# Private Cloud by VMware ESXi (hypervisor)
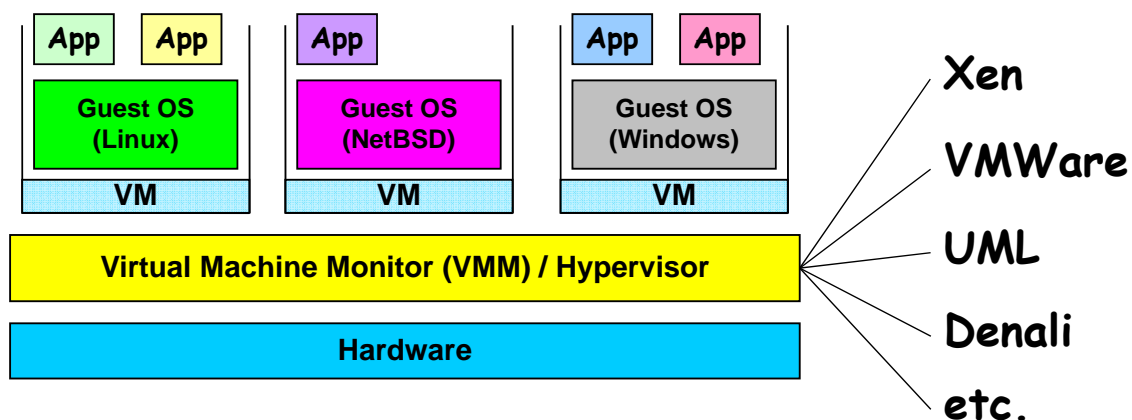
**Tien-Fu Chen**

**Dept. of Computer Science National Chiao Tung Univ.**

---

# Virtual Machine Management

- ❑ VM technology allows multiple virtual machines to run on a single physical machine.

| App | App | | App | | App | App |
|-----|-----|---|-----|---|-----|-----|
| Guest OS (Linux) | | | Guest OS (NetBSD) | | Guest OS (Windows) | |
| VM | | | VM | | VM | |

**Virtual Machine Monitor (VMM) / Hypervisor**

**Hardware**

Xen

VMWare

UML

Denali

etc.

*Performance*: Para-virtualization (e.g. Xen) is very close to raw physical performance!

# Virtualization in General

❑ Advantages of virtual machines:
- – Run operating systems where the physical hardware is unavailable,
- – Easier to create new machines, backup machines, etc.,
- – Software testing using "clean" installs of operating systems and software,
- – Emulate more machines than are physically available,
- – Timeshare lightly loaded systems on one host,
- – Debug problems (suspend and resume the problem machine),
- – Easy migration of virtual machines (shutdown needed or not).
- – Run legacy systems!

---

# Inside Look

| Citrix Xenserver | QEMU-KVM | Hyper-V (Azure*) | XenSource |
|---|---|---|---|
| ✓Free and Enterprise Versions | ✓FOSS | ✓Available on any Windows 2008R2 Server | ✓FOSS |
| ✓Unified Console | ✓Common API | ✓MMC Console | ✓Common API |
| •Hardware Virtualization | •Hardware Virtualization | •Hardware Virtualization | •Software Virtualization |
| ✓Easy Desktop Virtualization | ✓Plethora of Tools | ✓Easy Install and Usage | ✓Plethora of Tools |
| ✓Citrix Support and Application Integration | ✓API and Custom Application Friendly | ✓Integrates with most provided windows tools | ✓API and Custom Application Friendly |
| •Red Hat Based (Linux) Custom OS | •All Linux & Unix | •Windows Proprietary | •All Linux & Unix |
| -Strict Hardware Requirements | -Harder to Administer | -Hit and Miss Performance and Support Outside of Windows | -Performance Concerns |

**Table 3.6** VI Managers and Operating Systems for Virtualizing Data Centers [9]

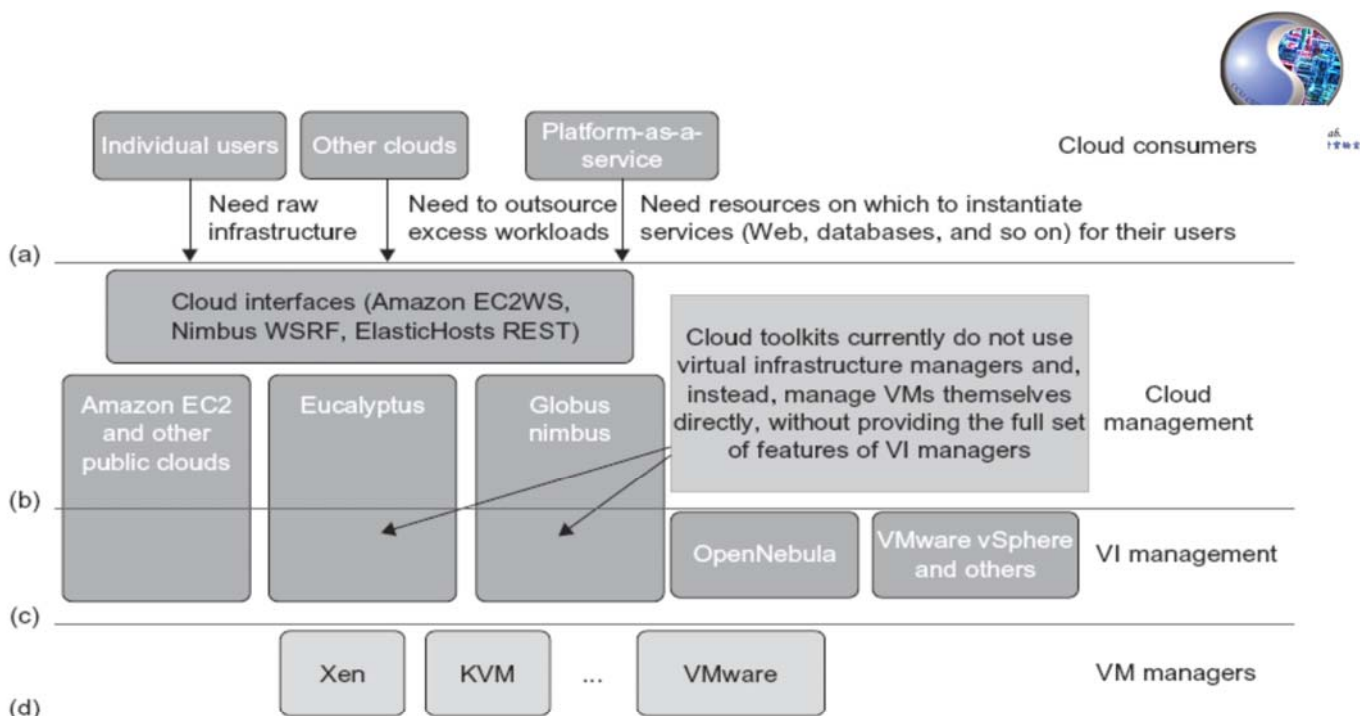| Manager/ OS, Platforms, License | Resources Being Virtualized, Web Link | Client API, Language | Hypervisors Used | Public Cloud Interface | Special Features |
|---|---|---|---|---|---|
| **Nimbus** Linux, Apache v2 | VM creation, virtual cluster, www.nimbusproject.org/ | EC2 WS, WSRF, CLI | Xen, KVM | EC2 | Virtual networks |
| **Eucalyptus** Linux, BSD | Virtual networking (Example 3.12 and [41]), www.eucalyptus.com/ | EC2 WS, CLI | Xen, KVM | EC2 | Virtual networks |
| **OpenNebula** Linux, Apache v2 | Management of VM, host, virtual network, and scheduling tools, www.opennebula.org/ | XML-RPC, CLI, Java | Xen, KVM | EC2, Elastic Host | Virtual networks, dynamic provisioning |
| **vSphere 4** Linux, Windows, proprietary | Virtualizing OS for data centers (Example 3.13), www.vmware.com/products/vsphere/ [66] | CLI, GUI, Portal, WS | VMware ESX, ESXi | VMware vCloud partners | Data protection, vStorage, VMFS, DRM, HA |

