

CSE 487/587

Data Intensive Computing

Lecture 17: SAN/NAS

Part 1 of Storage

Lot of info borrowed from EMC, MS

Vipin Chaudhary
vipin@buffalo.edu

716.645.4740
305 Davis Hall

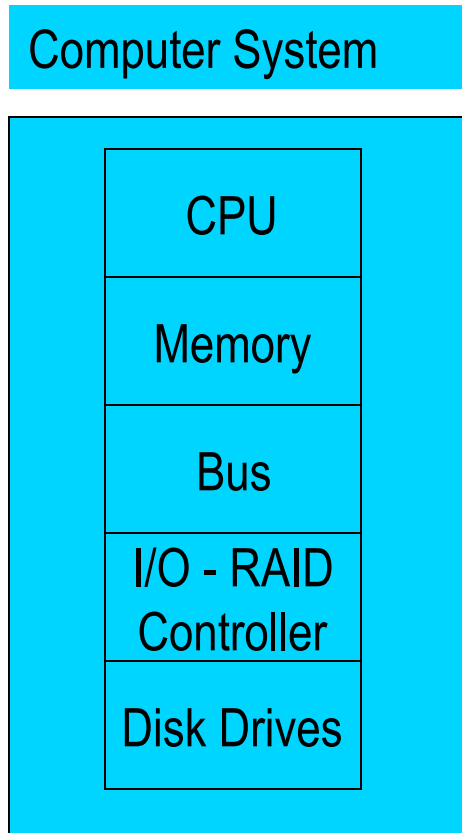
Overview of Lecture Series

- Basics – DAS, SAN, NAS
- RAID
- Erasure codes
- SSDs
- Newer Technologies

A Few Storage Basics....

- Where will data finally end up?
- How will it get there?
- What will it pass through?

Direct Attached Storage (Internal)

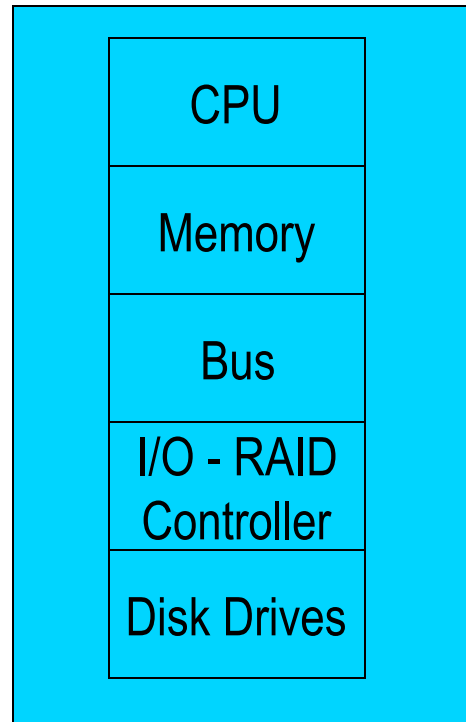


Direct Attached Storage (Internal)

305 Davis Hall
SUNY
Data Intensive
cse 487/587

Data

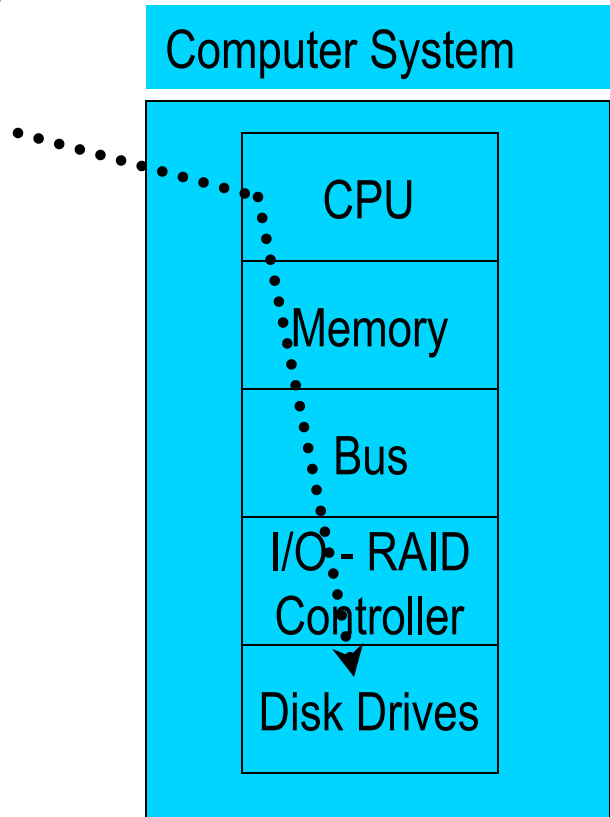
Computer System



Direct Attached Storage (Internal)

305 Davis Hall

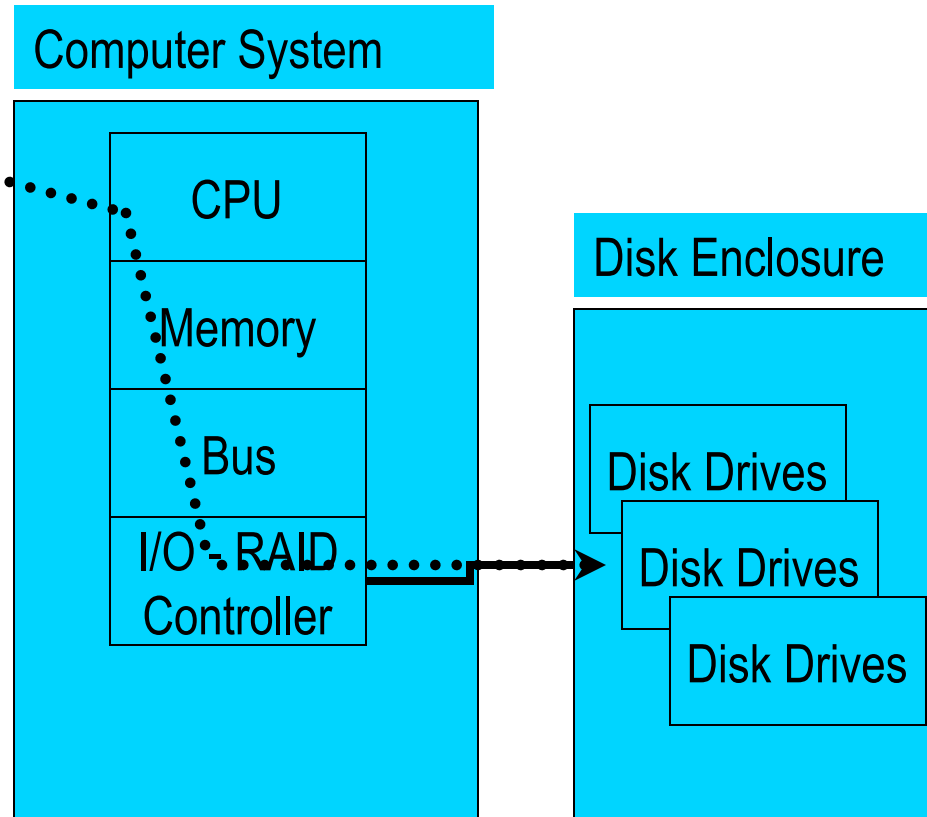
SUNY
Data Intensive
CSE 487/587



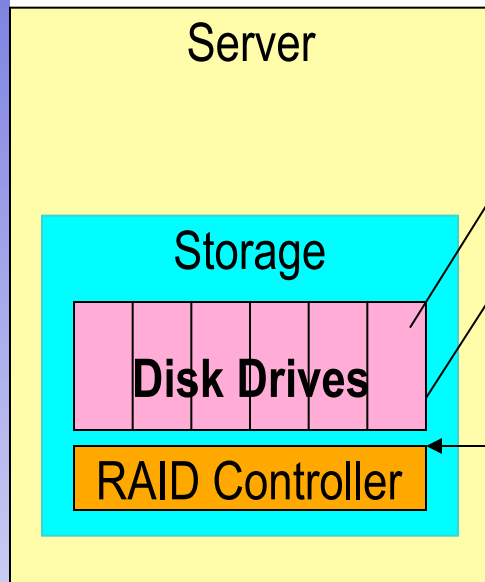
DAS w/ internal controller and external storage

305 Davis Hall

SUNY
Data Intensive
cse 487/587



Comparing Internal and External Storage



Internal Storage

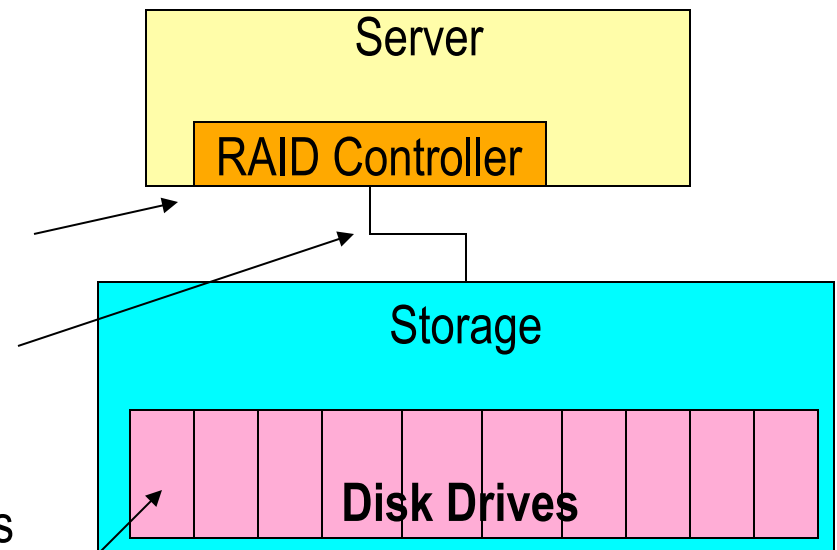
RAID controllers and disk drives are internal to the server

SCSI, ATA, or SATA protocol between controller and disks

RAID controller is internal

SCSI or SATA protocol between controller and disks

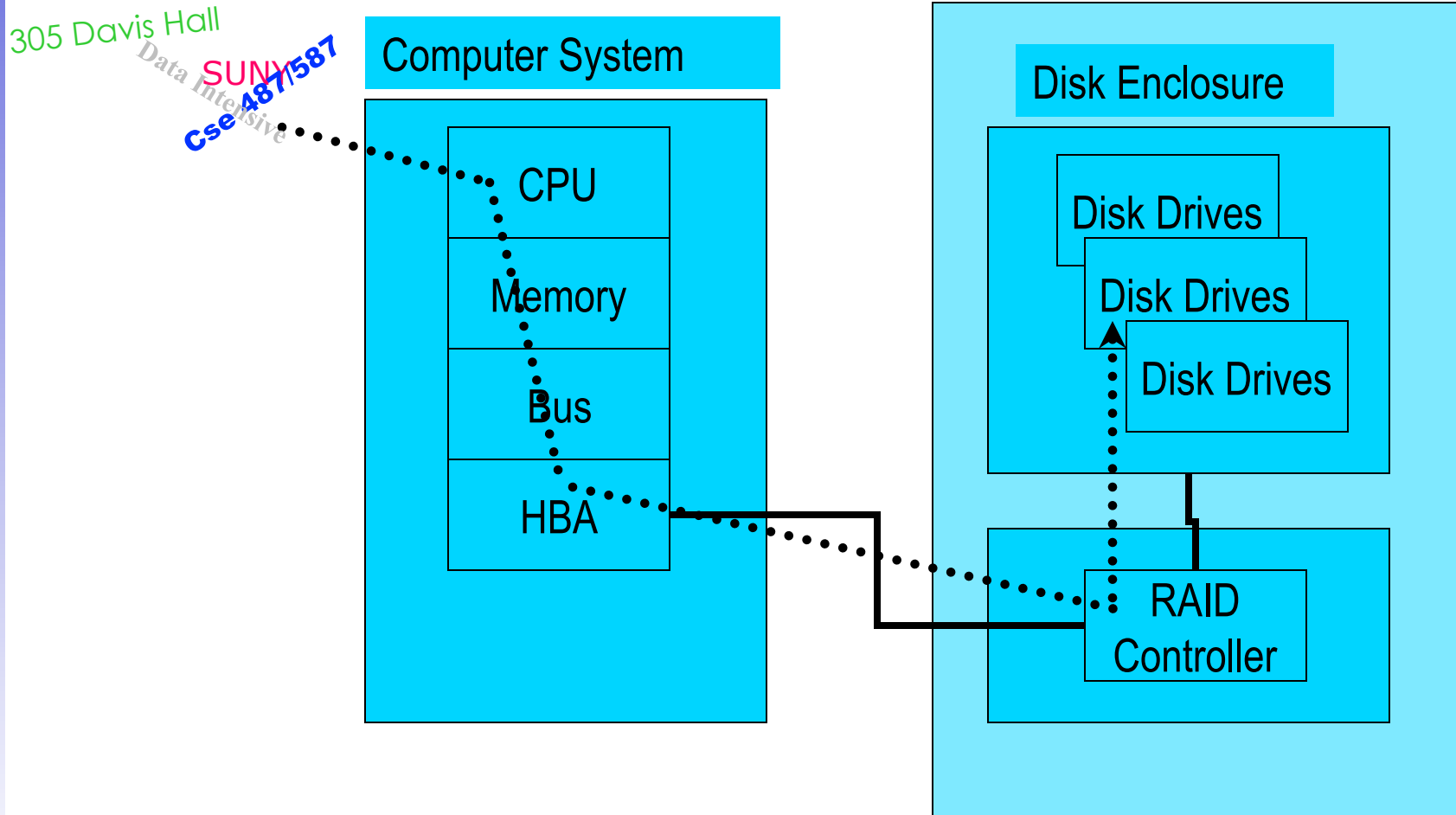
Disk drives are external



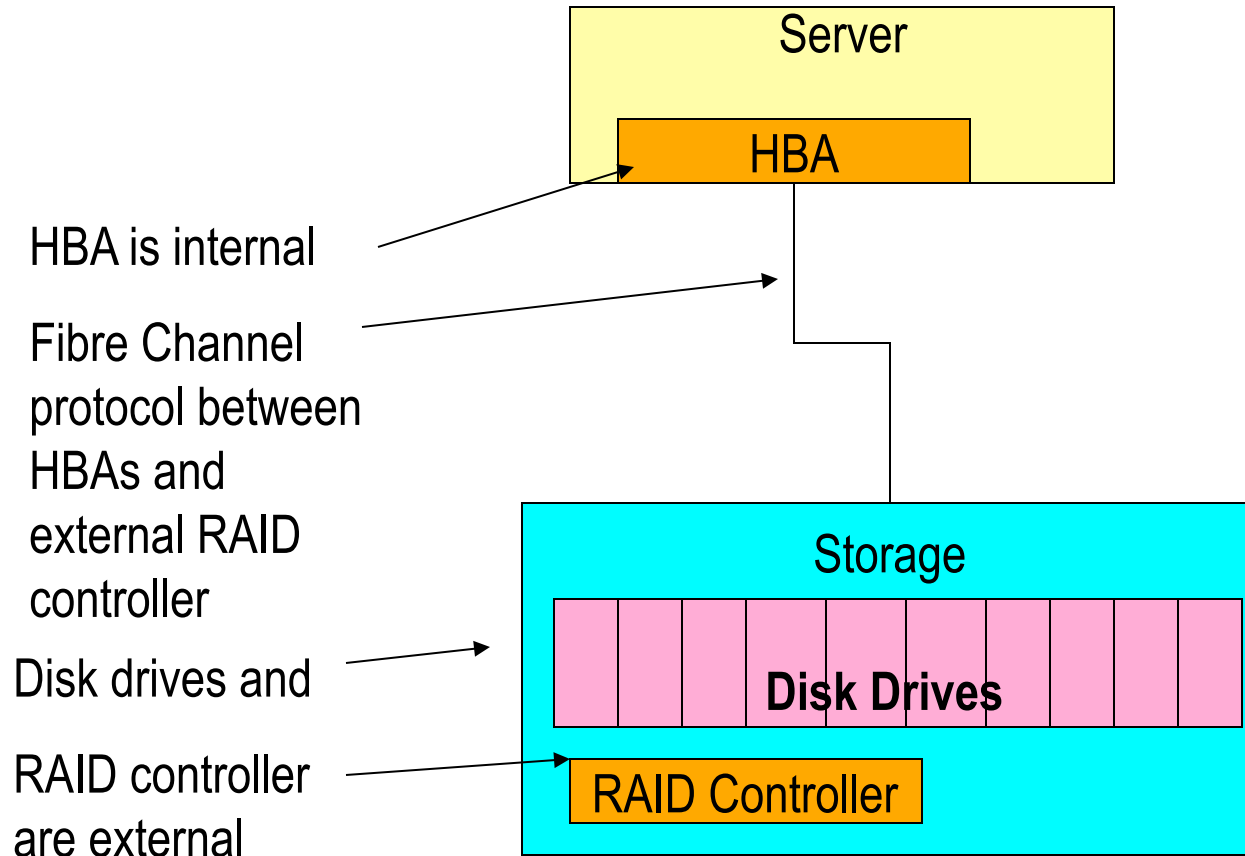
SCSI Bus w/ external storage

DAS w/ external controller and external storage

Storage System



DAS over Fibre Channel



External SAN Array

I/O Transfer

- RAID Controller
 - Contains the “smarts”
 - Determines how the data will be written (striping, mirroring, RAID 10, RAID 5, etc.)
- Host Bus Adapter (HBA)
 - Simply transfers the data to the RAID controller.
 - Doesn't do any RAID or striping calculations.
 - “Dumb” for speed.
 - Required for external storage.



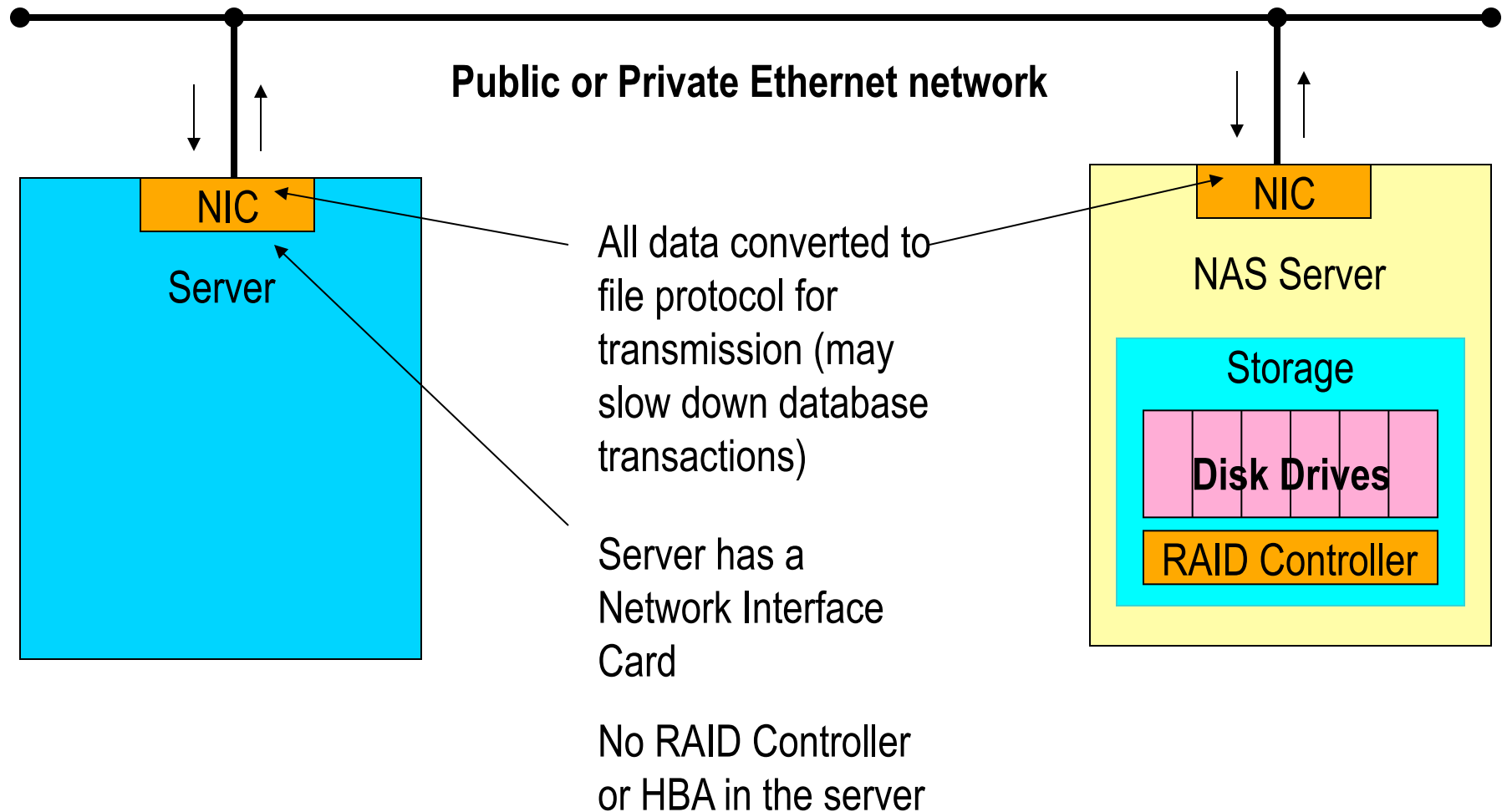
Storage types

- Single Disk Drive
- JBOD
- Volume
- Storage Array
- DAS
- NAS
- SAN

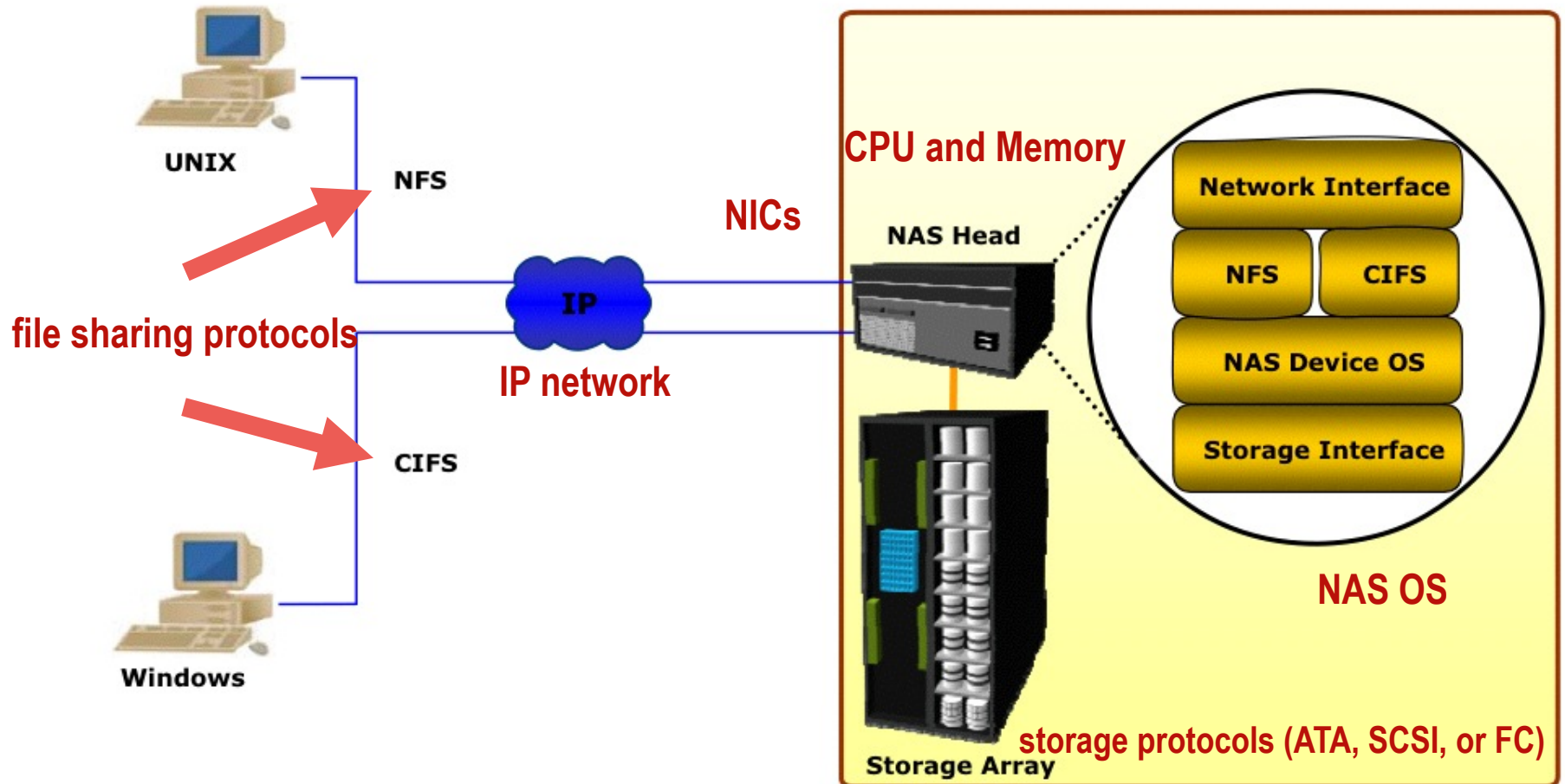
NAS: What is it?

