# NFS: What's Next

## David L. Black, Ph.D.

Senior Technologist

EMC Corporation

black_david@emc.com

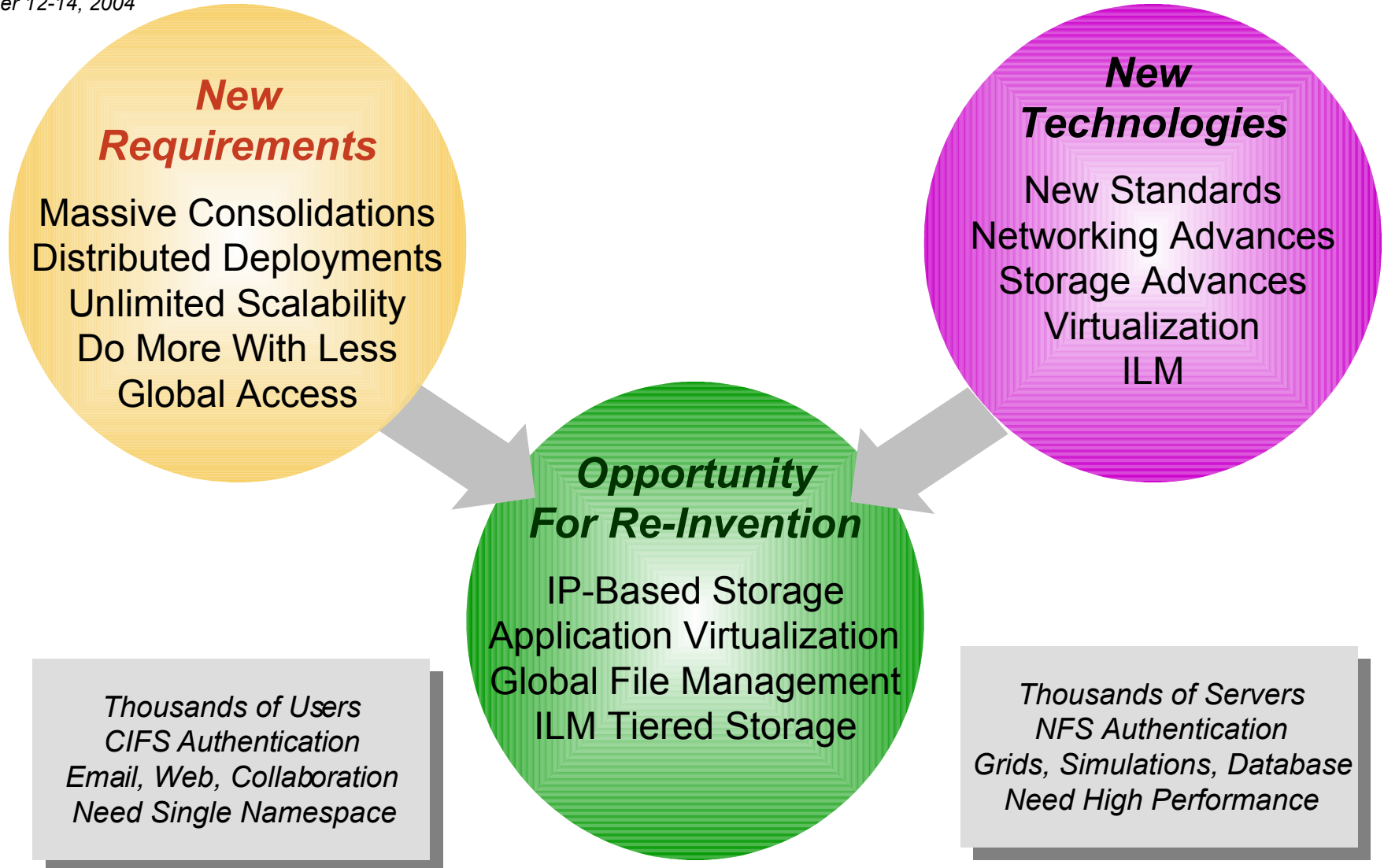# We Briefly Interrupt this Presentation ...

*EMC makes no representation and undertakes no obligations with regard to product planning information, anticipated product characteristics, performance specifications, or anticipated release dates (collectively, "Roadmap Information"). Roadmap Information is provided by EMC as an accommodation to the recipient solely for purposes of discussion and without intending to be bound thereby.*

### ... We Now Return to the Regularly Scheduled Presentation

# Networked Storage is Changing….

## New Requirements

Massive Consolidations
Distributed Deployments
Unlimited Scalability
Do More With Less
Global Access

## New Technologies

New Standards
Networking Advances
Storage Advances
Virtualization
ILM

## Opportunity For Re-Invention

IP-Based Storage
Application Virtualization
Global File Management
ILM Tiered Storage

*Thousands of Users*
*CIFS Authentication*
*Email, Web, Collaboration*
*Need Single Namespace*

*Thousands of Servers*
*NFS Authentication*
*Grids, Simulations, Database*
*Need High Performance*

