



# IS02 – Virtualisation Options with DPM, z/VM and KVM on IBM z Systems

Arwed Tschoeke  
Client Center Böblingen



# Agenda

- Virtualization basics
- PR/SM and DPM
- z/VM
- KVM
- Docker

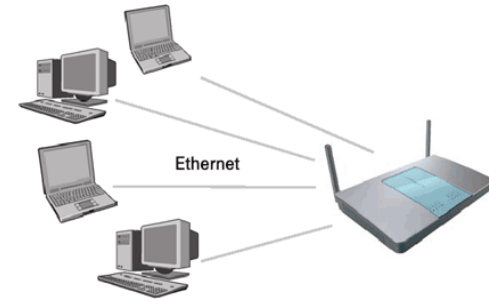


# Virtualization in **your** daily life

- Ideas!?



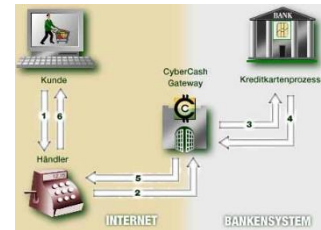
Car Sharing



Internet Access Sharing



Bathroom Sharing

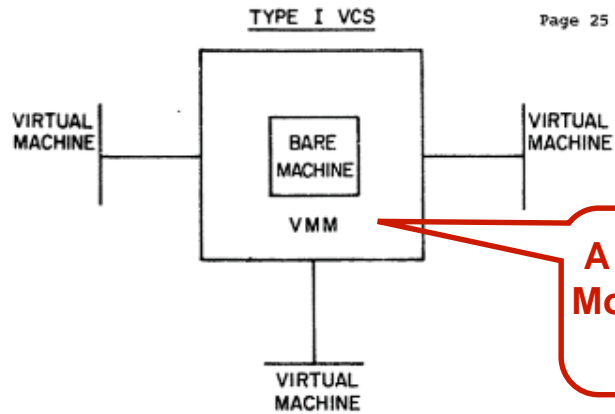


Money and virtual Money?

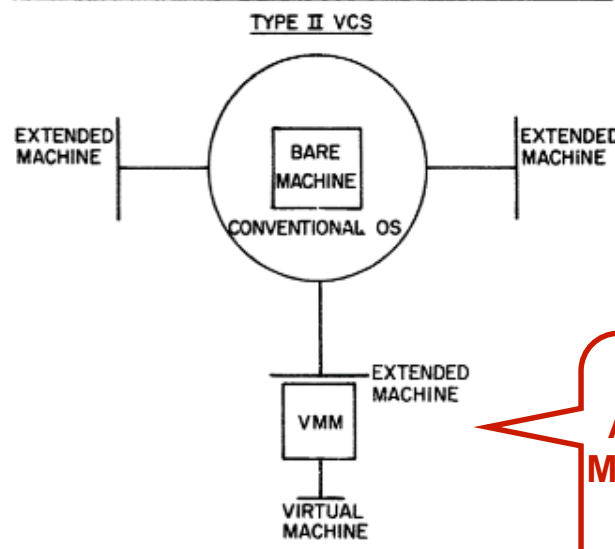
**=> Sharing „a few“ real resources with „many“ users**



# Type 1 Hypervisor



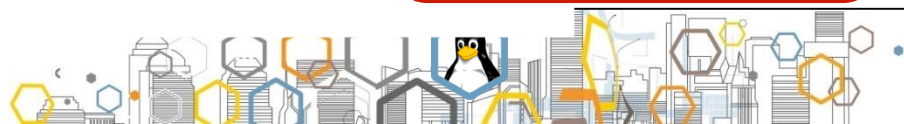
A Type 1 Virtual Machine Monitor runs as part of the Kernel



“Extended Machine” is the environment running a time-shared program. The EM evolved into a modern OS Process.

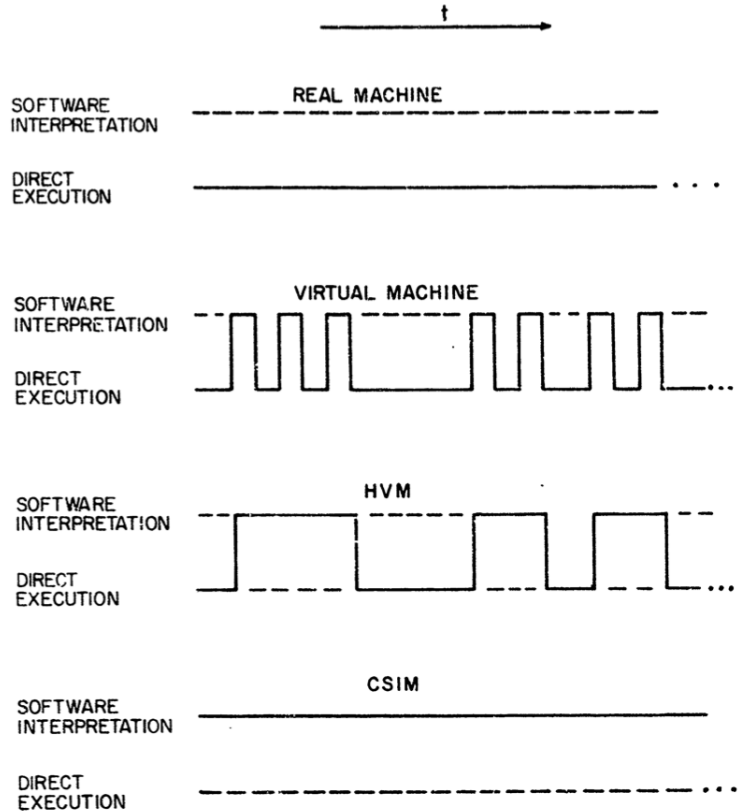
A Type II Virtual Machine Monitor runs as a standard OS Process

FIGURE 2-4 TYPE I vs TYPE II VCS



# Software Interpretation

Page 23



Non-Virtualized

KVM, VMware, Xen, PowerVM, z/VM

VMware Pre-VMX

Qemu

VIRTUAL MACHINE vs. OTHER CONSTRUCTS  
FIGURE 2-3



# Hypervisors and Virtualization for z Systems

## PR/SM-LPARs



- **Virtualization capabilities built into the system**
- PR/SM manages and virtualizes all the installed and enabled system resources as a single large SMP system
- **Full sharing of the installed resources** with high efficiency and very low overhead
- **High scalability** with support for up to 40 (for z13s) or 85 (for z13) logical partitions
- IBM Dynamic Partition Manager **simplifies management experience**
- Ensured **workload separation** based on highest EAL5+ security certification

## z/VM v6.4 (preview)



- **Enables extreme scalability, security and efficiency** creating cost savings opportunities
- **Ease Migration** with upgrade in place infrastructure provides a seamless migration path from previous z/VM releases (z/VM 6.2 and z/VM 6.3) to the latest version
- **Operational improvements** by enhancing z/VM to provide ease of use
- **Improved SCSI support** for guest attachment of disk and other peripherals, and hypervisor attachment of disk drives
- IBM Wave for z/VM **simplifies the management** of virtual Linux servers from a single user interface
- Provides the foundation for cognitive computing on z Systems