- Network Attached Storage
- Utilizes a TCP/IP network to “share” data
- Uses file sharing protocols like Unix NFS and Windows CIFS
- Storage “Appliances” utilize a stripped-down OS that optimizes file protocol performance

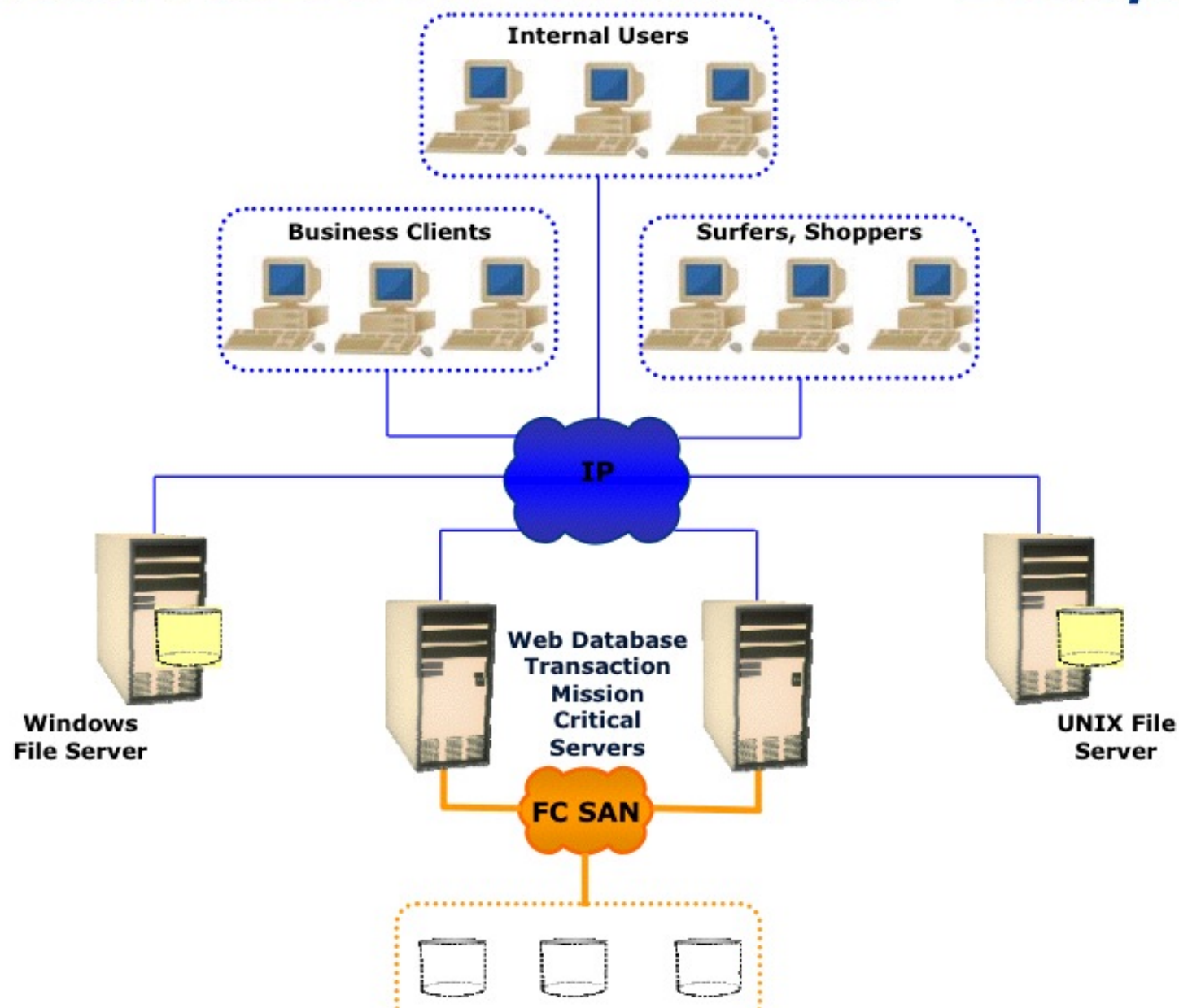
Networked Attached Storage



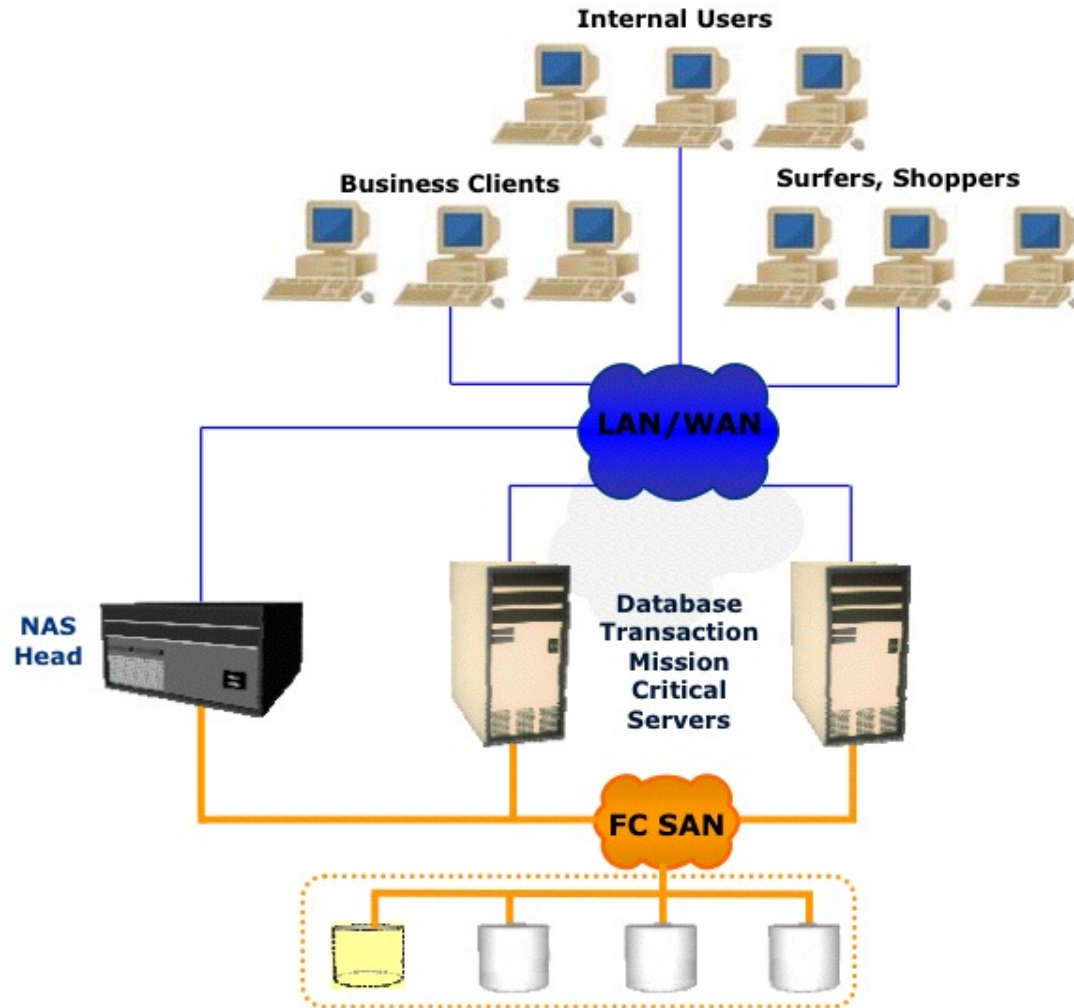
Components of NAS



Traditional File Server Environment – Example 1



Storage Consolidation with NAS

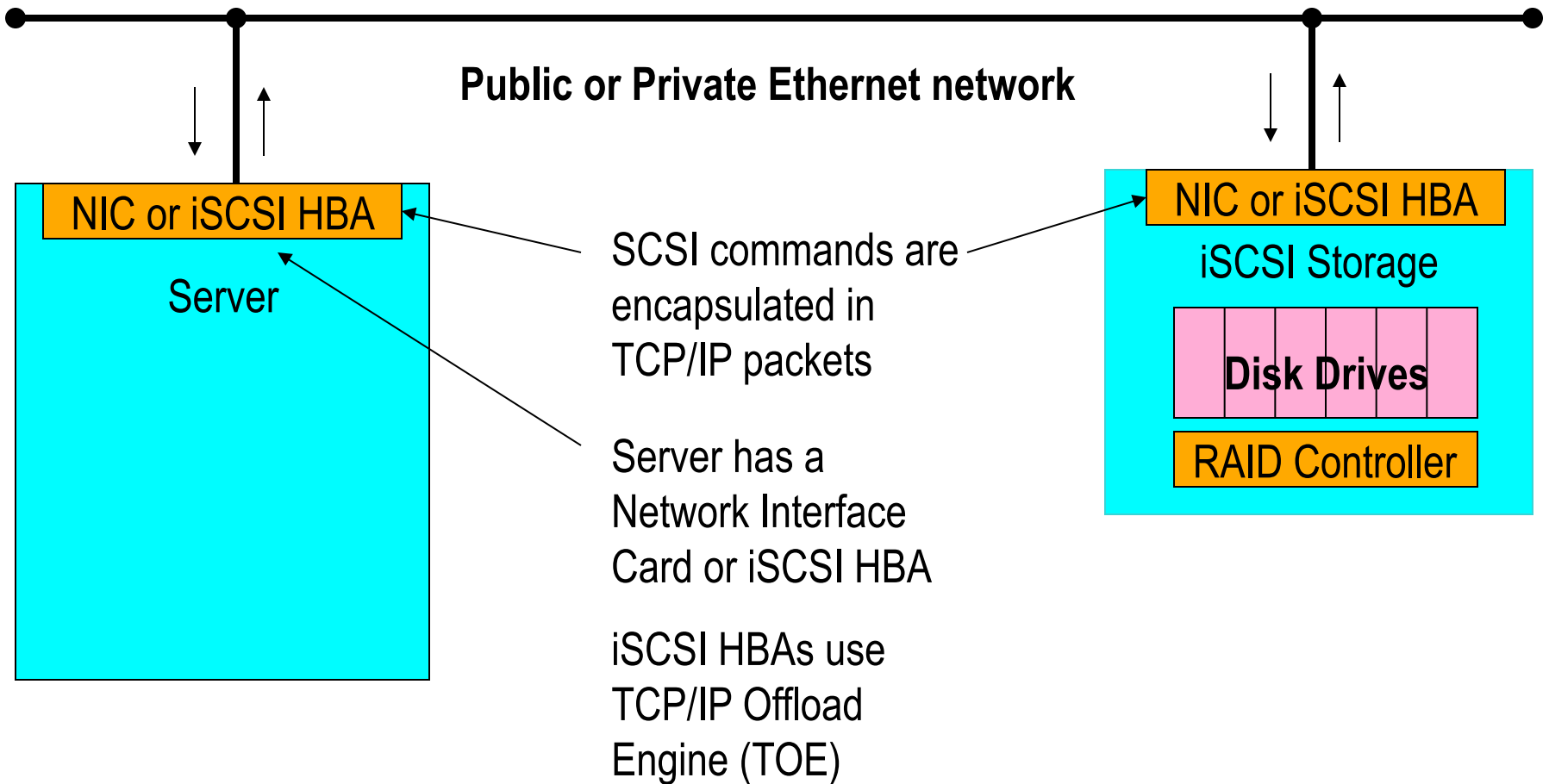


- Benefits:
- Increases performance throughput (service level) to end users
- Minimizes investment in additional servers
- Provides storage pooling
- Provides heterogeneous file servings
- Uses existing infrastructure, tools, and processes

iSCSI: What is it?

- An alternate form of networked storage
- Like NAS, also utilizes a TCP/IP network
- Encapsulates native SCSI commands in TCP/IP packets
- Supported in Windows 2003 Server and Linux
- TCP/IP Offload Engines (TOEs) on NICs speed up packet encapsulation

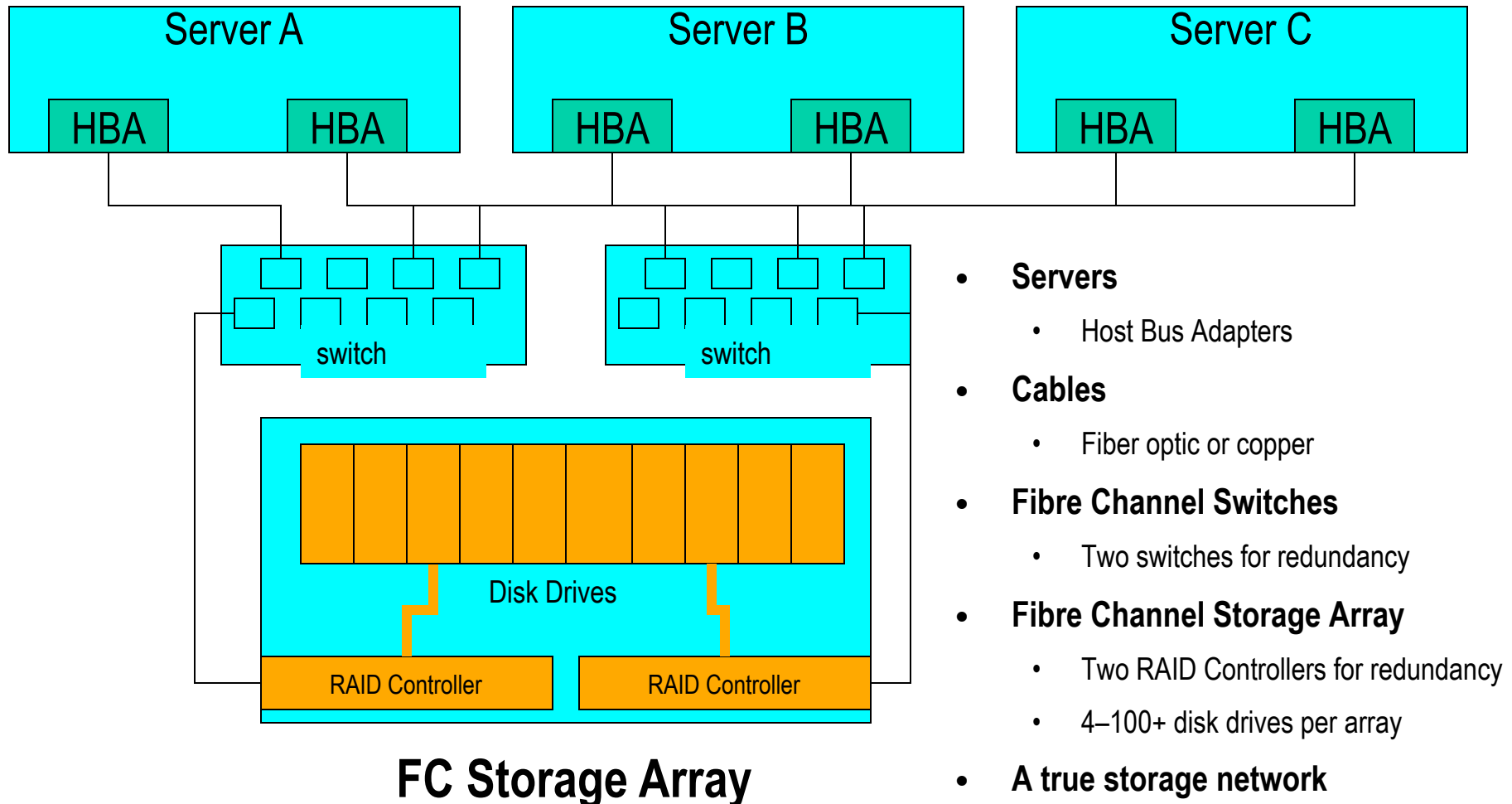
iSCSI Storage



Fibre Channel: What is it?

- Fibre Channel is a network protocol implemented specifically for dedicated storage networks
- Fibre Channel utilizes specialized
 - Switches
 - Host Bus Adapters
 - RAID controllers
 - Cables

Fibre Channel Components



- **Servers**
 - Host Bus Adapters
- **Cables**
 - Fiber optic or copper
- **Fibre Channel Switches**
 - Two switches for redundancy
- **Fibre Channel Storage Array**
 - Two RAID Controllers for redundancy
 - 4–100+ disk drives per array
- **A true storage network**
 - Multiple servers
 - Multiple switches
 - Multiple Storage Arrays

SAN: What is it?

- Storage Area Network
- A network whose primary purpose is the transfer of data between storage systems and computer systems
- Fibre Channel is the primary technology utilized for SANs
- Recently, SANs have been implemented with dedicated iSCSI networks

Benefits of SAN/Consolidated Storage

- Reduce cost of external storage (?)
- Increase performance
- Centralized and improved tape backup
- LAN-less backup
- High-speed, no single-point-of-failure clustering solutions
- Consolidation

Fibre Channel Technology

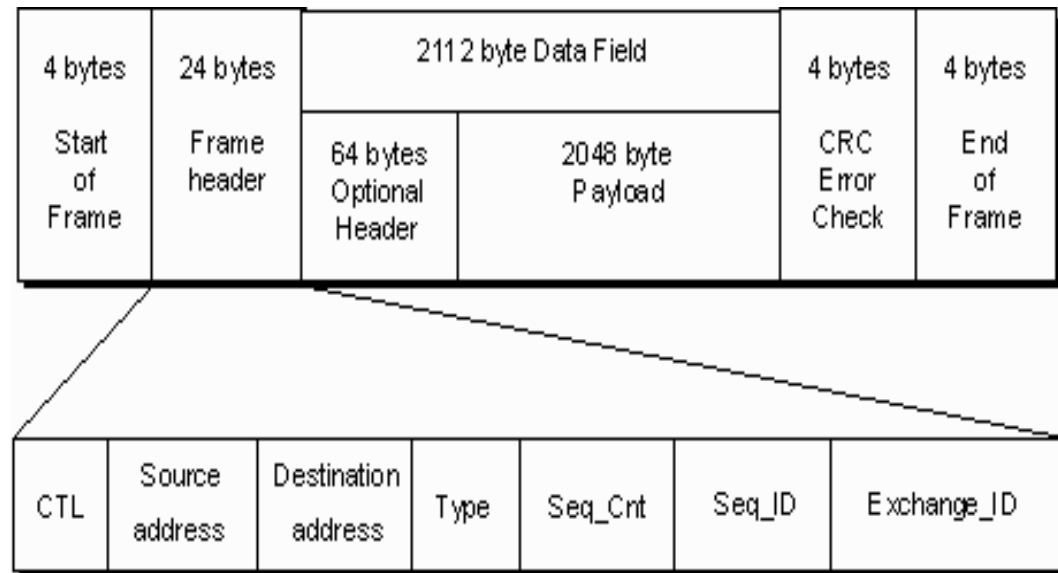
- Provides concurrent communications between servers, storage devices, and other peripherals
- A gigabit interconnect technology
 - FC1: Over 1,000,000,000 bits per second
 - FC2: Over 2,000,000,000 bits per second
 - FC16: now
 - FC128: projected in 2016?
- A highly reliable interconnect
- Up to 127 devices (SCSI: 15)
- Up to 10 km of cabling (3-15 ft. for SCSI)
- Physical interconnect can be copper or fiber optic

Fibre Channel - (continued)

- Hot-pluggable - Devices can be removed or added at will with no ill effects to data communications
- Provides a data link layer above the physical interconnect, analogous to Ethernet
- Sophisticated error detection at the frame level
- Data is checked and resent if necessary

Fibre Channel - Frame Dissection

- Up to 2048 byte payload
- 4 byte checksum for each frame



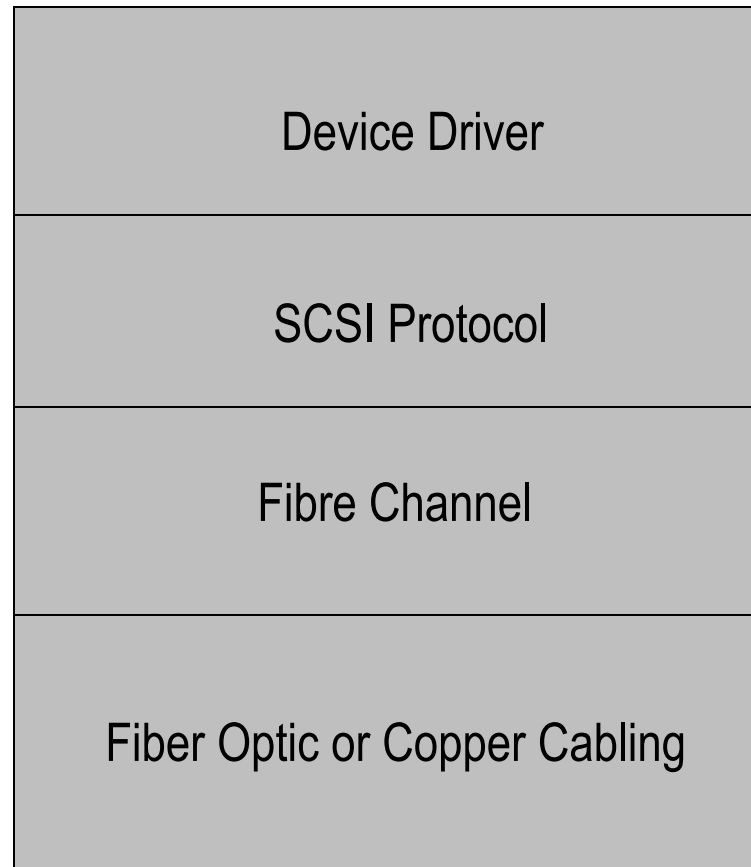
Fibre Channel

- What's with the funny name?
 - Some background history required
 - Originally developed to only support fiber optic cabling
 - When copper cabling support was added, ISO decided not to rename the technology
 - ISO changed to the French spelling to reduce association with fiber optics only medium

Fibre Channel

- How does it work?
 - Serial interface
 - Data is transferred across a single piece of medium at the fastest speed supported
 - No complex signaling required

