Designing Your VMware Virtual Infrastructure for Optimal Performance, Resilience and Availability – Straight from the Source Deji Akomolafe – VMware

David Klee – Heraflux Technologies Cody Chapman – Heraflux Technologies

usenix LISA16

December 4–9, 2016 | Boston, MA www.usenix.org/lisa16 #lisa16





Staff Solutions Architect, VMware Global Technical and Professional Services

- Microsoft Applications Virtualization Lead
- Member of VMware CTO Ambassador Program
- 20+ years IT experience, specializing in Microsoft technologies
- Former Microsoft MVP in multiple designations
 - Exchange Server
 - Directory Services
 - Windows Security
- Speaker at:
 - VMworld | EMCworld | VMUG | SQL Saturday



https://blogs.vmware.com/apps

http://www.dejify.com



usenix

LISA16

http://bit.ly/2h3Rf53



mware

About David Klee



- Performance Tuning & Troubleshooting
- Virtualization
- Cloud Enablement
- Infrastructure Architecture
- Health & Efficiency
- Capacity Management





About Cody Chapman



@codyrchapman
 heraflux.com
 linkedin.com/in/codyrchapman

- Performance Tuning & Troubleshooting
- Virtualization
- Infrastructure Architecture
- Scripting and Automation
- Health & Efficiency





Things we can all agree on

Virtualization is mainstream

You want to virtualize your applications

You care about the outcome

Your applications are **Important**

That is <u>WHY</u> we are here



Is the Application "Critical"?

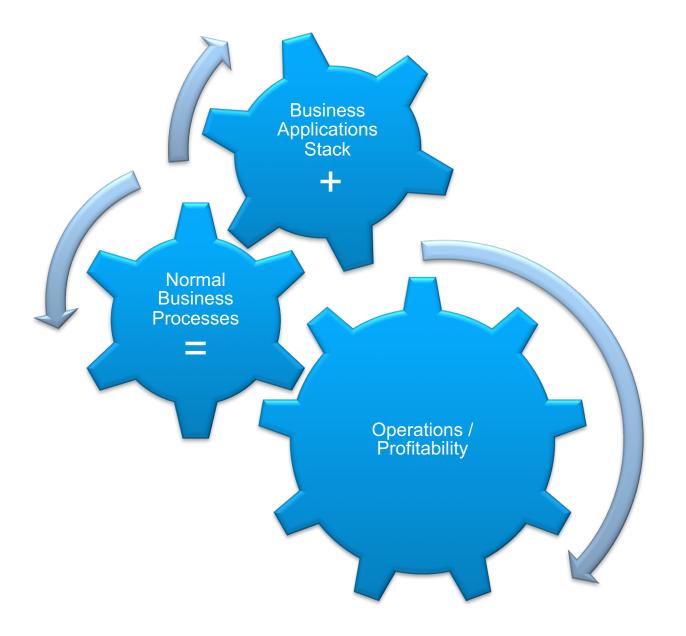


Is the outage impactful?

Outage <u>NOT</u> easily survivable?

Outage <u>NOT</u> easily recoverable?

Will it/you be missed?





Business Critical Applications Characteristics



LISA16

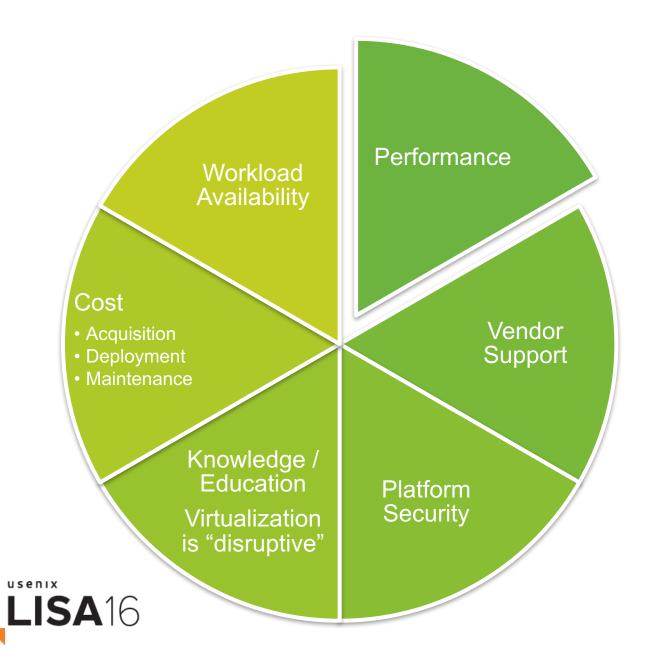
Why Virtualize Critical Applications

Resource Maximization	 Server resources increase too much for one application instance Virtualization improves resource utilization Reduces wastage
Enhanced Availability	 Native application HA features incomplete for most critical applications vSphere HA features complement native App HA features Result is improved availability
Dev Testing	 Virtualization improves adaptivity and elasticity Lifecycle management easier in virtual (provisioning/de-provisioning)
Rapid Provisioning And Scaling	 All the known and latent benefits of virtualization Project lifecycle considerably reduced
Job Security	It's 2016, and all the cool kids have done it You can't get to the "Cloud" without virtualizing
Lower TCO	Significant savings in power, cooling, and datacenter space, and administration
LISA16	

Common Objections To Virtualizing Critical Applications



Common Objections to Virtualizing BCA





Common Objections to Virtualization - Vendor Support

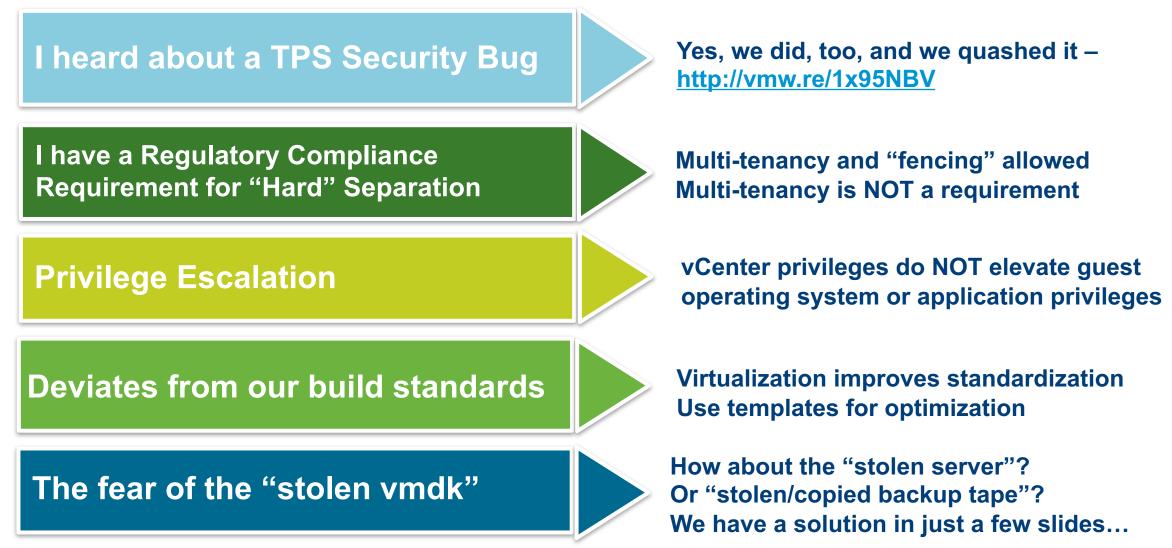
Vendor	Reference
Everything Business Critical Applications on VMware vSphere <u>http://www.vmware.com/business-critical-apps/</u> 	http://vmw.re/15MO7oL
Microsoft Supports Virtualization of ALL its Critical Applications	http://bit.ly/1uvVRkk
Exchange Server	http://bit.ly/1H1xYfu
SQL Server	http://bit.ly/15MrBMy
Oracle mySupport (Note 249212.1)	http://bit.ly/15DrLW3
SAP General Support Statement for Virtual Environments (Note 1492000)	http://bit.ly/1Ctkd4T
SAP on VMware	http://bit.ly/15NEiH4
SAP Notes Related to VMware	http://bit.ly/1wyohKe

For when you are in a jam http://www.tsanet.org





Common Objections to Virtualization - Security





Stolen VMDK? Meet VM Encryption

The "Dye Pack" of Enterprise Virtualization

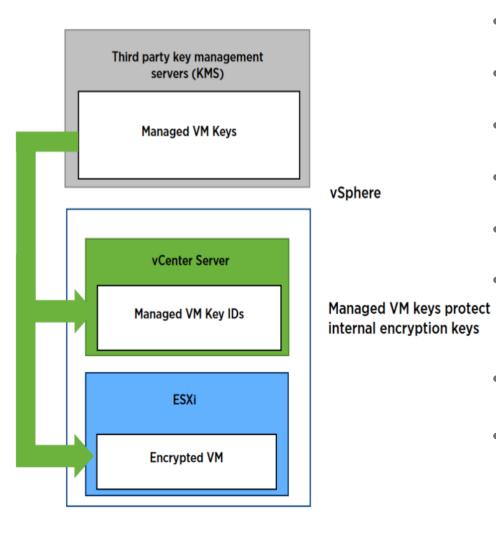




- Introduced in vSphere 6.5
- Secures Data in a VM's VMDK
- Uses vSphere APIs for I/O filtering (VAIO)
- VM Possesses Decryption Key
- vCenter Serves as Broker/Facilitator Only
- Data Meaningless to Unauthorized Entities
- No SPECIAL Hardware Required *

AES-NI Capable Server Hardware Improves Performance

VM Encryption – How it Works



usenix

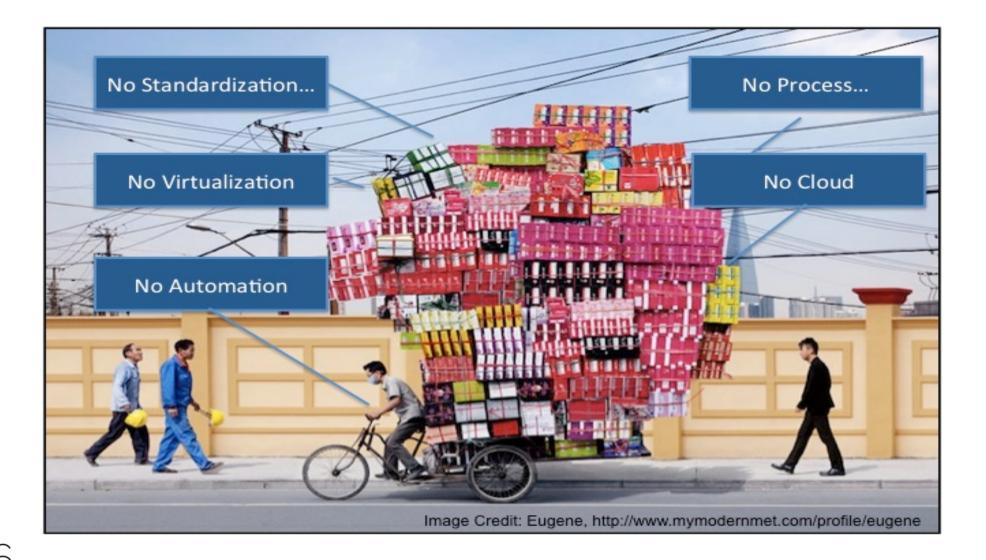
LISA16

- Customer-Supplied Key Management Server (KMS)
 - Customer-owned and Operated
- Centralized Repository for Crypto Keys
 - No Special Requirement KMIP 1.1-compliant
- KMS Clusters can be created
 - For Redundancy and Availability
- vCenter is Manually Enrolled in KMS
 - Establishing Trust
- vCenter Obtains Crypto KEKs from KMS
 - Distributes KEKs to ESXi
- ESXi Uses KEK to Generate DEK
 - Used for Encrypting/Decrypting VM Files
 - Encrypted DEKs Stored in VM Config Files
 - KEK for VMs Resides in ESXi's Memory
- IF ESXi Powered-Cycled (or Otherwise Unavailable), vCenter <u>Must Request</u> New KEK for Host
- If Encrypted VM Unregistered, vCenter <u>Must Request</u> KEK During Re-Registration

VM Unable to Power-On if Request Fails

Common Objections to Virtualization - Knowledge / Education

The Fear of Change.... Leads to inertia



LISA16

Virtualizing Applications for Performance and Scale



Can vSphere handle the load?