**FIGURE 4.4**

Cloud ecosystem for building private clouds: (a) Consumers demand a flexible platform; (b) Cloud manager provides virtualized resources over an IaaS platform; (c) VI manager allocates VMs; (d) VM managers handle VMs installed on servers.
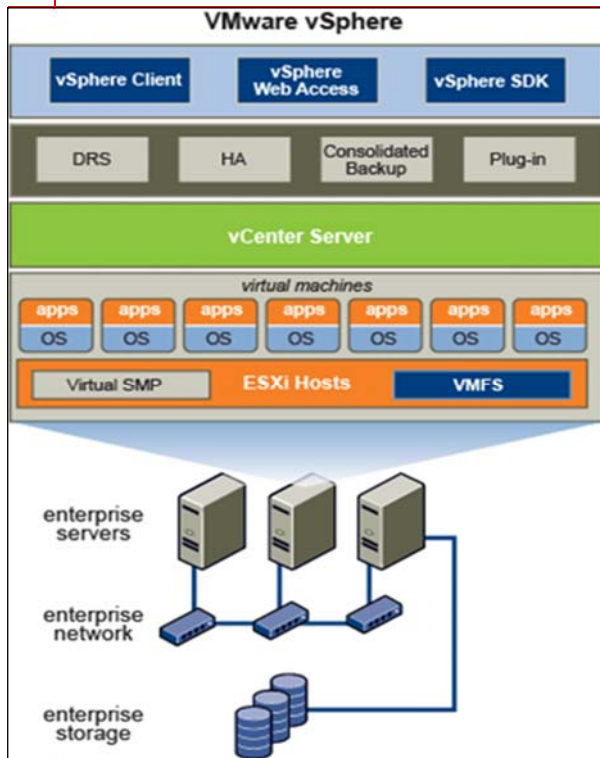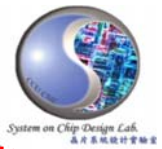
(Courtesy of Sotomayor, et al. [68])

# What Is VMware vSphere?



An infrastructure virtualization suite that provides virtualization, management, resource optimization, application availability, and operational automation capabilities

It consists of the following components:

- VMware ESXi
- VMware vCenter Server™
- VMware vSphere® Client™
- VMware vSphere® VMFS
- VMware vSphere® Virtual Symmetric Multiprocessing

