the number of online processors

What I hope I covered.....

- What are dispatchable units of work on z/OS
- How WLM manages dispatchable units of work
- The role of HiperDispatch
- Dispatching work to zIIP and zAAP engines
- RMF Work Unit Analysis Report



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 SEs 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 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.

SmarterComputing



The World of z/OS Dispatching on the zEC12 Processor

Glenn Anderson, IBM Lab Services and Training



Summer SHARE 2013
Session 14040



© 2013 IBM Corporation