

V01

Virtualization Basics

Brian K. Wade, Ph.D. - bkw@us.ibm.com

IBM zSeries Expo
April 16-20, 2007
Munich, Germany

Trademarks

IBM @server zSeries

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

CICS*	IBM logo	Virtual Image Facility
DB2	MQSeries*	VM/ESA*
DB2 Connect	Multiprise*	VSE/ESA
DB2 Universal Database	OS/390	WebSphere
e-business logo*	RISC	z/OS
FICON	S/390	z/VM
HiperSockets	S/390 Parallel Enterprise Server*	zSeries
IBM*		

* Registered trademarks of the IBM Corporation

The following are trademarks or registered trademarks of other companies.

Lotus, Notes, and Domino are trademarks or registered trademarks of Lotus Development Corporation.

Tivoli is a trademark of Tivoli Systems Inc.

Linux is a registered trademark of Linus Torvalds.

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.

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.

IBM considers a product "Year 2000 ready" if the product, when used in accordance with its associated documentation, is capable of correctly processing, providing and/or receiving date data within and between the 20th and 21st centuries, provided that all products (for example, hardware, software and firmware) used with the product properly exchange accurate date data with it. Any statements concerning the Year 2000 readiness of any IBM products contained in this presentation are Year 2000 Readiness Disclosures, subject to the Year 2000 Information and Readiness Disclosure Act of 1998.

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.

Credits

IBM server zSeries

People who contributed ideas and charts:

- Alan Altmark
- Bill Bitner
- John Franciscovich
- Reed Mullen
- Brian Wade
- Romney White

Thanks to everyone who contributed!

Introduction

IBM  zSeries

We'll explain basic concepts of zSeries:

- Terminology
- Processors
- Memory
- I/O
- Networking

We'll see that z/VM *virtualizes* a zSeries machine:

- Virtual processors
- Virtual memory
- ... and so on

Where appropriate, we'll compare or contrast:

- PR/SM or LPAR
- z/OS
- Linux

Terminology

zSeries Architecture

IBM @server zSeries

Every computer system has an *architecture*.

- Formal definition of how the hardware operates
- It's the hardware's functional specification
- What the software can expect from the hardware
- *What it does*, not how it does it

IBM's book [z/Architecture Principles of Operation](#) defines zSeries architecture

- Instruction set
- Processor features (registers, timers, interruption management)
- Arrangement of memory
- How I/O is to be done

Different *models* implement the architecture in different ways.

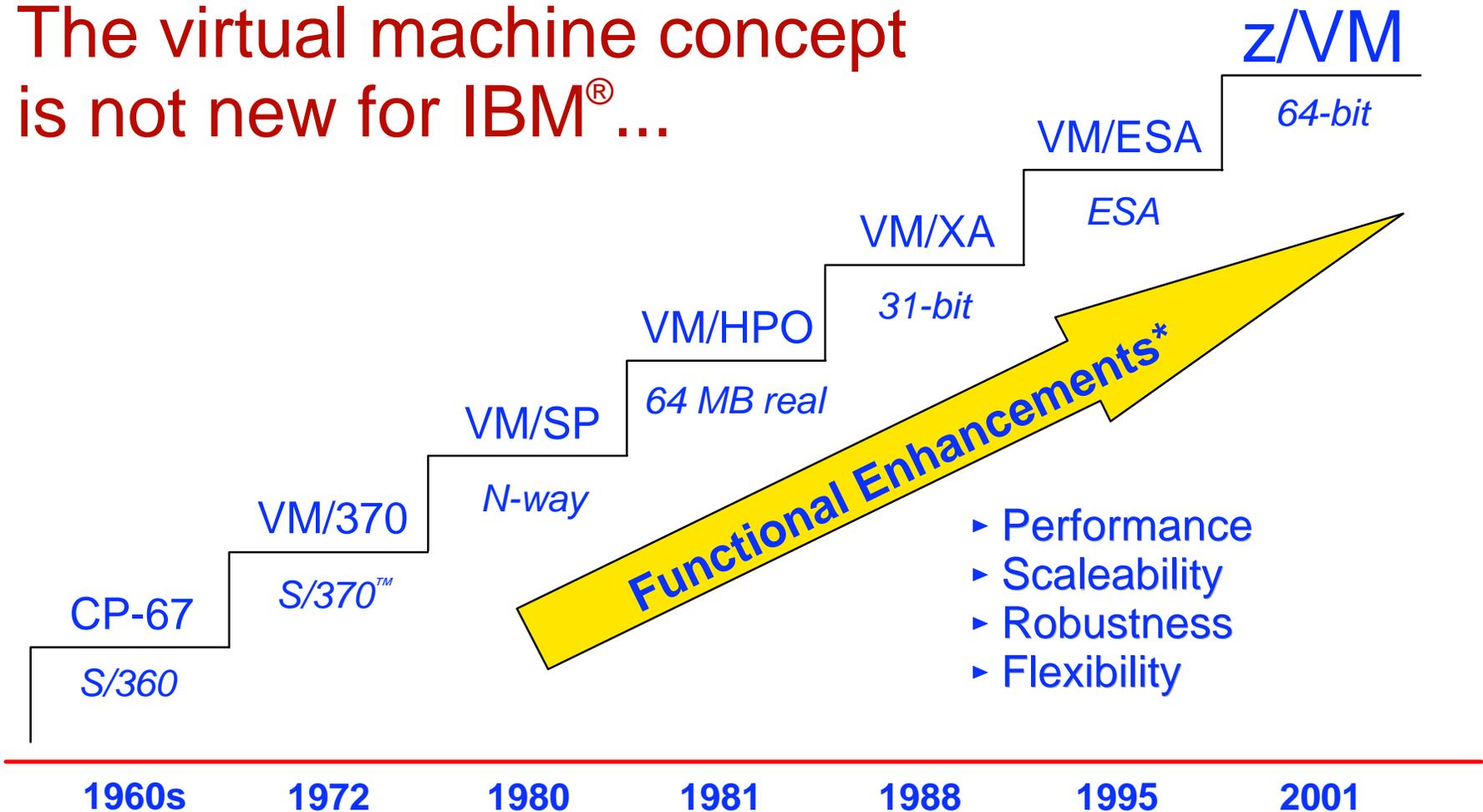
- How many processors there are
- How the processors connect to the memory bus
- How the cache is arranged
- How much physical memory there is
- How much I/O capability there is

z900, z990, and z890 are all *models* implementing z/Architecture.

