

## **Storage Industry Trends and Storage Clustering Keynote Address IEEE/NASA Mass Storage Conference April 2002**

Matthew T. O'Keefe, Ph.D.  
Founder and Chief Technical Officer  
Sistina Software, Inc.  
Minneapolis, Minnesota  
[okeefe@sistina.com](mailto:okeefe@sistina.com)

# Agenda

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

- Merrill Lynch/McKinsey Enterprise Storage Study
- Storage Clustering — Definition and Advantages
- Storage Clustering Applications
  - Compute Clusters, Edge Serving, Parallel Databases
- Questions?

# The Storage Report – Customer Perspectives & Industry Evolution

## Merrill Lynch

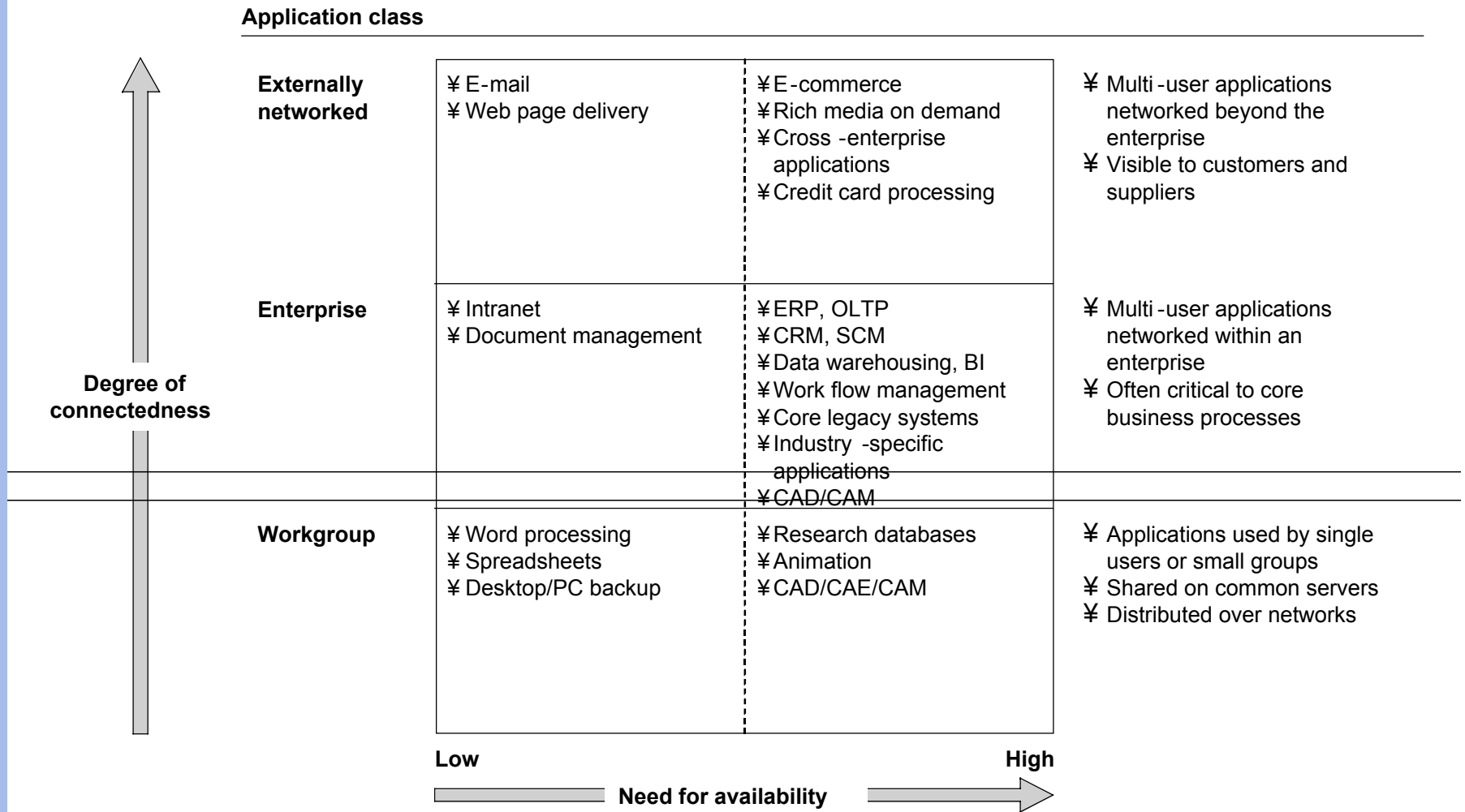
- Thomas Kraemer, James Berlino

## McKinsey & Company

- John Griffin, Doug Haynes, Thomas Herbig, Peter Stern, Alberto Torres
- Market surveys of several hundred storage industry customers on their storage needs
- In-depth interviews with storage customers across multiple industry sectors to understand capacity and performance needs, customer economics, purchasing processes, and decision criteria.

# Applications Classes

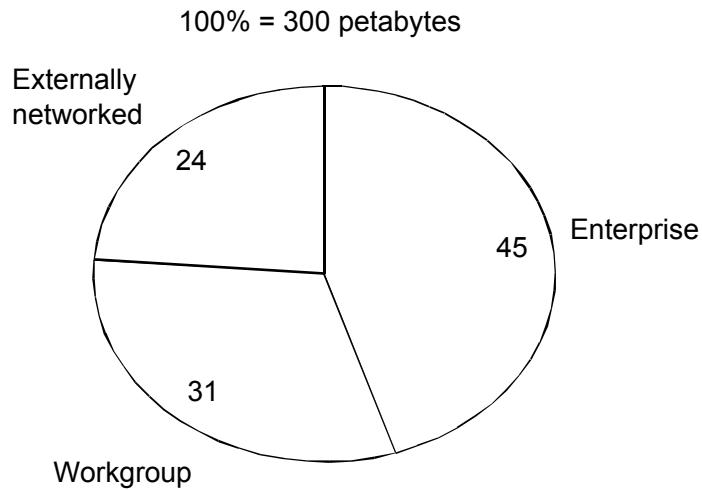
A NEW DIRECTION IN DATA STORAGE & MANAGEMENT



# Where is the Fastest Growth?

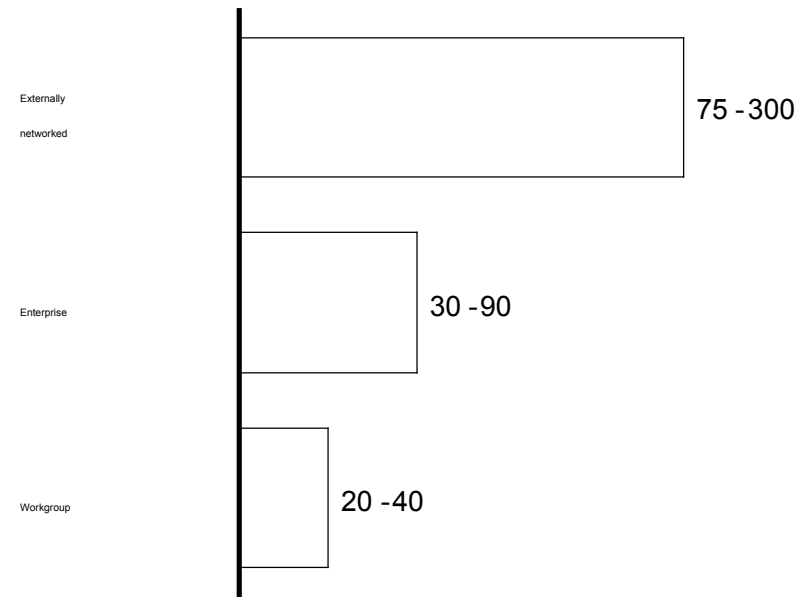
A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

**Breakdown of 2000 capacity shipped by application class**  
Percent



Source: Customer interviews; Merrill Lynch Reality Check survey of 110 C

**Annual capacity growth through 2003**  
Percent



IOs; IDC; Forrester; McKinsey and Merrill Lynch

# People costs, capacity, and NAS/SAN

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

When these increase

Storage people activities increase as follows in response\*

● Most activities increase  
○ Few activities increase

System component

Direct attached storage

Networked storage

Terabytes of capacity

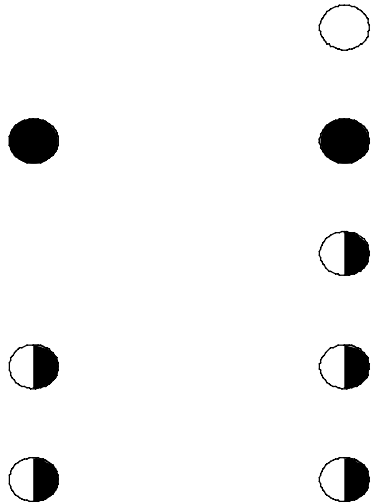
Storage subsystems

Servers

Applications

Server platforms

Capacity, subsystems, and servers are linked



In **DAS**, servers, subsystems, and capacity are linked, and increases in one often require increases in the other

**Networked storage** decouples these, and people costs no longer depend on raw capacity\*\*

\* Main activities are scaling, backing up, archiving, recovering, training, performance tuning. People costs often scale less steeply

\*\* When capacity growth exceeds 100% a year, they link terabytes to

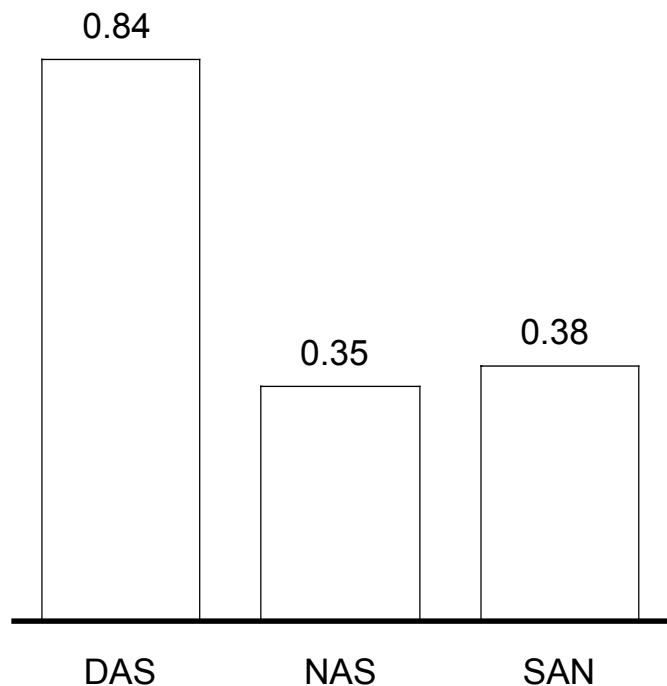
managing access, device and capacity monitoring, especially than system components due to economies of scale  
subsystems and again increase people costs

Source: Gartner; customer interviews; McKinsey and Merrill Lynch

# Network Storage has Better TCO than DAS

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

**3-year TCO by storage architecture\***  
\$ per megabyte of user data



¥ Cost saving of SAN and NAS driven by:

- Ĝ Improved disk utilization
- Ĝ Centralized management
- Ĝ Tape drive consolidation

¥ NAS does not have SAN network and installation charges

\* Based on 2TB of user data; as of March 2001

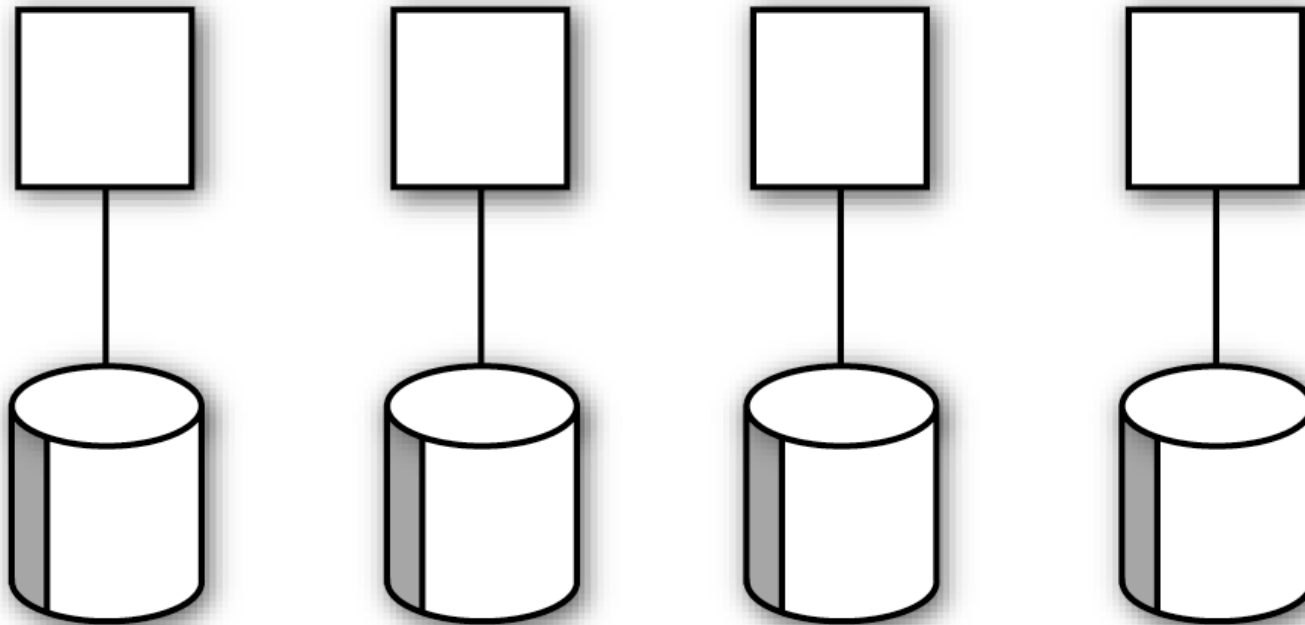
Source: Customer interviews; expert interviews; McKinsey and Merrill Lynch