2004 NAS Industry Conference

## *IP-Based Storage Delivering*

➲ Infinite Scalability

**SCALE UP**

# EMC NAS Vision

## *IP-Based Storage Delivering*

➲ Infinite
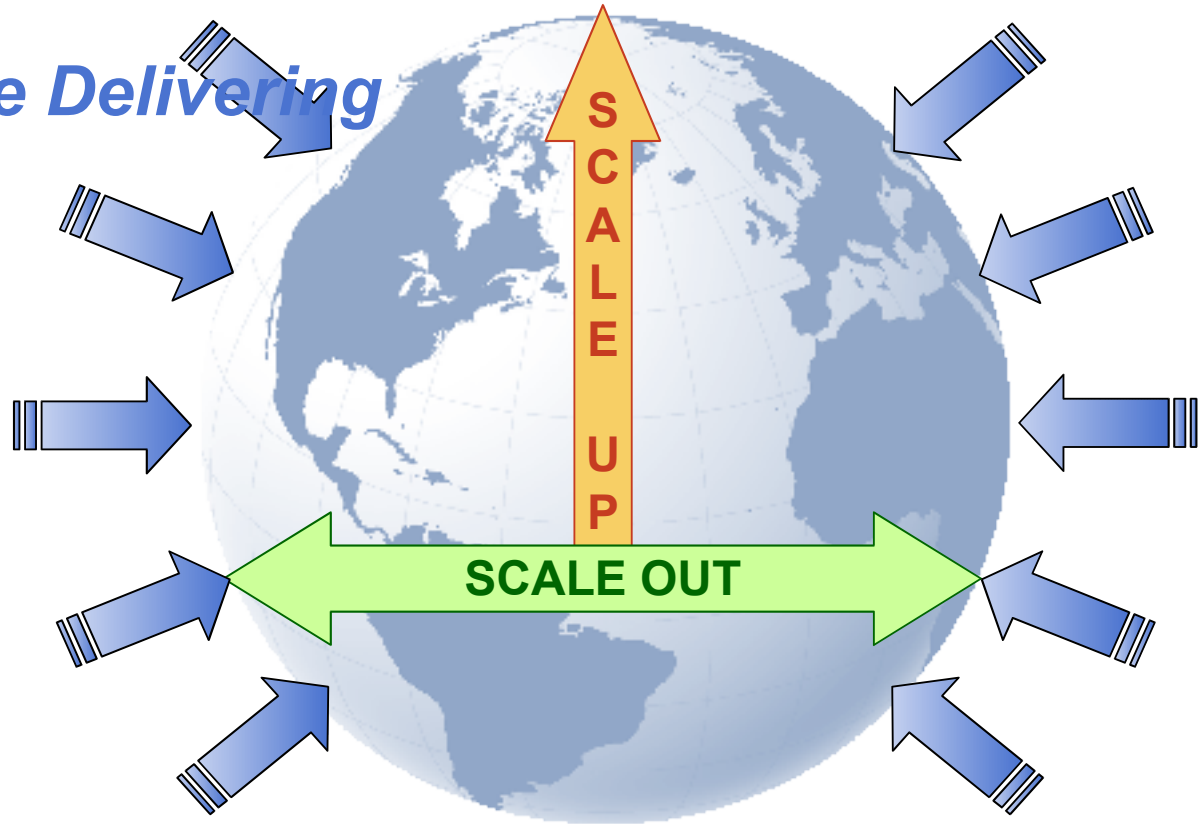   Scalability

➲ Optimized
   Data Placement

**SCALE UP**

**SCALE OUT**

# EMC NAS Vision

## *IP-Based Storage Delivering*

- ➔ Infinite Scalability

- ➔ Optimized Data Placement

- ➔ Global Accessibility

**SCALE UP**

**SCALE OUT**

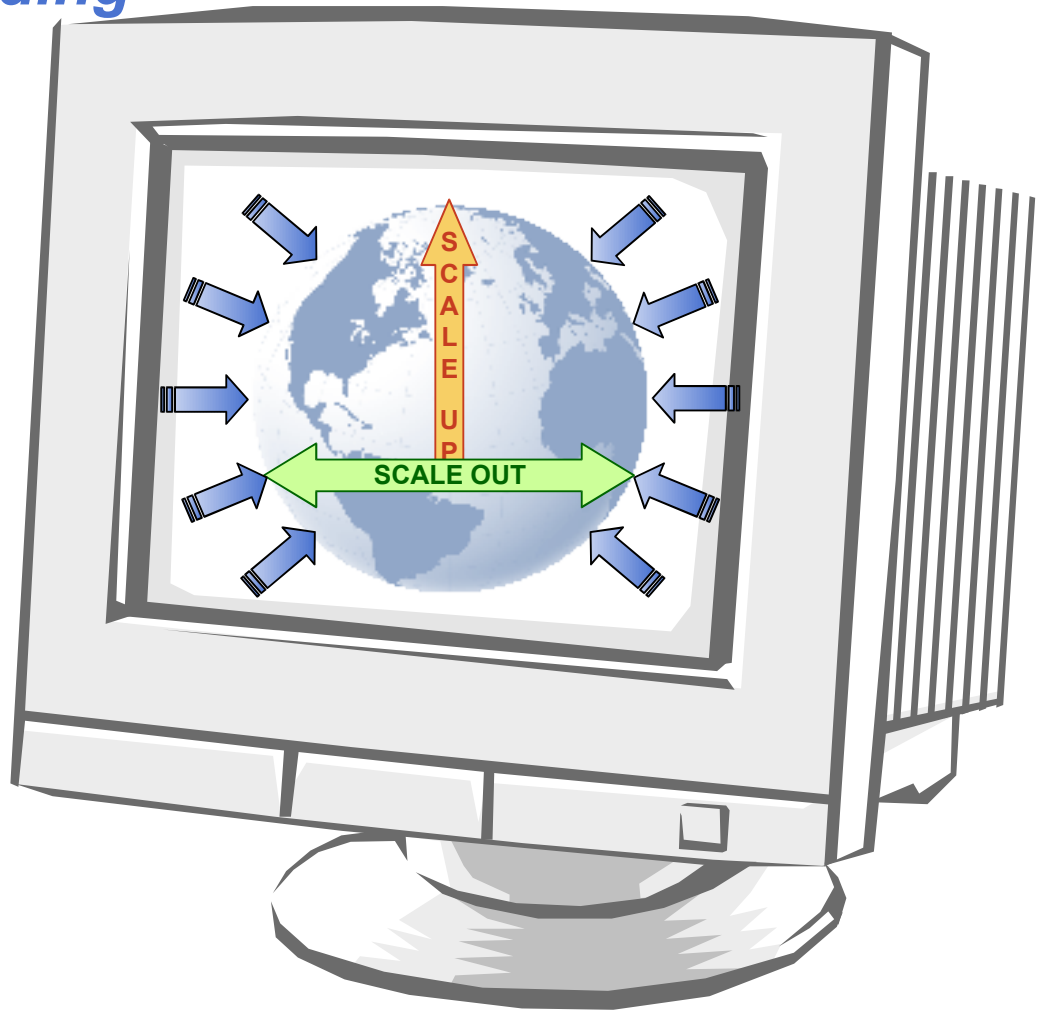# EMC NAS Vision

## *IP-Based Storage Providing*

➔ Infinite Scalability

➔ Optimized Data Placement

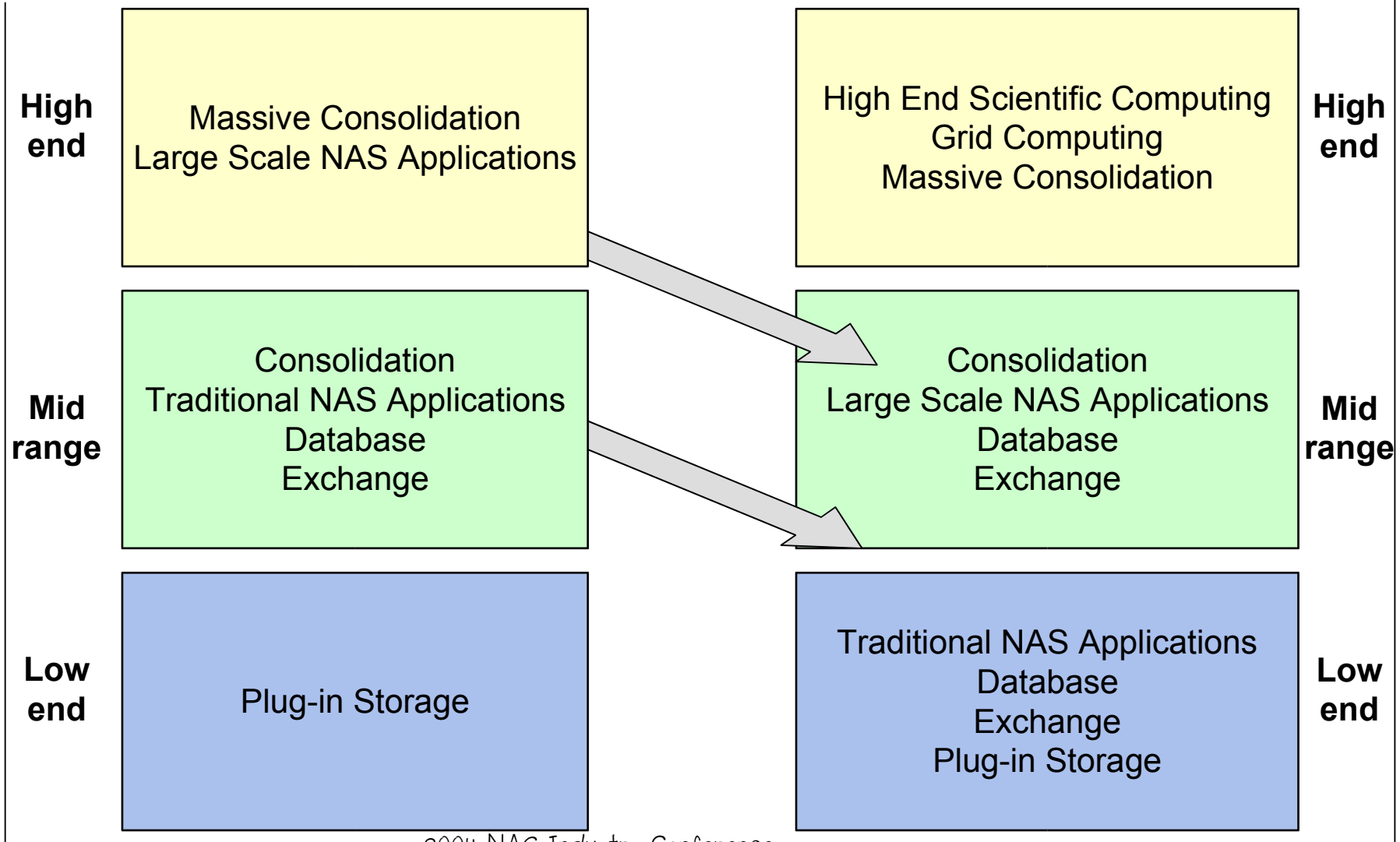➔ Global Accessibility

➔ Centralized Management

# NAS Use and Applications

**2004**

**2007**

| | 2004 | 2007 | |
|---|---|---|---|
| **High end** | Massive Consolidation<br>Large Scale NAS Applications | High End Scientific Computing<br>Grid Computing<br>Massive Consolidation | **High end** |
| **Mid range** | Consolidation<br>Traditional NAS Applications<br>Database<br>Exchange | Consolidation<br>Large Scale NAS Applications<br>Database<br>Exchange | **Mid range** |
| **Low end** | Plug-in Storage | Traditional NAS Applications<br>Database<br>Exchange<br>Plug-in Storage | **Low end** |

