

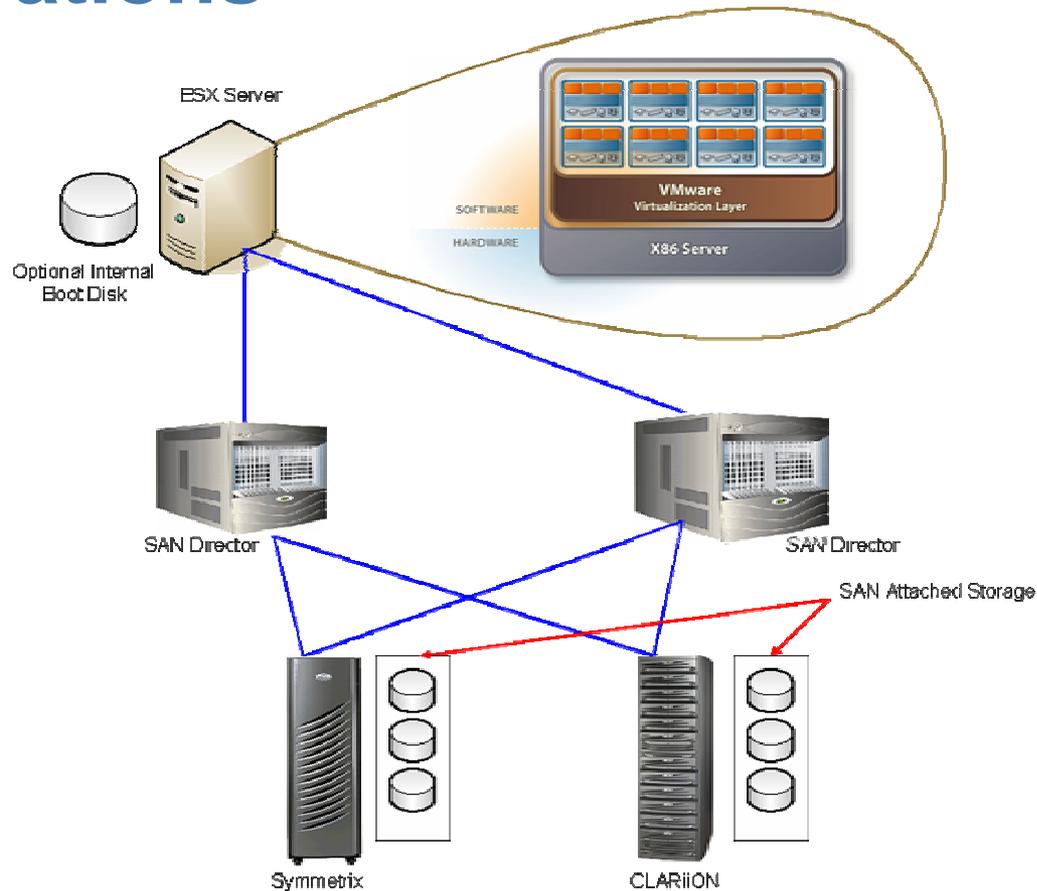
Storage Best Practices for VMware ESX Server 2.x

Bala Ganeshan
Corporate System Engineer
EMC Database Solutions Team

Agenda

- Introduction
- ESX Servers on SAN
- Booting ESX Servers of EMC Storage Arrays
- Using ESX Server v2.x on EMC Storage Arrays
- Using ESX Server v2.x with EMC Symmetrix Arrays
- Backup and recovery considerations
- Related sessions
- Question and answer session

SAN Attached ESX Server Configurations



Best Practices For Using EMC Storage With VMware ESX Server

- Recommendations to allow scalability
- Address potential performance limiting characteristics in the future
- May not be appropriate for all environments
- Recommendations enable optimal balance between performance, management and functionality

SAN Boot of ESX Servers

Booting ESX Server Off the SAN

- ESX Server 2.5 and later supports booting off the SAN
- Refer to “EMC Host Connectivity Guide for ESX Server 2.x” for configuration details
- ESX Servers can be booted off both Symmetrix and CLARiiON storage arrays
 - Advantages
 - Boot device protected on enterprise storage
 - x86 Servers become appliances (can be replaced in case of failures)
 - Leverage best practices for backing and restoring off the SAN

Booting ESX Server Off the SAN

- Disadvantages
 - Dedicated HBAs needed if RDM are to be used
 - Two different type of HBAs would be needed
 - Impacts VMkernel heap usage
 - ***Not supported by EMC***
 - Can impact management functions if HBAs are shared with service console
 - Potential performance impact if swap files are on the SAN disk
 - If internal disks are available, the swap file can be on protected internal disks.

Using ESX Servers with EMC Storage Arrays

Storage Layout Considerations

Storage Layout Considerations

- Do not present SAN storage to ESX Server farm as one big SCSI disk
 - Use multiple meta volumes or meta LUNs to present the storage
 - Plan on presenting storage as 8-10 meta volumes or meta LUNs
- Both meta volumes and meta LUNs can be grown non-disruptively
- Use of spanned VMFS with dynamic LUN growth mitigates risk

Storage Layout Considerations, cont.