## IBM Wave for z/VM



## KVM on z Systems v1.1.1



- **Support new analytics workloads** with Single Instruction Multiple Data (SIMD) for competitive advantage
- **Deliver higher compute capacity** with support for Simultaneous Multithreading (SMT) to meet new business requirements
- **RAS support enhanced for problem determination and high availability** setup to reduce down time and quickly react to business needs
- **Secure and protect** business data with Crypto exploitation



# PR/SM and DPM



# PR/SM or LPAR Hypervisor

- 'Processor Resource/System Manager' (PR/SM) and 'LPAR hypervisor' are commonly used synonymously.
- However the 'LPAR Hypervisor' is the program itself and 'PR/SM' is the facility of the whole
- So PR/SM aka LPAR hypervisor is a Type-1 Hypervisor that manages logical partitions:
  - Each partition owns a defined amount of physical storage
  - Strictly no storage shared across partitions
  - No virtual storage management / paging done by LPAR hypervisor
  - Zone relocation lets each partition start at address 0
  - CPUs may be dedicated to a partition or may be shared by multiple partitions
  - I/O channels may be dedicated to a partition or may be shared by multiple partitions (Multiple image facility, MIF)
  - Each LPAR has its own architecture mode (ESA/390 or z/Architecture)
- PR/SM is shipped with z Systems (considered as part of the firmware)
- PR/SM was initially introduced in 1988 with the IBM 3090 processors
- Beginning with z990, the PR/SM is always loaded (no Basic Mode anymore)
- Separation of logical partitions is considered as good as having each partition on a separate physical machine (Evaluation Assurance Level 5)





# Dynamic Partition Manager (DPM) – At a Glance



**Simplification**  
Provide simplified, consumable, enhanced Partition life-cycle and integrated dynamic I/O management capabilities.



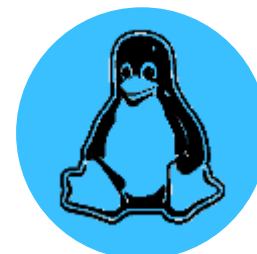
**DPM Mode**  
A CPC can be in non-DPM mode or DPM mode. Enable DPM mode with first IML.



**Cloud**  
Provides the technology foundation that enables IaaS and secure, private Clouds.



**FIE**  
Initial focus on First In Enterprise (FIE) customers with support for existing clients in a later stage.



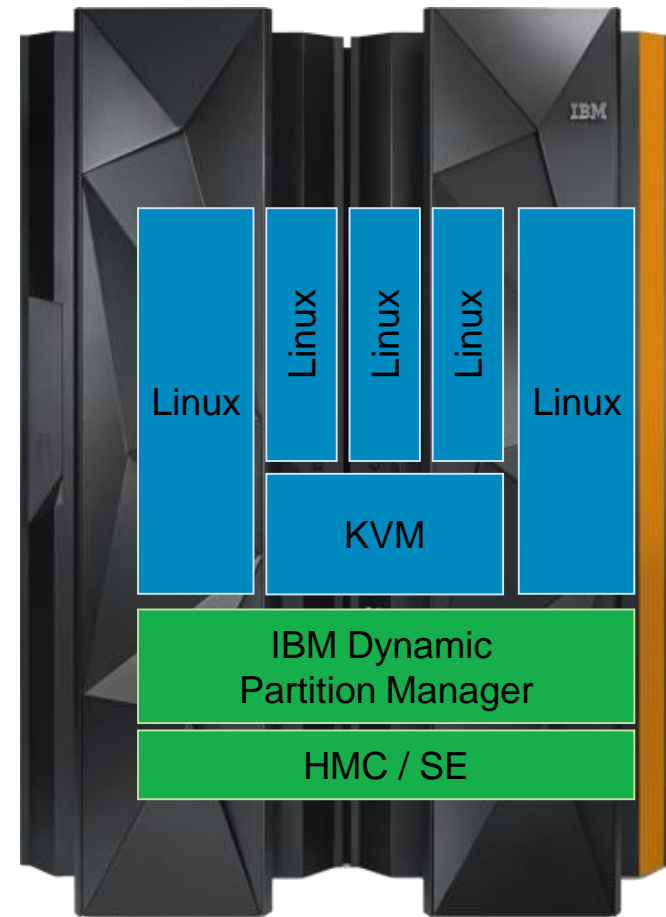
**Linux only**  
A CPC running in DPM mode is Linux only. No z/OS, z/VM, zVSE, zTPF support in Stage 1. FCP Storage only.

*“DPM provides simplified z Systems hardware and virtual infrastructure management including integrated dynamic I/O management for FIE customers that run KVM on z as a hypervisor and/or Linux on z as a Partition-hosted operating system.”*



# Partition and I/O device management at the HMC

- Linux partitions or KVM partitions correspond to LPARs under standard PR/SM
- Supports only Linux and Linux based hypervisors
- Creation of I/O Configuration Data Set (IOCDs) is “under the covers”
  - Supports dynamic updates of I/O
- Hardware and operating system message displays are unchanged
- Problem determination and maintenance continues to exist on the System Element (SE)
- On/Off Capacity on Demand (OoCoD) and Customer Initiated Upgrade (CIU) supported for Linux



# How DPM helps in a new Linux Environment

- z Systems and PR/SM require a HW definition
- Dynamic IO – one of the key differentiator of the platform would be nice
- Having the option to have a GUI-based administration
- Overcome the prejudice: z is old school and complicated
- Everything is scripted (that's why we need GUIs ;-)

```

ICP ICP070I SEARCH FOR '*ICP' TO FIND EACH IOCP MESSAGE
ID MSG1='IODF00',MSG2='SYS1.IODF00 - 2014-07-22 16:19', *
SYSTEM=(2828,1),LSYSTEM=P008B857, *
RESOURCE PARTITION=((CSS(0),(ZOS1,2),(ZVM1,1),(*,3),(*,4),(*,5*
),(*,6),(*,7),(*,8),(*,9),(*,A),(*,B),(*,C),(*,D),(*,E),*
(*,F)),(CSS(1),(*,1),(*,2),(*,3),(*,4),(*,5),(*,6),(*,7)*
,(*,8),(*,9),(*,A),(*,B),(*,C),(*,D),(*,E),(*,F)))
CHPID PATH=(CSS(0),25),SHARED,PARTITION=((ZOS1,ZVM1),(=)), *
PCHID=160,TYPE=FC
CHPID PATH=(CSS(0),26),SHARED,PARTITION=((ZOS1,ZVM1),(=)), *
PCHID=11C,TYPE=FC
CHPID PATH=(CSS(0),27),SHARED,PARTITION=((ZOS1,ZVM1),(=)), *
PCHID=161,TYPE=FC
CHPID PATH=(CSS(0),28),SHARED,PARTITION=((ZOS1,ZVM1),(=)), *
PCHID=11D,TYPE=FC
CHPID PATH=(CSS(0),F8),SHARED,PARTITION=((ZOS1,ZVM1),(=)), *
PCHID=17C,TYPE=OSD
CHPID PATH=(CSS(0),F9),SHARED,PARTITION=((ZOS1,ZVM1),(=)), *
PCHID=104,TYPE=OSD
CNTLUNIT CUNUMBR=0030,PATH=((CSS(0),F8)),UNIT=OSA
IODEVICE ADDRESS=(030,064),UNITADD=00,CUNUMBR=(0030),UNIT=OSA
CNTLUNIT CUNUMBR=0070,PATH=((CSS(0),F9)),UNIT=OSA
IODEVICE ADDRESS=(070,064),UNITADD=00,CUNUMBR=(0070),UNIT=OSA
CNTLUNIT CUNUMBR=0300,PATH=((CSS(0),25,26)), *
UNITADD=((00,032)),CUADD=8,UNIT=2105
IODEVICE ADDRESS=(300,032),CUNUMBR=(0300),STADET=Y,UNIT=3390
CNTLUNIT CUNUMBR=0320,PATH=((CSS(0),25,26)), *
UNITADD=((00,032)),CUADD=2,UNIT=2105
IODEVICE ADDRESS=(320,032),UNITADD=00,CUNUMBR=(0320),STADET=Y, *
UNIT=3390
CNTLUNIT CUNUMBR=0340,PATH=((CSS(0),25,26)), *
UNITADD=((00,032)),CUADD=4,UNIT=2105
IODEVICE ADDRESS=(340,032),UNITADD=00,CUNUMBR=(0340),STADET=Y, *
UNIT=3390
CNTLUNIT CUNUMBR=0360,PATH=((CSS(0),25,26)), *
UNITADD=((00,032)),CUADD=6,UNIT=2105
IODEVICE ADDRESS=(360,032),UNITADD=00,CUNUMBR=(0360),STADET=Y, *
UNIT=3390
CNTLUNIT CUNUMBR=0380,PATH=((CSS(0),25,26)), *
UNITADD=((00,032)),CUADD=8,UNIT=2105
IODEVICE ADDRESS=(380,032),UNITADD=00,CUNUMBR=(0380),STADET=Y, *
UNIT=3390
CNTLUNIT CUNUMBR=0400,PATH=((CSS(0),25,26)), *
UNITADD=((00,032)),CUADD=A,UNIT=2105
IODEVICE ADDRESS=(400,032),CUNUMBR=(0400),STADET=Y,UNIT=3390
  