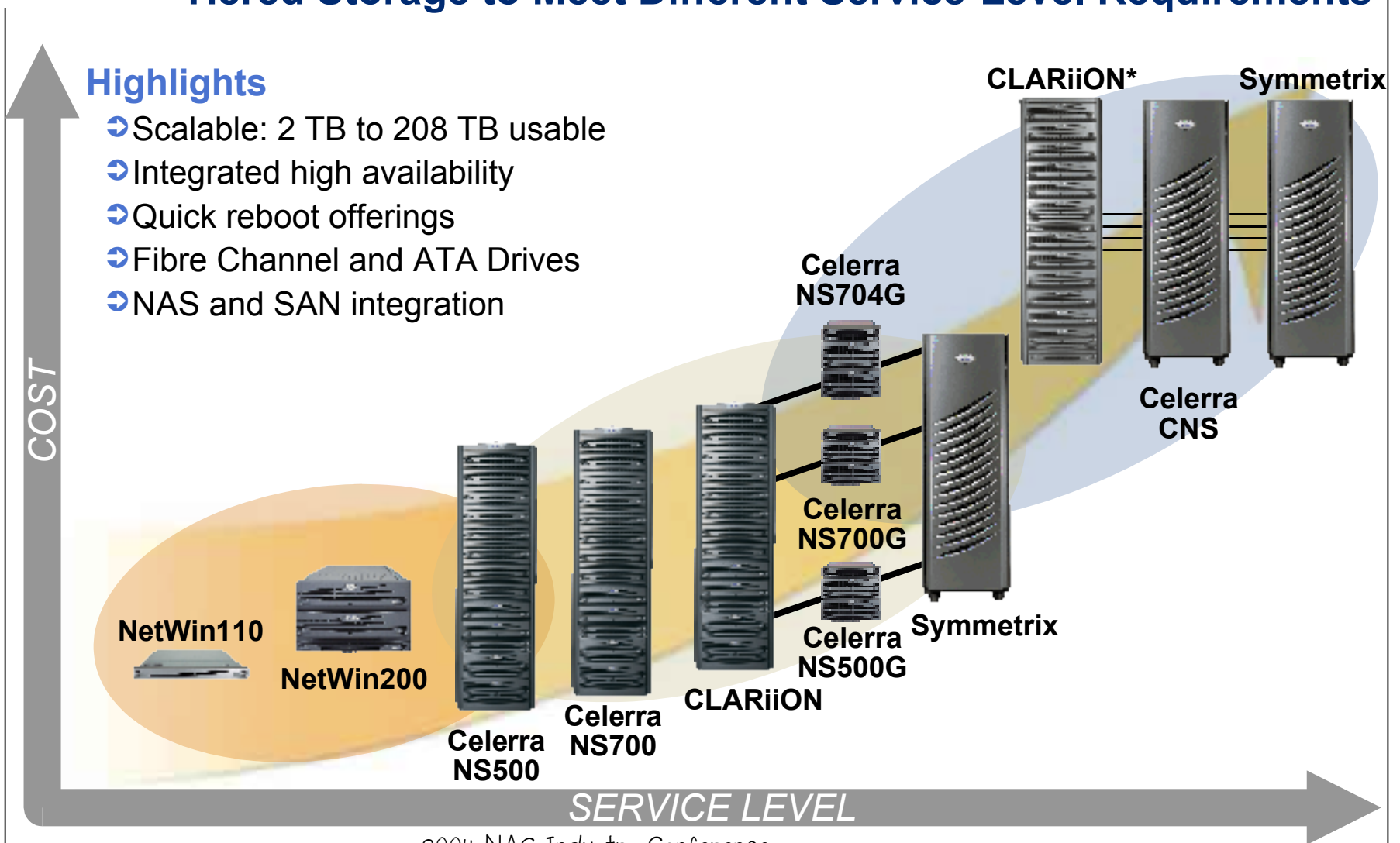# EMC NAS Family

## Tiered Storage to Meet Different Service-Level Requirements

### Highlights

➲ Scalable: 2 TB to 208 TB usable
➲ Integrated high availability
➲ Quick reboot offerings
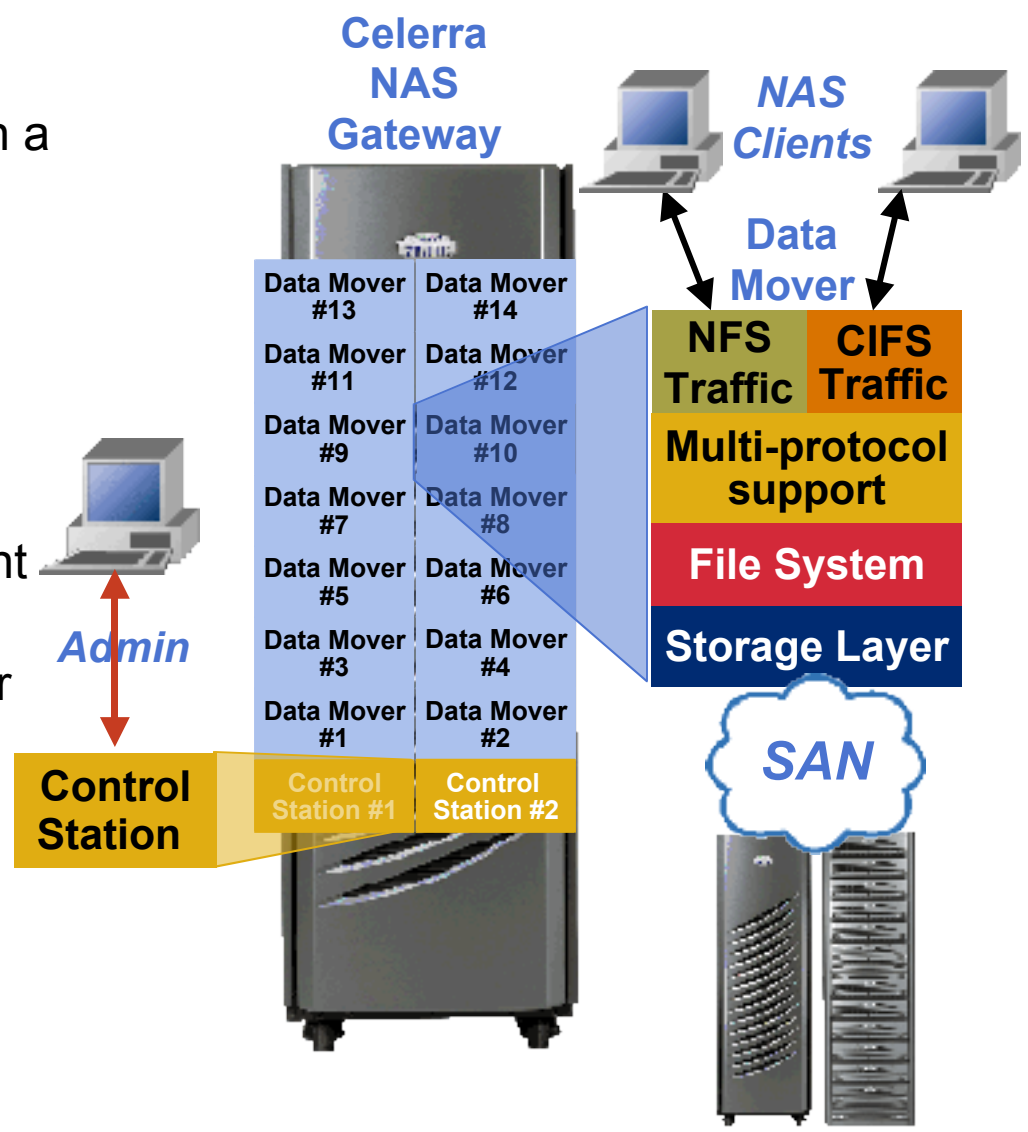➲ Fibre Channel and ATA Drives
➲ NAS and SAN integration

**COST**

**CLARiiON*** **Symmetrix**

**Celerra NS704G**

**Celerra CNS**

**Celerra NS700G**

**Symmetrix**

**Celerra NS500G**

**NetWin110**

**NetWin200**

**CLARiiON**

**Celerra NS500**

**Celerra NS700**

**SERVICE LEVEL**

2004 NAS Industry Conference

# EMC Celerra

- Up to 14 file servers contained in a single clustered system
- Managed as a single server
- NAS front-end scales independently of SAN back-end
- N to 1 failover options

- Control Station
  - Administration & management
  - Web-based GUI
  - Manages Data Mover failover

