



E45

z/VSE V3.1 and SCSI Performance Update

Ingo Franzki (ifranzki@de.ibm.com)

zSeries® EXPO

**FEATURING Z/OS, Z/VM, Z/VSE
AND LINUX ON ZSERIES**

September 19 - 23, 2005

San Francisco, CA

Trademarks

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

CICS*	IBM*	Virtual Image
DB2*	IBM logo*	Facility
DB2 Connect	IMS	VM/ESA*
DB2 Universal	Intelligent	VSE/ESA
Database	Miner	VisualAge*
e-business logo*	Multiprise*	VTAM*
Enterprise Storage	MQSeries*	WebSphere*
Server	OS/390*	xSeries
HiperSockets	S/390*	z/Architecture
	SNAP/SHOT	z/VM
	*	z/VSE
		zSeries

* Registered trademarks of IBM Corporation

The following are trademarks or registered trademarks of other companies.

LINUX is a registered trademark of Linus Torvalds

Tivoli is a trademark of Tivoli Systems Inc.

Java and all Java-related trademarks and logos are trademarks of Sun Microsystems, Inc., in the United States and other countries

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.

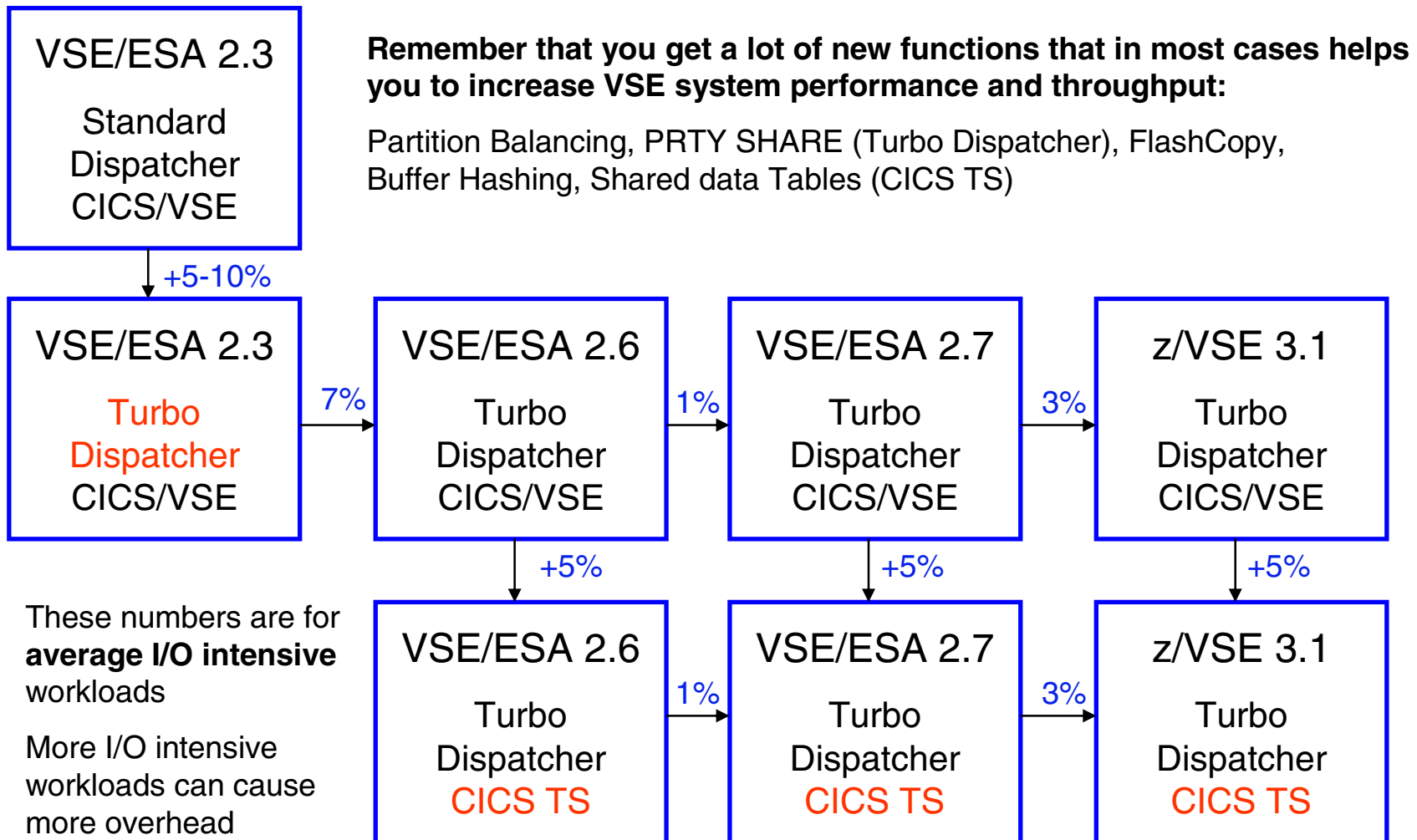
SET and Secure Electronic Transaction are trademarks owned by SET Secure Electronic Transaction LLC.

Intel is a registered trademark of Intel Corporation.

Agenda

- **z/VSE V3.1 Performance considerations**
 - Release overhead
- **SCSI Performance considerations**
 - Overhead
 - Basics
 - Migration aspects
- **Hardware support**
- **Turbo dispatcher**

Overhead Deltas for VSE Releases



Agenda

- **z/VSE V3.1**
 - Release overhead
- **SCSI Performance considerations**
 - Overhead
 - Basics
 - Migration aspects
- **Hardware support**
- **Turbo dispatcher**

Hardware and software requirements for SCSI

- **IBM eServer zSeries 800, 900, 890 or 990**
- **IBM zSeries FCP Adapter**
 - Microcode Level:
 - z800 und z900: J11233.015 or higher
 - z890 und z990: J13471.004 or higher
- **FCP Switch (e.g. IBM 2109)**
- **IBM TotalStorage Enterprise Storage Server (ESS)**
 - Microcode Level: 2.3.1 or higher
- **IBM TotalStorage DS6000 or DS8000**
- **z/VSE Version 3 Release 1**
- **z/VM 4.4. or higher (only if VSE runs under VM)**

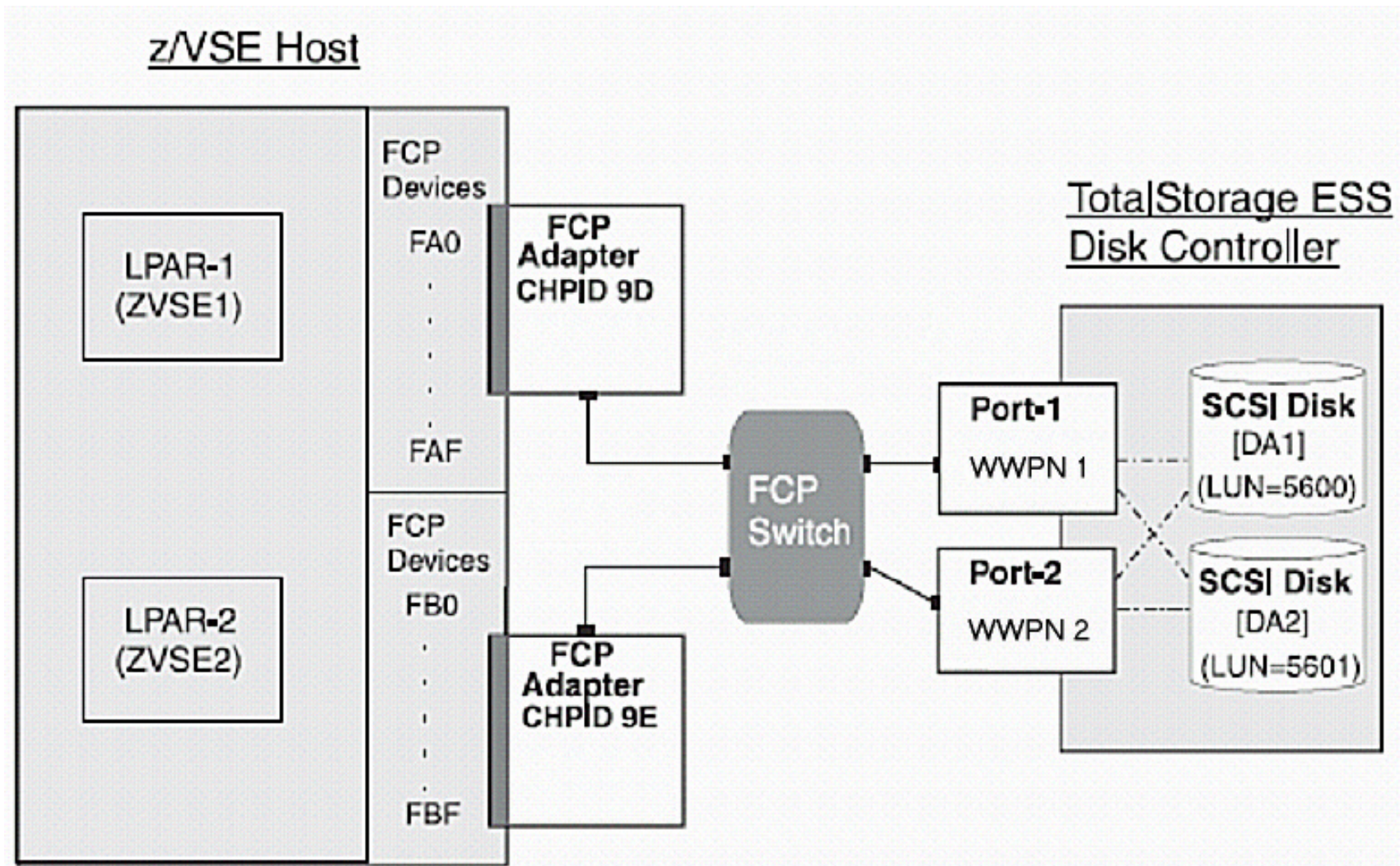
Hardware and software requirements for SCSI (2)