# Local Disk Approach is Inefficient

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

## UNDERUTILIZATION

Inefficient space allocation resulting from local disk approach



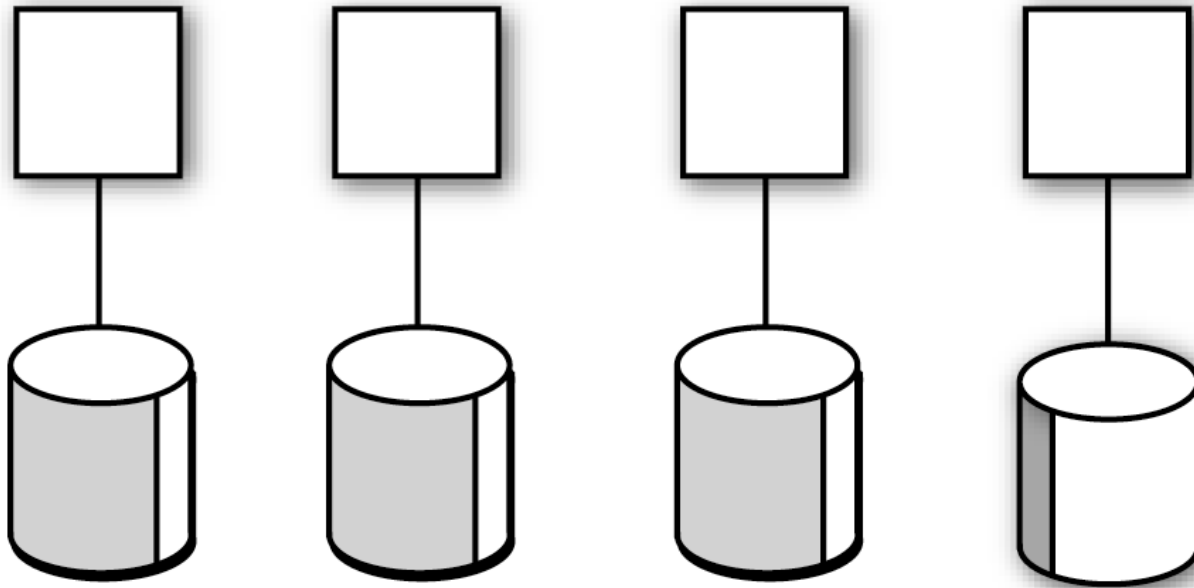
*Each local disk only 10% utilized*



# Lack of Balanced Usage Across Local Disks

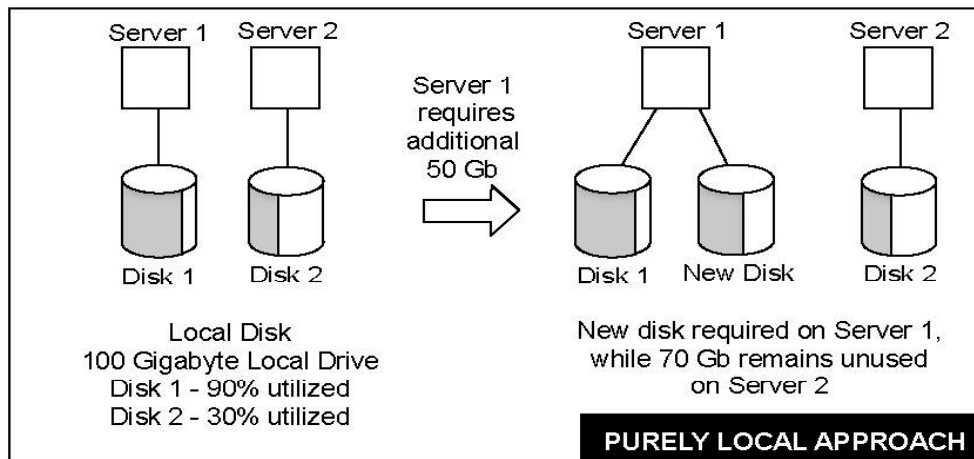
A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

## MIXED UTILIZATION RATIO

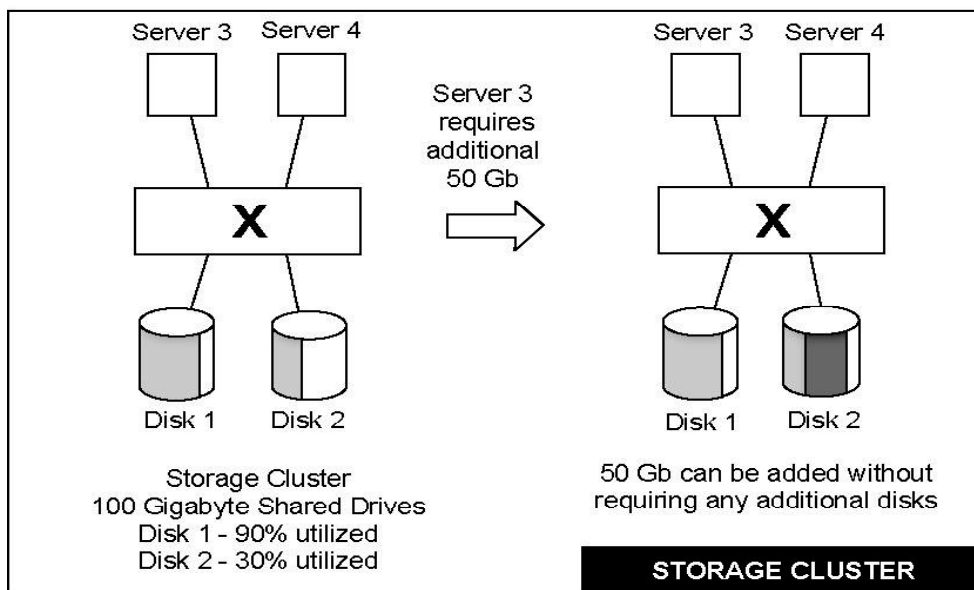


*Three local disks 90% utilized, one local disk 10% utilized*

**DYNAMIC SPACE UTILIZATION WITH STORAGE CLUSTERING  
VS.  
A PURELY LOCAL DISK APPROACH**

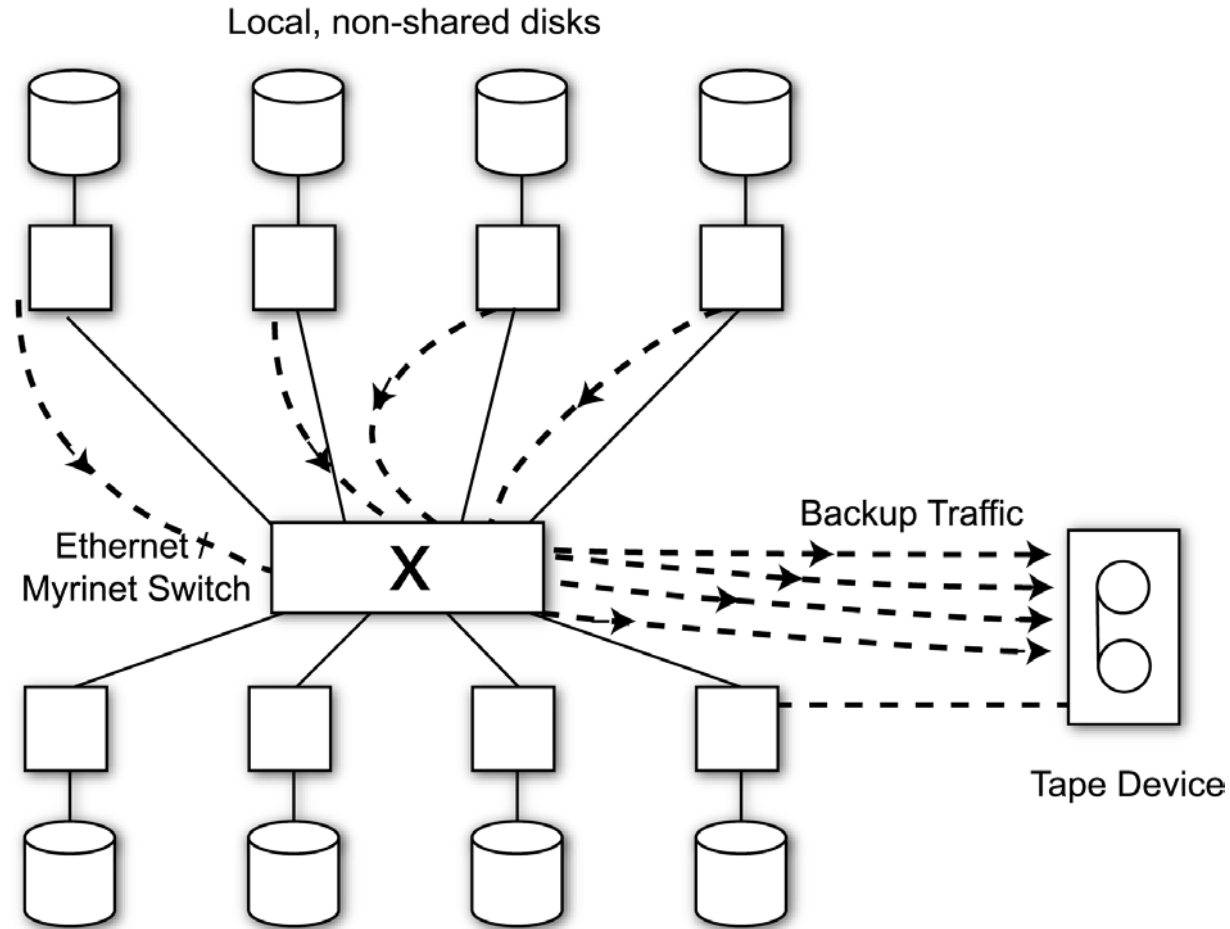


- Without shared storage, in this example a new server-attached disk is required
- With shared storage, no new disk is required, can efficiently utilize existing storage



