

VCAP-DCD

8/14/2014

Module 2: Design Process

Design Methodology: An iterative, high level process used to create a design.

4x Phases in vSphere Design Methodology:

- ① Architectural Vision
- ② Architectural Analysis
- ③ Technology Architecture
- ④ Migration Planning

① Architectural Vision: A high level vision for the project. Includes the following aspects

(1) Scope - Identify scope in detail!

- defined scope prevents unintended expansion. i.e. this project includes the United States production servers only. No dev/test and no Canada servers.
- prevents the need to renegotiate cost of a project or work for free.

(2) Goals - Set specific goals for a project.

- goals need to be specific, measurable, and actionable.
- without goals it is difficult to determine success and value of a project.
- i.e. the organization wants to achieve a 50% reduction in production server equipment by the end of the year.

(3) Requirements: Identify key business and technical requirements for the project.

- example: SOX compliance, physical separation from production + dev/test, uptime requirements, etc.

(4) Assumptions: Design components that are assumed valid without proof.

- i.e. the organization has sufficient bandwidth between sites for replication

(5) Constraints: Constraints limit design choices - could be a policy, process, or technical constraint

- i.e. due to existing relationships, all hardware is Dell.

(6) Risks: Identify risks that might prevent achieving project goals.

- i.e. lack of core redundancy introduces risk of 99.99% uptime.
- discussing risks eliminates surprises.

VCAP-DCD

8/14/2011

Module 2: Design Process (Continued..)

(1) Architectural Analysis: Deeper dive into current infrastructure.

Perform 2x Activities:

① Current State Analysis:

- gather and analyze info in current servers, storage, network,
- gather and analyze applications and operating systems.
- gather resource consumption data.

② Define target state required to meet organizational goals ↳ envision.

while factoring in requirements, assumptions, and constraints.

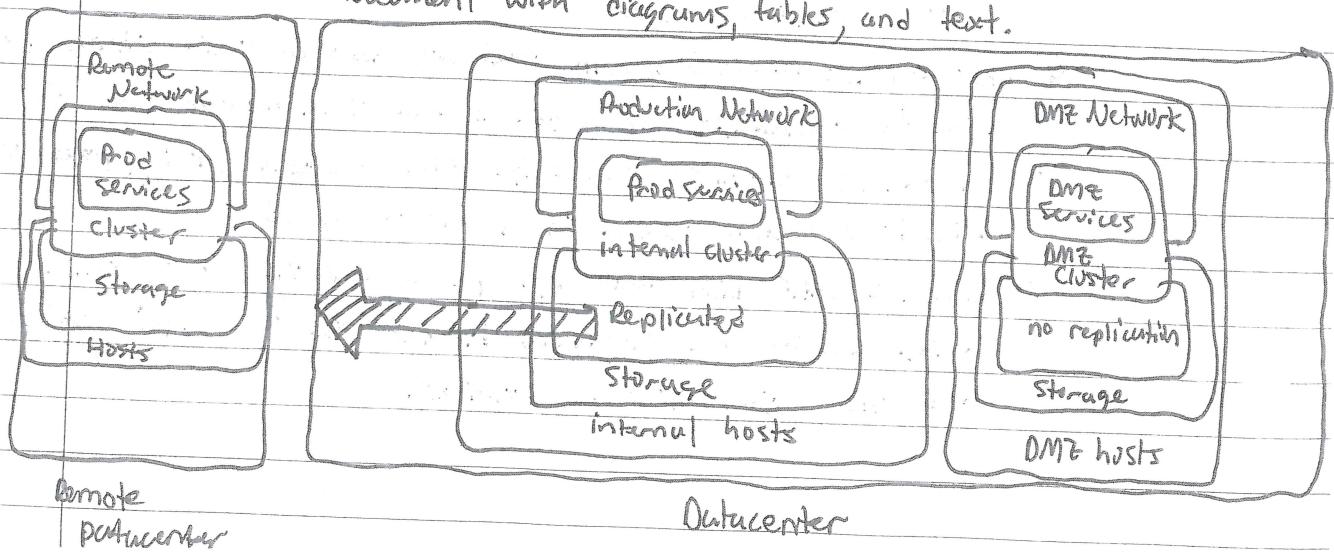
③ Technology Architecture: Develop 3x designs..

① Conceptual Design: focus on Achieving goals + requirements.

- use info from SME/stakeholder interviews (scope/goals/requirements/assumptions/constraints) and from info gathered in current state analysis
- Determine entities affected by project. (LOB, users, applications, processes, physical machines, etc.).
- Determine how goals map to each entity.
- Design infrastructure that achieves each entity's goals + requirements, but stay within constraints. i.e.: Where do you need availability, scalability, performance, security, manageability, etc.

• Document with diagrams, tables, and text.

Example Conceptual design diagram.



VCAP - DCD

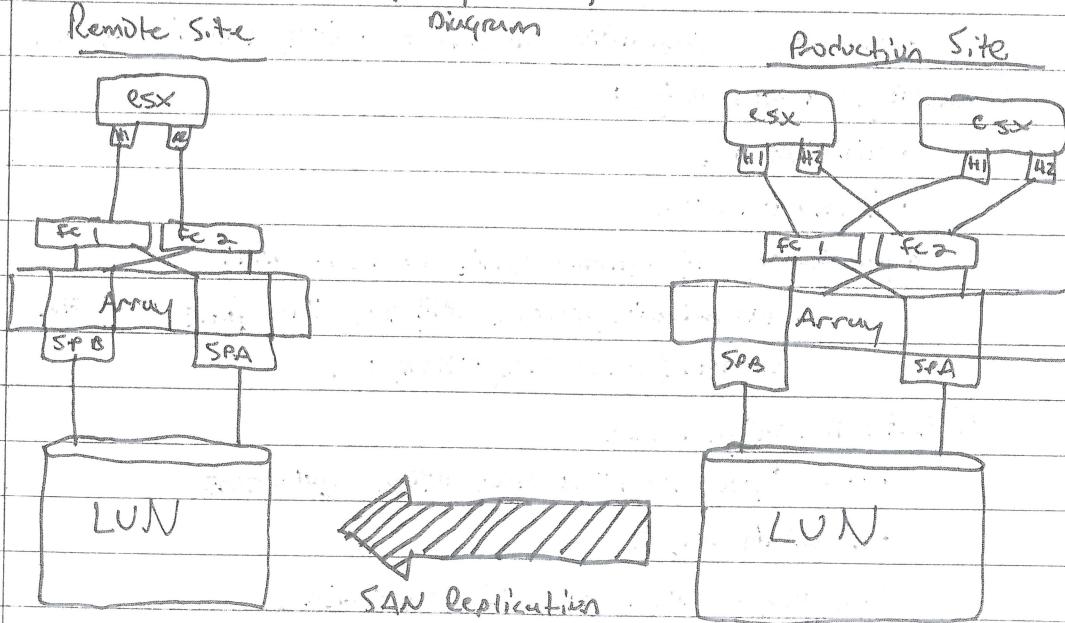
8/04/2011

Module 2: Design Process Continued...

② Logical Design: Design includes relationships between all major components of the infrastructure. Considers the conceptual design, constraints, and risks.

- Useful for understanding and evaluating the design of the infrastructure:
 - does it meet the requirements but stay within constraints?
 - Does NOT include physical details like port assignments, hardware vendor, IPs, etc.
 - illustrate how to arrange infrastructure components
 - don't get lost in configuration details.
 - be aware of capacity analysis, but don't include things like LUN sizing, CPUs, etc.
 - document in diagrams, tables, and text.