Configuration Item	ESXi 6.0	ESXi 6.5
Virtual CPUs per virtual machine (Virtual SMP)	128	128
RAM per virtual machine	4TB	6TB
Virtual machine swapfile size	4TB	6TB
Logical CPUs per host	480	576
Virtual CPUs per host	4096	4096
Virtual machines per host	1024	1024
Virtual CPUs per core	32	32
Virtual CPUs per FT virtual machine	4	4
FT Virtual machines per host	4	4
RAM per host	4TB	6TB
Hosts per cluster	64	64
Virtual Machines per cluster	8000	8000
LUNs per cluster/host	254	512
Paths per cluster/host	1024	2048
LUN / VMDK Size	62 TB	62 TB
Virtual NICs per virtual machine	10	10

LISA16

Ensuring Application Performance on vSphere

Physical Hardware

- VMware HCL
- BIOS / Firmware
- Power / C-States
- Hyper-threading
- NUMA

usenıx

LISA16

ESXi Host

- Power
- Virtual Switches/Portgroups
- vMotion Portgroups

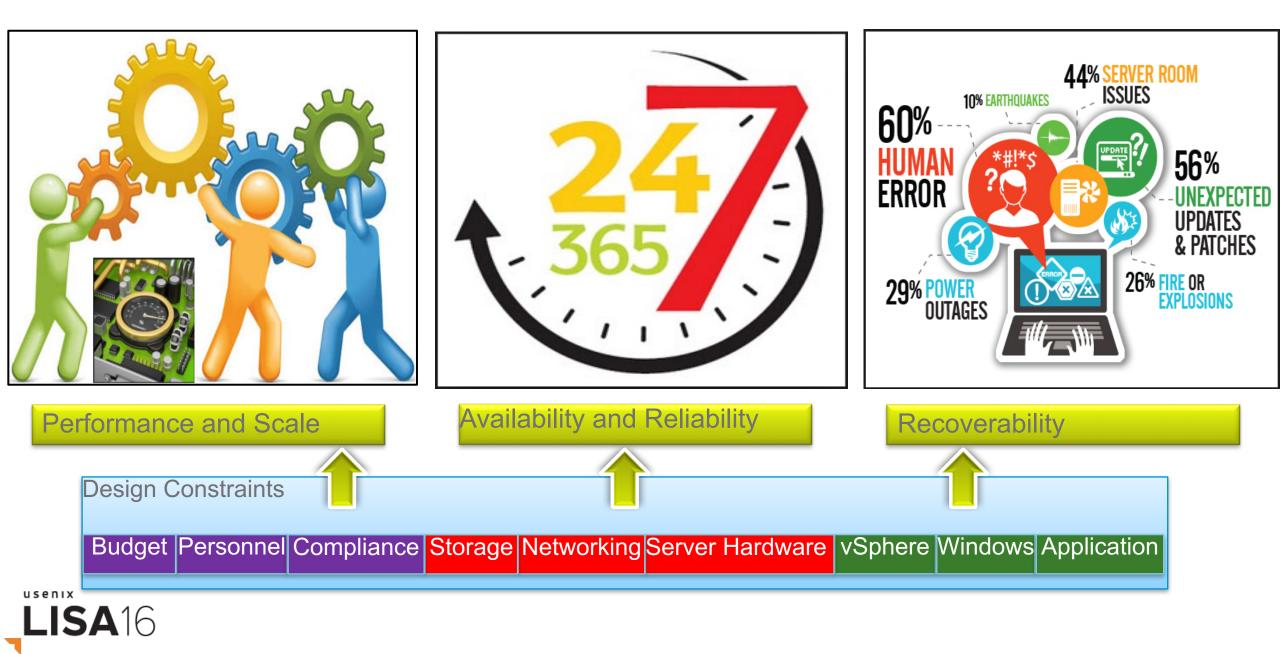
Virtual Machine

- Resource Allocation
- Storage
- Memory
- CPU / vNUMA
- Networking
- vSCSI Controller

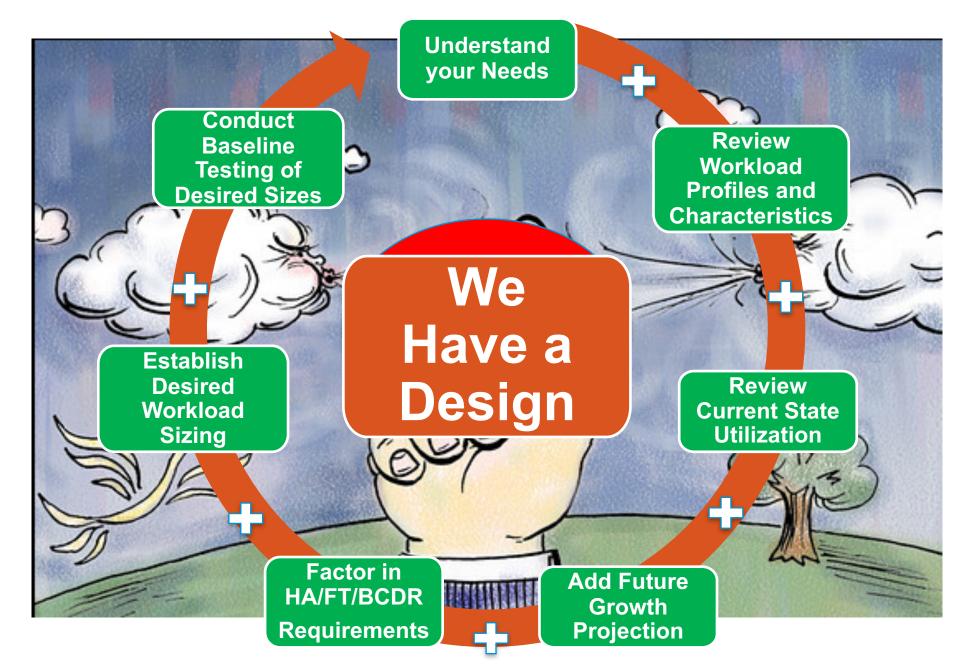
Guest Operating System

- Power
- CPU
- Networking
- Storage IO

Designing to Requirements – Know the Constraints



Performance-based Designing Tenets

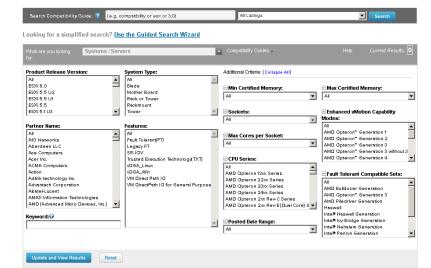


LISA16

Everything rides on the physical hardware – E.V.E.R.Y.T.H.I.N.G

- Physical Hardware
 - Hardware MUST Be On VMware's HCL
 - · Outdated drivers, firmware and BIOS Revs adversely impact virtualization
 - Always Disable unused physical hardware devices
 - Leave memory scrubbing rate in BIOS at default
 - Incorrect firmware, BIOS and Drivers Revs adversely impact virtualization
 - Default hardware Power Scheme unsuitable for virtualization
 - Change Power setting to "OS controlled"
 - Set ESXi Power Management Policy to "High Performance"
 - Enable Turbo Boost (or Equivalent)
 - Disable Processor C-states / C1E halt State
 - Enable All Cores Don't let hardware turn off cores dynamically

VMware Compatibility Guide



Looking for information on VMware product compatibility by version? <u>See the Product Interoperability Matrix</u> Looking for products verified and supported by partners? <u>Partner Verified and Supported Products</u>

WRONG BIOS, FIRMWARE, AND DRIVERS REVS ADVERSELY IMPACT VIRTUALIZATION



Time-Keeping in your vSphere Infrastructure

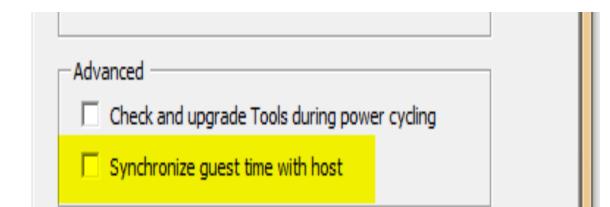


Back in the Days.....

lardware Options Resources	S VServices	Virtual Machine Version: 10
Settings	Summary	Power Controls
General Options	Deji-Lat-SensVM	Shut Down Guest
App Options	Disabled	
/Mware Tools	Shut Down	Suspend 🔽
Power Management	Standby	Power on / Resume virtual machine
dvanced		
General	Normal	Restart Guest
CPUID Mask	Expose Nx flag to	Due Uthurus Teels Carista
Memory/CPU Hotplug	Disabled/Disabled	Run VMware Tools Scripts
Boot Options	Normal Boot	After powering on
Fibre Channel NPIV	None	
CPU/MMU Virtualization	Automatic	After resuming
Swapfile Location	Use default settings	
		Before suspending
		Before shutting down Guest
		Advanced
		Check and upgrade Tools during power cycling
		Synchronize guest time with host



That was Problematic





But, That, Too, Is Insufficient

Because Even When You Do THAT, We Still Do THIS

Disabling Time Synchronization (1189)

Details

In the VMware Tools control panel, the time synchronization checkbox is unselected, but you may experience these symptoms:

When you suspend a virtual machine, the next time you resume that virtual machine it synchronizes the time to adjust it to the host.

Time is resynchronized when you migrate the virtual machine using vMotion, take a snapshot, restore to a snapshot, shrink the virtual
disk, or restart the VMware Tools service in the virtual machine (including rebooting the virtual machine).

Reference: <u>http://kb.vmware.com/kb/1189</u>



Preventing Bad Time Sync

- Ensure Hardware Clock on ESXi Hosts is CORRECT
- Configure Reliable NTP on ALL ESXi Hosts
- Configure in-Guest NTP Source
- IF Internal Authoritative Time Source Virtualized
 - (e.g.) Windows Active Directory PDC
 - Disable DRS for the VM
 - Use Host-Guest Affinity Rule for the VM
 - Helps you find it in Emergency



Completely Disabling Time Sync

Add the Following VM's Advanced Configuration Options to your VMs/Templates

```
tools.syncTime = "0"
time.synchronize.continue = "0"
time.synchronize.restore = "0"
time.synchronize.resume.disk = "0"
time.synchronize.shrink = "0"
time.synchronize.tools.startup = "0"
time.synchronize.tools.enable = "0"
```

To add these settings across multiple VMs at once, use <u>VMware vRealize Orchestrator</u>:

http://blogs.vmware.com/apps/2016/01/completely-disable-time-synchronization-for-your-vm.html



Designing for Performance

• NUMA

usenix

_ISA16

- To enable or to not enable? Depends on the Workloads
- More on NUMA later

Sockets, Cores and Threads

- Enable Hyper-threading
- Size to physical cores, not logical hyper-threaded cores.

Reservation, Limits, Shares and Resource Pools

- Use reservation to guarantee resources IF mixing workloads in clusters
- Use limits CAREFULLY for non-critical workloads
 - Limits <u>must</u> never be less than Allocated Values *
- Use Shares on Resource Pools
 - Only to contain non-critical Workload's consumption rate
 - Resource Pools must be continuously managed and reviewed
 - Avoid nesting Resource Pools complicates capacity planning

*Only possible with scripted deployment



Designing for Performance

Network

- Use VMXNET3 Drivers
 - VMXNET3 Template Issues in Windows 2008 R2 <u>kb.vmware.com\kb\1020078</u>
 - Hotfix for Windows 2008 R2 VMs <u>http://support.microsoft.com/kb/2344941</u>
 - Hotfix for Windows 2008 R2 SP1 VMs <u>http://support.microsoft.com/kb/2550978</u>
 - Remember Microsoft's "Convenience Update"? <u>https://support.microsoft.com/en-us/kb/3125574</u>
- Disable interrupt coalescing at vNIC level
- On 1GB network, use dedicated physical NIC for different traffic type

Storage

usenix

LISA16

- Latency is king Queue Depths exist at multiple paths (Datastore, vSCSI, HBA, Array)
- Adhere to storage vendor's recommended multi-pathing policy
- Use multiple vSCSI controllers, distribute VMDKS evenly
- Disk format and snapshot
- Smaller or larger datastores?
 - Determined by storage platform and workload characteristics (VVOL is the future)
- IP Storage? Jumbo Frames, if supported by physical network devices



The more you know...

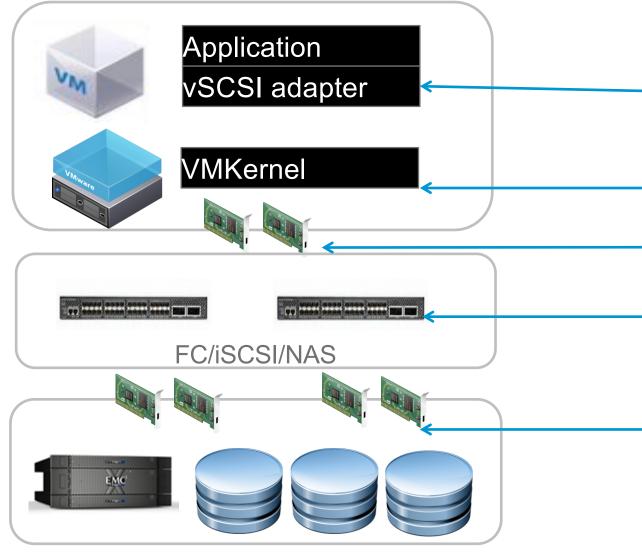
It's the Storage, Stupid	 There is ALWAYS a Queue One-lane highway vs 4-Lane highway. More is better PVSCSI for all data ask volumes Ask Your Storage Vendor about multi-pathing policy
More is NOT Better	 Know your hardware NUMA boundary. Use it to guide your sizing Beware of the memory tax Beware of CPU fairness There is no place like 127.0.0.1 (VM's Home Node)
Don't Blame the vNIC	 VMXNET3 is NOT the problem Outdated VMware Tools <u>MAY</u> be the problem Check in-guest network tuning options – e.g. RSS Consider Disabling Interrupt Coalescing
Use Your Tools	 Virtualizing does NOT change OS/App administrative tasks ESXTop – Native to ESXi Visualesxtop - <u>https://labs.vmware.com/flings/visualesxtop</u> Esxplot - <u>https://labs.vmware.com/flings/esxplot</u>



Storage Optimization



Factors affecting storage performance



Adapter type Number of virtual disks Virtual adapter queue depth

VMKernel admittance (Disk.SchedNumReqOutstanding)

Per path queue depth Adapter queue depth

Storage network (link speed, zoning, subnetting)

LUN queue depth Array SPs HBA target queues Number of disks (spindles)