	Storage as Single LUNs	Storage as Multiple LUNs
Management	<ul style="list-style-type: none"> ▪ Easier management ▪ Storage can be over provisioned ▪ One VMFS to manage 	<ul style="list-style-type: none"> ▪ Slightly harder management. ▪ Storage provisioning has to be on demand ▪ One VMFS to manage (spanned)
Performance	<ul style="list-style-type: none"> ▪ Can result in poor response time ▪ No manual load balancing 	<ul style="list-style-type: none"> ▪ Multiple queues ensure minimal response times ▪ Manual load balancing
Scalability	<ul style="list-style-type: none"> ▪ Limited # of virtual machines due to response time issue ▪ Limited # of IO intensive virtual machines since one VMFS 	<ul style="list-style-type: none"> ▪ Multiple VMFS allows more virtual machines per ESX Server ▪ Response time of limited concern (can optimize)
Functionality	<ul style="list-style-type: none"> ▪ All virtual machines share one LUN ▪ Cannot leverage ALL available storage functionality 	<ul style="list-style-type: none"> ▪ Use VMFS when storage functionality not needed ▪ Judicious use of RDM vs. VMFS

Storage Layout Considerations, cont.

- Use RDM (raw disk mapping) instead of VMFS if advanced storage functionality is desired
 - Required if using CLARiiON storage array application
- Use physical disk compatibility mode for RDMs
 - Enables VMotion while allowing virtual machines to recognize storage disk characteristics

Storage Layout Considerations, cont.

- If practical present multiple VMFS per ESX Server
 - Ideally 4-5 VMFS per ESX Server
 - Separate boot disk and application data
 - Separate log and data volumes
 - Separate classes of applications
- With ESX Server 2.x, VMFS can span multiple physical disks and/or partitions
 - Can be done using “vmkfstools” or MUI

Storage Layout Considerations, cont.

- Span VMFS across SAN volumes if needed
 - Advantages
 - Allows ESX Server administrator to still manage storage as one entity
 - Allows VMkernel to schedule multiple IOs to underlying disks
 - Provides flexibility to both ESX Server and storage administrator in assigning storage

Storage Layout Considerations, cont.

- Disadvantages:
 - Limits the number of IO intensive virtual machines on a ESX Server farm
 - Approximately 32 IO intensive virtual machines per VMFS
 - **Potential availability issue. Loss of one member of spanned VMFS will cause disruption and loss of all data on that VMFS**
 - Not a major issue with enterprise storage
 - Ensure VMFS does not span different storage class
 - Filesystem is spanned. IOs not balanced across VMFS “physical extents”

Using ESX Servers with EMC Storage Arrays

Path Management Considerations

Path Management Considerations

- PowerPath is not supported
 - If applicable, uninstall PowerPath from virtual machine after P2V is run
- ESX Server 2.x does not support active load balancing. Only path failover is available
- Set preferred path to enable some level of load balancing across HBAs
- Statistics will continue to be assigned on the first path discovered by Vmkernel
 - However, IOs will be directed to the active path

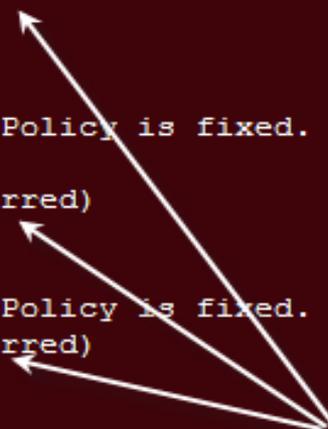
Path Management Considerations, cont.

```
root@losap152:~  
[root@losap152 root]# let i=1; while (( $i <= 3 )); do vmkmultipath -q vmhba1:0:  
$i; let i+=1; done  
Disk and multipath information follows:  
  
Disk vmhba1:0:1 (8,631 MB) has 2 paths. Policy is fixed.  
    vmhba1:0:1      on (active, preferred)  
    vmhba2:0:1      on  
Disk and multipath information follows:  
  
Disk vmhba1:0:2 (8,631 MB) has 2 paths. Policy is fixed.  
    vmhba1:0:2      on (active, preferred)  
    vmhba2:0:2      on  
Disk and multipath information follows:  
  
Disk vmhba1:0:3 (8,631 MB) has 2 paths. Policy is fixed.  
    vmhba1:0:3      on (active, preferred)  
    vmhba2:0:3      on  
[root@losap152 root]#
```