Example Logical Design



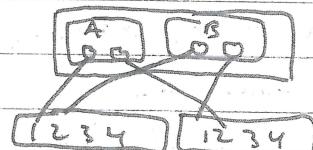
③ Physical Design: use the logical design; include specific hardware details and

Example Physical

design diagram

implementation information.

- includes port assignments, pci slots, etc.



VCAP-DCD

8/14/2011

Design Criteria

① Usability

1. Performance - may be measured by throughput, latency, transactions, etc.

2. Availability - defined as access to resources when needed.

achieved using redundancy

defined in SLA.

3. Scalability - support future growth while maintaining acceptable performance.

② Manageability

- easy to deploy.

- easy to administer + maintain

- easy to update + upgrade.

⇒ Simplification is key.

• unnecessary complexities can lead to failures, rises in costs, etc.

③ Security - a good design:

- minimizes risks

- is easy to secure.

⇒ good security design typically applies defence in depth:

- authentication, firewalls, DMZ, IPS, IDS, filters, VAN, etc.

- biometrics, smart cards, tokens, etc.

④ Cost: the design needs:

- to be "good enough" to meet business requirements

- fit within budget.

⇒ All designs are a series of compromises.

Design Approach: Sound Approaches include:

① USE A METHODOLOGY!

② Conduct interviews with SME's and knowledgeable stakeholders.

⇒ discuss design alternatives + effects on performance, availability, scalability, manageability, compliance, etc.

③ Thoroughly DOCUMENT all aspects of a design. Two engineers should be able to build the same infrastructure from the doc.

VCAP - DCD

8/14/2011

End Goals: Good Designs

① Involve SME's + stakeholders.

② balances business requirements with technical considerations and best practices.

③ is not unnecessarily complex

④ is REPEATABLE

⑤ is not ambiguous.

⑥ has documented rationales + considerations.

⌘ clearly document all design decisions.

⌘ Always involve SME's + stakeholders

⌘ create design documents like construction blueprints.

⌘ Simple designs are easier to understand, perform better, cost less to manage etc.

VCAP-DCD

8/20/2014

Module 3: ESX/ESXi Host Design.

CPU Capacity:

- Determine Requirements by using Capacity Planner, OS and app vendor guidelines/documentation, talking to SME + Stakeholders.
- Add needed growth to current usage for CPU capacity.
- DO NOT Plan to fully utilize CPU.
 - ↳ Plan for avg utilization of 60-70%

Leave headroom for → short term util spikes. (business cycle)

- Patching / Maintenance
- VMware Ht.
- Datacenter failover
- VMK + SC Overhead.
- Growth.

Number of Hosts: CPU Capacity.

- No matter what, host config must meet CPU Req.
- How many CPUs per host?
 - Cost of dense vs non-dense servers?
 - ↳ Consider not just the initial purchase, but the operational costs over time:
 - i.e. Power, Cooling, man hour overhead for maintenance, service contracts, etc, licensing

- ↳ Can enough memory be installed to keep the CPUs busy
- ↳ Are there enough storage + network ports available to provide bandwidth to support VMs that could run on a dense host

★ Consider VMware Ht and the # of VMs per host limits.
Isolation policies may drive more smaller hosts.

Cores per CPU:

- ★ The # of cores per host must match or exceed the # of vCPUs in the largest VSMP VM.

READ!!

★ CPU

Schedular
policy

VCAP - DCD

8/20/2011

Module 3 → ESX/ESXi Host Design.

of vCPUs per Core:

Check vCPU per core maximum.

VM Performance drives this limit

Goal is to efficiently use CPU capacity without reducing performance
maximum limit: Config mux:

effective limit: based on effective acceptable performance.

More cores = more sharing of memory & data busses.

* 4-6 w/ dual core processors

* 3-5 w/ Quad core processors. ← lower bc of bus sharing.

1-2 vCPUs higher w/ VDI.

features (CPU)

Buy the fastest CPUs and the most cache you can afford.

64-bit

Intel VT-x + EPT support

AMD -V + RVI. (Rapid Virtualization Index)

Licensing:

Must support # of sockets

Must support # of cores per pCPU.

Memory Capacity:

Same As CPU. See GPU.

DO NOT Plan to fully utilize memory resources.

70-90% util.

Leave headroom for util, maintenance, etc., ← Same as CPU.

VCAP - DCD

8/20/11

Module 3 Continued...

Memory Overcommitment:

2 types: total → does not affect performance

Active. → DOES →

→ AVOID ← Add more memory to design
↓ ↓ Reduce consolidation Ratio

Active Memory = Active Guest Memory.

Amount of memory in use by guest OS and Applications.

If active memory is oversubscribed:

- ① Add memory to each host.
- ② Reduce # of vms on each host.
- ③ Reduce vRAM allocation.

With large memory pages → Reduce memory overcommit.

Large pages are not shared.

Service Console Memory: (ESX only)

* If the host has more than 6GB of memory, configure the SC with 800MB of memory

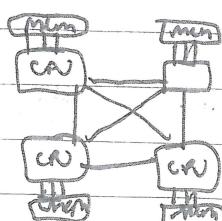
* Always configure MAX SC Allocation. (800MB)

* Configure SWAP w/ (1,600 MB) Partition

SWAP always bubble SC Memory config.

Cost of extra space is less than cost of reassign to resize.

NUMA - Nonuniform Memory Access.



Processor affinity may not be manageable by NUMA.

→ NUMA. → dense VMs that have ^{more} vCPUs than one socket cannot be managed

VCAP - DCD

8/20/2011

Module 3 (Continued)

NUMA

- ★ Evenly Distribute memory among CPUs.
- ★ Disable node interleaving in BIOS
- ★ By default only enabled on hosts with a total of at least 4 CPU cores and with at least 2 CPU cores per NUMA Node.

Node interleaving is like RAID 1 disks. All data sent to physical memory is written across all physical host memory. = data written to both Local AND Remote Memory. = BAD PERFORMANCE!

Host Hardware Type:

Blades or fixed boxes?

- ★ Choice is not always technical. *

Blades:

Good:

- flexible
- easy to install + replace
- simpler cabling.
↳ fewer I/O ports.
- ~~less power~~ (overall)
- management. (Open)
- more DC engineering.
- smaller footprint.

Bad

- locks into vendor (proprietary)
- hot spots
- more power + cooling per rack
- shared chassis + components could be single point, or could lead to larger outages.
- not as much expansion.

Design A Consistent Platform.

Same servers - CPU, memory, ports, PCI assignments.

↳ Simplifies troubleshooting, planning, automated installation, cloning, HA, etc.

VCAP - DCI

8/20/2011

Module 3 Continued.

ESX or ESXi:

ESXi: Embedded wherever possible.

Embedded or installable

- less attack surface

- easier to patch

- no SC consuming resources

What would force ESX over ESXi?

(1) Check Agents.

(2) vSphere Web Access not included in ESXi.

BIOS Settings:

Enable VT.

Disable node interleaving.

Disable unused USB ports, serial ports, network ports, etc.

Disable power saving options for Power + CPU & FSP w/ FT.

PCI Slot Design:

SAME CONFIG!!!

↳ Consistent HBA Numbering: Simplify troubleshooting, auto install + config, easier administration.

1Gb cards should go in PCIe RX slots to reduce bus limitation

Power & Cooling:

Design for Redundancy:

check vendor limits (power) per rack.

VCAP - DCD

8/20/2011

Before placing Hosts into Production:

- ① Burn in all equipment.
- ② Test memory.
- ③ Run HW Diags.

* Assign Static IP Addresses.

If IPs change; vCenter needs to know.

Configure hostnames in DNS!!!

Standardize naming conventions!!!

Host Security:

GSXi no SC = more secure.

GSXi has lockdown mode. (prevents remote access as root)

* Use Centralized Management!

→ vCenter

vMA.

