



Education

File Systems for Object Storage Devices

Craig Harmer
Independent Data Storage Consultant

- The material contained in this tutorial is copyrighted by the SNIA.
- Member companies and individual members may use this material in presentations and literature under the following conditions:
 - ◆ Any slide or slides used must be reproduced in their entirety without modification
 - ◆ The SNIA must be acknowledged as the source of any material used in the body of any document containing material from these presentations.
- This presentation is a project of the SNIA Education Committee.
- Neither the author nor the presenter is an attorney and nothing in this presentation is intended to be, or should be construed as legal advice or an opinion of counsel. If you need legal advice or a legal opinion please contact your attorney.
- The information presented herein represents the author's personal opinion and current understanding of the relevant issues involved. The author, the presenter, and the SNIA do not assume any responsibility or liability for damages arising out of any reliance on or use of this information.

NO WARRANTIES, EXPRESS OR IMPLIED. USE AT YOUR OWN RISK.

➤ File Systems for Object Storage Devices

Object-based storage devices (OSDs) may well be the “next big thing” in file-oriented data storage. Already popular in the high-performance computing arena, they are poised to enter general enterprise computing environments. By distributing storage management and enabling secure data transfer between storage devices and clients, OSDs promise significant improvements in scaling and administrative simplicity. But making effective use of OSDs requires a new breed of file system—one that makes effective use of the new devices to deliver the promised benefits. This tutorial will describe the salient properties of OSDs, explain how file system technology is evolving to exploit the scalability and administrative simplicity they offer, identify the mature and emerging segments of the OSD-based file storage market, and show how technology that has been successful in HPC can be beneficially employed in general enterprise data center environments. Standardization activities, most notably the parallel NFS (pNFS) protocol for addressing OSDs will be discussed.

- Limitations of current file storage system designs

- Object Storage Devices (OSDs)
 - ◆ What they are
 - ◆ How they help overcome the limitations

- File systems for OSDs
 - ◆ Basic architecture
 - ◆ Adding value via performance, availability, and manageability

- So, you want an object-based file system...

- The ideal enterprise storage system would enable storage administrators to:
 - ◆ Share storage capacity across the data center
 - › Storage in a shared storage pool could be used by any host
 - ◆ Share data across the data center
 - › Data could be shared between any set of hosts at file or file system level
 - ◆ Share securely
 - › Data and storage would be shared only with authorized clients, preventing unauthorized or unintended access
 - ◆ Centralize administration
 - › Storage and storage access could be managed from a single console
 - ◆ Thin provision
 - › Only storage actually used by hosts would need to be installed

The Ideal Enterprise Storage System

- To satisfy those goals, a storage system would need:
 - ◆ Extreme capacity
 - › Billions of files, petabytes of data, support for thousands of clients
 - ◆ Scalable performance
 - Low latency and high bandwidth that scales nearly linearly with capacity
 - ◆ Universal access
 - › Data sharing among all data center computing platforms
 - ◆ File based storage
 - › Data shared at the file level
 - › File-based storage & dynamically sized file systems implicitly thin provision
 - ◆ Security
 - › Robust protection against unauthorized or unintended access
 - ◆ Incremental growth
 - › Scaling from Tbyte to Pbyte in bite-size increments

Teasing the problem apart ...

➤ People deal with “business objects”

- ◆ AKA “files”
- ◆ Names, sizes, access rights, lifetimes...



Teasing the problem apart ...

➤ People deal with “business objects”

- ◆ AKA “files”
- ◆ Names, sizes, access rights, lifetimes...



➤ Computers deal with storage objects

- ◆ Typically “blocks”
- ◆ Which have addresses ...and not much else



Teasing the problem apart ...

➤ People deal with “business objects”

- ◆ AKA “files”
- ◆ Names, sizes, access rights, lifetimes...



➤ In between: “file systems”

- ◆ Where is the free space for new files?
- ◆ Who is allowed to access `/year2010/march/.../results`?
- ◆ Where is the data for `/year2010/march/.../results` stored ?



- ◆ Namespace management
- ◆ Storage management



➤ Computers deal with storage objects

- ◆ Typically “blocks”
- ◆ Which have addresses ... and not much else



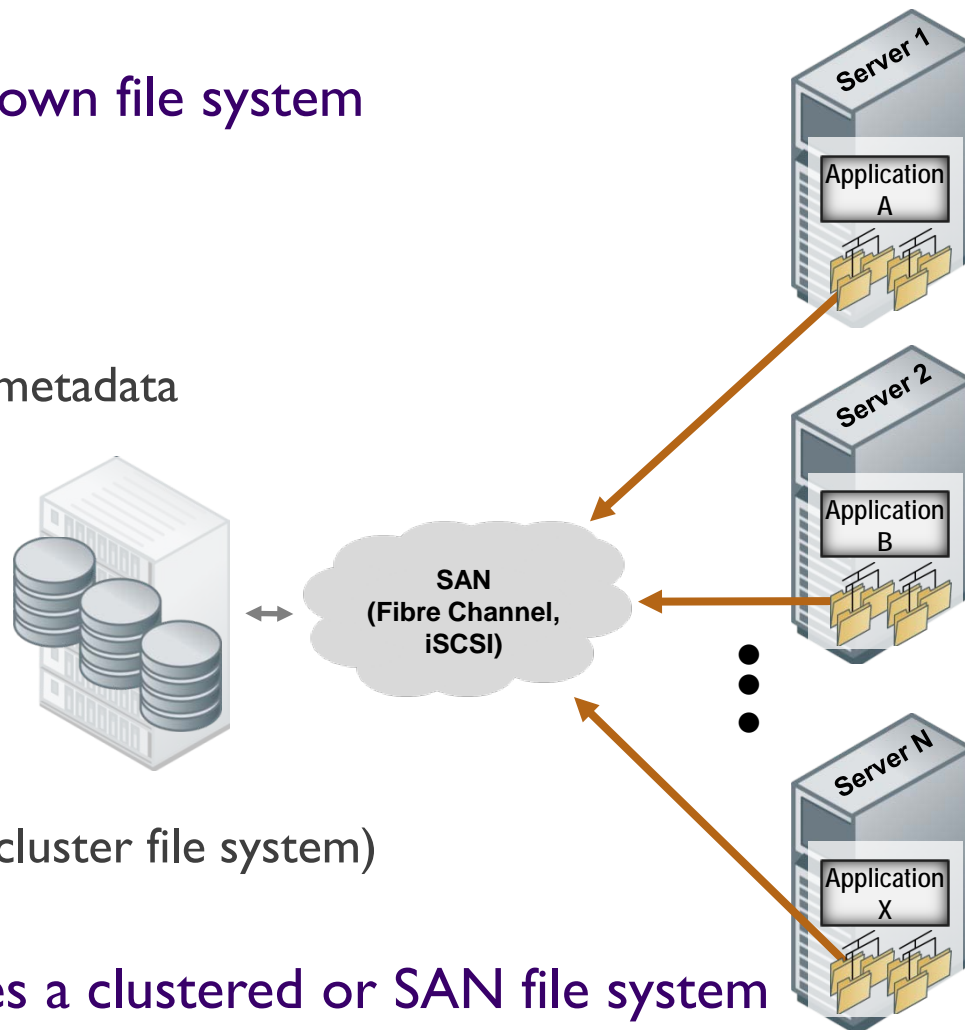
➤ Every client computer has its own file system