IBM Virtualization Technology Evolution

IBM @server zSeries

The virtual machine concept
is not new for IBM® ...



* Investments made in hardware, architecture, microcode, software

IBM @server. For the next generation of e-business.

zSeries Parts Nomenclature

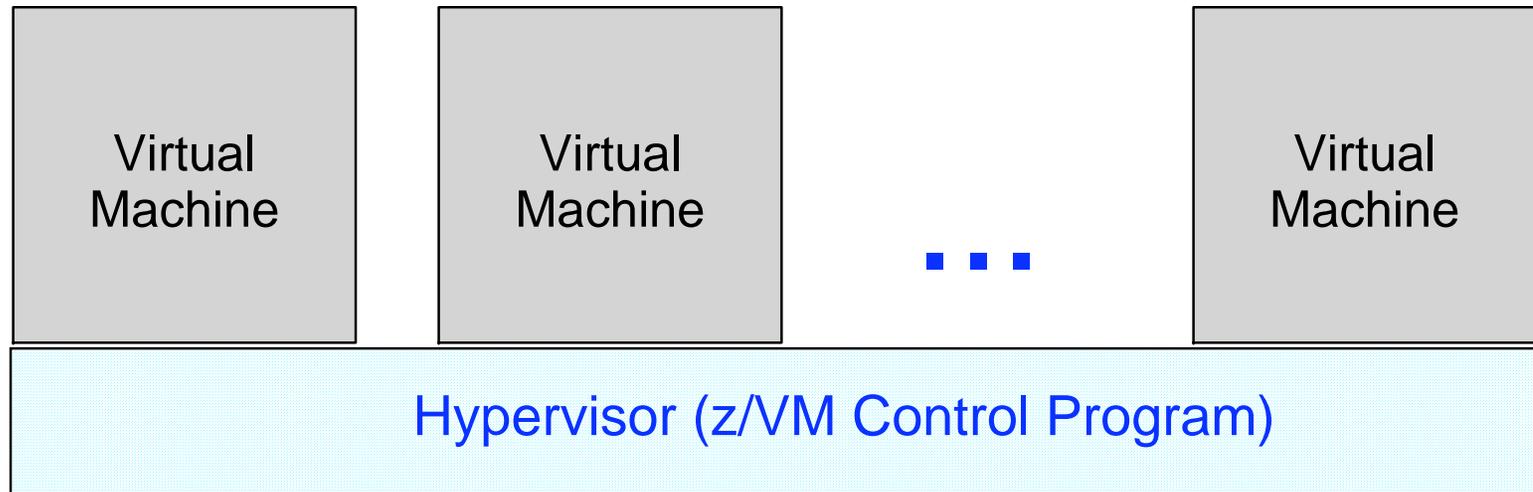
IBM @server zSeries

Intel, pSeries, etc.	zSeries
Memory	Storage (though we are moving toward "memory")
Disk, storage	DASD- Direct Access Storage Device
Processor	Processor, CPU (central processing unit), engine, IFL (Integrated Facility for Linux), IOP (I/O processor), SAP (system assist processor), CP (central processor), PU (processing unit), zAAP (zSeries Application Assist Processor)
Computer	CEC (central electronics complex)

Virtual Machines

What: Virtual Machines

IBM @server zSeries



A **virtual machine** is an execution context that obeys the architecture.

The purpose of z/VM is to **virtualize** the real hardware:

- Faithfully replicate the z/Architecture Principles of Operation
- Permit any virtual configuration that could legitimately exist in real hardware
- Let many virtual machines operate simultaneously
- Allow overcommitment of the real hardware (processors, for example)
- Designed for many thousands of virtual machines per z/VM image (I have seen 40,000)
- Your limits will depend on the size of your physical zSeries computer

What: A Virtual Machine

IBM @server zSeries

Virtual
machine

z/Architecture

512 MB of memory

2 processors

Basic I/O devices:

A console

A card reader

A card punch

A printer

Some read-only disks

Some read-write disks

Some networking devices

We permit any configuration that a real zSeries machine could have.

In other words, we completely implement the z/Architecture Principles of Operation.

There is no "standard virtual machine configuration".

How: VM User Directory

IBM @server zSeries

Definitions of:	USER LINUX01 MYPASS 512M 1024M G
	MACHINE ESA 2
- memory	IPL 190 PARM AUTOOCR
- architecture	CONSOLE 01F 3270 A
	SPOOL 00C 2540 READER *
- processors	SPOOL 00D 2540 PUNCH A
	SPOOL 00E 1403 A
- spool devices	SPECIAL 500 QDIO 3 SYSTEM MYLAN
- network device	LINK MAINT 190 190 RR
	LINK MAINT 19D 19D RR
- disk devices	LINK MAINT 19E 19E RR
	MDISK 191 3390 012 001 ONEBIT MW
- other attributes	MDISK 200 3390 050 100 TWOBIT MR

How: CP Commands

IBM  zSeries

CP DEFINE

- Adds to the virtual configuration somehow
- CP DEFINE STORAGE
- CP DEFINE PROC
- CP DEFINE *{device} {device_specific_attributes}*

CP ATTACH

- Gives an entire real device to a virtual machine

CP DETACH

- Removes a device from the virtual configuration

CP LINK

- Lets one machine's disk device also belong to another's configuration

Changing the virtual configuration after logon is considered normal.
Usually the guest operating system detects and responds to the change.