PowerCLI.

— Limit local user accounts.

— Limit root access.

Limit installed agents.

install agents in vMA where possible.

VCAP - DCD

8/20/2011

Module 4: vSphere Virtual Datacenter Design.

vCenter w/ 300 hosts or 3000 VMs. Require 64-bit hosts. (200P 64-bit)
Best Practice 200 hosts or 2000 VMs. ↗

Requires 2CPU.

3GB Memory. → increase if SQL runs there.

2 GB Storage.

Host	VMs	CPU	Memory
50 host	250 VMs	2	4
200	2000	4	4
300	3000	4	8

vpxa.cfg → IP of vCenter system on each host.

Virtual or Physical?

Best Practice: Use a virtual machine.

If key stakeholders insist on a physical machine, list as a risk!

- | VM | Physical |
|---|---|
| <ul style="list-style-type: none"> Protect with HA or VM HB. Dynamic Sizing easy image backups. allows HW maintenance of vCenter root server. | <ul style="list-style-type: none"> Protect w/ VM HB only. Static Buildups / Restore harder |

Manual Failover: initial setup cheap → operating expense expensive.

Auto failover: vCenter HB.

VCAP-DCD

8/20/2014

Module 4

vCenter Protection

- ① Manual - initial cheap, long term expensive. complex, high overhead. failure difficult.
- ② Cluster w/ HB - expensive, but automatic + fast.
- ③ Backup/Restore - too long RTD for most.
- ④ VMware HA. - always use.

↳ disable DRS Migration.

MSCS supported in vCenter 4. Not 4.1.

vCenter becoming Tier 1.

vCenter Database:

4.0 vCenter requires 32-bit DSN. 4.1?

Best Practice: vCenter existing database server. if VM, better.

PSA: existing db machine.

Only small infrastructures should colocate vCenter, db, + vCenter server.

SQL express 5 hosts 50 VMs.

High availability of Database system: use vendor tools if possible. otherwise use
VMware HA. ↳ does not protect db info.

Certificates:

- if concerned about man in the middle attacks from management servers, use certs.
- FT and linked mode require vCenter to have certificate checking enabled.

VCAP-DCD

8/21/2011

Module 4: vSphere Virtual Datacenter Design.

Clusters and Resource Pools.

Critical workloads should be protected by VMware HA:

Scale up or Scale Out?

- Cost of hardware

- Cost of licensing.

- Cost of operational overhead.

- ↳ maintenance

- ↳ patching

- ↳ power, cooling, rackspace (especially in Datacenter)

Scale up clusters	Scale out clusters
<ul style="list-style-type: none"> • easier to manage fewer hosts, reduces costs. • Host failure affects more virtual machines • HA failover takes longer. • Certain failover policies might reserve more resources. Typically # of Host Failures + specify failover host. • need to be mindful of virtual machine per host limits. • DRS has fewer migration choices, so workloads could be less balanced. 	<ul style="list-style-type: none"> • host failure affects fewer VMs. • failover is faster. • fewer resources reserved for failover. • more DRS migration choices • may affect size of virtual machines. • more reestablish required. • more power + cooling. • higher operational overhead • expensive to manage • more licensing.

Mixed Clusters:

- if new features in VM HW version 7 are required, do not use mixed clusters.
- if they are not needed, mixed mode clusters are fine.

VCAP-DCD

8/21/2011

Module 4.

Cluster Failure Policies.

Smaller clusters w/ # of hosts failed tolerates vs designated failover node policies configured, more resources will be reserved for failover.

Ex. 5 hosts equally sized, if 2 hosts tolerated, 40% of resources in CLV or

→ Avoid above situation, configure percentage based policy. Risk:

↳ downside is that guarantees cannot be made on restarting all VMs.

explain risks to SME + stakeholders.

VMware recommendation is to protect the workload.

Demand Failing recommends % based.

Reservations:

- Powered on VM reservations determine the slot size used by Host failures cluster tolerates policy.
- Even a single VM w/ abnormally large reservations will reduce the # of slots available to run VMs, because the slots are bigger.
- Host failures cluster tolerates works best with large clusters where VMs have similarly sized reservation settings.

In clusters w/ VMs that have highly variable reservation settings: 3 choices:

Best WL ← • Audit the usage slot size + understand fewer VMs can power on.

- Configure das.slotCpuInMHz to manually set CLV slot size
- Configure das.slotMemInMB to manually set memory slot size.

→ Risk here is that any VM w/ reservations larger than the slot size may not power on in a failover.

- Configure a failover policy that is not based on slot sizes.

↳ either % based.

↳ Specify ^{failover} host policy: → Risk is if host fails while another is in maintenance mode, VMs may not have resources to start.

VCAP-DCD

8/21/2011

Module 4: vSphere Virtual Datacenter Design

Cluster Failure Policies:

Host failures Cluster Tolerates:

- Default.
- Always reserve enough failure capacity to accommodate host failures while others are offline or in maintenance mode.
 - ↳ does increase costs, however should still be getting cost savings from consolidation
- = Always add 1 to the ^{# desired} host failure cluster tolerates to accommodate host failure during maintenance.

* Always configure strict Admission control to protect the workloads

Only turn it off in extreme circumstances.

make Admission control (strict) a part of Change Control policies.

↳ powering on unprotected vms should require change control pr
↳ makes this a management choice, instead of an admin choice.

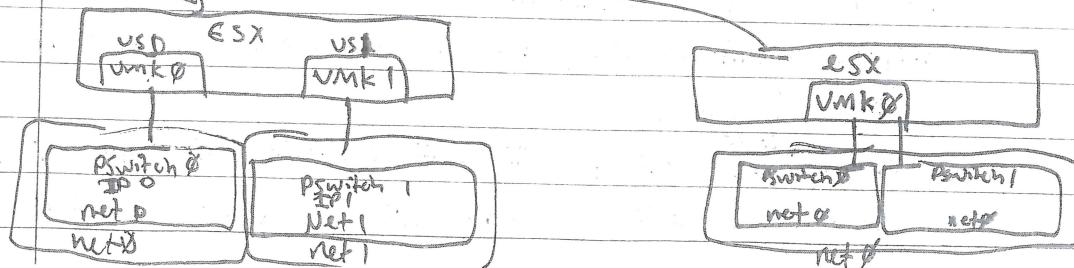
HA Network Redundancy:

* Always configure it.

A switch failure could trigger datacenter wide isolation response.

2 choices:

- o 2x separate networks, separate subnets, etc.
- o Single mgmt network w/ NSX teaming to multiple switches.



VCAP-DCD

8/21/2011

Module 4:

* HA will use all mgmt interfaces to send heartbeats.

↳ in esxi on all vmkernel networks (except for those marked as vMotion).

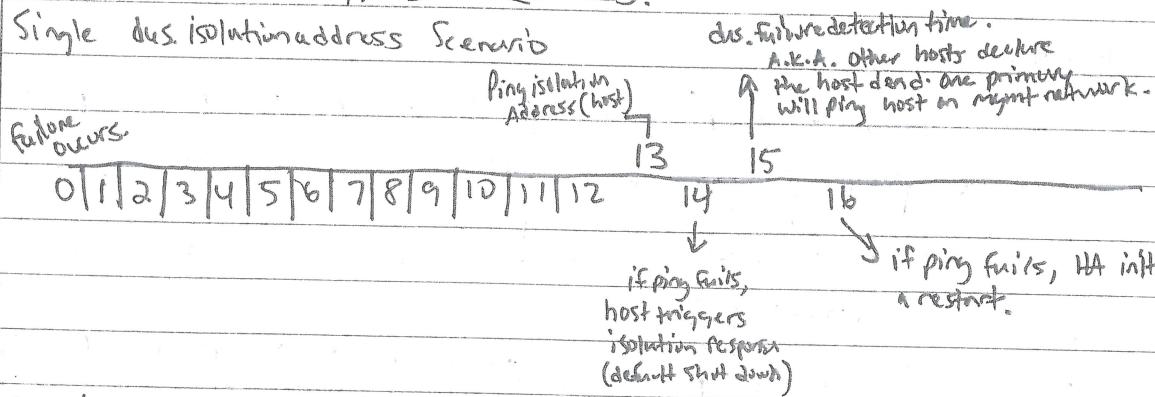
If there are 2x management networks, use dvs.isolationaddress to add an isolation test address for each additional management network.