- Data Mover
  - Optimized real-time OS
  - Concurrent NFS and CIFS file access
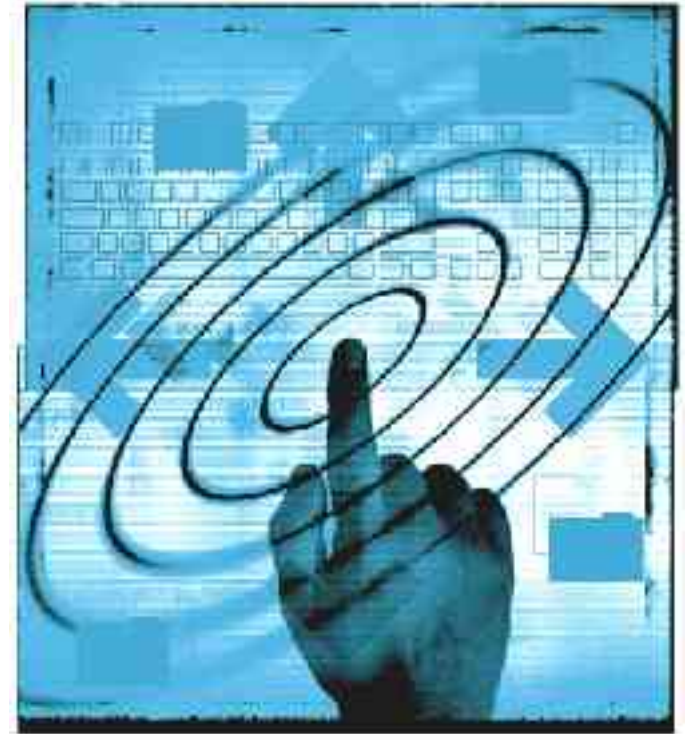  - No performance impact after failover

**Celerra NAS Gateway**

*NAS Clients*

**Data Mover**

| Data Mover #13 | Data Mover #14 |
| Data Mover #11 | Data Mover #12 |
| Data Mover #9 | Data Mover #10 |
| Data Mover #7 | Data Mover #8 |
| Data Mover #5 | Data Mover #6 |
| Data Mover #3 | Data Mover #4 |
| Data Mover #1 | Data Mover #2 |

**NFS Traffic**  **CIFS Traffic**

**Multi-protocol support**

**File System**

**Storage Layer**

*Admin*

**Control Station**

Control Station #1  **Control Station #2**

*SAN*

2004 NAS Industry Conference

# NAS Usage Scenarios

- Massive Consolidation Workloads
  - Cluster FS, Single Namespace

- Tiered Storage
  - Celerra FileMover API

- High Performance Computing
  - Multi-path IP SAN Filesystem

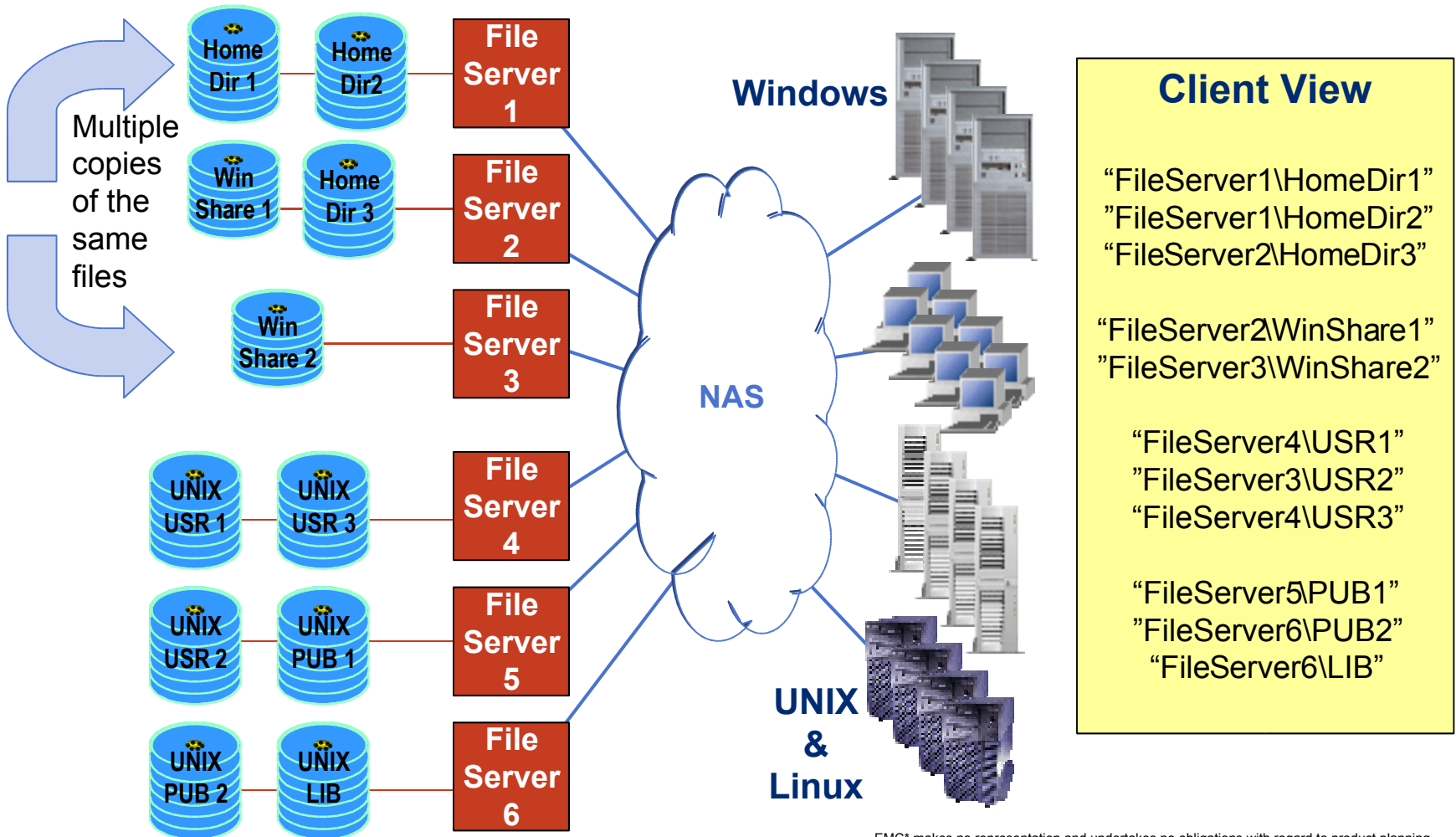- Integrated Block and File
  - iSCSI Target and Initiator

# Massive Consolidation Workloads: Before

Multiple copies of the same files

**Home Dir 1** — **Home Dir2** — **File Server 1**

**Win Share 1** — **Home Dir 3** — **File Server 2**

**Win Share 2** — **File Server 3**

**UNIX USR 1** — **UNIX USR 3** — **File Server 4**

**UNIX USR 2** — **UNIX PUB 1** — **File Server 5**

**UNIX PUB 2** — **UNIX LIB** — **File Server 6**

**Windows**

**NAS**

**UNIX & Linux**

## Client View

"FileServer1\HomeDir1"
"FileServer1\HomeDir2"
"FileServer2\HomeDir3"

"FileServer2\WinShare1"
"FileServer3\WinShare2"

"FileServer4\USR1"
"FileServer3\USR2"
"FileServer4\USR3"