Processors

What: Processors

IBM server zSeries

Configuration

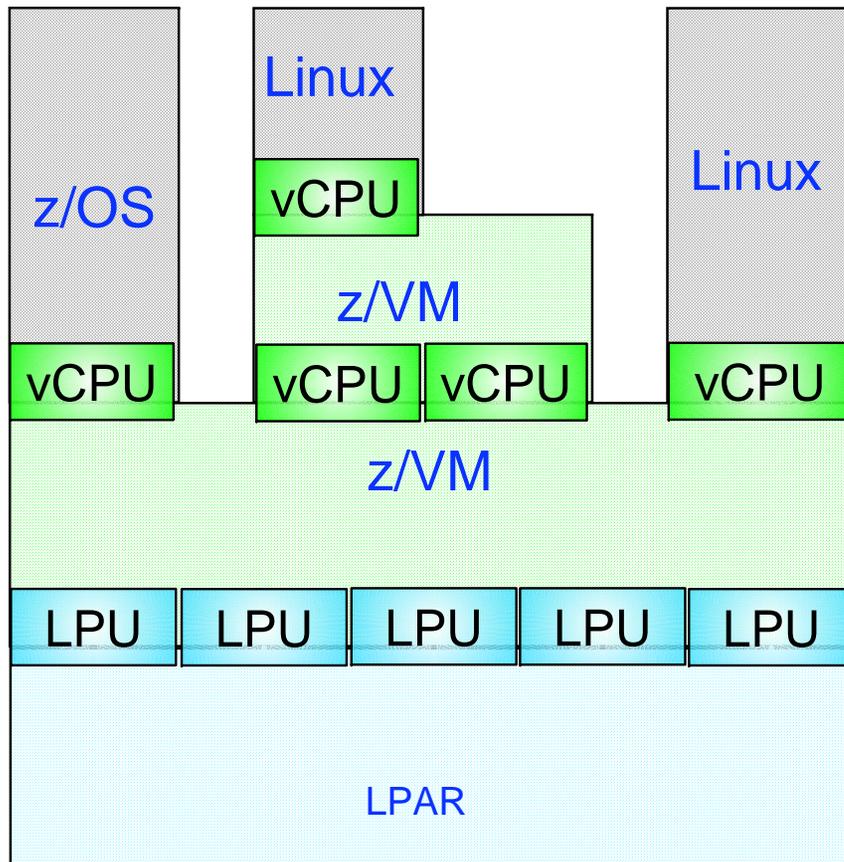
- Virtual 1- to 64-way
 - Defined in user directory, or
 - Defined by CP command
- A real processor can be dedicated to a virtual machine

Control and Limits

- Scheduler selects virtual processors according to apparent CPU need
- "Share" setting - prioritizes real CPU consumption
 - Absolute or relative
 - Target minimum and maximum values
 - Maximum values (limit shares) either hard or soft
- "Share" for virtual machine is divided among its virtual processors

What: Logical and Virtual Processors

IBM @server zSeries



How: Start Interpretive Execution (SIE)

IBM @server zSeries

- SIE = "Start Interpretive Execution", an instruction
- z/VM (like the LPAR hypervisor) uses the SIE instruction to "run" virtual processors for a given virtual machine.
- Our processor chips contain special hardware (registers, etc.) to make SIE fast
- SIE has access to:
 - A control block that describes the virtual processor state (registers, etc.)
 - The Dynamic Address Translation (DAT) tables for the virtual machine
- z/VM gets control back from SIE for various reasons:
 - Page faults
 - I/O channel program translation
 - Privileged instructions (including CP system service calls)
 - CPU timer expiration (dispatch slice)
 - Other, including CP asking to get control for special cases
- CP can also shoulder tap SIE from another processor to remove virtual processor from SIE (perhaps to reflect an interrupt)

How: Scheduling and Dispatching

IBM @server zSeries

VM

- *Scheduler* determines priorities based on *share* setting and other factors
- *Dispatcher* runs a virtual processor on a real processor
- Virtual processor runs for (up to) a *minor time slice*
- Virtual processor keeps competing for (up to) an *elapsed time slice*

LPAR hypervisor

- Uses *weight* settings for partitions, similar to share settings for virtual machines
- Dispatches logical processors on real engines

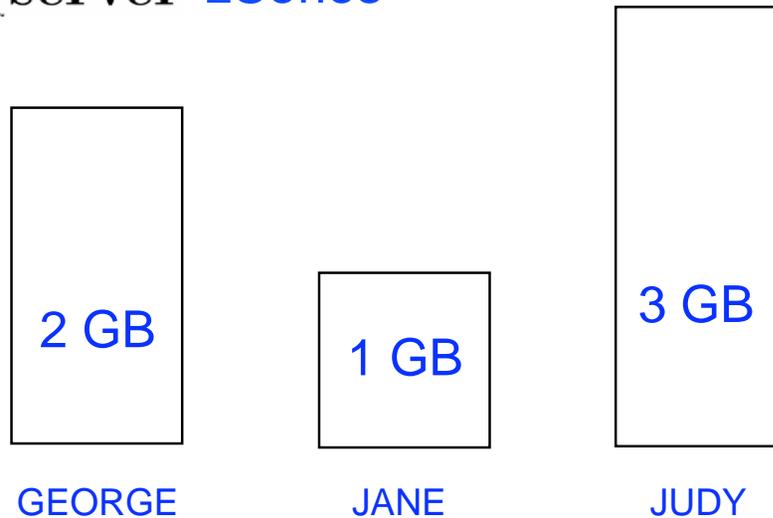
Linux

- *Scheduler* handles prioritization and dispatching processes for a time slice or *quantum*