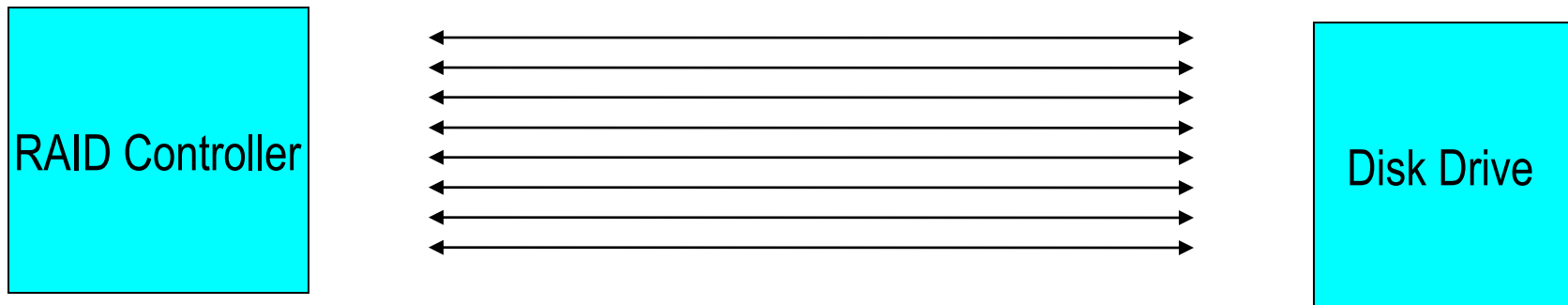
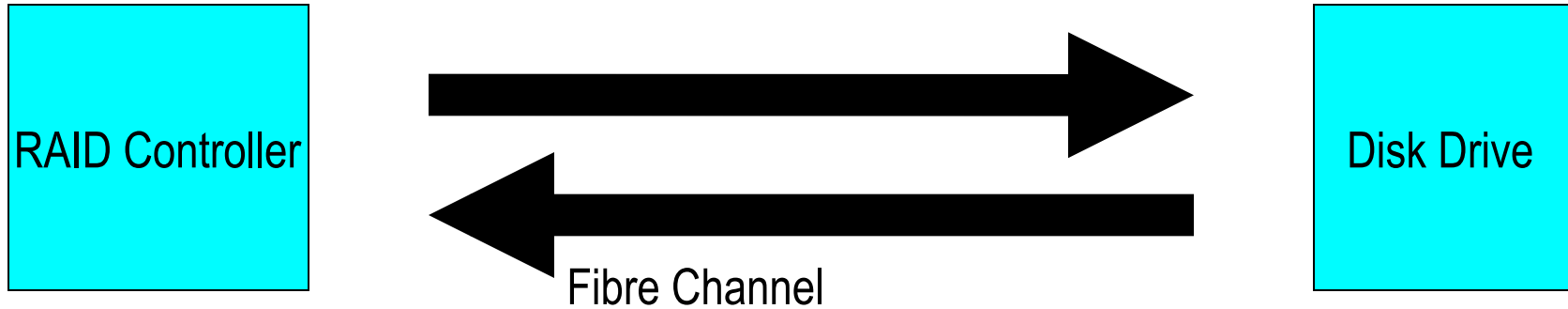
Fibre Channel Interface Layers



SCSI vs. Fibre Channel *Protocol*

- SCSI
 - SCSI protocol vs. SCSI device
 - SCSI is an established, tried and true protocol
 - Provides services analogous to TCP/IP
 - Supported in every major OS on market
- Fibre Channel
 - Fibre Channel runs on top of SCSI
 - No re-inventing the wheel
 - Immediate OS support

SCSI vs. FC Transmission



SCSI

SCSI vs. Fibre Channel

- Interface for internal storage to external disks
- Potential down time w/ SCSI
- Single bus
- RAID controller is SCSI hardware
- Standards:
 - Ultra2 (80 MB/sec)
 - Ultra 160 (160 MB/sec)
 - Ultra 320 (320 MB/sec)
- Media specific (copper only)
- SCSI Limitations:
 - Cables can't be any longer than 3 feet for single ended; 15 feet for LVD (low voltage differential)
 - No more than 15 devices on a SCSI bus
 - # of disk drives

- Used with SAN
- Lots of built-in redundancy with connections
- Redundant network
- HBA is fibre channel hardware
- Standards:
 - FC16: 16 GB/sec
 - FC128: 128 GB/sec
- Provides a data link layer above the physical interconnect
 - Analogous to Ethernet
 - FC is a network of devices
 - It can be media independent- copper or fibre optic
- Fibre Channel limitations:
 - Cable length: Up to 10 kilometers (more a limitation of cable than FC itself)
 - Up to 127 devices
 - # of disk drives

Fibre Channel vs. iSCSI

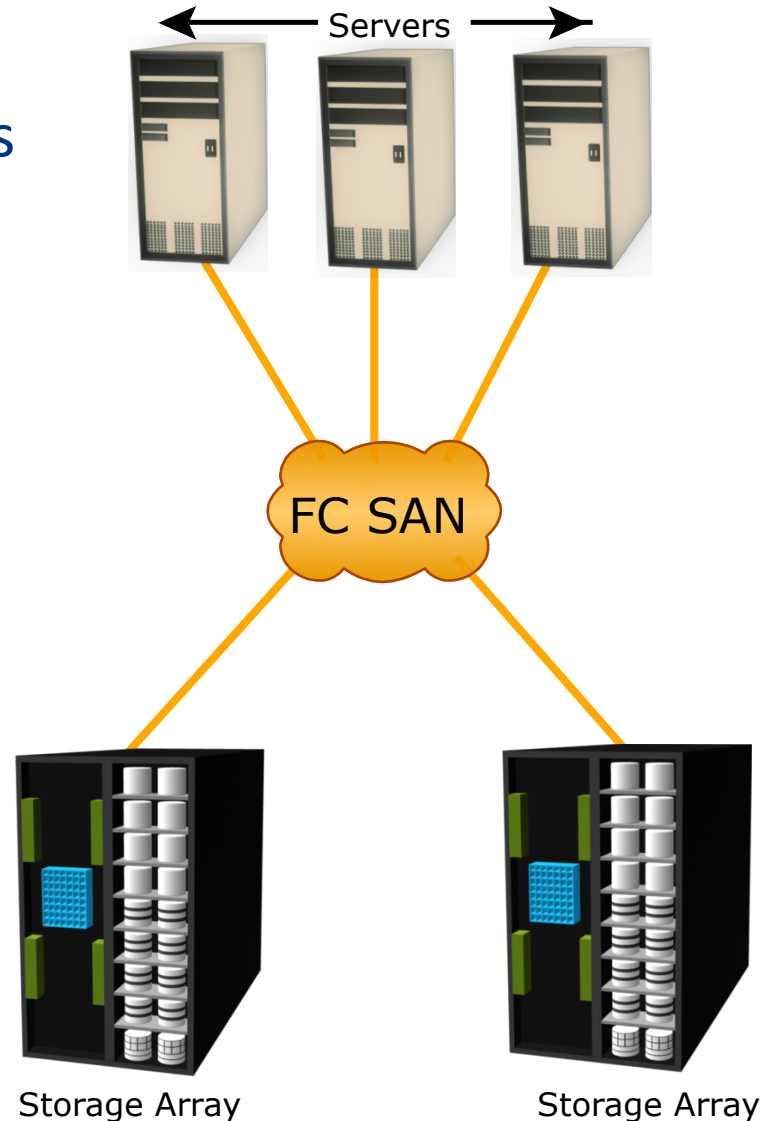
- Fibre Channel
 - The current market leader for shared storage technologies
 - Provides the highest performance levels
 - Designed for mission-critical applications
 - Cost of components is relatively high, particularly per server HBA costs
 - Relatively difficult to implement and manage
- iSCSI
 - Relatively new, but usage is increasing rapidly
 - Performance can approach Fibre Channel speeds
 - A better fit for databases than NAS
 - A good fit for Small to Medium Size Businesses
 - Relatively inexpensive, compared to Fibre Channel
 - Relatively easy to implement and manage

Summary

- How data is routed through a server to I/O
- Types of storage
 - DAS
 - NAS
 - iSCSI
 - SAN
- Benefits of SAN technology
 - Storage consolidation
 - Reduced costs
 - Centralized, LAN-free backup and restore
- The Fibre Channel protocol
 - How it works
 - Fibre Channel protocol vs. SCSI protocol
- Comparing Fibre Channel SANs and iSCSI SANs
 - Fibre Channel SANs offer mission-critical performance, with relatively high costs and high complexity
 - iSCSI SANs offer moderate to high performance at an attractive price/performance ratio and are relatively easy to administer

What is a SAN ?

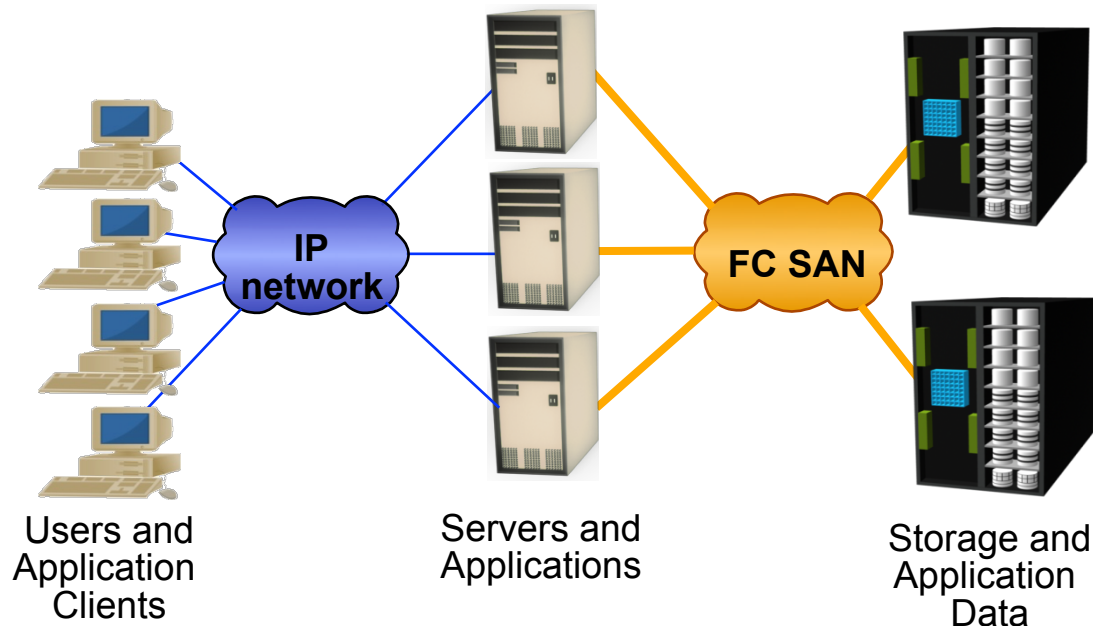
- Dedicated high speed network of servers and shared storage devices
- Provide block level data access
- Resource Consolidation
 - Centralized storage and management
- Scalability
 - Theoretical limit: Appx. 15 million devices
- Secure Access



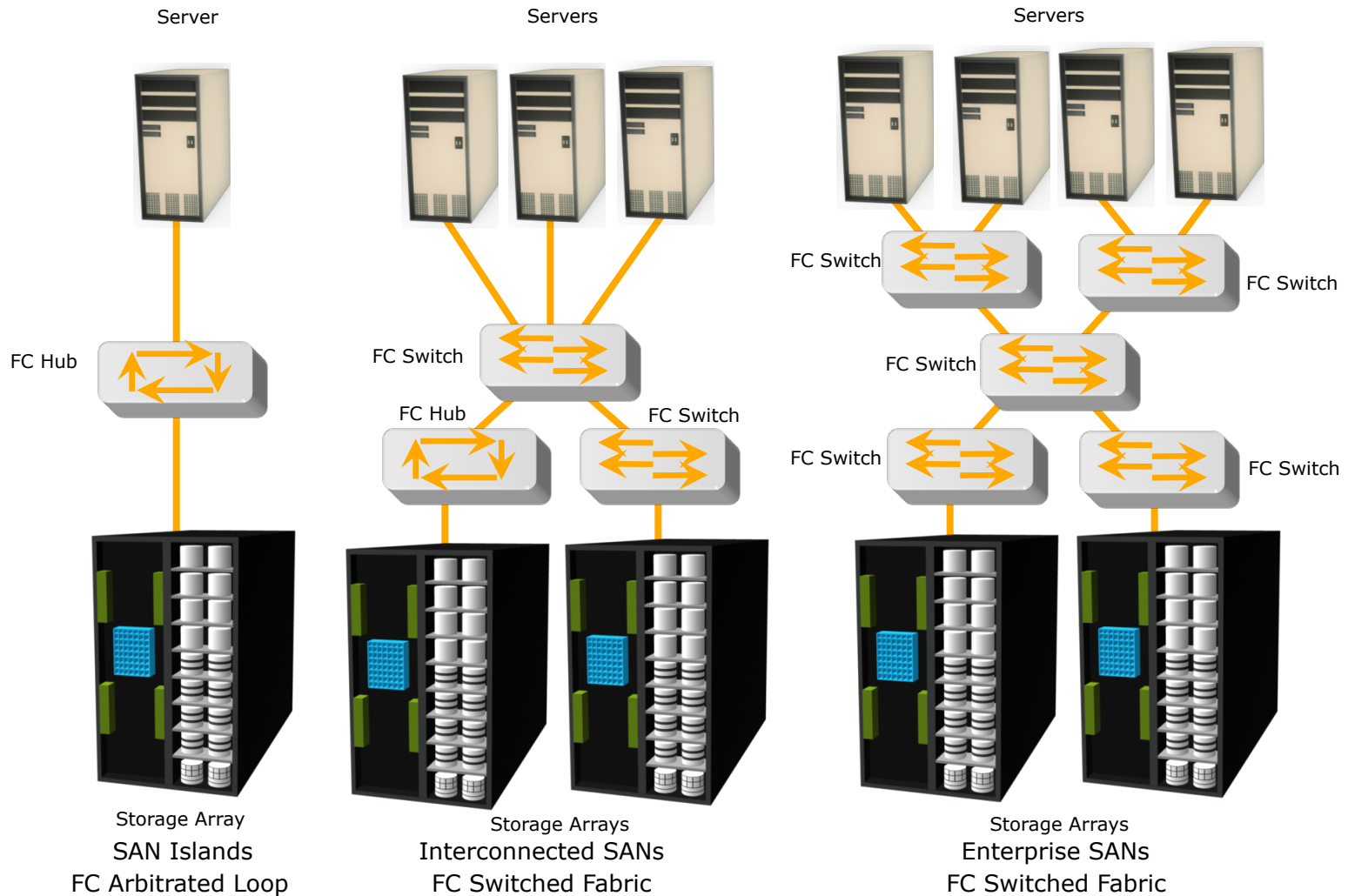
Understanding Fibre Channel

- Fibre Channel is a high-speed network technology that uses:
 - Optical fiber cables (for front end connectivity)
 - Serial copper cables (for back end connectivity)
- Latest FC implementations support 16Gb/s

o Servers are attached to 2 distinct networks



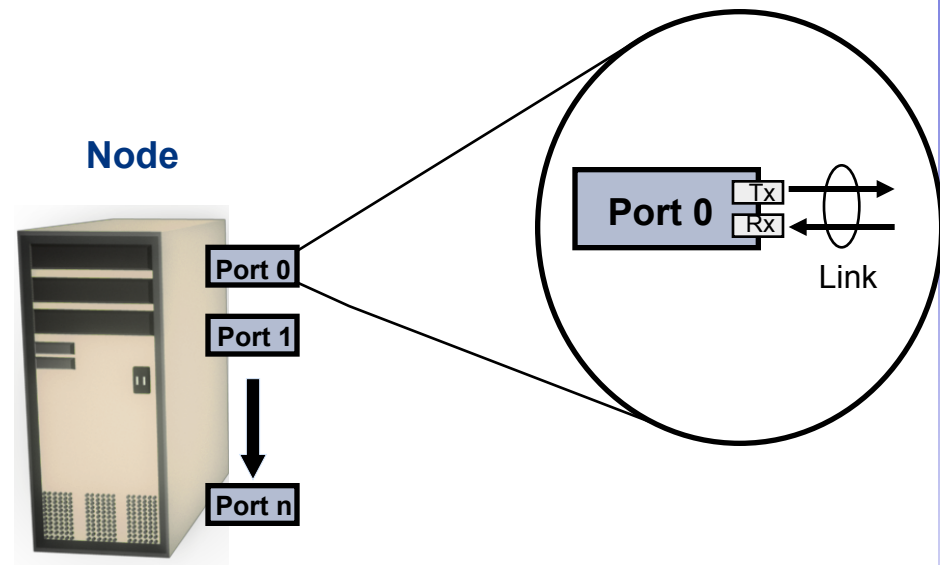
FC SAN Evolution



Fibre Channel SAN Evolution

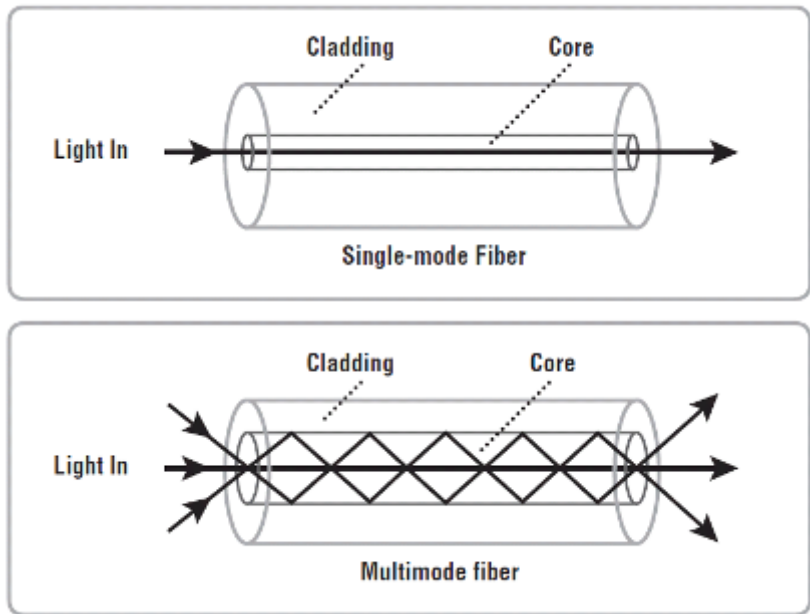
Components of SAN: Node ports

- Examples of nodes
 - Hosts, storage and tape library
- Ports are available on:
 - HBA in host
 - Front-end adapters in storage
 - Each port has transmit (Tx) link and receive (Rx) link
- HBAs perform low-level interface functions automatically to minimize impact on host performance



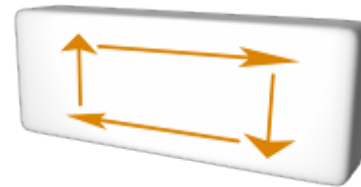
Components of SAN: Cabling

- SAN implementation uses:
 - Copper cables for short distance
 - Optical fiber cables for long distance
- Two types of optical cables
 - Single-mode
 - Can carry single beams of light
 - Distance up to 10 KM
 - Multi-mode
 - Can carry multiple beams of light simultaneously
 - Distance up to 500 meters

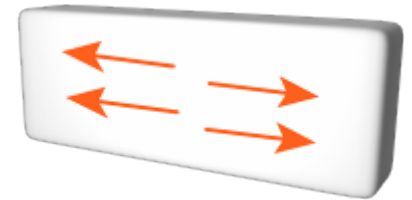


Components of SAN: Interconnecting devices

- Basis for SAN communication
 - Hubs
 - Switches and
 - Directors



FC HUB



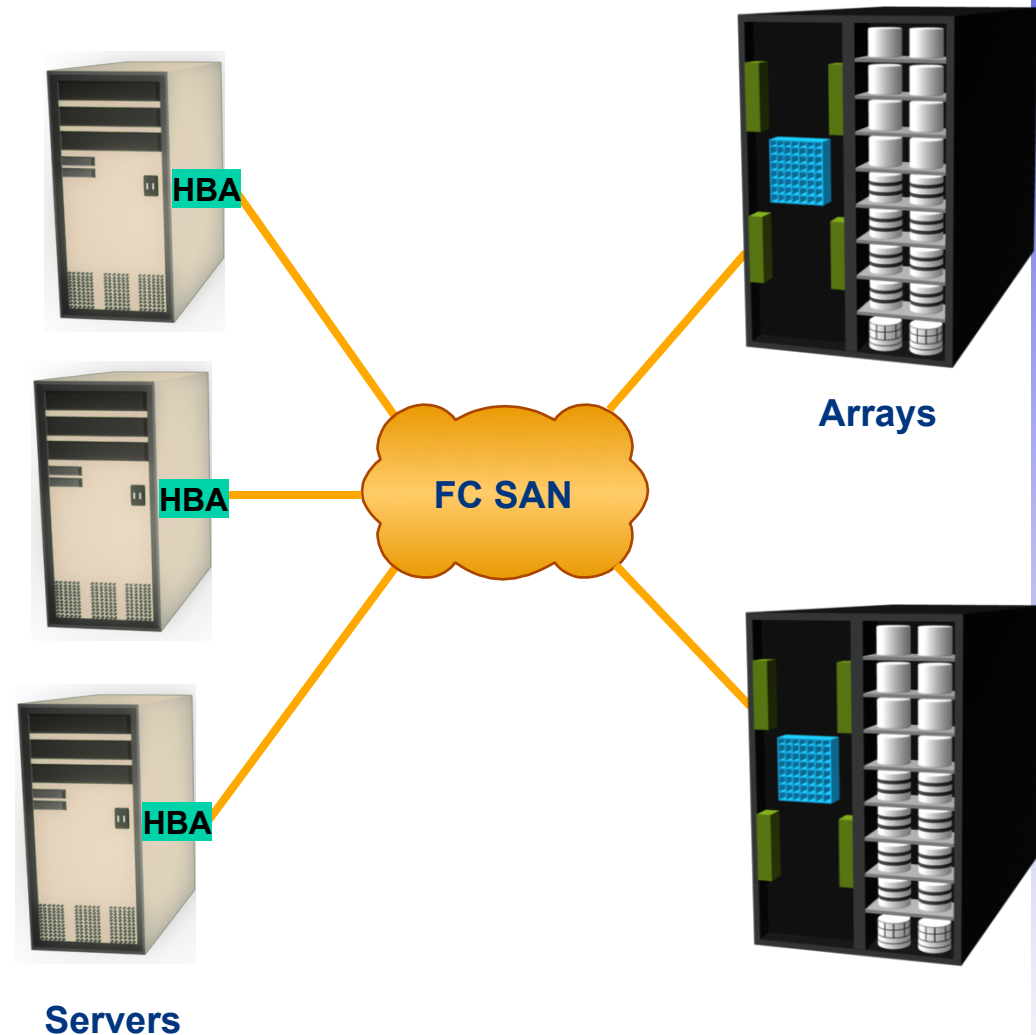
FC Switch



Director

Components of SAN: Storage array

- Provides storage consolidation and centralization
- Features of an array
 - High Availability/Redundancy
 - Performance
 - Business Continuity
 - Multiple host connect

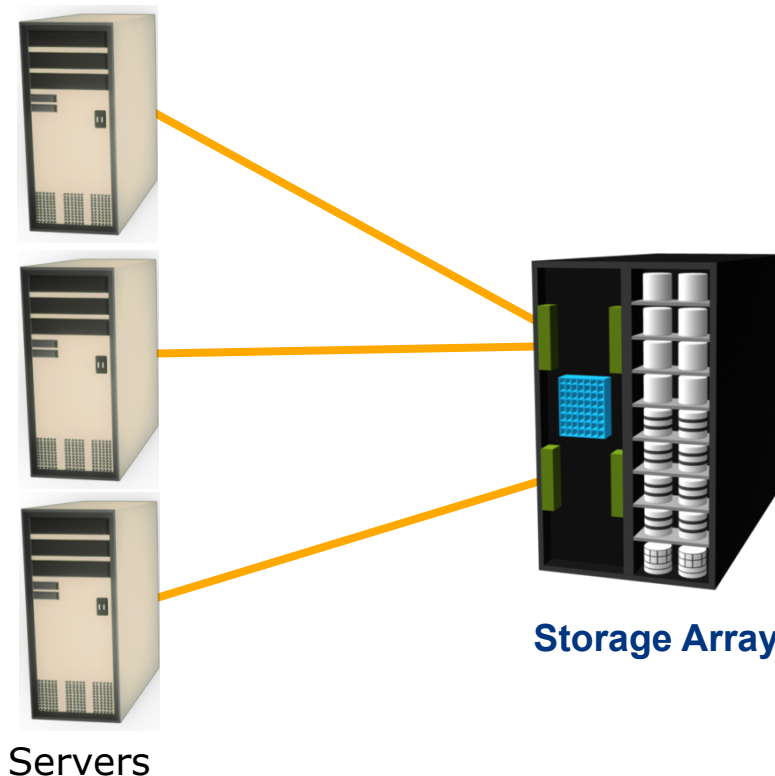


SAN Interconnectivity Options:

Point to Point

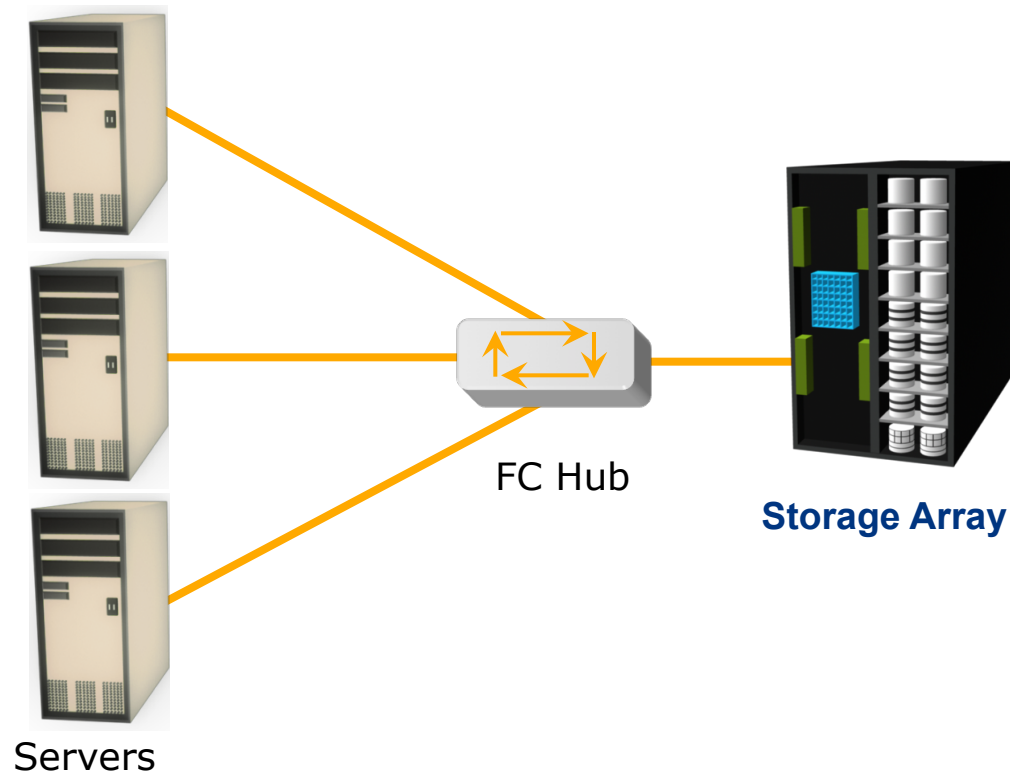
Point to point (Pt-to-Pt)

- Direct connection between devices
- Limited connectivity

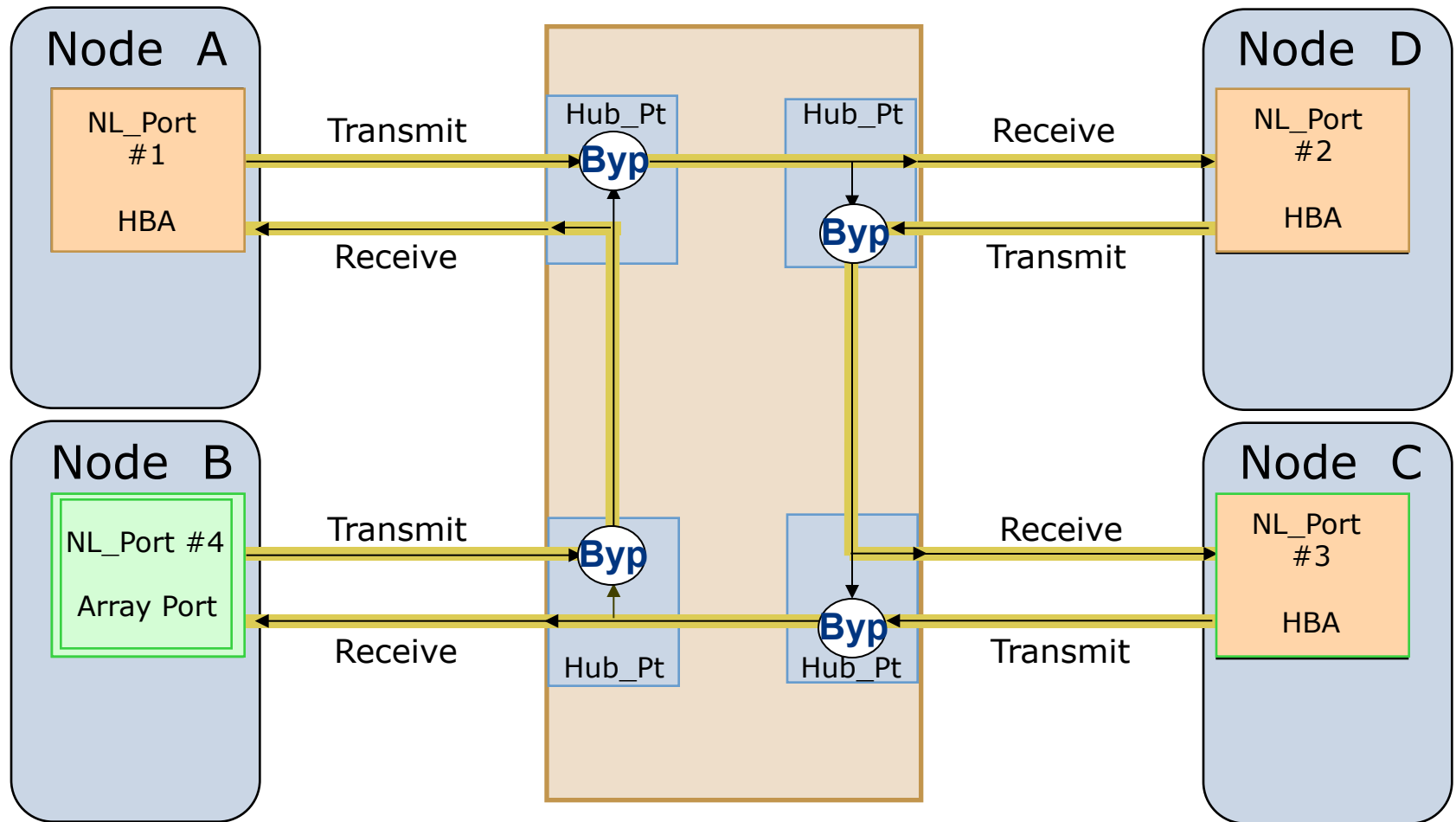


SAN Interconnectivity Options: FC-AL

- Fibre Channel Arbitrated Loop (FC-AL)
 - Devices must arbitrate to gain control
 - Devices are connected via hubs
 - Supports up to 127 devices

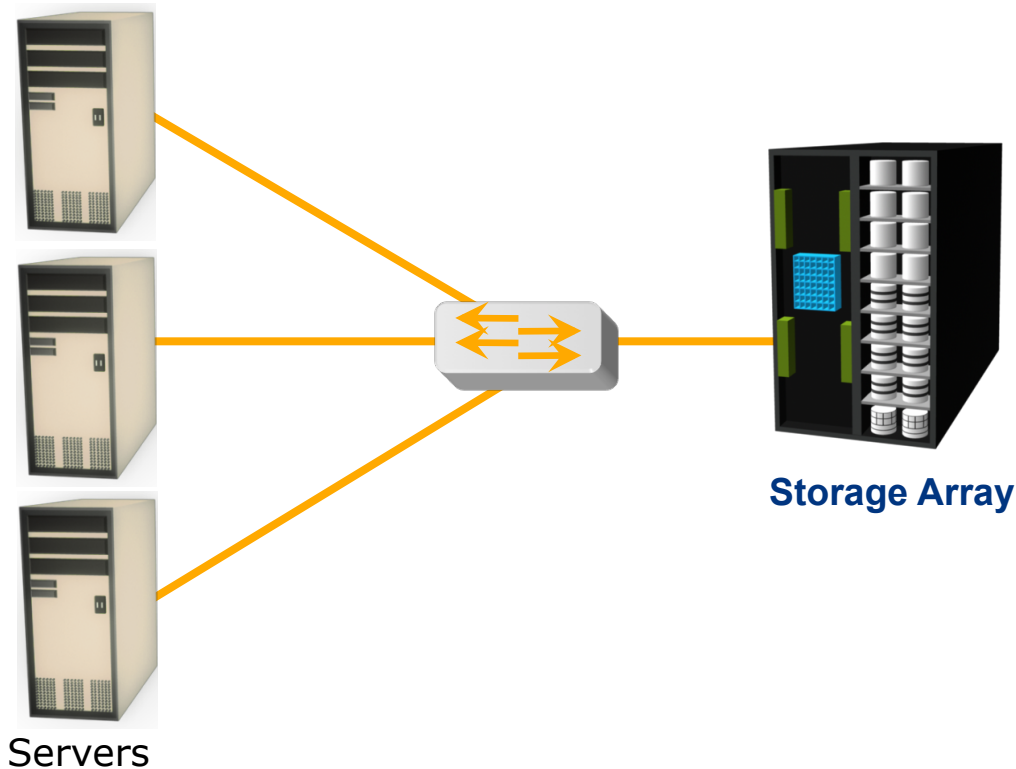


FC-AL Transmission

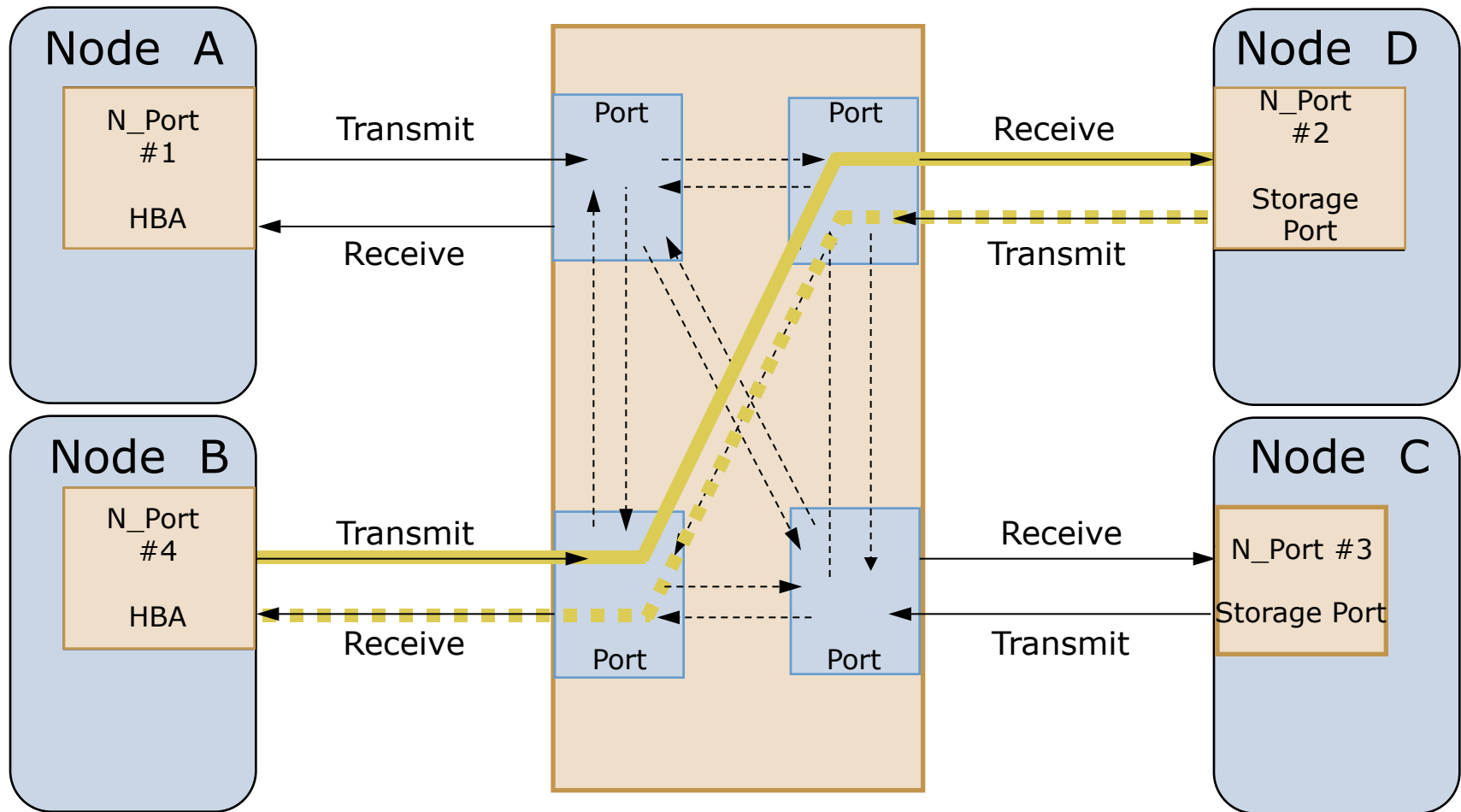


SAN Interconnectivity Options: FC-SW

- Fabric connect (FC-SW)
 - Dedicated bandwidth between devices
 - Support up to 15 million devices
 - Higher availability than hubs



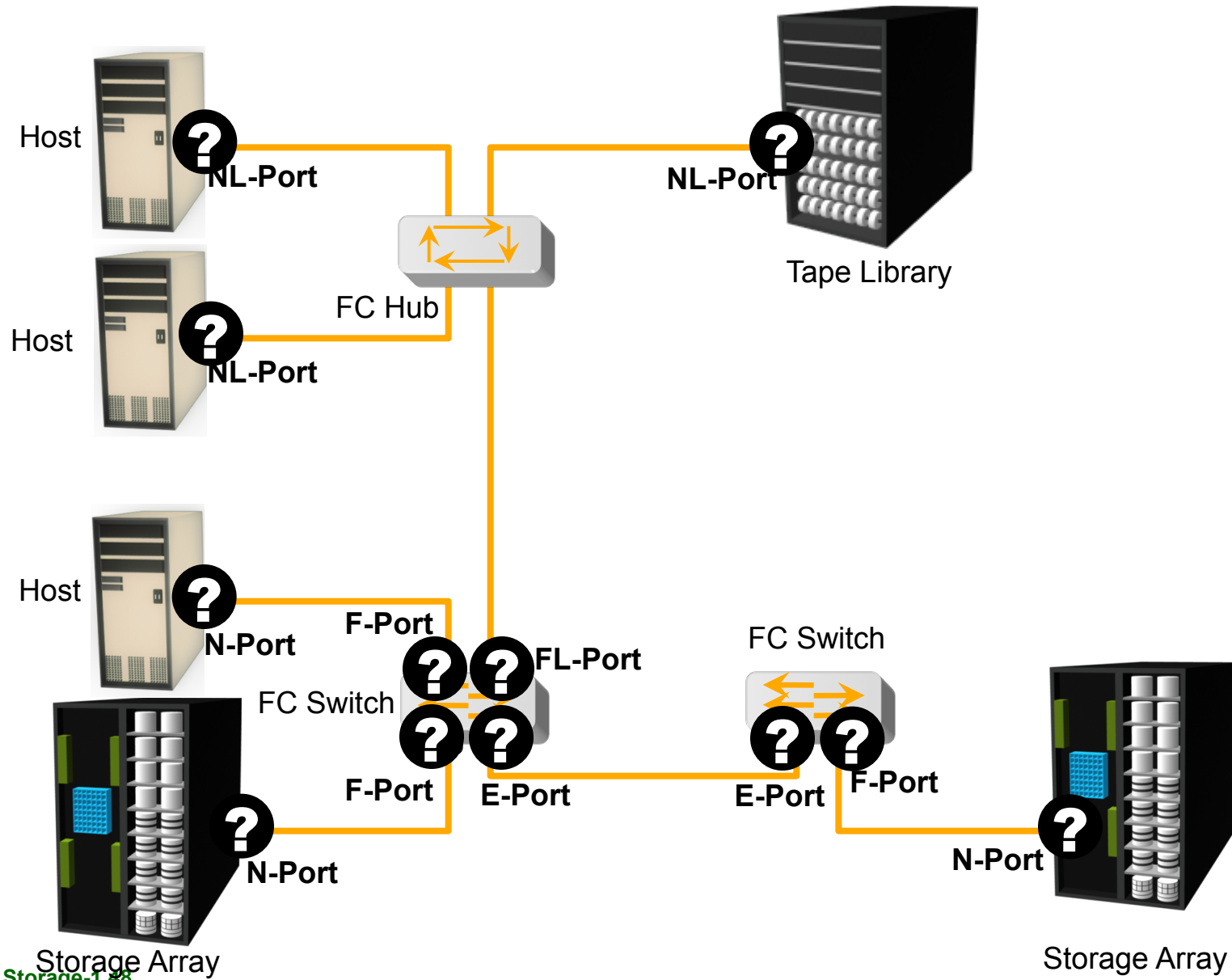
FC-SW Transmission



Port Types

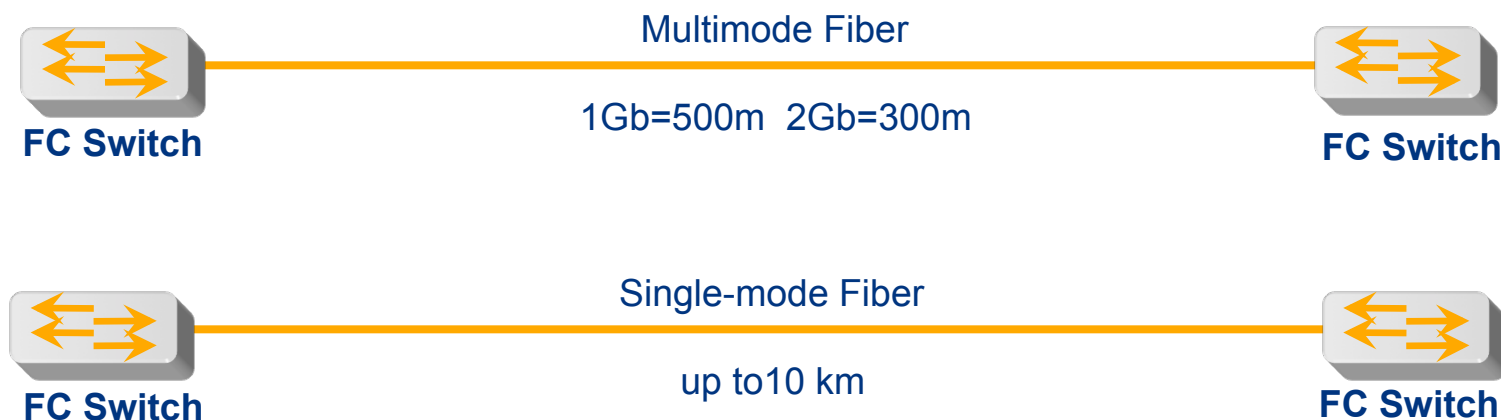
- **N_port** is a port on the node (e.g. host or storage device) used with both FC-P2P or FC-SW topologies. Also known as **node port**.
- **NL_port** is a port on the node used with an FC-AL topology. Also known as **Node Loop port**.
- **F_port** is a port on the switch that connects to a node point-to-point (i.e. connects to an N_port). Also known as **fabric port**. An F_port is not loop capable.
- **FL_port** is a port on the switch that connects to a FC-AL loop (i.e. to NL_ports). Also known as **fabric loop port**.
- **E_port** is the connection between two fibre channel switches. Also known as an **Expansion port**. When E_ports between two switches form a link, that link is referred to as an inter-switch link (**ISL**).
- **B_port** A Bridge Port is a Fabric inter-element port used to connect Bridge devices with E_Ports on a Switch. The B_Port provides a subset of the E_port functionality
- **D_port** is a diagnostic port, used solely for the purpose of running link-level diagnostics between two switches and to isolate link level fault on the port, in the SFP, or in the cable.
- **EX_port** is the connection between a fibre channel router and a fibre channel switch. On the side of the switch it looks like a normal E_port, but on the side of the router it is an EX_port.
- **TE_port** * Is an extended ISL or **EISL**. The TE_port provides not only standard E_port functions but allows for routing of multiple **VSANs** (Virtual SANs). This is accomplished by modifying the standard Fibre Channel frame (vsan tagging) upon ingress/egress of the VSAN environment. Also known as **Trunking E_port**.
- **VE_Port** an **INCITS** T11 addition, FCIP interconnected E-Port/ISL, i.e. fabrics will merge.
- **VEX_Port** an **INCITS** T11 addition, is a FCIP interconnected EX-Port, routing needed via Isan zoning to connect initiator to a target.