↳ eliminates isolation test address as a single point of failure.

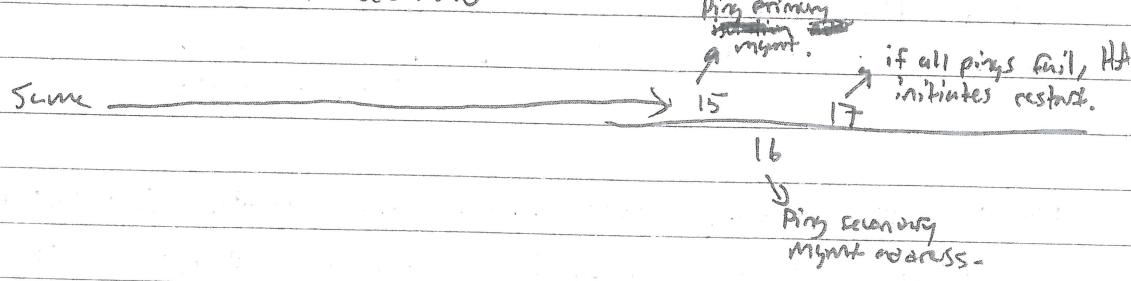
* When you specify an additional test address, increase the setting for dvs.failuredetectiontime to 20,000 milliseconds or greater.

~~Failure occurs 10 seconds.~~

Single dvs.isolationaddress Scenario



2x dvs.isolationaddress Scenario



2 (or 3) second time limit where isolation response could initiate a power off at second 14. If the heartbeat comes back at second 16 or before, the VMs will already have been powered down but HA will do nothing because the heartbeat returned.

VCAP-DCD

8/21/2011

Module 4 VSphere Virtual Datacenter Design.

Isolation Response:

Default is Shut down VM. - must have VMware tools. 300 second timeout until VMs are powered off.

If the mgmt network has redundancy, shut down VM is appropriate

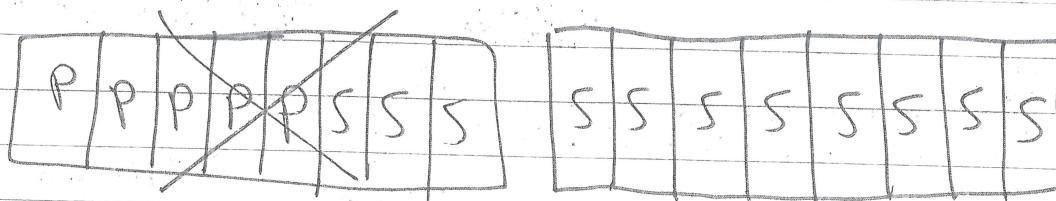
If not, select "leave powered on."

VM Monitoring:

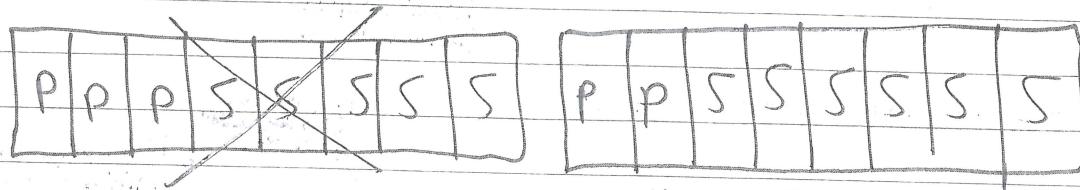
UX it.

FT and MSCS will respond faster than VM monitoring.

Blades:



Cluster fails



Cluster Survives.

Primary hosts maintain cluster consistency and handle failover events. If a primary is lost, a secondary is promoted to primary as long as there is at least 1 surviving primary.

* Never add more than 4 hosts in the same blade enclosure to a single VMware HA cluster. ← guarantees that all primaries will not hit the same chassis.

VCAP-DCD

8/21/2011

Mod 4 → VMware HA Cont..

VMware FT.

- Only single CPU Workloads.
- FT VMs cannot have snapshots, DRS, or Storage vMotion.
- if using FT, clusters should not have both esx / esxi hosts. One or the other.
- FT VMs need to be in clusters with at least 3 hosts.
↳ ensures redundancy if a host fails.
- FT must satisfy HA cluster requirements. ← only in HA cluster.

Use hosts with the same CPU speeds.

Disable CPU power saving in BIOS on hosts running FT workload.

Configure Minimum 1Gb ethernet logging.

w/ 1Gb net, no more than 4x VMs protected by FT on that host.
10Gb net is better.

Distribute primary VMs among hosts to reduce contention on the logging NIC
→ logging is asymmetric. Primary ⇒ secondary.

Create redundant logging networks.

An FT VM automatically reserves all of its memory. Configure resource pools ~~with~~ memory to have memory above configured memory for the VM for overhead.

VCAP - DCD

8/21/2011

Module 4 Continued.

DRS Clusters:

- vMotion Requirements must be met.
- Enabling Enhanced vMotion Compatibility to future proof cluster.
- Configure DRS for full Automation, w/ Default 3 Star Rating
- DRS can balance workloads better in scale out scenarios.
 ↗ gives the algorithms more choices.

Use Affinity + anti-Affinity as the exception, not the norm.

↳ many rules could limit migration choices.
 has negative effect on the scheduler.

Affinity can be useful:

- if VMs have a lot of network traffic
- if VMs share similar memory pages - TPS can be more efficient.

anti-affinity is used to increase availability for service workflow

Multiple Clusters:

Reasons for multiple clusters:

- Cluster may hit host limit (32)
 ↳ VMs per host limits.
 ↳ Hosts per LUN limit.
- desktop cluster vs server cluster.
- Organization requires hardware isolation.
 ↳ security
 ↳ dev/test/prod
 ↳ departmental for cost accounting.



VCAP - DCD

8/21/2011

Module 4 Continued --

Cluster Size:

SMBs will typically deploy 1 cluster

Larger orgs with 20+ hosts may choose multiple clusters

↳ Physical separation

↳ avoid limits.

⊗ In general, medium sized clusters perform better than large clusters ↳ cluster operations will be faster.

↳ consider 4x 8-host clusters instead of 2x 16-host Cls because of sluggish vCenter performance.

⊗ Avoid Extremes.

Resource Pools.

Deploy resource pools with limit and reservation settings ONLY

For virtual machine workloads that require DEDICATED and ISOLATED Resources.

Easier to allocate reservations and limits at the RP level, rather than at the individual VM level.

||| ⊗ Resource Pool CPU and Memory share settings do NOT get pushed down to the VM level. If VMs in a specific RP need higher shares, they must be configured manually.

If common VMs are created + need share settings, create a VM template for it.

VCAP-DCD

8/21/2011

Module 4 → vSphere virtual Datacenter Design.

VMware DPM. - Distributed Power Management

- Enable DPM in environments where workloads vary significantly over time.
- Buy hosts with iLO or IPMI (intelligent Platform mgmt); (interf)
- if hosts don't support above, buy WOL NICs.
- Configure DPM for Automatic. Use the default threshold.
 - ↳ decreases power/cooling costs.
 - ↳ decreases mgmt overhead.