# Backup Traffic Can Saturate a Compute Network

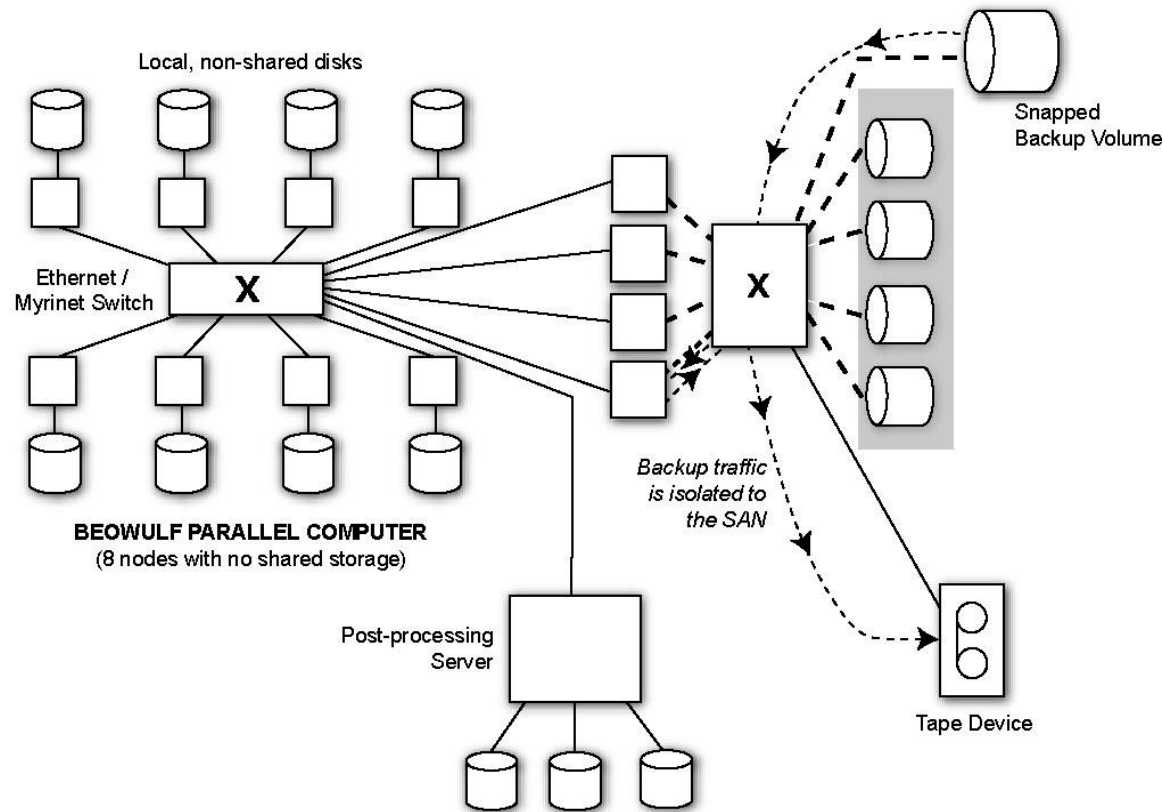
A NEW DIRECTION IN DATA STORAGE & MANAGEMENT



**BEOWULF PARALLEL COMPUTER**  
(8 nodes with no shared storage)

# Backup Traffic Off-Loaded to Storage Cluster

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT



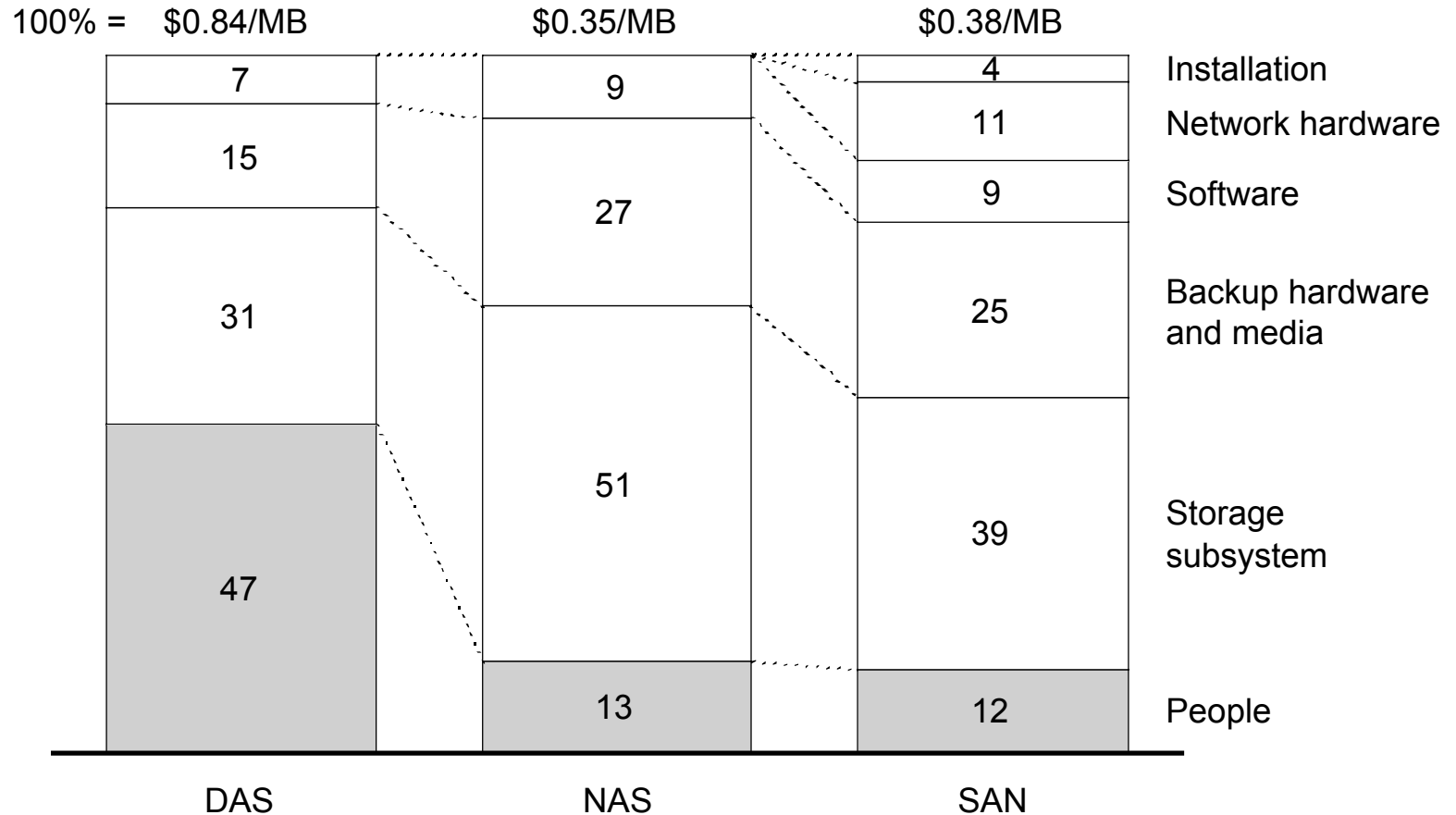
- Volume snapshot of GFS file system made first
- Read-only GFS mount then backed up across SAN
- Can be done on-line without disrupting computations
- Off-line storage consolidation

# Relative Cost Components

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

## Cost breakdown of storage architectures\*

Percent



\* Based on 2TB of user data and 10 servers; as of March 2001; see Source: Customer interviews; expert interviews; McKinsey and Merrill Lynch

Appendix E for assumptions

# Storage Decision Drivers

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

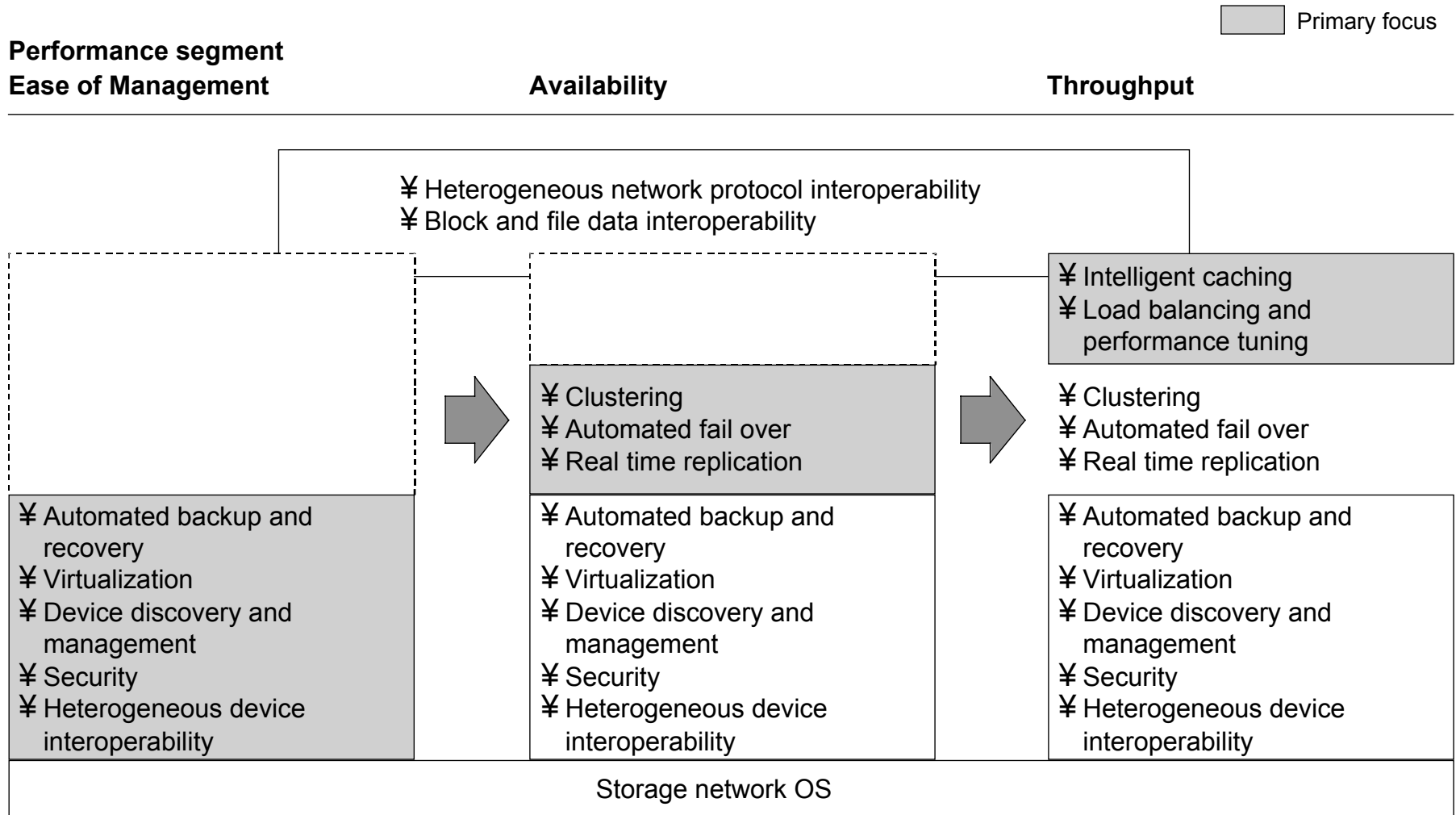
- Three performance characteristics are the primary drivers of storage decisions:
  - ease of management
  - availability
  - throughput
- Reliability and scalability are not the primary differentiators.
  - reliability is important, but today's storage systems are sufficiently reliable that this is no longer a differentiator

# Storage Decision Drives

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

- Scalability is no longer a major concern for most customers for two reasons:
  - improvements in storage density are accelerating more rapidly than the growth in demand for storage
  - networked storage greatly simplifies scaling
- Vendors should focus on the next set of customer challenges:
  - ease of management: NAS
  - availability: SAN
  - throughput: storage clustering (SAN + CFS + NAS)