Memory

What: Virtual Memory

IBM @server zSeries



Configuration

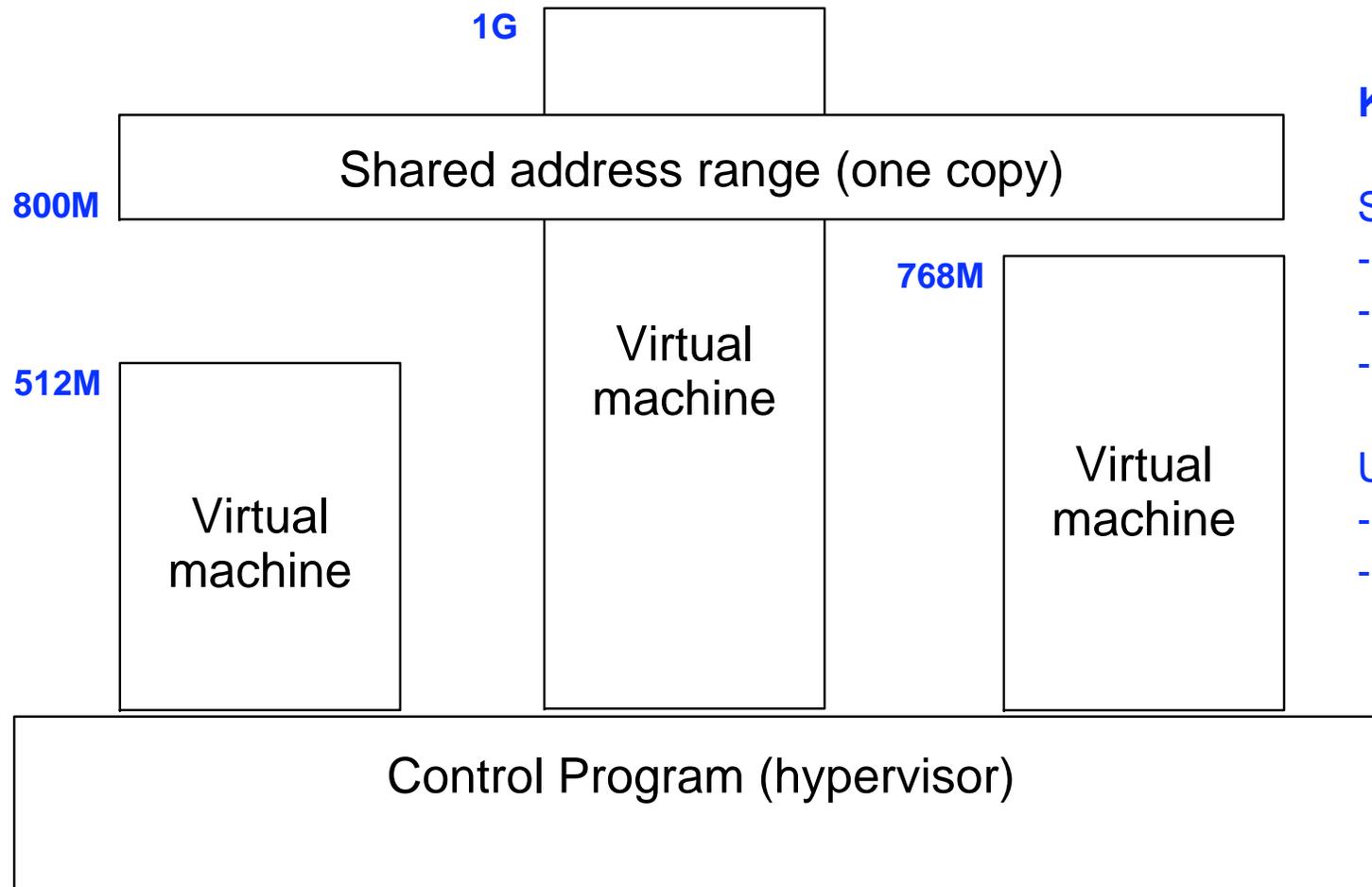
- Defined in CP directory entry or via CP command
- Can define storage with gaps (useful for testing)
- Can attach expanded storage to virtual machine

Control and Limits

- Scheduler selects virtual machines according to apparent need for storage and paging capacity
- Virtual machines that do not fit criteria are placed in the *eligible list*
- Can reserve an amount of real storage for a guest's pages
- Can lock certain specific guest pages into real storage

What: Shared Memory

IBM @server zSeries



Key Points:

Sharing:

- Read-only
- Read-write
- Security knobs

Uses:

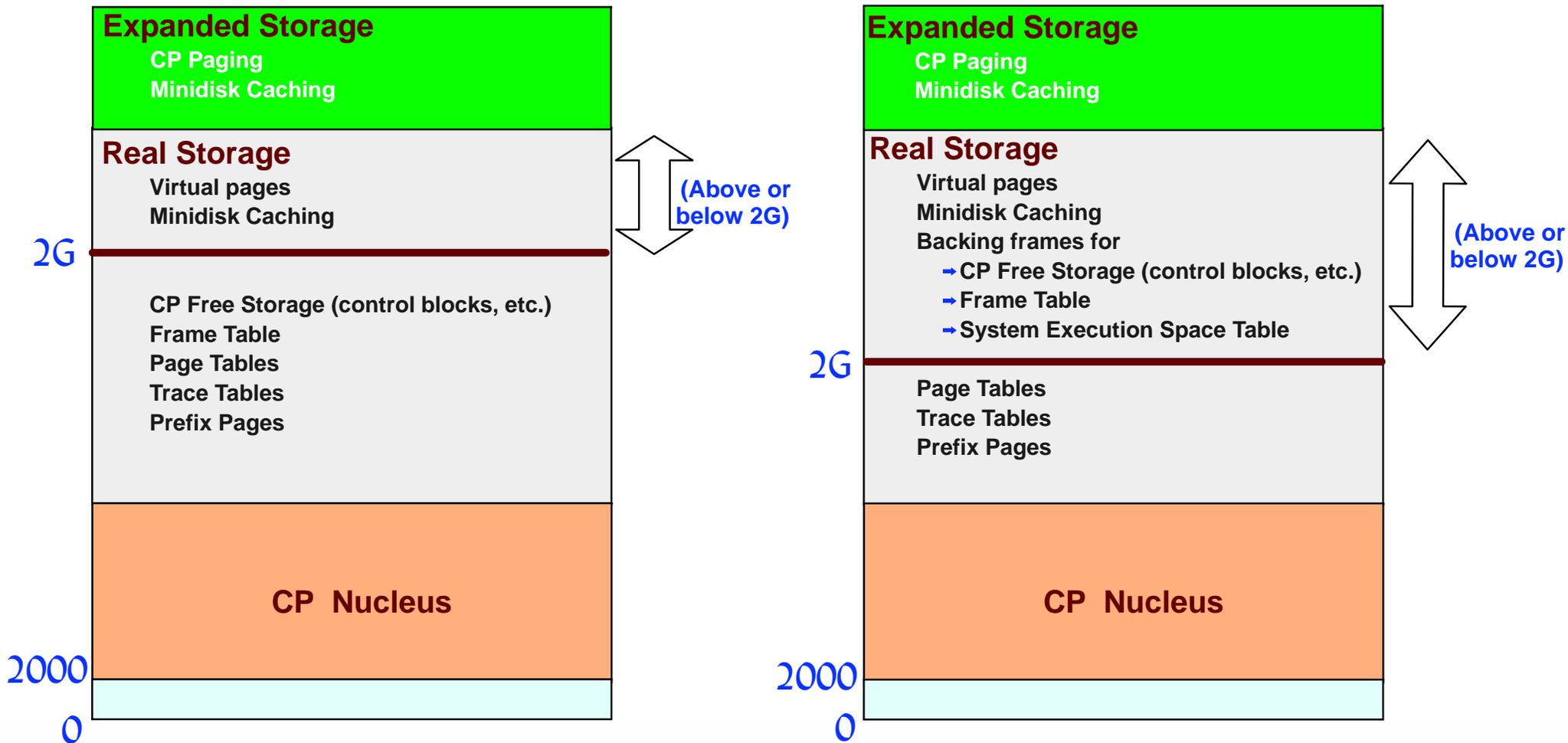
- Common kernel
- Shared programs

More: Layout of Real Storage

IBM @server zSeries

z/VM 5.1.0

z/VM 5.2.0



IBM @server. For the next generation of e-business.

How: Memory Management

IBM @server zSeries

VM

- Demand paging between central and expanded
- Block paging with DASD (disk)
- Steal from central based on LRU with reference bits
- Steal from expanded based on LRU with timestamps
- Paging activity is traditionally considered normal

LPAR

- Dedicated storage, no paging

Linux

- Paging on per-page basis to swap disks
- No longer swaps entire processes
- Traditionally considered bad

I/O Resources

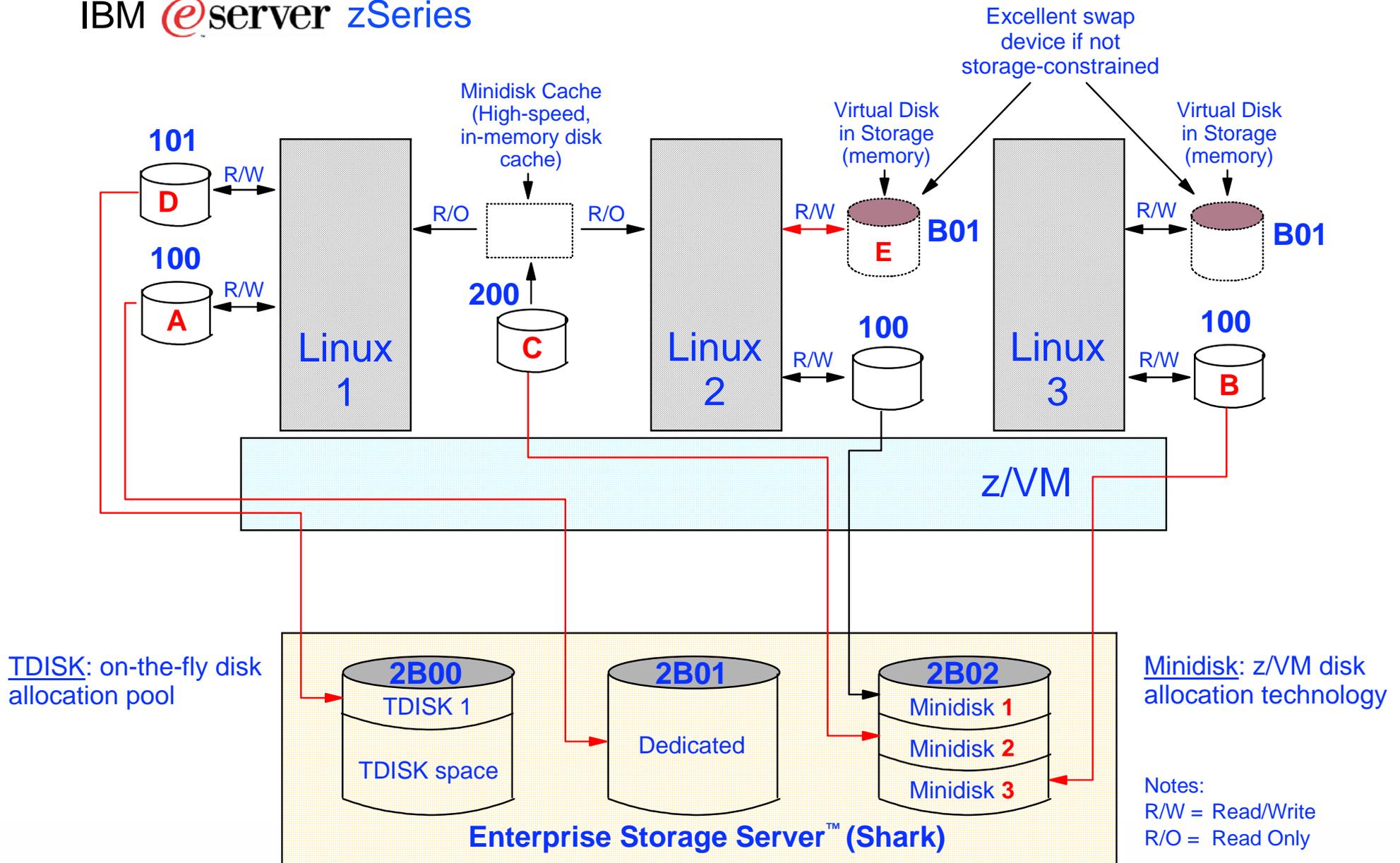
What: Device Management Concepts

IBM @server zSeries

- *Dedicated* or *attached*
 - The guest has exclusive use of the entire real device.
- *Virtualized*
 - Present a slice of a real device to multiple virtual machines
 - Slice in time or slice in space
 - E.g., DASD, crypto devices
- *Simulated*
 - Provide a device to a virtual machine without the help of real hardware
 - Virtual CTCAs, virtual disks, guest LANs, spool devices
- *Control and Limits*
 - Indirect control through "share" setting
 - Real devices can be "throttled" at device level
 - Channel priority can be set for virtual machine
 - MDC fair share limits (can be overridden)

What: Virtualization of Disks

IBM @server zSeries



IBM @server. For the next generation of e-business.