➤ That's good

- ◆ Short I/O paths
- ◆ Coherent caching of data and metadata

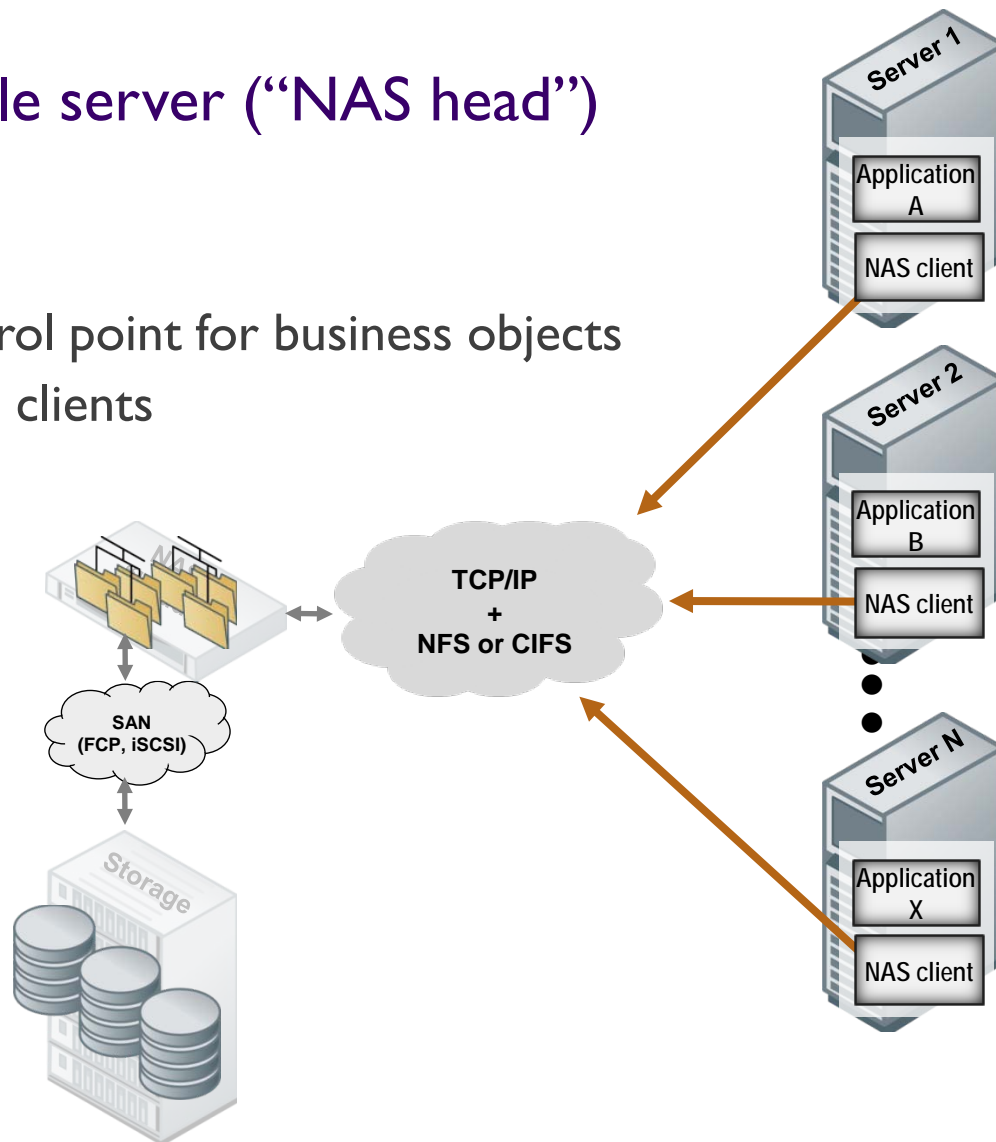
➤ ... and not so good

- ◆ Lots of file systems to manage
- ◆ Lots of storage access rights to coordinate
- ◆ No sharing of data (without a cluster file system)



➤ Sharing data or scaling requires a clustered or SAN file system

- Shared file system in the file server (“NAS head”)
- That’s good...
 - ◆ Single name space and control point for business objects
 - ◆ Consistent semantics for all clients
- And not so good...
 - ◆ High latency protocols
 - ◆ “NAS head” bottleneck
- Scalability via
 - ◆ Clustered NAS
 - ◆ NAS aggregators



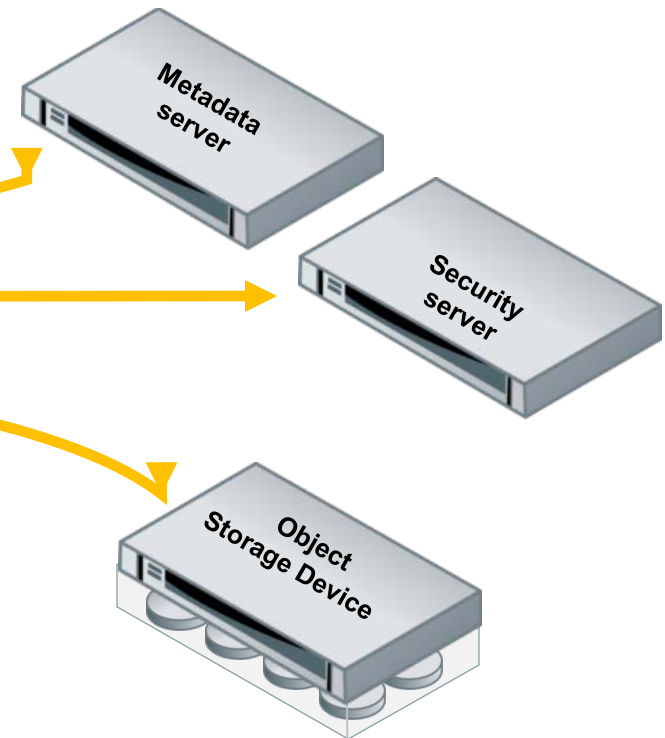
OSDs: A step closer to The Ideal

➤ The Object Storage Device (OSD)

- ◆ Hybrid between a disk and a file server

Allows the problem to be divided

- ◆ Namespace management
- ◆ Access security
- ◆ Storage management



➤ Enables

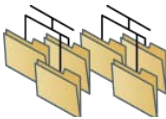
- ◆ Extreme scaling
- ◆ Robustness

OSDs: A step closer to The Ideal

➤ The Object Storage Device (OSD)

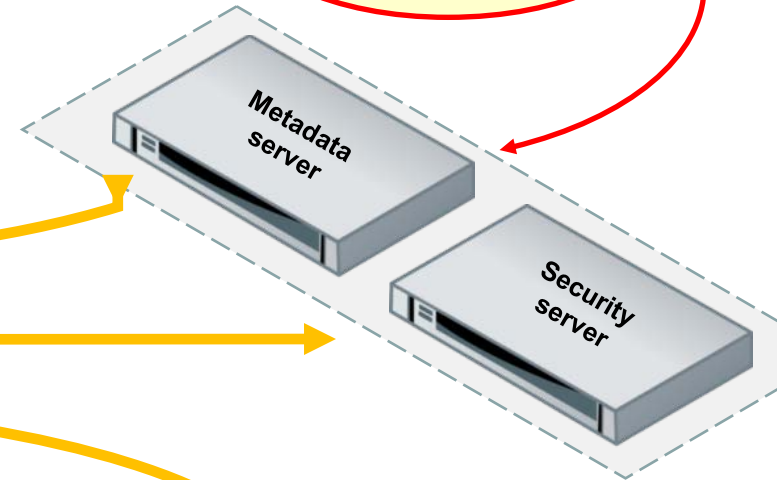
- ◆ Hybrid between a disk and a file server