"FileServer5\PUB1"
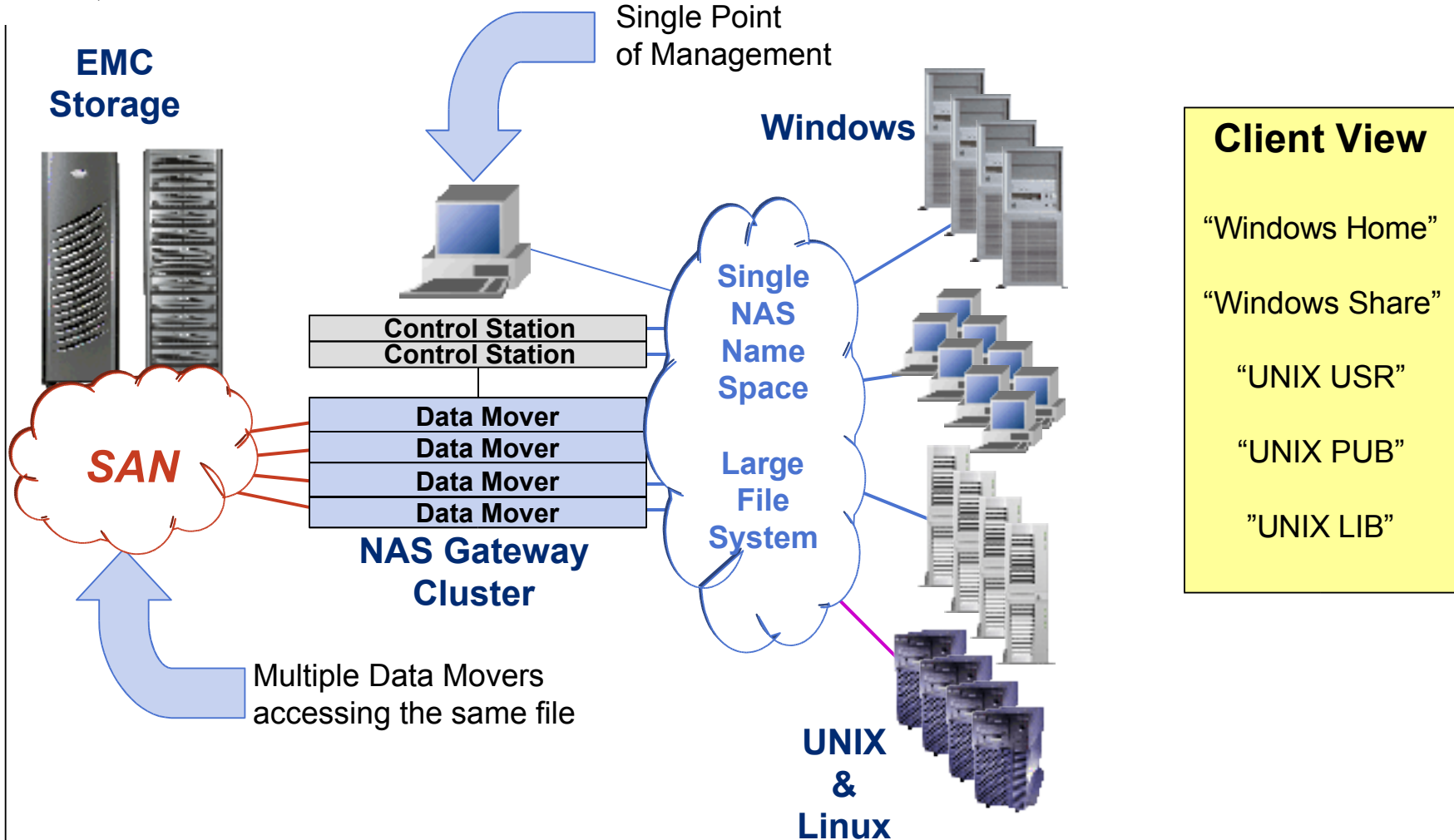"FileServer6\PUB2"
"FileServer6\LIB"

EMC* makes no representation and undertakes no obligations with regard to product planning information, anticipated product characteristics, performance specifications, or anticipated release dates (collectively, "Roadmap Information"). Roadmap Information is provided by EMC as an accommodation to the recipient solely for purposes of discussion and without intending to be bound thereby.

# Massive Consolidation Workloads: After

**EMC Storage**

Single Point of Management

**Windows**

Control Station
Control Station

Data Mover
Data Mover
Data Mover
Data Mover

**NAS Gateway Cluster**

*SAN*

**Single NAS Name Space**

**Large File System**

Multiple Data Movers accessing the same file

**UNIX & Linux**

**Client View**

"Windows Home"

"Windows Share"

"UNIX USR"

"UNIX PUB"

"UNIX LIB"

EMC* makes no representation and undertakes no obligations with regard to product planning information, anticipated product characteristics, performance specifications, or anticipated release dates (collectively, "Roadmap Information"). Roadmap Information is provided by EMC as an accommodation to the recipient solely for purposes of discussion and without intending to be bound thereby.

2004 NAS Industry Conference

# Tiered Storage: Celerra FileMover API

**Primary Storage**
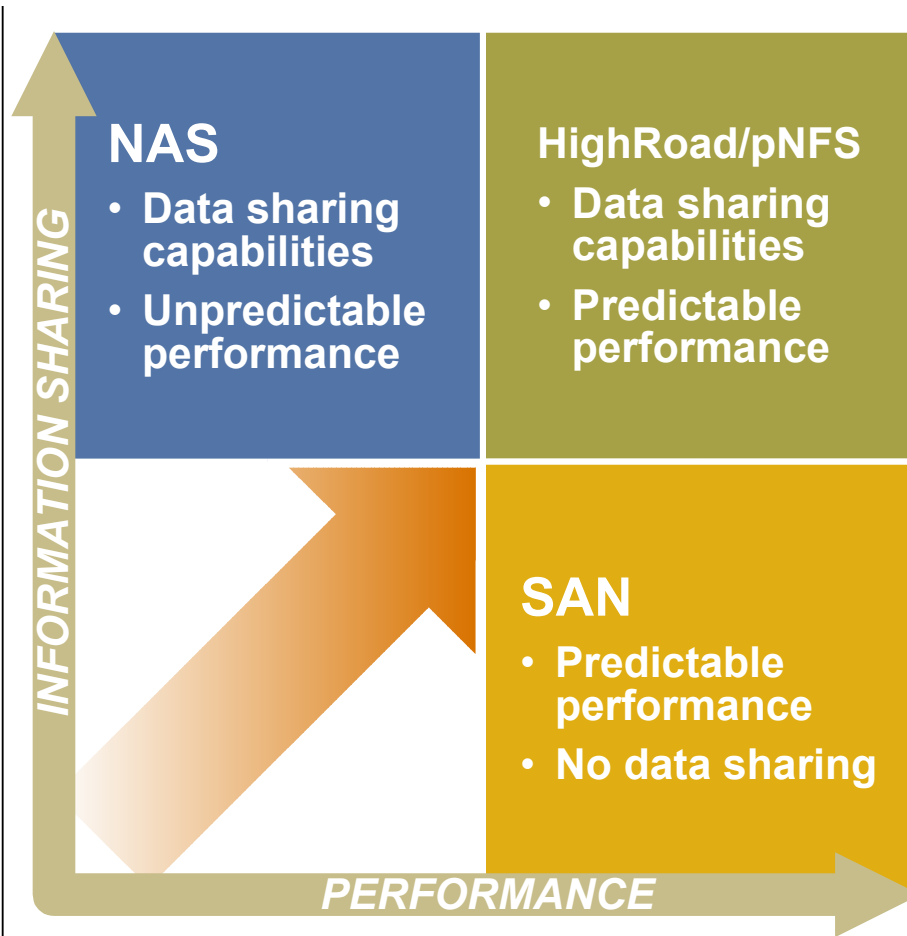
File I/O

**NAS Clients**

1. Policy Engine controls data migration

2. Metadata remains on fileserver

3. On NAS client access, fileserver retrieves data:
   - Pass through
   - Migration back

**ILM Policy Engine**

Retrieve File

Migrate File

**Celerra**

**Secondary Storage**

**Centera**     **ATA**     **Tape/Optical**

# Celerra HighRoad and Parallel NFS (pNFS)

**INFORMATION SHARING** →

**NAS**
- Data sharing capabilities
- Unpredictable performance

**HighRoad/pNFS**
- Data sharing capabilities
- Predictable performance

**SAN**
- Predictable performance
- No data sharing

**PERFORMANCE** →

## Applications

- Media
  - Post production
  - Television finishing
  - Streaming video
  - Advertising

- Large image processing
  - Seismic
  - Medical
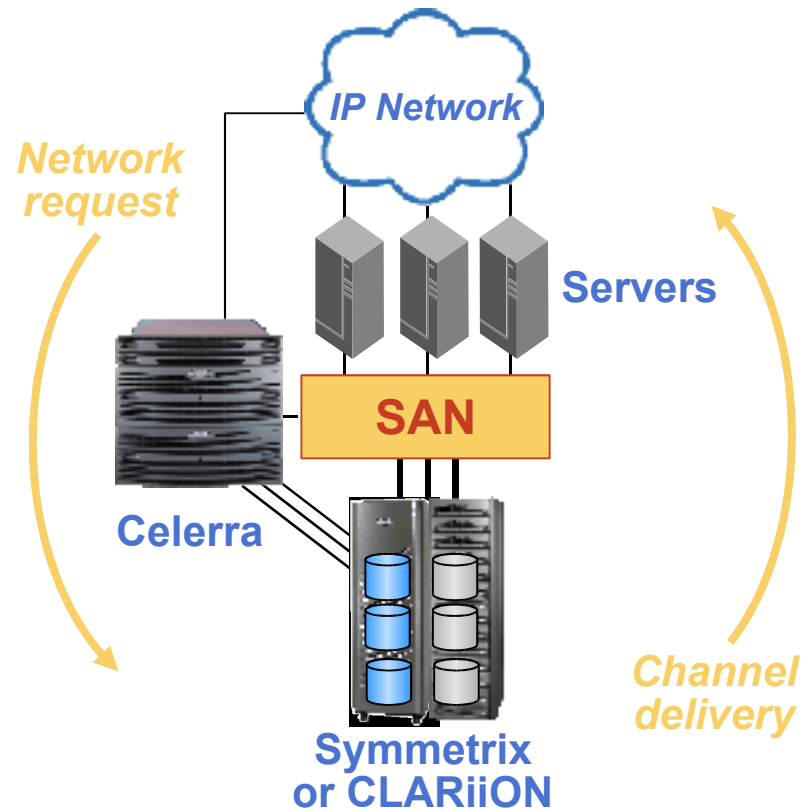  - CAD / CAM
  - Scientific simulations