# Requirements by Performance Segment

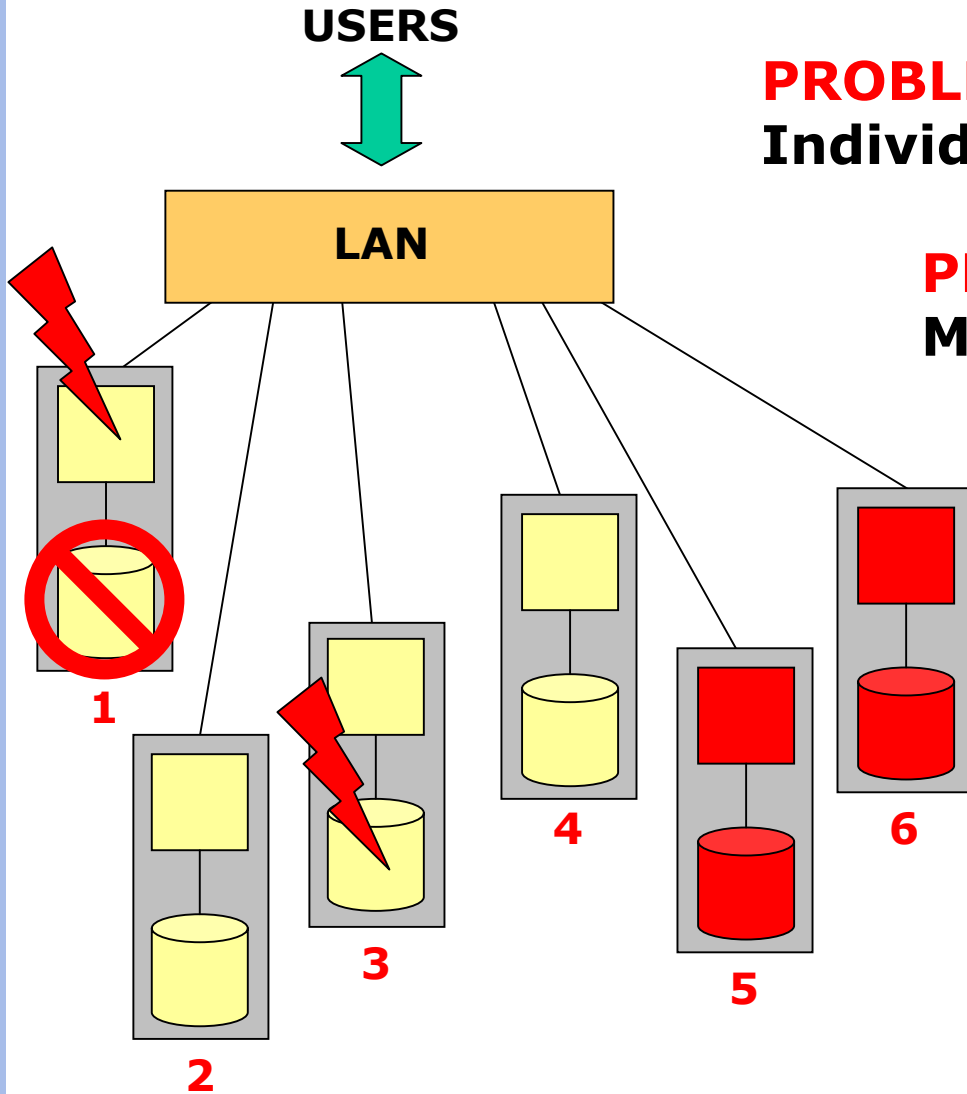


Source: Customer interviews; McKinsey and Merrill Lynch



# Today's Infrastructure

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT



**PROBLEM:** Servers Attach Individually to Disparate Disks

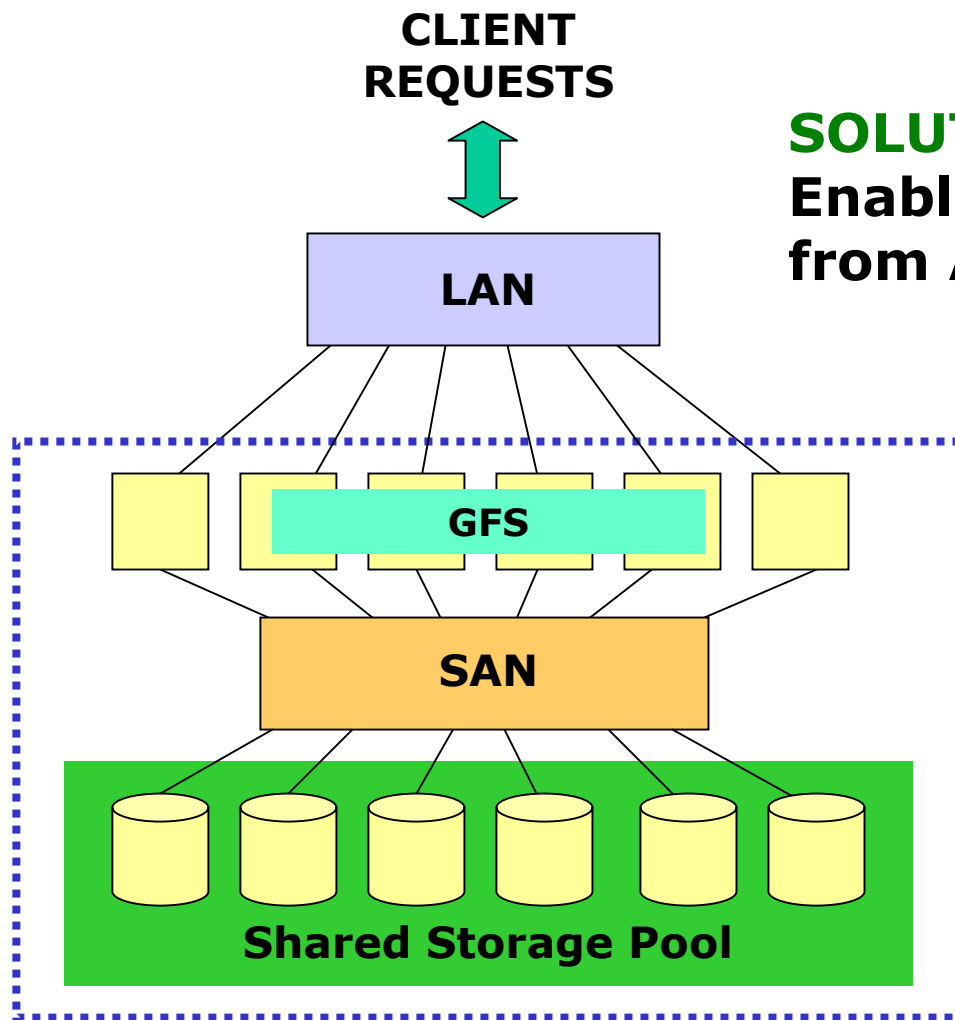
**PROBLEM:** Many Separate Management Domains

**PROBLEM:** Inefficient Scalability

**PROBLEM:** Data Accessibility

# Storage Clustering

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT



**SOLUTION:** Storage Clustering Enables a SAN to Create a Path from All Servers to All Disks...

**SOLUTION:** One Management Domain with Shared Storage

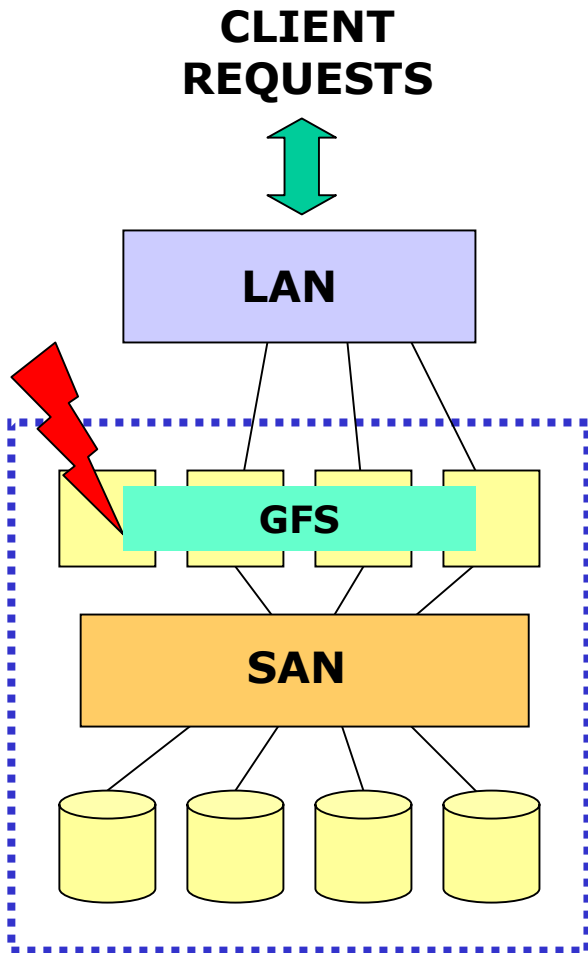
**SOLUTION:** Efficient Server Scalability

**SOLUTION:** Efficient Storage Scalability

1 Management Domain

# Storage Clustering

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT



Major Advantage:

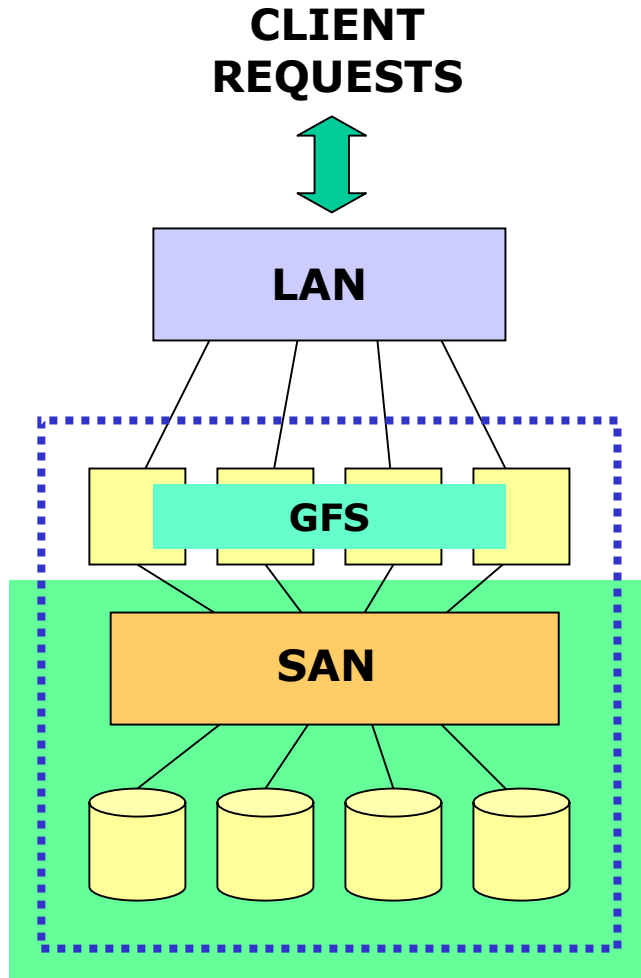
## Constant Accessibility

**In the case of a server failure the load is balanced across the remaining servers without access interruption.**

**The failed server can be brought back on-line while the system remains running.**

# Storage Clustering

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT



Major Advantage:

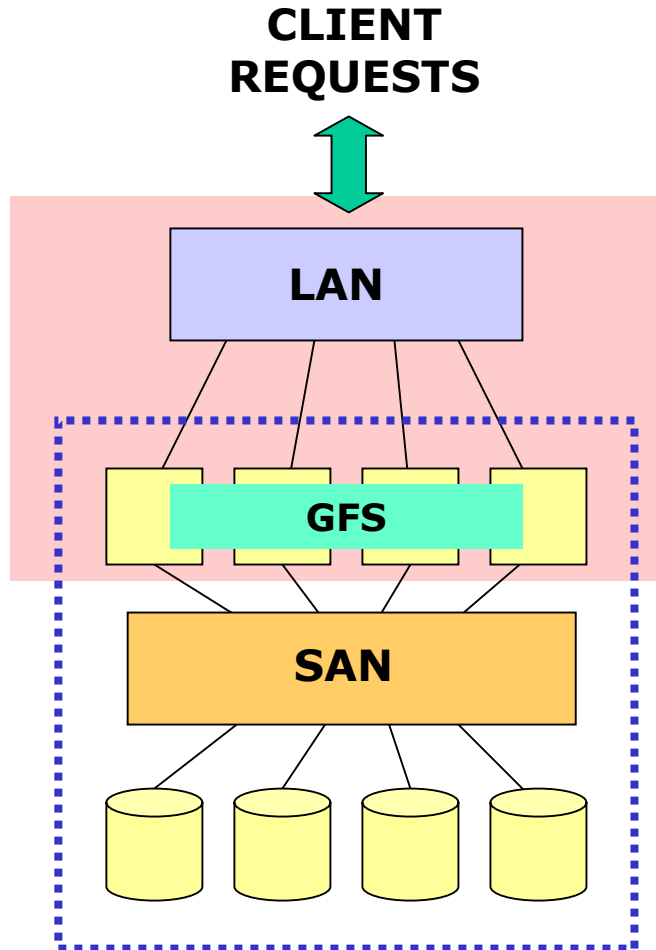
## Supports Storage Virtualization

Virtualization allows the aggregation and management of blocks.

A cluster file system completes the virtualization concept by enabling data sharing.