- **IPL von SCSI**
 - CPU Feature Code 9904
 - z800 and z900:
 - Microcode Level EC-Number J12811 or higher
 - z890 and z990:
 - Microcode Level EC-Number J12221 or higher
- **IPL von SCSI under z/VM 4.4**
 - z/VM Service Level:
 - UM31181 (English)
 - UM31180 (German)
 - UM31179 (Kanji)
- **Emulated FBA Disks:**
 - z/VM 5.1

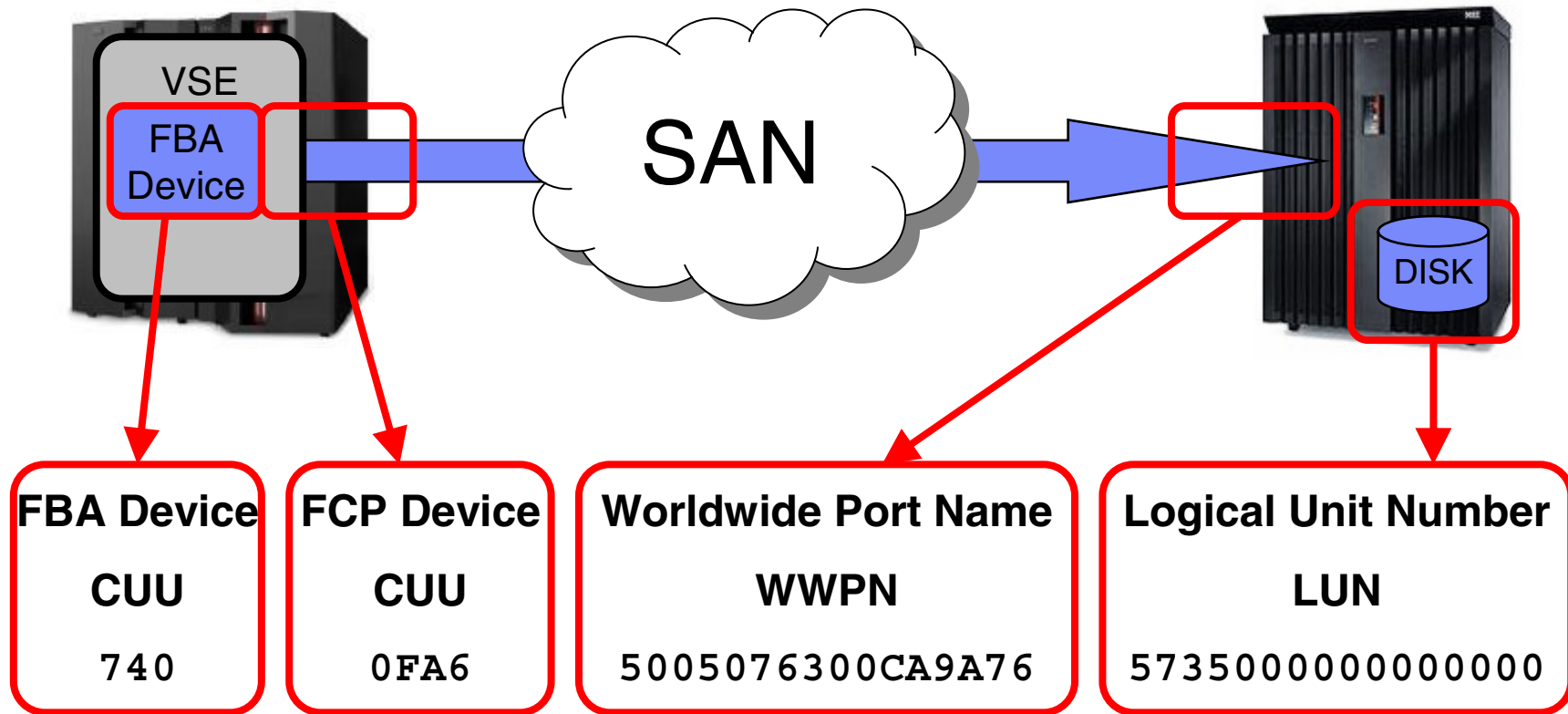
SCSI disk characteristics

- **Size of a SCSI Disk**
 - Minimum 8 MB
 - Maximum about. 24 GB
 - 4 MB are used internally from z/VSE
 - Usable size = size – 4 MB
 - VSAM can only use the first 16 GB
- **Model**
 - SCSI Disks are defined as FBA Devices
 - 9336 Model 20
- **Block sizes**
 - z/VSE supports only SCSI Disks with a block size of 512 Bytes
- **Standards**
 - SCSI Disks must support ANSI SCSI Version 3
- **z/VM SCSI to FBA emulation**
 - Maximum size of 2 GB

SCSI components



SCSI Addressing



SCSI setup for z/VSE



- **FCP Devices:**
 - ADD 4A7:4A9,FCP
- **FBA Devices:**
 - ADD 608:61B,FBA
- **Define SCSI:**
 - DEF SCSI,FBA=608,FCP=4A7,WWPN=5005076300CA9A76,LUN=5710
 - DEF SCSI,FBA=609,FCP=4A7,WWPN=5005076300CA9A76,LUN=5711
- **IPL from SCSI (VM)**
 - Minimum 32M Memory
 - SET LOADDEV PORT 50050763 00CE9A76 LUN 57350000 00000000
 - IPL 4B8 (IPL from FCP device)

SCSI commands

▪SCSI definitions during IPL:

```
DEF SCSI, FBA=cuu, FCP=cuu, WWPN=nnnnnnnnnnnnnnnnnn, LUN=nnnn
```

▪SCSI definitions online:

```
SYSDEF SCSI, FBA=cuu, FCP=cuu, WWPN=nnnnnnnnnnnnnnnnnn, LUN=nnnn
```

- FBA Device and FCP Device must have been already defined during IPL

▪Delete SCSI definitions:

```
SYSDEF SCSI, DELETE, FBA=cuu, FCP=cuu, WWPN=nnnnnnnnnnnnnnnnnn, LUN=nnnn
```

▪Display SCSI definitions:

```
QUERY SCSI
```

AR	FBA-CUU	FCP-CUU	WORLDWIDE PORTNAME	LOGICAL UNIT NUMBER
AR 0015	608	4A7	5005076300CA9A76	5710000000000000
AR 0015	609	4A7	5005076300CA9A76	5711000000000000
AR 0015	60A	4A7	5005076300CA9A76	5712000000000000
AR 0015	60B	4A7	5005076300CA9A76	5713000000000000
AR 0015	60C	4A7	5005076300CA9A76	5714000000000000
AR 0015	60D	4A8	5005076300CA9A76	5715000000000000

```
QUERY SCSI,608 (FBA Device)
```

AR	FBA-CUU	FCP-CUU	WORLDWIDE PORTNAME	LOGICAL UNIT NUMBER
AR 0015	608	4A7	5005076300CA9A76	5710000000000000

Interactive Interface Dialog

```
TAS$ICME          HARDWARE CONFIGURATION AND IPL: DEF SCSI
Enter the required data and press ENTER.

FBA .....      DA1          cuu of the FBA-SCSI device
FCP .....      FA0          cuu of the FCP device
WWPN .....     5005076300CA9A76 World wide port name of the
                                remote controller
LUN .....      5600         Logical unit number of the SCSI

PF1=HELP      2=REDISPLAY  3=END
```