- Backup

*Address both sharing and bandwidth challenges*

# Celerra HighRoad and pNFS



INTEGRATED NETWORK INFRASTRUCTURE

IP Network

Network request

Servers

Celerra

SAN

Channel delivery

Symmetrix or CLARiiON

- Network Request and Channel Delivery
  - Servers connected to storage over SAN
  - Servers connected to an out-of-band "meta data" cluster via IP
  - Servers send file requests to cluster via NFS/CIFS
  - Data access is direct via SAN (performance)
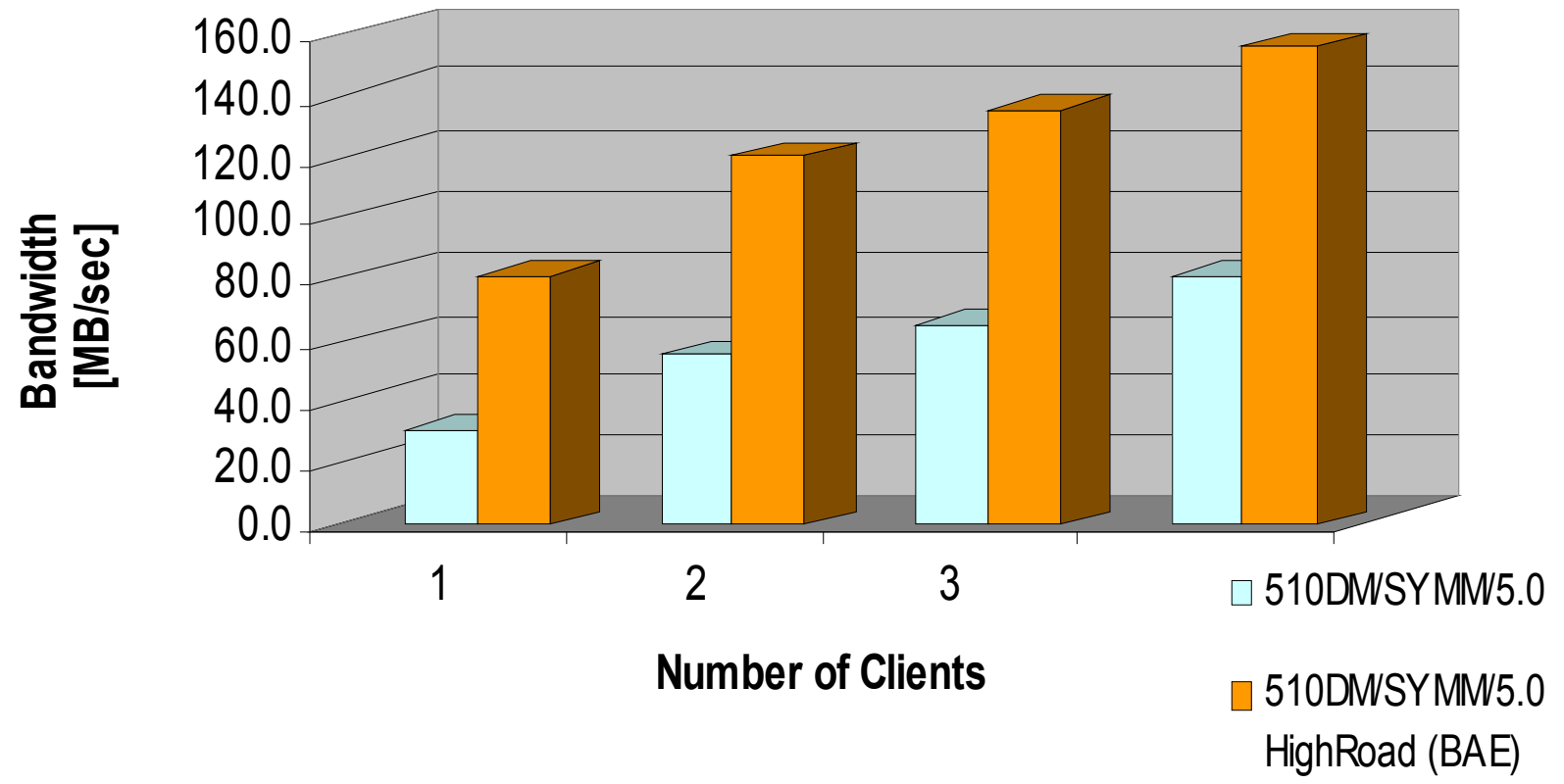  - Meta data cluster scaling improved by data bypass

# HighRoad Client Driver

**NFS operations**

**FileSystem**

**getaddr readdir etc.**

**NFS**

**Pass Through**

**HighRoad Driver**

**FMP**

**read write commit**

**Intercept**

**File mappings**

Small I/Os sent directly via NFS to avoid expense of obtaining mapping data

**SCSI / Fibre Channel**

**Data**

# HighRoad Metadata protocol: FMP

- Client asks fileserver "where is this file?"

- Fileserver provides map as answer
  - And grants read/write access permissions to client
  - Client can now read/write file blocks (via SAN)

- Update server state as needed
  - Hole filling (zeros or data)
    - Client does writes, tells server what it did
  - File extension or truncation (new EOF)
    - Server propagates new EOF to other clients

# FMP and pNFS

- ## FMP: Stand alone protocol
  - Can be used with NFSv2, NFSv3, NFSv4 and CIFS
  - Data caching and consistency are independent of protocol

- ## pNFS: NFSv4 extensions for mapping
  - Allow use of compound operations
  - Similar functional behavior to FMP
  - To be standardized in IETF

- ## FMP $\rightarrow$ pNFS
  - EMC has provided FMP specification to the pNFS effort
  - Enable pNFS effort to learn from HighRoad experience
  - Celerra HighRoad product will evolve to support pNFS
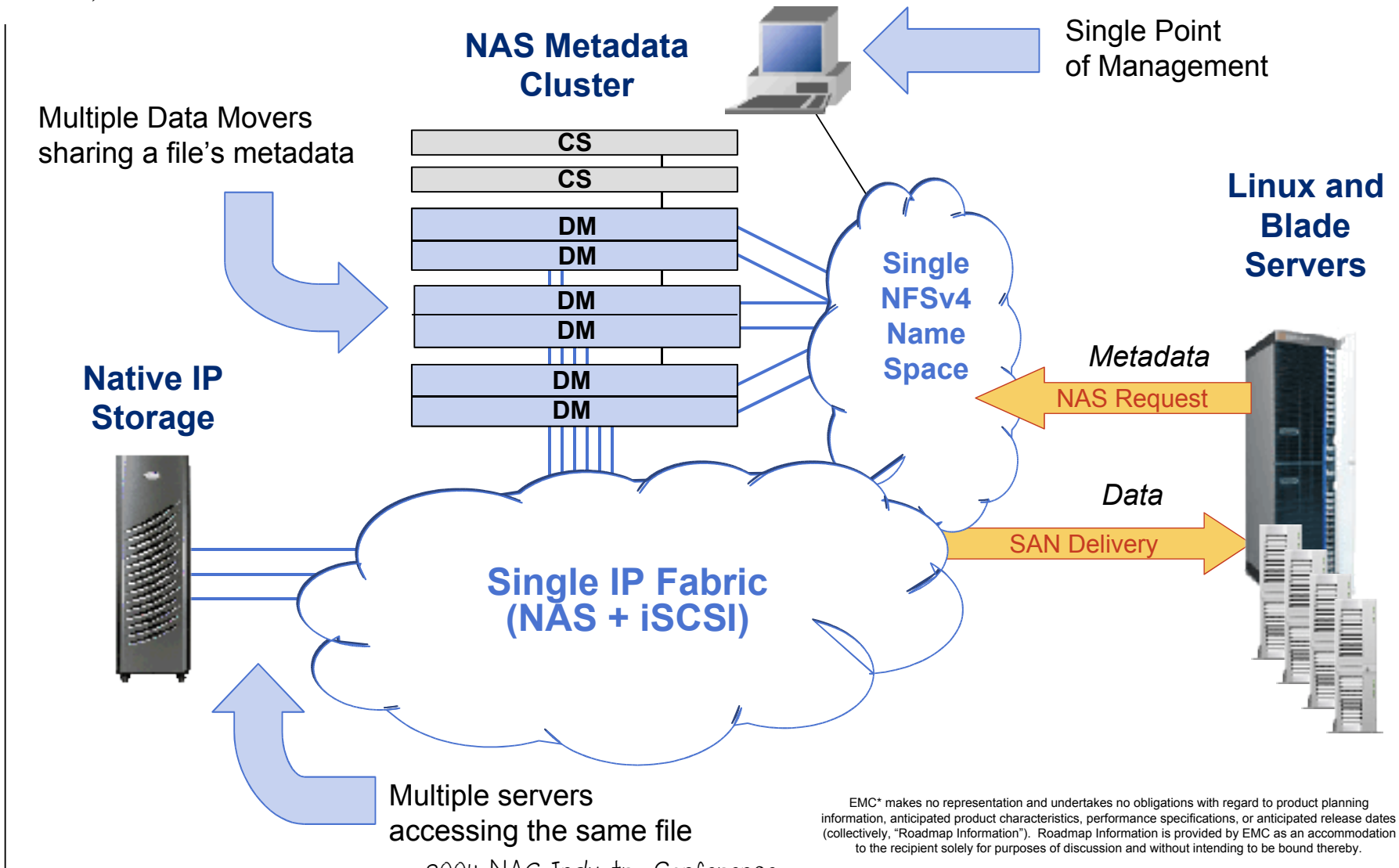
# Some Lessons Learned from HighRoad