```

short version: No Texteditor  
required to get started, dynamic IO available



# How DPM looks like

New Partition - new\_server\_1

General  
 Status  
 Processors  
 Memory  
 Network  
 Storage  
 Accelerators  
 Cryptos  
 Boot  
 Controls

General
--- General

Name:

Description:

Short name:

Partition ID:   Generate automatically

Reserve resources to ensure they are available when the partition is started ?

Status
--- Status

Acceptable statuses: ?

Active  
 Starting  
 Terminated  
 Status check

Stopped  
 Stopping  
 Paused  
 Communications not active

Degraded  
 Reservation error

Processors
--- Processors

Processor type:  Central Processor (CP)  Integrated Facility for Linux (IFL)

Processor mode:  Shared  Dedicated

Processors:  ?

Processing weight:  ?

Enforce weight capping ?

Enforce absolute processor capping ?

Number of processors (0.01-255.0):

Processors

Shared Processors

Virtual:Physical-400.00% ?

Related Tasks  
[System Details](#)  
[Manage Adapters](#)  
[Monitor System](#)

Basic
OK Cancel Help



Systems Management > Z13S

Partitions Topology Monitor

Tasks Views: Puritscher

S...	Name	Status	Proc...	Memory (GB)	Processor Utilization	Network Utilization	OS Na...	OS Type	OS Level	Description
<input type="checkbox"/>	tmcc40	Active	4	20.0	1%	0%	TMCC40	z/VM	6.3.0 - 1601	z/VM 6.3 + RACF + Openstack (CMA)
<input type="checkbox"/>	waveserv	Active	1	8.0	8%	0%	ZLIN104	Linux	3.12.62	sles12sp1 with IBM Wave for z/VM
<input type="checkbox"/>	zcloudb	Active	4	30.0	100%	0%	ZVM640	z/VM	6.4.0	z/VM IESP + RACF + Openstack
<input type="checkbox"/>	zkvm105	Active	1	32.0	0%	0%	ZKVM105	KVMIBM	1.1.1	KVM tests with crypto (EP11), TSM client. RoCE ...
<input type="checkbox"/>	zkvm230	Active	2	8.0	0%	0%				test KVM automated install (master workshop system)
<input type="checkbox"/>	zkvm231	Stopped	2	8.0						ZKVM2xx created by Web API.
<input type="checkbox"/>	zkvm232	Stopped	2	8.0						ZKVM2xx created by Web API.
<input type="checkbox"/>	zkvm233	Stopped	2	8.0						ZKVM2xx created by Web API.
<input type="checkbox"/>	zkvm234	Stopped	2	8.0						ZKVM2xx created by Web API.
<input type="checkbox"/>	zkvm235	Stopped	2	8.0						ZKVM2xx created by Web API.
<input type="checkbox"/>	zkvm236	Stopped	2	8.0						ZKVM2xx created by Web API.
<input type="checkbox"/>	zkvm237	Stopped	2	8.0						ZKVM2xx created by Web API.
<input type="checkbox"/>	zkvm238	Stopped	2	8.0						ZKVM2xx created by Web API.
<input type="checkbox"/>	zkvm241	Active	2	8.0	0%	0%				ZKVM2xx created by Web API.
<input type="checkbox"/>	zkvm242	Active	2	8.0	0%	0%				ZKVM2xx created by Web API.
<input type="checkbox"/>	zkvm250	Stopped	2	8.0						ZKVM2xx created by Web API.
<input type="checkbox"/>	zlin019	Active	2	8.0	1%	0%		Linux	3.12.62	sles12sp1
<input type="checkbox"/>	zlin101	Active	1	8.0	0%	0%				ubuntu 16.10
<input type="checkbox"/>	zlin102	Active	1	8.0	0%	0%		Linux	3.10.0	rhel7u2 (demolan only)
<input type="checkbox"/>	zbxwork	Stopped	1	8.0						z install server VLAN 1747
<input type="checkbox"/>	zTSMserver	Active	1	32.0	1%	0%	ZKVM106	Linux	3.0.101	TSM server sles11sp3
<input type="checkbox"/>	zvm640	Stopped	4	20.0			ZVM640	z/VM	6.4.0	z/VM IESP install test

Max Page Size: 500 Total: 22 Filtered: 22 Selected: 0



# Partition Details - zTSMserver

- General
- Status
- Controls
- Processors
- Memory
- Network
- Storage**
- Accelerators
- Cryptos
- Boot

Storage

### HBA's

+  -    Actions

Search

<input type="checkbox"/>	Name	WWPN	Type	Adapter Name	Device Number	Card Type	Description
<input type="checkbox"/>	sana	C05076D7D2000019	FCP	SAN-A	7000	FICON Express16s	
<input type="checkbox"/>	sanb	C05076D7D200001A	FCP	SAN-B	7100	FICON Express16s	

Total: 2 Selected: 0

Accelerators

### Accelerator Virtual Functions

- Related Tasks
- [Stop](#)
- [System Details](#)
- [Manage Adapters](#)
- [Monitor System](#)



# Partition Details - zTSMserver

- General
- Status
- Controls
- Processors
- Memory
- Network**
- Storage
- Accelerators
- Cryptos
- Boot