Interactive Interface Dialog

```
TAS$ICMD          HARDWARE CONFIGURATION AND IPL: DEF SCSI

Enter the required data and press ENTER.

OPTIONS: 1 = ADD          2 = ALTER
         5 = DELETE

  OPT   FBA      FCP      WWPN          LUN
  ---   ---     ---     ---             ---
  -     233     C01     5005076300C693CB  5176
  -     DA1     FA0     5005076300CA9A76  5600
  -
  -
  -
  -
  -
  -
  -
  -
  -
  -

PF1=HELP          2=REDISPLAY  3=END          5=PROCESS
```

SCSI messages

- **AR 0033 0S45I SCSI DEVICE 618 CONSISTS OF 03906304 BLOCKS, 03897432 BLOCKS ARE AVAILABLE, 680 BLOCKS ARE UNUSED**
 - Multiple of 777 Blocks
 - Internal Model: 1 „Cylinder“ = 777 Blocks
- **AR 0033 0S40I SCSI PROCESSING EVENT: REASON=0060
FUNCTION=INIT-SCSI FBA=609 FCP=4A7 WWPN=5005076300CA9A76
LUN=5711000000000000**
 - Message description shows the reason based on the Reason Codes
- **0S46I I/O ERROR ON FBA=600 FCP=C00 RC=01, REASON==052500**
 - SCSI I/O Error has been mapped to an FBA Error
 - Message description shows the reason based on the Reason Codes
- **SCSI I/O errors are mapped in S/390 I/O errors for user programs (e.g. Unit check)**
 - In addition, message 0S40I message is issued, to inform about the exact reason

SCSI multipathing

- **One or more alternative paths to the same SCSI Disk**
 - Increases the availability
 - NOT: Workload-Balancing
- **Each path must be defined over a different FCP adapter**
 - One FCP card can contain multiple FCP adapters (CHIPID)
 - To increase availability, you should use different FCP adapters on different physical FCP cards
- **As best, even over different switches and/or ports**
- **Example:**

```
DEF SCSI , FBA=DA1 , FCP=FA0 , WWPN=5005076300CA9A76 , LUN=5600
```

```
DEF SCSI , FBA=DA1 , FCP=FB0 , WWPN=5005076300C29A76 , LUN=5600
```

- **QUERY SCSI**

```
AR 0015 FBA-CUU FCP-CUU WORLDWIDE PORTNAME LOGICAL UNIT NUMBER
AR 0015 DA1      FA0      5005076300CA9A76 5600000000000000
AR 0015 DA1      FB0      5005076300C29A76 5600000000000000
```

- The first path is currently used to access the SCSI disk

Sharing SCSI disks

- **ADD cuu,FBA,SHR (FBA Device)**
- **Lockfile is used (see DLF)**
 - Using Reserve/Release SCSI Command (internal)
 - Reserve is based on FCP Adapter
 - Release must be done from same FCP Adapter
 - Possible problem:
 - Hardwait during disk is reserved
 - Disk stays reserved
 - The system tries to release the disk during hardwait processing, but this may fail
- **Restrictions**
 - Lockfile can not reside on DOSRES/SYSWK1 (only for SCSI)
 - No Multipathing possible for Lockfile-Disks
 - Each VSE System must access the Lockfile using its own FCK CHPID
- **Suggestion**
 - Use separate Disk for Lockfile

Base installation on SCSI disks

- **Only base installation possible**
 - FSU from ECKD- to SCSI-Disks not possible
- **Automatic installation**
 - IPL from Tape

```
BG 0000 SI70D IF YOU WANT TO USE SCSI DEVICES SPECIFY YES, ELSE NO
```

```
0 YES
```

```
BG 0000 SI75I ENTER SCSI COMMAND FOR DOSRES
```

```
BG 0000 SA80D SCSI, FBA=CUU, FCP=CUU, WWPN=PORTNAME, LUN=LUN
```

```
0 SCSI, FBA=608, FCP=C00, WWPN=5005076300C69A76, LUN=5745
```

```
AR 0033 0S45I SCSI DEVICE 608 CONSISTS OF 09765632 AVAILABLE, 651 BLOCKS ARE UNUSED
```

```
BG 0000 SA76I ENTER SCSI COMMAND FOR SYSWK1
```

```
BG 0000 SA80D SCSI, FBA=CUU, FCP=CUU, WWPN=PORTNAME, LUN=LUN
```

```
0 SCSI, FBA=609, FCP=D00, WWPN=5005076300C29A76, LUN=5746
```

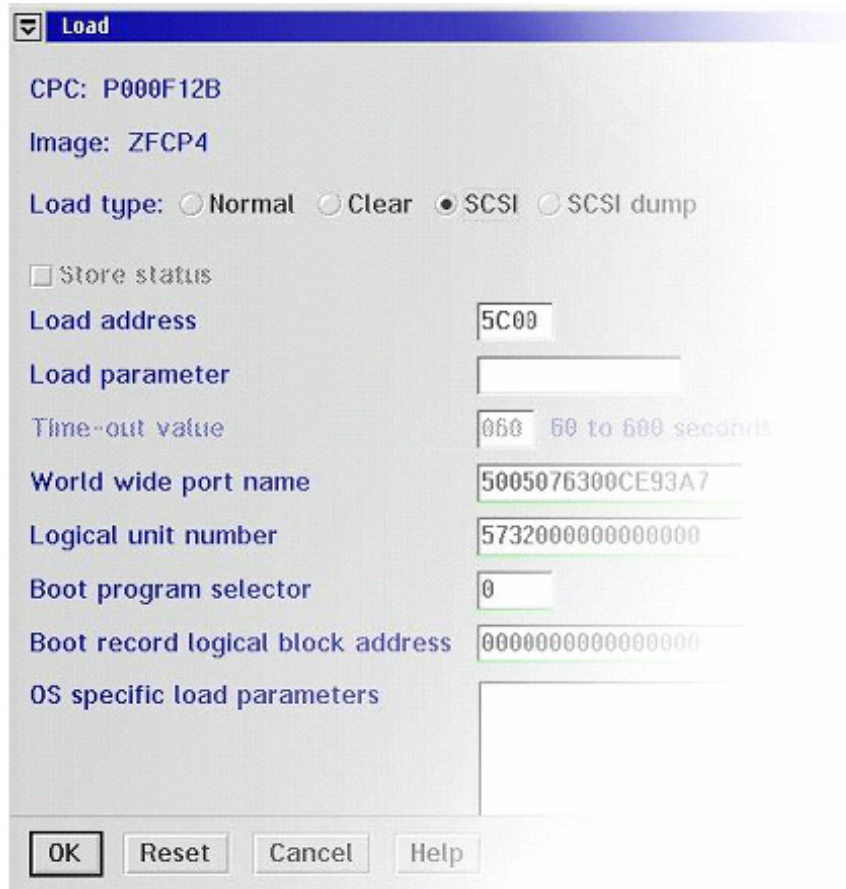
```
AR 0033 0S45I SCSI DEVICE 609 CONSISTS OF 09765632 AVAILABLE, 651 BLOCKS ARE UNUSED
```

```
BG 0000 SI08I DOSRES IS 608, DEVICE TYPE FBA
```

```
BG 0000 SI09I SYSWK1 IS 609, DEVICE TYPE FBA
```

- **Hardware Configuration Dialog**
 - Press PF5 to catalog the IPLPROC
 - Otherwise the next IPL will fail because it does not find SYSWK1

IPL from SCSI



Load

CPC: P000F12B
Image: ZFCP4

Load type: Normal Clear SCSI SCSI dump

Store status

Load address: 5C00

Load parameter:

Time-out value: 060 60 to 600 seconds

World wide port name: 5005076300CE93A7

Logical unit number: 5732000000000000

Boot program selector: 0

Boot record logical block address: 0000000000000000

OS specific load parameters:

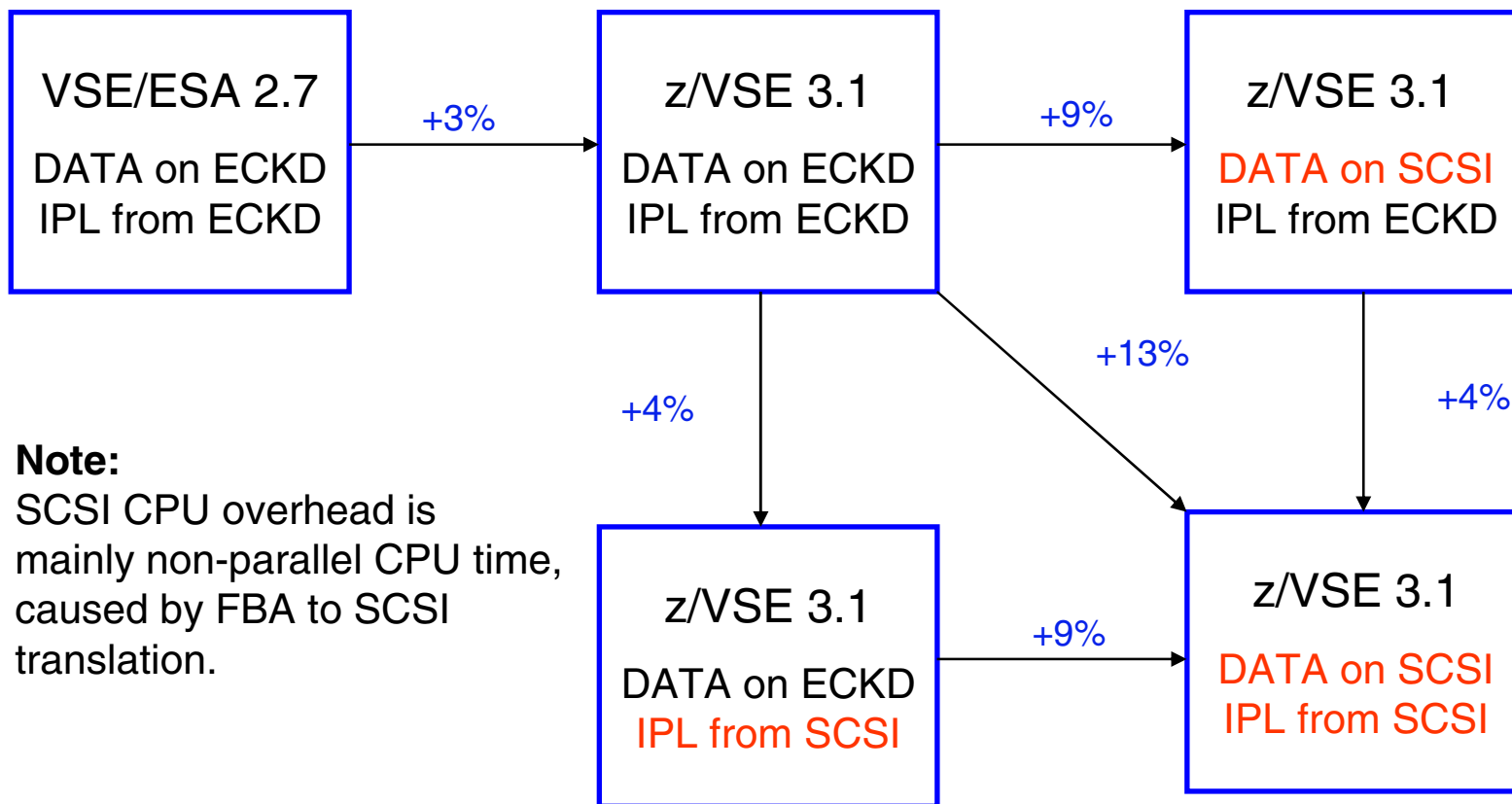
OK Reset Cancel Help

- **Uses the „Machine Loader“**
 - Platform independent Hardware-Tool
- **Native or LPAR**
 - Perform a Load using Hardware Management Console (HMC)
 - Load Address = FCP Device
 - WWPN
 - LUN number
- **Under z/VM**
 - SET LOADDEV PORTNAME 5005076300C29A76 LUN 56010000 00000000
 - IPL cuu (FCP Device)

Migration from ECKD to SCSI

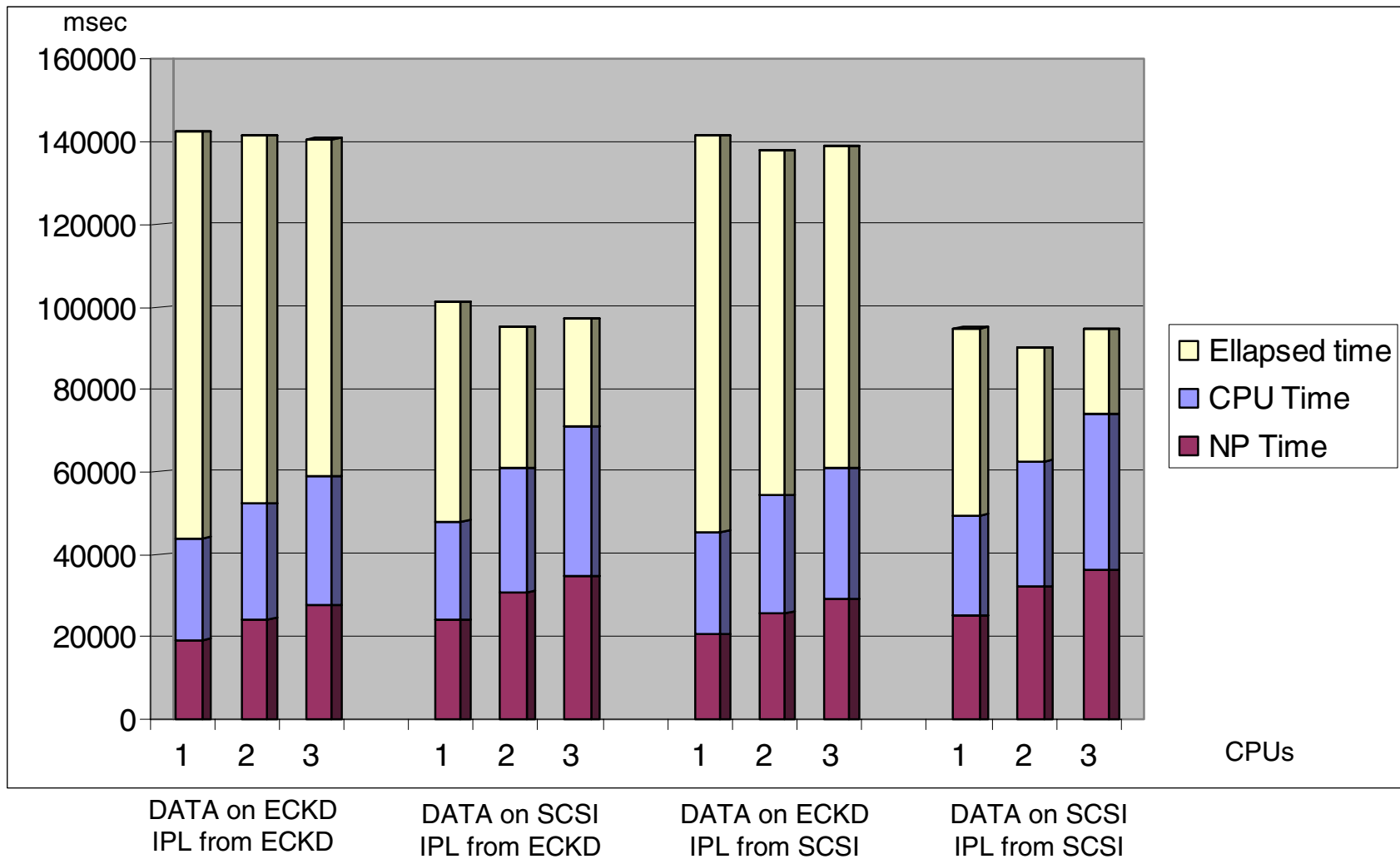
- **FSU from ECKD to SCSI not possible**
 - Only base installation
- **File allocations must be adapted**
 - Tracks/Cylinder into Blocks (for 3390):
 - 1 Track = about 112 Blocks
 - 1 Cylinder = about 1680 Blocks
 - VSAM Space
 - Hint: Specify cluster sizes in RECORDS, not Tracks
 - Sequential files
 - VSE Libraries
 - Hint: 1 LIBR Block = 1024 Bytes = 2 SCSI Blocks
- **Programs must be able to work with FBA disks**
 - As best, implement it device independent

Overhead Deltas for SCSI



Note:
SCSI CPU overhead is mainly non-parallel CPU time, caused by FBA to SCSI translation.

SCSI Overhead



SCSI I/O count considerations

- **For PACEX16 workload**
 - ECKD:
 - 18000 ECKD I/Os per disk
 - SCSI:
 - 20000 FBA I/Os per disk
 - ECKD → FBA:
 - 11% increased I/O counts
- **In general 1 FBA I/O is translated into 1 SCSI I/O**
- **Except for**
 - Long CCW chains
 - Overlapping addresses (e.g. PHASE loading)

Comparison ESCON – FICON/FCP

	ESCON	FICON/FCP
Max # of channels	256	4 x 256
Max # device addresses per link	4096	65536
Max # logical CU-paths per port	64	256
Device addresses per channel	1024	16384
Link rate	20 MB/sec	200 MB/sec
Max achievable transfer rate	17 MB/sec	170 MB/sec
Full duplex	No	Yes
Concurrent I/O operations	1	Up to 32
CCW execution	Synchrony	Asynchrony (FICON)

When should I (not) use SCSI?