**All Active (and Preferred)
Path on vmhba1**

Path Management Considerations, cont.

```
root@losap152:~  
[root@losap152 root]# vmkmultipath -s vmhba1:0:2 -r vmhba2:0:2  
VMware::ExtHelpers::System[788]: '/sbin/fdisk' -l /dev/sdo  
[root@losap152 root]# r let  
let i=1; while (( $i <= 3 )); do vmkmultipath -q vmhba1:0:$i; let i+=1; done  
Disk and multipath information follows:  
  
Disk vmhba1:0:1 (8,631 MB) has 2 paths. Policy is fixed.  
    vmhba1:0:1      on (active, preferred)  
    vmhba2:0:1      on  
Disk and multipath information follows:  
  
Disk vmhba1:0:2 (8,631 MB) has 2 paths. Policy is fixed.  
    vmhba1:0:2      on  
    vmhba2:0:2      on (active, preferred)  
Disk and multipath information follows:  
  
Disk vmhba1:0:3 (8,631 MB) has 2 paths. Policy is fixed.  
    vmhba1:0:3      on (active, preferred)  
    vmhba2:0:3      on  
[root@losap152 root]#
```



Active (and Preferred) Paths distributed between HBA's

Using ESX Servers with EMC Storage Arrays

Partition Alignment

Alignment Issue on ESX Server

- ESX Server exhibits the same alignment problems seen on ALL Intel based platforms
- Legacy BIOS code from IBM PC
 - Used Cylinder, Head and Sector addressing instead of LBA addressing
 - Cylinder 0 – 1023, head 0-254, sector 1-63
 - Our reported geometry is 15 heads and 64 sectors
 - HBA BIOS code use mapping: e.g. 30 heads and 32 sectors

LBA	BIOS Mapping
Cyl 0, head 0, sec 1-64	Cyl 0, head 0, sec 1-63 AND cyl 0 head 1, sec 1
Cyl 0, head 1 , sec 1-64	Cyl 0, head 1, sec 2-63 AND cyl 0, head 2, sec 1-2
Cyl 0, head 2, sec 1-64	Cyl 0, head 2, sec 3-63 AND cyl 0, head 3, sec 1-3

Alignment Issue on ESX Server

- First track is reserved for boot code
- Means first partition starts at cyl 0 head 1 sector 1
- This is LBA 63 and is therefore, unaligned
- **As environments grow it can become performance limiting**
- To prevent this misalignment:
 - Offset the starting point of the partition to specified block
 - Create an aligned dummy partition
- Presentation shows how this can be done
 - Offset method on Linux
 - Creating dummy partition on Windows
- Either methodology can be used in both operating systems

Alignment Issue on ESX Server

- Default block size for VMFS is 1 MB
 - Use fdisk on service console to ensure virtual disks are aligned
 - Guest OS can misalign the IOs. Use diskpart or fdisk to align on the guest
- Use diskpart if RDM is presented to Windows as guest OS
- Use fdisk to align partitions on RDM presented to Linux guest OS

Using fdisk on Service Console to Align VMFS

- By default, ESX Server will create VMFS that are misaligned
- Execute the following steps to align VMFS
 1. On service console, execute “fdisk /dev/sd<x>”, where sd<x> is the device on which you would like to create the VMFS
 2. Type “n” to create a new partition
 3. Type “p” to create a primary partition
 4. Type “1” to create partition #1
 5. Select the defaults to use the complete disk
 6. Type “x” to get into expert mode
 7. Type “b” to specify the starting block for partitions
 8. Type “1” to select partition #1
 9. Type “128” to make partition #1 to align on 64KB boundary
 10. Type “r” to return to main menu
 11. Type “t” to change partition type

Continued on next page

Using fdisk on Service Console to Align VMFS

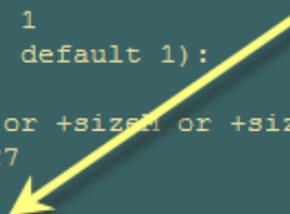
12. Type "1" to select partition 1
13. Type "fb" to set type to fb (VMFS volume)
14. Type "w" to write label and the partition information to disk

- You are now ready to create VMFS on the partition
- Use MUI or vmkfstools to create the VMFS
- The virtual disks on the VMFS will now be aligned
- However, guest OS will misalign IOs inside the virtual disk
- For Linux guest OS follow the procedure listed above
- Procedures for Windows hosts is presented after service console screen shots

Using fdisk on Service Console to Align VMFS

```
root@esx-node1:~  
[root@esx-node1 modules]# cd  
[root@esx-node1 root]# fdisk /dev/sdb  
  
The number of cylinders for this disk is set to 4427.  
There is nothing wrong with that, but this is larger than 1024,  
and could in certain setups cause problems with:  
1) software that runs at boot time (e.g., old versions of LILO)  
2) booting and partitioning software from other OSs  
   (e.g., DOS FDISK, OS/2 FDISK)  
  
Command (m for help): n  
Command action  
   e   extended  
   p   primary partition (1-4)  
p  
Partition number (1-4): 1  
First cylinder (1-4427, default 1):  
Using default value 1  
Last cylinder or +size or +sizeM or +sizeK (1-4427, default 4427):  
Using default value 4427  
  
Command (m for help): x  
  
Expert command (m for help): █
```