LISA16

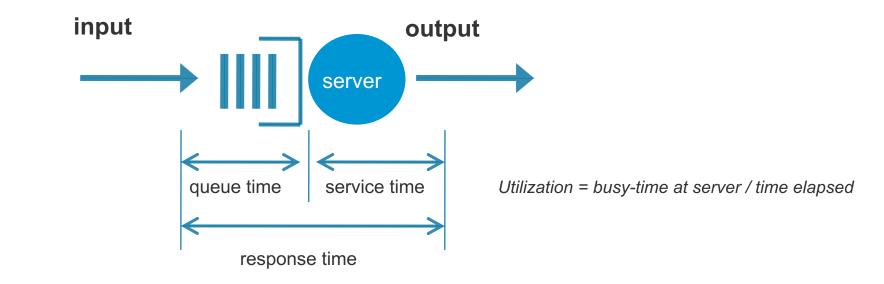
Nobody Likes Long Queues



Arriving Customers

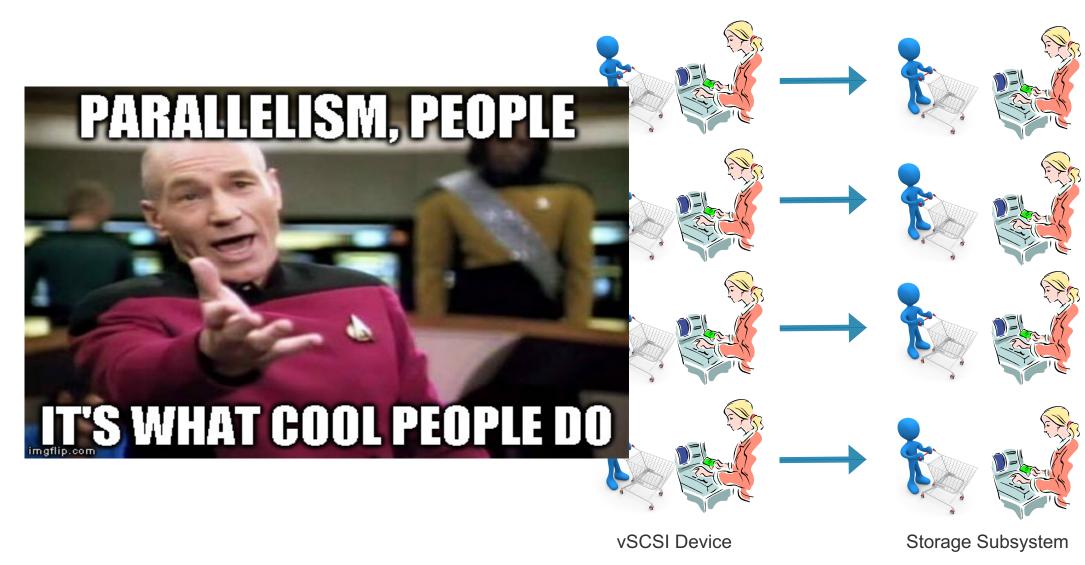
Queue

Checkout





Additional vSCSI controllers improves concurrency



LISA16

Optimize for Performance – Queue Depth

- vSCSI Adapter
 - Be aware of per device/adapter queue depth maximums (KB 1267)
 - Use multiple PVSCSI adapters
- VMKernel Admittance
 - VMKernel admittance policy affecting shared datastore (<u>KB 1268</u>), use dedicated datastores for DB and Log Volumes
 - VMKernel admittance changes dynamically when SIOC is enabled (may be used to control IOs for lower-tiered VMs)
- Physical HBAs
 - Follow vendor recommendation on max queue depth per LUN (<u>http://kb.vmware.com/kb/1267</u>)
 - Follow vendor recommendation on HBA execution throttle
 - Be aware settings are global if host is connected to multiple storage arrays
 - Ensure cards are installed in slots with enough bandwidth to support their expected throughput
 - Pick the right multi-pathing policy based on vendor storage array design (ask your storage vendor)



Increase PVSCSI Queue Depth

- Just increasing LUN, HBA queue depths is <u>NOT ENOUGH</u>
- PVSCSI <u>http://KB.vmware.com/kb/2053145</u>
- Increase PVSCSI Default Queue Depth (after consultation with array vendor)

• Linux:

- Add following line to /etc/modprobe.d/ or /etc/modprobe.conf file:
- options vmw_pvscsi cmd_per_lun=254 ring_pages=32
- OR, append these to the appropriate kernel boot arguments (grub.conf or grub.cfg)
 - vmw_pvscsi.cmd_per_lun=254
 - vmw_pvscsi.ring_pages=32
- Windows:
 - Key: HKLM\SYSTEM\CurrentControlSet\services\pvscsi\Parameters\Device
 - Value: DriverParameter | Value Data: "RequestRingPages=32,MaxQueueDepth=254"



Optimize for Performance – Storage Network

- Link Type/Speed
 - FC vs. iSCSI vs. NAS
 - · Latency suffers when bandwidth is saturated
- Zoning and Subnetting
 - Place hosts and storage on the same switch, minimize Inter-Switch Links
 - Use 1:1 initiator to target zoning or follow vendor recommendation
 - Enable jumbo frame for IP based storage (MTU needs to be set on all connected physical and virtual devices)
 - Make sure different iSCSI IP subnets cannot transmit traffic between them

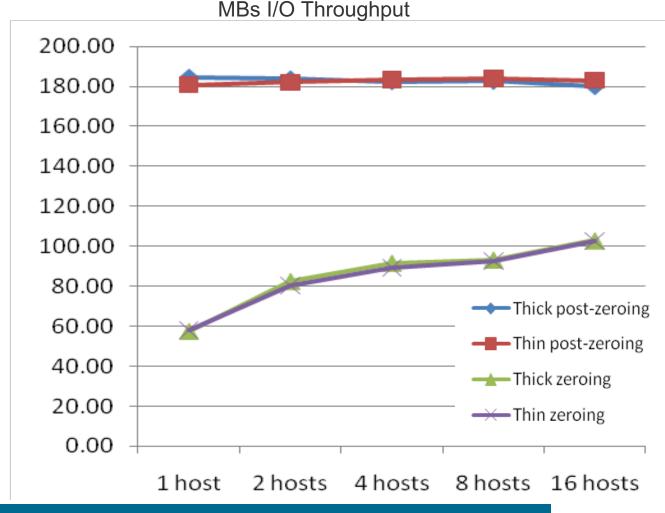


"Thick" vs "Thin"

usenix

SΔ

- Thin (Fully Inflated and Zeroed) Disk Performance = Thick Eager Zero Disk
- Performance impact due to zeroing, not result of allocation of new blocks
- To get maximum performance from the <u>start</u>, must use Thick <u>Eager Zero</u> Disks (think Business Critical Apps)
- Maximum Performance happens eventually, but when using lazy zeroing, zeroing needs to occur before you can get maximum performance



Choose Storage which supports VMware vStorage APIs for Array Integration (VAAI)

VMFS or RDM?

- Generally similar performance http://www.vmware.com/files/pdf/performance_char_vmfs_rdm.pdf
- vSphere 5.5 and later support up to 62TB VMDK files
- Disk size no longer a limitation of VMFS

VMFS	RDM
Better storage consolidation – multiple virtual disks/virtual machines per VMFS LUN. But still can assign one virtual machine per LUN	Enforces 1:1 mapping between virtual machine and LUN
Consolidating virtual machines in LUN – less likely to reach vSphere LUN Limit of 255	More likely to hit vSphere LUN limit of 255
Manage performance – combined IOPS of all virtual machines in LUN < IOPS rating of LUN	Not impacted by IOPS of other virtual machines

- When to use raw device mapping (RDM)
 - Required for shared-disk failover clustering
 - Required by storage vendor for SAN management tools such as backup and snapshots
- Otherwise use VMFS



Example Best Practices for VM Disk Layout (Microsoft SQL Server)

Characteristics:

- OS on shared DataStore/LUN
- 1 database; 4 equally-sized data files across 4 LUNs
- 1 TempDB; 4 (1/vCPU) equally-sized tempdb files across 4 LUNs
- Data, TempDB, and Log files spread across 3 PVSCSI adapters
 - Data and TempDB files share PVSCSI adapters
- Virtual Disks could be RDMs

Advantages:

usenix

LISA16

- Optimal performance; each Data, TempDB, and Log file has a dedicated VMDK/Data Store/LUN
- I/O spread evenly across PVSCSI adapters
- Log traffic does not contend with random Data/TempDB traffic

SQL Server OS Can be Mount Points under a drive as well. **C:**\ D:\ H:\ E:\ **I:**\ **F:**\ J:\ G:\ K:\ L:\ T:\ OS TmpFile2 DataFile3 LogFile1. DataFile1 TmpFile1 DataFile5 TmpFile3 DataFile7 TmpFile4 TmpLog1 NTFS Partition: ldf ldf .mdf .mdf .ndf .ndf .ndf .ndf .ndf .ndf 64K cluster size Can also be a shared LUN since TempDB is usually in Simple ESX Host Recovery Mode LSI1 PVSCSI1 PVSCSI2 PVSCSI3 VMDK4 VMDK5 VMDK6 OS VMDK VMDK1 VMDK2 VMDK3 VMDK6 VMDK5 VMDK6 VMDK5 Can be placed on a DataStore/LUN with other OS Data Store 2 Data Store 5 Data Store 1 Data Store 3 Data Store 4 Data Store 6 Data Store 5 Data Store 6 Data Store 5 Data Store 6 VMDKs y a y a ys gg LUN5 LUN6 LUN1 LUN2 LUN3 LUN4 LUN6 LUN5 LUN6 LUN5

Disadvantages:

- You can quickly run out of Windows driver letters!
- More complicated storage management

Realistic VM Disk Layout (Microsoft SQL Server)

Characteristics:

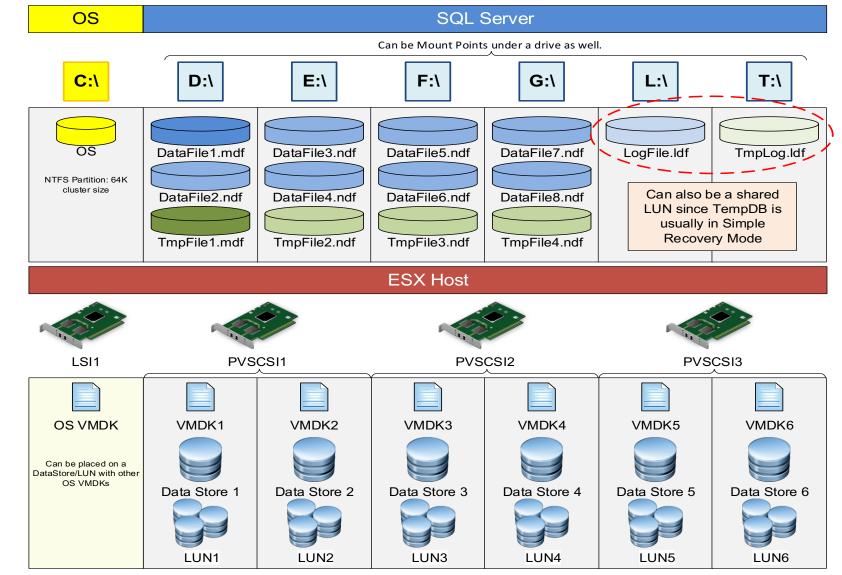
- OS on shared DataStore/LUN
- 1 database; 8 Equally-sized data files across 4 LUNs
- 1 TempDB; 4 files (1/vCPU) evenly distributed and mixed with data files to avoid "hot spots"
- Data, TempDB, and Log files spread across 3 PVSCSI adapters
- Virtual Disks could be RDMs

Advantages:

usenix

LISA16

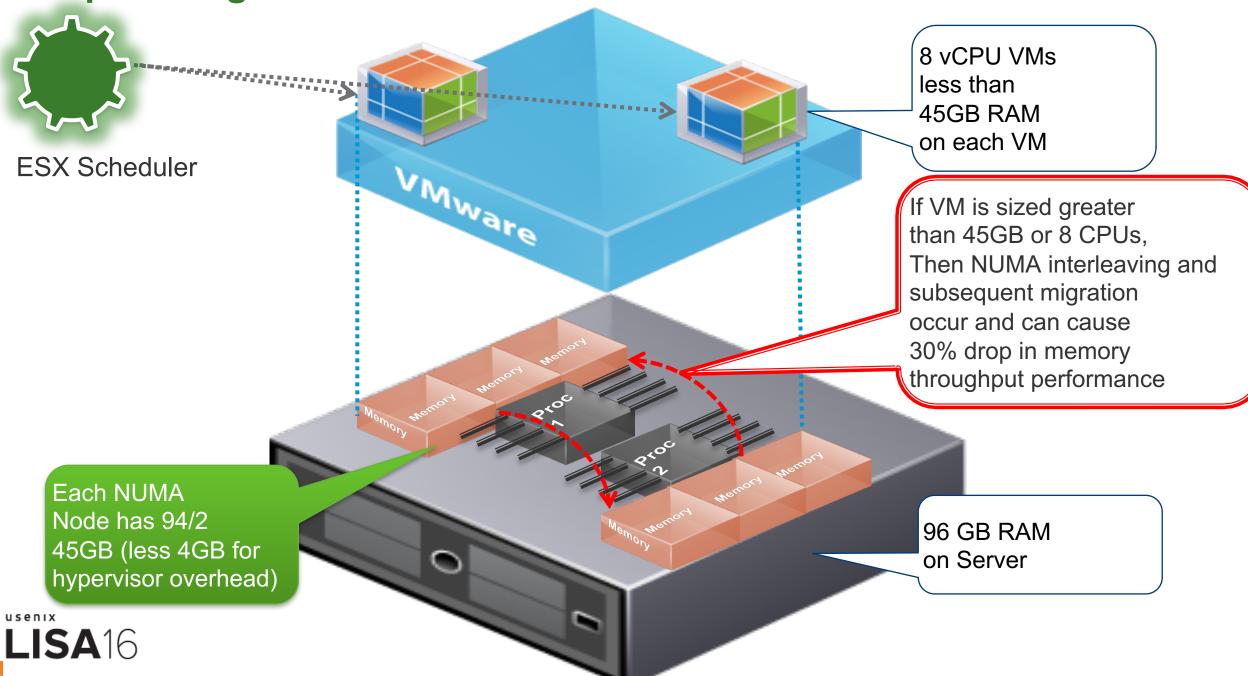
- Fewer drive letters used
- I/O spread evenly/TempDB hot spots avoided
- Log traffic does not contend with random Data/TempDB traffic