**vmware**

Source: VMware vSphere: Overview

**7**

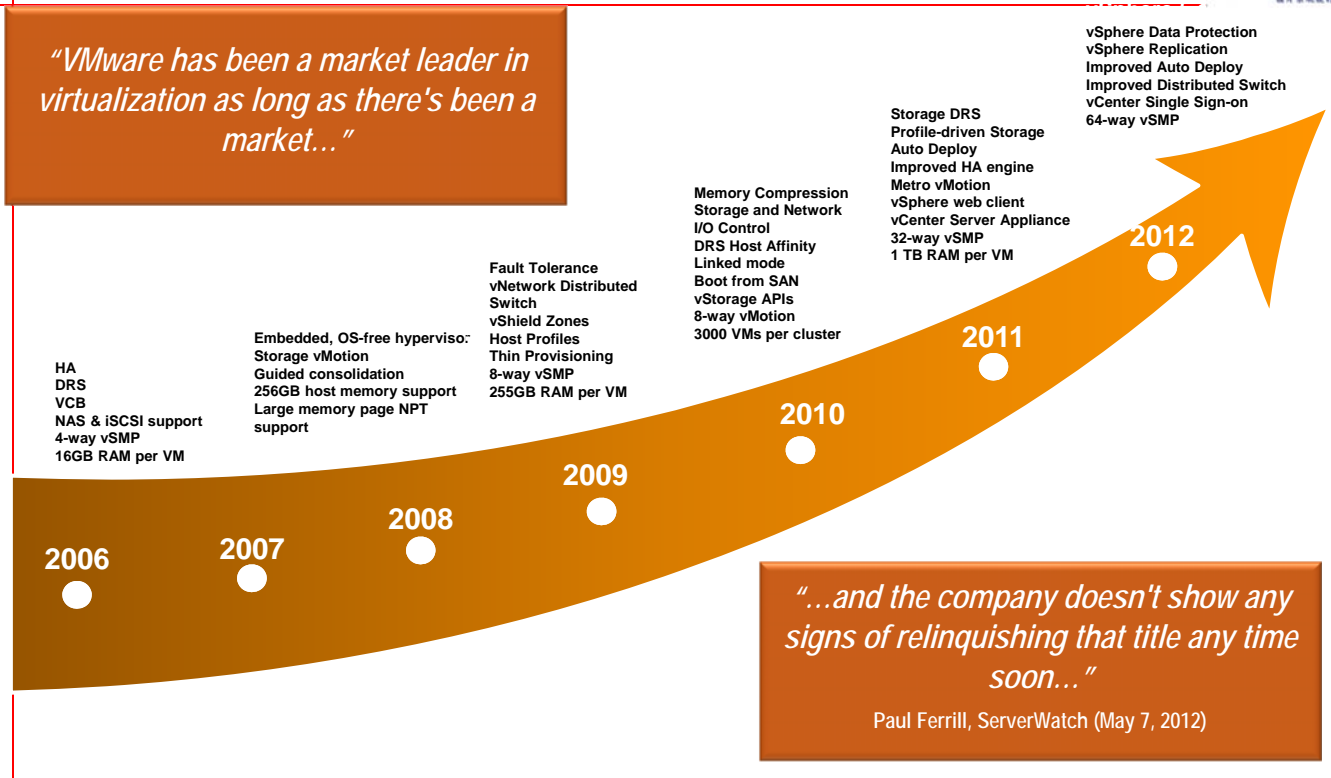T.-F. Chen@NCTU CSIE

---

# VMware vSphere: Most Comprehensive OS Support

| MS Hyper-V 3 | Citrix XenServer 6.1 | VMware vSphere 5.1 | |
|---|---|---|---|
| Windows Server 2003 (32/64) | Windows Server 2003 (32/64) | MS-DOS 6.22 | Oracle Linux 4 (32/64) |
| Windows Server 2008 (32/64) | Windows Server 2008 (32/64) | Windows 3.1 | Oracle Linux 5 (32/64) |
| Windows Server 2012 | Windows 7 (32/64) | Windows 95 | Oracle Linux 6 (32/64) |
| Windows Home Server 2011 | Windows Vista | Windows 98 | Asianux 3 (32/64) |
| Windows Small Business Server 2011 | Windows XP | Windows NT | Asianux 4 (32/64) |
| Windows 8 (32/64) | RHEL 4 | Windows XP (32/64) | Ubuntu 8 (32/64) |
| Windows 7 (32/64) | RHEL 5 (32/64) | Windows Vista (32/64) | Ubuntu 9 (32/64) |
| Windows Vista (32/64) | RHEL 6 (32/64) | Windows 7 (32/64) | Ubuntu 10 (32/64) |
| Windows XP (32/64) | SLES10 (32/64) | Windows 8 (32/64) | Ubuntu 11 (32/64) |
| RHEL 5 (32/64) | SLES11 (32/64) | Windows 2000 | Ubuntu 12 (32/64) |
| RHEL 6 (32/64) | Debian Squeeze 6 (32/64) | WinServer 2003 (32/64) | FreeBSD 6 (32/64) |
| SLES 11 (32/64) | CentOS 4 | WinServer 2008 (32/64) | FreeBSD 7 (32/64) |
| CentOS 5 (32/64) | CentOS 5 (32/64) | WinServer 2012 | FreeBSD 8 (32/64) |
| CentOS 6 (32/64) | CentOS 6 (32/64) | RHEL 2.1 | FreeBSD 9 (32/64) |
| Open SUSE 12 (32/64) | Oracle Linux 5 (32/64) | RHEL 3 (32/64) | Solaris 10 (32/64) |
| Ubuntu 12 (32/64) | Oracle Linux 6 (32/64) | RHEL 4 (32/64) | Solaris 11 |
| | Ubuntu 10 (32/64) | RHEL 5 (32/64) | IBM OS/2 Warp 4 |
| | Ubuntu 12 (32/64) | RHEL 6 (32/64) | NetWare 5 |
| | | SLES 8 | NetWare 6 |
| | | SLES 9 (32/64) | eComStation 1 |
| | | SLES 10 (32/64) | eComStation 2 |
| | | SLES 11 (32/64) | SCO UnixWare 7 |
| | | SLED 10 (32/64) | SCO OpenServer 5 |
| | | SLED 11 (32/64) | Mac OS X 10.6 (32/64) |
| | | Debian 4 (32/64) | Mac OS X 10.7 (32/64) |
| | | Debian 5 (32/64) | Mac OS X 10.8 (32/64) |
| | | Debian 6 (32/64) | |
| | | CentOS 4 (32/64) | |
| | | CentOS 5 (32/64) | |
| | | CentOS 6 (32/64) | |
| **Total: 29** | **Total: 32** | **Total: 95** | |

Data collected Dec 4, 2012

**vSphere = Most guest OSs; More versions of Windows**

See http://www.vmware.com/technical-resources/advantages/guest-os.html

# VMware: Consistent Delivery, Continuous Innovation

"VMware has been a market leader in virtualization as long as there's been a market…"

**2006**
HA
DRS
VCB
NAS & iSCSI support
4-way vSMP
16GB RAM per VM

**2007**
Embedded, OS-free hypervisor:
Storage vMotion
Guided consolidation
256GB host memory support
Large memory page NPT support

**2008**

**2009**
Fault Tolerance
vNetwork Distributed Switch
vShield Zones
Host Profiles
Thin Provisioning
8-way vSMP
255GB RAM per VM

**2010**
Memory Compression
Storage and Network I/O Control
DRS Host Affinity
Linked mode
Boot from SAN
vStorage APIs
8-way vMotion
3000 VMs per cluster

**2011**
Storage DRS
Profile-driven Storage
Auto Deploy
Improved HA engine
Metro vMotion
vSphere web client
vCenter Server Appliance
32-way vSMP
1 TB RAM per VM

**2012**
vSphere Data Protection
vSphere Replication
Improved Auto Deploy
Improved Distributed Switch
vCenter Single Sign-on
64-way vSMP

"…and the company doesn't show any signs of relinquishing that title any time soon…"

Paul Ferrill, ServerWatch (May 7, 2012)

---

# VMware ESXi: 3rd Generation Hypervisor Architecture

**VMware GSX (VMware Server)**
- Installs as an application
- Runs on a host OS
- Depends on OS for resource management

VM VM VM VM VM
VMware Server
Windows or Linux OS

**VMware ESX Architecture**
- Installs "bare metal"
- Relies on a Linux OS (Service Console) for running partner agents and scripting

VM VM VM VM VM
VMware ESX

**VMware ESXi Architecture**
- Installs "bare metal"
- Management tasks are moved outside of the hypervisor

VM VM VM VM VM
VMware ESXi

**2001**          **2003**                    **2007**

The ESXi architecture runs independently of a general purpose OS, simplifying hypervisor management and improving security.
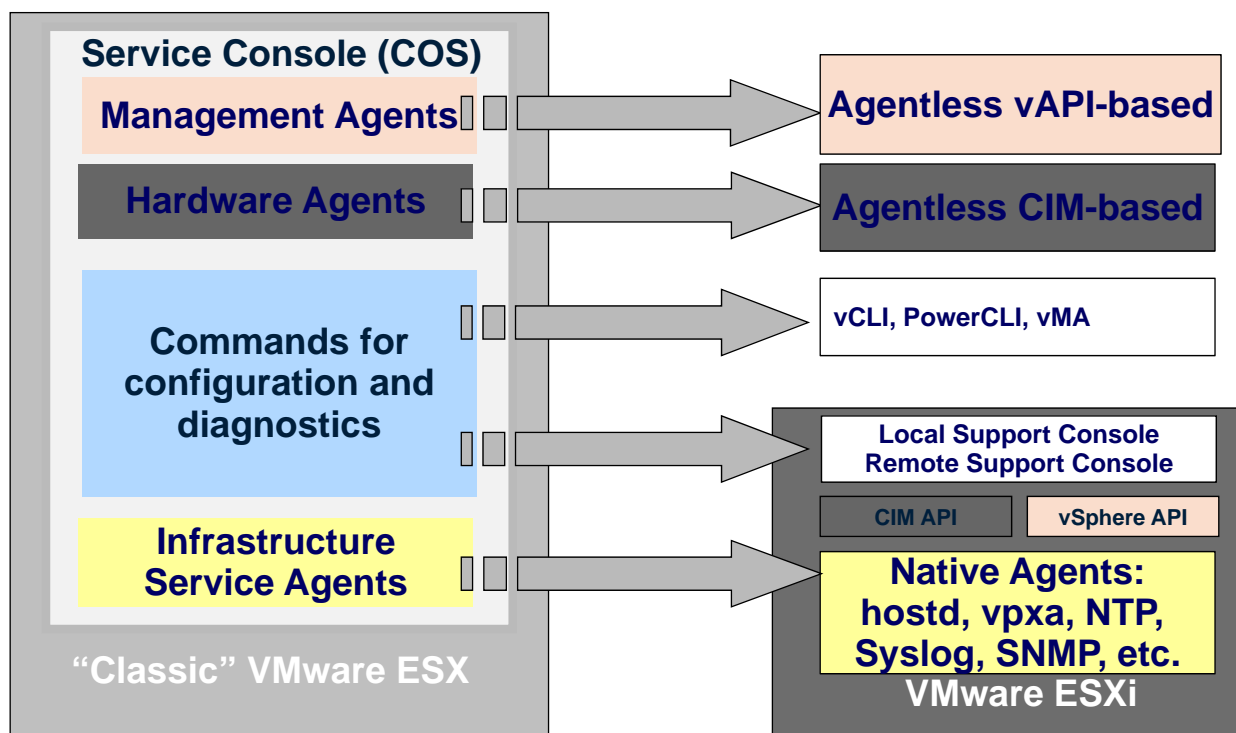
# VMware vSphere 5.0: *What's New?*