Usually the same physical server



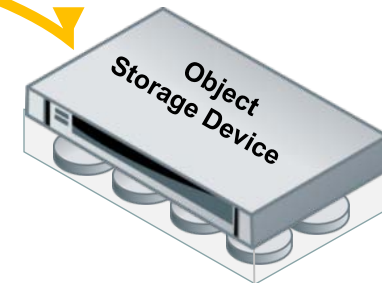
Allows the problem to be divided

- ◆ Namespace management
- ◆ Access security
- ◆ Storage management



➤ Enables

- ◆ Extreme scaling
- ◆ Robustness



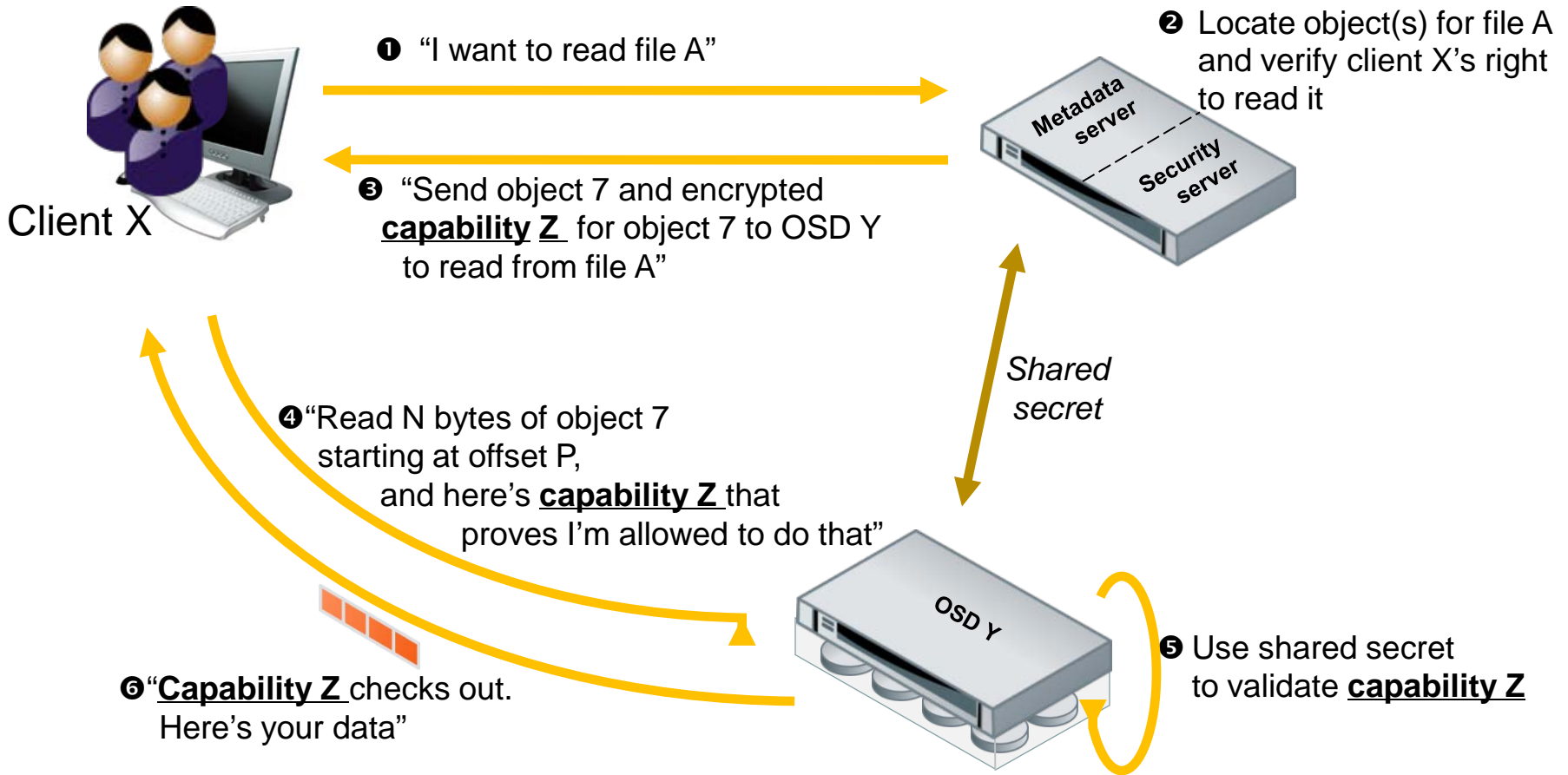
OSDs: Like disks, only different

	Disk	OSD
Model	Array of blocks <ul style="list-style-type: none"> • Number never changes • Size never changes 	Objects <ul style="list-style-type: none"> • Created and deleted • Grow and shrink
Operations	Read/write disk blocks	Create/delete object Read/write object blocks
Security	Zoning, LUN masking <ul style="list-style-type: none"> •Applies to entire device 	Capability-based <ul style="list-style-type: none"> •Applies to each object and op
Typical transports	Fibre Channel, SCSI, iSCSI	iSCSI, <u>ONC-RPC over TCP/IP</u>

OSDs: Like a file server, only different

	File server	OSD
Model	Files	Objects
Naming	Human-readable names in a hierarchical directory tree	Two level name space: 64 bit object “name” in a 64 bit partition “name”
Operations	File: create, delete, rename File byte range: read, write, append, truncate	Object: create, delete, Object block range: read, write, append, truncate
Security	User group world × rwX or access control lists • Checked at initial file access	Digitally signed capabilities • Checked for every I/O request

OSD Capabilities



File Systems for OSDs

➤ Scaling

- ◆ What if there's more data than the biggest OSD can hold?
- ◆ What if too many clients access an OSD at the same time?
- ◆ What if there's a file bigger than the biggest OSD can hold?

➤ Robustness

- ◆ What happens to data if the OSD fails?
- ◆ What happens to data if the metadata server fails?

➤ Performance

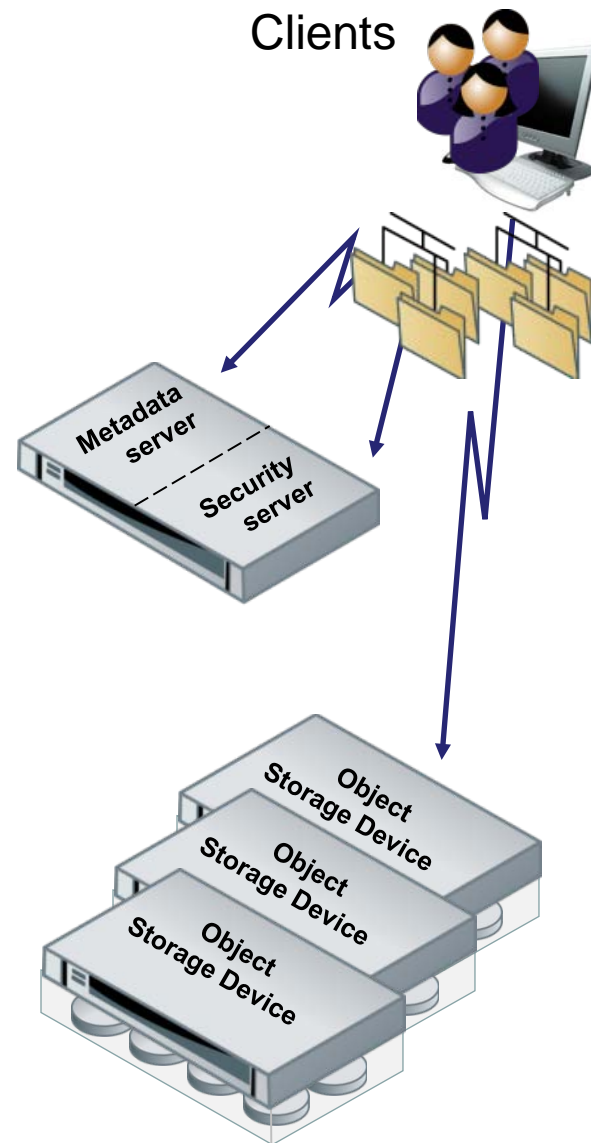
- ◆ What if thousands of objects are accessed concurrently?
- ◆ What if big files have to be transferred really fast?