Port Types



Inter Switch Links (ISL)

- ISL connects two or more FC switches to each other using E-Ports
- ISLs are used to transfer host-to-storage data as well as the fabric management traffic from one switch to another
- ISL is also one of the scaling mechanisms in SAN connectivity



Login Types in a Switched Network

Extended Link Services that are defined in the standards:

- FLOGI - Fabric login
 - Between N_Port to F_Port
- PLOGI - Port login
 - Between N_Port to N_Port
 - N_Port establishes a session with another N_Port
- PRLI - Process login
 - Between N_Port to N_Port
 - To share information about the upper layer protocol type in use
 - And recognizing device as the SCSI initiator, or target

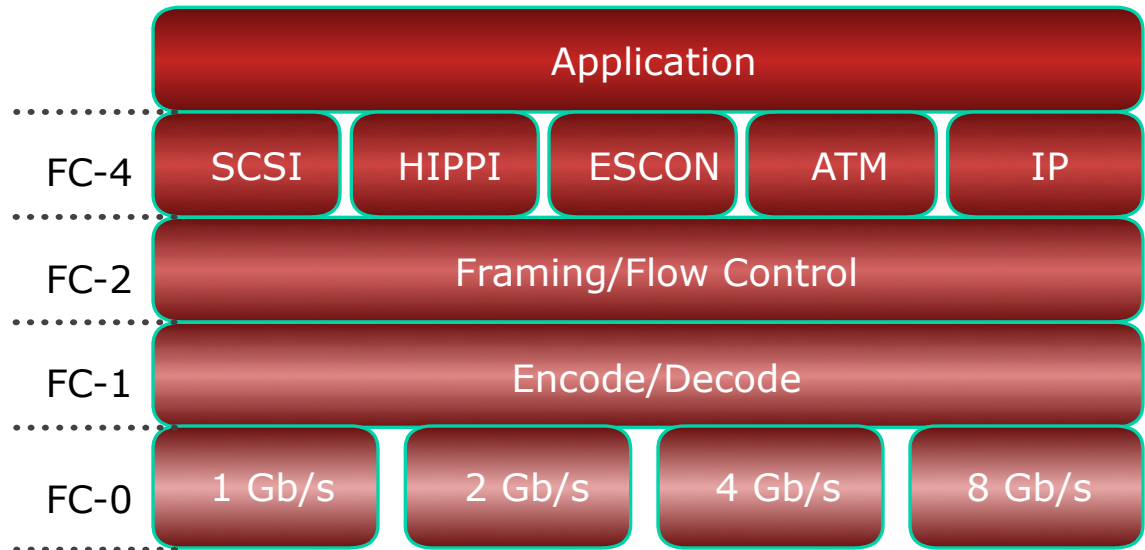
FC Architecture Overview

- FC uses channel technology
- Provide high performance with low protocol overheads
- FCP is SCSI-3 over FC network
 - Sustained transmission bandwidth over long distances
 - Provides speeds up to 8 Gb/s (8 GFC)

- FCP has five layers:

- FC-4
- FC-2
- FC-1
- FC-0

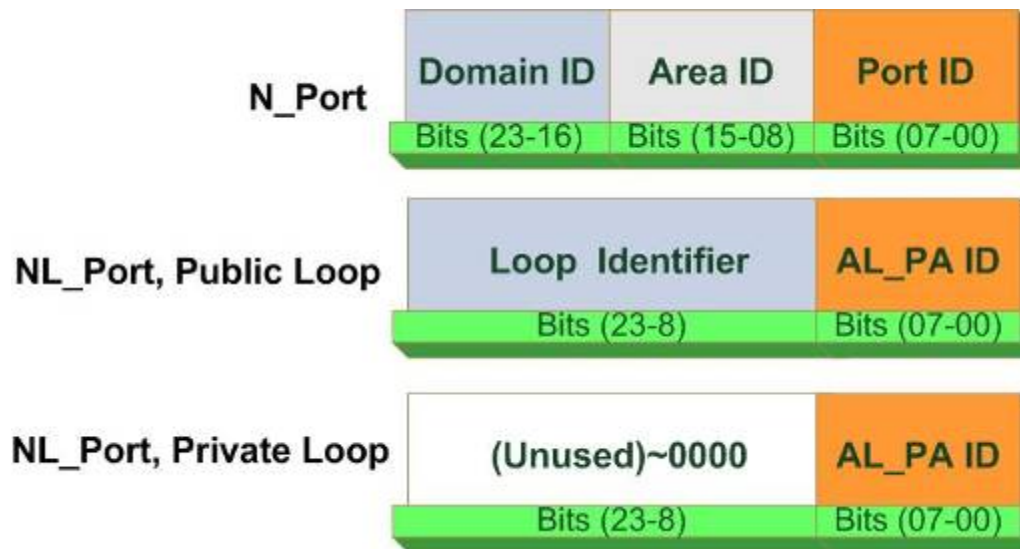
*FC-3 is not yet implemented



Fibre Channel Protocol Stack

FC layer	Function	SAN relevant features specified by FC layer
FC-4	Mapping interface	Mapping upper layer protocol (e.g. SCSI-3 to FC transport)
FC-3	Common services	Not implemented
FC-2	Routing, flow control	Frame structure, ports, FC addressing, buffer credits
FC-1	Encode/decode	8b/10b encoding, bit and frame synchronization
FC-0	Physical layer	Media, cables, connector

Fibre Channel Addressing



- FC Address is assigned during Fabric Login
 - Used to communicate between nodes within SAN
 - Similar in functionality to an IP address on NICs
- Address Format:
 - 24 bit address, dynamically assigned
 - Contents of the three bytes depend on the type of N-Port
 - For an N_Port or a public NL_Port:
 - switch maintains mapping of WWN to FC-Address via the Name Server

World Wide Names

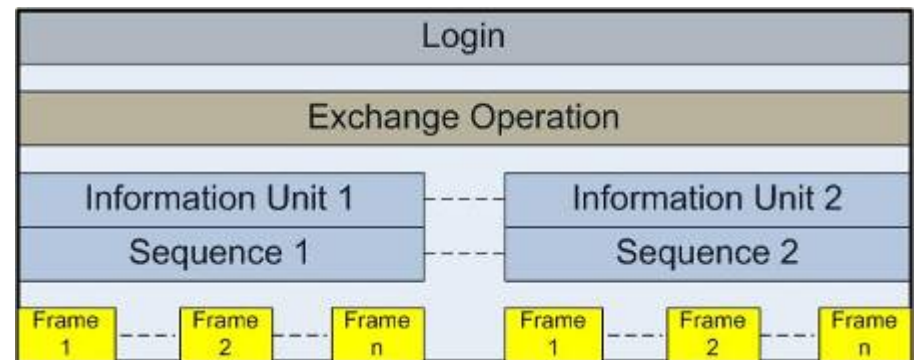
- Unique 64 bit identifier
- Static to the port
 - Used to physically identify ports or nodes within SAN
 - Similar to NIC's MAC address

World Wide Name - Array															
5	0	0	6	0	1	6	0	0	0	6	0	0	1	B	2
0101	0000	0000	0110	0000	0001	0110	0000	0000	0000	0110	0000	0000	0001	1011	0010
Company ID 24 bits							Port	Model Seed 32 bits							

World Wide Name - HBA															
1	0	0	0	0	0	0	0	c	9	2	0	d	c	4	0
Reserved 12 bits				Company ID 24 bits						Company Specific 24 bits					

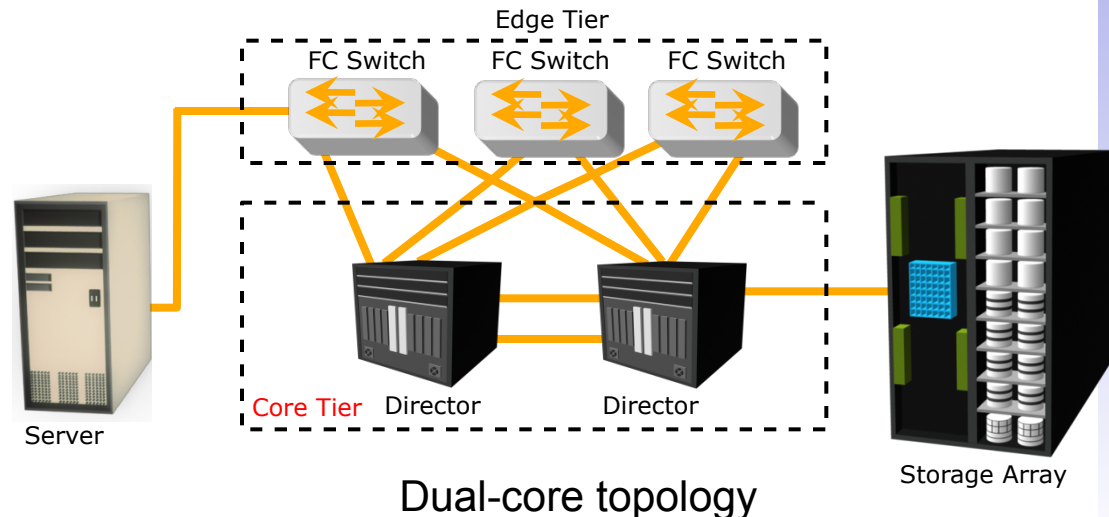
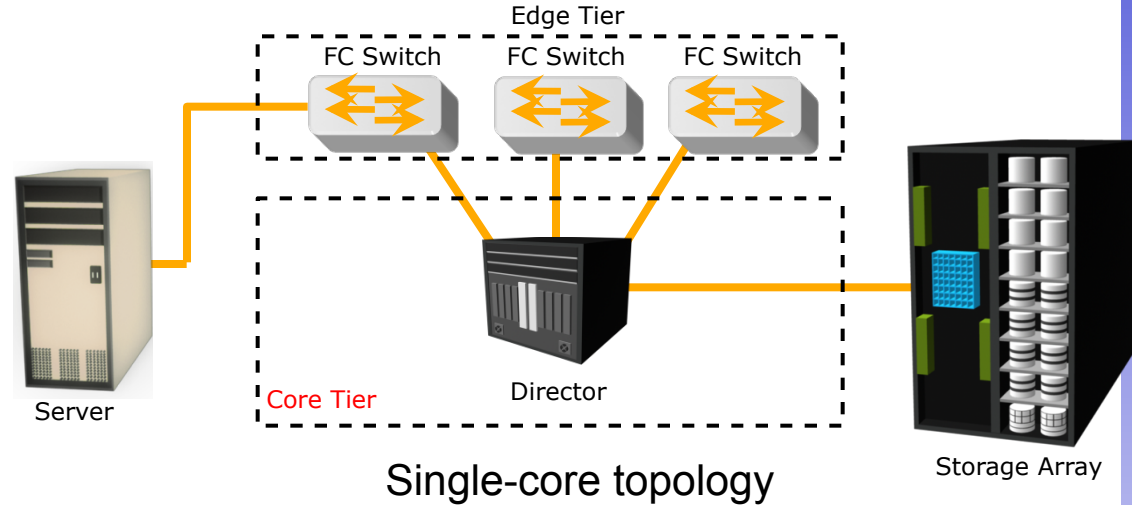
Structure and Organization of FC Data

- FC data is organized as:
 - Exchange operations
 - Enables two N_ports to identify and manage a set of information units
 - Maps to sequence
 - Sequence
 - Contiguous set of frames sent from one port to another
 - Frames
 - Fundamental unit of data transfer
 - Each frame can contain up to 2112 bytes of payload



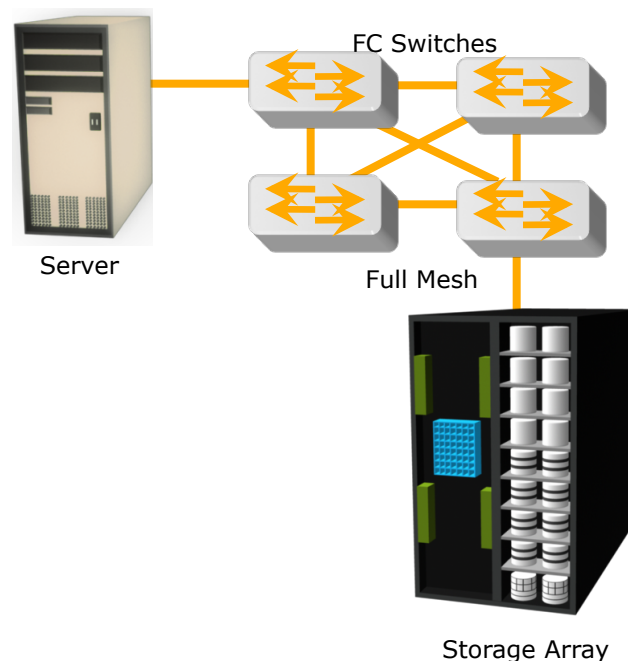
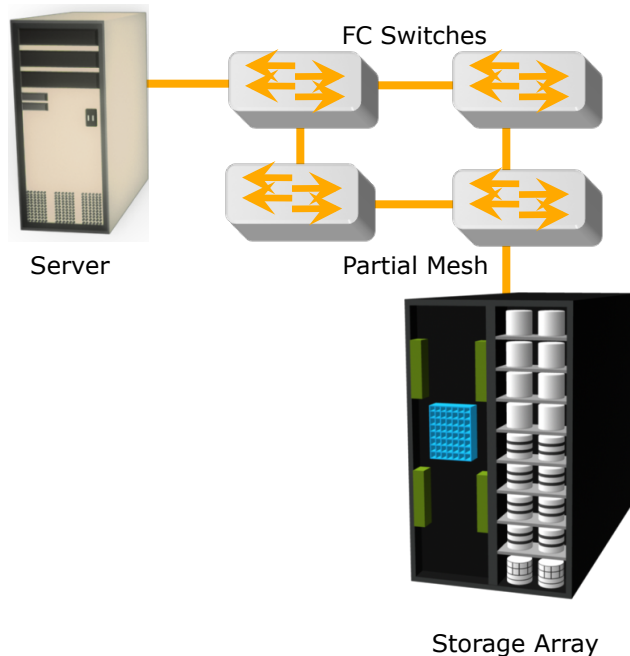
Fabric Topology: Core-Edge Fabric

- Can be two or three tiers
 - Single Core Tier
 - One or two Edge Tiers
- In a two tier topology, storage is usually connected to the Core
- Benefits
 - High Availability
 - Medium Scalability
 - Medium to maximum Connectivity

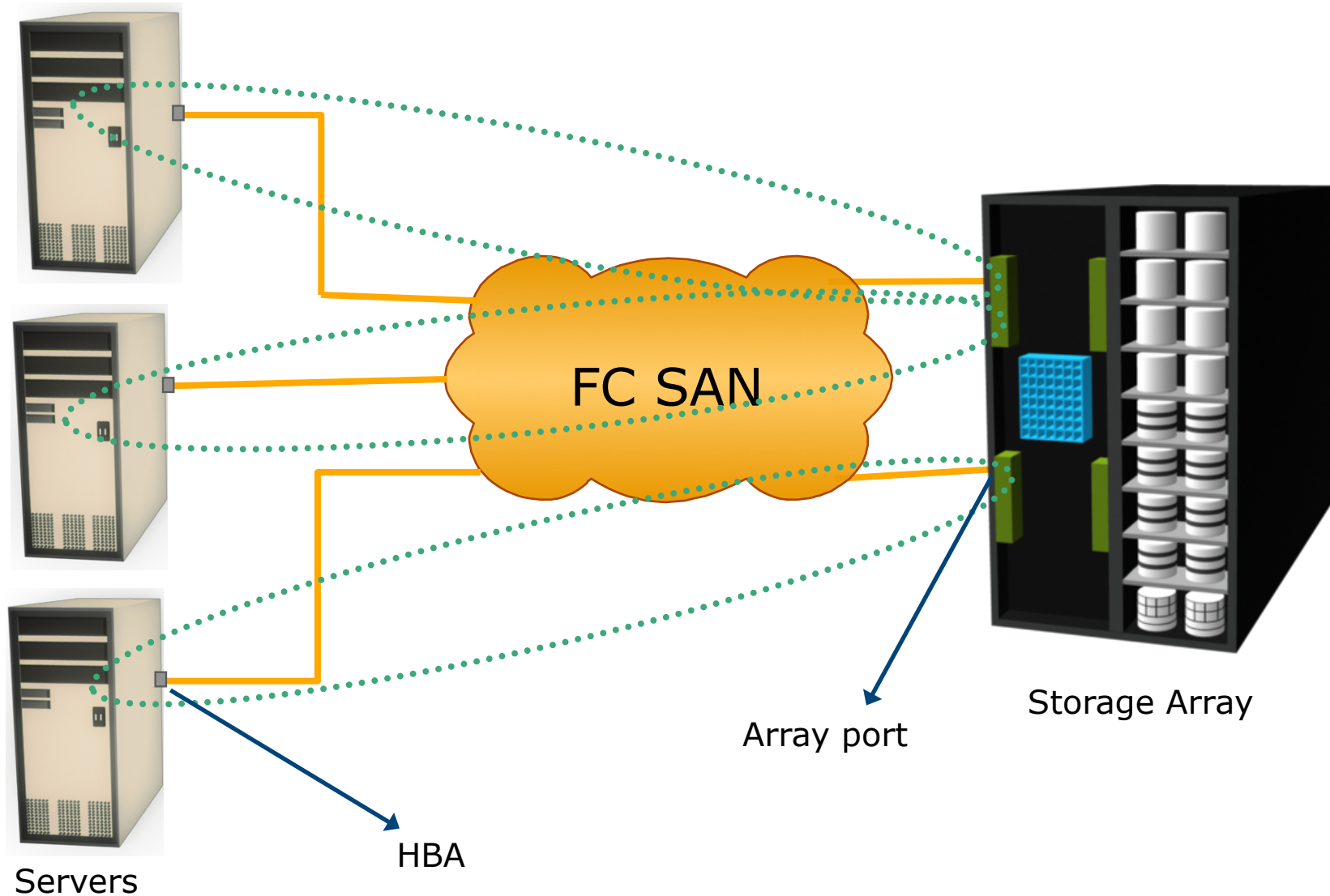


Fabric Topology: Mesh

- Can be either partial or full mesh
- All switches are connected to each other
- Host and Storage can be located anywhere in the fabric
- Host and Storage can be localized to a single switch



Fabric Management: Zoning

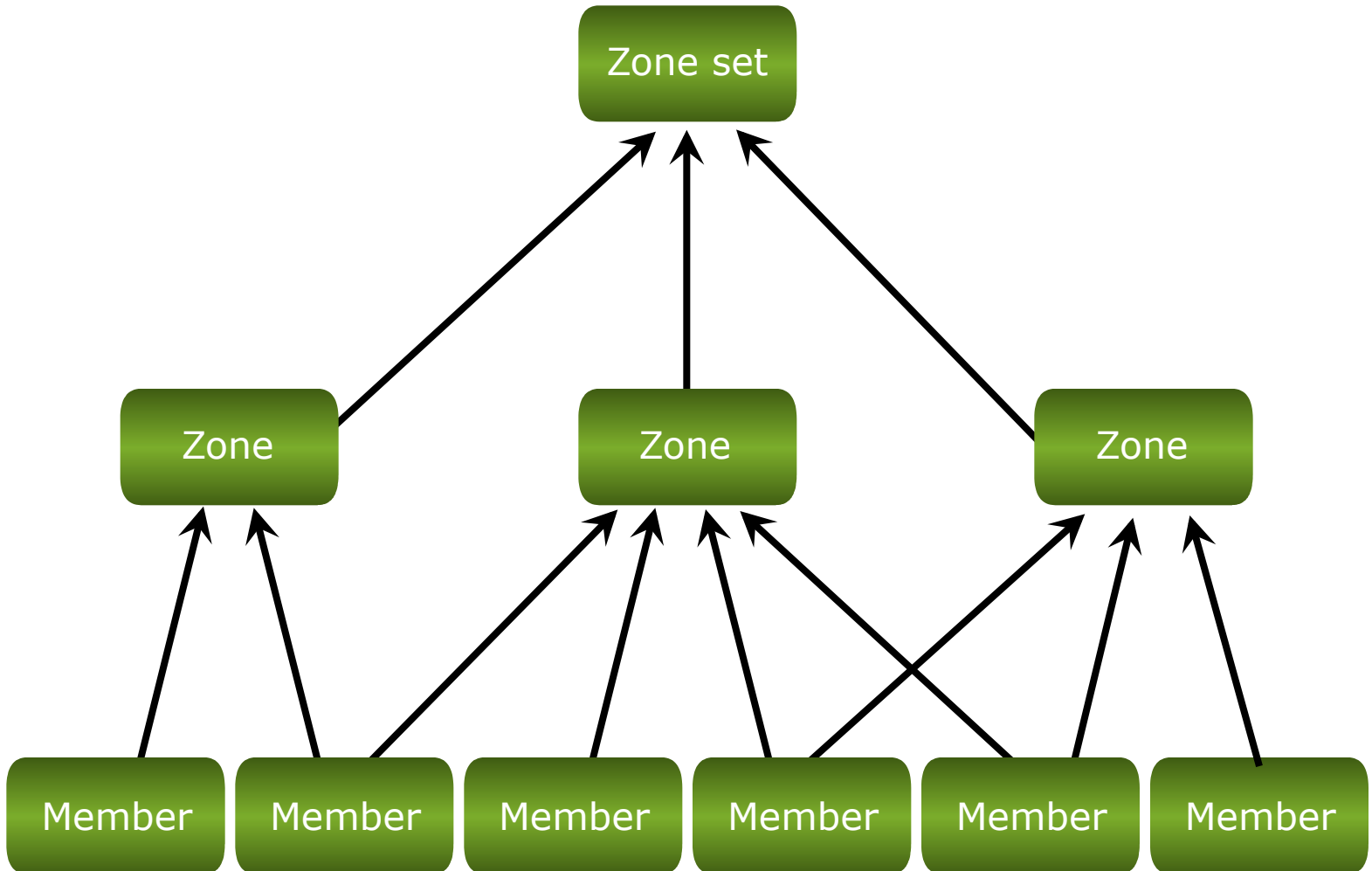


Zoning Components

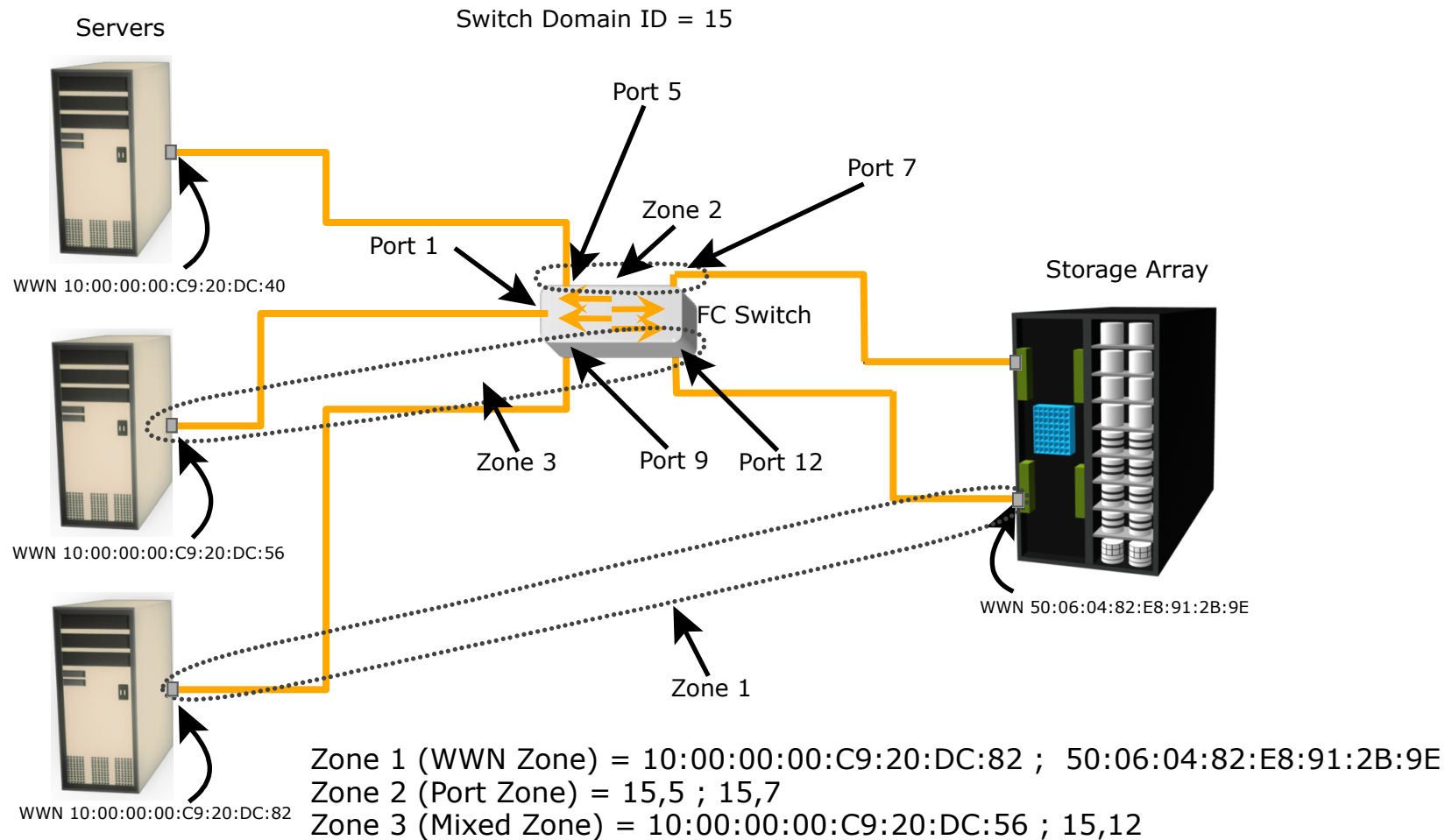
Zone sets
(Library)

Zone
(Library)