- **When should I use SCSI?**

- When enough CPU power is available to handle the additional SCSI overhead.
 - SCSI overhead is mostly non-parallel code.

- **When should I not use SCSI?**

- When you are already CPU constraint.
 - If you are today running at 80% CPU utilization, SCSI would fill up your CPU up to 100 %

Agenda

- **z/VSE V3.1**
- **SCSI Performance considerations**
 - Overhead
 - Basics
 - Migration aspects
- **Hardware support**
- **Turbo dispatcher**

z/VSE 3.1 Hardware support

- **z/VSE 3.1 and VSE/ESA 2.7 runs on the following machines**
 - IBM System z9 109
 - zSeries: z800, z900, z990, z890
 - 9672 Parallel Enterprise Server (G5/G6)
 - Multiprice 3000 (7060)
 - equivalent emulators (Flex-ES)
- **z/VSE 3.1 and VSE/ESA 2.7 is based on the hardware instruction set described in the manual 'ESA/390 Principles of Operation' (SA22-7201).**
 - It is assumed that all the ESA/390 instructions and facilities described in that manual can be used.

z/VSE 3.1 Hardware support - continuation

- **z/VSE 3.1 is designed to support:**
 - IBM System z9 109
 - IBM eServer zSeries 800, 900, 890 and 990
 - SCSI disks attached to zSeries FCP channels
 - OSA-Express2 and FICON Express2 adapters
 - Crypto Express2 and CP Assist for Cryptographic Function (CPACF)
 - IBM TotalStorage 3494 Virtual Tape Server
 - improved support for IBM 3494 Tape Library
 - IBM TotalStorage DS8000 and DS6000 series Storage Servers
 - IBM TotalStorage Enterprise Storage Server (ESS)

Supported VSE Releases

VSE Release	Available	End of Marketing	End of Service
z/VSE 3.1	03/04/2005		
VSE/ESA 2.7	03/14/2003	09/30/2005	02/28/2007
VSE/ESA 2.6	12/14/2001	03/14/2003 (no longer orderable)	03/31/2006
VSE/ESA 2.5	09/29/2000	12/14/2001	12/31/2003 (out of service)
VSE/ESA 2.4	06/25/1999	09/29/2000	06/30/2002 (out of service)
VSE/ESA 2.3	07/12/1997	06/30/2000	12/31/2001 (out of service)

VSE Server Support

IBM Server	z/VSE 3.1	VSE/ESA 2.7	VSE/ESA 2.6	VSE/ESA 2.5	VSE/ESA 2.4/2.3
IBM System z9 109	Yes	Yes	Yes (PTF required)	Yes (PTF required)	No
zSeries 890, 990	Yes	Yes	Yes (PTF required)	Yes (PTF required)	No
zSeries 800, 900	Yes	Yes	Yes	Yes	Yes
S/390 Parallel Enterprise Server G5/G6	Yes	Yes	Yes	Yes	Yes
S/390 Multiprise 3000	Yes	Yes	Yes	Yes	Yes
S/390 Parallel Enterprise Server G3/G4	No	No	Yes	Yes	Yes
S/390 Multiprise 2000	No	No	Yes	Yes	Yes
S/390 Integrated Server	No	No	Yes	Yes	Yes
S/390 Parallel Enterprise Server G2 / G1 (out of Service)	No	No	Yes	Yes	Yes
ES/9000 – 9221, 9121, 9021 (out of Service)	No	No	Yes	Yes	Yes
P/390 and R/390 (out of Service)	No	No	Yes	Yes	Yes

VSE Hardware Support

VSE Release	HiperSockets	OSA Express (QDIO mode)	Hardware Crypto
z/VSE 3.1	Yes	Yes	Yes (PCICA, CEX2C, CPACF)
VSE/ESA 2.7	Yes	Yes	Yes (PCICA, CPACF)
VSE/ESA 2.6	No	Yes	No
VSE/ESA 2.5	No	No	No
VSE/ESA 2.4	No	No	No
VSE/ESA 2.3	No	No	No

Crypto Card	z800	z900	z890	z990
PCICA	No	Yes	Yes	Yes
CEX2C	No	No	Yes	Yes
CPACF	No	No	Yes	Yes
CEX2A	No	No	No	No

zSeries Remarks

- **Prior to zSeries there is one cache for data and instructions**
- **zSeries has splitted data and instruction cache**
- **Performance implications:**
 - If **program variables** and **code that updates** these program variables are **in the same cache line** (256 byte)
 - Update of program variable invalidates instruction cache
 - Performance decrease if update is done in a loop
 - See APAR PQ66981 for FORTRAN compiler

zSeries Remarks - example

Not causing a problem:

```

LA      R1,PHASNAME    POINT AT PHASE NAME
CDDELETE (1)
+*     SUPERVISOR - CDDELETE - 5686-032-06
+      CNOP  0,4
+      BAL   15,*+8
+      DC    A(B'00010010')
+      L     15,0(,15)
+      SVC   65          ISSUE SVC FOR CDDELETE
+      DS    0H

```

CDDELETE uses an inline flag byte,
but does not modify it

Can cause a problem:

```

WTO TEXT=DATA
+      CNOP  0,
+      BAL   1,IHB0003A  BRANCH AROUND MESSAGE
+      DC    AL2(8)      TEXT LENGTH
+      DC    B'00000000000010000'  MCSFLAGS
+      DC    AL4(0)      MESSAGE TEXT ADDR
...
+IHB0003A DS    0H
+      LR    14,1        FIRST BYTE OF PARM LIST
+      SR    15,15       CLEAR REGISTER 15
+      AH    15,0(1,0)   ADD LENGTH OF TEXT + 4
+      AR    14,15       FIRST BYTE AFTER TEXT
+      LA    15,DATA     LOAD TEXT VALUE
+      ST    15,4(0,1)   STORE ADDR INTO PLIST
+*     SUPERVISOR - SIMSVC - 5686-032
...
+      SVC   35          ISSUE SVC 35
@GE00016 DS    0H

```

WTO uses an inline parameter list,
but modifies the parameter list

Note: WTO can be coded with an external
parameter list: WTO ...,MF=(E,addr)

Z890, z990 and z9-109 Considerations

- **The z890, z990 and z9-109 are LPAR-only machines**
 - No basic mode any more
 - Even if you run just one VSE system, it now runs in an LPAR
 - Running VSE systems under z/VM means
 - running VSE in z/VM in an LPAR
 - No I/O Assist in LPARs
 - Only available if z/VM runs in basic mode, but no basic mode available on z890, z990, z9-109

z/VM 5.1 considerations

- **z/VM 5.1 no longer supports V=R and V=F guests**
- **z/VM 5.1 no longer support I/O Assist**
 - If you currently run with preferred guests, you will need to estimate and plan for a likely increase in processor requirements as those preferred guests become V=V guests as part of the migration.
 - Refer to Preferred Guest Migration Considerations at <http://www.vm.ibm.com/perf/tips/z890.html> for assistance and background information
- **How to size the impact (on your current system)**
 - **Loss of I/O Assist:** Run your workload with CP SET IOASSIST OFF and measure the increase
 - **Loss of V=R/F:** Run your workload with V=V and use the CP Monitor to watch for increased CPU consumption
- **How to tune**
 - **Dedicated processors:** CP SET SHARE ABSOLUTE
 - **Dedicated memory:** CP SET RESERVED
 - **I/O Assist:** Use minidisks, turn minidisk caching on (MDC)

Possible performance problem with PPRC

- **Problem can occurs if**
 - PPRC is used
 - VSE runs in native or in LPAR
 - Not all devices that are defined in IOCP are also defined in VSE ADD statements
- **In case there is an PPRC state change, interrupts are sent to all LPARs where the related device are defined in IOCP.**
 - If the device is defined in VSE ADD, no problem occurs: VSE will process the interrupt correctly.
 - If the device is NOT defined in VSE ADD, the interrupt is ignored by VSE and the interrupt is resent very quickly to that LPAR
 - Results in very high channel activity (up to 100%)
- **Solution:**
 - Define ALL devices in VSE ADD that are defined in IOCP

VSE/POWER POFFLOAD Performance Problems

- **Caused by ‘incompatibility’ between VSE/POWER tape format and new tape drives**
- **3490F empties cache for FSF used by POFFLOAD LOAD**
 - Install **DY46164/DY46245** for VSE/ESA 2.7/2.6
- **3590 synchronizes cache with tape for each WTM**
 - Install microcode **FC0520** on A60 controller + VSE/AF **APAR DY45817** + AR command **TAPE WTM=NOSYNC**
 - Unfortunately controller A50 is too small to install FC0520