➤ Architecture

- ◆ File = one or more *groups of objects* (usually on different OSDs)
- ◆ Clients access metadata servers to locate data
- ◆ Clients transfer data directly to & from OSDs

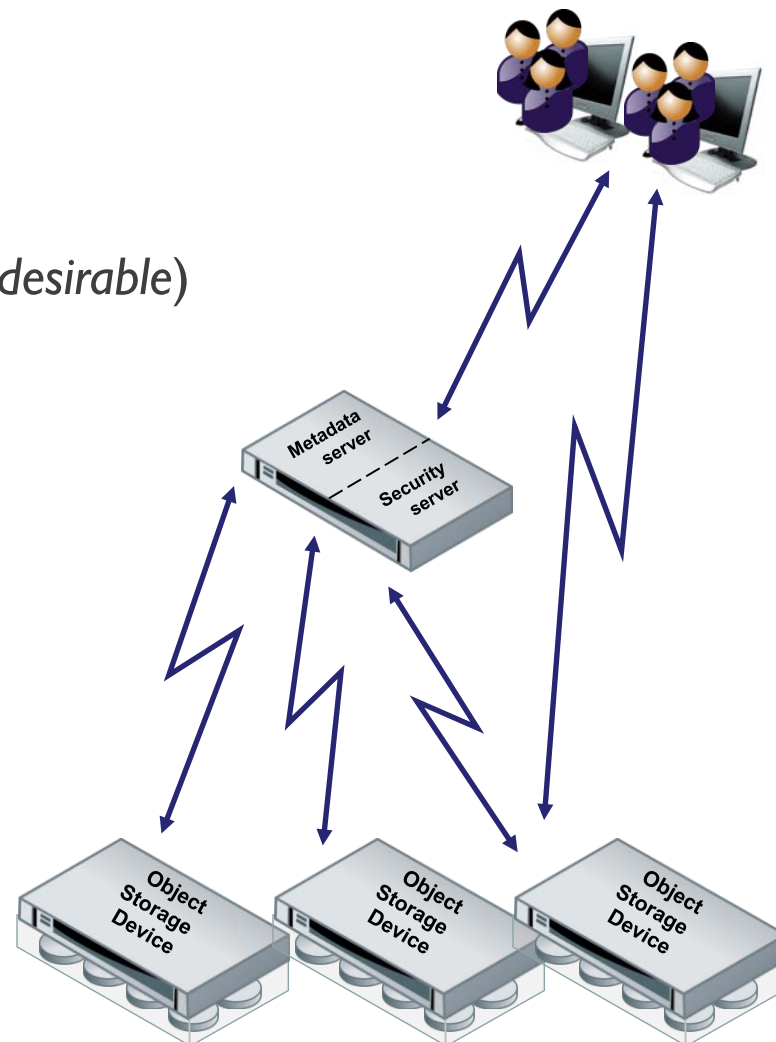
➤ Solves

- ◆ Capacity scaling
- ◆ Robustness
- ◆ Performance scaling



➤ Add OSDs

- ◆ Increase total system capacity
- ◆ Support bigger files:
(files can span OSDs if necessary or desirable)