- Network

NICs

Actions

<input type="checkbox"/>	Name	Type	Adapter Name	Adapter Port	Device Number	Card Type	Description
<input type="checkbox"/>	demolan	OSA	DemoLAN	0	EA00	OSA-Express5s 1000Base-T	
<input type="checkbox"/>	intranet	OSA	Intranet	0	EC00	OSA-Express5s 1000Base-T	

Total: 2 Selected: 0

- Storage

HBA's

- Related Tasks
- [Stop](#)
- [System Details](#)
- [Manage Adapters](#)
- [Monitor System](#)



# Partition Details - zkvm230

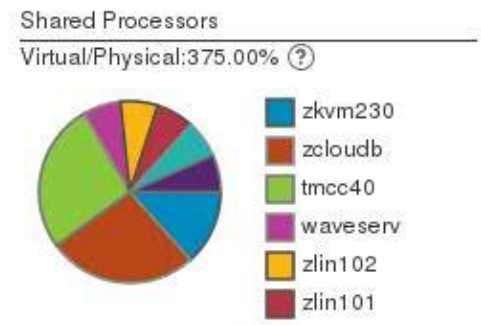
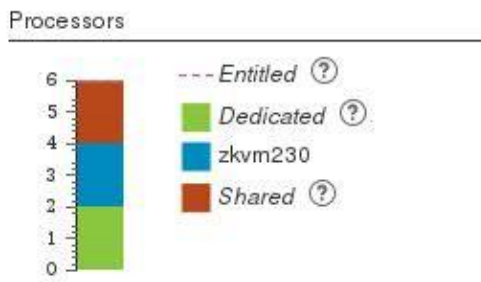
- General
- Status
- Controls
- Processors
- Memory
- Network
- Storage
- Accelerators
- Cryptos
- Boot

## Processors

Processor type:  Central Processor (CP)  Integrated Facility for Linux (IFL)  
Processor mode:  Shared  Dedicated

\* Processors:

1 2 4 5 6 2



\* Processing weight: ?

999 900 700 500 300 100 1

Very High  
High  
Medium  
Low  
Very Low

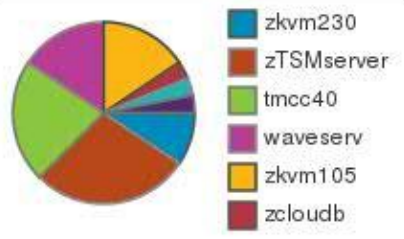
300

- Enforce weight capping ?
- Enforce absolute processor capping ?

Number of processors (0.01-255.0): 1

[Manage Processor Sharing](#)

## Active Processing Weights



- Related Tasks
- [Stop](#)
- [System Details](#)
- [Manage Adapters](#)
- [Monitor System](#)

## Memory



# Partition Details - zTSMserver

General

Status

Controls

Processors

Memory

Network

Storage

Accelerators

Cryptos

Boot

Boot

Boot from: Storage Device(SAN)

HBA



Storage Device(SAN)

Network Server(PXE)

FTP server

Hardware Management Console removable media

ISO Image

None

Name	Description
<input checked="" type="radio"/> sana	
<input type="radio"/> sanb	C05076D7D200001A 7100

Total: 2 Selected: 1

Target WWPN: 500507630300D5AA

Target LUN: 4000400B00000000

Boot program selector (0-30): 0

Boot record logical block address:

OS load parameters:

Related Tasks

[Stop](#)

[System Details](#)

[Manage Adapters](#)

[Monitor System](#)

OK Apply Cancel Help



## New Partition

General

Status

Processors

Memory

Network

Storage

Accelerators

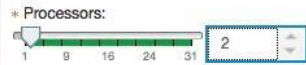
Cryptos

Boot

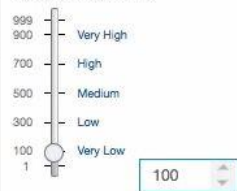
Controls

Processors

Processor type:  Central Processor (CP)  Integrated Facility for Linux (IFL)  
 Processor mode:  Shared  Dedicated



\* Processing weight: ?



Enforce weight capping ?  
 Enforce absolute processor capping ?

Number of processors (0.01-255.0):

Processors

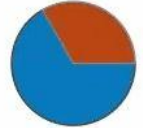


Active Processing Weights



Shared Processors

Virtual/Physical: 4.76% ?



Related Tasks  
[System Details](#)  
[Manage Adapters](#)  
[Monitor System](#)



# z/VM



# IBM z/VM Hypervisor



- z/VM is the product name of a Type-1 Hypervisor
- z/VM
  - virtualizes the architecture:
    - Guests definitions are completely virtual (and do not necessarily be consistent with physical HW)
  - support DASD and FCP
    - Offers the possibility to choose the solution with the largest convenience factor
  - SSI Clustering for increased availability
  - Integration into GDPS
- Since z990 (with the removal of the Basic Mode), z/VM always runs either in an LPAR or nested on another z/VM systems



# z/VM Version 6 Release 3 Making Room to Grow Your Business

Product General Availability

z/VM support for zEDC Express and 10GbE RoCE Express features Available, CPU Pooling

**January 14**  
z13 and z/VM Enhancements Announcement

**Feb 13**  
Base z13 & Crypto support Available

**March 13**  
SMT and Scalability Support Available

**June 26**  
Multi-VSwitch Link Aggregation Available

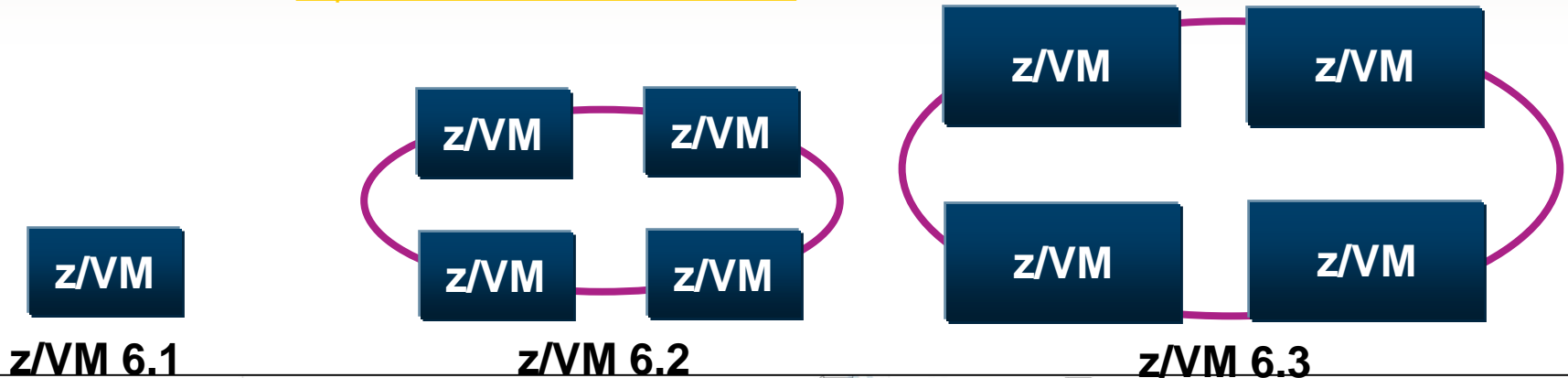
**September 15**  
RACf enhancements, Prorated Core time