Member
WWN's



Types of Zoning

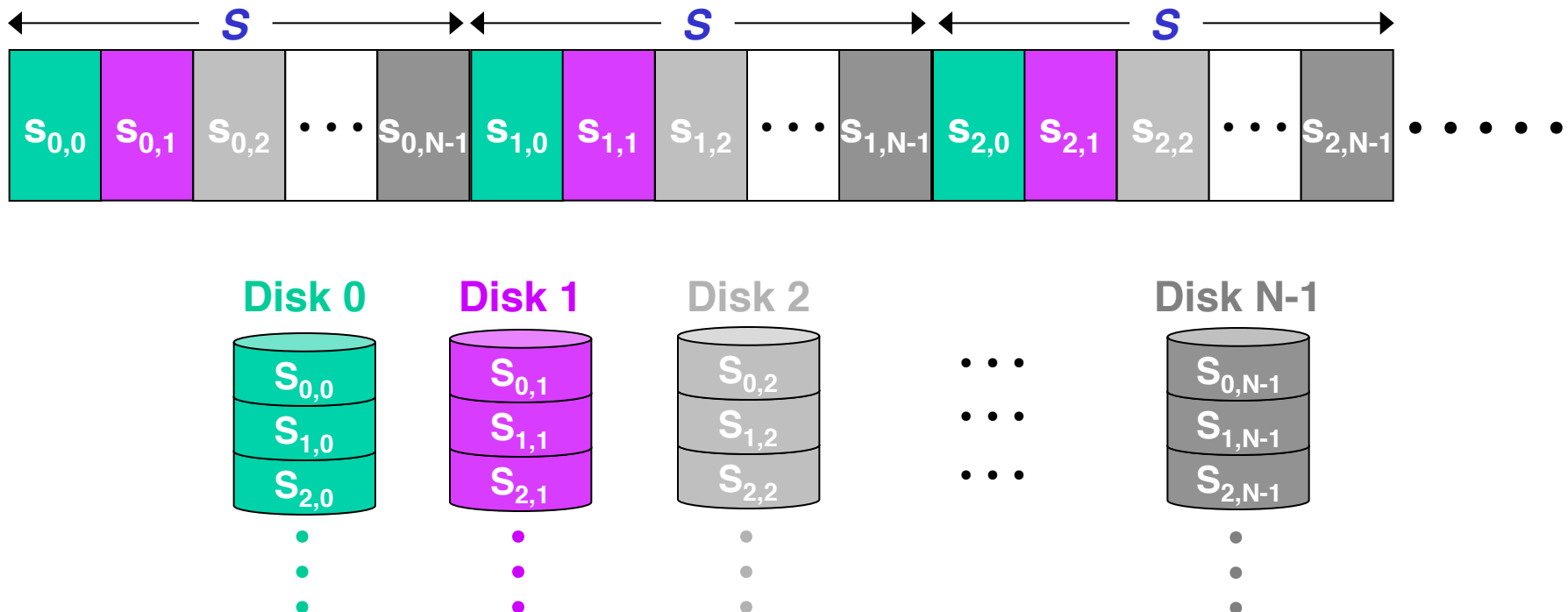


Motivation for Disk Arrays

- Typical memory bandwidths \approx 150 MB/sec
- Typical disk bandwidths \approx 10 MB/sec
- Result: *I/O-bound applications limited by disk bandwidth*
 - (not just by disk latency!)
- Two common disk bandwidth-limited scenarios
 - Scientific applications: one huge I/O-hungry app
 - Servers: many concurrent requests with modest I/O demands

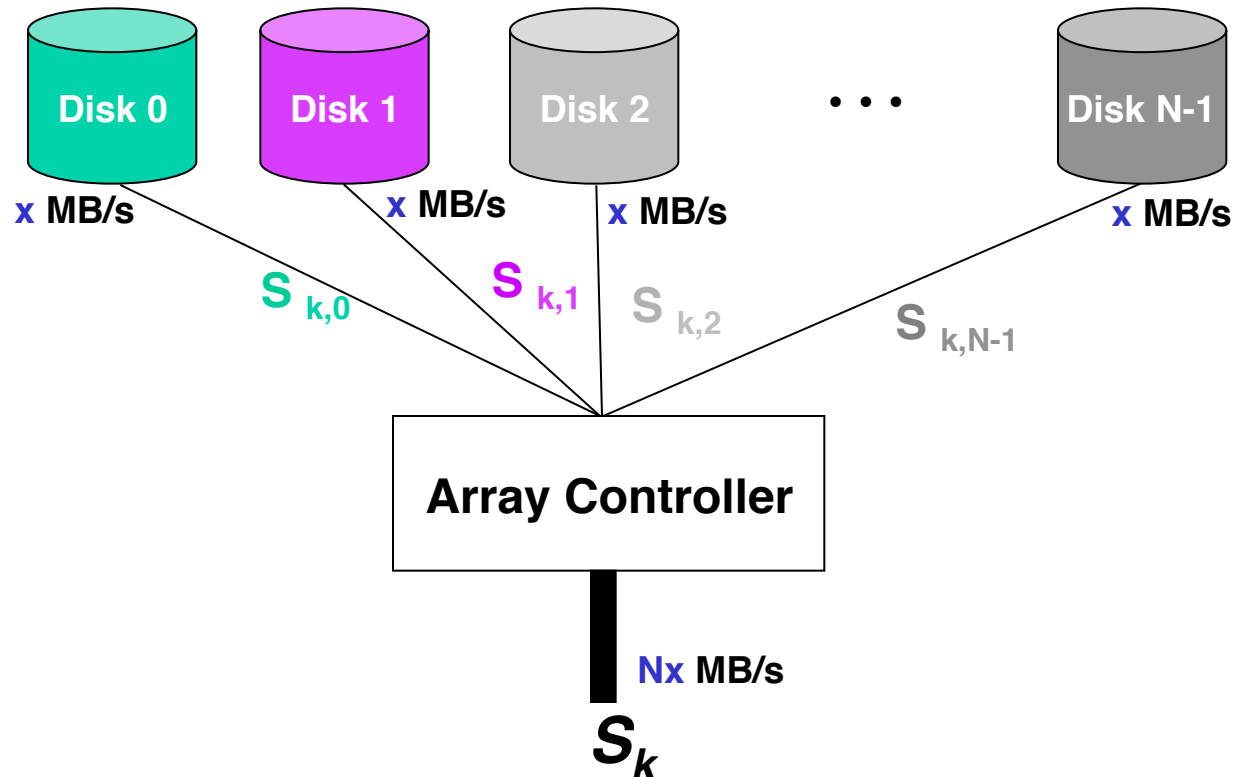
Solution: Exploit Parallelism

- *Stripe* the data across an array of disks
 - many alternative striping strategies possible
- Example: consider a big file striped across N disks
 - *stripe width* is S bytes
 - hence each *stripe unit* is S/N bytes
 - sequential read of S bytes at a time

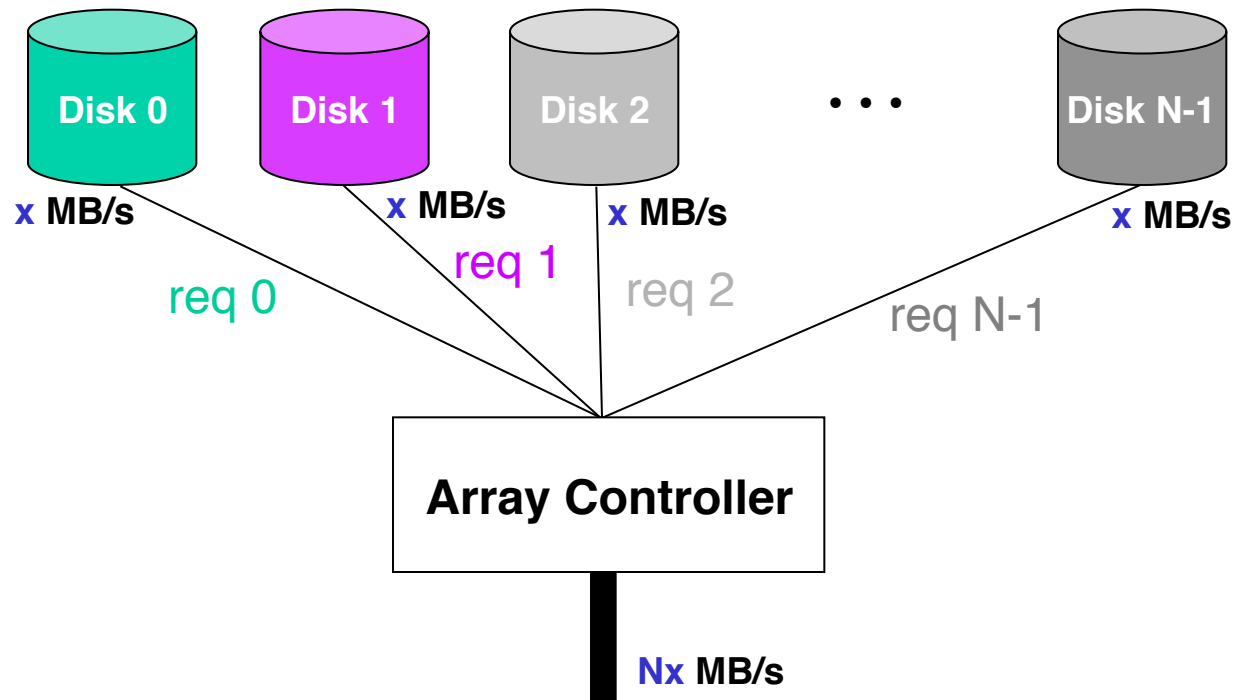


Performance Benefit

- *Sequential read or write of large file*
 - application (or I/O buffer cache) reads in multiples of S bytes
 - controller performs parallel access of N disks
 - aggregate bandwidth is N times individual disk bandwidth
 - (assumes that disk is the bottleneck)



- *N concurrent small read or write requests*
 - randomly distributed across N drives (we hope!)
 - common in database and Web server environments



N independent requests

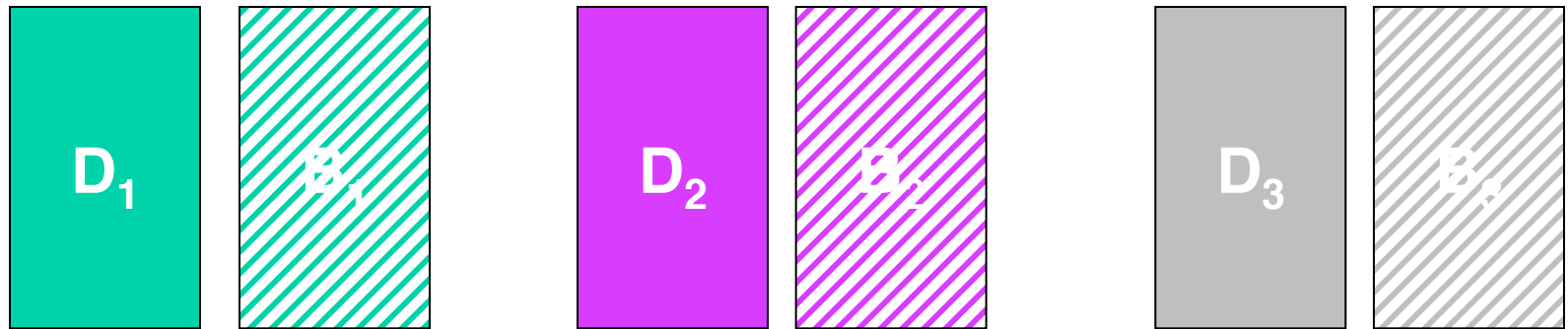
Reliability of Disk Arrays

- As number of disks grows, chances of at least one failing increases
- Reliability of N disks = (reliability of 1 disk) / N
 - suppose each disk has MTTF of 50,000 hours
 - (roughly 6 years before any given disk fails)
 - then some disk in a 70-disk array will fail in (50,000 / 70) hours
 - (roughly once a month!)
- Large arrays without redundancy too unreliable to be useful
 - *“Redundant Arrays of Independent Disks” (RAID)*

RAID Approaches

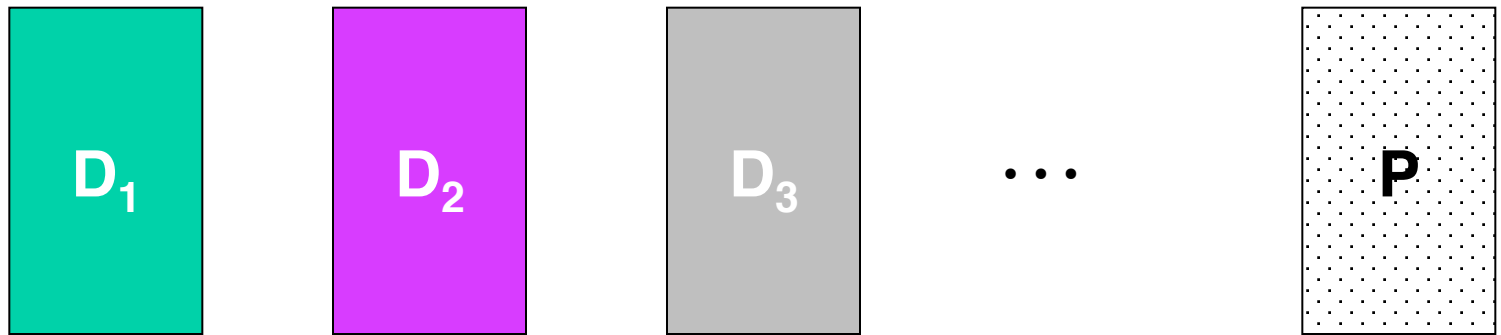
- Many alternative approaches to achieving this redundancy
 - *RAID levels 1 through 5*
 - *hot sparing* allows reconstruction concurrently with accesses
- Key metrics to evaluate alternatives
 - wasted space due to redundancy
 - likelihood of “hot spots” during heavy loads
 - degradation of performance during repair

RAID Level 1



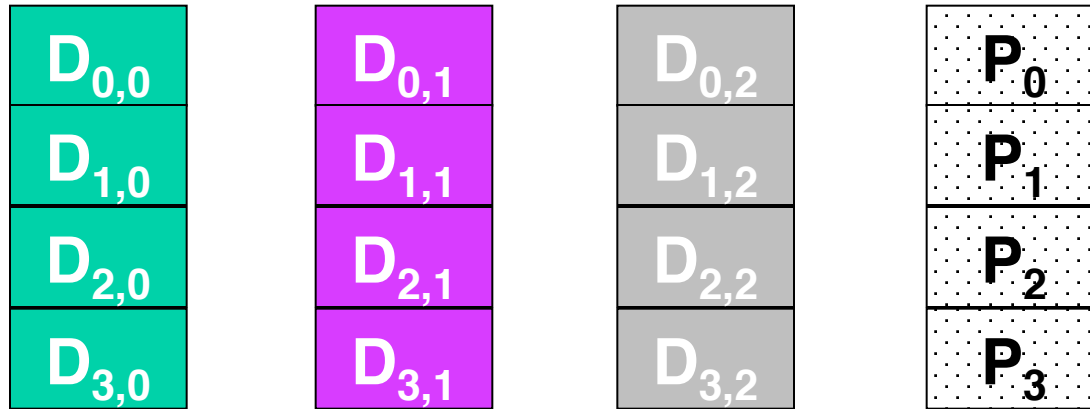
- Also known as “*mirroring*”
- To read a block:
 - read from either data disk or backup
- To write a block:
 - write both data and backup disks
 - failure model determines whether writes can occur in parallel
- Backups can be located far way: safeguard against site failure

RAID Levels 2 & 3



- These are *bit-interleaved* schemes
- In Raid Level 2, P contains memory-style ECC
- In Rail Level 3, P contains simple parity
- Rarely used today

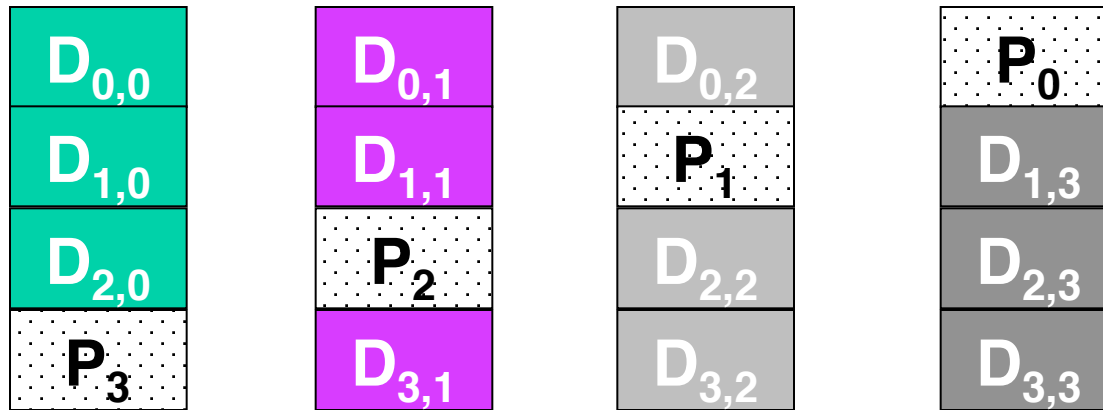
RAID Level 4



$$D_{0,0} \oplus D_{0,1} \oplus D_{0,2} = P_0$$

- *Block-interleaved parity*
- Wasted storage is small: one parity block for N data blocks
- Key problem:
 - parity disk becomes a hot spot
 - write access to parity disk on every write to any block

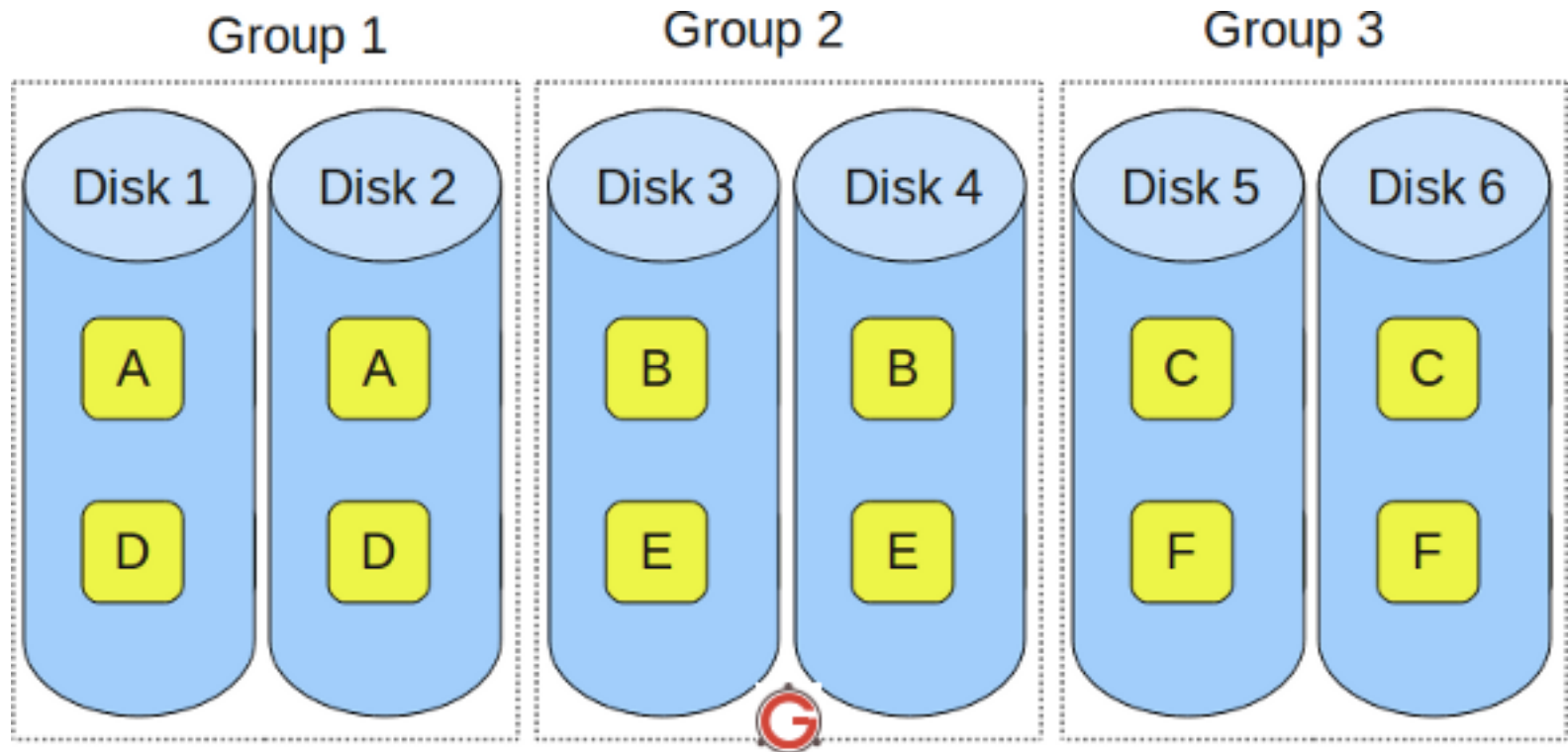
RAID Level 5



$$D_{0,0} \oplus D_{0,1} \oplus D_{0,2} = P_0$$

- *Rotated parity*
- Wastage is small: same as in Raid 4
- Parity update traffic is distributed across disks

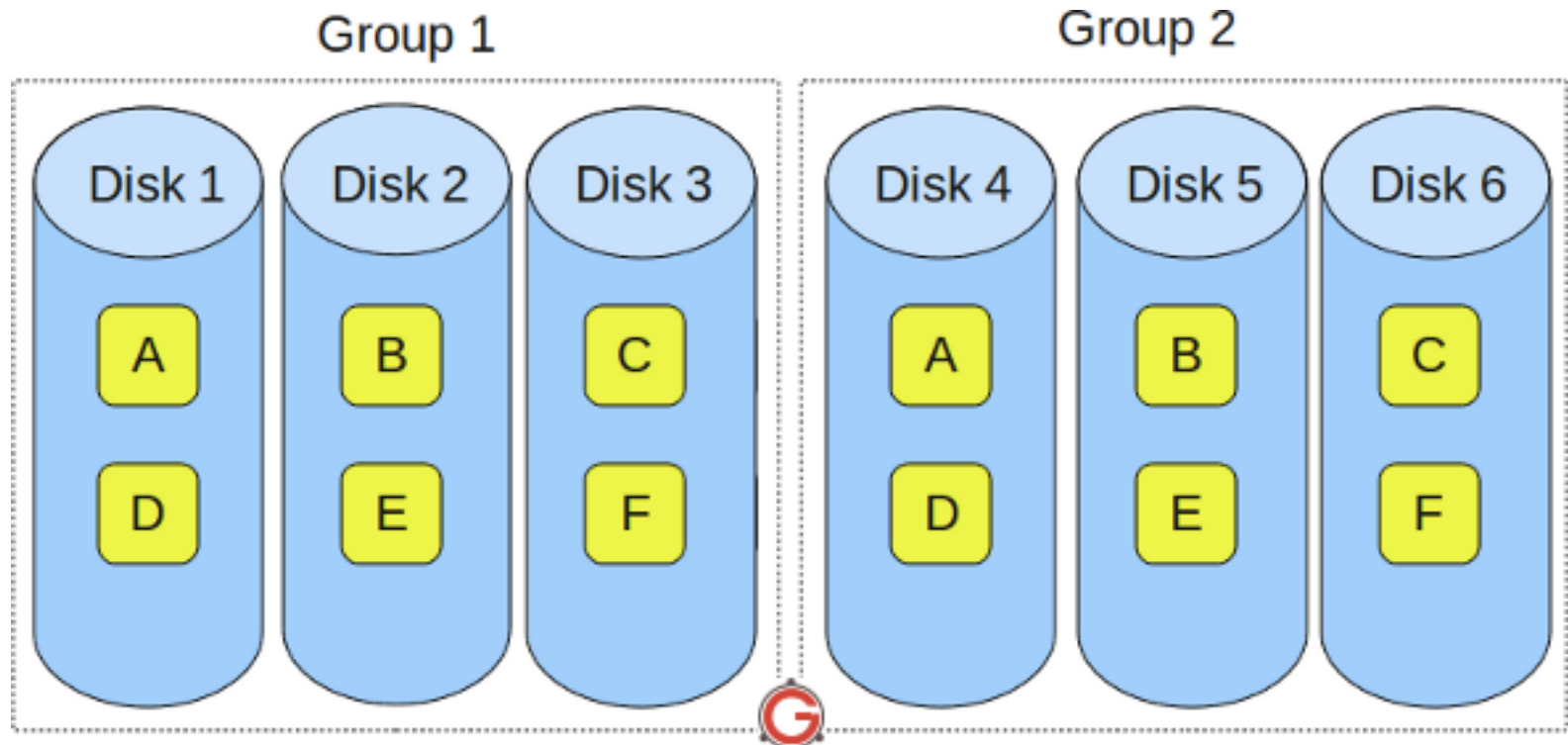
RAID Level 10



RAID 10 – Blocks Mirrored. (and Blocks Striped)