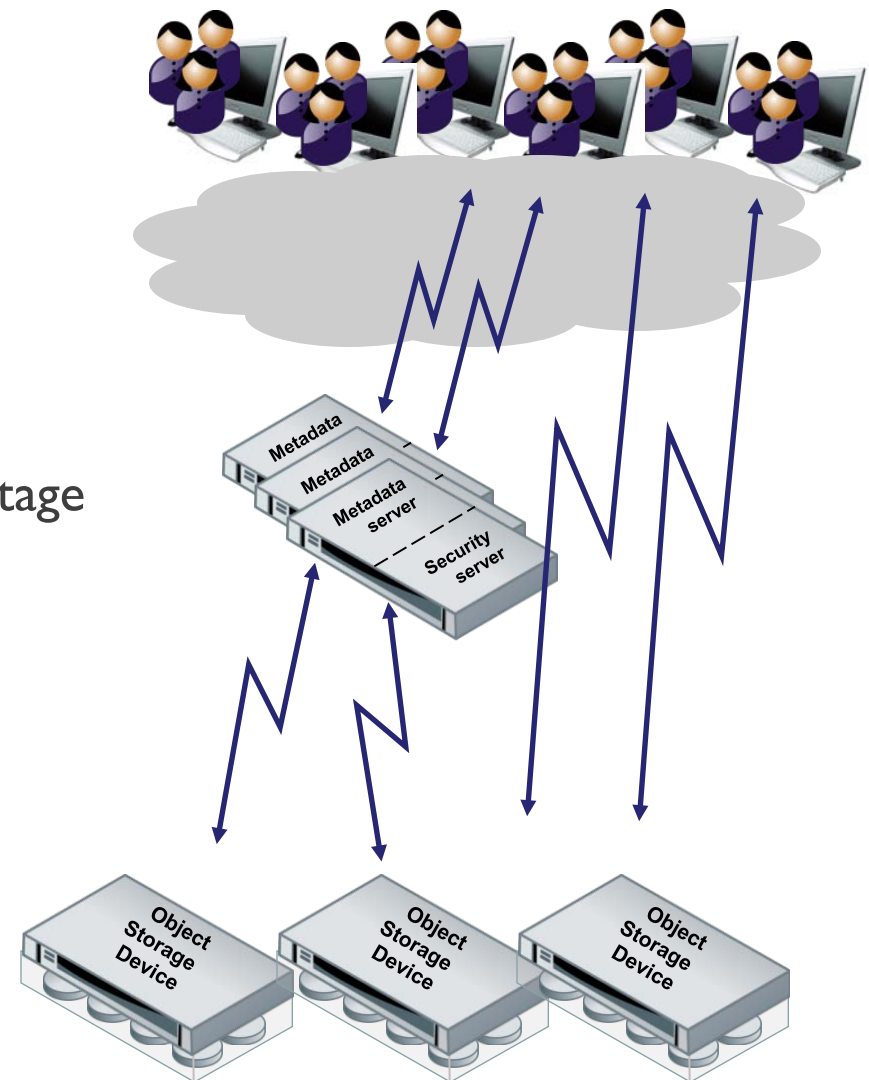


➤ Add metadata servers

- ◆ Resilient metadata services
- ◆ Resilient security services

➤ Add OSDs

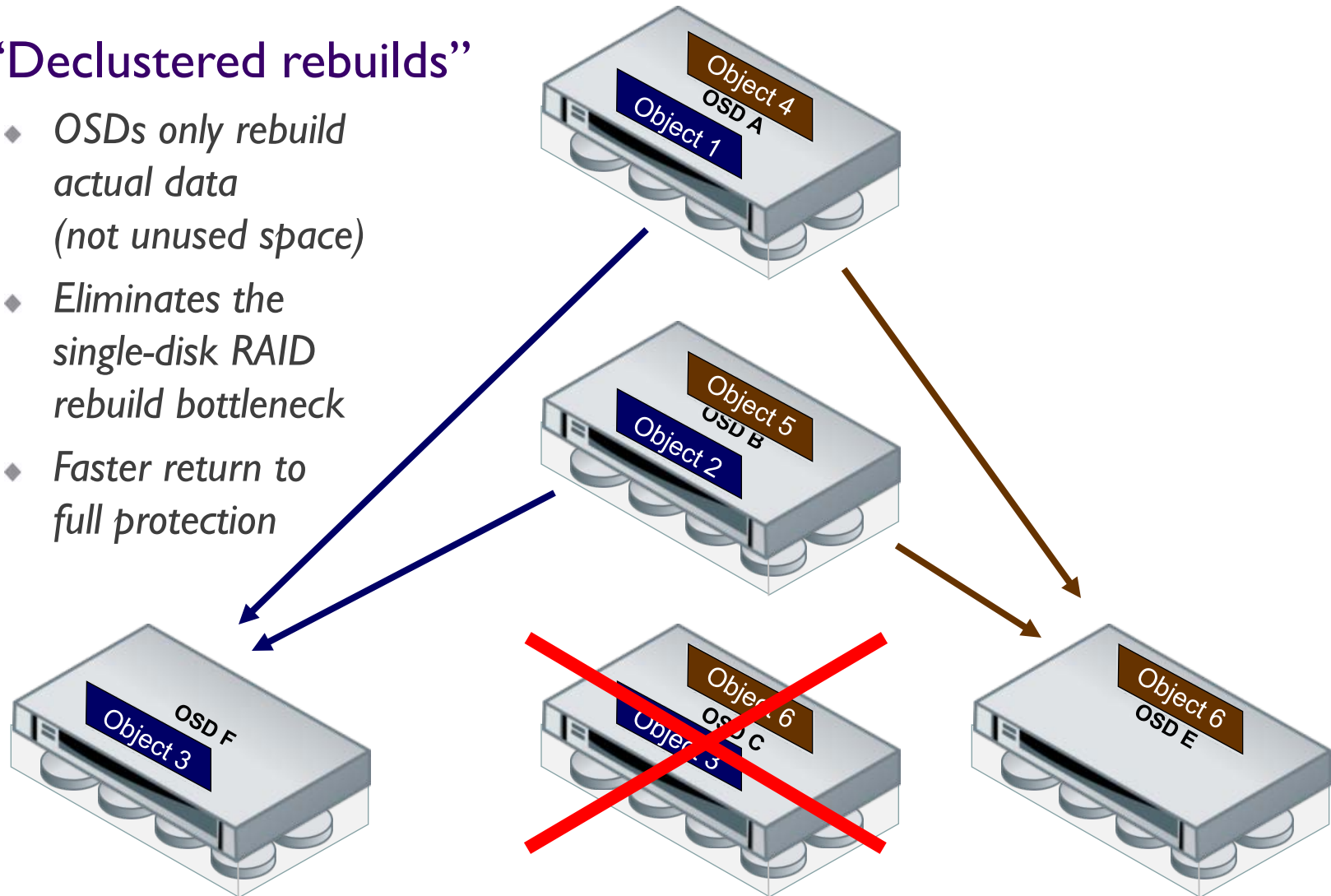
- ◆ Failed OSD affects smaller percentage of system resources
- ◆ Inter-OSD mirroring and RAID
- ◆ Near-online file system checking



An important advantage

➤ “Declustered rebuilds”

- ◆ *OSDs only rebuild actual data (not unused space)*
- ◆ *Eliminates the single-disk RAID rebuild bottleneck*
- ◆ *Faster return to full protection*

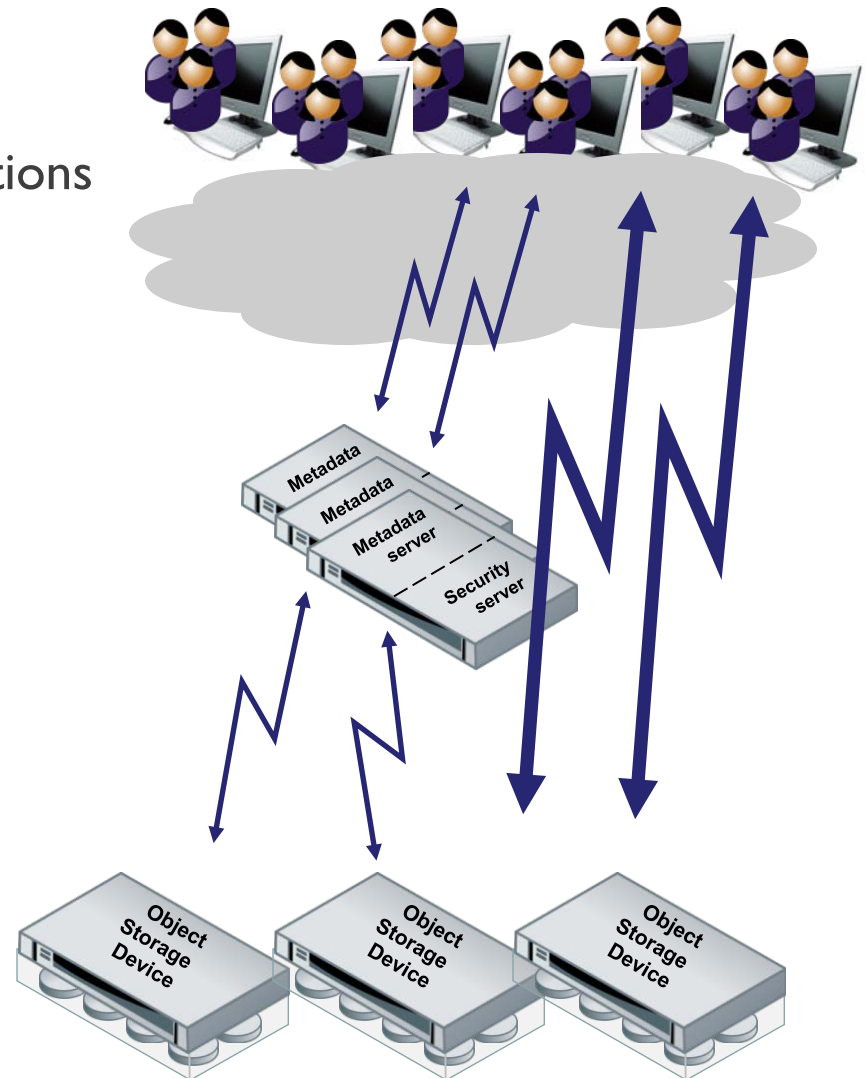


➤ Add metadata servers

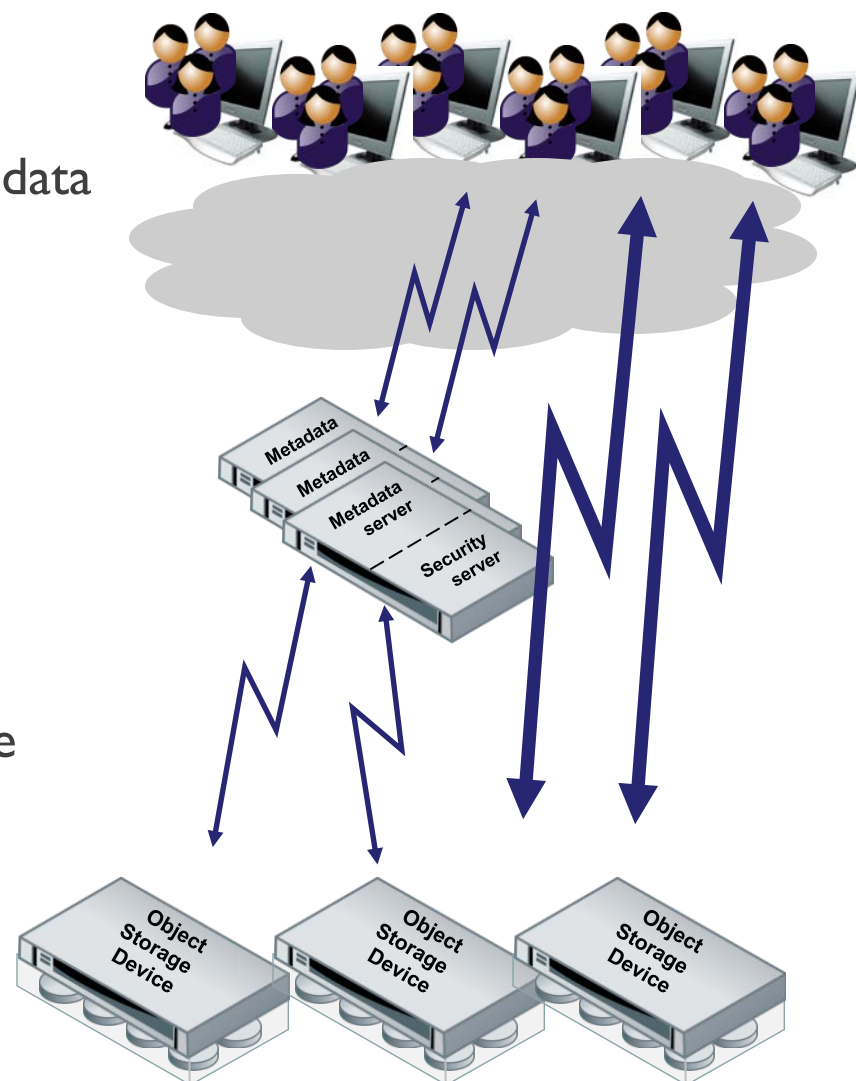
- ◆ More concurrent metadata operations (Getattr, Readdir, Create, Open,...)