What: Data-in-Memory

IBM @server zSeries

Minidisk Cache

- Write-through cache for non-dedicated disks
- Cached in central or expanded storage
- Psuedo-track cache
- Great performance - exploits access registers
- Lots of tuning knobs

Virtual Disk in Storage

- Like a RAM disk that is pageable
- Volatile
- Appears like an FBA disk
- Can be shared with other virtual machines
- Plenty of knobs here too

Networking

What: Virtual Networks

IBM @server zSeries

One Linux guest (or z/VM TCP/IP stack) connects to the external network

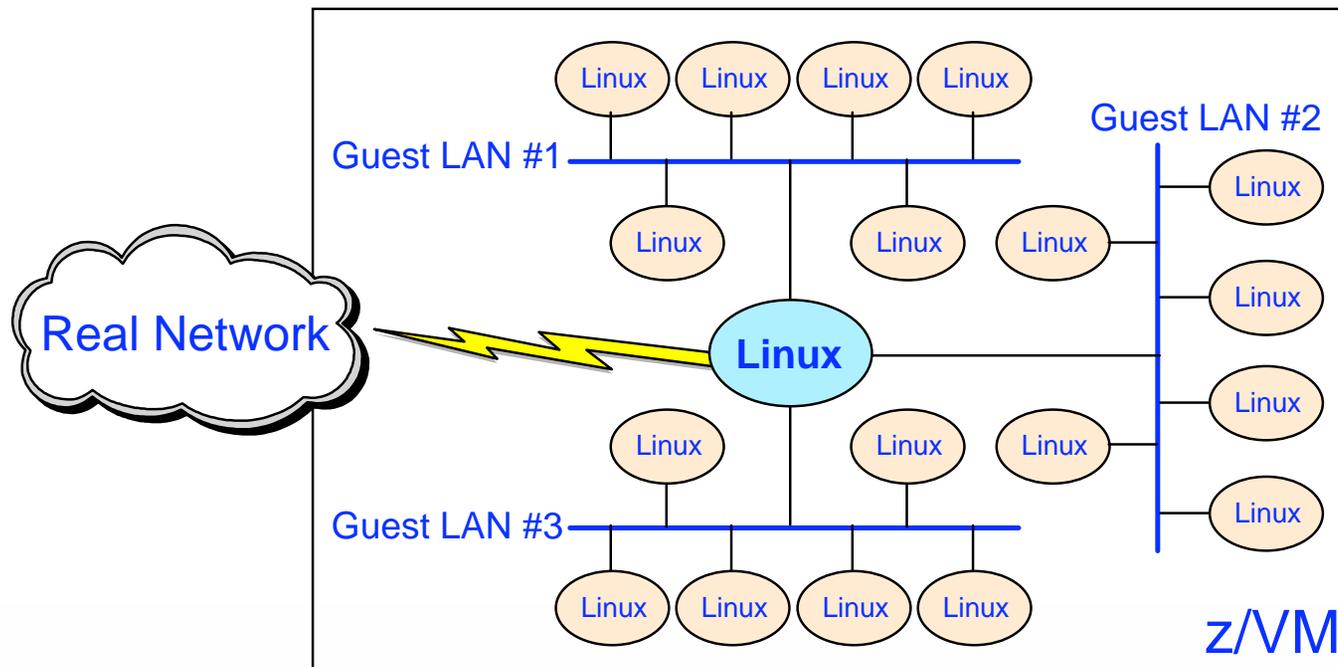
- Owns the physical OSA (to real LAN) or HiperSockets device (to another LPAR)
- Also connected to multiple guest LANs (each guest LAN is a distinct IP subnet)
- Provides routing services for guests

Another choice is the z/VM Virtual Switch

- z/VM CP itself owns the physical OSA
- Guests' virtual network adapters seem to be on the external IP subnet