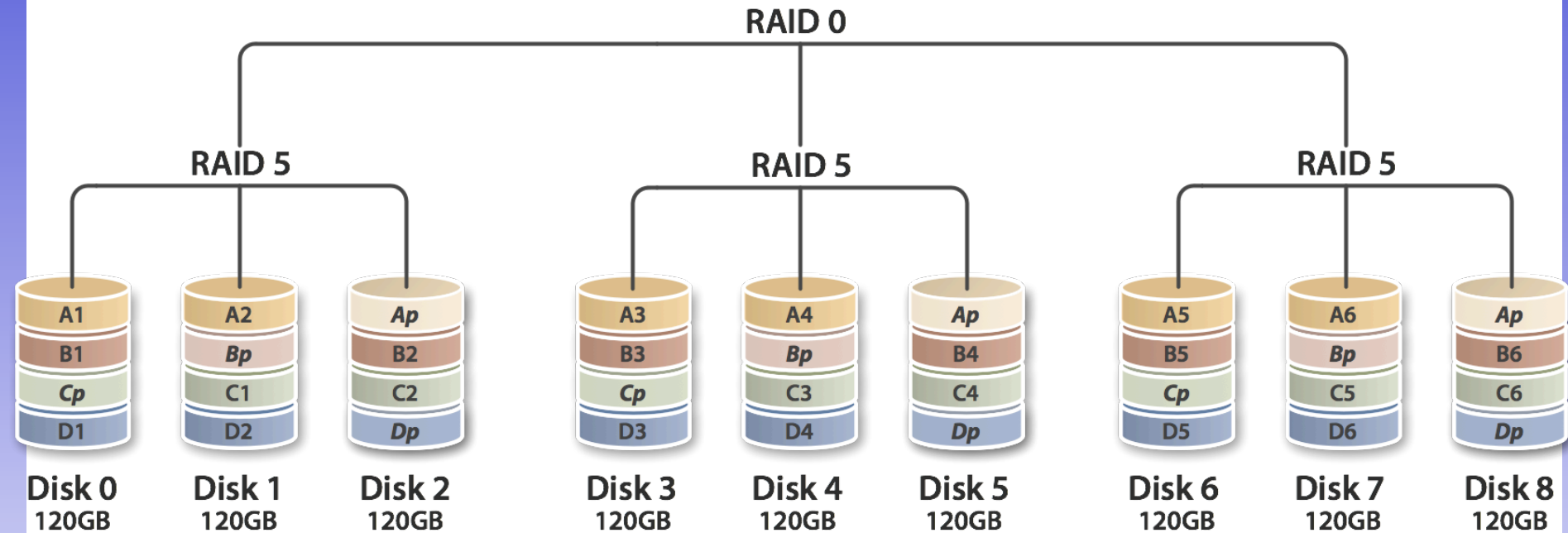
- *Difference between RAID 0+1 and RAID 1+0*

RAID Level 0+1



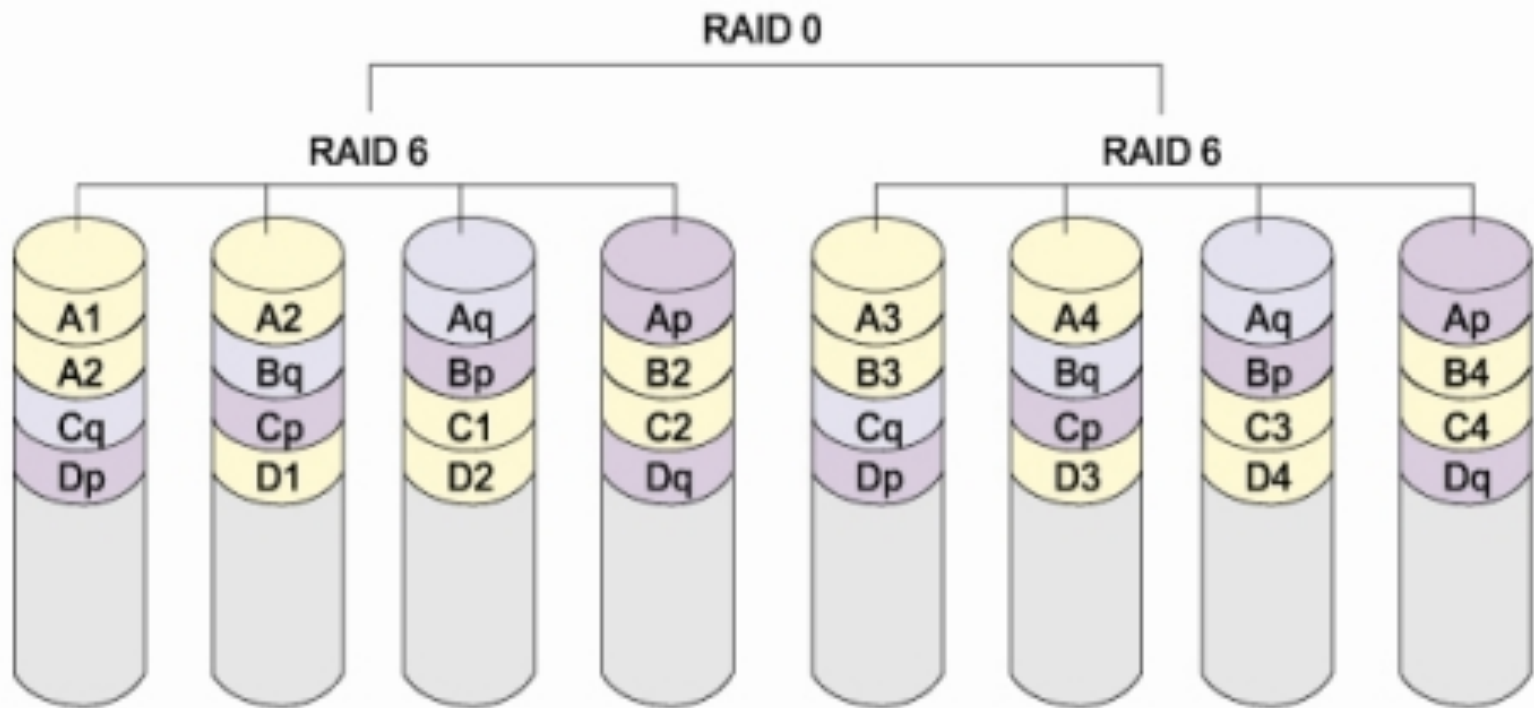
RAID 01 – Blocks Striped. (and Blocks Mirrored)

RAID Level 50



RAID Level 60

RAID 60



RAID Level 100

RAID-0

RAID-0

RAID-0

RAID-0

RAID-1

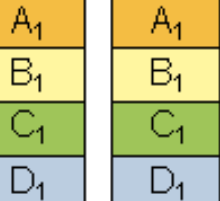
RAID-1

RAID-1

RAID-1

RAID-1

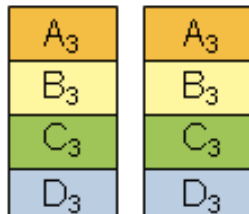
RAID-1



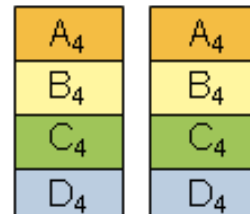
Disk 1 Disk 2



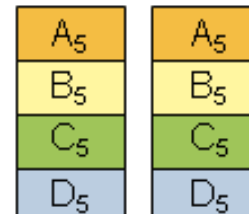
Disk 3 Disk 4



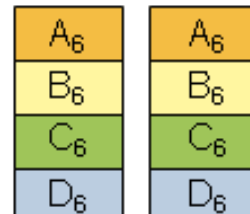
Disk 5 Disk 6



Disk 7 Disk 8



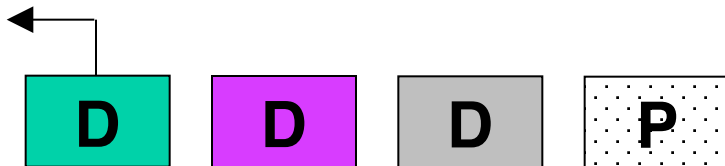
Disk 9 Disk 10



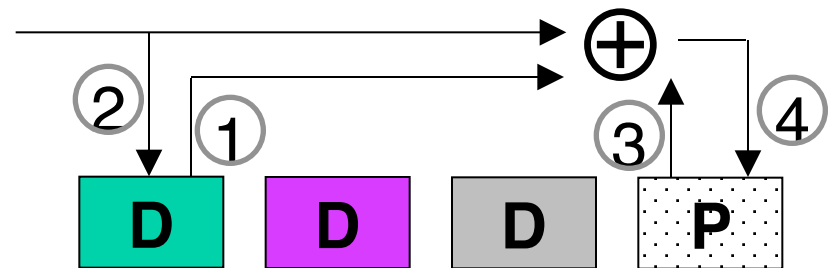
Disk 11 Disk 12

RAID 5 Actions

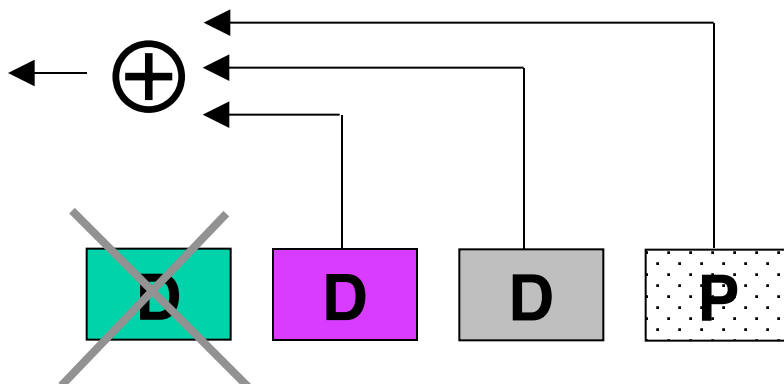
Fault-free Read



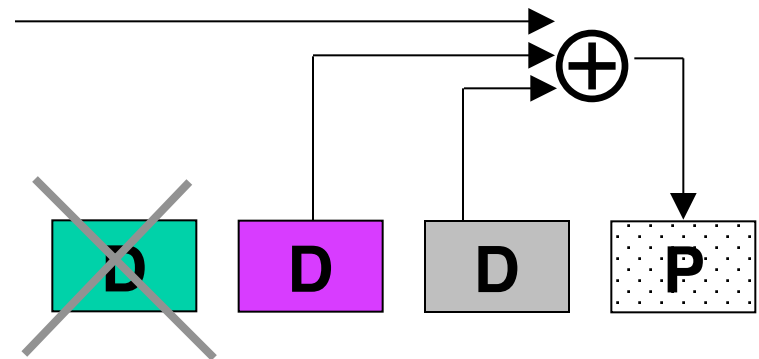
Fault-free Write



Degraded Read



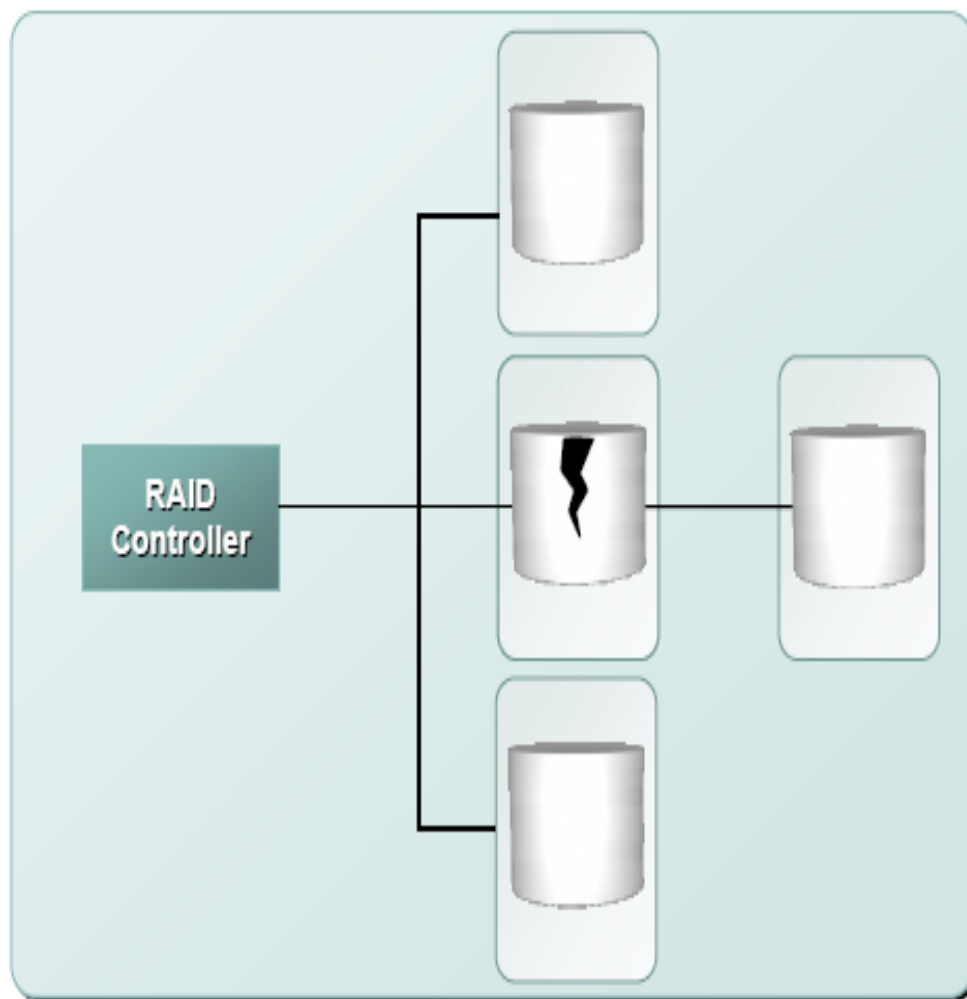
Degraded Write



Hot Spare

- A hot spare refers to a spare HDD in a RAID array that temporarily replaces a failed HDD of a RAID set.

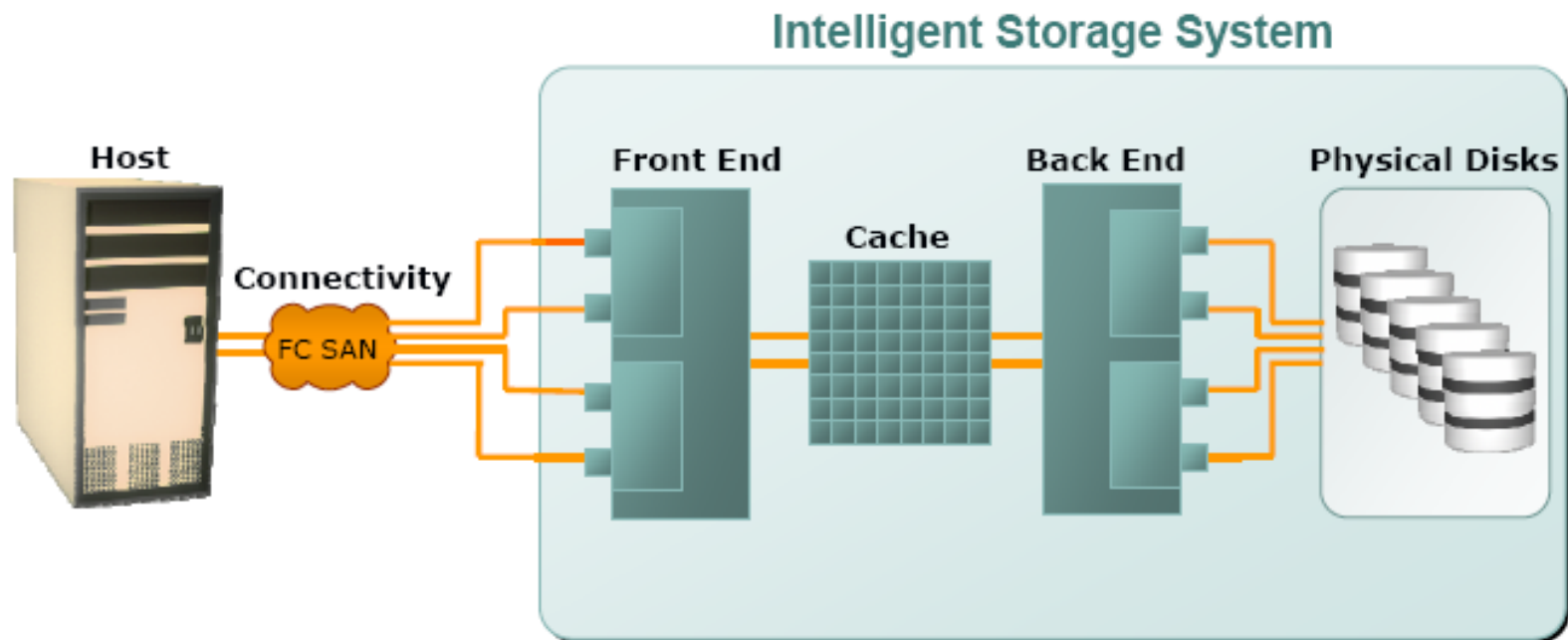
Hot Spares



- When the failed HDD is replaced with a new HDD, The hot spare replaces the new HDD permanently, and a new hot spare must be configured on the array, or data from the hot spare is copied to it, and the hot spare returns to its idle state, ready to replace the next failed drive.
- A hot spare should be large enough to accommodate data from a failed drive.
- Some systems implement multiple hot spares to improve data availability.
- A hot spare can be configured as automatic or user initiated, which specifies how it will be used in the event of disk failure

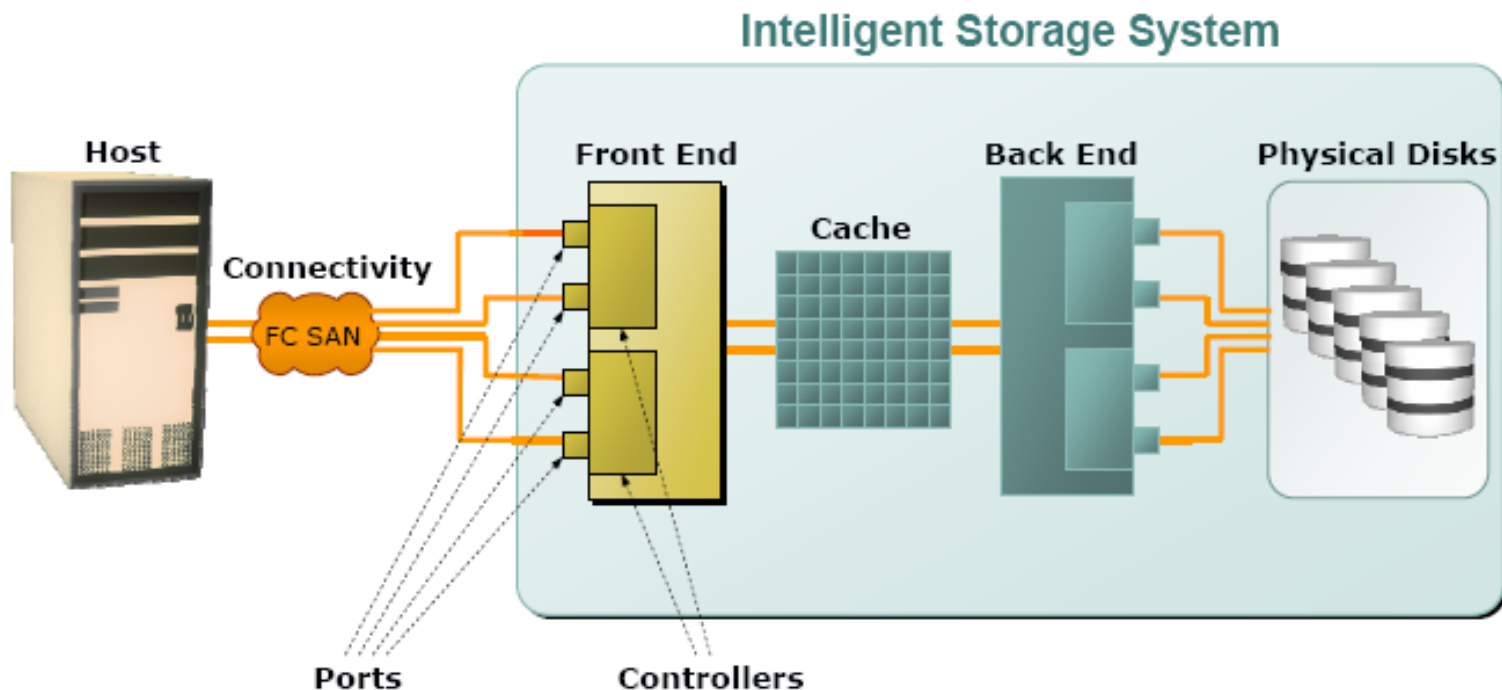
Intelligent Storage System

- An intelligent storage system consists of four key components: *front end*, *cache*, *back end*, and *physical disks*.