➤ Add OSDs

- ◆ More concurrent I/O operations
- ◆ More bandwidth directly between clients and data

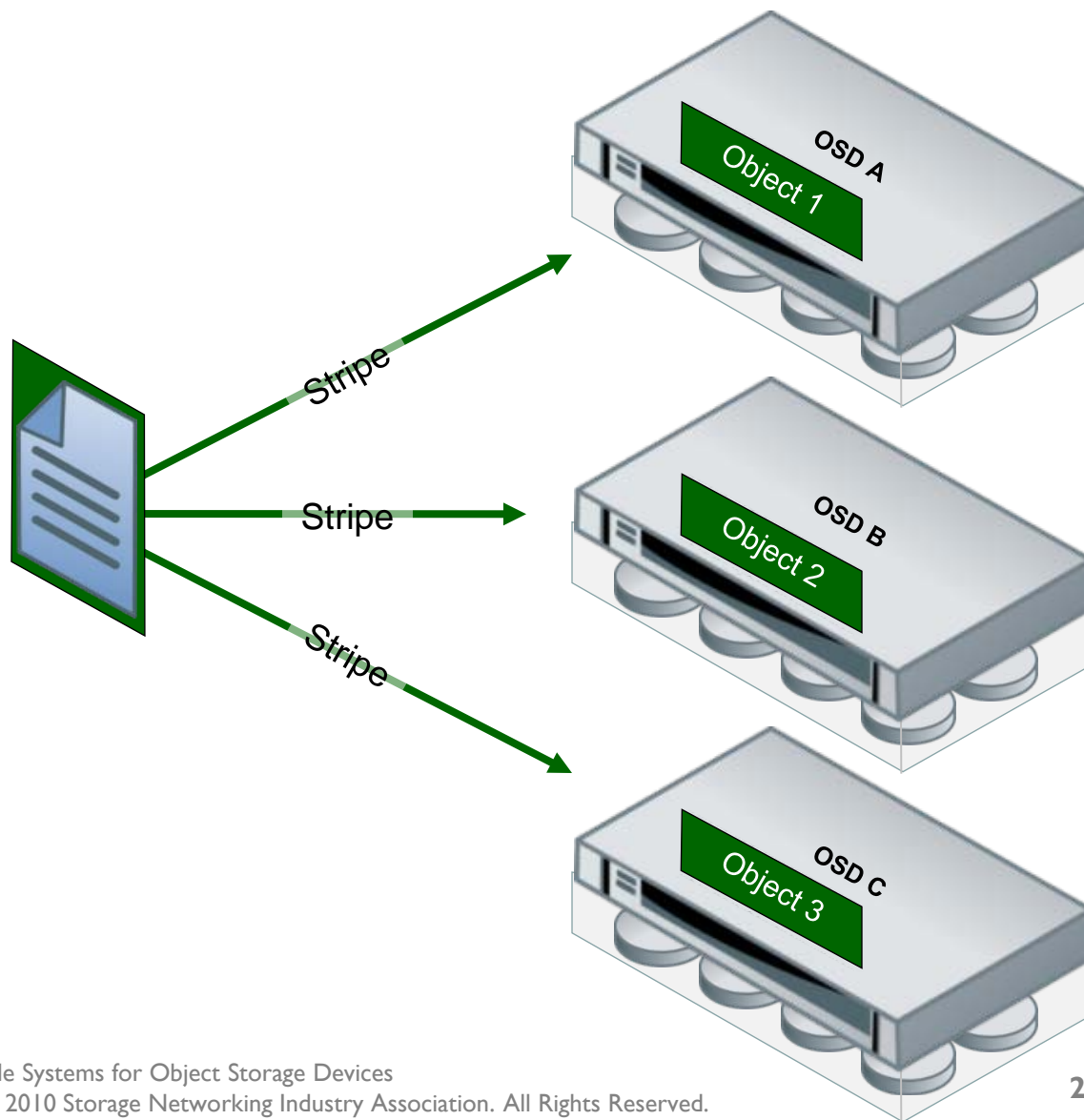


- Optimal data placement
 - ◆ Within OSD: proximity of related data
 - ◆ Load balancing across OSDs
- System-wide storage pooling
 - ◆ Across multiple file systems
- Storage tiering
 - ◆ **Per-file control** over performance and resiliency