Agenda

- **z/VSE V3.1**
- **SCSI Performance considerations**
 - Overhead
 - Basics
 - Migration aspects
- **Hardware support**
- **Turbo dispatcher**

Turbo Dispatcher - Overview

- **Turbo Dispatcher**
 - available since 1995
 - VSE/ESA 2.1-2.3 Standard and Turbo Dispatcher
 - since VSE/ESA 2.4 only Turbo Dispatcher
 - last changes:
 - VSE/ESA 2.6.2 (APAR DY45869)
 - VSE/ESA 2.7.0 (APAR DY45926)
 - Supports basic (native), LPAR and VM mode
 - Runs on Uni- and n-Way-procercssors
 - CPUs have "equal" rights
 - more than 3 CPUs are not recommended

Turbo Dispatcher - Overview (2)

- **IPL is done on 1 CPU only**
 - after IPL other CPUs can be started
 - CPUs can be started or stopped without re-IPL
 - at least 1 CPU (IPL CPU) must always be active

`SYSDEF TD,START=n|ALL`

`SYSDEF TD,STOP=n|ALL`

`SYSDEF TD,STOPQ=n|ALL`

`QUERY TD`

Turbo Dispatcher - Quiesced CPUs

- **SYSDEF TD,STOPQ=n to set a CPU in quiesced mode**
 - Implemented for z/VM guest systems
 - Not started guest CPUs stop IOASSIST
 - STOPQ remains IOASSIST active, and avoids TD Overhead, (CPU will no longer participate in work unit selection)
 - quiesced CPUs will not process any work units
 - quiesced CPUs will not handle any interrupt
 - quiesced CPUs can be started with SYSDEF TD,START

Turbo Dispatcher - Design

- **TD dynamically assigns partitions to CPUs**
 - Work unit = from assignment to one CPU until next interrupt/SVC
 - If one task (subtask) of a partition is active, no other task of the same partition will be selected
 - TD dispatches on partition-basis, not on task-basis
 - A job running in a partition is processed in several work units.

Turbo Dispatcher - Design (2)

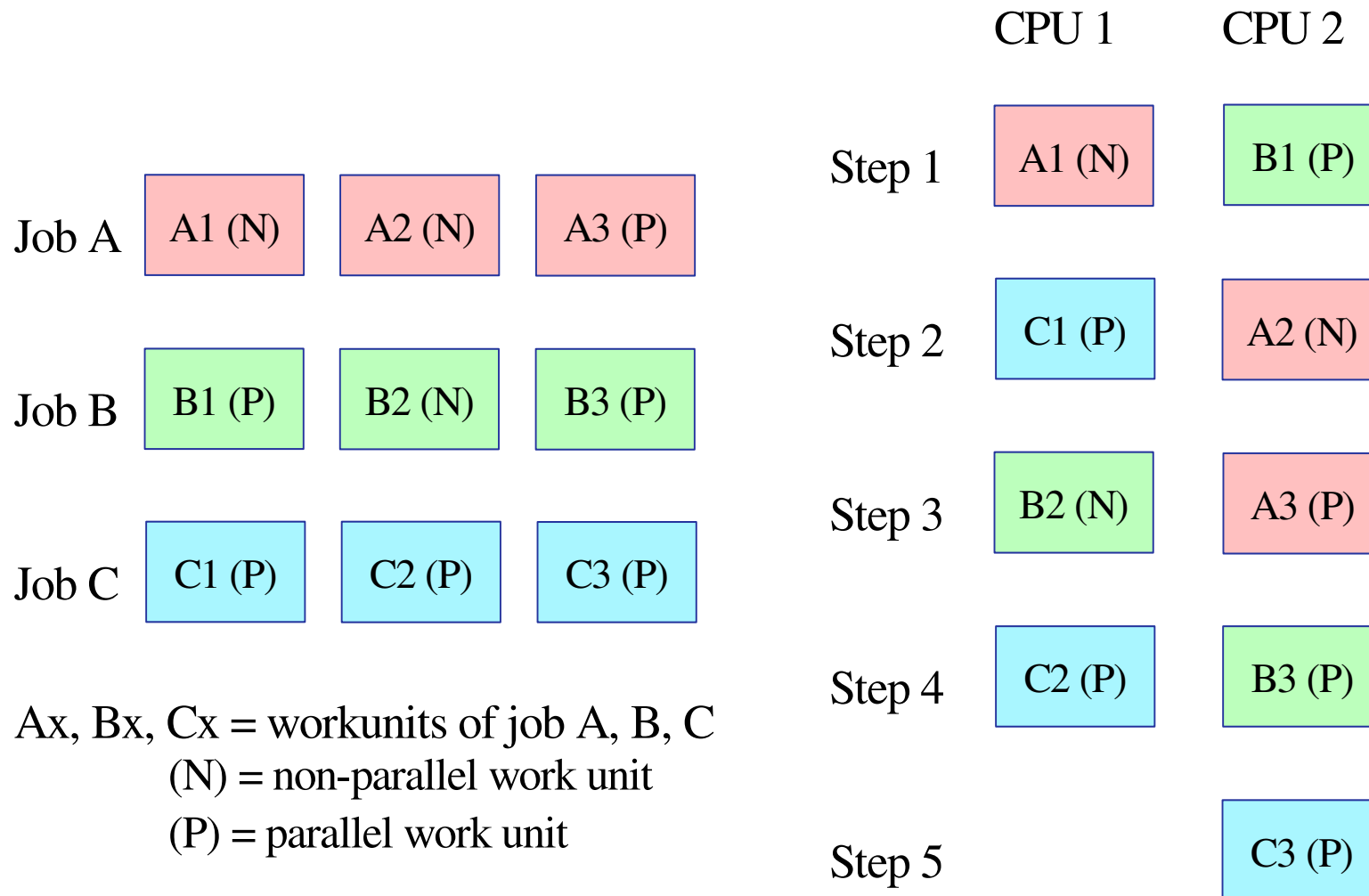
- **parallel work units**

- application code (CICS, Batch)
- may run on any CPU concurrently with other parallel or non-parallel work units.

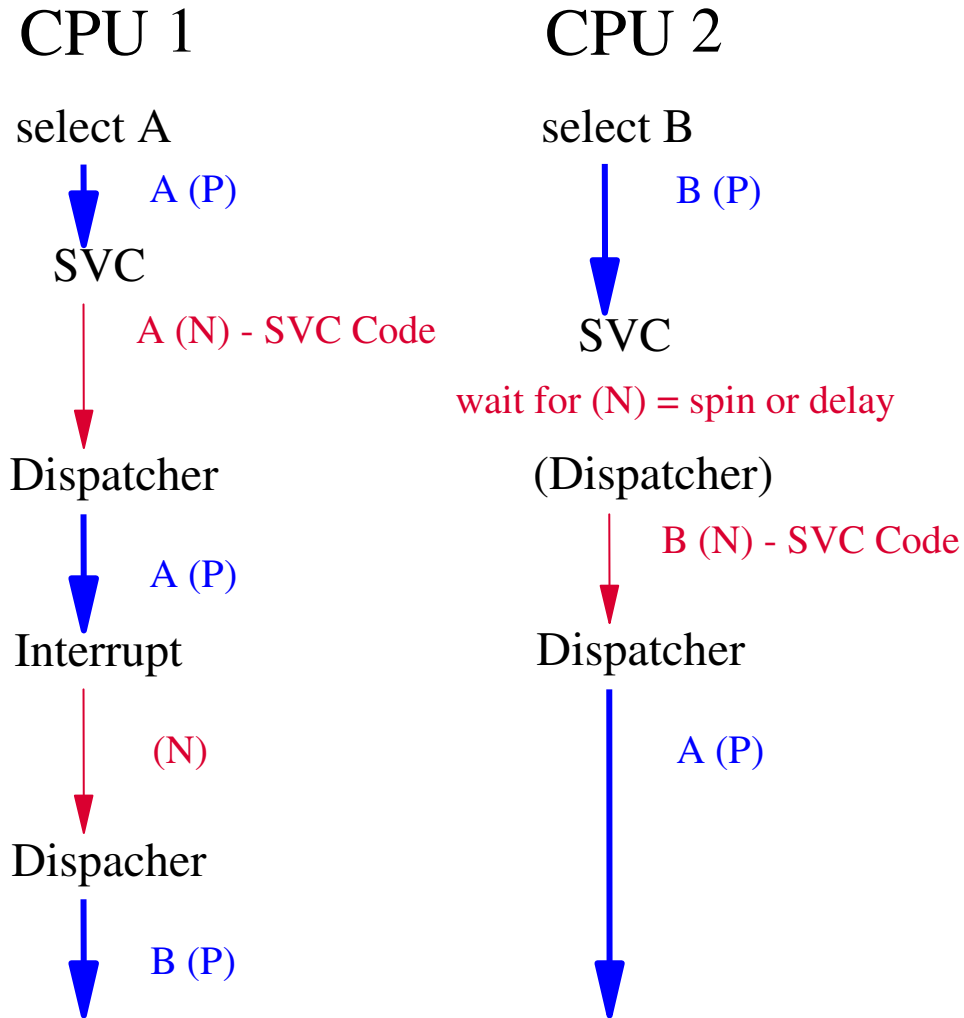
- **non-parallel work units**

- system code (Services, VTAM, Vendor code)
- As long as one non-parallel work unit is active on one CPU, no other non-parallel work unit can execute on any other CPU.