# Storage Clustering

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT



Major Advantage:

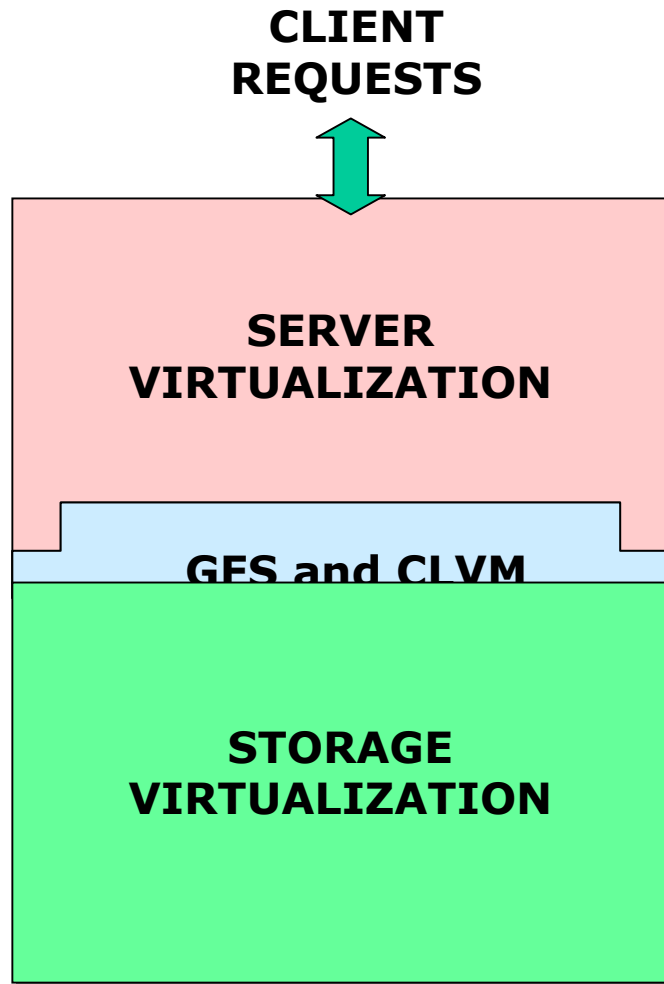
## Supports Server Virtualization

A cluster file system and SAN enables the evolution of clustered applications beyond pair-based, failover software.

Storage clustering technology reduces the expense and complication of HA software.

# A Cluster File System as Middleware

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

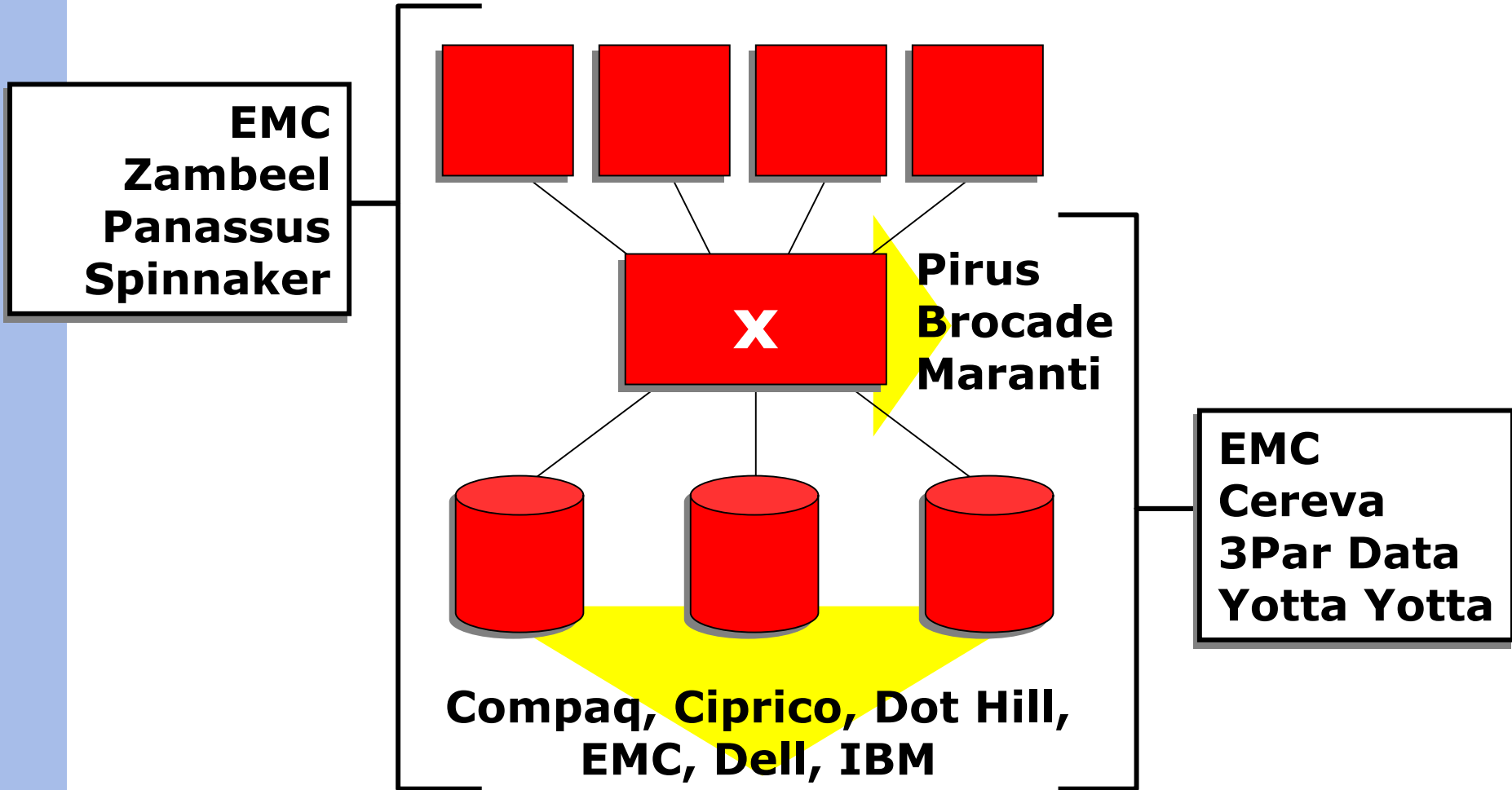


Major Advantage:

**A Cluster File System  
and Cluster LVM  
Enables Server and  
Storage Virtualization to  
Operate Together**

# Lots of Innovation

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT



# Storage Clustering Applications

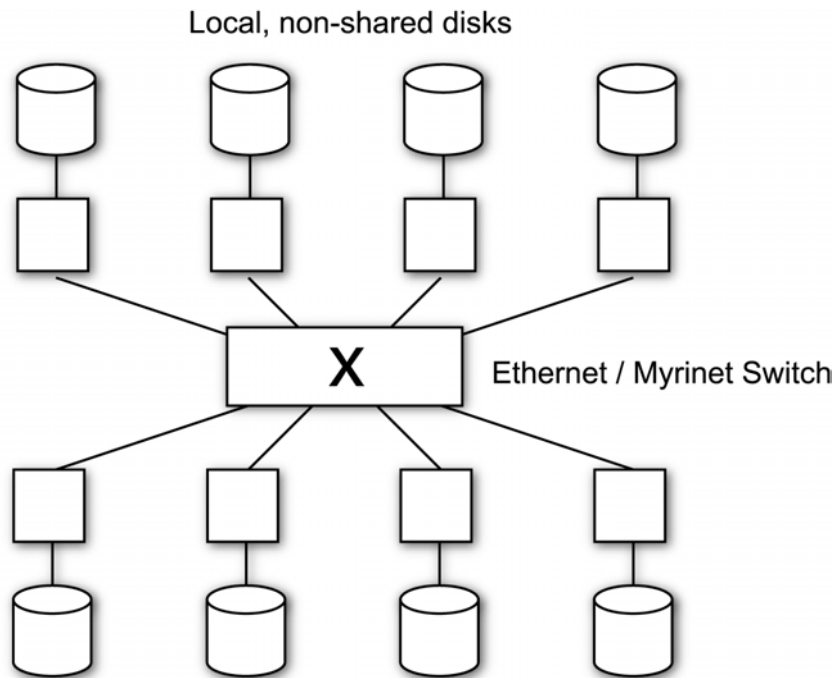
A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

- Compute Clusters
- Edge Serving
- Parallel Database Serving



# Compute Clusters

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

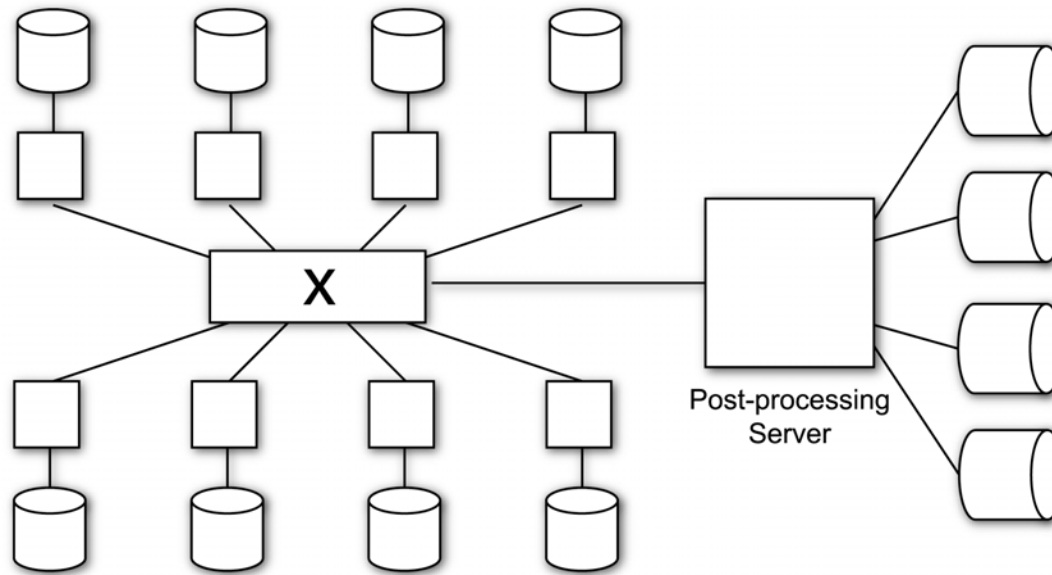


**BEOWULF PARALLEL COMPUTER**  
(8 nodes with no shared storage)

- Compute clusters becoming increasingly popular for HPC
- Basic concept: commodity processors, network, and disks make HPC affordable and scalable
- Hard to program, no fault-tolerance, but cheap

# Pre- and Post-Processing

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

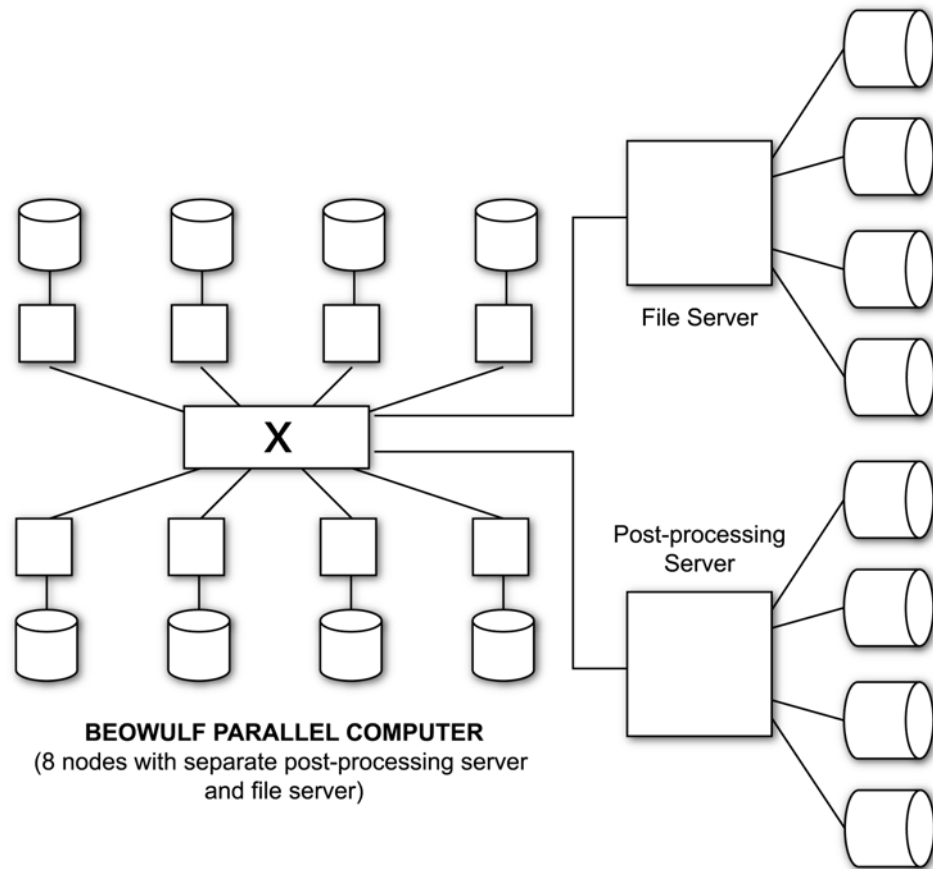


**BEOWULF PARALLEL COMPUTER**  
(8 nodes with separate post-processing server)

- Pre- and Post-processing servers often used to process computed data
- ftp or NFS used to get data from each compute node
- Or the data is just ignored