**vCenter Server**
- Virtual Appliance
- Web Client

**Application Services**

**VMware vSphere 5**

**Availability**
- New HA Architecture
- vMotion over higher latency links

**Security**
- ESXi Firewall

**Scalability**
- 32 way SMP
- 1 TB VMs

**Infrastructure Services**

**Compute**
- ESXi Convergence
- Auto Deploy
- HW version 8

**Storage**
- Storage DRS
- Profile-Driven Storage
- VMFS 5
- Storage I/O Control (NFS)

**Network**
- Network I/O Control (per VM controls)
- Distributed Switch (Netflow, SPAN, LLDP)

---

# ESX vs ESXi

**Service Console (COS)**

**Management Agents** → **Agentless vAPI-based**

**Hardware Agents** → **Agentless CIM-based**

**Commands for configuration and diagnostics** → vCLI, PowerCLI, vMA

→ **Local Support Console Remote Support Console**
CIM API | vSphere API

**Infrastructure Service Agents** → **Native Agents: hostd, vpxa, NTP, Syslog, SNMP, etc.**

**"Classic" VMware ESX**

**VMware ESXi**

# Install ESXi

**Tien-Fu Chen**

**Dept. of Computer Science
National Chiao Tung Univ.**

---

# Burn a VMware ESXi CD or bootable USB

- Download the VMware ESXi 5.5 ISO file from the VMware Download Center..

- Burn ISO onto your USB by Unetbootin or Rufus

- Enable virtualization VT-x in your bios

- Enable VT-d if your bios supports for Directed I/O

- If fail at "initializing IOV", enter noIOMMU (after shift-O)



```
ESXi-5.5.0-20140302001-standard Boot Menu

ESXi-5.5.0-20140302001-standard Installer
Boot from local disk




Press [Tab] to edit options
Automatic boot in 6 seconds...
```

Andy Barnes
VMadmin.co.uk

# Select a disk for VMFS datastore

- ❑ Select the correct storage device to install ESXi on and press "Enter"

- ❑ VMware VMFS (Virtual Machine File System)



15

# Set up initial network configuration

- ❑ Pressing F2 and entering your Root accounts password.

- ❑ Select Configure Management Network then IP configuration.

- ❑ Give your ESXI host a static address on your network.

- ❑ Once done select DNS configuration,

- ❑ Enter your DNS servers IP address

- ❑ Give your ESXI host a valid host name on your network.

# User Interfaces



**vSphere Client**

**Your desktop**

**Web Client**

**vCenter Server**

**ESXi host**

*vm*ware®

---

# Web Client



## Overview

- Run and manage vSphere from any web browser anywhere in the world

## Benefits

- Platform independence

- Replaces Web Access GUI

- Building block for cloud based administration

# vSphere Web Client Install

# Install VMware vSphere Client

❑ Download the files available at VMware Downloads.
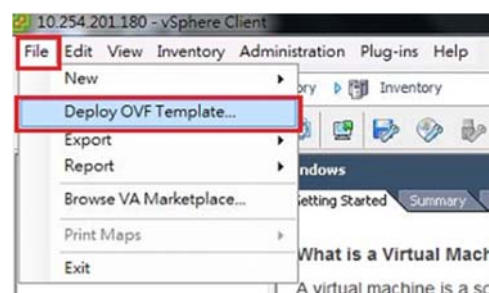
Create a VM

---

# Export existing VM to OVF/OVA

❑ File Formats for Virtual Machines

– Open Virtualization Format (OVF)

XML-based describing the properties of a virtual system. has generous allowances for extensibility

– Open Virtual Appliance (OVA)

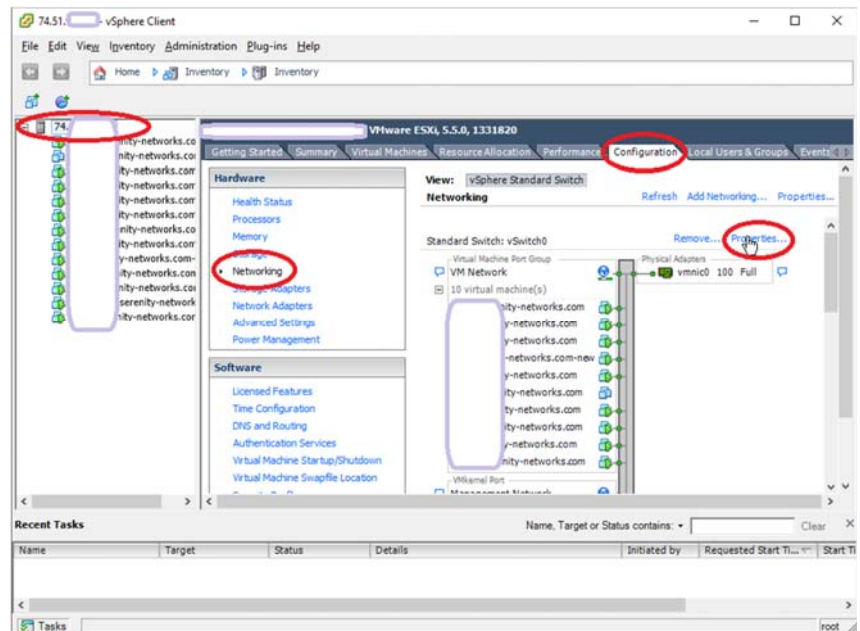An OVA is an OVF file packaged together with all of its supporting files (disk images, etc.).



22

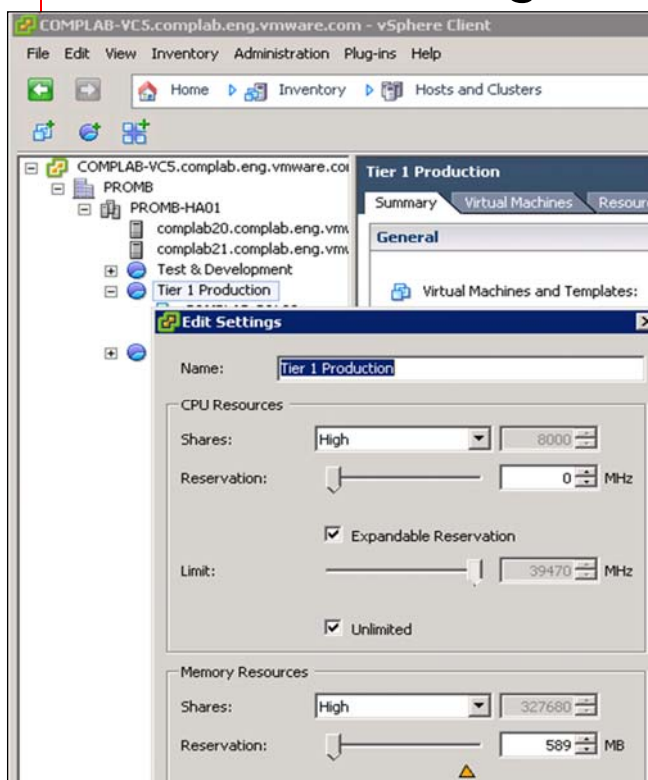# VM vswitch and port group

❑ Configuration > Networking> Properties of the vSwitch

❑ add a port group exclusive to the vLAN

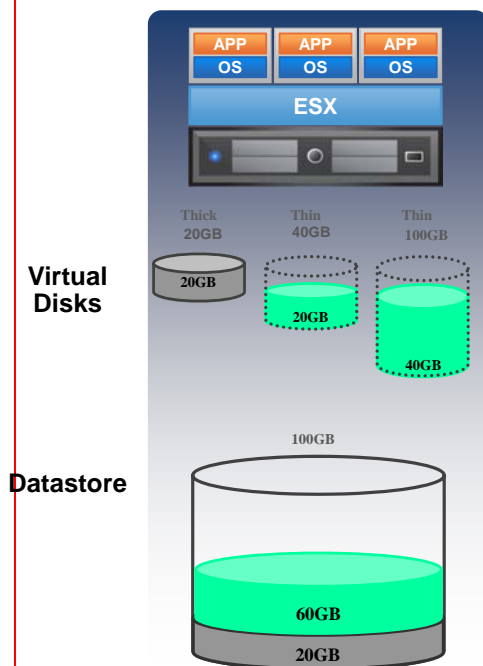❑ select a connection type.  Select Virtual Machine