Expert mode (x option) is needed to align partitions



Using fdisk on Service Console to Align VMFS

```
root@esx-node1:~
Last cylinder or +size or +sizeM or +sizeK (1-4427, default 4427):
Using default value 4427

Command (m for help): x

Expert command (m for help): b
Partition number (1-4): 1
New beginning of data (63-71119754, default 63): 128

Expert command (m for help): p

Disk /dev/sdb: 255 heads, 63 sectors, 4427 cylinders

Nr AF Hd Sec Cyl Hd Sec Cyl Start Size ID
 1 00  1  1   0 254  63 1023   128 71119754 83
 2 00  0  0   0  0  0   0     0  0 00
 3 00  0  0   0  0  0   0     0  0 00
 4 00  0  0   0  0  0   0     0  0 00

Expert command (m for help): r

Command (m for help): t
Partition number (1-4): 1
Hex code (type L to list codes): █
```

Shows the starting sector as 63

This is the geometry of the disk as presented by the BIOS

Using fdisk on Service Console to Align VMFS

```
root@esx-node1:~  
3 00 0 0 0 0 0 0 0 0 00  
4 00 0 0 0 0 0 0 0 0 00  
  
Expert command (m for help): r  
  
Command (m for help): t  
Partition number (1-4): 1  
Hex code (type L to list codes): fb  
Changed system type of partition 1 to fb (Unknown)  
  
Command (m for help): p  
  
Disk /dev/sdb: 255 heads, 63 sectors, 4427 cylinders  
Units = cylinders of 16065 * 512 bytes  
  
   Device Boot      Start         End      Blocks   Id  System  
/dev/sdb1          1         4427    35559813+  fb  Unknown  
  
Command (m for help): w  
The partition table has been altered!  
  
Calling ioctl() to re-read partition table.  
Syncing disks.  
[root@esx-node1 root]#
```

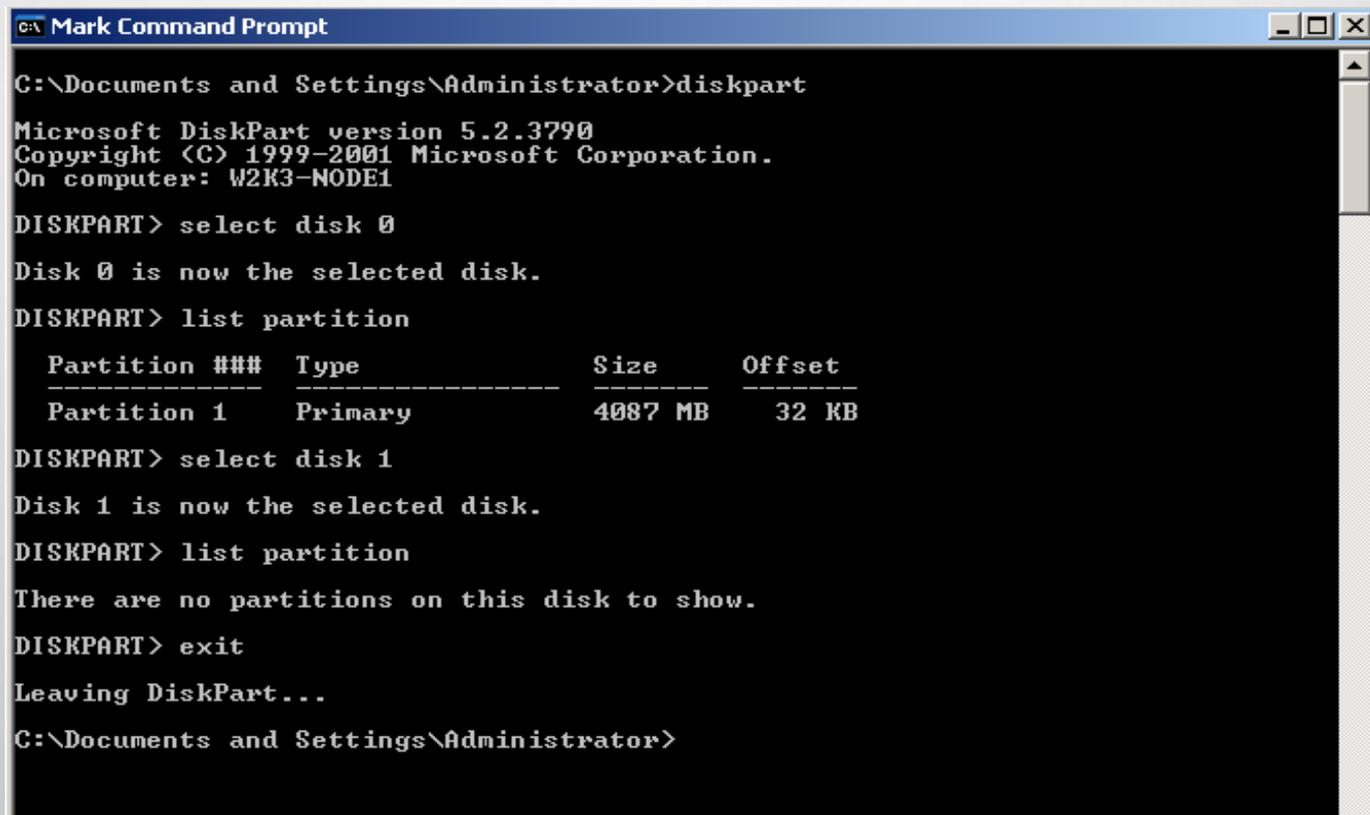
fdisk tells you that the size of the filesystem does not "align" with the geometry