# Sharing with One File Server

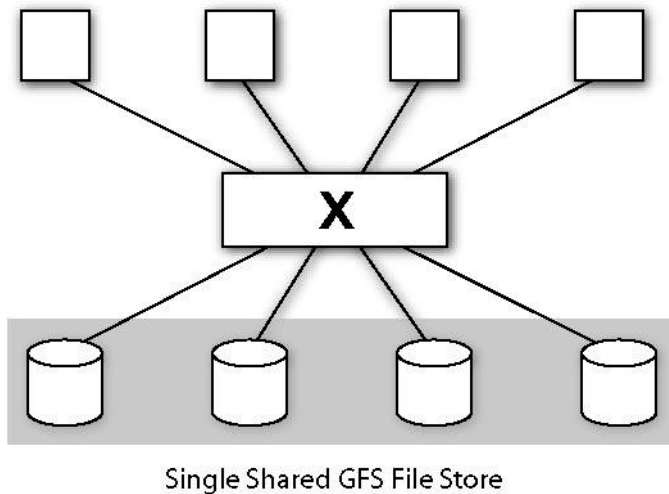
A NEW DIRECTION IN DATA STORAGE & MANAGEMENT



- NFS file server can be used to share data between compute cluster and outside processing nodes
- Results in replicated data: on both cluster and the file server
- Slow, does not scale

# GFS-Based Storage Clusters

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

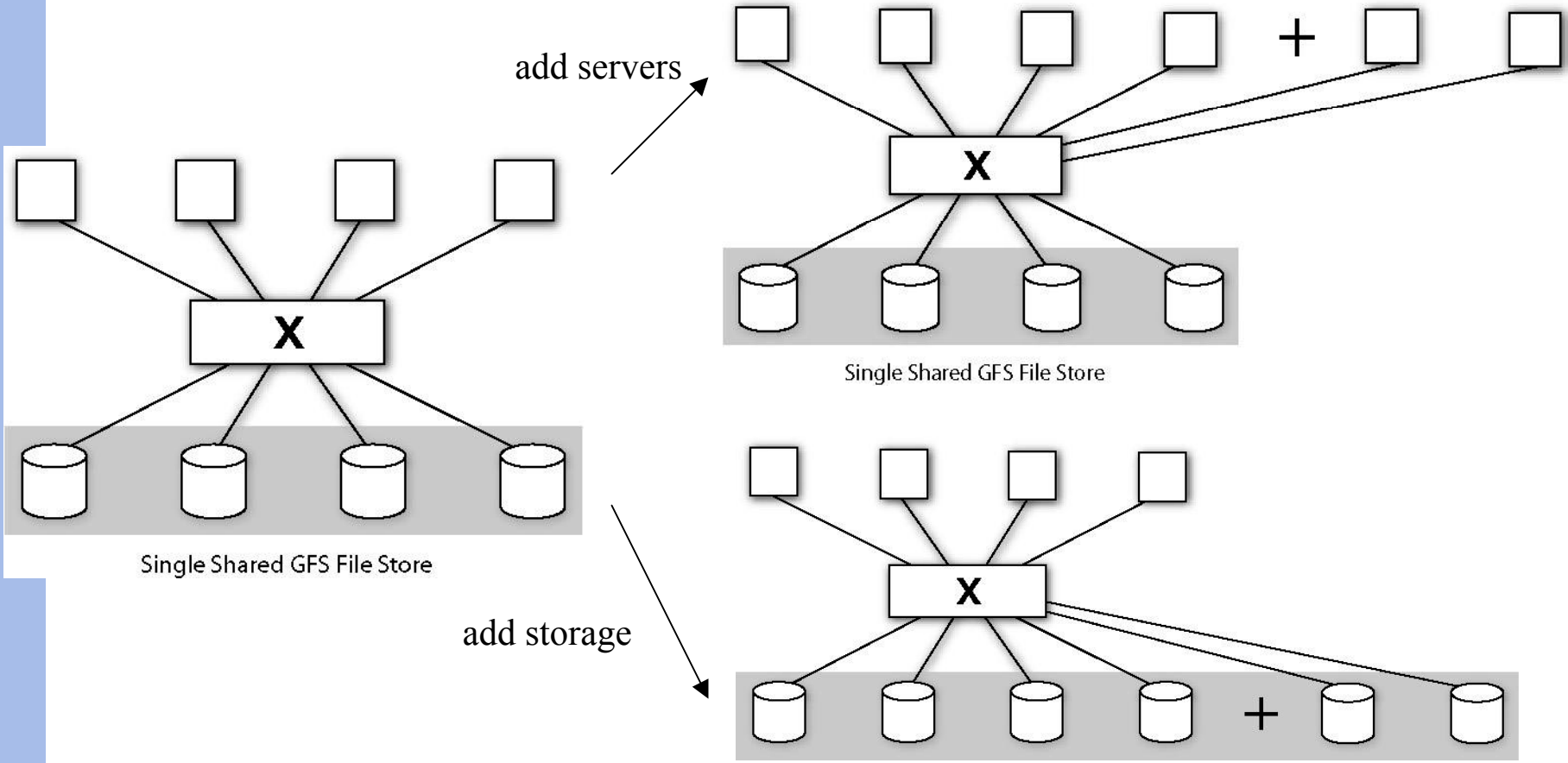


- Storage cluster consists of
  - storage area network
  - shared storage devices,
  - cluster file system/volume manager running on the servers
- Sistina's GFS and CLVM: cluster file system for Linux: can be used to build a Linux-based storage cluster

# Storage Cluster Scalability

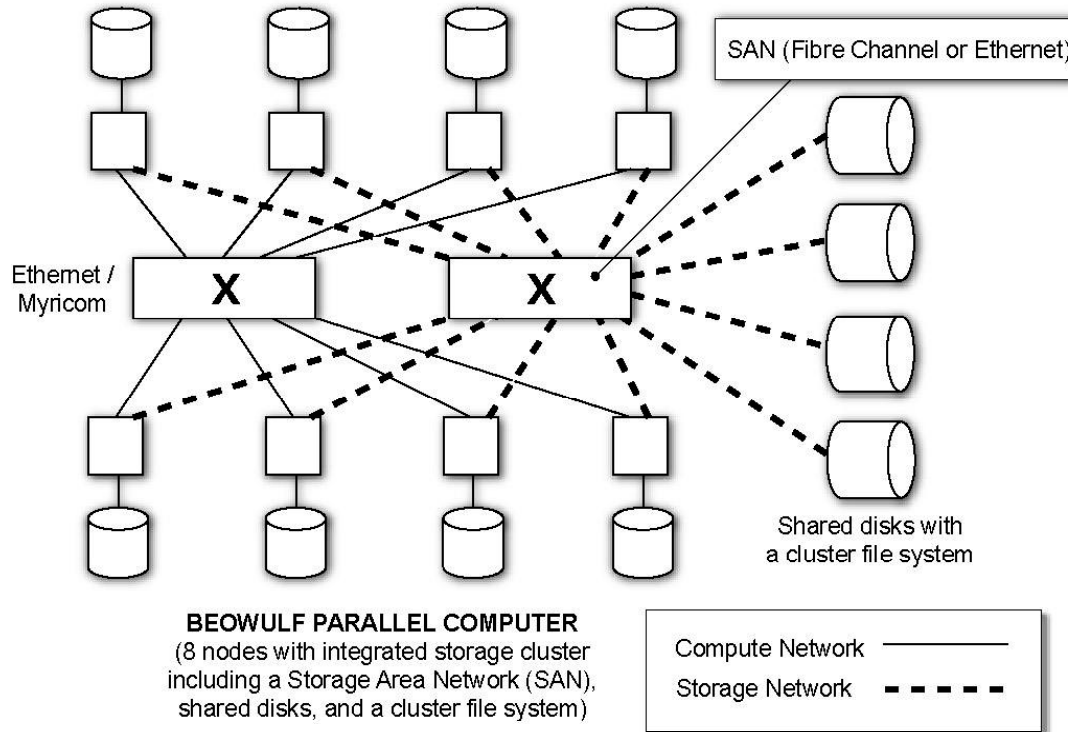
A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

- Scale either servers or storage



# Compute Cluster with Integrated Storage Cluster

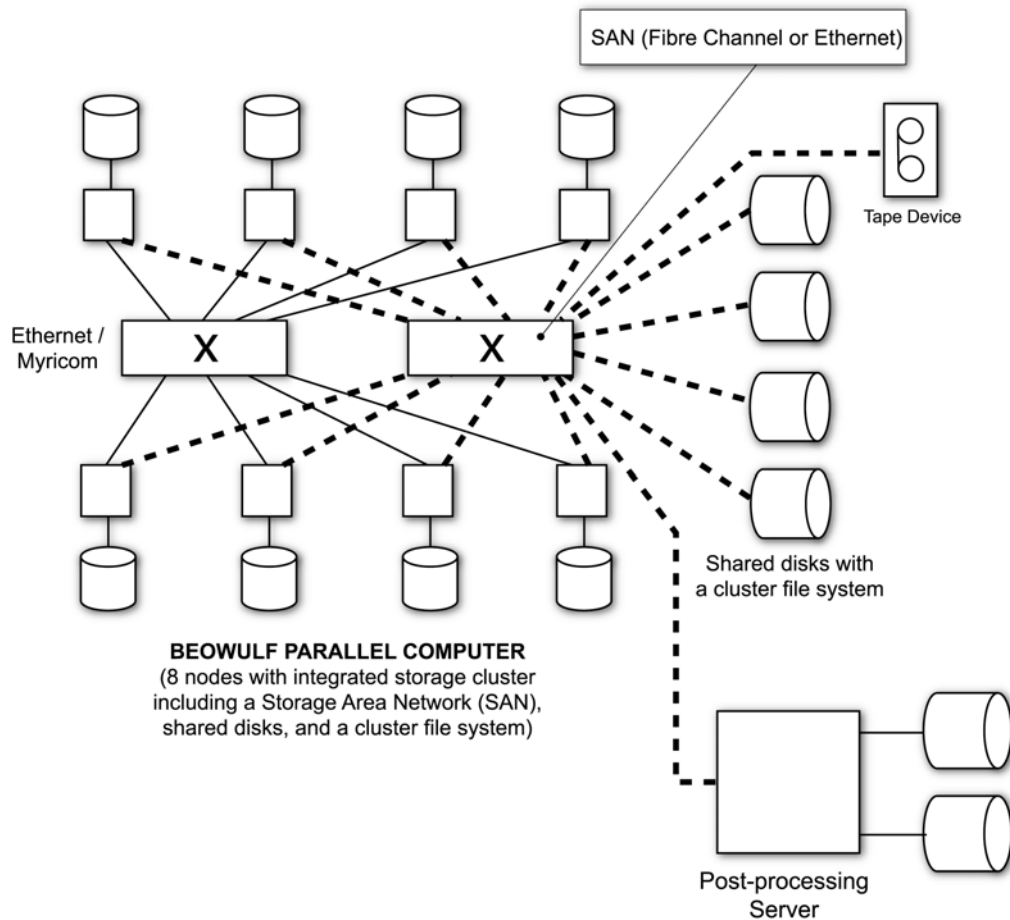
A NEW DIRECTION IN DATA STORAGE & MANAGEMENT



- Storage cluster: servers sharing storage with a cluster file system over a SAN
- Implies set of shared storage devices
- SAN with FC or IP

# Post-Processing Server Can Use Shared Storage

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT



- Pre- or Post-processing hardware can be attached to the storage cluster
- One image of data between the compute cluster and other servers
- Consolidates storage into single, easily managed pool — avoids replication