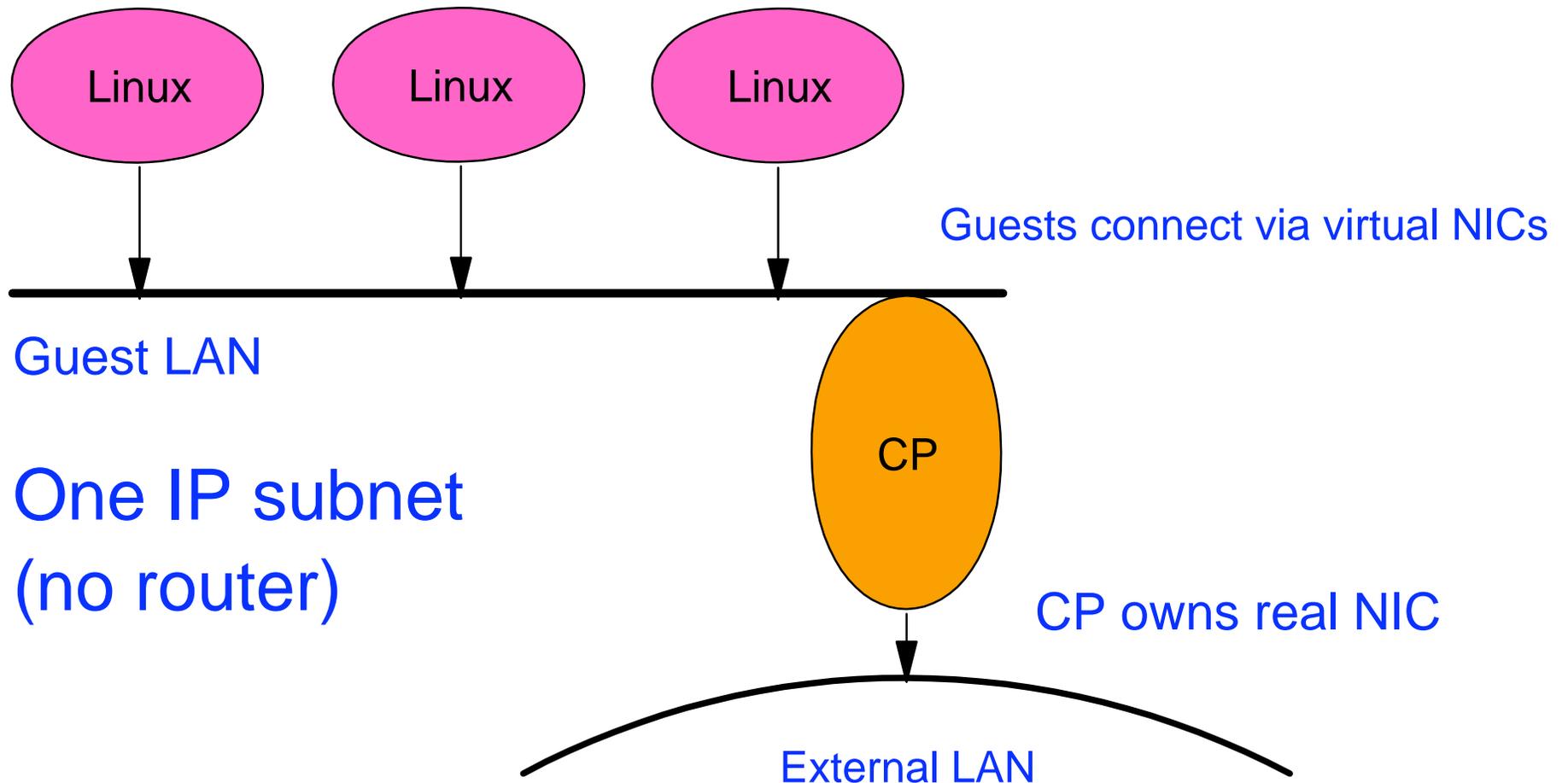
Other Linux guests connect to individual guest LAN(s)

- Virtual HiperSockets and OSA Express connections supported
- Point-to-point, Multicast, and Broadcast (QDIO) supported



What: Virtual Switch

IBM @server zSeries



Beyond Virtualization

What: Other Control Program (CP) Interfaces

IBM @server zSeries

Commands

- Query or change virtual machine configuration
- Debug and tracing
- Commands fall into different privilege classes
- Some commands affect entire system

Inter-virtual-machine communication

- Connectionless or connection-oriented protocols
- Most pre-date TCP/IP

System Services

- Enduring connection to hypervisor via a connection-oriented program-to-program API
- Various services: Monitor (performance data), Accounting, Security

Diagnose Instructions

- These are really programming APIs (semantically, procedure calls)
- Operands communicate with hardware (or in this case the virtual hardware) in various ways

What: Debugging a Virtual Machine

IBM @server zSeries

Tracing of virtual machine

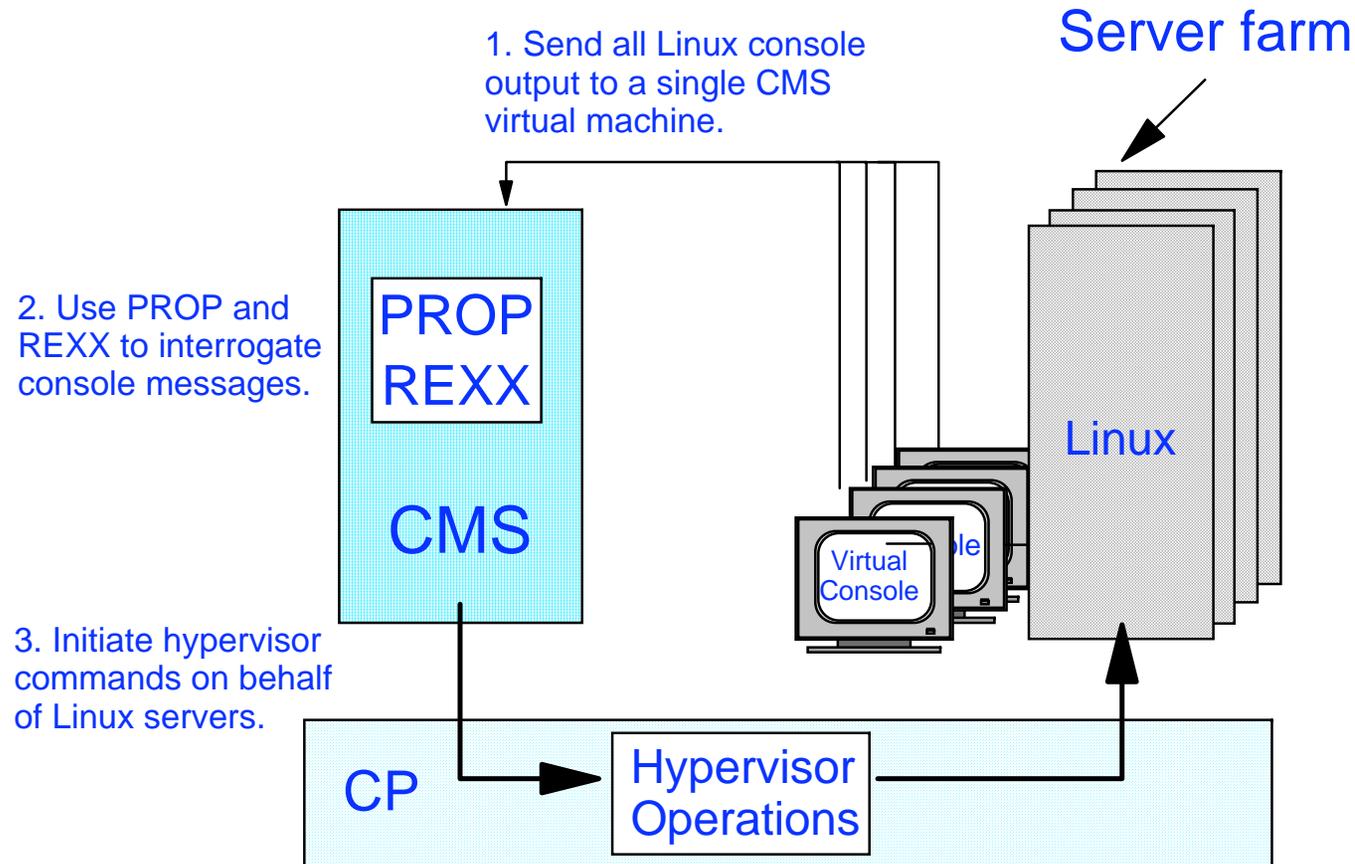
- CP TRACE command has >40 pages of documentation on tracing of:
 - instructions
 - storage references
 - some specific opcodes or privileged instructions
 - branches
 - various address space usage
 - registers
 - etc
- Step through execution or run and collect information to spool
- Trace points can trigger other commands