**Cloud system**

---

# VMware vSphere Provides Advanced Resource Management



❑ Granular control of CPU and memory allocation with Resource Pools

❑ Easy to configure and view allocations

❑ Apply resource priority for multiple virtual machines
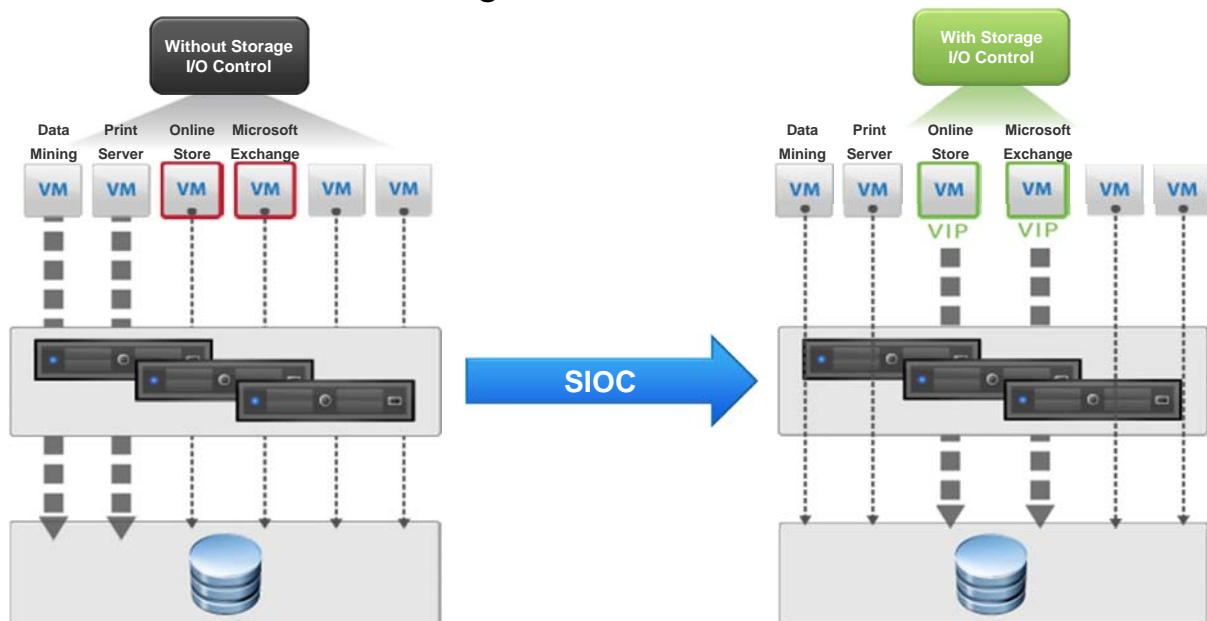
❑ Supports nesting

# vStorage Thin Provisioning

- Virtual machine disks consume only the amount of physical space in use
  - Virtual machine sees full logical disk size at all times
  - Full reporting and alerting on allocation and consumption
- Significantly improve storage utilization
- Eliminate need to over-provision virtual disks
- Reduce storage costs by up to 50%

---

# VMware vSphere Provides Advanced Resource Management
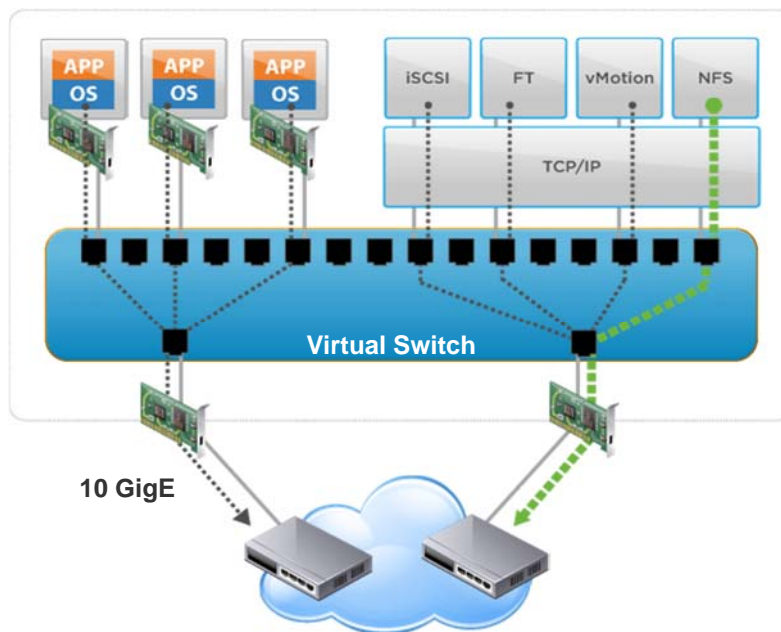
## Storage I/O Prioritization

During high I/O from non-critical application

# VMware vSphere Provides Advanced Resource Management

**Network I/O Prioritization**

---

# The Best of the Rest

❑ **Platform**
- Hardware Version 8 – EFI virtual BIOS

❑ **Network**
- Distributed Switch (Netflow, SPAN support, LLDP)
- Network I/O Controls (per VM), ESXi firewall

❑ **Storage**
- VMFS 5
- iSCSI UI
- Storage I/O Control (NFS)
- Array Integration for Thin Provisioning
- Swap to SSD, 2TB+ VMFS datastores
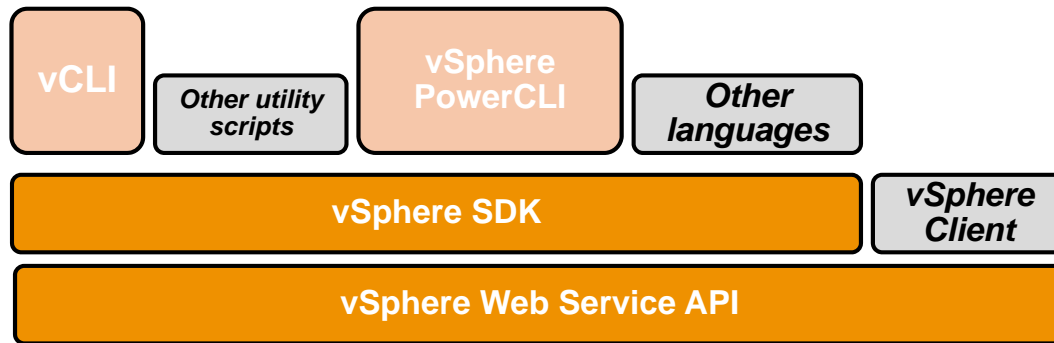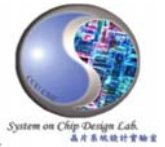- Storage vMotion Snapshot Support