Naming Conventions:

USE Standards.

Pick meaningful names.

VCAP-DCD

8/24/2011

Module 5: vSphere Network Design.

Design Infrastructure that reduces costs, boosts performance, improves availability, provides security, and enhances functionality.

Configuring 2x 1Gb NICs could be more cost effective than 1x 1Gb NICs on separate physical networks.

High level Design Requirements:

- Connectivity, Bandwidth, Latency, Availability, cost.

Build a modular network solution. → helps control future costs.

Components that ~~should~~ be on separate networks:

- VMs
- VMware vMotion
- VMware FT
- IP Storage
- VMware HA
- Management

2 Reasons to separate Networks:

- Reduce latency & improve performance (esp IP Storage) ← Reduce contention
- Improve security → FT, vMotion, etc.

Architect enough Physical NIC ports for bandwidth.

- Use current state analysis (capacity) to determine requirements.
- Consider failure. Will there still be enough BW left after device failure?

Segmentation!

VLANs or Physical Separation?

Depends on # of NICs, # of switchports,

The preferred choice is VLANs, however some orgs may have policies against them (constraint)

VCAP - DCD

8/24/2011

Module 5: vSphere Network Design.

Physical segmentation will likely require bigger hosts with greater I/O expansion. More NICs, more switch ports, more cables, more rack space, etc.

Other factors.

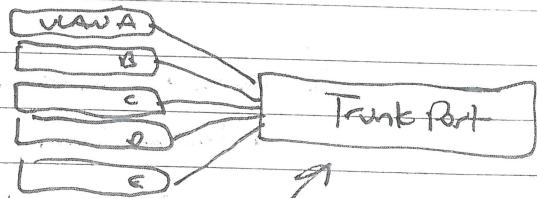
- 1Gb Ethernet will not be able to handle as many VLANs than 10Gb.
- Physical infrastructure may not support VLANs.
- Security policies may require physical separation.
- Organization may not currently use VLANs.

VLAN Benefits

- Reduced physical ports.
- Reduced cabling.
- Flexibility.

VLAN Risks

- port failures can be more serious
- easy to oversubscribe & create contention.



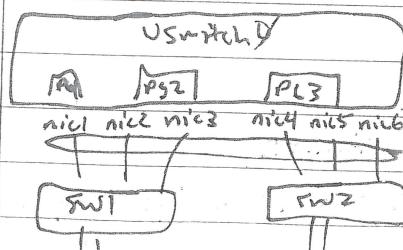
Implementing VLANs:

- 802.1Q Switches.
- Configure switch ports as trunks.
- Configure vSwitch trunk ports
- Configure STP PortFast on all ESX-facing ports!

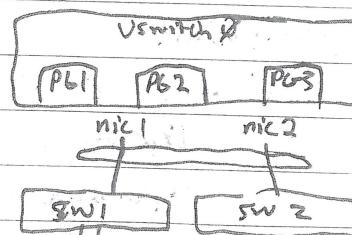
enough bandwidth?

Configure learning!

1Gb Ethernet



10Gb Ethernet



VCAP-DCD

8/24/20 11

Module 5, VSphere Network Design.

Private VLANs (PVLANs)

Provide layer 2 isolation between servers on the same layer 2 network without needing to do MAC ACLs.



This can reduce the # of VLANs required and still achieve isolation.

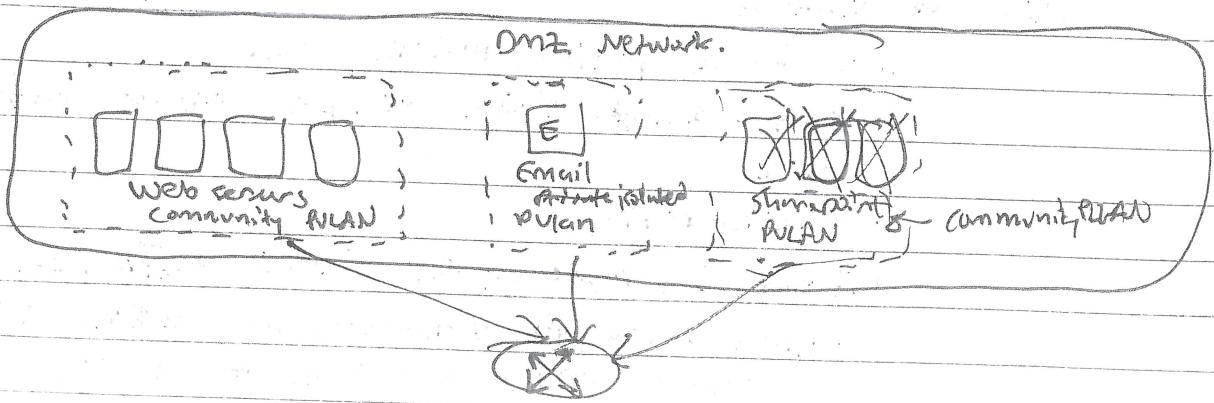
Configure PVLANs to protect the servers from one another if a machine becomes compromised.

Community PVLANs allow servers in the same community PVLAN to communicate with one another, but are isolated from the rest of the machines on the same layer 2 subnet.

Isolated PVLANs mean each server is isolated from other servers on the same subnet.

Important for DMZ/public networks.

Example:



Router in Promiscuous VLAN,

VCAP-DCD

8/24/2011

Module 5: vSphere Network Design

Type of vSwitches:

Distributed vSwitch	Standard vSwitch
<ul style="list-style-type: none"> Centralized management. consistent configuration. offer PVLANs, load based teaming, persistent network state, bi-directional traffic shaping. vCloud + Labmgr cross host fencing require dswitches. not manageable when vCenter is down. Requires ent + licenses. 	<ul style="list-style-type: none"> required if blending with vIs.

Nexus 1000v
<ul style="list-style-type: none"> consistent mgmt interface extends virtual network feature set to a level consistent w/ circuit switches. satisfies mgmt requirements by allowing the network team to manage the network infrastructure Ent + <u>nexus</u> license.

Q. How Many vSwitches?

- A. As few as possible. Preferably one.
- Configure a single vSwitch with port groups for each type of traffic.
 - Simplifies config.
 - If there is a requirement that VM-to-VM traffic passes through a physical firewall, the infrastructure needs multiple vSwitches.
→ could use vShield zones if acceptable by org.

NIC Teaming:

Increases BW, Redundancy.

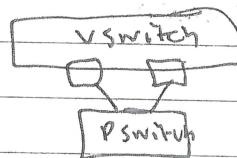
Try to team Across Onboard + Riser NICs.

VCAP-DCD

8/24/2011

Mod 5: Network Design.

- Create a single virtual switch with teamed Nics across separate physical switches:

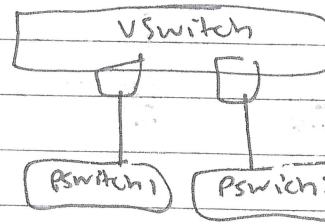


OK.

Source port ID

Source MAC Addr

IP hash.



Source port ID

Source MAC Addr

IP hash (with stacked switches)

Use originating port ID as default, IP hash where hardware supports it.

Use Active / Standby port configuration to load balance between port groups.