- ## Volume identification (where's the file/FS?)
  - Have to use volume and/or filesystem labels
  - Addresses don't work (what's a SCSI address?)

- ## Access permission recall is essential and subtle
  - Server recall may conflict with pending client request
  - Out-of-order delivery can create race conditions

- ## Block permission granularity makes a difference
  - Whole file granularity creates false sharing conflicts

- ## Keep-alive needed to detect client death
  - And clean up any access permissions it held

# High Performance Computing Workloads

**NAS Metadata Cluster**

**Single Point of Management**

Multiple Data Movers sharing a file's metadata

| CS |
| CS |
| DM |
| DM |
| DM |
| DM |
| DM |
| DM |

**Linux and Blade Servers**

**Single NFSv4 Name Space**

*Metadata*

NAS Request

**Native IP Storage**

*Data*

SAN Delivery

**Single IP Fabric (NAS + iSCSI)**

Multiple servers accessing the same file

EMC* makes no representation and undertakes no obligations with regard to product planning information, anticipated product characteristics, performance specifications, or anticipated release dates (collectively, "Roadmap Information"). Roadmap Information is provided by EMC as an accommodation to the recipient solely for purposes of discussion and without intending to be bound thereby.
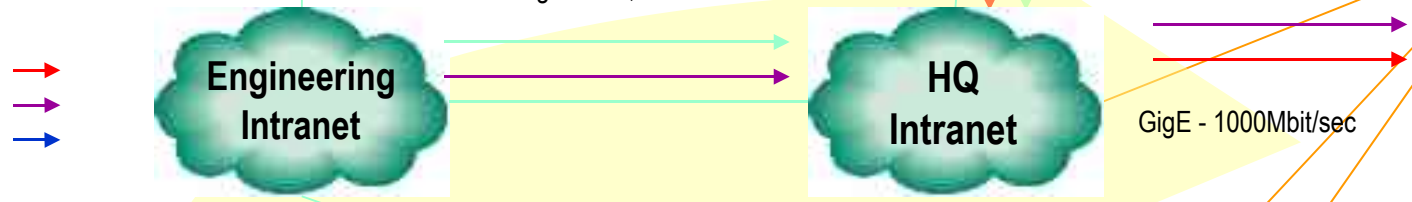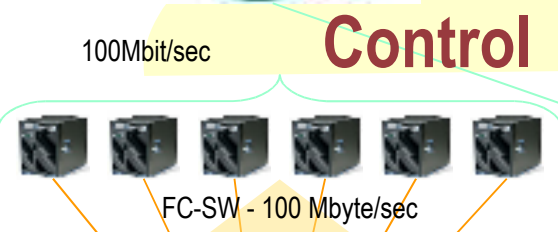
2004 NAS Industry Conference
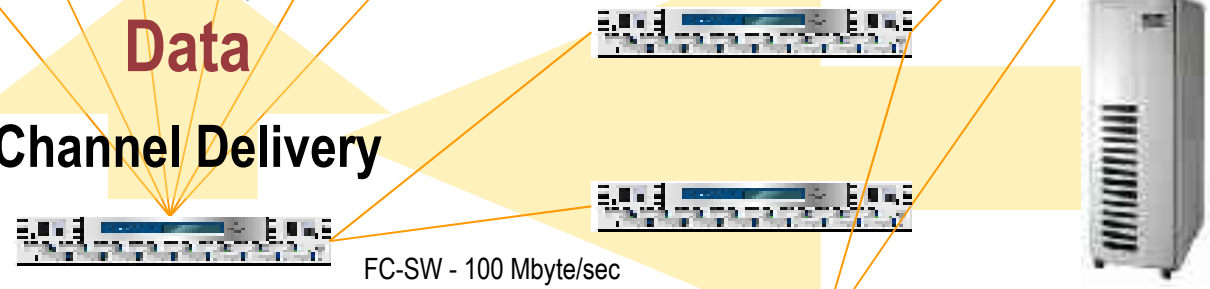
# Example: Consumer Electronics Manufacturer

- 150-200 UNIX workstations

- Sun, SGI, HP, IBM, Etc.

- Engineering users running CAD/CAM/CAE

Windows Clients

100Mbit/sec

ATM OC3, 155 Mbit/sec

GigE Soon, 1000 Mbit/sec

100Mbit/sec

**Engineering Intranet**

**HQ Intranet**

GigE - 1000Mbit/sec

**Control**

6 HP J6000 for Engineering Simulation runs on LSF Computer Engines. File sharing at channel speeds. 100Mbyte/sec data access to common file system.

100Mbit/sec

FC-SW - 100 Mbyte/sec

**Data**

FC-SW - 100 Mbyte/sec

**Channel Delivery**

CIFS Data Flow
NFS Data Flow
FTP Data Flow

FC-SW - 100 Mbyte/sec

Engineering NFS File
Windows CIFS File S
6 Active Data Movers
1 Standby Data Mov
1 Control Station
~4TB Usable Storage