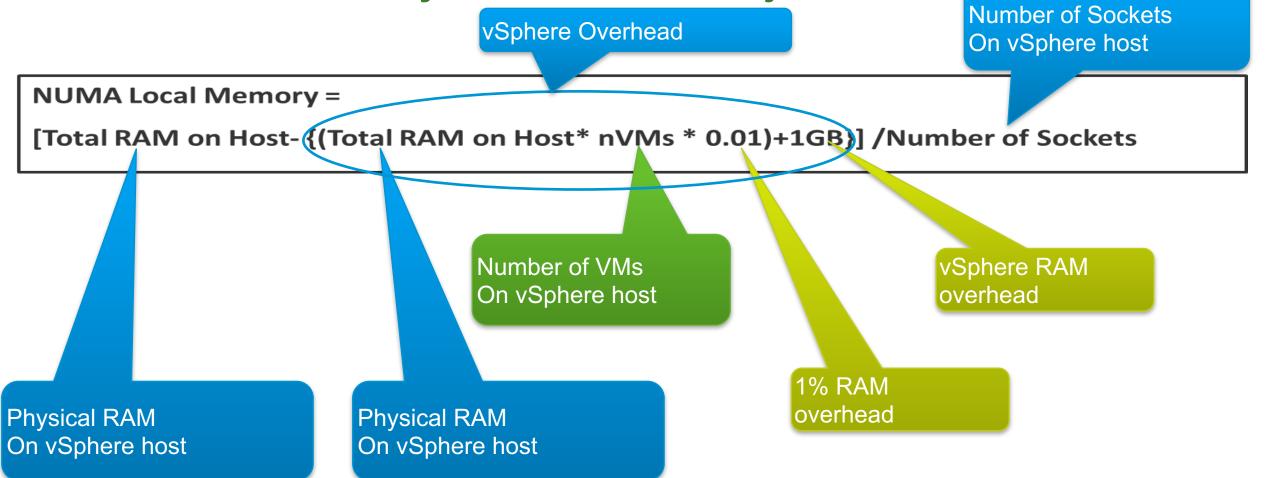
Lets talk about CPU, vCPUs and other Things



Optimizing Performance – Know Your NUMA



NUMA Local Memory with Overhead Adjustment





NUMA and vNUMA FAQ!

- Shall we Define NUMA Again? Nah.....
- Why VMware Recommends Enabling NUMA
 - Modern Operating Systems are NUMA-aware
 - Some applications are NUMA-aware (some are not)
 - vSphere Benefits from NUMA
- Use it, People
 - Enable Host-Level NUMA
 - Disable "Node Inter-leaving" in BIOS on HP Systems
 - Consult Hardware Vendor for SPECIFIC Configuration
- Virtual NUMA
 - Auto-enabled on vSphere for Any VM with 9 or more vCPUs
 - Want to use it on Smaller VMs?
 - Set <u>"numa.vcpu.min"</u> to # of vCPUs on the VM
- CPU Hot-Plug <u>DISABLES</u> Virtual NUMA
- vSphere 6.5 changes vNUMA config



vSphere 6.5 vCPU Allocation Guidance

Dharaing L Town Lower	vCPUs	VM Config	uration	Resulting vNUMA			
Physical Topology	Required	vSockets	vCores	Nodes Presented			
	1	1	1	1			
	2	1	2	1			
	3	1	3	1			
	4	1	4	1			
	5	1	5	1			
	6	1	6	1			
	7	1	7	1			
	8	1	8	1			
	9	1	9	1			
	10	1	10	1			
Hardware: Intel	11	1	11	1			
	12	1	12	1			
 2 Sockets 	13	1	13	1			
 2 pNUMA Nodes 	14	1	14	1			
	15	-	Sub-opt				
 14 Cores Per Socket 	16	2	8	2			
28 Logical Threads	17 18	2	Sub-opt				
 28 Logical Threads 	18	2					
	20	2	Sub-opt	2			
	20	2	Sub-opti				
	22	2	11	2			
	23	~ ~	Sub-opti				
	24	2	12	2			
	25	Sub-optimal					
	26	2	13	2			
	27		Sub-opti				
	28	2	14	2			



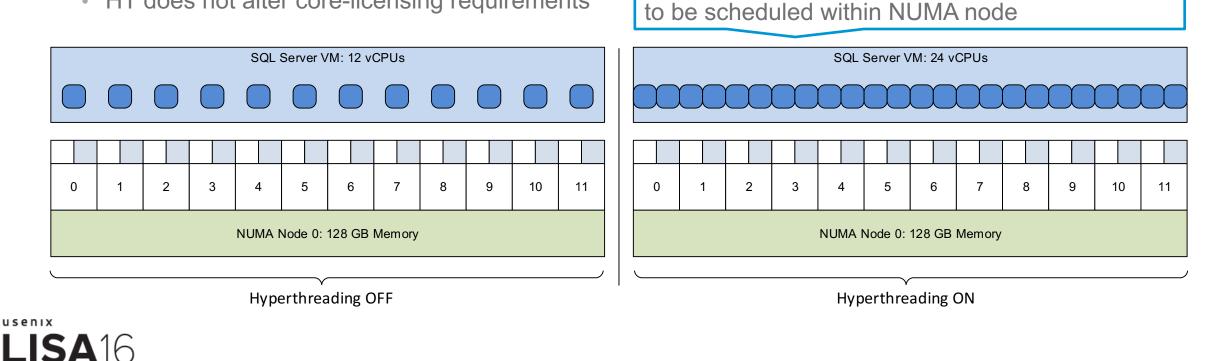
NUMA Best Practices

- Avoid Remote NUMA access
 - Size # of vCPUs to be <= the # of cores on a NUMA node (processor socket)
 - Where possible, align VMs with physical NUMA boundaries
 - For wide VMs, use a multiple or even divisor of NUMA boundaries
 - http://www.vmware.com/files/pdf/techpaper/VMware-vSphere-CPU-Sched-Perf.pdf
- Hyper-threading
 - Initial conservative sizing: set vCPUs equal to # of physical cores
 - HT benefit around 30-50%, < for CPU intensive batch jobs (based on OLTP workload tests)
- Allocate vCPUs by socket count
 - Default "Cores Per Socket" is set to "1"
 - Applicable to vSphere versions prior to 6.5. Not as relevant in 6.5
- ESXTOP to monitor NUMA performance in vSphere
 - Coreinfo.exe to see NUMA topology in Windows Guest
- vMotioning VMs between hosts with dissimilar NUMA topology leads to performance issues



Non-Wide VM Sizing Example (VM fits within NUMA Node)

- 1 vCPU per core with hyper-threading OFF
 - Must license each core for SQL Server
- 1 vCPU per thread with hyper-threading ON
 - 10%-25% gain in processing power •
 - Same licensing consideration
 - HT does not alter core-licensing requirements



"numa.vcpu.preferHT" to true to force 24-way VM

Wide VM Sizing Example (VM crosses NUMA Node)

- Extends NUMA awareness to the guest OS
- Enabled through multicore UI
 - On by default for 9+ vCPU multicore VM
 - Existing VMs are not affected through upgrade
 - For smaller VMs, enable by setting numa.vcpu.min=4
- Do NOT turn on CPU Hot-Add
- For wide virtual machines, confirm feature is on for best performance

🔗 sql_test_srv1 - Virtual Machine	Properties							
Hardware Options Resources Profiles VServices Virtual Machine Version: 8								
Show All Devices	Add Remove	Number of virtual sockets:	4 💌					
Hardware	Summary	Number of cores per socket:	4 🔻					
Memory	8192 MB							
🔲 CPUs (edited)	16	Total number of cores:	16					
	(n							

Virtu	Virtual NUMA Node 0 SQL Server VM: 24 vCPUs Virtual NUMA Note 0 Vi											Node 1											
																							\bigcirc
0	1	2	3	4	5	6	7	8	9	10	11	0	1	2	3	4	5	6	7	8	9	10	11
	NUMA Node 0: 128 GB Memory NUMA Node 1: 128 GB Memory																						
												·											



Designing for Performance

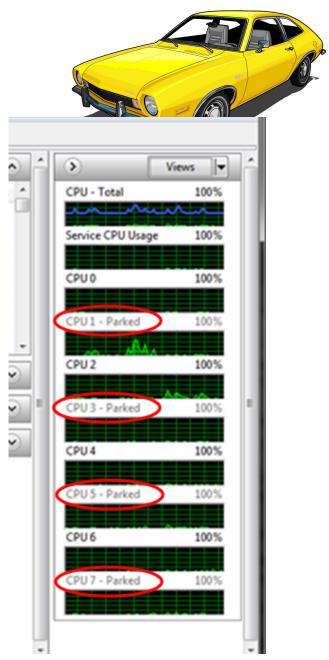
- The VM itself matters In-guest optimization
 - Windows CPU Core Parking = BAD
 - Set Power to "High Performance" to avoid core parking
 - Relevant IF ESXi Host Power Setting NOT "High Performance"
 - Windows Receive Side Scaling settings impact CPU utilization
 - Must be enabled at NIC and Windows Kernel level
 - Use "netsh int tcp show global" to verify
 - More on this later
- Application-level tuning
 - Follow vendor's recommendation
 - Virtualization does not change the consideration





Why Your Windows App Server Lamborghini Runs Like a Pinto

File Monitor Help	lav. I						
		Network	_		_		A Views V
Processes	0% CPU Us	age		59% Maximum	Freque	ngy 🔿	Views V
Image	PID	Descrip	Status	Threads	CPU	Averag ^	CPU - Total 100%
perfmon.exe	5496	Resour	Runni	19	0	0.99	
PaintDotNet.exe	5708	Paint.N	Runni	30	0	0.96	Service CPU Usage 100%
dwm.exe	2728	Deskto	Runni	5	0	0.23	Service CPO Osage 100%
System Interrupts	-	Deferr	Runni	-	0	0.10	
chrome.exe	4296	Googi	Runni	6	0	0.05	CPU 0 100%
System	4	NT Ker	Runni	149	0	0.05	
explorer.exe	3080	Windo		35	0	0.04	
csrss.exe	576	Client		12	0	0.04	CPU1 - Parked 100%
MSIAfterburner.exe	3588	MSIAft		6	0	0.04	- AA.
for the first of the second			O	,	-		CPU 2 100%
Services	0% CPU Us	age				•	
Associated Handles			Sea	rch Handles		P 47 👻	E CPU 3 - Parked 100%
Associated Modules						$\overline{\mathbf{v}}$	CPU 4 100%
							100%
							CPU 5 - Parked 100%
							CPU 6 100%
							CP06 100%
							CPU 7 - Parked 100%



LISA16

Memory Optimization



Memory Reservations

- Guarantees allocated memory for a VM
- The VM is only allowed to power on if the CPU and memory reservation is available (strict admission)
- If Allocated RAM = Reserved RAM, you avoid swapping
- Do NOT set memory limits for Mission-Critical VMs
- If using Resource Pools, Put Lower-tiered VMs in Resource Pools
- Some Applications Don't Support "Memory Hot-add"
 - E.g. Microsoft Exchange Server CANNOT use Hot-added RAM
- Don't use it on ESXi versions lower than 6.0
- Virtual: Physical memory allocation ratio <u>should</u> not exceed 2:1
- Remember NUMA? It's not just about CPU
 - Fetching remote memory is VERY expensive
 - Use "numa.vcpu.maxPerVirtualNode" to control memory locality

What about Dynamic Memory?

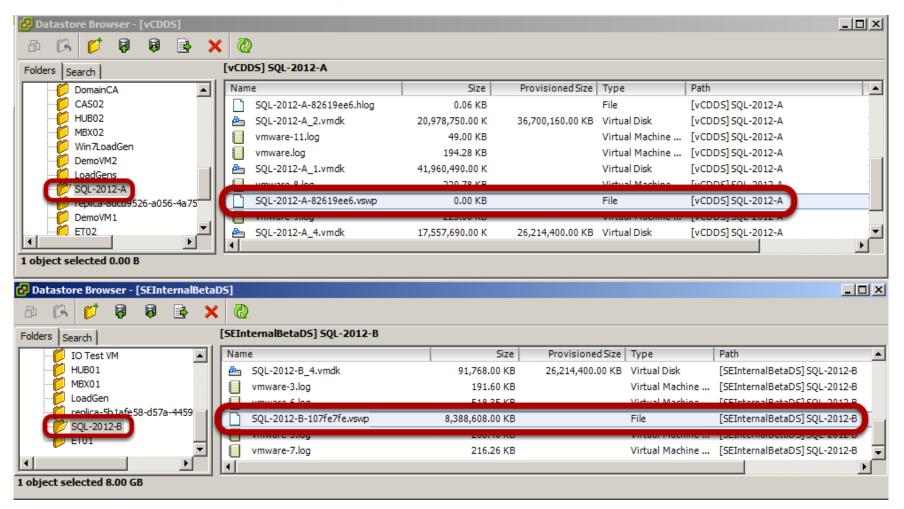
Not Supported by Most Microsoft's Critical Applications