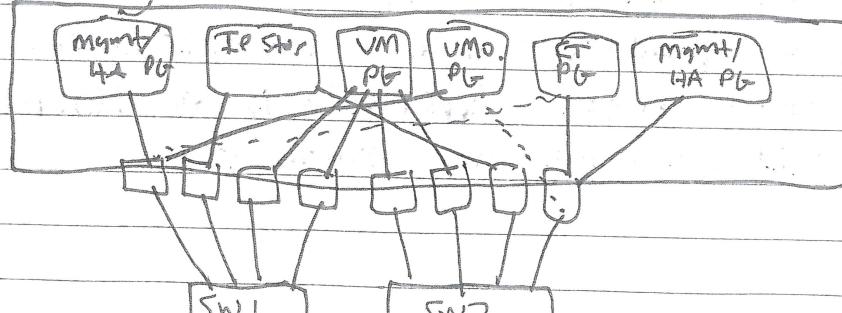
* Management Network needs to be redundant for HA isolation testing + heartbeats.

* FT logging network should be redundant.

* Beacon probing only works w/ teams of at least 3 ports, connected to separate physical switches.

* Beacon probing is unable to detect upstream network failures, that requires link state tracking or dependency ports in DELL.

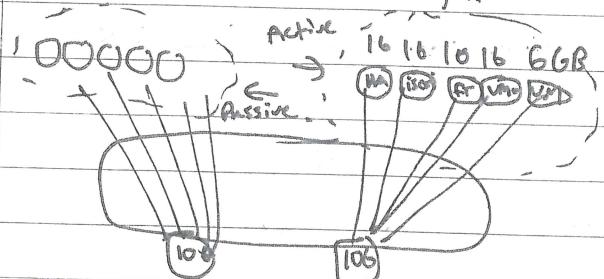
16B Sample Config



VCAP-DCD

8/24/2011

Mod 5 : Network Design.



1GbE example Active/Passive w/ traffic shaping.

← ingress traffic shaping

put 1GbE Adapters in PCIe 8x slots.

dual port 1GbE adapters in PCIe 16x slots.

Buckplanes can be a bottleneck (esp. riser cards)

Design book says to put SC/VMK nics in vSwitch. Real world suggests otherwise.

vSwitch Security Design:

- Change default settings for forged transmits and MAC address changes to REJECT, unless otherwise required not to.
↳ prevents compromised VM from using MAC spoofing to impersonate another ~~host~~ server in the network.
- * Need to allow → NLB.
- Testing IOS,
when App licensing depends on MAC Addresses.
- If MAC changes + forged transmits are needed, enable them on a specific port group.
- Leave promiscuous mode disabled unless needed. Enable only on required port group. Apps needing it are IOS, packet capture utilities, and performance monitoring tools.

VCAP - DCD

8/26/2011

Mod5 Network Designs.

Physical NICs.

- Configure at least 2x pNICs for each network.

Checksum offload

TCP segmentation offload (TSO)

64-bit direct DMA.

Jumbo frames.

- Leave Auto negotiation enabled.

- Use 1Gb Adapters that support NetQueue

Multiple queues
each queue associated w/ MACs
managed by the hypervisor

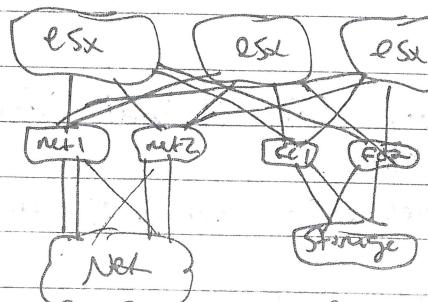
FCoE

- If the infrastructure supports it, use FCoE to converge FC/IP traffic, reduces # of HBAs.
 → CNA not supported w/ iSCSI or NFS in VS4.0.

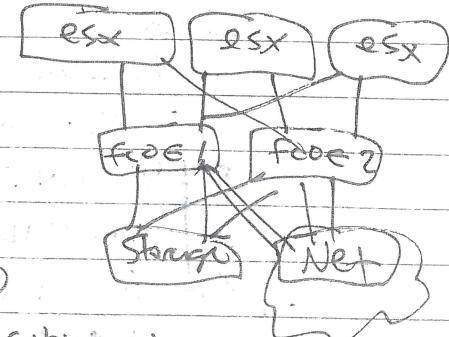
CNA appears to ESX as 1Gb Ethernet NIC + FC HBA: use priority queue bandwidth reservations to slice traffic up. t tool on the CNA.

Converged Design Comparison

Traditional:



Converged:



Result = Fewer ID

Ports, switches, cables, etc.

VCAP - DCI

8/24/2014

Mod 5: Network Design:

Physical Switches:

- Always redundant.
- Manually configure trunk ports if using VLANs.
- STP Portfast enabled.

Jumbo Frames:

- Boot says use them. Community Disagrees.
 - ↳ for IP Storage
- for VMs, leave Standard MTU unless specific app requires it

VM DirectPath ID.

- Don't use unless there is a clear use case.
 - ↳ are there examples of Apps with high transactional workloads that require the lowest possible latency?
 - ↳ esp effective on high speed devices like 10GbE
- What is it?
 - ↳ provides access to a network device while bypassing the VM layer.

Limitations:

- Limited Hardware Support: Only w/ Intel VT-d, only on a few 10Gb network controllers.
- ✗ You cannot assign devices in the same PCI slot to both VMdirect and VMkernel.
 - ↳ i.e. on a 10Gb NIC (Dual port) you can't split it so 1 goes to VMDPID and the other to another Port group.
- NOT compatible with vMotion, HA, DRS, hotAdd,

VCAP-DCD

8/24/2011

Module 5: Network Design.

IPv6: Enable it only if Required:

IPv6 increases mgmt overhead.

No compelling reason to enable.

Limitations:

IP Storage is experimental on IPv6.

HA/FT not supported

Requires vswitches.

DNS.

Resolve long ANY short Names.

Records for the following, both forward & Reverse:

- ESX/ESXi hosts
- vCenter
- vCenter Module Systems
- IP Storage.

Always Highly Available.

Configure each host w/ redundant DNS servers., sum for VMs.

* Configure common port group names.

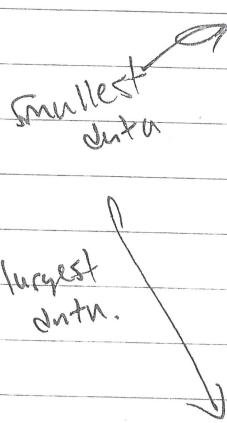
VCAP-DCD

8/21/2011

Mod 6: Storage Design.

Designing Storage Tiers is cost efficient.

Tiers: Mission Critical Data - Highest performance, 5 9's, Highest Cost, near 0 downtime

Smallest data
Business Critical Data - Good perf + 5 9's, less than 8 hr recovery, moderate cost

accessible-online data - Online performance, 4 9's, less 8hr recovery, lower cost

near-line data - Large volumes, protected data 3.9's, less than 1hr auto retrieval, cost-sensitive

Offline data - Archive data for backup or compliance, 72 hr retrieval, very cost sensitive.

Determine what tier data falls into:

- IOPS
- Capacity
- latency.
- Throughput
- availability

Consider SCA's.

Data may move b/w tiers during information lifecycle.

Tiering Example:

Tier 1	FC	15k	R10	8 Disks	10 VMs per disk
Tier 2	SAS	10k	R5	6 Disks	15 VMs / DS.
Tier 3	SATA	7.2k	R6	6 Disks	15 VMs / DS.

VCAP - DCD

8/24/11

MoC6 Storage.

What Protocol?

iSCSI, NFS, FC, FCoE?

- Factor in existing expertise.
- cost
- vendor relationships
- existing infrastructure.
- Performance (throughput + latency) requirements.

IP Storage Design.

- Use a dedicated storage network.
→ increases security.

If you can, physical isolation is preferred.

- limit the # of switches to traverse. more switches = more latency
- routing iSCSI/NFS is not supported.
- Configure Jumbo frames.
- use TSO NICs.
- Storage Design should be modular + scalable.
↳ be able to add spindles
front end ports
network infrastructure.