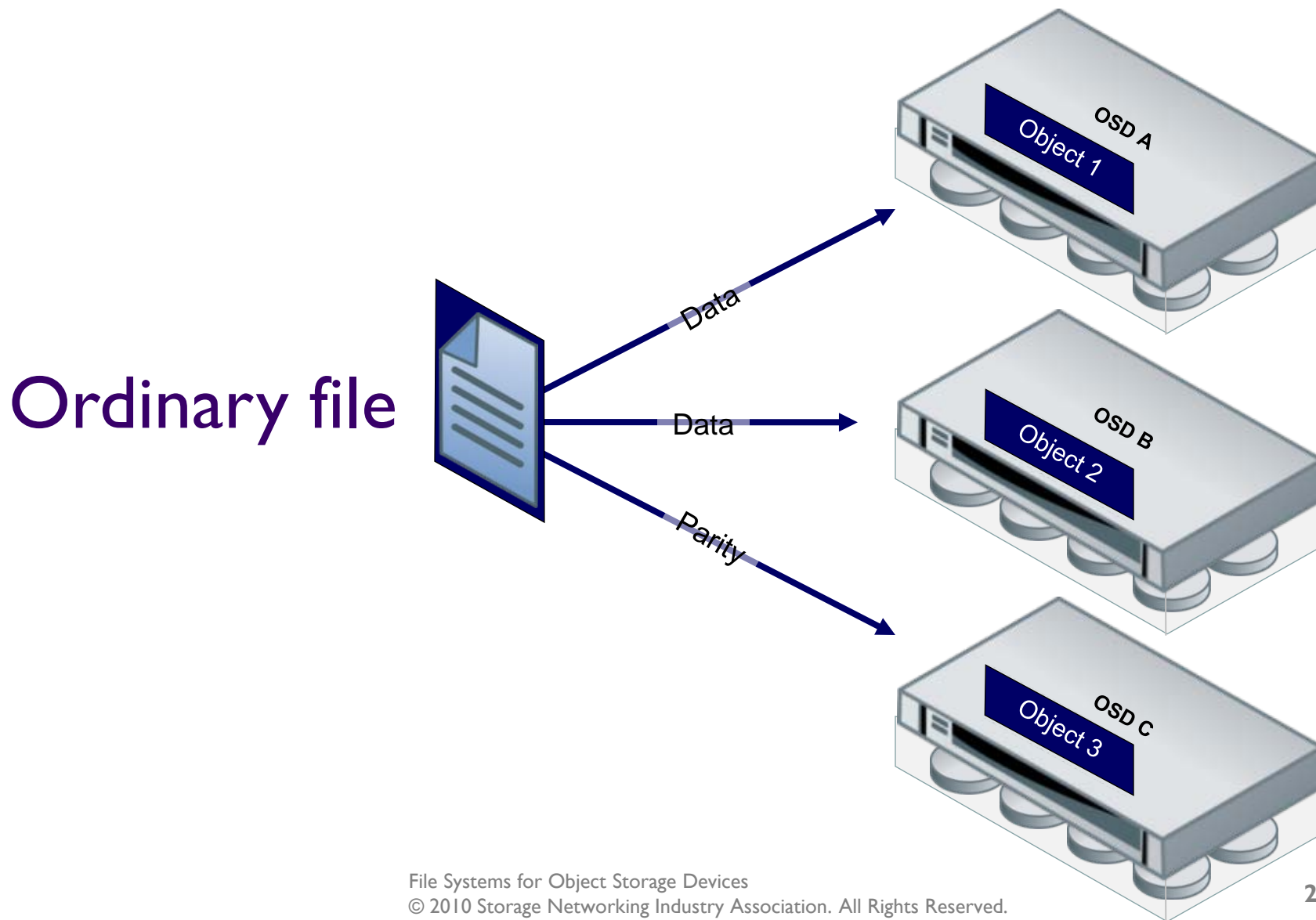


Per-file tiering: striping

Scratch file

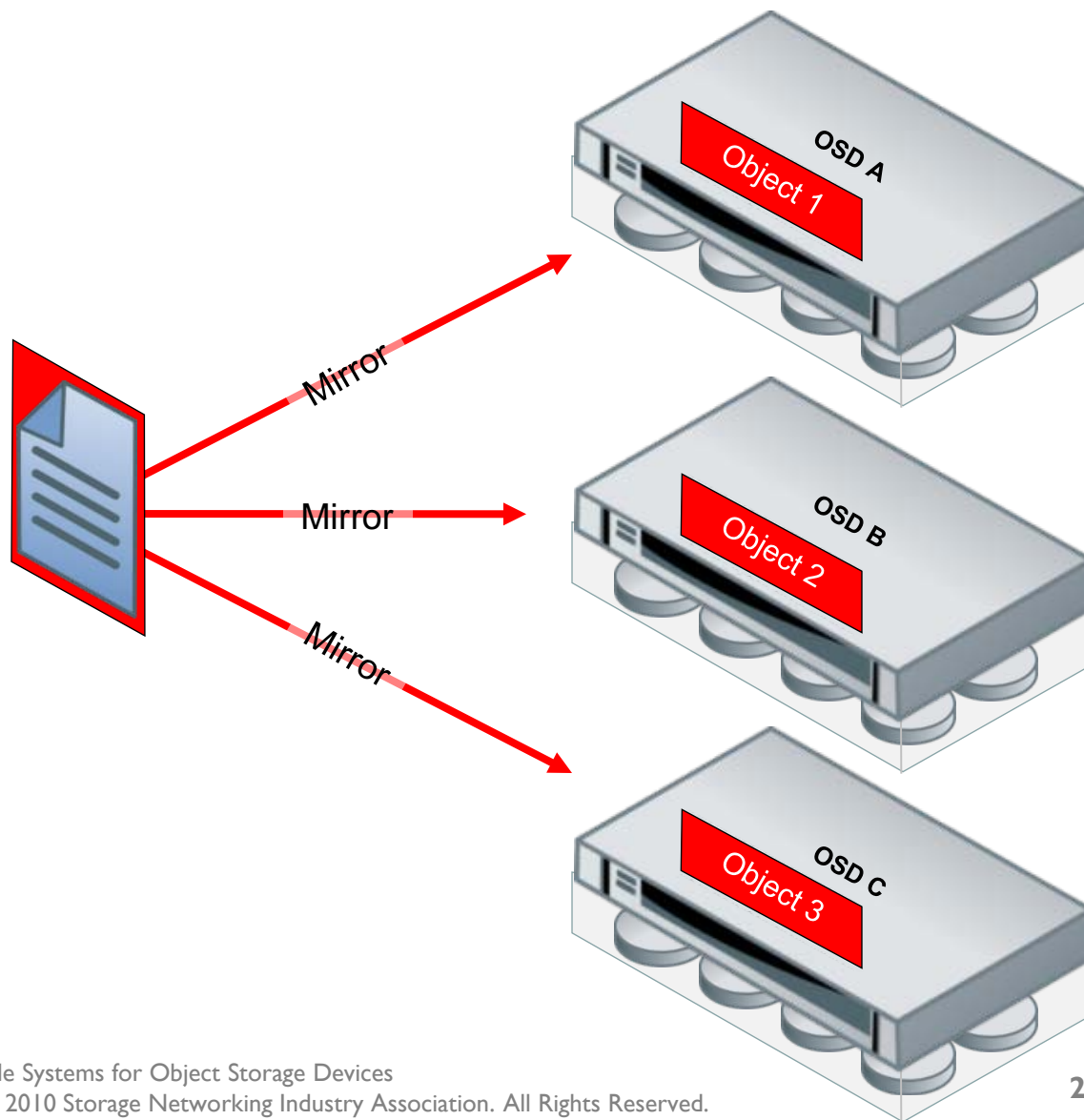


Per-file tiering: RAID-5



Per-file tiering: mirroring

Critical file



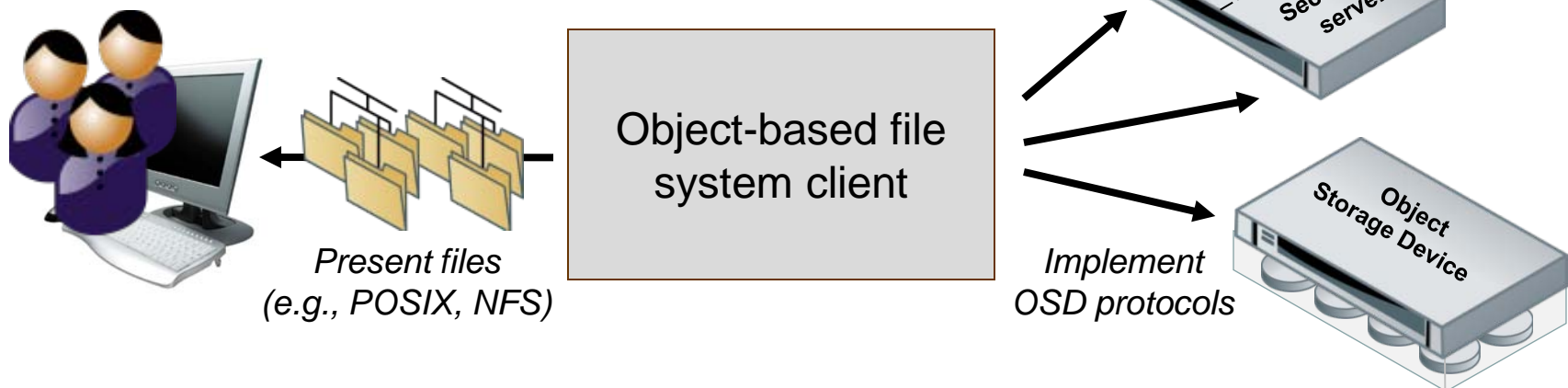
Accessing OSD-based file systems

➤ It's not

- ◆ SCSI
- ◆ NFS/CIFS

➤ Needs a “client component”