Turbo Dispatcher - Design - Example 1



Turbo Dispatcher - Design - Example 2



Turbo Dispatcher - Exploitation

■ **Uni-Processor**

- new Partition Balancing Concept
 - Helps to set priorities of partitions
- Determination of non-parallel share, to find out if a 2. or 3. CPU would be of use

■ **n-Way Processors (2-3 CPUs)**

- System tuning required for exploitation
- Increased Capacity (dependent on workload)
 - Exploitation increases by reduction of non-parallel work units

Turbo Dispatcher - CPU time measurement

- **CPU time measurement (overall system)**
 - SYSDEF TD,RESETCNT
 - Workload (e.g. run a job)
 - QUERY TD (QUERY TD,INTERNAL)

CPU	STATUS	SPIN_TIME	NP_TIME	TOTAL_TIME	NP/TOT
00	ACTIVE	0	237100	416698	0.568
01	ACTIVE	0	157556	415229	0.379
02	QUIESCED	0	0	0	*.***
03	INACTIVE				

TOTAL		0	394656	831927	0.474
			NP/TOT: 0.474	SPIN/ (SPIN+TOT) :	0.000
			OVERALL UTILIZATION: 179%	NP UTILIZATION:	85%
			ELAPSED TIME SINCE LAST RESET:	463433	

NP/TOT = non-parallel share (NPS)
 SPIN_TIME = CPU time waiting for NP

Display System Activity Dialog

```

Session C - [32 x 80]
File Edit View Communication Actions Window Help
IESADMDA DISPLAY SYSTEM ACTIVITY 15 Seconds 13:55:26
*---- SYSTEM (CPUs: 1 / 0) ----* *----- CICS : DBDCCICS -----*
| CPU : 0% I/O/Sec: 1 | | No. Tasks: 7,018 Per Second : *
| Pages In : 0 Per Sec: * | | Dispatchable: 0 Suspended : 3
| Pages Out: 0 Per Sec: * | | Peak Active : 7 MXT reached: 0
*-----* *-----*
Priority: Z,Y,S,R,P,C,BG,FA,F9,F8,F6,F5,F4,F2,F7,FB,F3,F1

ID S JOB NAME PHASE NAME ELAPSED CPU TIME OVERHEAD %CPU I/O
F1 1 POWSTART IPWPOWER 29:23:33 1.23 .37 6,000
F3 3 VTAMSTRT ISTINCVT 29:23:28 18.13 5.65 304,230
FB 8 SECSERV BSTPSTS 29:23:33 .03 .01 213
*F7 7 TCPIP00 IPNET 29:23:28 1.61 .77 814
F2 2 CICSICCF DFHSIP 29:23:28 597.71 169.82 8,718
F4 4 <=WAITING FOR WORK=> .00 .00 2
F5 5 <=WAITING FOR WORK=> .00 .00 2
F6 6 <=WAITING FOR WORK=> .00 .00 2
F8 8 <=WAITING FOR WORK=> .00 .00 2
F9 9 <=WAITING FOR WORK=> .00 .00 2
FA A <=WAITING FOR WORK=> .00 .00 2
BG 0 <=WAITING FOR WORK=> .00 .00 2
PF1=HELP 2=PART.BAL. 3=END 4=RETURN 5=DYN.PART 6=CPU
    
```



Migration aspects

- **Consider hard-/software requirements:**
 - Does my largest partition still fit into a single CPU of the target processor?
 - Note: a partition can only run on 1 CPU at a time!
 - Is the processor capacity and speed still sufficient to run the workload?
 - Does multiprocessing help to run the workload?
 - What about non-parallel share (on 1-Way)?
 - Are there many parallel batch jobs?
 - A large CICS partition does not benefit of a 2. CPU

Migration overhead

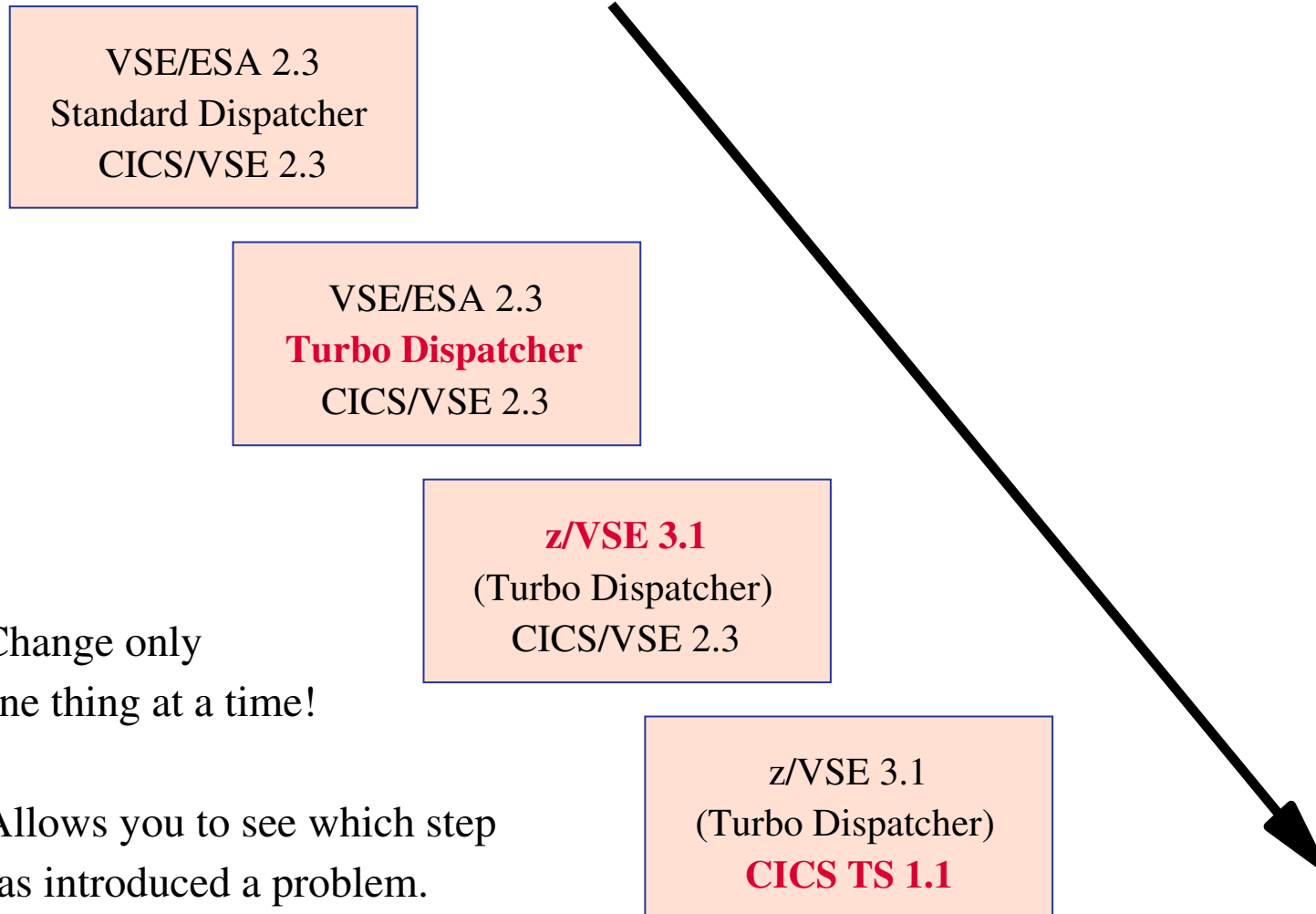
■ **Uni-Processor**

- increased overhead because of
 - Release migration (VSE/ESA 2.6/2.7 vs. z/VSE 3.1)
 - TD overhead (Standard Dispatcher vs. TD)
 - CICS/VSE vs. CICS TS

■ **N-Way Processor**

- CPU time increases when migrating from uni to n-Way Processor (for the same workload)
 - For PACEX Workload: Factor 1.4 (2 CPUs)
 - TD overhead for multiprocessor exploitation
 - z/VM Overhead

Migration path



Change only
one thing at a time!

Allows you to see which step
has introduced a problem.