❑ **Availability**
- vMotion with higher latency links
- Data Recovery Enhancements

❑ **Management**
- Inventory Extensibility
- Solution Installation and Management
- iPad client

# vCLI and PowerCLI: primary scripting interfaces

| vCLI | Other utility scripts | vSphere PowerCLI | Other languages |
|---|---|---|---|

| vSphere SDK | vSphere Client |
|---|---|

| vSphere Web Service API |
|---|

- ❏ vCLI and PowerCLI built on same API as vSphere Client
  - – Same authentication (e.g. Active Directory), roles and privileges, event logging
  - – API is secure, optimized for remote environments, firewall-friendly, standards-based
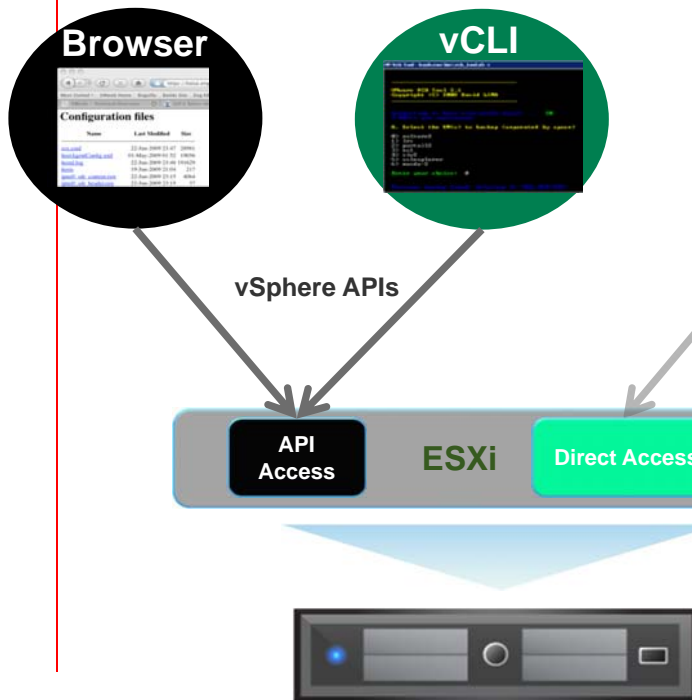
---

# Diagnostic Commands for ESXi: vCLI

- ❏ Familiar set of 'esxcfg-*' commands available in vCLI
  - – Names mapped to '**vicfg-*'**
  - – Also includes
    - ❏ **vmkfstools**
    - ❏ **vmware-cmd**
    - ❏ **resxtop**
    - ❏ **esxcli**: suite of diagnostic tools

# Summary of ESXi Diagnostics and Troubleshooting

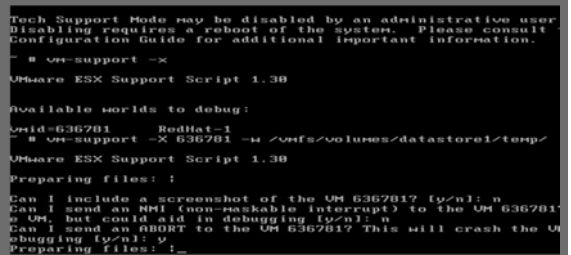**• During normal operations:**

**• If things go wrong:**

**Browser**

**vCLI**

vSphere APIs

**API Access**  **ESXi**  **Direct Access**

DCUI: misconfigs / restart mgmt agents

TSM: Advanced troubleshooting

---

# Virtualization Performance I

# CPU

# Performance Metrics

- CPU
  - Throughput: MIPS (%used), Goodput: useful instructions
  - Latency: Instruction Latency (cache latency, cache miss)
- Memory
  - Throughput: MB/Sec, Goodput: useful data
  - Latency: nanosecs
- Storage
  - Throughput: MB/Sec, IOPS/Sec, Goodput: useful data
  - Latency: Seek time
- Networking
  - Throughput: MB/Sec, IO/Sec, Goodput: useful traffic
  - Latency: microseconds

---

# CPU – Overview

- ❑ Raw processing power of a given host or VM
  - Hosts provide CPU resources
  - VMs and Resource Pools consume CPU resources

- ❑ CPU cores/threads need to be shared between VMs

- ❑ Fair scheduling vCPU time
  - Hardware interrupts for a VM
  - Parallel processing for SMP VMs
  - I/O

# CPU – esxtop

```
10:10:36am up 28 days  3:28, 321 worlds, 5 VMs, 7 vCPUs; CPU load average: 0.01, 0.01, 0.01
PCPU USED(%): 6.0 1.2 0.8 0.9 0.2 0.2 2.4 1.9 0.4 1.3 0.3 0.9 AVG: 1.4
PCPU UTIL(%): 9.4 3.7 2.4 2.7 0.8 0.6 5.2 6.2 1.5 4.4 1.1 2.9 AVG: 3.4

     ID      GID NAME              NWLD   %USED    %RUN   %SYS    %WAIT %VMWAIT    %RDY   %IDLE  %OVRLP   %CSTP  %MLMTD  %SWPWT
      1        1 idle                12 1127.07 1200.00   0.01     0.00       - 1200.00    0.00    1.94    0.00    0.00    0.00
 697664   697664 DC                   5    4.90    6.18   0.05   476.03    0.25    0.33   90.14    0.03    0.00    0.00    0.00
 744427   744427 RedHat 5.5           5    3.16    8.32   0.19   474.13    0.49    0.10   87.86    0.01    0.00    0.00    0.00
1324719  1324719 vIN                  6    1.62    3.99   0.15   574.55    0.00    0.52  189.10    0.02    0.00    0.00    0.00
1073009  1073009 UI VM               6    1.55    3.80   0.14   574.76    0.00    0.49  189.27    0.02    0.00    0.00    0.00
  17742    17742 vCOPs standalon      5    1.42    3.67   0.06   478.58    0.00    0.30   92.88    0.01    0.00    0.00    0.00
1369428  1369428 esxtop.1681008       1    0.96    1.10   0.00    95.41       -    0.00    0.00    0.01    0.00    0.00    0.00
    756      756 hostd.2825          20    0.48    0.92   0.00  1929.09       -    0.18    0.00    0.00    0.00    0.00    0.00
1069135  1069135 vpxa.948012         19    0.28    0.58   0.01  1832.94       -    0.17    0.00    0.00    0.00    0.00    0.00
1069450  1069450 fdm.1310934         18    0.08    0.21   0.01  1736.84       -    0.12    0.00    0.00    0.00    0.00    0.00
      2        2 system              10    0.04    0.10   0.00   964.98       -    0.04    0.00    0.00    0.00    0.00    0.00
      8        8 helper              75    0.03    0.09   0.00  7238.18       -    0.06    0.00    0.00    0.00    0.00    0.00
    606      606 vmsyslogd.2659       3    0.02    0.04   0.00   289.48       -    0.00    0.00    0.00    0.00    0.00    0.00
1369424  1369424 sshd.1683052         1    0.01    0.03   0.00    96.48       -    0.00    0.00    0.00    0.00    0.00    0.00
    713      713 vmware-usbarbit      2    0.01    0.03   0.00   192.99       -    0.01    0.00    0.00    0.00    0.00    0.00
    645      645 vmkiscsid.2703       2    0.01    0.02   0.00   192.98       -    0.02    0.00    0.00    0.00    0.00    0.00
      9        9 drivers             11    0.01    0.02   0.00  1061.58       -    0.02    0.00    0.00    0.00    0.00    0.00
    732      732 net-lbt.2803         1    0.01    0.02   0.00    96.49       -    0.00    0.00    0.00    0.00    0.00    0.00
    679      679 ntpd.2748            2    0.01    0.02   0.00   192.99       -    0.01    0.00    0.00    0.00    0.00    0.00
   1090     1090 openwsmand.3207      3    0.01    0.02   0.00   289.50       -    0.01    0.00    0.00    0.00    0.00    0.00
    978      978 dcbd.3062            1    0.00    0.01   0.00    96.49       -    0.01    0.00    0.00    0.00    0.00    0.00
   1463     1463 sfcb-ProviderMa     10    0.00    0.01   0.00   965.08       -    0.01    0.00    0.00    0.00    0.00    0.00
    776      776 vprobed.2849         3    0.00    0.01   0.00   289.52       -    0.00    0.00    0.00    0.00    0.00    0.00
    853      853 storageRM.2931       2    0.00    0.00   0.00   193.01       -    0.00    0.00    0.00    0.00    0.00    0.00
   1016     1016 vobd.3101           15    0.00    0.00   0.00  1447.63       -    0.00    0.00    0.00    0.00    0.00    0.00
   1461     1461 sfcb-ProviderMa      8    0.00    0.00   0.00   772.07       -    0.00    0.00    0.00    0.00    0.00    0.00
```
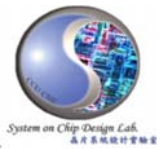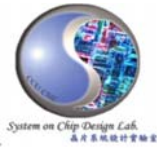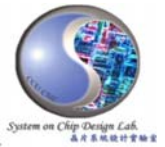
---