Aligning Partitions on Windows Hosts (Guest OS)

- The procedure applies to both disks that are presented as System LUN/Disk and virtual disks (vmdk files on VMFS)
- Windows 2000 and later “lies” about alignment when using diskpart
 - Partitions are offset by 32256 and not 32768 bytes
- First Partition by default is at LBA 63. This misaligns IOs
- Use diskpar.exe (and not diskpart.exe) to align the usable partitions

Output From DiskPart

- Diskpart will show partitions as aligned. In reality they are not



```
C:\> Mark Command Prompt
C:\Documents and Settings\Administrator>diskpart
Microsoft DiskPart version 5.2.3790
Copyright (C) 1999-2001 Microsoft Corporation.
On computer: W2K3-NODE1

DISKPART> select disk 0
Disk 0 is now the selected disk.
DISKPART> list partition

   Partition ###  Type              Size              Offset
-----
   Partition 1    Primary           4087 MB           32 KB

DISKPART> select disk 1
Disk 1 is now the selected disk.
DISKPART> list partition
There are no partitions on this disk to show.
DISKPART> exit
Leaving DiskPart...
C:\Documents and Settings\Administrator>
```

How to Use Diskpart to Create Aligned Partitions

- Create a 1 MB “dummy” partition
- The partition is aligned on 32 KB boundary
- All subsequent partitions will be aligned

```
ca Command Prompt
C:\Documents and Settings\Administrator>diskpart -i 1
---- Drive 1 Geometry Information ----
Cylinders = 522
TracksPerCylinder = 255
SectorsPerTrack = 63
BytesPerSector = 512
DiskSize = 4293596160 (Bytes) = 4094 (MB)

End of partition information. Total existing partitions: 0

C:\Documents and Settings\Administrator>diskpart -s 1
Set partition can only be done on a raw drive.
You can use Disk Manager to delete all existing partitions
Are you sure drive 1 is a raw device without any partition? (Y/N) y

---- Drive 1 Geometry Information ----
Cylinders = 522
TracksPerCylinder = 255
SectorsPerTrack = 63
BytesPerSector = 512
DiskSize = 4293596160 (Bytes) = 4094 (MB)

We are going to set the new disk partition.
All data on this drive will be lost. continue (Y/N)? y

Please specify starting offset (in sectors): 64
Please specify partition length (in MB) (Max = 4094): 1

Done setting partition.
---- New Partition Information ----
StartingOffset = 32768
PartitionLength = 1048576
HiddenSectors = 64
PartitionNumber = 1
PartitionType = 7

You now should use Disk Manager to format this partition
C:\Documents and Settings\Administrator>
```

How to Use Diskpart to Create Aligned Partitions

- 1 MB Partition that was created in the previous step allows rest of the disk to have aligned IOs

The screenshot shows the Windows Computer Management console. The left pane shows the navigation tree with 'Disk Management' selected. The main pane displays a table of disks and their partitions.

Volume	Layout	Type	File System	Status
Boot Disk 1 (C:)	Partition	Basic	NTFS	Healthy

Disk	Layout	Type	Size	Health
Disk 0	Basic	3.99 GB	Online	
Disk 1	Basic	3.99 GB	Online	