**January 14**  
dynamic PDR migration, SIMD support

**Q4/14**  
z/VM6.4 GA



See <http://www.vm.ibm.com/zvm630/>



# System Programmer & Management Capability

- Upgrade In Place migration enhancements
  - Upgrade In Place migration was introduced in z/VM 6.3
  - Enhanced to allow migration to z/VM 6.4 from
    - z/VM 6.2 or z/VM 6.3 (but not both at same time in cluster)
    - Supports migration for clustered or non-clustered systems



# KVM



# KVM Overview

- KVM (Kernel Virtual Machine) is a Linux kernel-based hypervisor
- Developed and maintained by Avi Kivity / Qumranet, recently acquired by Red Hat
- KVM turns the Kernel into a hypervisor by loading a kernel module and opening a device node. The main parts of KVM are:
  - Kernel module `kvm.ko`
  - Hardware specific modules
  - Device node `/dev/kvm` (to create/run VMs from userspace with a set of `ioctl(s)`)
- Virtual machines (or guests or domains) appear as normal Linux processes and integrate seamlessly into the rest of Linux
- A VM has its own memory, that is separated from the user space process
- Virtual CPUs are not scheduled on it's own (vCPUs are realized as Linux threads, and are still scheduled by the Linux Kernel process scheduler)
- In full virtualization mode it's possible to run multiple unmodified guest OSes in parallel, with each having private virtual hardware (network, disk, graphics etc.)
- Exploits 'SIE' hardware instruction on z Systems

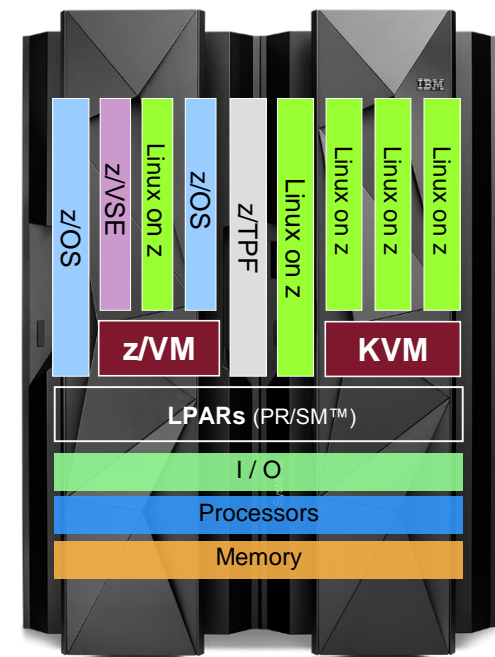




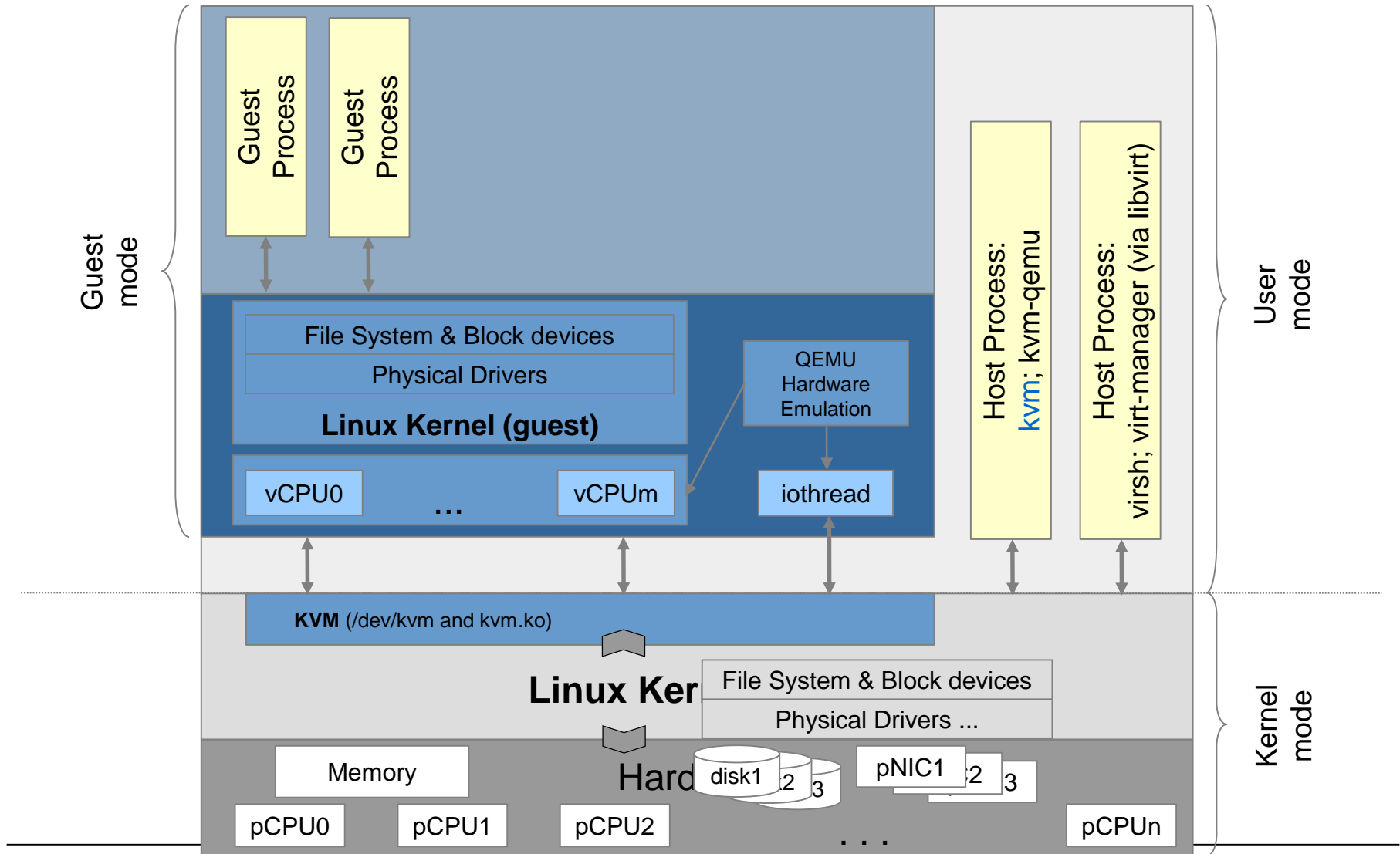
# KVM for z Systems

**In addition to z/VM, IBM supports a Kernel-based Virtual Machine (KVM) offering for z Systems that hosts Linux on z Systems guest virtual machines.**

- ◆ The KVM can be installed on z Systems processors.
- ◆ The KVM offering co-exists with z/VM virtualization environments, z/OS, Linux on z Systems, z/VSE and z/TPF.
- ◆ The KVM offering is optimized for the z Systems architecture and provides standard Linux and KVM interfaces for operational control of the environment.
  - Enterprises will be enabled to easily integrate Linux servers into their existing infrastructure and cloud offerings.