Not a feature of VMware vSphere



Memory Reservations and Swapping on vSphere

• Setting a reservation creates zero (or near-zero) swap file

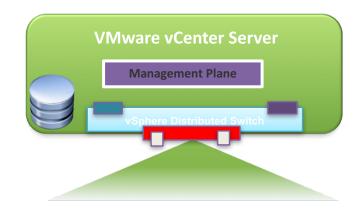


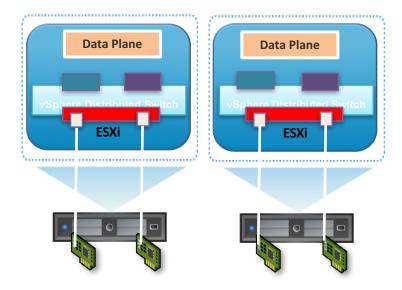
LISA16

Network Optimization



vSphere Distributed Switch (VDS) Overview



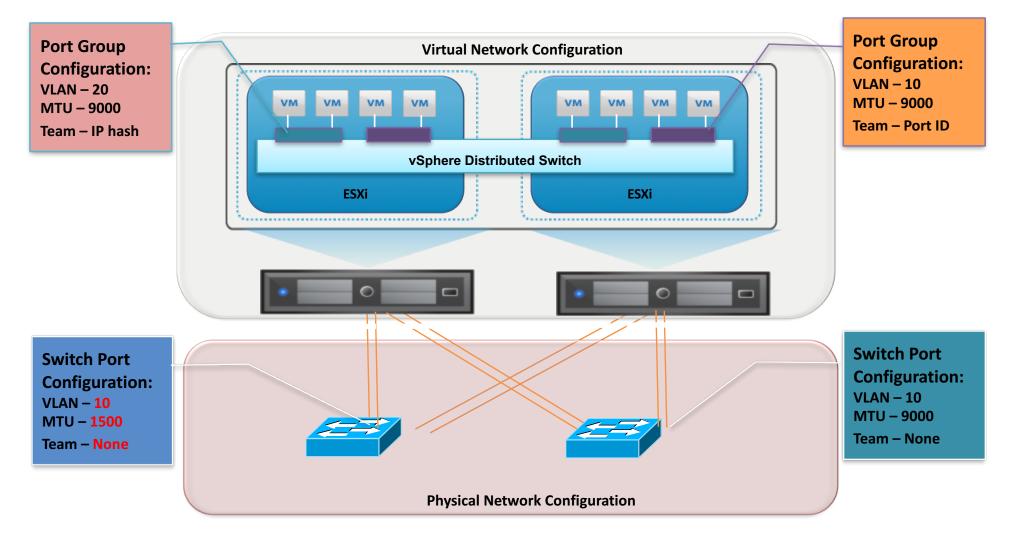


usenix

LISA16

- Unified network virtualization management
 - Independent of physical fabric
- vMotion aware : Statistics and policies follow the VM
- vCenter management plane independent of data plane
- Advanced Traffic Management features
 - Load Based Teaming (LBT)
 - Network IO Control (NIOC)
- Monitoring and Troubleshooting features
 - NetFlow
 - Port Mirroring

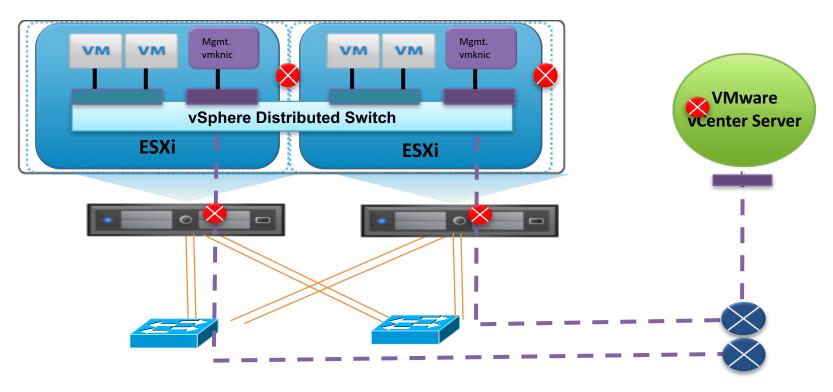
Common Network Misconfiguration



The network health check feature sends a probe packet every 2 mins



Misconfiguration of Management Network



Two different updates that triggers rollback

- Host level Rollback gets triggered when there is change in the host networking configurations such as: Physical NIC speed change, Change in MTU configuration, Change in IP settings etc..
- VDS level rollback can happen after the user updates some VDS related objects such as port group or dvports.



Network Best Practices

- Allocate separate NICs for different traffic type
 - Can be connected to same uplink/physical NIC on 10GB network
- vSphere versions 5.0 and newer support multi-NIC, concurrent vMotion operations
- Use NIC load-based teaming (route based on physical NIC load)
 - For redundancy, load balancing, and improved vMotion speeds
- Have minimum 4 NICs per host to ensure performance and redundancy of network
- Recommend the use of NICs that support:
 - Checksum offload, TCP segmentation offload (TSO)
 - Jumbo frames (JF), Large receive offload (LRO)
 - Ability to handle high-memory DMA (i.e. 64-bit DMA addresses)
 - Ability to handle multiple Scatter Gather elements per Tx frame
 - NICs should support offload of encapsulated packets (with VXLAN)
- ALWAYS Check and Update Physical NIC Drivers

• Keep VMware Tools Up-to-Date - ALWAYS LISA16



Network Best Practices (continued)

- Use Virtual Distributed Switches for cross-ESX network convenience
- Optimize IP-based storage (iSCSI and NFS)
 - Enable Jumbo Frames
 - Use dedicated VLAN for ESXi host's vmknic & iSCSI/NFS server to minimize network interference from other packet sources
 - Exclude in-Guest iSCSI NICs from WSFC use
 - Be mindful of converged networks; storage load can affect network and vice versa as they use the same physical hardware; ensure no bottlenecks in the network between the source and destination
- Use VMXNET3 Para-virtualized adapter drivers to increase performance
 - NEVER use any other vNIC type, unless for legacy OSes and applications
 - Reduces overhead versus vlance or E1000 emulation
 - Must have VMware Tools to enable VMXNET3
- Tune Guest OS network buffers, maximum ports



Network Best Practices (continued)

- VMXNET3 can bite but only if you let it
 - ALWAYS keep VMware Tools up-to-date
 - ALWAYS keep ESXi Host Firmware and Drivers up-to-date
 - Choose your physical NICs wisely
 - Windows Issues with VMXNET3
 - Older Windows versions



- VMXNET3 template issues in Windows 2008 R2 <u>kb.vmware.com\kb\1020078</u>
- Hotfix for Windows 2008 R2 VMs <u>http://support.microsoft.com/kb/2344941</u>
- Hotfix for Windows 2008 R2 SP1 VMs <u>http://support.microsoft.com/kb/2550978</u>
- Disable interrupt coalescing at vNIC level
 - ONLY if ALL other options fail to remedy network-related performance Issue



:(

Your PC ran into a problem and needs to restart. We're just collecting some error info, and then we'll restart for you.

If you'd like to know more, you can search online later for this error: ALWAYS_LOOK_ON_THE_BRIGHT_SIDE_OF_LIFE



What is NSX?

- Network Overlay
- Logical networks
- Logical Routing
- Logical Firewal
- Logical Load Balancing
- Addition a Netw services (NAT, more)
- Programmatically Controlled

usenix

LISA16

production src,dest,port,protocol database tier allow<=application tier> customer Data allow<appid=3456> pci data allow<appid=6789> quarantine cvss=2

VM

Network Hypervisor

General Purpose IP Hardware

L3

Logical Switch (L2)

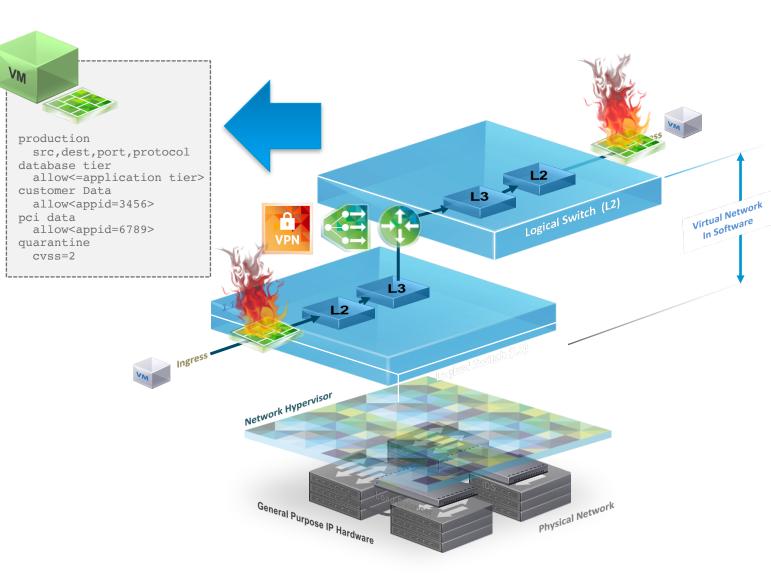
Physical Network

Virtual Network

In Software

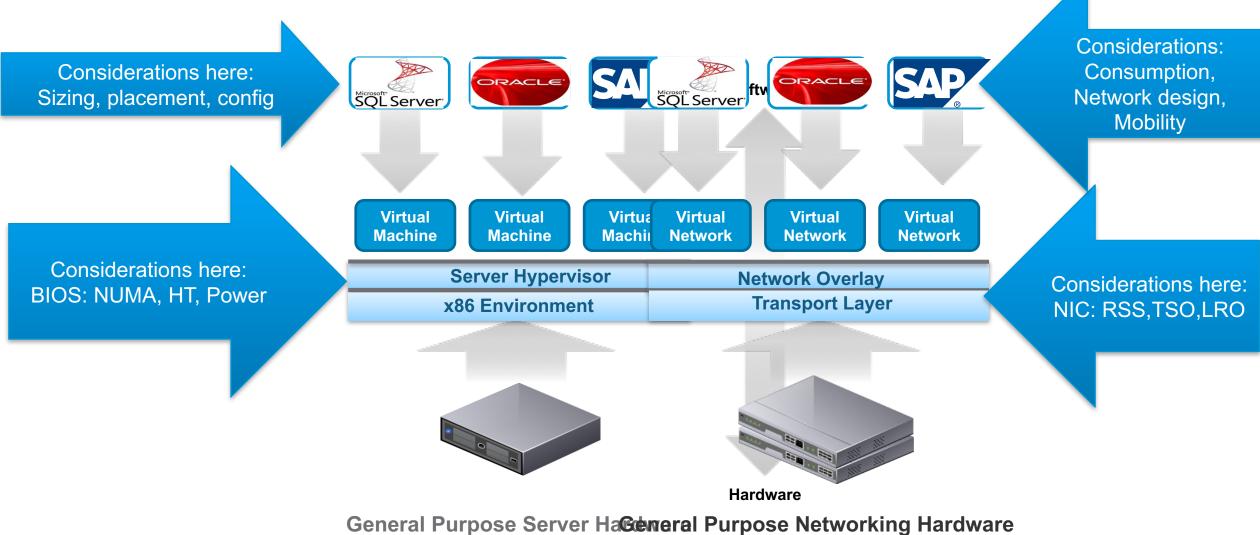
What is NSX?

- Network Overlay
- Logical networks
- Logical Routing
- Logical Firewall
- Logical Load Balancing
- Additional Networking services (NAT, VPN, more)
- Programmatically Controlled





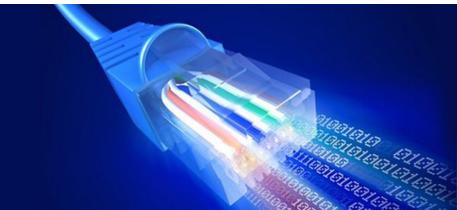
What do app owners care about?



LISA16

Performance Considerations

- All you need is IP connectivity between ESXi hosts
- The physical NIC and the NIC driver should support:
 - TSO TCP Segmentation Offload = NIC divides larger data chunks into TCP segments
 - VXLAN offload NIC encapsulates VXLAN instead of ESXi
 - RSS Receive side scaling, allows the NIC to distribute received traffic to multiple CPU
 - LRO (Large Receive Offload) NIC reassembles incoming network packets





App owners say...

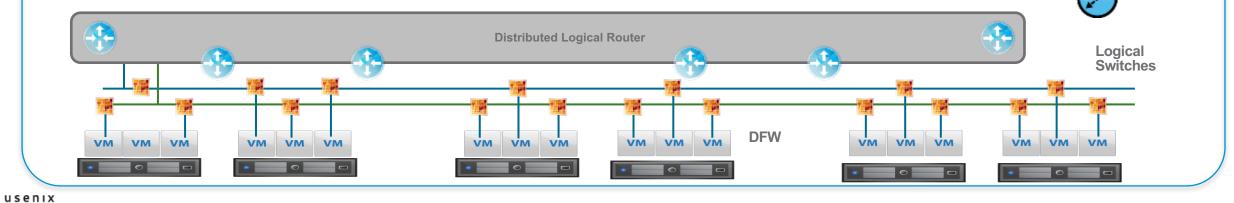
LISA16

- So if the "Network hypervisor" fail does my app fail?
- What about NSX components dependencies?



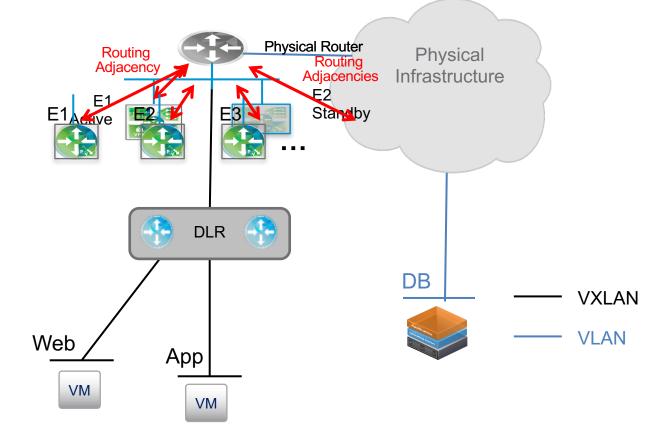
vCenter & NSX Manager A	Management plane: UI, API access Not in the data path
Controller Cluster	Control plane: Decouples virtual networks form physical topology Not in Data Path Highly Available

Data plane: Logical switches, Distributed Routers, Distributed Firewall, Edge devices



Connecting to the physical network

- Typical use case: 3-tier application, Web/App/DB, with non-virtualized DB tier.
- Option 21-- Routesing and Edge device on AA Anotacle:



usenix

LISA16

Altonos SNO TS Addrew for state fulse hvarses blat A Ti, eL @dg. (P\$Nuch as NAT, LB, VPN.