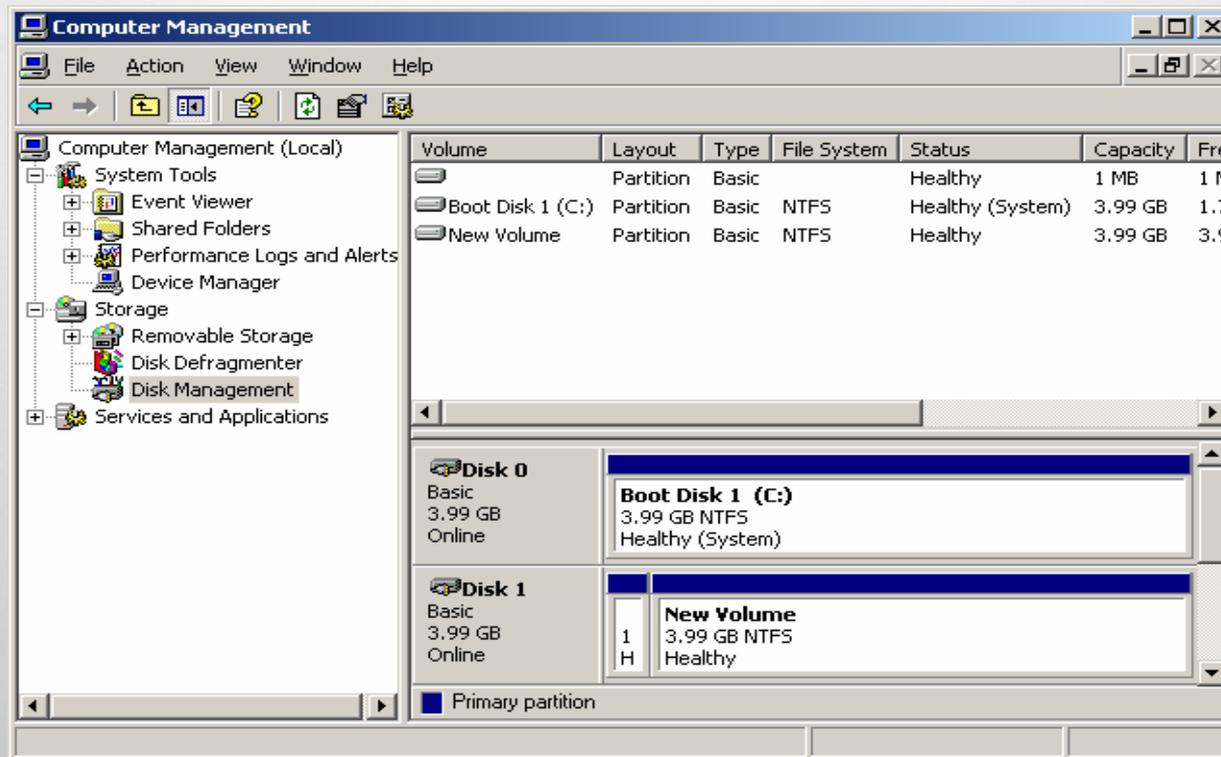
Partition	Layout	Type	Size	File System
1	H	3.99 GB	Unallocated	

Legend: ■ Unallocated ■ Primary partition

Overlaid on the right is the 'VMware, VMware Virtual S SCSI Disk Device Properties' dialog box. The 'General' tab is active, showing the device name, type, manufacturer, and location. The 'Device status' section indicates the device is working properly. The 'Device usage' dropdown is set to 'Use this device (enable)'.

How to Use Diskpart to Create Aligned Partitions

- Create partition and format NTFS with allocation unit of 32KB
- This will align all IOs on 32KB boundaries



Using ESX Servers With Symmetrix Arrays

Best Practices For Using EMC Symmetrix Arrays With VMware ESX Servers

- Use striped meta volumes
 - Follow EMC best practices when creating meta volumes
 - Use the standard split size if it exists
 - Ideally 4, 8 or 16 way meta volumes
 - Meta volumes should not have members on the same physical disks
 - Stripe size should be two cylinders
- For optimal performance and scalability meta volume should not be much larger than 200 GB
 - Balance between performance versus management
 - With 8 to 10 LUNs per ESX Server farm, this provides approximately 2 TB of storage
 - Can be grown as appropriate

Best Practices For Using EMC Symmetrix Arrays With VMware ESX Servers

- Use multi-target zoning if possible

```
root@L82AP128:~  
[root@L82AP128 root]# exitScript done, file is typescript  
vVMware::ExtHelpers::System[788]: '/sbin/fdisk' -l /dev/sdw  
[root@L82AP128 root]# vi typescript  
[root@L82AP128 root]# more /tmp/output.1  
Script started on Tue May 24 20:54:46 2005  
root@L82AP128 root]# vmkmultipath -q  
Disk and multipath information follows:  
  
Disk vmhba0:0:1 (9 MB) has 2 paths. Policy is fixed.  
    vmhba0:0:1      on (active, preferred)  
    vmhba1:0:1      on  
  
Disk vmhba0:0:41 (4.814 MB) has 4 paths. Policy is fixed.  
    vmhba0:0:41     on (active, preferred)  
    vmhba0:1:41     on  
    vmhba1:0:41     on  
    vmhba1:3:41     on
```

The server has two HBA's

but has 4 paths to this LUN

Best Practices For Using EMC Symmetrix Arrays With VMware ESX Servers

- Solutions Enabler 5.4 and higher is supported
 - Do not install Solutions Enabler under guest operating system unless RDM is used (no value add)
 - Installation of extraneous software on service console is NOT recommended. SE 6.0 will NOT run on service console
- Use GNS to enable propagation of information across multiple ESX Servers (farms) and guest operating systems