# Qemu/KVM Component Diagram



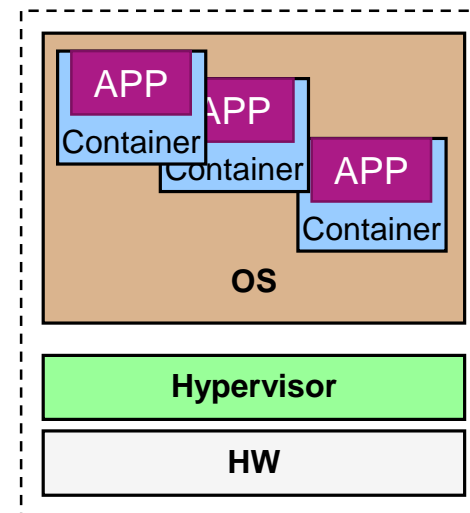
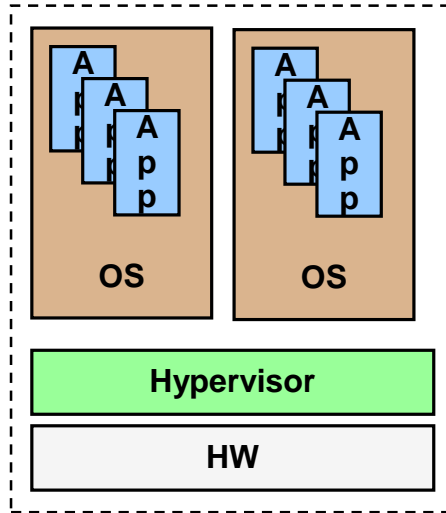
# Others



## Virtualization

vs.

## Containers



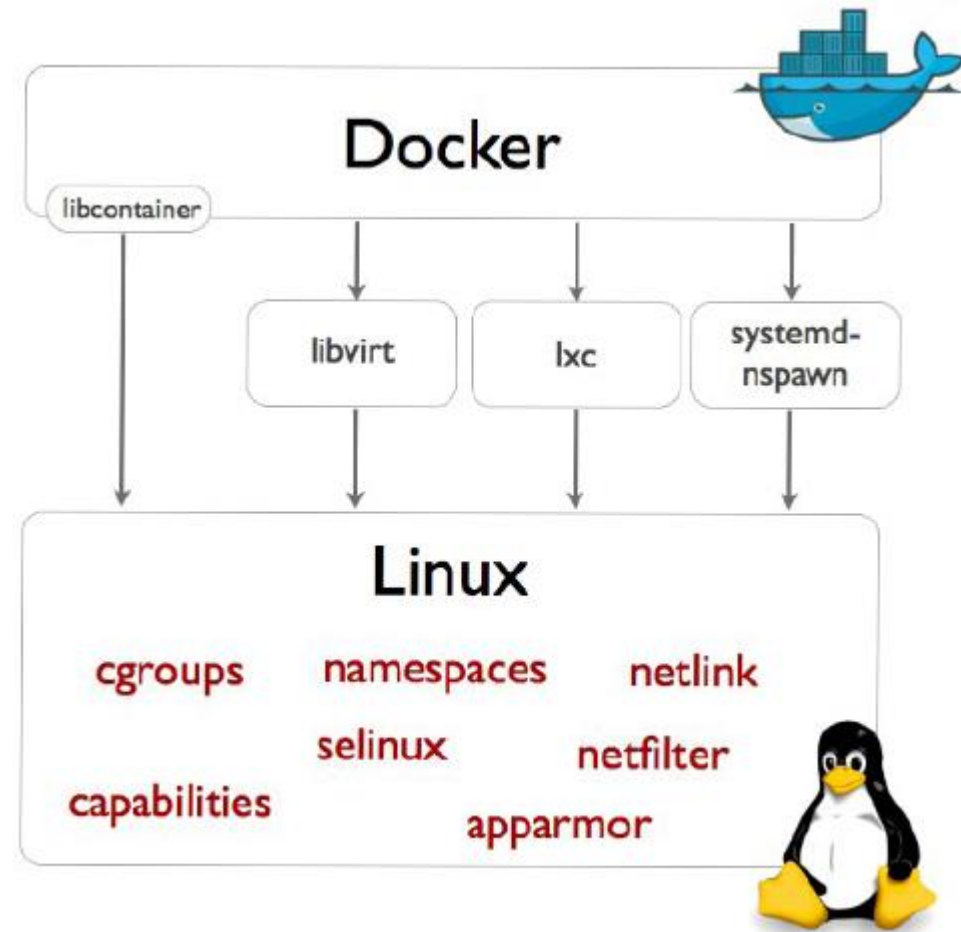
- Infrastructure oriented:
  - coming from servers, now virtualized
  - several applications per server
  - isolation
  - Separation between tenants

- Service oriented:
  - application-centric
  - solution decomposed
  - DevOps
  - separations between the apps of a tenant