Display or store into virtual memory

- Helpful, especially when used with tracing
- Valid for various virtual address spaces
- Options for translation as EBCDIC, ASCII, or 390 opcode
- Locate strings in storage
- Store into virtual memory (code, data, etc.)

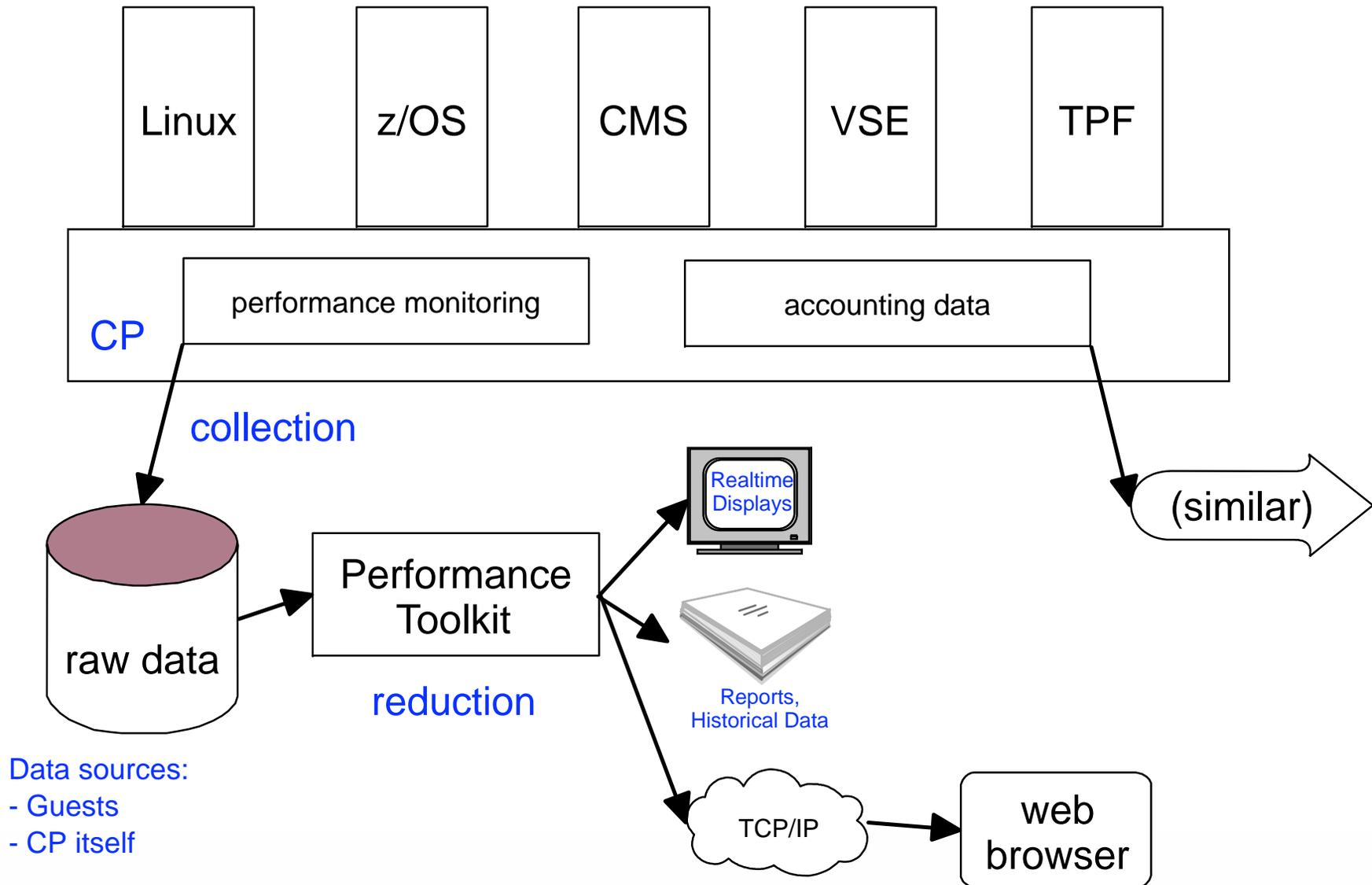
What: Programmable Operator

IBM @server zSeries



What: Performance and Accounting Data

IBM @server zSeries



IBM @server. For the next generation of e-business.

References

IBM @server zSeries

VM web site: www.vm.ibm.com

Publications on VM Web Site

- <http://www.vm.ibm.com/pubs/>
- Follow the links to the latest z/VM library
- Of particular interest:
 - z/VM CP Command and Utility Reference
 - z/VM CP Planning and Administration
 - z/VM CP Programming Services
 - z/VM Performance

IBM Systems Journal Vol. 30, No. 1, 1991

- Good article on SIE

End of Presentation

Question and Answer Time