LimiteB icath satilightep provideG britone (singtern N 160) de Firewall can be service by the FailoDerVtakes a few seconds

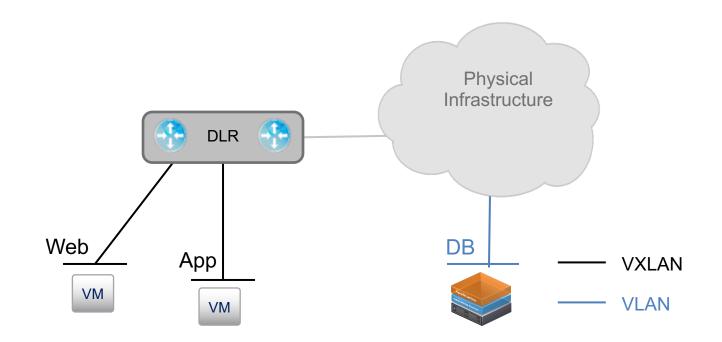
High throughput of upto 80Gbit

Provides highest redundancy with multipath

NSX Edge

Connecting to the physical network

- Typical use case: 3-tier application, Web/App/DB, with non-virtualized DB tier.
- Option 3 Bridging L2 network using software or hardware GW:



Straight from the ESXi kernel to the VLAN backed network

Lowest Latency L2 adjacency between the tiers Design complexity Redundancy limitations



Designing for Availability



vSphere Native Availability Features

vSphere vMotion

- Can reduce virtual machine planned downtime
- Relocates VMs without end-user interruption
 - Behavior COMPLETELY Configurable
- Enables Admin to perform on-demand host maintenance without service interruption

• vSphere DRS

usenix

LISA16

- Monitors state of virtual machine resource usage
- Can automatically and intelligently locate virtual machine
- Can create a dynamically balanced Exchange Server deployment
- Uses vMotion. Behavior COMPLETELY Configurable
- vSphere High Availability (HA)
 - HA Evaluates DRS Rules **<u>BEFORE</u>** Recovery Just a checkbox operati
 - * Now DEFAULT BEHAVIOR is vSphere 6.5
 - Does not require Vendor-specific clustering solutions
 - **NOT** a replacement for app-specific native HA features
 - **<u>COMPLEMENTS</u>** and **<u>ENHANCES</u>** app-specific HA features
 - Automatically restarts failed virtual machine in minutes

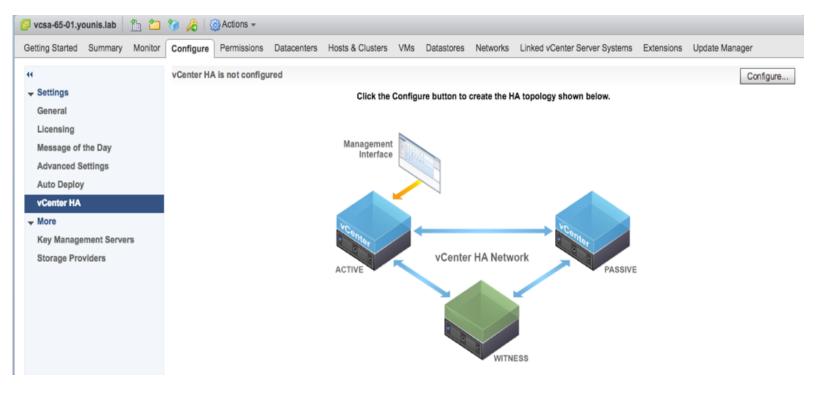
vSphere HA Rule Settings

vSphere HA can enforce VM/Host rules when restarting virtual machines.

VM anti-affinity rules	Ignore rules	
VM to Host affinity rules	Ignore rules	

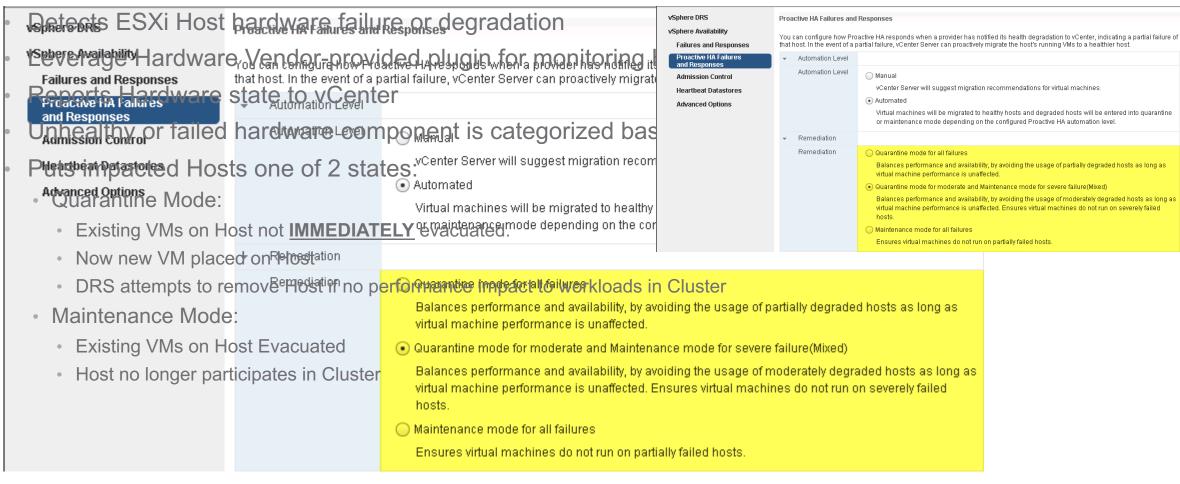


- vCenter High Availability
 - vCenter Server Appliance <u>ONLY</u>
 - Active, Passive, and Witness nodes Exact Clones of existing vCenter Server.
 - Protects vCenter against Host, Appliance or Component Failures
 - 5-minute RTO at release



LISA16

Proactive High Availability





Continuous VM Availability

For when VMs MUST be up, even at the expense of <u>PERFORMANCE</u>

vSphere DRS	Admission Control					
vSphere Availability Failures and Responses Proactive HA Failures	Admission control is a policy used by vSphere HA to ensure failover capacity within a cluster. Increasing the value of host failures cluster tolerates will increase the availability constraints and capacity reserved.					
Admission Control	Host failures cluster tolerates	1 🚔 Maximum is one less than number of hosts in cluster.				
Heartbeat Datastores	Define host failover capacity by	Cluster resource percentage				
Advanced Options		Override calculated failover capacity. CPU 33 * % Memory 33 * %				
	Performance degradation VMs tolerate	 Percentage of performance degradation the VMs in the cluster are allowed to tolerate during a failure. 0% - Raises a warning if there is insufficient failover capacity to guarantee the same performance after VMs restart. 100% - Warning is disabled. 				



vSphere DRS Rules

Rules now includes "VM Dependencies"

•• → S

► Vi - Co

Allows VMs to be recovered in order of PRIORITIES

	VM/Host Ru	les				
ervices	Add	Edit Delete				
phere DRS	Name		Туре			
phere Availability		ate VMs	Separate Virtual Machines			
rtual SAN		riority-1-Rule	Run VMs based on group dependent	×		
onfiguration		riority-2-Rule	Run VMs based on group dependent			
eneral						
censing	🚯 TSALA	AB-65-Cluster - Edit VM/Host Rule	(?)			
Mware EVC	Name:	Boot-Priority-1-Rule				
M/Host Groups		Enable rule.				
M/Host Rules	Type:	Virtual Machines to Virtual Machines	•			
M Overrides	Description:					
ost Options						
ofiles		achines in the VM group Boot-Priority-1-VM a in the VM group Boot-Priority-2-VMs will k				
) Filters		pendency restart condition has been met				
				uster deper		
	First resta	rt VMs in VM group:				
	Boot-Pri	ority-1-VMs	•	Add		
	Than ract	art VMs in VM group:		ot-Priority-2-V VMW-SQL		
		prity-2-VMs	•	VMW-SQL		
	BOOLET III	7my 2 mmo	·]	41010-3GL		



Predictive DRS

usenix

Integrated with VMware's vRealize Operations Monitoring Capabilities

vSphere DRS	Predictive DRS	×
vSphere Availability		Enable Predictive DRS
Failures and Responses		In addition to realtime metrics, DRS will respond to forecasted metrics provided by
Proactive HA Failures and Responses		vRealize Operations server. You must also configure Predictive DRS in a version of vRealize Operations that supports this feature.
Admission Control		

- Network-Aware DRS
 - Considers Host's Network Bandwith Utilization for VM Placement
 - Does NOT Evacuate VMs Based on Utilization

Simplified Advanced DRS Configuration Tasks

Now just <u>Checkbox</u> options

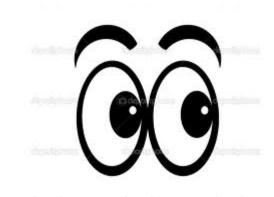
vSphere DRS	☑ Turn ON ∨Sphere DRS			
vSphere Availability Failures and Responses	DRS Automation Fully Automated The second sec			
Proactive HA Failures and Responses Admission Control	VM Distribution	For availability, distribute a more even number of virtual machines across hosts.		
Heartbeat Datastores Advanced Options	Memory Metric for Load Balancing	Load balance based on consumed memory of virtual machines rather than active memory. This setting is only recommended for clusters where host memory is not over-committed.		
	CPU Over-Commitment	Control CPU over-commitment in the cluster Over-commitment ratio (% of cluster capacity): 0 🐳 Min: 0 Max: 500		

Combining Windows Applications HA with vSphere HA Features – The Caveats



Are You Going to Cluster THAT?

- Do you <u>NEED</u> App-level Clustering?
 - Purely business and administrative decision
 - Virtualization does not preclude you from doing so
- Share-nothing Application Clustering?
 - No "Special" requirements on vSphere
- Shared-Disk Application Clustering (e.g. FCI / MSCS)
 - You MUST use Raw Device Mapping (RDM) Disks Type for Shared Disks
 - MUST be connected to vSCSI controllers in PHYSICAL Mode Bus Sharing
 - Wonder why it's called "Physical Mode RDM", eh?
 - In Pre-vSphere 6.0, FCI/MSCS nodes CANNOT be vMotioned. Period
 - In vSphere 6.0 and above, you have vMotions capabilities under following conditions
 - Clustered VMs are at Hardware Version > 10
 - vMotion VMKernel Portgroup Connected to 10GB Network





vMotioning Clustered Windows Nodes – Avoid the Pitfall

- Clustered Windows Applications Use Windows Server Failover Clustering (WSFC)
 - WSFC has a Default 5 Seconds Heartbeat Timeout Threshold
 - vMotion Operations MAY Exceed 5 Seconds (During VM Quiescing)
 - Leading to Unintended and Disruptive Clustered Resource Failover Events
- SOLUTION
 - Use MULTIPLE vMotion Portgroups, where possible
 - Enable jumbo frames on all vmkernel ports, IF PHYSICAL Network Supports it
 - If jumbo frames is not supported, consider modifying default WSFC behaviors:
 - (get-cluster).SameSubnetThreshold = 10
 - (get-cluster).CrossSubnetThreshold = 20
 - (get-cluster).RouteHistoryLength = 40
 - NOTES:

usenix

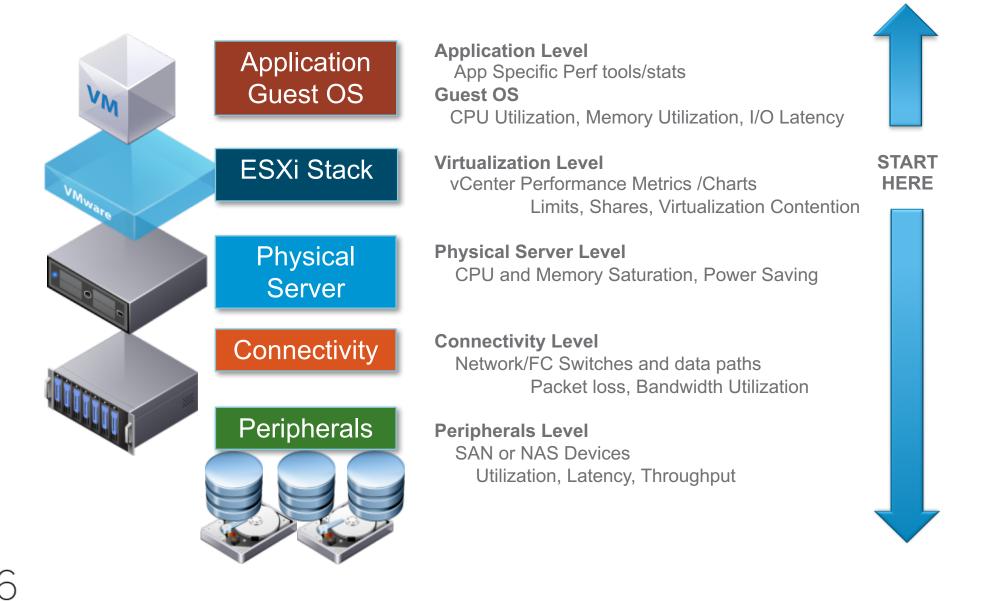
LISA16

- You may need to <u>"Import-Module FailoverClusters</u>" first
- Behavior NOT Unique to VMware or Virtualization
 - If Your Backup Software Quiesces Exchange, You Experience Symptom
 - See Microsoft's "Tuning Failover Cluster Network Thresholds" <u>http://bit.ly/1nJRPs3</u>