Best Practices For Using EMC Symmetrix Arrays With VMware ESX Servers

- Use of RDM is recommended if advanced functionality is needed
 - Leverage consistency technology
 - Storage functionality can be coordinated across multiple OS and applications
- Gatekeepers need to be provided and mapped to guest OS to leverage storage functionality

Backup and Recovery Considerations

Backup and Recovery Considerations

- ESX Server and Virtual Machines/Applications
 - Native tools provided by VMware
 - vmsnap, vmsnap_all, vmres etc.
 - Symmetrix Layered Applications
 - TimeFinder/Mirror, TimeFinder/Clones, TimeFinder/Snap, SRDF/S, SRDF/A and Open Replicator
 - CLARiiON Layered applications
 - SnapView, MirrorView, SAN Copy
 - CLARiiON Disk Library

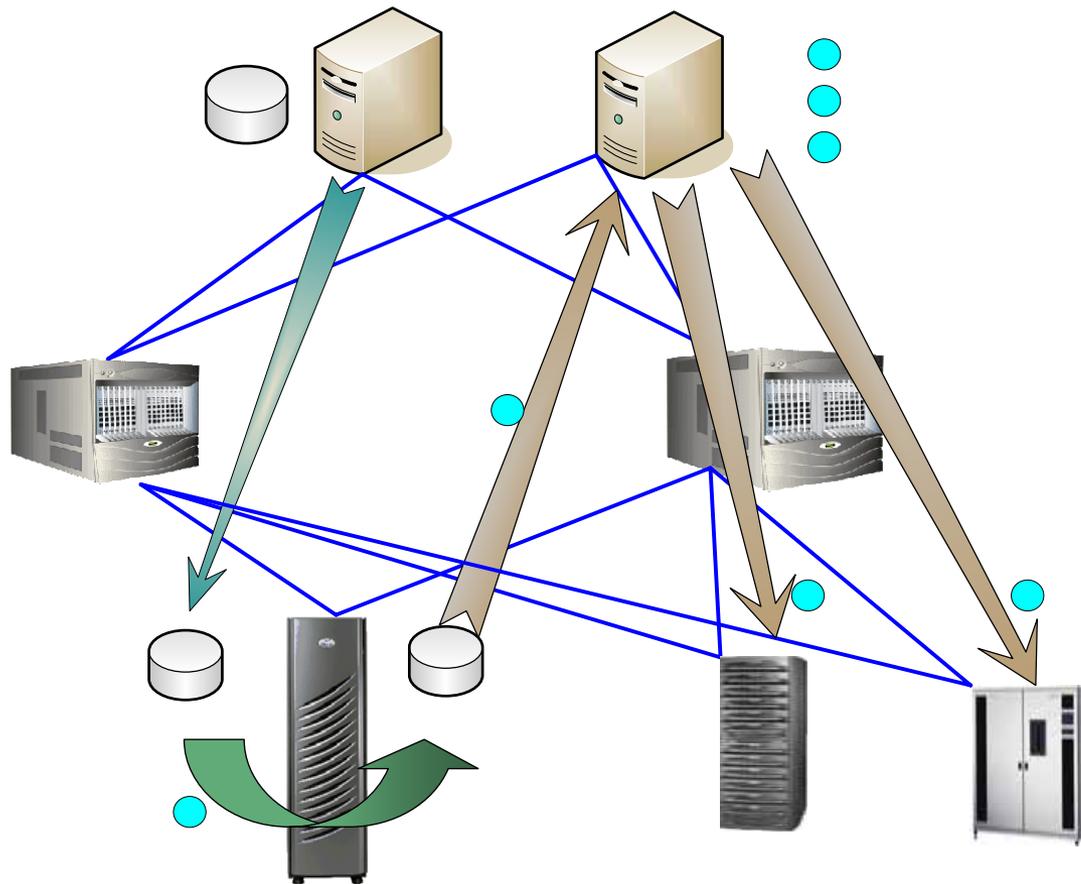
Backup and Recovery Considerations

- Backup Software
 - VMware tests major backup software. Refer to matrix http://www.vmware.com/pdf/esx2_backup_guide.pdf
 - Legato NetWorker has been tested.
 - “Storage Node” in Virtual Machines (Linux and Windows)
 - “NetWorker Client” in Virtual Machines
 - “Storage Node” or “NetWorker Client” for Linux on Service Console
 - EMC recommends configuring “Storage Node” or “NetWorker Client” in Virtual Machines
 - Provides granular backup and restores in this configuration
 - Add on modules (Oracle, Exchange etc.) can be used
- VMware add-on tools (vmware-mount.pl & vmware-loop) maybe needed for restores of individual files