# CPU – esxtop

- Interpret the esxtop columns correctly

- %RDY - The percentage of time a VM is ready to run, but no physical processor is ready to run it which may result in decreased performance
- %USED – Physical CPU usage
- %SYS – Percentage of time in the VMkernel
- %RUN – Percentage of total scheduled time to run
- %WAIT – Percentage of time in blocked or busy wait states
- %IDLE – %WAIT- %IDLE can be used to estimate I/O wait time

# CPU – Performance Overhead & Utilization

❑ Different workloads have different overhead costs (%SYS) even for the same utilization (%USED)

❑ CPU virtualization adds varying amounts of system overhead

- Direct execution vs. privileged execution
- Non-paravirtual adapters vs. emulated adaptors
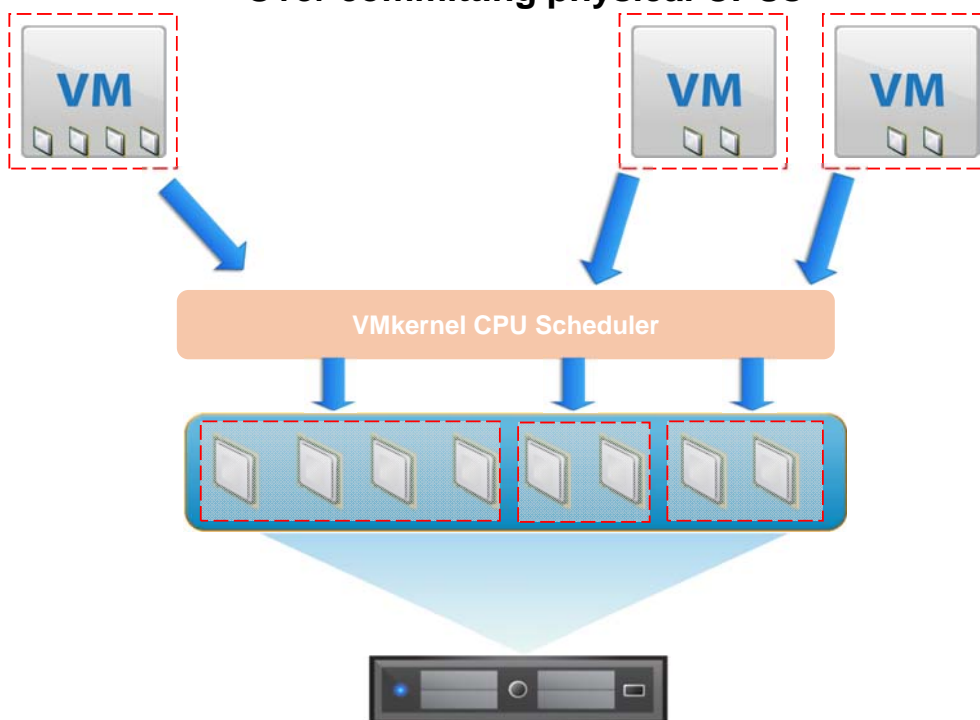- Virtual hardware (Interrupts!)
- Network and storage I/O

---

# CPU – vSMP

❑ Relaxed Co-Scheduling: vCPUs can run out-of-sync

❑ Idle vCPUs incur a scheduling penalty
- configure only as many vCPUs as needed
- Imposes unnecessary scheduling constraints
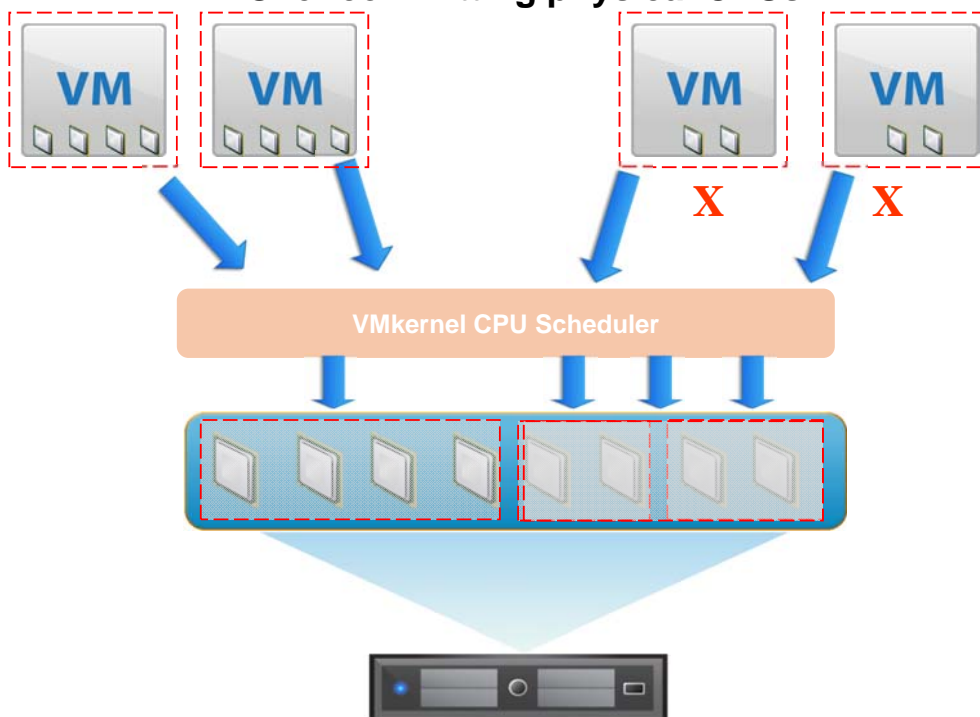
❑ Use Uniprocessor VMs for single-threaded applications