Monitoring and Identifying Performance Bottlenecks



Performance Needs Monitoring at Every Level

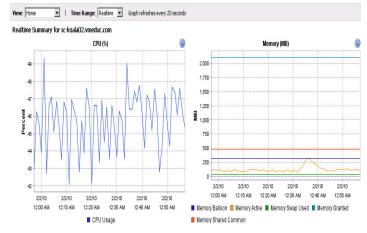


LISA16

Host Level Monitoring

- VMware vSphere Client[™]
 - GUI interface, primary tool for observing performance and configuration data for one or more vSphere hosts
 - Does not require high levels of privilege to access the data
- Resxtop/ESXTop
 - Gives access to detailed performance data of a single vSphere host
 - Provides fast access to a large number of performance metrics
 - Runs in interactive, batch, or replay mode
 - ESXTop Cheat Sheet <u>http://www.running-system.com/vsphere-6-esxtop-quick-overview-for-troubleshooting/</u>

Overview Advanced



00		root	@bk09-	h380-11	:∼ — ssh	$-118 \times$	21				
8:10:46am up 4 days	4:27, 143 wo	rlds; CPU	load ave	erage: 0.8	1, 0.01,	0.01					ŕ
ADAPTR PATH	NPTH	CMDS/s	READS/s	WRITES/s	MBREAD/s	MBWRTN/s	DAVG/cmd	KAVG/cmd	GAVG/cmd	QAVG/cmd	
vmhba0 -	1	0.59	0.00	0.59	0.00	0.00	0.08	0.01	0.09	0.00	
vmhbal -	6	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
vmhba2 -	6	1741.29	1740.29	0.99	54.38	0.01	17.31	1.08	18.39	0.00	
vmhba3 -	1	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
vmhba32 -	8	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
vmhba33 -	8	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
vmhba34 -	0	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
vmhba35 -	0	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
vmhba36 -	8	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	



Key Metrics to Monitor for vSphere

Resource	Metric	Host / VM	Description
	%USED	Both	CPU used over the collection interval (%)
CPU	%RDY	VM	CPU time spent in ready state
	%SYS	Both	Percentage of time spent in the ESX Server VMKernel
Momony	Swapin, Swapout	Both	Memory ESX host swaps in/out from/to disk (per VM, or cumulative over host)
Memory	MCTLSZ (MB)	Both	Amount of memory reclaimed from resource pool by way of ballooning
	READs/s, WRITEs/s	Both	Reads and Writes issued in the collection interval
Dist	DAVG/cmd	Both	Average latency (ms) of the device (LUN)
Disk	KAVG/cmd	Both	Average latency (ms) in the VMkernel, also known as "queuing time"
	GAVG/cmd	Both	Average latency (ms) in the guest. GAVG = DAVG + KAVG
	MbRX/s, MbTX/s	Both	Amount of data transmitted per second
Network	PKTRX/s, PKTTX/s	Both	Packets transmitted per second
	%DRPRX, %DRPTX	Both	Drop packets per second



Key Indicators

CPU

Ready (%RDY)

- % time a vCPU was ready to be scheduled on a physical processor but couldn't' t due to processor contention
- Investigation Threshold: 10% per vCPU

Co-Stop (%CSTP)

- % time a vCPU in an SMP virtual machine is "stopped" from executing, so that another vCPU in the same virtual machine could be run to "catch-up" and make sure the skew between the two virtual processors doesn't grow too large
- Investigation Threshold: 3%
- Max Limited (%MLMTD)
 - % time VM was ready to run but wasn't scheduled because it violated the CPU Limit set ; added to %RDY time
 - Virtual machine level processor queue length



Key Performance Indicators

Memory

Balloon driver size (MCTLSZ)

the total amount of guest physical memory reclaimed by the balloon driver

Investigation Threshold: 1

Swapping (SWCUR)

the current amount of guest physical memory that is swapped out to the ESX kernel VM swap file.

Investigation Threshold: 1

Swap Reads/sec (SWR/s)

the rate at which machine memory is swapped in from disk. **Investigation Threshold: 1**

Swap Writes/sec (SWW/s)

The rate at which machine memory is swapped out to disk. **Investigation Threshold: 1**

Network

Transmit Dropped Packets (%DRPTX)

The percentage of transmit packets dropped. Investigation Threshold: 1

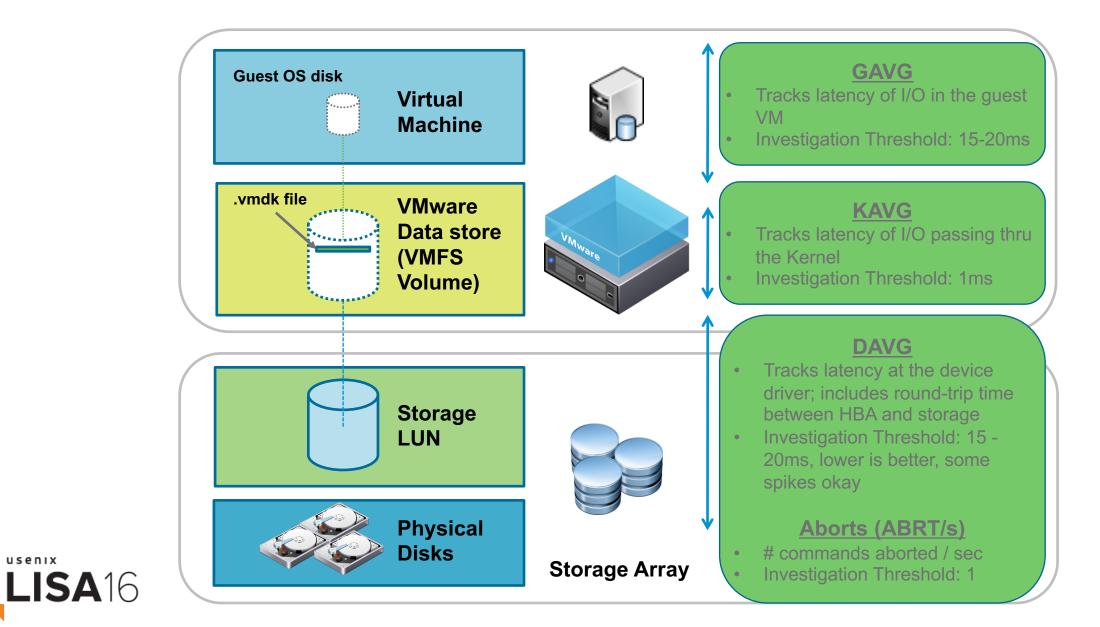
Receive Dropped Packets (%DRPRX)

The percentage of receive packets dropped. Investigation Threshold: 1



Logical Storage Layers: from Physical Disks to vmdks

usenıx



Key Indicators

Storage

- Kernel Latency Average (KAVG)
 - This counter tracks the latencies of IO passing thru the Kernel
 - Investigation Threshold: 1ms
- Device Latency Average (DAVG)
 - This is the latency seen at the device driver level. It includes the roundtrip time between the HBA and the storage.
 - Investigation Threshold: 15-20ms, lower is better, some spikes okay
- Aborts (ABRT/s)
 - The number of commands aborted per second.
 - Investigation Threshold: 1
- Size Storage Arrays appropriately for Total VM usage
 - > 15-20ms Disk Latency could be a performance problem
 - > 1ms Kernel Latency could be a performance problem or a undersized ESX device queue



Storage Performance Troubleshooting Tools



Storage Profiling Tips and Tricks

- Common IO Profiles (database, web, etc): http://blogs.msdn.com/b/tvoellm/archive/2009/05/07/useful-io-profiles-for-simulating-various-workloads.aspx
- Make Sure to Check / Try:
 - Load balancing / multi-pathing
 - Queue depth & outstanding I/Os
 - pvSCSI Device Driver
- Look out for:
 - I/O contention
 - Disk Shares
 - SIOC & SDRS
 - IOP Limits





vscsiStats – DEEP Storage Diagnostics

- vscsiStats characterizes IO for each virtual disk
 - Allows us to separate out each different type of workload into its own container and observe trends
- Histograms only collected if enabled; no overhead otherwise
- Metrics

usenıx

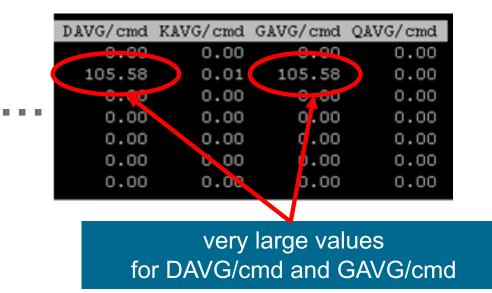
SA16

- I/O Size
- Seek Distance
- Outstanding I/Os
- I/O Interarrival Times
- Latency

	Virtual Machine isk C Virtual SCSI HB LSI Logic or BUS L ESX Serve	Disk D A, ogic	
LUN 1	LUN 2		
	/1:1- /		- U ×
<mark>≰ root@bk09-h380-11:/usr</mark> [root@bk09-h380-11 b		s -w 4968 -i 8208 -p seekDistance	
Histogram: distance rtual machine worldG { min : -16771839 max : 807168 mean : -36 count : 429263 { 1 0 0 0 429261 0 0	<pre>compID : 4968, v</pre>	-128) -64) -32) -1) 0) 1) 32)	
	(<= (>	64) 128)	-

Monitoring Disk Performance with esxtop

ADAPTR	CMDS/s	READS/s	WRITES/s	
vmhba0	0.00	0.00	0.00	
vmhba1	9.16	7.44	1.72	
vmhba2	0.00	0.00	0.00	
vmhba3	0.00	0.00	0.00	
vmhba32	0.00	0.00	0.00	
vmhba33	0.00	0.00	0.00	
vmhba34	0.00	0.00	0.00	



- Rule of thumb
 - GAVG/cmd > 20ms = high latency!
- What does this mean?
 - When command reaches device, latency is high
 - Latency as seen by the guest is high
 - Low KAVG/cmd means command is not queuing in VMkernel

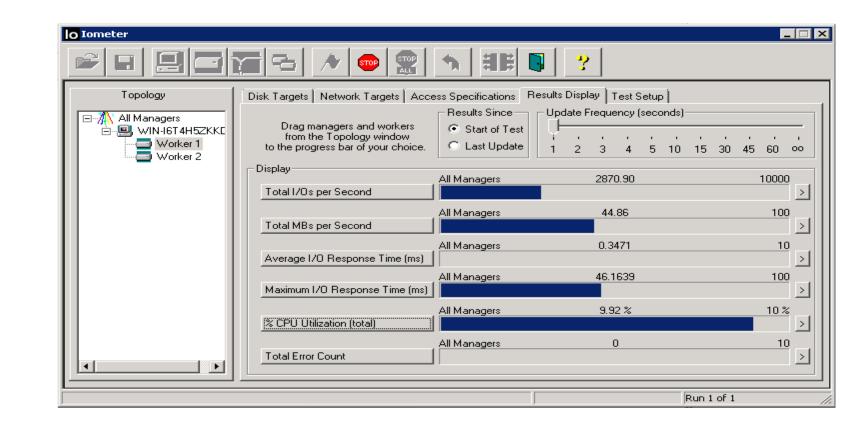


lometer

An I/O subsystem measurement and characterization tool for single and clustered systems.

Supports Windows and Linux

- Windows and Linux
- Free (Open Source)
- Single or Multi-server capable
- Multi-threaded
- Metrics Collected
 - Total I/Os per Sec.
 - Throughput (MB)
 - CPU Utilization
 - Latency (avg. & max)





DiskSpd Utility: A Robust Storage Testing Tool (SQLIO)

https://gallery.technet.microsoft.com/DiskSpd-a-robust-storage-6cd2f223 http://hfxte.ch/diskspd

- Windows-based feature-rich synthetic storage testing and validation tool
- Replaces SQLIO and effective for baselining storage for MS SQL Server workloads
- Fine-grained IO workload characteristics definition
- Configurable runtime and output options

usenix

LISA16

 Intelligent and easy-to-understand tabular summary in text-based output

thread	bytes	I/Os	MB/s	I/O per s	AvgLat	LatStdDev	file
	9101312	1111	4.34	555.48	7.193	5.687	s:\diskspd\io.dat (500
1	9256960	1130	4.41	564.98	7.073	6.121	s:\diskspd\io.dat (500
2	8814592	1076	4.20 j	537.98	7.437	5.880	s:\diskspd\io.dat (500
3	9101312	1111	4.34	555.48	7.181	5.555	s:\diskspd\io.dat (500
total:	36274176	4428	17.30	2213.90	7.219	5.817	
Read IO							
thread	bytes	I/Os	MB/s	I/O per s	AvgLat	LatStdDev	file
0	7397376	903	3.53	451.48	8.614	5.384	s:\diskspd\io.dat (500
1	7528448	919	3.59	459.48	8.480	5.954	s:\diskspd\io.dat (500
2	7086080	865	3.38	432.48	9.018	5.498	s:\diskspd\io.dat (500
3	7282688	889	3.47	444.48	8.740	5.136	s:\diskspd\io.dat (500
total:	29294592	3576	13.97	1787.92	8.708	5.508	
write IO							
thread	bytes	I/Os	MB/s	I/O per s	AvgLat	LatStdDev	file
0	1703936	208	0.81	104.00	1.026	0.370	s:\diskspd\io.dat (500
1	1728512	211	0.82	105.50	0.948	0.335	s:\diskspd\io.dat (500
2	1728512	211	0.82	105.50	0.955	0.328	s:\diskspd\io.dat (500
3	1818624	222	0.87	111.00	0.938	0.284	s:\diskspd\io.dat (500
total:	6979584	852	3.33	425.98	0.966	0.332	

I/O Analyzer

A virtual appliance solution

Provides a simple and standardized way of measuring storage performance.

http://labs.vmware.com/flings/io-analyzer

- Readily deployable virtual appliance
- Easy configuration and launch of I/O tests on one or more hosts
- I/O trace replay as an additional workload generator
- Ability to upload I/O traces for automatic extraction of vital metrics
- Graphical visualization



IO Blazer

usenıx

LISA16

Multi-platform storage stack micro-benchmark.