Backup and Recovery Considerations

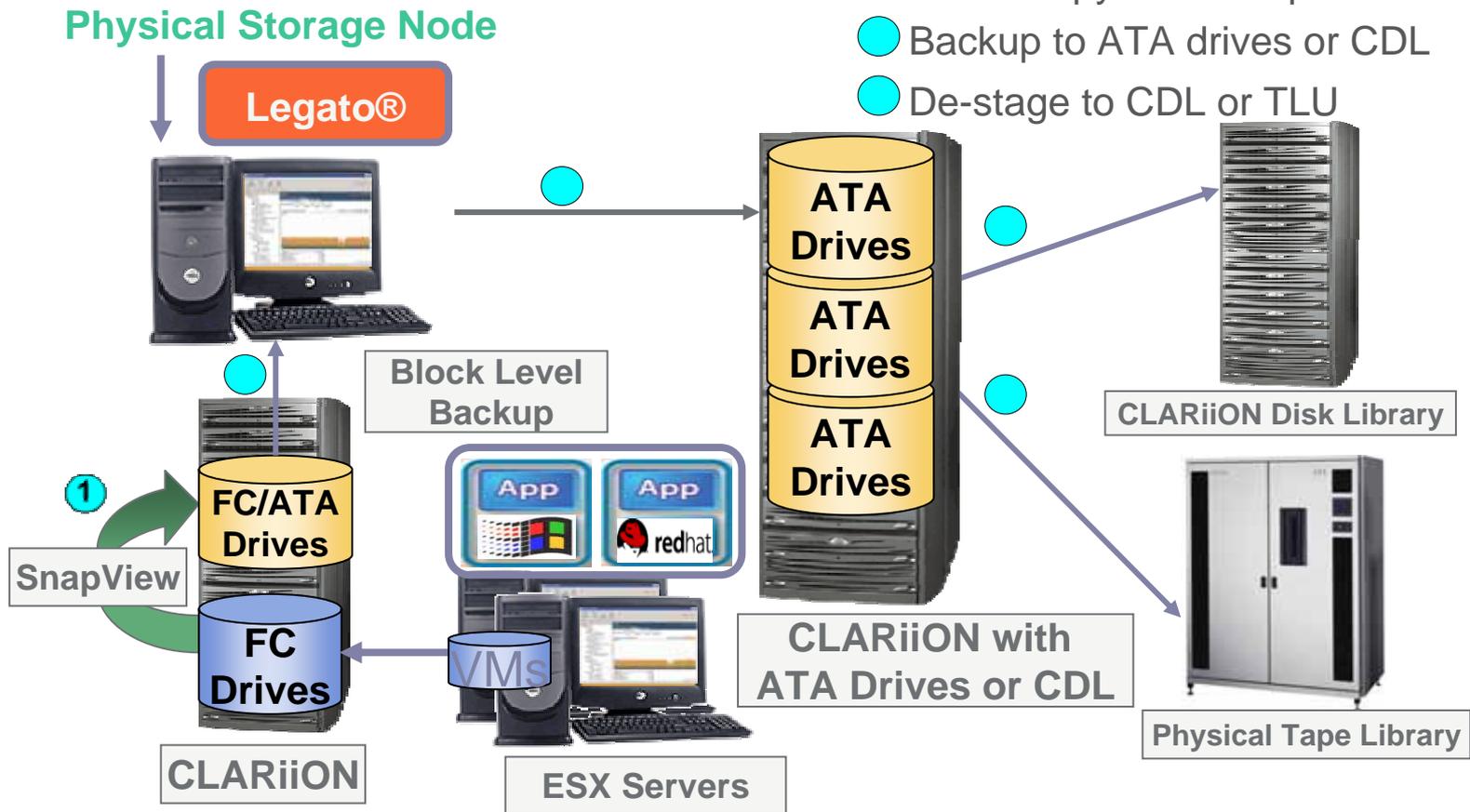
- On the Symmetrix side, extensive testing has been done with TimeFinder/M and SRDF/S
 - Refer to “Integrating VMware ESX Server with EMC Symmetrix Remote Data Facility” white paper
 - Solution Guides discussing integration of ESX Server with Symmetrix are also available
- All CLARiiON layered applications listed above are fully supported
 - Some of the restrictions are discussed in the previous slides
 - Refer to the white paper “CLARiiON Integration with VMware ESX Server” for further details

Backup and Restore Using TimeFinder



Backup using CLARiON Disk Library

- 1 Create copy of data
- Mount copy on backup server
- Backup to ATA drives or CDL
- De-stage to CDL or TLU



Summary

- Storage should be presented as multiple disks
- If using VMFS separate log and application data on separate file system
- Use static load balancing
- Use multi target zoning if possible
- Ensure partitions and file systems are aligned
- Follow best practices for creating metavolume and metaLUNs
- Use Symmetrix GNS if there are multiple ESX servers sharing storage
- Using CLARiiON layered applications require RDM
- When possible leverage layered applications to save precious CPU cycles

Related Sessions

- **SLN056** – VMware ESX Server Workload Analysis: How to Determine Good Candidates for Virtualization
- **SLN381** – VMware with CLARiiON

EMC²[®]

where information lives[®]

VMworld2005

virtualize^{now}

las vegas • october 18-20, 2005