# CPU– Scheduling
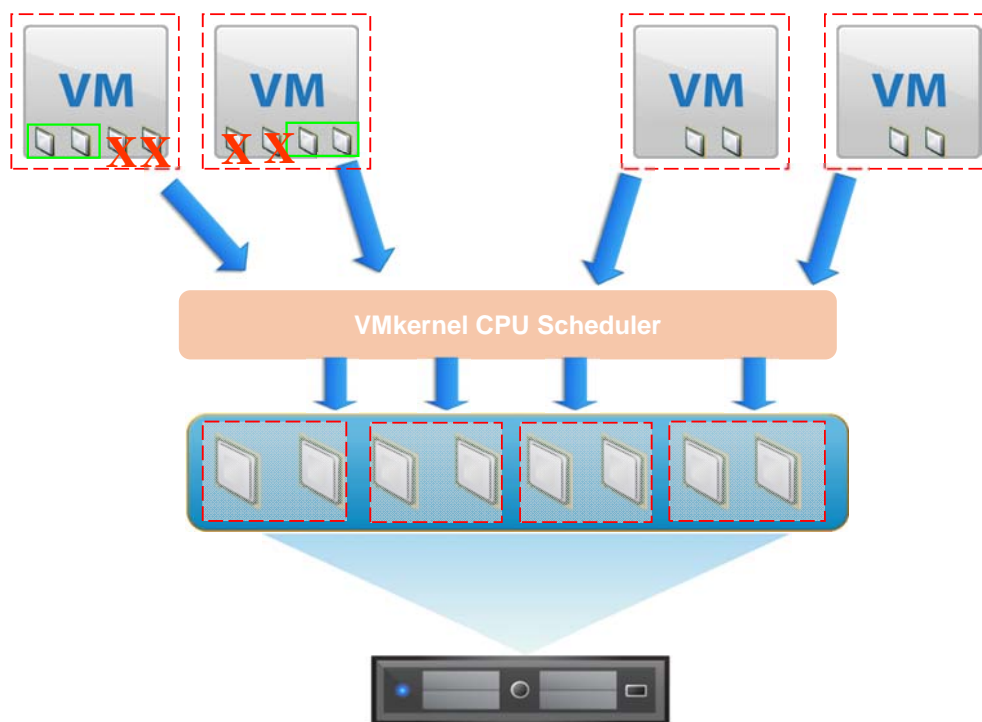
**Over committing physical CPUs**

# CPU– Scheduling

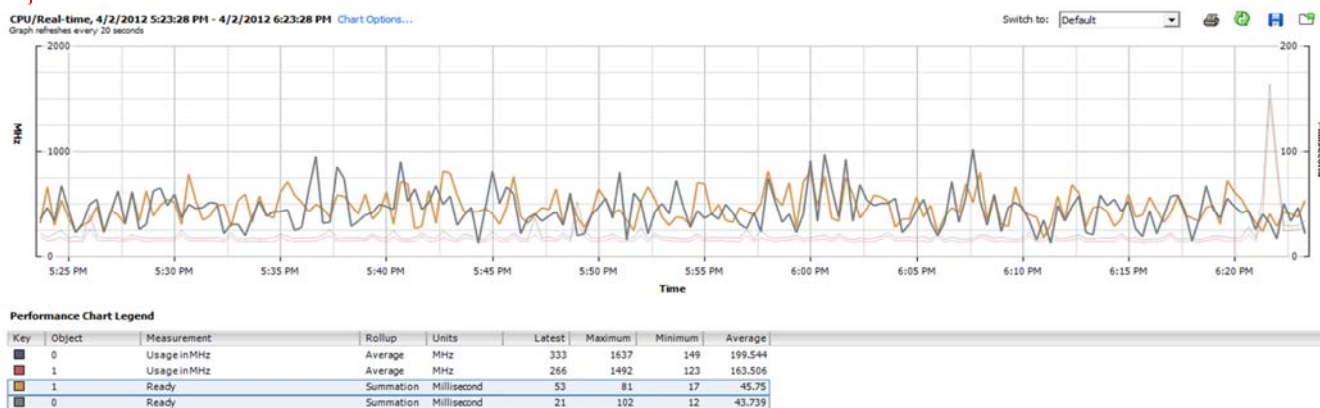**Over committing physical CPUs**

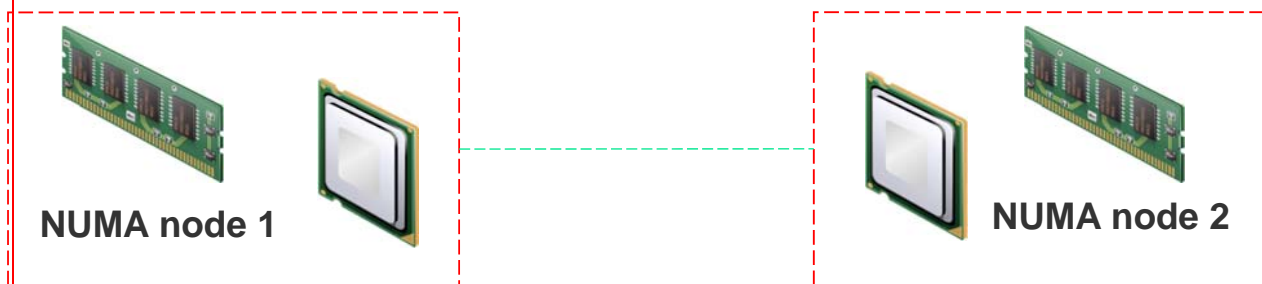# CPU– Scheduling

**Over committing physical CPUs**

---

# CPU – Ready Time

❑ The percentage of time that a vCPU is ready to execute, but waiting for physical CPU time

❑ Does not necessarily indicate a problem
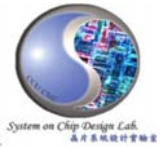 – Indicates possible CPU contention or limits

# CPU – NUMA nodes

❏ Non-Uniform Memory Access system architecture

❏ Each node consists of CPU cores and memory

❏ A CPU core in one NUMA node can access memory in another node, but at a small performance cost

**NUMA node 1**

**NUMA node 2**

---

# CPU – Troubleshooting

● vCPU to pCPU over allocation
   – HyperThreading does not double CPU capacity!

● Limits or too many reservations
   • can create artificial limits.

● Expecting the same consolidation ratios with different workloads
   • Virtualizing "easy" systems first, then expanding to heavier systems

● Compare Apples to Apples
   • Frequency, turbo, cache sizes, cache sharing, core count, instruction set…

# Demystifying "Ready" time

- ❑ Powered on VM could be either running, halted or in a ready state
- ❑ Ready time signifies the time spent by a VM on the run queue waiting to be scheduled
- ❑ Ready time accrues when more than one world wants to run at the same time on the same CPU
    - PCPU, VCPU over-commitment with CPU intensive workloads
    - Scheduler constraints - CPU affinity settings
- ❑ Higher ready time reduces response times or increases job completion time
- ❑ Total accrued ready time is not useful
    - VM could have accrued ready time during their runtime without incurring performance loss (for example during boot)
- ❑ %ready = ready time accrual rate

---

# Resource Over-Commitment

- ❑**CPU Over-Commitment**
    - Higher CPU utilization does not necessarily mean lesser performance.
        - ❑ Application's progress is not affected by higher CPU utilization
        - ❑ However if higher CPU utilization is due to monitor overheads then it may impact performance by increasing latency
        - ❑ When there is no headroom (100% CPU), performance degrades
    - 100% CPU utilization and %ready are almost identical – both delay application progress
    - CPU Over-Commitment could lead to other performance problems
        - ❑ Dropped network packets
        - ❑ Poor I/O throughput
        - ❑ Higher latency, poor response time