Supports Linux, Windows and OSX.

http://labs.vmware.com/flings/ioblazer

- Capable of generating a highly customizable workloads
- Parameters like: IO size, number of outstanding los, interarrival time, read vs. write mix, buffered vs. direct IO
- IOBlazer is also capable of playing back VSCSI traces captured using vscsiStats.
- Metrics reported are throughput and IO latency.

usage: ioblazer [options]

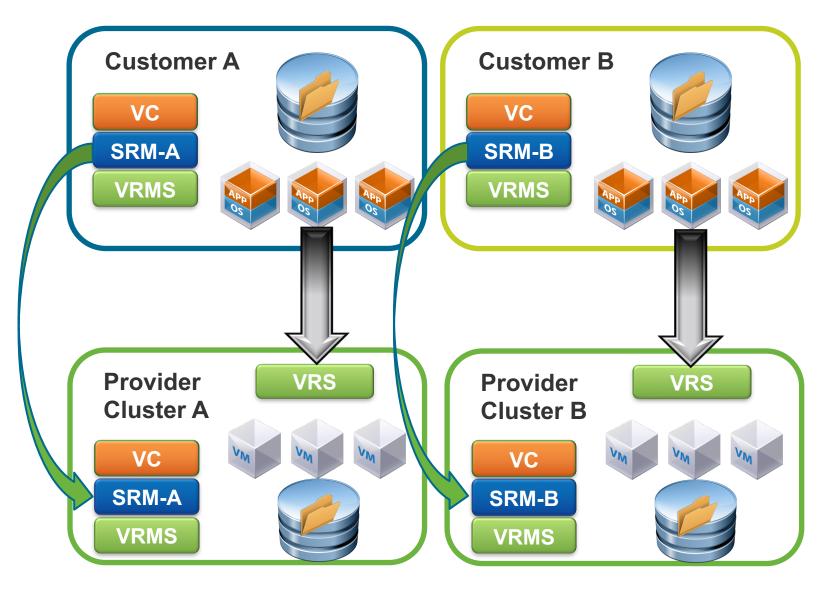
IOBlazer, the ultime IO benchmark

Options	Description	Default
-a <s r></s r>	IO access pattern: sequential or random	r
-b <size></size>	Memory buffer size in MB	1 MB
-B	Buffered IOs (go thru FS buffer cache)	0
-c	Computes pseudo-checksum on data	0
-d <path></path>	Device or File path	d:\tmp
-f <size></size>	File or Device size in MB	1 MB
-F	Fill the test file before starting IOs	0
-g <gap></gap>	Inter-burst time in microseconds (Avg)	0 us
-G <f u=""></f>	Inter-burst time pattern (fixed or uniform)	f
-h	Print this message	
-i <size></size>	IO size in B (Avg)	8192 B
-I <f b="" u=""></f>	IO size pattern (fixed, uniform or bimodal)	f
-1 <thresh></thresh>	Latency Alarm Threshold in ms	5000 ms
-o <#0I0s>	Burst size [aka Otstanding IOs] (Avg)	32
-0 <f u=""></f>	Burst size (fixed or uniform)	f
-p <c f="" z=""></c>	, ,	f
-P <file></file>	Playback VSCSI trace from 'file'	NULL
-r <ratio></ratio>	-	1.0
-R	Raw Device Access	0
-t <time></time>		180 s
-w <num></num>	Number of worker threads	1
- W Shulle	Autor of worker chiedus	-

Disaster Recovery with VMware Site Recovery Manager (SRM)



Architectural model #1 – Dedicated 1 to 1 Architecture



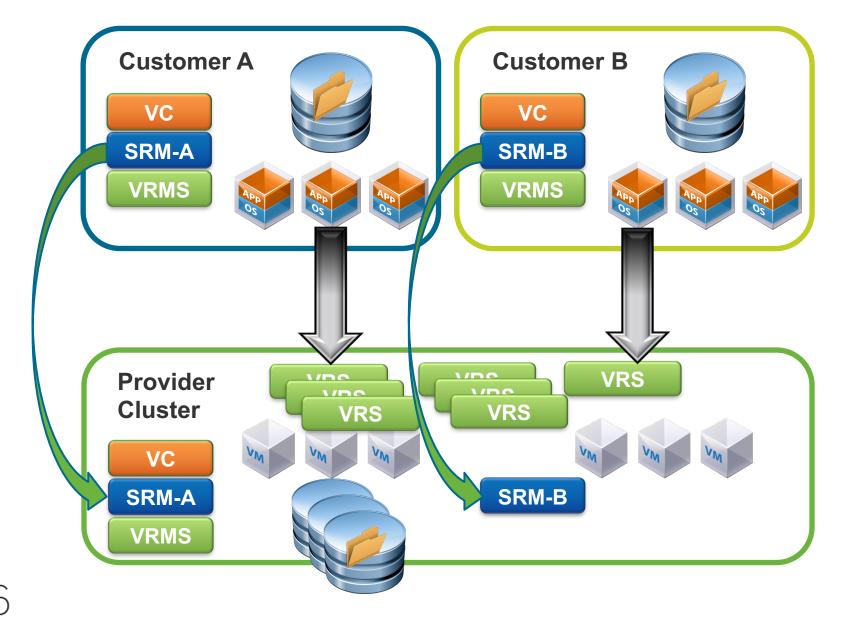
LISA16

Pros and Cons of 1 to 1 paired architecture

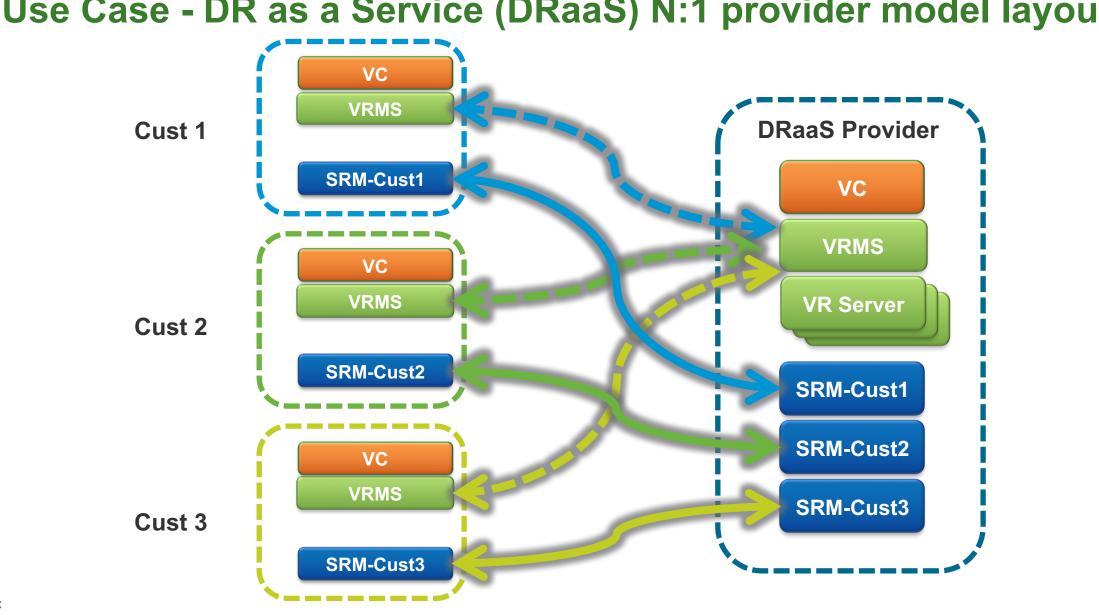
Pros	Cons
Ensures customer isolation	Highest cost model
Dedicated resources per consumer	High level of ongoing management
Can provide full admin rights to consumers	Wasted resources during non-failover times
Easy self-service for consumers	
Well known and traditional model for configuration	
Easy upgrades	
Custom options allowable per consumer	



Use Case – Shared N to 1 Architecture



LISA16



Use Case - DR as a Service (DRaaS) N:1 provider model layout

usenıx **LISA**16

Use Case – DR as a Service (DRaaS) provider model

Minimum Component Requirements

- Same as site-to-site requirements
- Remote customer site SRM "pairs" installed using SRM shared site option
- Remote customer site VRMS connection paired to recovery site VRMS as per default VRMS setup
- Typically provider runs whole solution as a managed service
- Provider usually own / administer all component VM's (SRM servers etc.) to reduce security complexities (i.e. user accounts / credentials)
- Current targeted N:1 limit is 10:1 meaning for each vCenter at the provider site there can be up to 10 inbound customers. To go beyond this scale out by adding additional clusters with own dedicated vCenter/VRMS/VR components
- Up to 500 VMs can be protected by VR under a single framework

Host Requirements

- Remote site VMs protected with vSphere Replication
 - You WILL need ESXi hosts to run those VMs on
 - Typically, provider will configure VR at customer site



Pros and Cons of Shared N to 1 architecture

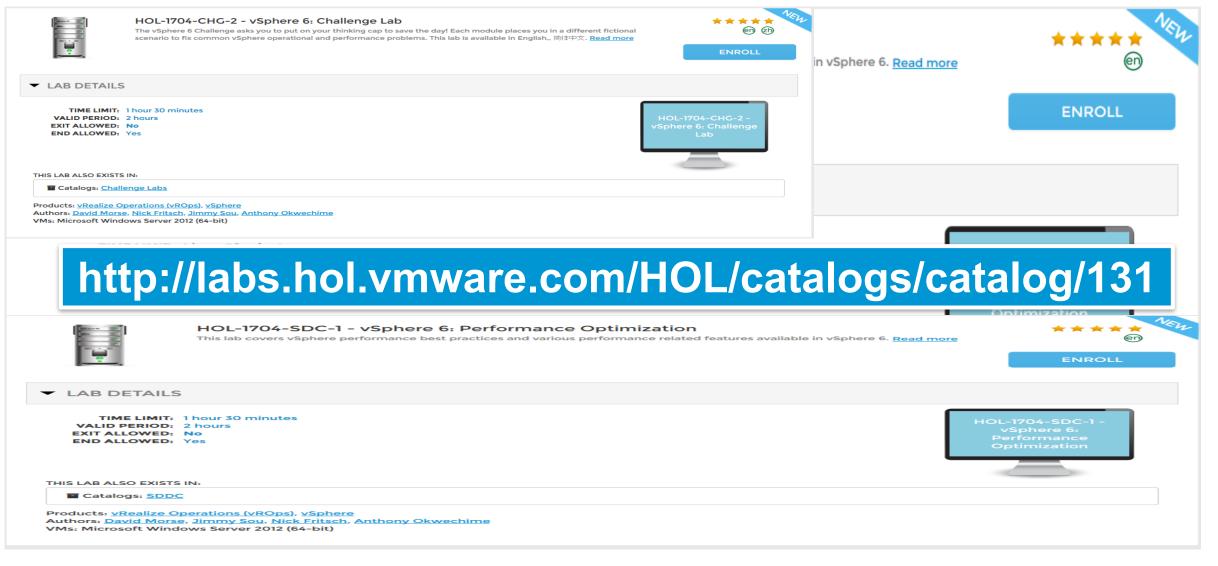
Pros	Cons
Lower cost of infrastructure	Difficult coordinated upgrades
Ease of management	Difficult to isolate customer environments
Ease of scaling	Cluster-wide events affect every customer
Central management of customer environments	More difficult to provide self-service
	Requires extensive role and permission management
	Scalability limits of 10:1



Resources



VMware Hands-on Labs



LISA16

The Links are Free. Really

Virtualizing Business Critical Applications

- <u>http://www.vmware.com/solutions/business-critical-apps/</u>
- <u>http://blogs.vmware.com/apps</u>

VMware vSphere 6.5 Document

- https://www.vmware.com/support/pubs/vsphere-esxi-vcenter-server-6-pubs.html
- https://pubs.vmware.com/vsphere-65/index.jsp
- <u>http://pubs.vmware.com/vsphere-65/topic/com.vmware.ICbase/PDF/vsphere-esxi-vcenter-server-65-setup-mscs.pdf</u>

VMware's Performance – Technical Papers

- https://www.vmware.com/pdf/vsphere6/r65/vsphere-65-configuration-maximums.pdf
- <u>http://www.vmware.com/files/pdf/techpaper/VMW-Tuning-Latency-Sensitive-Workloads.pdf</u>
- <u>http://pubs.vmware.com/vsphere-65/topic/com.vmware.ICbase/PDF/vsphere-esxi-vcenter-server-65-monitoring-performance-guide.pdf</u>
- <u>http://www.vmware.com/files/pdf/techpaper/VMware-PerfBest-Practices-vSphere6-0.pdf</u>
- http://www.vmware.com/pdf/Perf_Best_Practices_vSphere5.5.pdf
- http://www.running-system.com/vsphere-6-esxtop-quick-overview-for-troubleshooting/ ESXTop Cheat Sheet
- <u>VMware vSphere Data Protection Documentation page</u>



Questions? #rtfm

usenix LISA16

December 4–9, 2016 | Boston, MA www.usenix.org/lisa16 #lisa16

Thanks for attending

usenix LISA16

December 4–9, 2016 | Boston, MA www.usenix.org/lisa16 #lisa16