Calculating Capacity:

use CP or vendor tools.

factor in swap files + snaps. ← 20-30% for this.

factor in growth.

← store sizes can be imputed by RTD limits.

VCAP - DCD

8/20/2011

Module 5: vStorage Design.

12 to 15 VMs per VMFS Volume.

20-30% for overhead.

VMFS Block Size: MAX Disk Size

1 MB 256 GB

2 MB 512 GB

4 MB 1024 GB

8 MB 2048 GB

Determine block size by largest required virtual disk.

Block size also used to extend thin LUNs. inconsequential for design.

Recommend always using MAX block size.

Reducing latency:

- Dedicated Networks • reduce queuing.
- jumbo frames

Command Queuing on Host:

- VMs on the same host share a LUN Queue depth and sometimes
- If more IOs are generated than the Queue depth, IOs are queued in the kernel and latency goes up.

Queuing @ the Array.

- IO can bottle up @ the LUN Queue or at the individual disk Queue, depending on the array.

Queuing = BAD

VCAP-DCD

8/26/2011

Module 5 vStorage Design.

- ensure IOs are not consistently hitting LUN Queues.

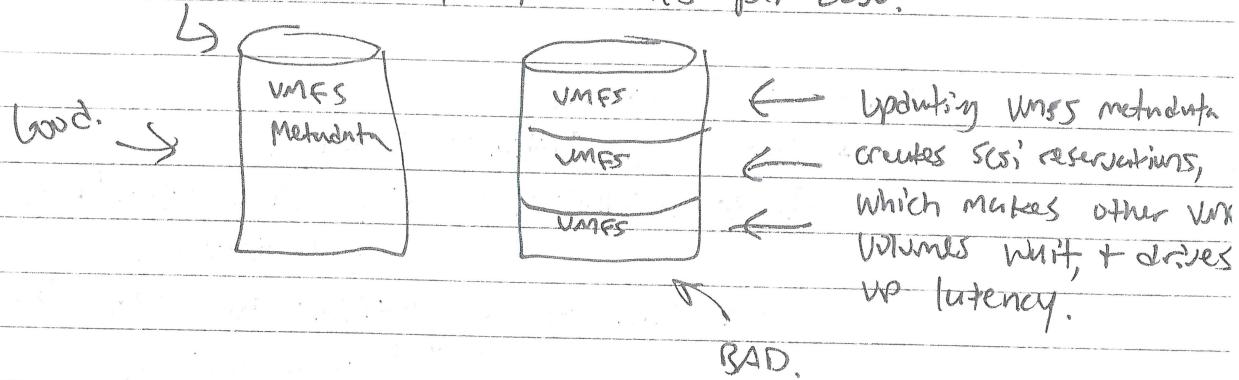
↳ either increase the LUN Queue depth or

↳ move VMs to another (less congested LUN)

use vSustains to see if queue depths are being stressed.

VMFS volumes per LUN:

Single Partition w/ Single VMFS per LUN.



Latency:

Below 15ms Good (green)

B/W 15ms + 30ms (yellow)

Above 30ms Bad. (red)

Improve latency? Reduce oversubscription!

Control Access to NFS:

- ① Network Segmentation
- ② Don't mount datastore
- ③ Datastore permissions.

ISCSI ACCESS:

- configure static IPs.
- physical + VLAN separation
- CHAP.
- Datastore permissions.

VCAP - DCD

8/21/2011

Module 5: vStorage Design.

Fiber Channel Access Control:

- ① Zoning
- ② Lun masking
- ③ Datastore Permissions.

Zoning is done for Performance + Access control.

↳ Controls which HBAs have access to which Storage Bits.

Zoning for Performance:

* Use Single initiator Zoning!!

Single initiator: → single HBA port w/ one or more target devices.

Multi-initiator: → multiple HBA ports w/ one or more target devices.

no need for initiators to talk to each other.

also risks ~~RSCN~~ RSCN (Registered State Change)

Notification messages from interrupting normal IO.

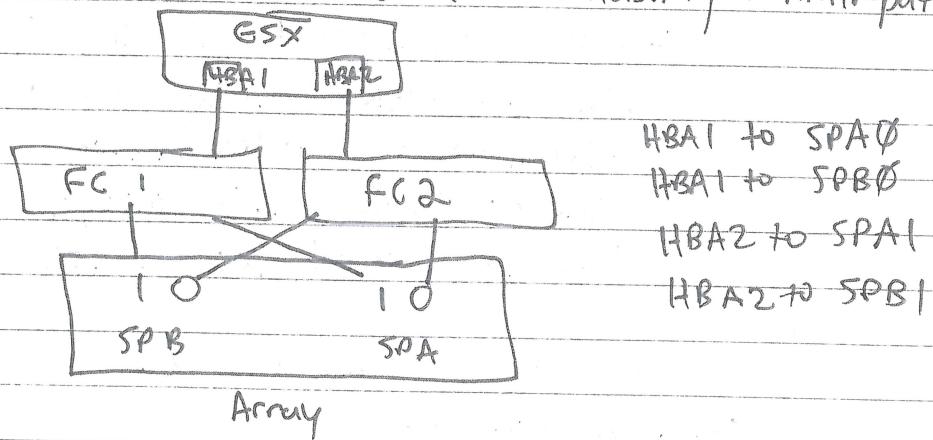
→ RSCN Messages are sent when an initiator enters or leaves a fabric. This initiates a Loop Initialization Primitive (LIP), which can interrupt IO from other initiators. Single initiator zones cannot interrupt other initiators w/ LIP's due to RSCN because the initiators cannot talk to each other.

VCAP - DCD

8/24/2011

Mod 5 : vStorage Design.

Zone for performance: Make sure HBA's are zoned to SP ports on different controllers for availability & multi-pathing.



Zoning for Access Control:

Security.

Performance. Limit LUNs per host

Cluster members need access to the same LUNs.

Use WWN Zoning rather than port zoning!

↳ can move cables to different ports w/o affecting zoning

- LUN Masking: Can be used when Access Zoning is not enough.
 - Masking on the Array is preferred, although it is possible in the Host.

Use WW Masking for Boot from SAN Requirements.

Use the same CUN #s across hosts, esp for arrays that do not support NAA numbers.

Multi-pathing : Recommends 4 paths per LUN for max performance + availability.

VCAP - DCD

8/24/2011

Mod 5: vStorage Design.

Use 2 single port HBA over 1 dual port HBA.

- Verify that the correct MP Policy is selected for the array
 - MRU for Active/Passive → avoids path thrashing.
 - fixed, RR for Active/Active
 - MRU or RR for ALUA arrays.
 - MRU for virtual port arrays.

Use the vSphere client to create volumes. Auto Aligned on 64K boundaries.

VMDK on boot disk cannot be aligned.

VMDK or RDM?

VMDK where possible. Performance is not a reason to use RDM.

RDM Required for:

MS Clustering.

NPIV

SAN Mgmt inside VM.

Anything needing access to HW specific SCSI command

Use virtual compatibility where possible.

Allows snapshots.

WZK8 clustering requires physical RDM Access.

Thick or Thin?