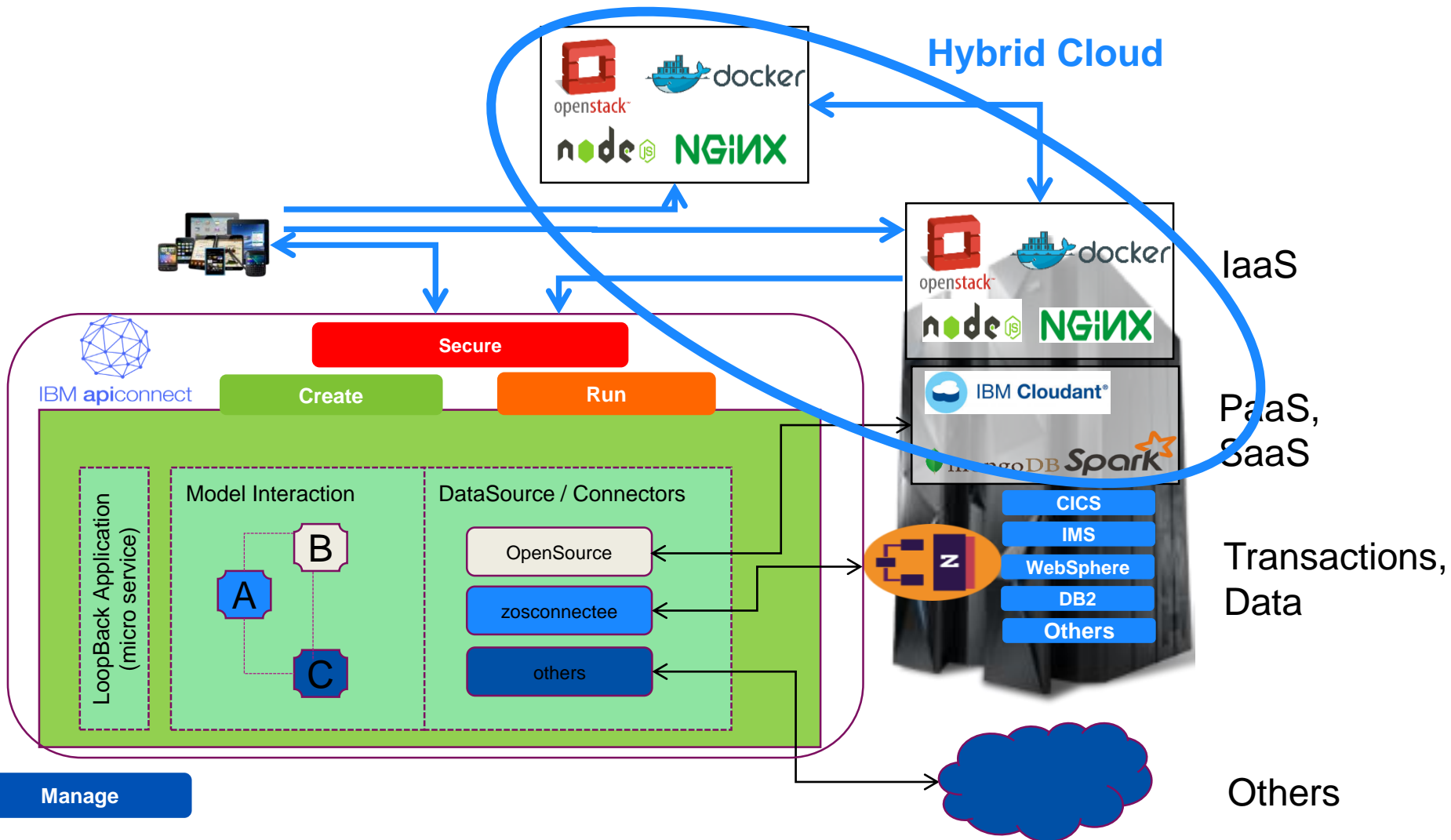


# Docker and Containers

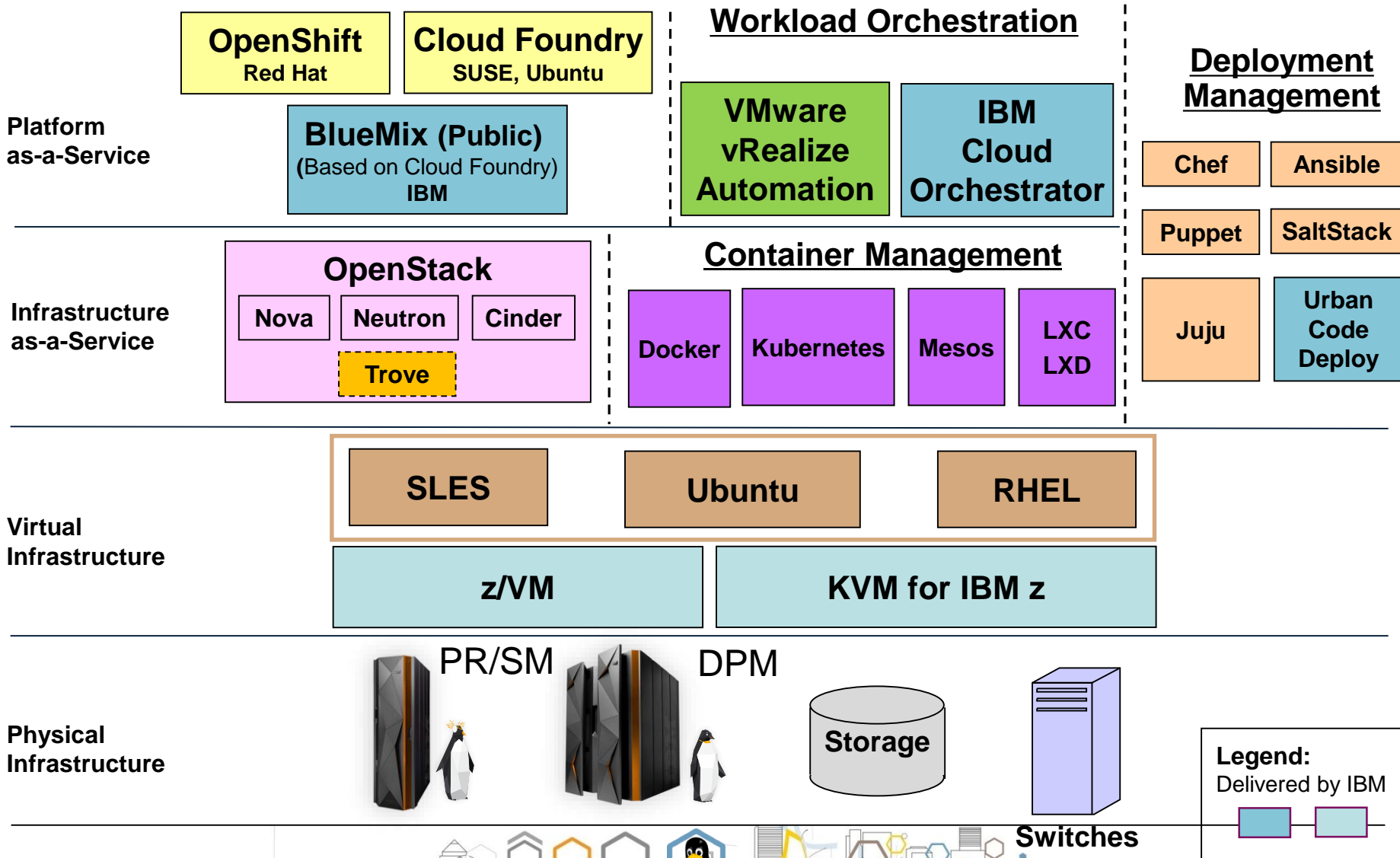
- Docker in general offers different ways to access the Linux Kernel and software resources that allows to constitute and form Containers:
  - libvirt
  - systemd-nspawn
  - lxc and
  - libcontainer
- The recently introduced 'libcontainer' library – Docker's own way to access these resources, like namespaces and cgroups – seems to prevail and to become accepted.