# GFS-Based Shared Storage

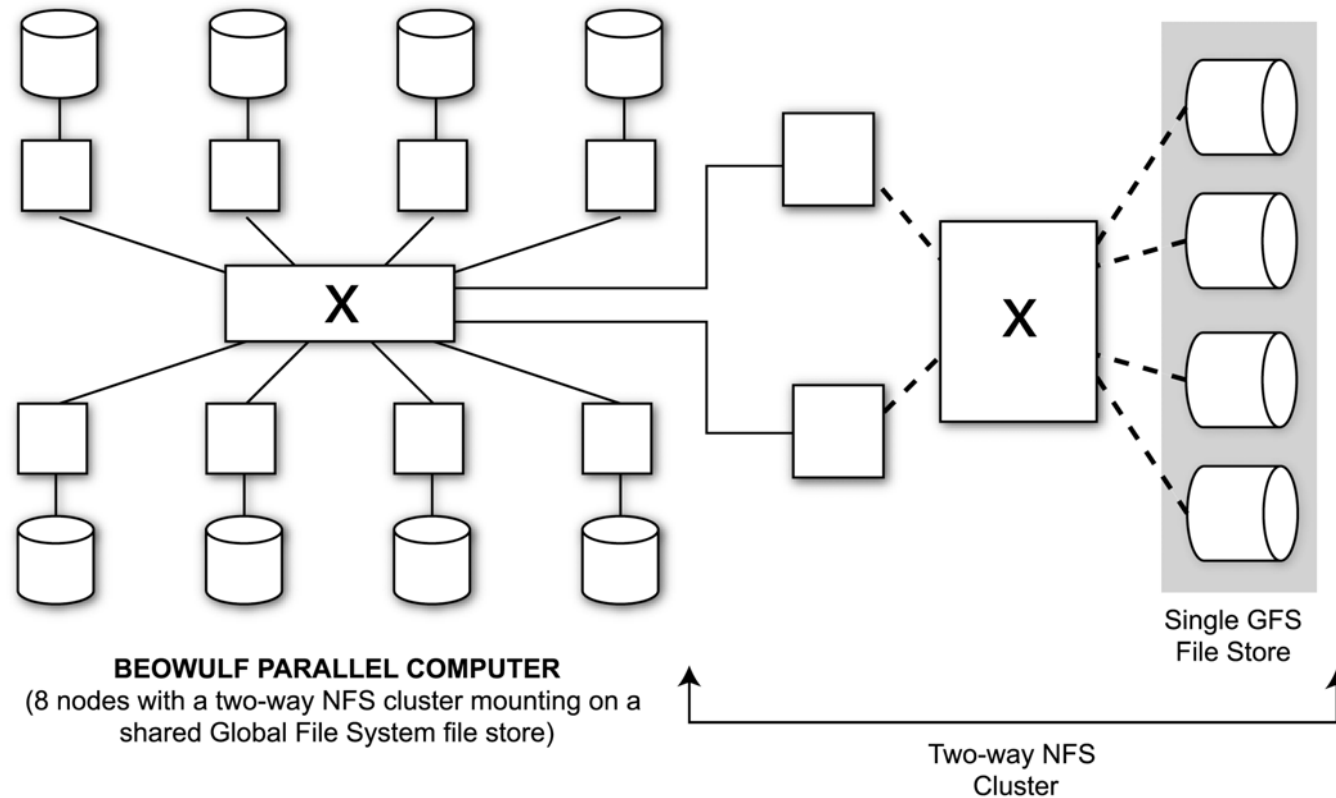
A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

- Shared storage is great, but FC-based SANs are still expensive relative to compute node cost
- Instead of a single large NFS server, system architects can use GFS to enable an NFS cluster to provide scalable storage service for a large compute cluster



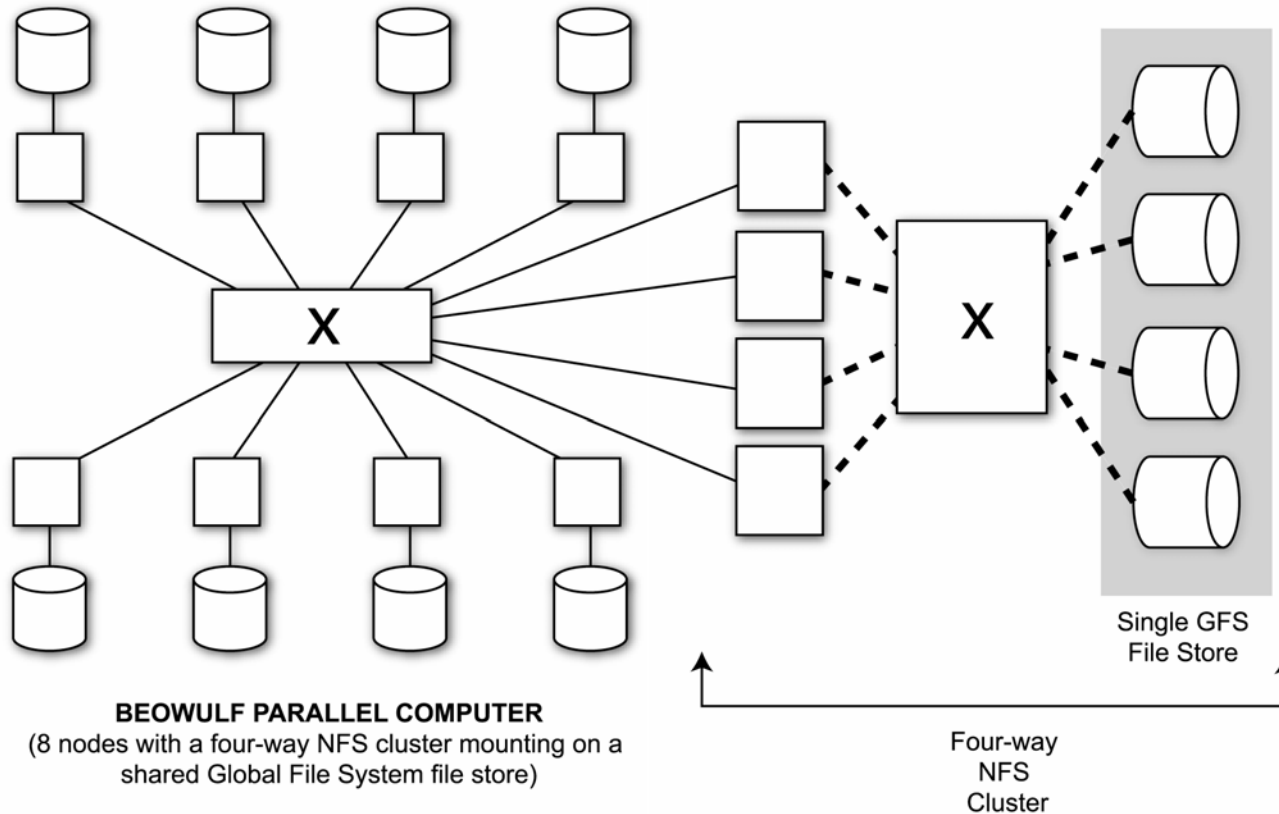
# NFS Clustering for Shared Storage: 2-way

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT



# NFS-Clustering for Shared Storage: 4-way

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT



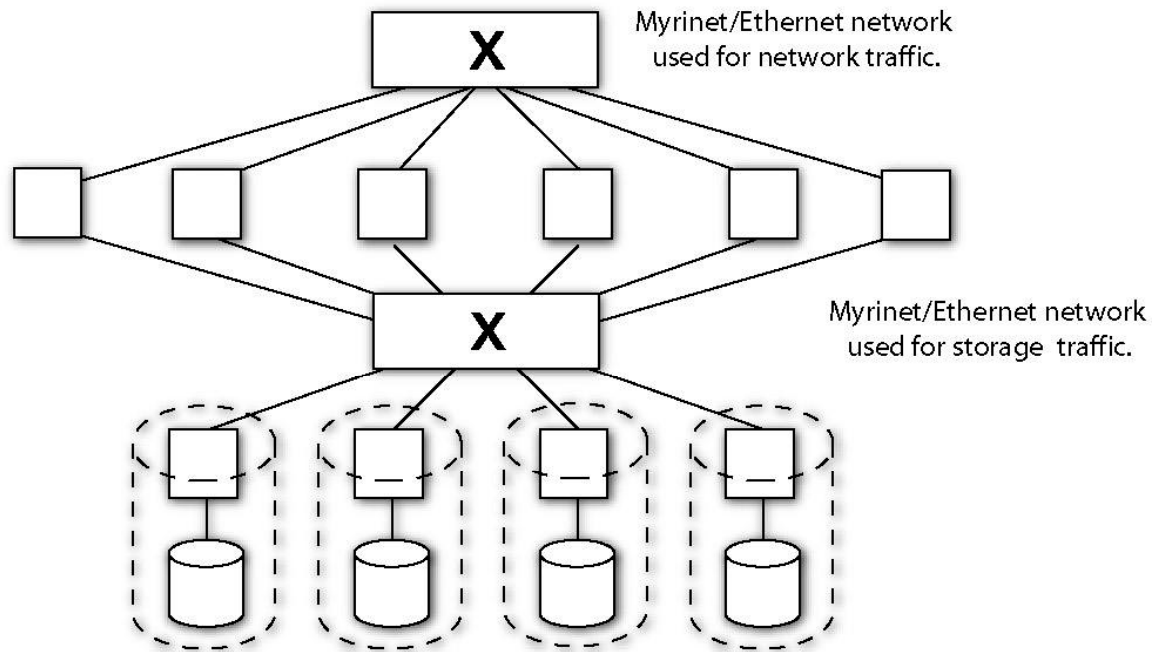
- Advantage: expensive SAN hardware not required for every node
- Disadvantage: not quite as scalable

**Table 1**  
**Storage Clustering Alternatives Comparison Matrix**

	<b>File Protocol</b>	<b>Hardware Cost</b>	<b>Management Cost</b>	<b>Performance IO</b>	<b>Scalability (Bandwidth &amp; Capacity)</b>	<b>Availability (Downtime &amp; Failover)</b>
<b>Fully-integrated Storage Cluster</b>	Local	Medium (IP) High (FC)	Low	High	High	High
<b>N-way NFS Cluster</b>	NFS	Medium	Medium	Medium - High	Medium	High
<b>Single Node Shared NFS Server</b>	NFS	Medium	Medium	Low	Low	Low
<b>No Shared Storage</b>	None	Low	Very High	Very Low	Very Low	Very Low

# Storage Clustering without Fibre Channel

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT



## SHARED VIRTUAL DISKS USING GNBD

6 compute nodes connected to 4 virtual disks  
over Ethernet or Myrinet (without Fibre Channel).

# Edge Serving — Definition

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

- At boundary between external Internet and internal data center operation, organizations typically deploy servers to implement a variety of client-server protocols
- These servers provide external client access to internal server data via web, email, file, or database protocols

# Edge Server Clustering Protocols

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

- Web Services
  - http and http proxy (e.g., squid cache)
- Mail Services
  - pop, imap, smtp, mta
- File Services
  - ftp, NFS, CIFS, Appletalk, Oracle IFS
- Client Login Services
  - ssh, telnet, etc.

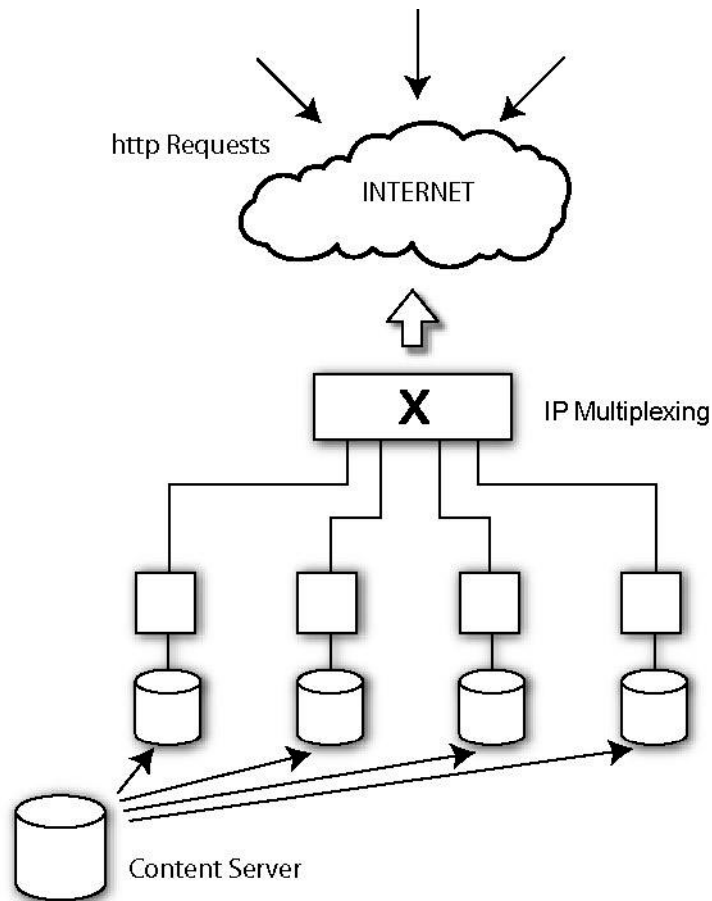
# Edge Server Farms

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

- A collection of edge servers will be referred to as an *edge server farm*
  - *multiple servers required to meet load requirements and to increase availability*
- Edge server farms are *homogeneous* if all servers in the farm are providing the same service or set of services