Thick

Thin

- | | |
|--|---|
| <ul style="list-style-type: none"> Support for FT. Simpler config & monitoring less efficient more expensive | <ul style="list-style-type: none"> decreases cost. pay as you grow. high risk. complex. |
|--|---|

- | | |
|--|---|
| <ul style="list-style-type: none"> Support for FT. Simpler config & monitoring less efficient more expensive | <ul style="list-style-type: none"> decreases cost. pay as you grow. high risk. complex. |
|--|---|

VCAP-DCD

8/24/2011

Mod 5 vStorage:

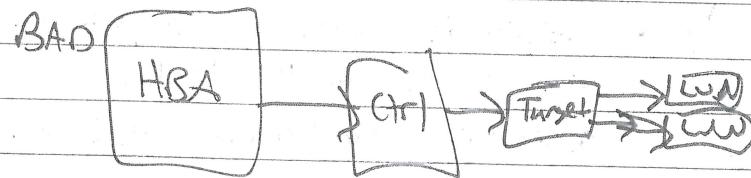
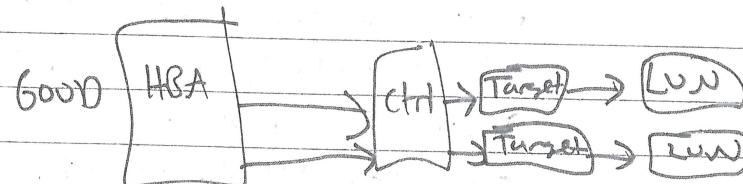
Software vs Hardware initiators.

- Software is fine.
- do not mix both types in the same host.
- Hw iSCSI uses less host CPU. About same performance.

One LUN per target is best for performance.

Initiator establishes a single connection to a target.

One target per LUN.



Configure ports consistently across hosts, switches, and arrays.

↳ consistent HBA + easier troubleshooting.

* Don't mix HBAs from different vendors. Different model from same vendor is ok.

Boot from SAN:

FC & Hw iSCSI only.

- can increase complexity, adds add dependency on SAN.
- HBA smaller.

Servers can be denser. Run cooler, NO disks to configure or install. Hosts less expensive.

VCAP-DCD

8/24/2011.

Boot from SAN.

- Exclusively must boot LUN to host.
 - use single initiator zoning.
 - create single 1600MB Disk partition
- BFS supported in ESX 4.0, ESXi 4.1.

NPIV:

NPIV assigns a WWN to a VM. Gives a VM an identity on the SAN.

- track storage traffic per VM.
 - Zone + Map LWNs per VM.
 - SAN QoS per VM.
 - Improve IO to VM via array caching.
- Configure if there is a mgmt requirement to monitor SAN LUN usage @ VM level.
- a security requirement that mandates zoning a LUN directly to a VM.

Requirements:

- VMs use RAMs.
 - FC HBA must support NPIV.
 - FC Switches must support NPIV.
 - ESX/ESXi must have access to all LUNs used by their V
- NPIV cannot be used on VMs protected w/ FT.

VCAP-DCD

8/24/2011

Module 7: Virtual Machine Design

Virtual Machine CPUs.

How many? Default to 1. Scale up if there is a need.

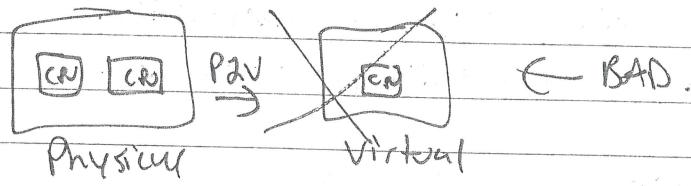
If using SMP VMs, make sure apps are multi-threaded so they can make use of 2x or 4x or 8x CPUs.

If the App is not multi-threaded, use multiple single vCPU VMs.

Use as few vCPUs as possible.

~~✓~~ # of vCPUs in a VM cannot exceed the core/hyperthread count in a host
it must be possible for all vCPUs to be scheduled.

For P2V's, create VM hardware that matches the physical machine's hardware abstraction layer.



Virtual Machine Memory:

- keep all guest active memory in physical RAM.
 - limit host memory overcommitment, or set reservations, or both.
 - Reservations should be set slightly higher than avg active memory size (for overhead)

Better to design around not overcommitting memory

~~✓~~ Always have Transparent Page Sharing Enabled.

Always load VMware tools and always enable ballooning.

VCAP-DCD

8/24/2011

Module 7 : Virtual Machine Design.

Shares, Reservations, Limits.

→ ✘ Deploy VMs w/ default settings unless clear reason to do otherwise.

? Are there Apps that NEED resources even during contention?
↳ Use reservations.

increases complexity + admin overhead.

VM DISKS:

Replay a system and data disk.

☒ Don't place ^{all} system disks on one datastore and data disks on another.

☒ Put system + data disks on the same datastore, unless they have widely varying IO characteristics.

☒ Configure one partition per virtual disk.

Simplifies SRM + snapshots.

↳ exception is if one VDisk is very large.

Swap location.

1st: Default, shared storage w/ VM files.

↳ con is that: vswap files are replicated w/ SRM.

2nd

Local storage: Reduces replication time. Local swap should not affect vm performance. Slows vMotion migration.

3rd

Dedicated vswap SAN storage. Good vMotion performance, non-replicated LUN, bad is admin overhead.

VCAP-DCD

8/24/2011

Module 7 : Virtual Machine Design.

Disk Adapters

MS Clustering Requires SAS Devices for Shared & Quorum disks.

Use Paravirtual SCSI adapters.

Cannot be used on boot disks in U.D.

PvSCSI not suited for local storage.

Snapshots can negate performance of PvSCSI.

Not supported w/ FT in 4.0.

NICs.

use vmnet3 where possible

use Standard Sizing for VMs

Ex:

item	small	medium	large	ex large
cpu	1	2	3	4
memory	1	2	4	8
disk	50	100	200	300

use vApps to guarantee power on order.

secure virtual machines as you would physical machines.

↳ AV, AS, ID, firewalls.

use vCenter server roles to limit access to VM console views.

Create VMs from templates that have been secured.

Add security products like vShield zones to the design.

VCAP - DCD

8/24/2011

Module 8: Management and Monitoring Design.

Design management around ESXi, even if using vSphere.

Limit the use of service console CLI and agents.

Deploy vCLI, PowerCLI, vMA

↳ install into VMs protected by HT.

Use ESXi embedded where possible.

Always automate! Use host profiles!

Create directory service users.

Create specific windows groups. Remove windows Administrators group from the vCenter Administrator role and create a specific vCenter Admins group.

Grant roles to groups, not users.

Use Folders to assign roles to objects that require similar access.

In large environments, consider a small management cluster.

F

FF Use thinapp to package the vSphere Client w/ all necessary plugins
↳ put that thinapp on a network share.

Use NTP for all hosts + Mgmt tools. Native time sync in VMs. Disable VMware tools time sync.

VCAP-DCD

8/24/2011

Module 8. Mgmt + Monitoring Design.

Create snapshot management policies that govern their use.

Snapshots are not backup solutions. They are back-out solutions.

Start level 1 for longest term retention.

Level 2 for long term performance monitoring when device

stats are not required but you want more than basics.

Level 3 for short term monitoring after encountering a problem
w/ a device.

Level 4 same as 3.

Levels 3 and 4 should only be used for short time.

Use best practices + SLAs to determine what to monitor.

Logging → if org has its own logging (syslog) server that can handle vSphere, use it.

Otherwise, use VMA + vlogger.

4.0

VCAP - DCI

Config maximums for HA.

VMs per host (~~host~~^{or fewer} 8 hosts in cluster) in 4.0 = 100
 $4.0 \text{ UI} = 160$

VMs per host for clusters w/ greater than 8 hosts = 40

Max failover hosts = 4

DCIcenter, powered on VMs = 2000 (32-bit), 3000 inventory.
 3000 (64-bit), 4500 inventory.

100 hosts per datacenter.

320 VMs per host.

256 kW/floor.

4.1

→ 10,000.

→ 320 VMs per host.

4.0 GA: Provides or smaller - 100 VMs / host
 4.0 Update / 160 VMs / host

All 4.0 9 nodes or higher - 40 VMs / host

4.1 → 320 VMs per host.