# z Systems: connecting data and transaction to Hybrid clouds



# No architectural limits for integration with other platforms

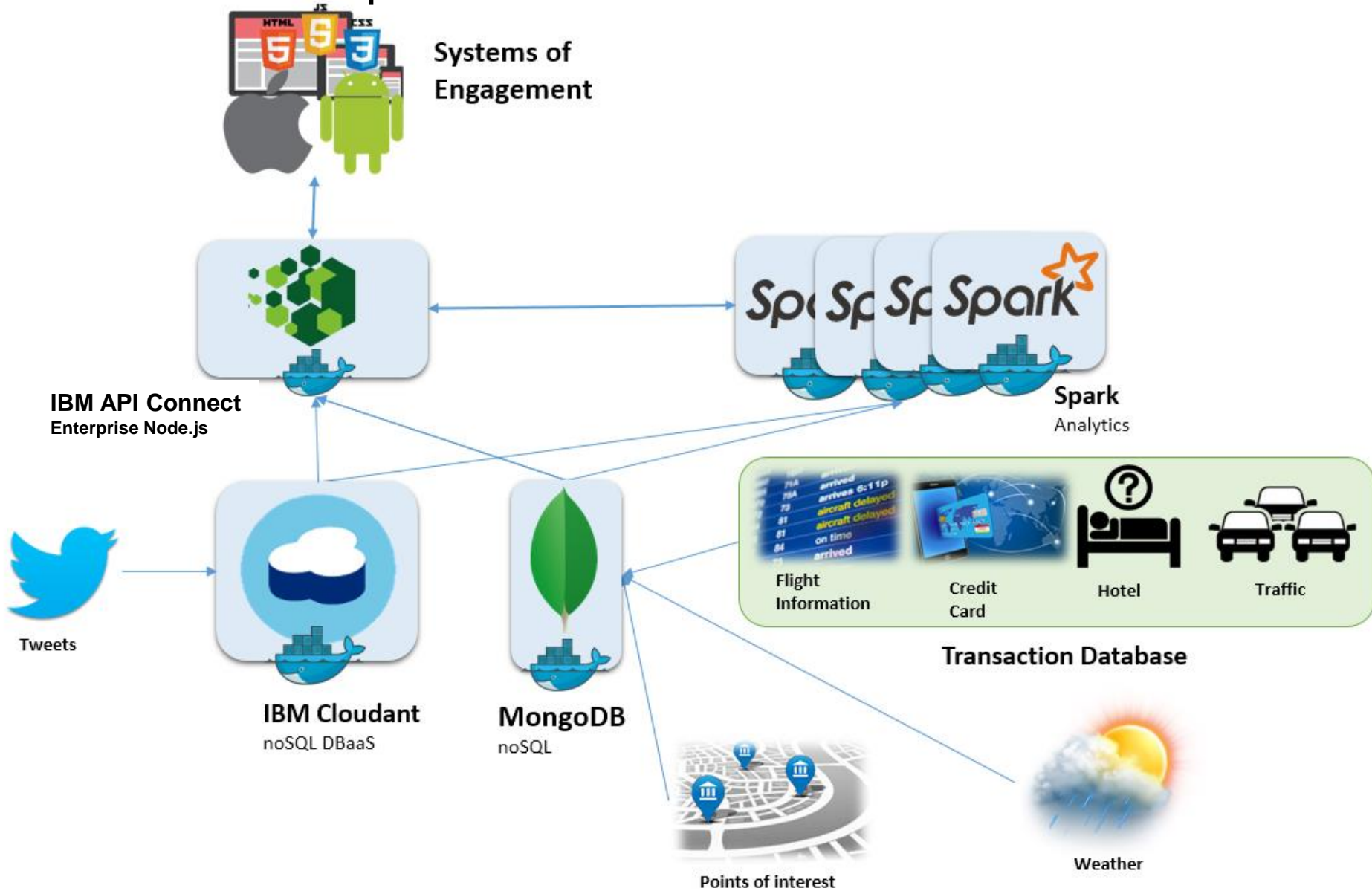


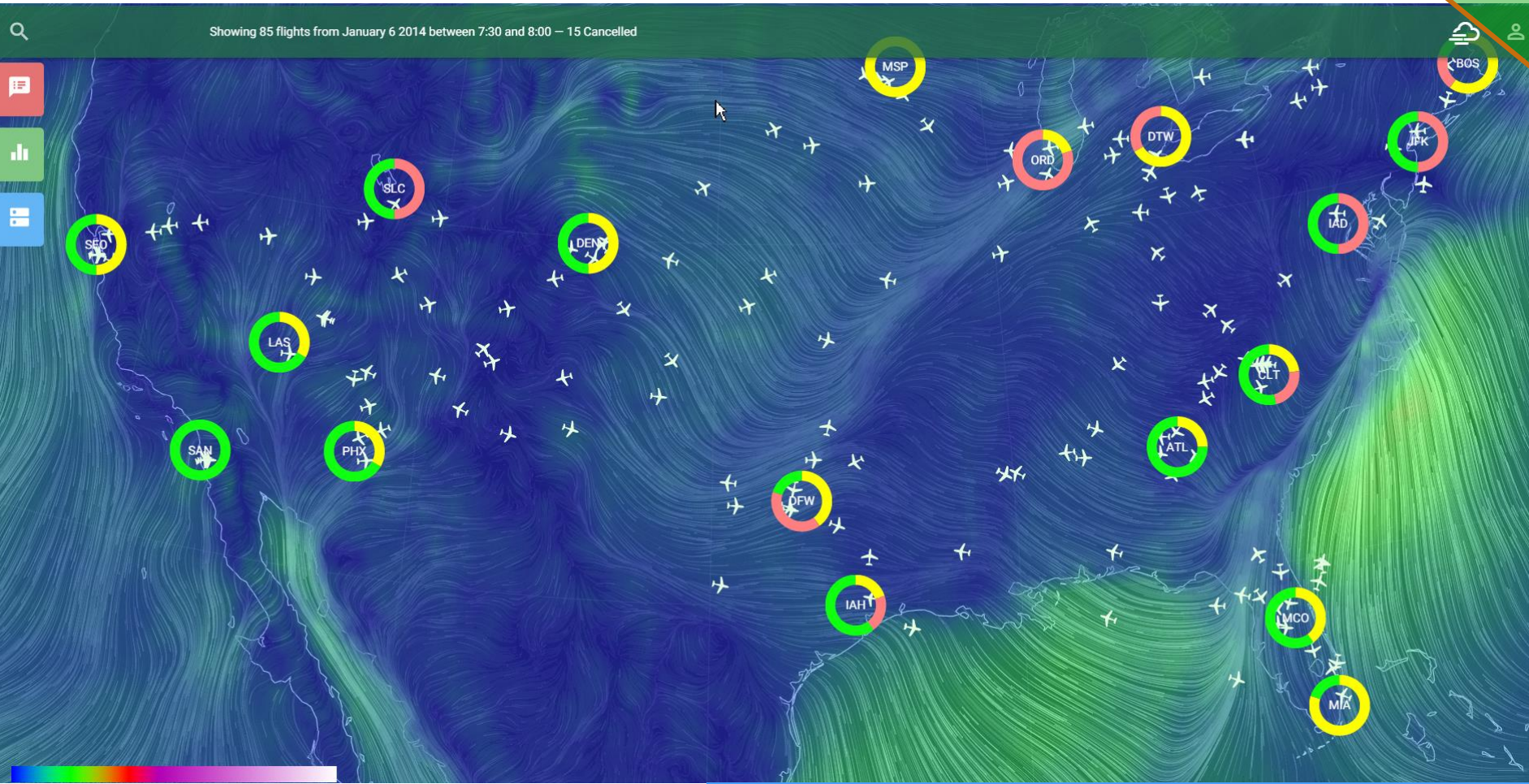
# Summary





# Architecture of a $\mu$ service Based Solution





## Why LinuxONE

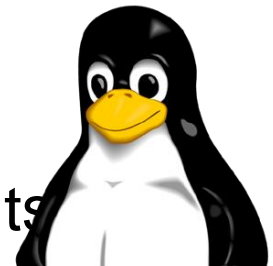
- In 2014, there were nearly 850 Million US Airline passengers or ~ 2.3 millions passengers per day.
- On a bad travel day, an average user could generate ~20 page loads with each page load generating ~100 web events.
- This drives a server volume of ~ 4.6 billion web events per day!

## Topology



# Things to Remember

- HW-accelerated virtualization
  - integrates into the existing world
  - efficient, secure, scalable
  - large benefits in current environments (consolidation, efficiency) and future ( $\mu$  services, highly scalable, on requirement)
- Choices:
  - take what ever is suited best
    - What is standard in your IT
    - What provides the best match to the requirements
    - Or both
  - Remember: OR, not XOR





***Arwed Tschoeke***

IBM Client Center –  
Systems and Software –  
z ATS  
IBM Germany Lab

*Schoenaicher Str. 220  
D-71032 Boeblingen*

*Phone +49 (0) 171 863 7780*

*[arwed.tschoeke@de.ibm.com](mailto:arwed.tschoeke@de.ibm.com)*