Performance Tips

- **A partition can only exploit 1 CPU at a time**
 - 2 CPUs do not have any benefit for a CICS partition
 - Use as many partitions as required for selected n-way
- **Use/define only as many CPUs as really needed**
 - additional CPUs create more overhead, but no benefit
- **Partitions setup**
 - Set up more batch and/or (independent) CICS partitions
 - Split CICS production partitions into multiple partitions

Performance Tips (2)

- **1 CPU must be able to handle all non-parallel workload**
- **Non-parallel code limits the n-Way exploitation**
 - QUERY TD: $NP/TOT = NPS$
 - Measure NPS before migration
 - **max CPUs = $0.9 / NPS$**

NPS	#CPUs	NPS	#CPUs
0.20	4.5 (4)	0.40	2.2 (2)
0.25	3.6 (3)	0.45	2.0 (2)
0.30	3.0 (3)	0.50	1.8 (1)
0.35	2.6 (2)	0.55	1.6 (1)

Performance Tips (3)

- **Non-parallel code limits the maximum MP exploitation**
- **System code (Key 0) increases non-parallel share**
 - Vendor code can have significant impact
- **Overhead increases when NP code limits throughput**
- **Data In Memory (DIM) reduces non-parallel code**
 - less system calls (I/Os)
 - may increase throughput
- **Change VSE/POWER startup to WORKUNIT=PA**
- **In general **ONE faster CPU** is better than multiple slower ones**
 - Even if sum of slower CPUs is higher than one faster CPU

CICS Implications

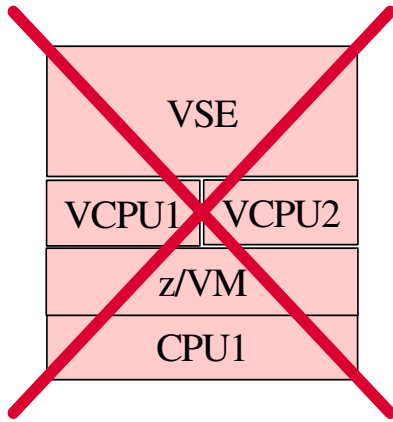
- **Single CICS**
 - Can consume processing power of one CPU only
 - parallel batch jobs may exploit 2. CPU
- **Multiple CICS partitions**
 - Number of CPUs depends on non-parallel share (NPS)
 - Function shipping and Transaction routing
 - AOR, TOR, FOR

Partition Balancing

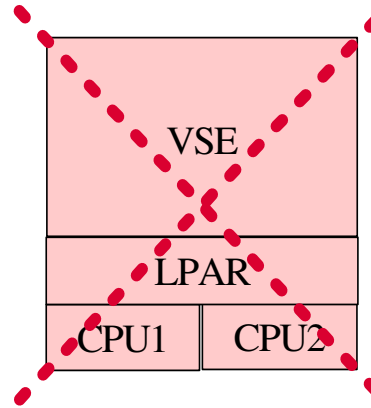
- **Balanced Group is defined with PRTY:**
 - PRTY BG,C=F5=F8,F2,F3,F1
 - Each partition/class of the group has a default-SHARE (100)
 - Dynamic partitions gets the SHARE of its class
- **To set a SHARE (1-1999)**
 - PRTY SHARE,F5=50
 - SHARE = 0 means the lowest priority within the group

```
PRTY
AR 0015 PRTY BG,C=F5=F8,F2,F3,F1
AR 0015
AR 0015 SHARE F5= 50, F8= 100, C= 100
MSECS
AR 0015 MSECS 976 <---- influences task selection
```

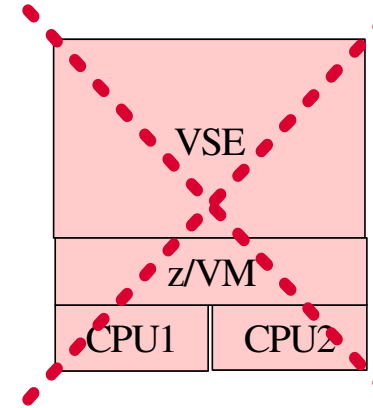
Do's and Don't Do's



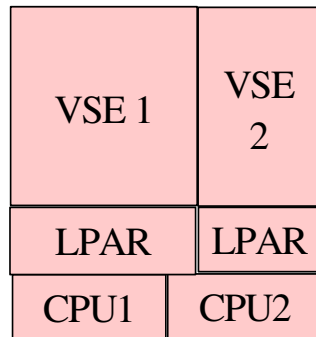
no virtual CPUs!
(creates overhead)



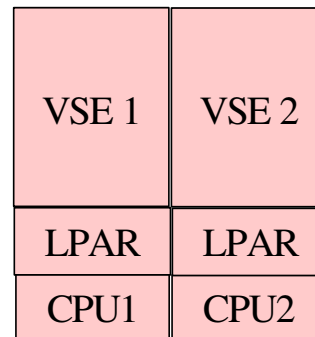
only if NPS < 4.5



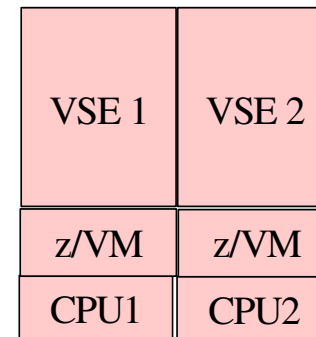
only if NPS < 4.5



VSE 1 = Production
VSE 2 = Test



dedicated CPU
per VSE



dedicated CPU
per VSE

Do's and Don't Do's (2)

**The fastest
uni-processor
is (almost always)
the best processor !**

VSE Health Check

- **Goals**
 - Recognize actual/upcoming problems
 - Optimize the system for new/current workload
- **A-B-C analysis**
 - A - concentrate on the essentials
 - 20 % work for 80 % results
 - B - more detailed analysis
 - 30 % work for 15 % results
 - C - analyze all details
 - 50 % work for 5 % results
- **A-B analysis takes about 2 days**
- **C analysis takes about 1 week**
- **Should be done about once a year**

VSE Health Check - continued

■ What should be checked?

- Processor (utilization, dispatching, z/VM, ...)
- DASD, Tapes (I/O rate, cache, ...)
- Network (network load, misrouted packets, ...)
- System software
 - Turbo Dispatcher (PRTY, PRTY SHARE, ...)
 - VSAM (CA/CI sizes, share options, buffers, ...)
 - CICS (MXT, DSA/EDSA sizes, SOS, ...)
 - Storage Layout (GETVIS 24, SVA, partitions, DSPACE, ...)
 - VTAM (buffer pool)
 - POWER (DBLK, DBLKGP, ...)
 - LE runtime options (Heap size, ...)
- Application software
- **New Tool: VSE Health Checker**
<http://www.ibm.com/servers/eserver/zseries/zvse/downloads/#healthchecker>

Hints and Tips for Performance

- **Try to exploit Turbo Dispatcher functions**
 - Priority settings
 - Partition balancing
 - Partition balancing groups
- **Use as much data in memory (DIM) as possible**
 - CICS Shared Data Tables
 - Large/many VSAM Buffers (with buffer hashing)
 - Virtual Disks
- **Switch tracing/DEBUG off for production**

Hints and Tips for Connector and TCP/IP-Performance

- **Reduce amount of data transferred**
 - Transfer only data that is needed
 - Issue only requests that are needed
- **Use connection pooling**
 - Reduce overhead of connection establishment
- **Performance of connectors depends on**
 - Network performance
 - Performance of "server"
 - Performance of "client" or middle tier
- **Reduce misrouted packets**
- **Use a packet filter**
 - Unwanted packets increases TCP/IP and CPU load

Documentation

- **z/VSE homepage:**
 - <http://www.ibm.com/servers/eserver/zseries/zvse/>
- **VSE Performance:**
 - <http://www.ibm.com/servers/eserver/zseries/zvse/documentation/performance.html>
- **z/VM homepage:**
 - <http://www.ibm.com/vm>
- **z/VM 5.1 Preferred Guest Migration Considerations**
 - <http://www.vm.ibm.com/perf/tips/z890.html>
- **IBM eServer zSeries 890 and 990:**
 - <http://www.ibm.com/servers/eserver/zseries/z890/>
 - <http://www.ibm.com/servers/eserver/zseries/z990/>
- **IBM TotalStorage DS8000 and DS 6000:**
 - <http://www.ibm.com/servers/storage/disk/ds8000/index.html>
 - <http://www.ibm.com/servers/storage/disk/ds6000/index.html>
- **IBM TotalStorage 3494 Virtual Tape Server:**
 - <http://www.ibm.com/servers/storage/tape/3494vts/index.html>
- **IBM 3494 Tape Library:**
 - <http://www.ibm.com/servers/storage/tape/3494/index.html>

Questions ?