- ◆ Proprietary
- ◆ Standard

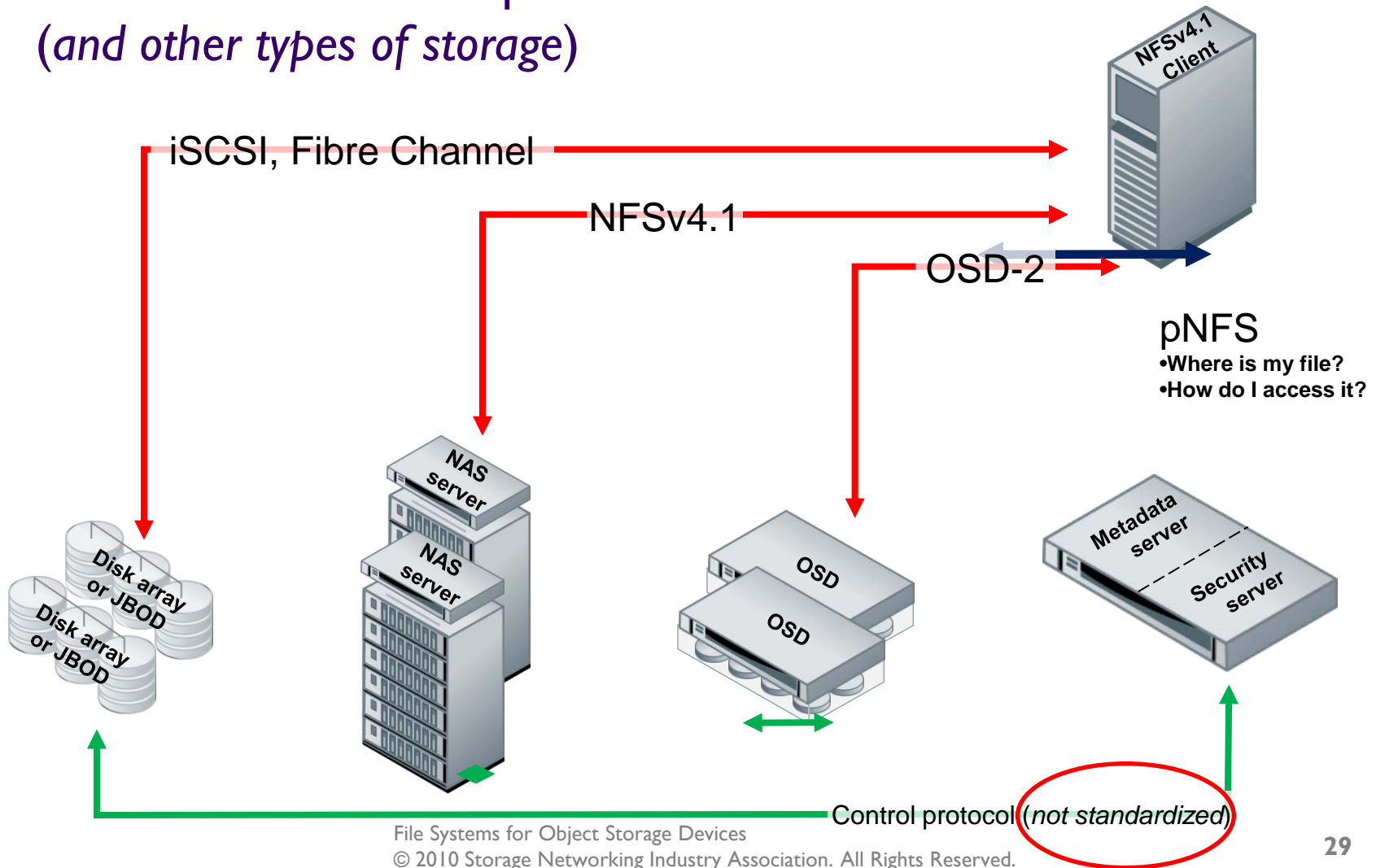


NFS v4.1 (aka pNFS)



Check out **SNIA Tutorial:**
pNFS

- A standard file access protocol for OSD
(and other types of storage)



*OK, so you buy
the concept*

➤ Two basic classes of OSD-based storage systems

- ◆ Server with embedded disks
- ◆ Software-based OSD functionality



- ◆ Purpose-built OSD “bricks”
- ◆ Highly-integrated system



Which type is optimal depends on what's important to you

- Be clear about the problem(s) you are trying to solve
 - ◆ Massive capacity?
 - ◆ Very large single name space?
 - ◆ Secure distributed access to files?
 - ◆ File integrity (AV, DLP,...)?
 - ◆ Extreme resiliency?
 - ◆ Network simplification?
 - ◆ Performance you can't get from your current architecture?
 - ◆ Administrative simplicity you can't get from your current architecture?

➤ Be clear about your IT organization's capabilities

- ◆ Favor turn-key solutions vs. integration of pieces?
- ◆ Conservative or open to new technologies and techniques?
- ◆ Experience with large-scale computing and data?
- ◆ SAN or NAS oriented?
- ◆ Network expertise?

The SNIA web site slide deck for this tutorial contains an OSD storage system checklist

- OSDs: new technology that enables file storage systems of extreme
 - ◆ Scale
 - ◆ Robustness
 - ◆ Performance

... based on custom or low-cost “commodity” components

- OSDs require specialized file systems
 - ◆ Metadata/security servers
 - ◆ Client software (drivers)

- OSD-based storage system user community
 - ◆ Initially high-performance computing
 - ◆ Today: being adopted for data-intensive applications throughout the enterprise storage market
 - › Financial services, telecom, biotech, oil & gas, aerospace, semiconductor

- As a promising new technology, it deserves a (careful) look

- Please send any questions or comments on this presentation to SNIA: trackfilemgmt@snia.org

**Many thanks to the following individuals
for their contributions to this tutorial.**

- SNIA Education Committee

**Tushar Tambay
Julian Satran
Rich Ramos
Erik Riedel
Mike Mesnier
Ralph Weber**

Appendix

Further Reference

➤ Academic research

- ◆ www.pdl.cmu.edu
- ◆ www.dtc.umn.edu

➤ Standards work

- ◆ www.snia.org/apps/org/workgroup/osd
 - › www.tl0.org/drafts.htm#OSD_Family
- ◆ datatracker.ietf.org/wg/nfsv4/charter
 - › www.rfc-editor.org/rfc/rfc5661.txt

➤ Industry research & development

- ◆ www.sun.com/lustre
- ◆ www.opensolaris.org/os/project/nfsv4/
- ◆ www.panasas.com
- ◆ www.open-osd.org
- ◆ ceph.newdream.net

OSD storage system checklist

Vendor / product maturity	How does the vendor rate ?	How important is it to me ?
Years of field experience		
Number of versions / updates / component refreshes		
Number of installations		
Largest (smallest) installation		
Geographic coverage		
Dominant applications / data access profiles		
Types of client platforms installed		
Types of backbone networks installed		

OSD storage system checklist

Flexibility	How does the vendor rate ?	How important is it to me ?
Largest (and smallest) supported configuration		
Increments of expansion (e.g., <i>capacity, cache, metadata processing, bandwidth</i>)		
Support for “mix-n-match” components (e.g., <i>different generations, disk capacities,...</i>)		
Standards compliance and alternate component sources		

OSD storage system checklist

Availability / data protection	How does the vendor rate ?	How important is it to me ?
Data protection models (e.g., mirror, RAID5-6,...)		
Fault protection domains (e.g., disks, OSDs, network links & switches...)		
Sustainable combinations of faults		
Monitoring, fault detection and notification mechanisms		
Self-healing		
Component hot-swap		
Online hardware upgrade		
Rolling software upgrade		

OSD storage system checklist

Functionality	How does the vendor rate ?	How important is it to me ?
Client access networks and features (e.g., <i>link aggregation</i>)		
Client access protocols (e.g., <i>NFS, REST,...</i>)		
Application semantics (e.g., <i>POSIX, NFS,...</i>)		
Automatic load / capacity balancing		
Advanced functions (e.g., <i>backup integration, archiving, duplicate elimination, anti-virus, data loss prevention, data classification,...</i>)		
Integration with system management frameworks		
Disaster recoverability (<i>remote replication and failover</i>)		

OSD storage system checklist

Security	How does the vendor rate ?	How important is it to me ?
Client authentication		
LDAP/AD support		
Per-client/per-user data access authorization		
Access control lists		
Digitally signed capabilities		
Protection against misbehaving clients		
Data encryption on the network		
Data encryption on media		