2004 NAS Industry Conference

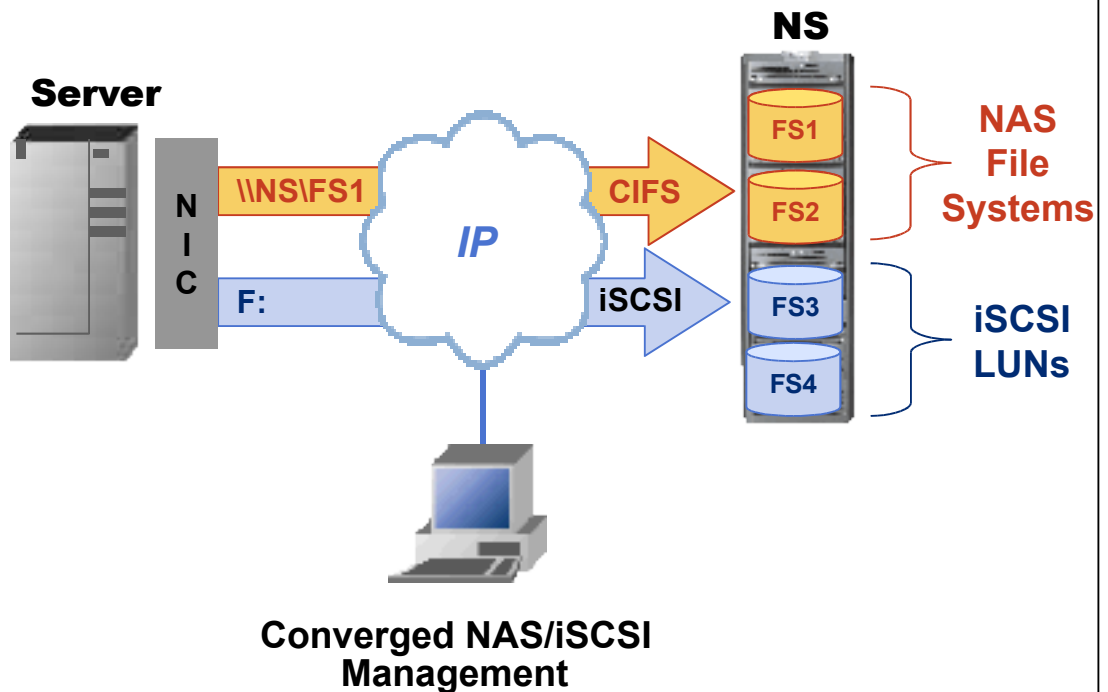# Block and File Workloads

## *iSCSI Target*

- Microsoft Logo Certified
  - iSNS Naming Service
  - CHAP Authentication

- Simple Management
  - Web-based GUI
  - Common toolset for NAS and IP SAN

- High Availability
  - Data Mover failover
  - Port/path failover



**Server**

**N I C**

**\\NS\FS1**

**F:**

**IP**

**CIFS**

**iSCSI**

**NS**

**FS1**

**FS2**

**FS3**

**FS4**

**NAS File Systems**

**iSCSI LUNs**

**Converged NAS/iSCSI Management**

EMC* makes no representation and undertakes no obligations with regard to product planning information, anticipated product characteristics, performance specifications, or anticipated release dates (collectively, "Roadmap Information"). Roadmap Information is provided by EMC as an accommodation to the recipient solely for purposes of discussion and without intending to be bound thereby.
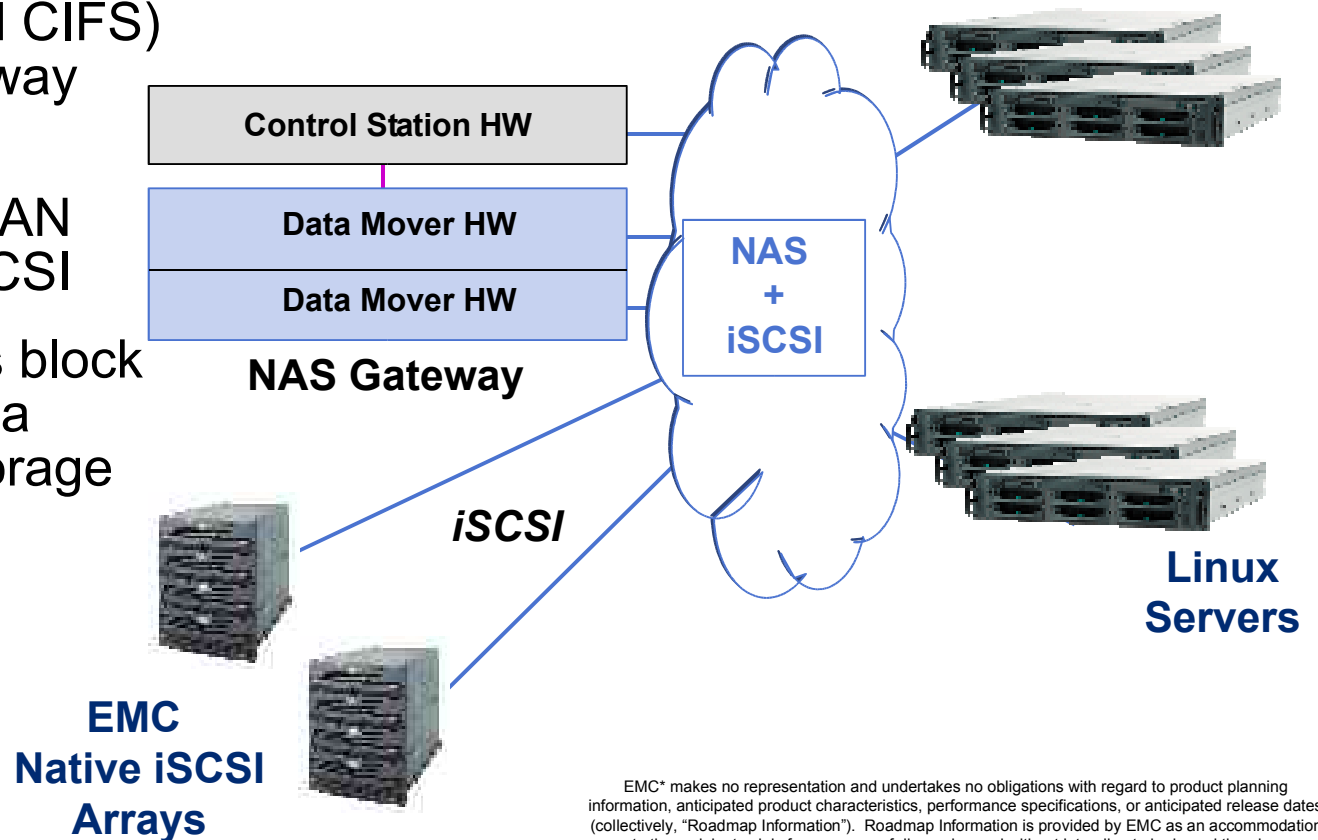
# NAS Gateway iSCSI Initiator

## *Adds NAS Services to IP SANs*

- Clients access file data (NFS and CIFS) via NAS Gateway

- NAS Gateway accesses IP SAN storage via iSCSI

- Clients access block data directly via iSCSI to IP storage

**Windows Servers**

| Control Station HW |
|---|
| Data Mover HW |
| Data Mover HW |

**NAS Gateway**

**NAS + iSCSI**

*iSCSI*

**EMC Native iSCSI Arrays**

**Linux Servers**

EMC* makes no representation and undertakes no obligations with regard to product planning information, anticipated product characteristics, performance specifications, or anticipated release dates (collectively, "Roadmap Information"). Roadmap Information is provided by EMC as an accommodation to the recipient solely for purposes of discussion and without intending to be bound thereby.
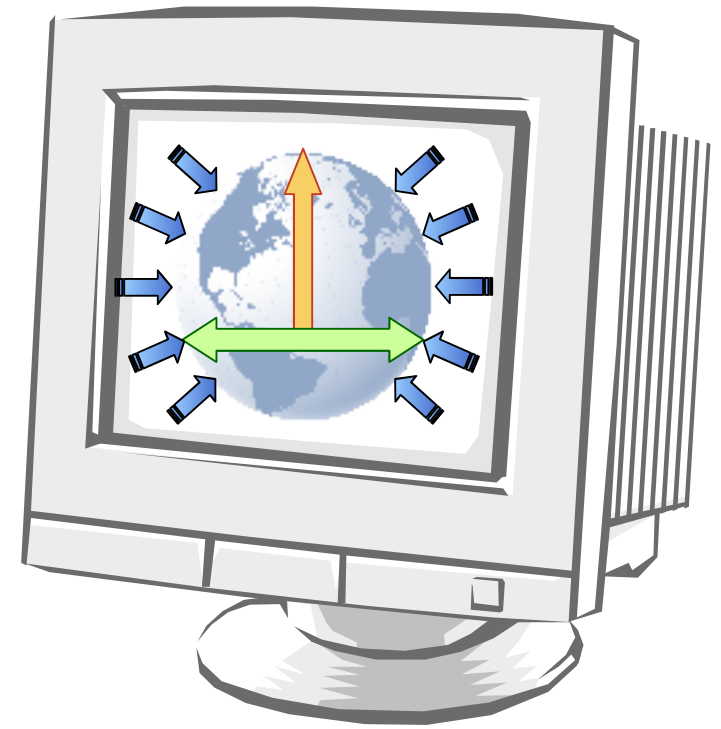
# Delivering on the EMC NAS Vision

## *Delivering on ILM*

✔ Infinite Scalability

- – Cluster File System and Large File Systems
- – Multi-path IP SAN File System

✔ Optimized Data Placement

- – FileMover API

✔ Global Accessibility

- – Single Name Space

✔ Centralized Management

- – Single Name Space
- – ISCSI Target and Initiator

EMC* makes no representation and undertakes no obligations with regard to product planning information, anticipated product characteristics, performance specifications, or anticipated release dates (collectively, "Roadmap Information").  Roadmap Information is provided by EMC as an accommodation to the recipient solely for purposes of discussion and without intending to be bound thereby.