# Edge Server Farms

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

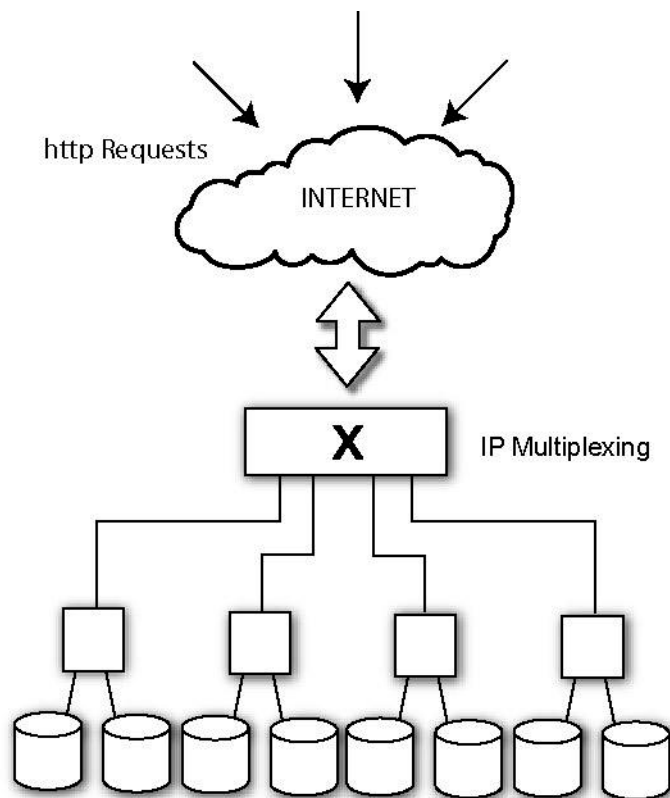


- Multiple servers provide content through IP load balancing switch to clients over Internet
- Content must be duplicated, no data sharing between servers, complex to manage and scale
- Scalable from hardware cost perspective, but complexity increases with each server added



# Edge Server Farms

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT



- If content and storage requirements double, storage on each node must be doubled
- Replicated content complex to manage and update
- Replicated content works only for simple, read-only protocols (http) not file serving (NFS,CIFS)

- Farms don't scale well from a management or efficiency standpoint (see Sun's latest ads about consolidation)
- Clustered, shared storage solves the problems of edge server farm management, scaling, performance, and availability
- GFS and CLVM provide clustered storage for Linux

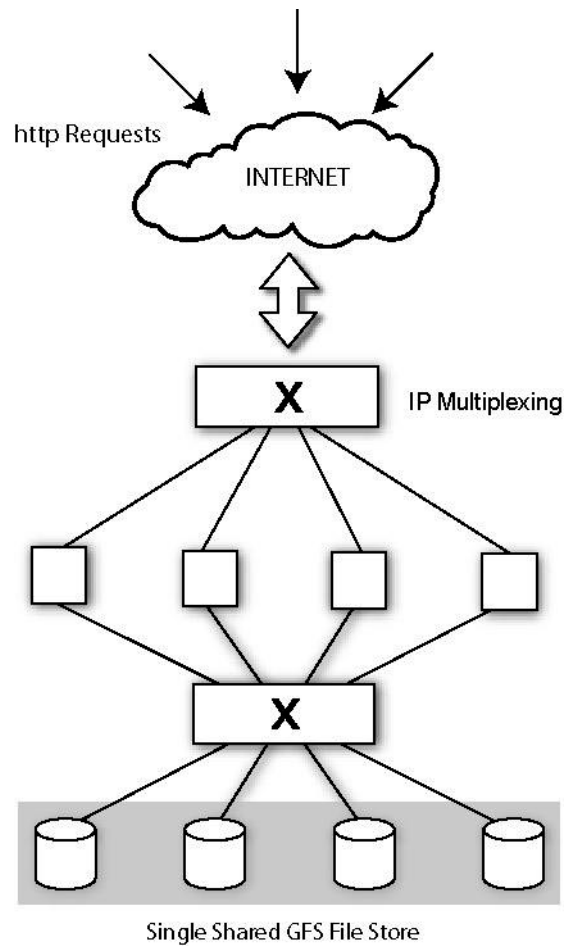
# Storage Clustering

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

- Storage cluster: servers sharing storage with a cluster file system and volume manager over a storage area network
- Integrating an edge server farm with a storage cluster creates an *edge server cluster*

# Edge Server Cluster

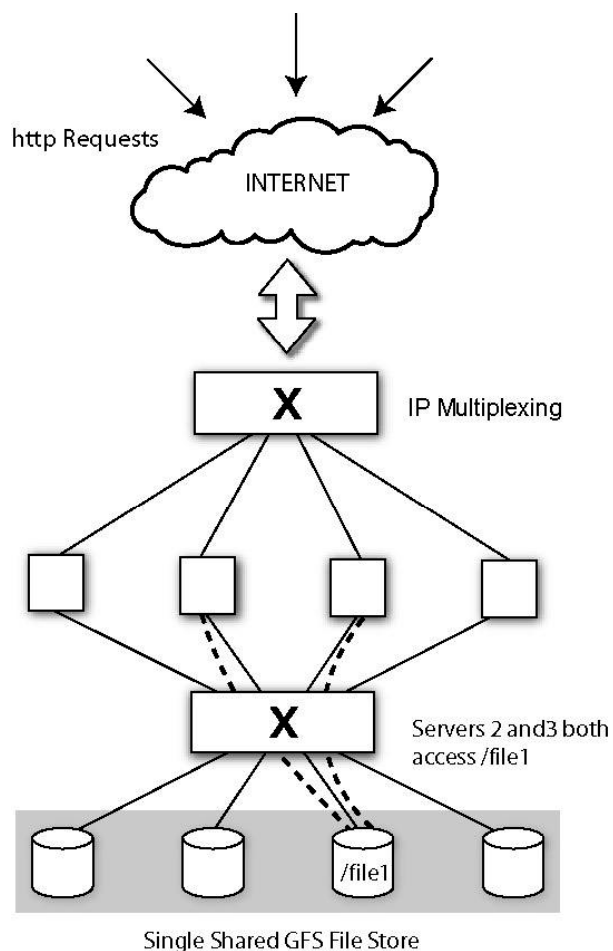
A NEW DIRECTION IN DATA STORAGE & MANAGEMENT



- Servers share storage: either storage or servers can be scaled independently, integrates fail-over and load balancing
- Edge server clusters exploit commodity technologies to achieve enterprise-class scalability, performance, and availability
- Approach works well for most client-server protocols

# Shared File Access in an Edge Server Cluster

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT



- Storage area networks are not good enough: cluster file system provides shared read/write access to applications across the edge server cluster
- IP Load Balancer is an important part of the architecture: helps virtualize the edge servers from the standpoint of outside clients

- About 40% of our users are building edge servers
  - applications include web serving, ftp, and email
  - availability is a big factor, especially for email
  - performance is also a factor
  - users are learning the storage management advantages of clustering as they deploy Sistina's storage cluster technology

# Edge Serving Applications

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

- We are collaborating with the Samba team to insure CIFS clustering support
- We are developing an NFS cluster server stack to integrate with GFS and CLVM
  - current Linux NFS stack works fine, but additional performance and capabilities can be provided
- We are beginning to work with load balancing switch vendors

# Summary on Edge Serving

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

- Edge server clustering can greatly simplify service and storage management complexity
- It can also increase ease of management, throughput, and availability
- GFS cluster file system and CLVM logical volume management technologies enable edge server clusters in Linux



# Oracle Night and Day

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

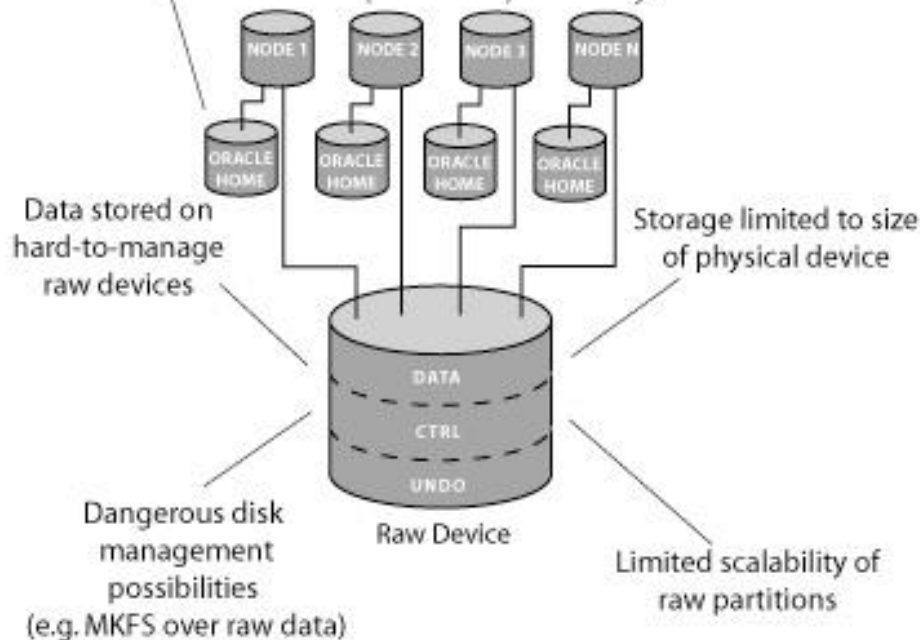
## TRADITIONAL ORACLE OPS ARCHITECTURE

Each node requires individual maintenance

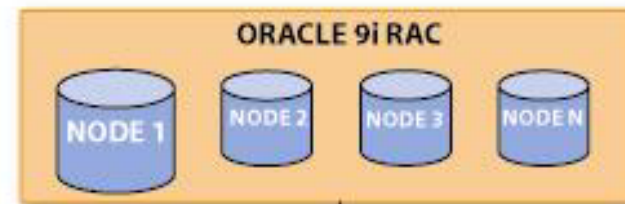
Each node must be identically configured

All 105,000 Oracle Home files must be installed directly to each node.

Adding nodes is difficult and complex



## ORACLE 9i RAC WITH SISTINA'S GFS



More Reliable

Higher Availability

Easier Management

Higher Performance

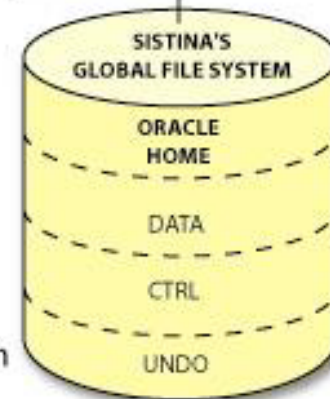
Completely Scalable

Shared Root File System

Safer Data Management

Data Files in File Systems

Lowest TCO via Commodity Hardware



# Cost Savings Example (Legacy)

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

Qty	Vendor	Model	Description
2	Sun	Enterprise 10000	64 466MHz Sparc II processors
2	Sun	Solaris 8	Included w/Sun System cost
40	Sun	Sbus HBA	Dual-Loop FC-AL Sbus HBA
6	Brocade	Silkworm 2040	8 port Fibre Switch
4	EMC	Clariion FC4500	Rack Mountable Raid 0+1 Disk Array
64	EMC	Clariion 18GB Drives	18GB drive for Clariion
2	Cisco	Catalyst 2924	24-port 10/100 Switch (Enterprise Edition)
2	Sun	Sun Clusters	HA Agent for Failover
2	Sun	SC Agent for Oracle	HA Agent for Oracle Database
<b>Total</b>			<b>\$3,414,194.00</b>

# Cost Savings Example (w/GFS)

A NEW DIRECTION IN DATA STORAGE & MANAGEMENT

Qty	Vendor	Model	Description
64	Dell	1U Rack Server	Dual Pentium III @ 933MHz
64	SuSE	SuSE 7.3	Linux Operating System
64	Qlogic	QLA2200f	64bit, 64MHz, Copper, HBA
5	Brocade	Silkworm 2040	8 port fibre switch
4	EMC	Clariion FC4500	Rack Mountable Raid 0+1 Disk Array
64	EMC	Clariion 18GB Drives	18GB drive for Clariion
4	Cisco	Catalyst 2924	24-port 10/100 Switch (Enterprise Edition)
64	Sistina	GFS 5.0	Cluster File System for Linux
<b>Total</b>			<b>\$225,654.00</b>

- Enterprise applications require ease of management, availability, throughput
- Storage clustering technology can provide all three — SAN and NAS alone do not
- Faster, cheaper SANs and virtualization, coupled with cluster file systems provide storage clustering foundations
- Thanks for your attention
- Questions?