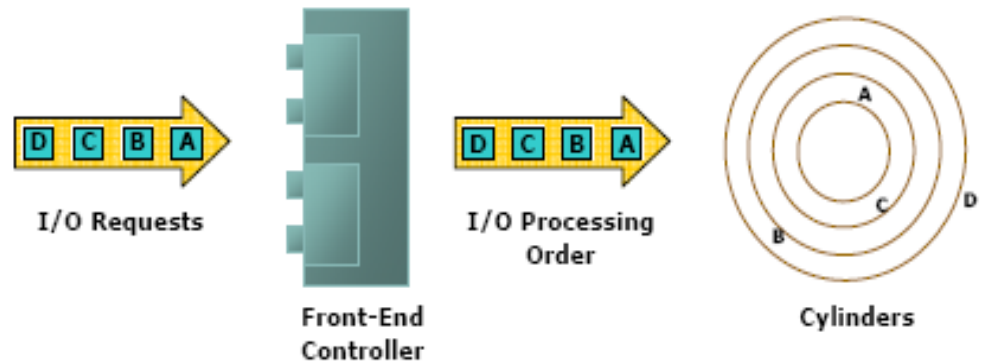
Intelligent Storage System

- The front end provides the interface between the storage system and the host.
 - two components: front-end ports and front-end controllers
- The *front-end ports* enable hosts to connect to the intelligent storage system, and has processing logic that executes the appropriate transport protocol, such as SCSI, Fibre Channel, IB, or iSCSI, for storage connections
- *Front-end controllers* route data to and from cache via the internal data bus. When cache receives write data, the controller sends an acknowledgment

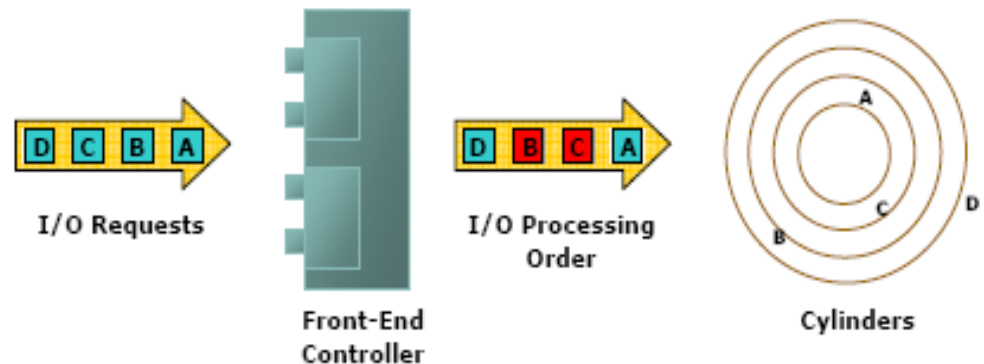


Intelligent Storage System

- Controllers optimize I/O processing by using *command queuing* algorithms
 - technique implemented on front-end controllers
 - determines the execution order of received commands and can reduce unnecessary drive head movements and improve disk performance



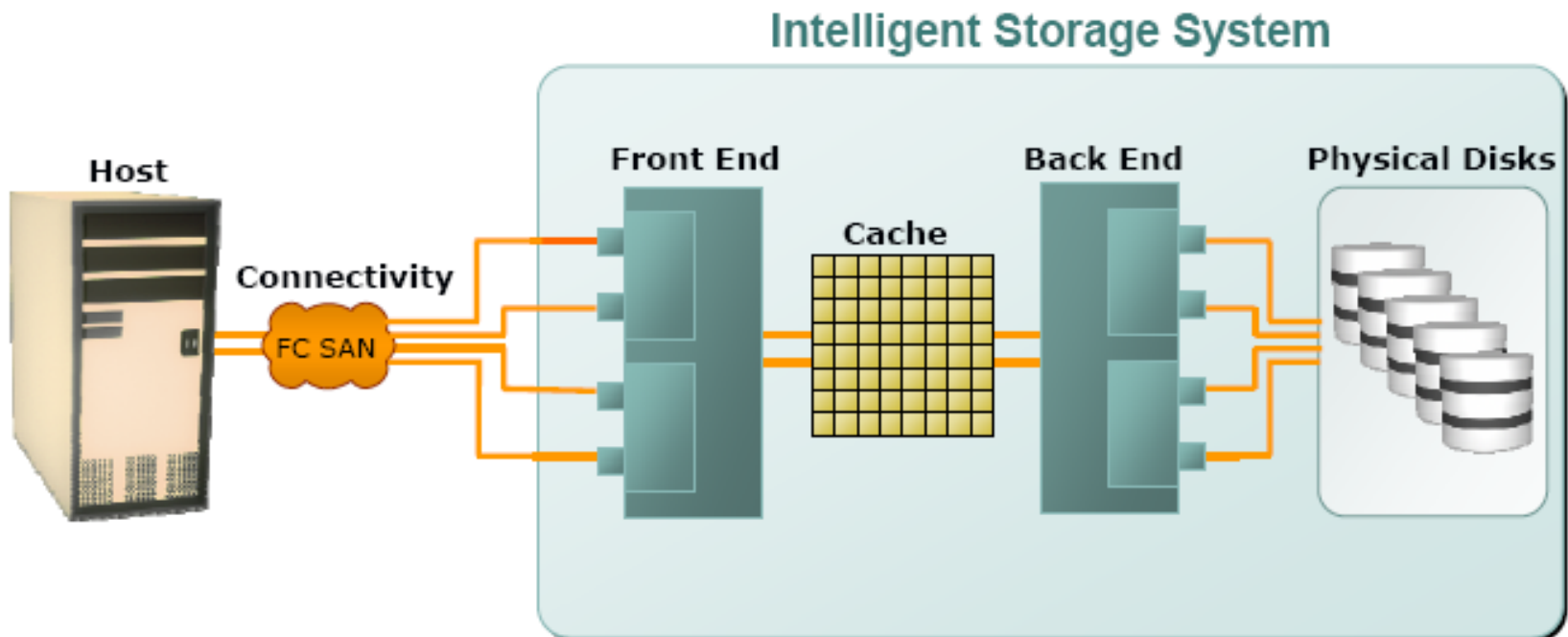
Without Optimization (FIFO)



With command queuing

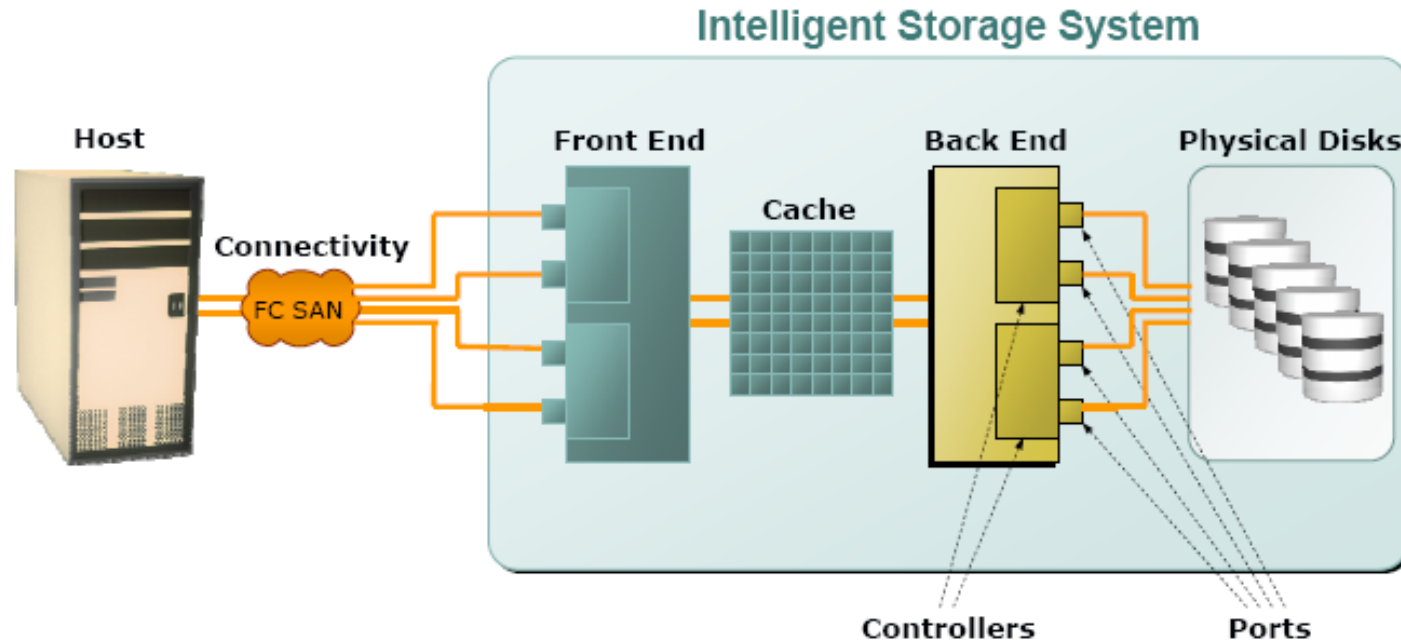
Intelligent Storage System: Cache

- Cache enhances the I/O performance in an intelligent storage system.
 - Accessing data from disk usually takes a few ms, from cache takes less than a ms.
 - Write data is placed in cache and then written to disk
- *Cache mirroring*: Each write to cache held in two different memory locations on two independent memory cards
- *Cache vaulting*: Cache exposed to risk of uncommitted data loss due to power failure
 - using battery power to write the cache content to the disk storage vendors use a set of physical disks to dump the contents of cache during power failure



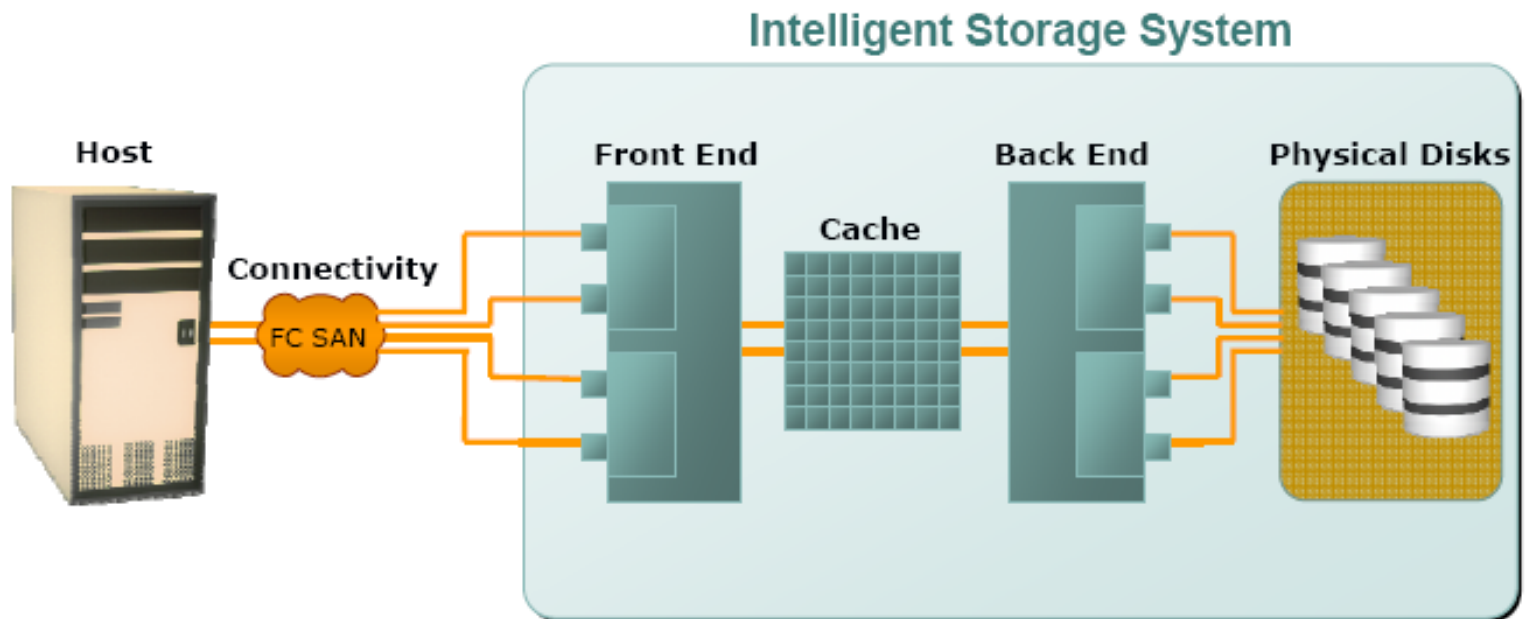
Intelligent Storage System: Back End

- Two components: back-end ports and back-end controllers
 - Physical disks are connected to ports on the back end.
- The back end controller communicates with the disks when performing reads and writes and also provides additional, but limited, temporary data storage.
- The algorithms implemented on back-end controllers provide error detection and correction, along with RAID functionality.
- Multiple controllers also facilitate load balancing



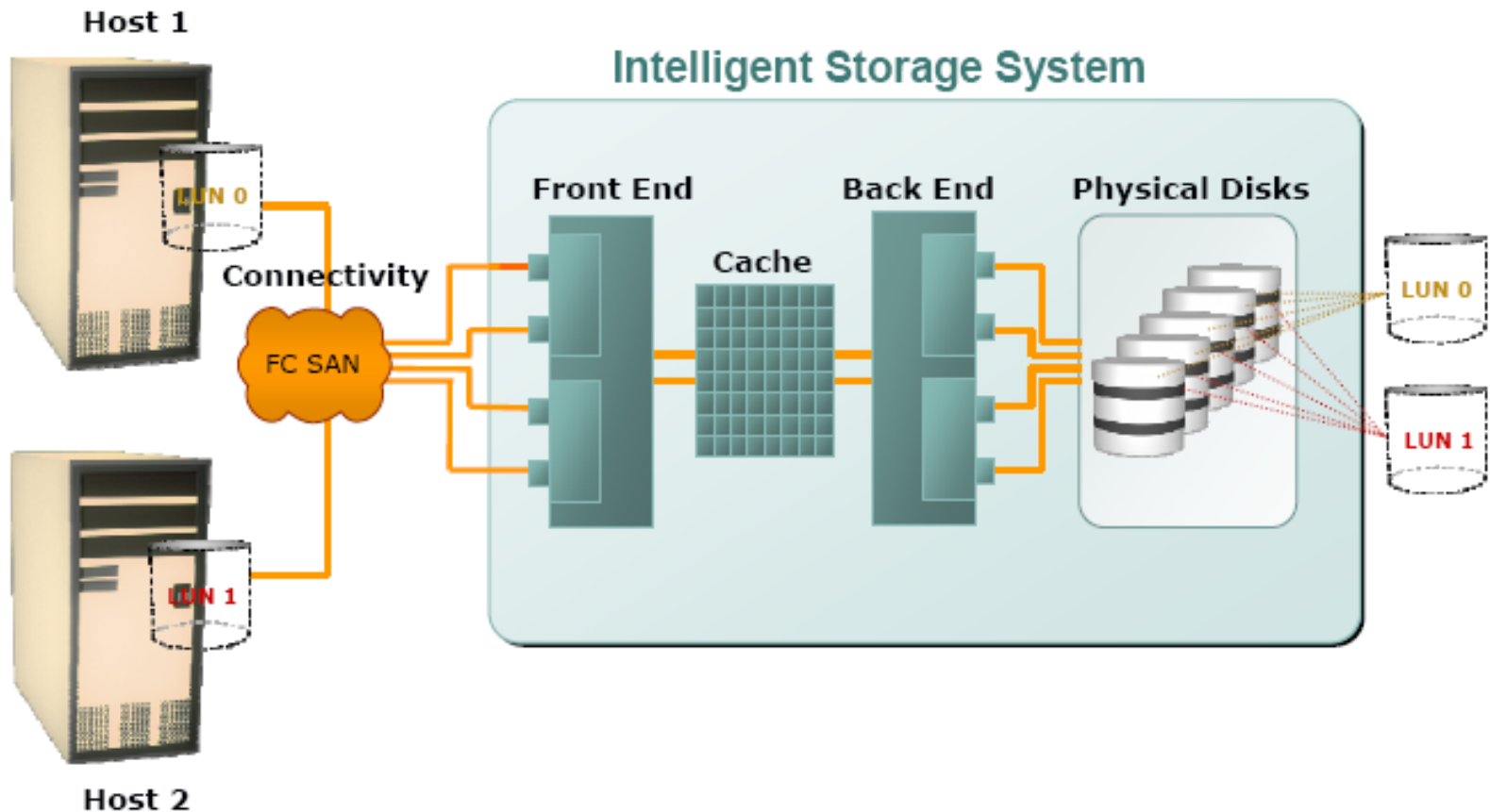
Physical Disks

- Disks are connected to the back-end with either SCSI, IB, Fibre Channel interface



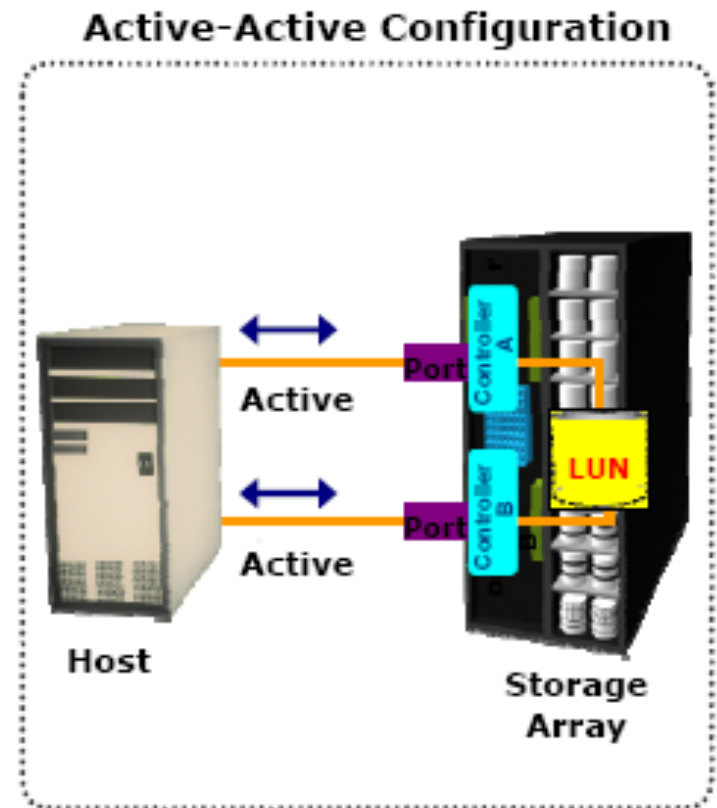
What is LUNs

- Physical drives or groups of RAID protected drives can be logically split into volumes known as logical volumes, commonly referred to as *Logical Unit Numbers* (LUNs)



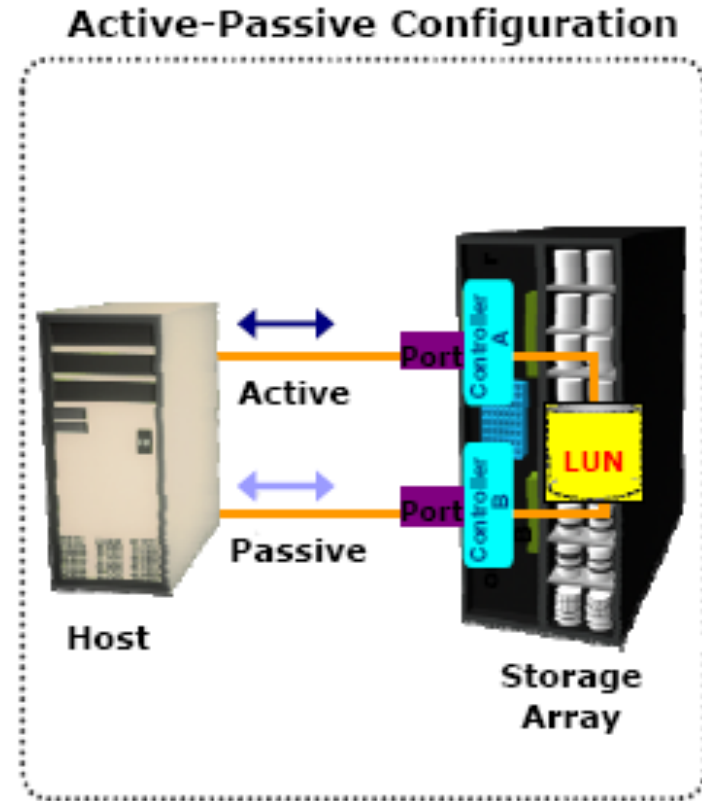
High-end Storage Systems

- High-end storage systems, referred to as *active-active arrays*, are generally aimed at large enterprises for centralizing corporate data
- These arrays are designed with a large number of controllers and cache memory
- An active-active array implies that the host can perform I/Os to its LUNs across any of the available Paths



Midrange Storage Systems

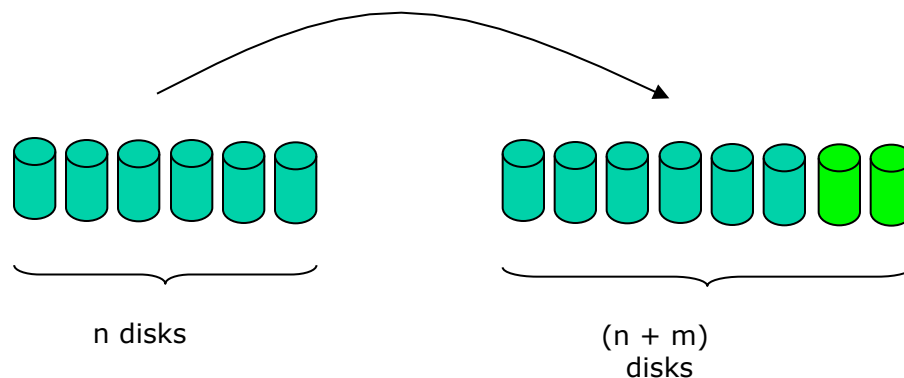
- Also referred as **Active-passive arrays**
- Host can perform I/Os to LUNs only through active paths
- Other paths remain passive till active path fails
- Midrange array have two controllers, each with cache, RAID controllers and disks drive interfaces
- Designed for small and medium enterprises
- Less scalable as compared to high-end array



Performance Metrics for Storage

- **Storage efficiency**: how much redundant information do you store?
- **Saturation throughput**: how many I/O requests can the system handle before it collapses (or delay increases to infinity)?
- **Rebuild time**: how fast can you replace information lost due to disk failure?
- **Mean time to data loss**: under assumptions on failure and usage models of the system, how long do you expect to run without any permanent loss of data?
- **Encoding/Decoding/Update/Rebuild complexity**: the computation power needed for all these operations; also, how many bytes of data on how many disks do you have to update if you just want to update 1 byte of user data?
- **Sequential read/write bandwidth**: bandwidth the system can provide for streaming data
- **Scale-out** (as opposed to “Scale-up”): Add low cost commodity hardware to increase capacity incrementally (scale-out) rather than make up-front large investment in more expensive complex hardware
- **High availability**: Cost of down time for businesses is large
- Systems cannot be taken down for backing up data
- Rebuild time should be small
- **Reliability vs. cost**: Replication based schemes for reliability are expensive

Erasure Codes



Encode data on n disks onto $(n+m)$ disks such that the whole system can tolerate up to m disk failures

- Reliability
 - Specified by avg # disk failures tolerated
 - If avg # disk failures tolerated = m , then the code is Maximum Distance Separable (MDS)
- Performance (encoding/decoding/update)
- Space usage/Rate: Rate = $n / (n+m)$
- Flexibility
 - can you arbitrarily add nodes? Change rate?
 - how does this affect failure coverage?

Examples: Reed-Solomon codes, Parity Array codes (EvenOdd coding, X code, Weaver), LDPC codes

Other Issues

- Failure Models:
 - What is a disk failure?
 - Completely unreadable? Mechanical or chip failures, for instance.
 - Latent sector errors?
 - How do you test a drive? Read all sectors and decide faulty if any one operation takes above threshold
 - Result depends on threshold ☹️
 - User's point of view: A disk has failed if the user feels it is no longer satisfactory for his/her needs
 - Measures of disk reliability
 - Manufacturers:
 - MTTF: mean time to failure (based on accelerated tests on a large number of disk)
 - AFR: annualized failure rate (percentage of disks expected to fail in a year)
 - More recent measures (from user's PoV)
 - ARR: annualized *replacement* rate (percentage of disks replaced by a user in a year)

Other Issues

- Failure Models (contd.):
 - Traditional assumptions
 - Disk failures are independent and is a Poisson process (that is, time between failures is exponentially distributed with parameter λ ; then the MTTF = $1/\lambda$)
 - Bathtub curve (plot of failure rate vs. time looks like a bathtub):
 - “infant mortality”: failure rates are high in the first few months to a year or so
 - “useful life period”: the failure rate is minimum and stays constant from a year to about 4 or 5 years
 - “wearout period”: failure rate again goes up after 4 or 5 years
 - New findings (Schroeder and Gibson, CMU, and Google, 2007) based on disk replacement (rates) rather than disk failure (rates)
 - Disk replacements are not independent and do not follow a Poisson process
 - Longer since the last disk failed => longer till the next disk will fail!
 - Possible explanation: environmental and other factors (temperature etc.) are more important than component specific factors
 - Disk replacement rates do not enter a steady state after 1 year (like a bathtub curve); instead, the replacement rates steadily increase over time
 - Different failure model => different values for performance metrics! => different design of codes for achieving performance!

Other Issues

- Usage models are as important as failure models in estimating performance metrics and designing good codes:
 - **Write once, (almost) never read:** eg. Archival storage
 - **Write once, read many times:** eg. Streaming applications like Youtube
 - **Random short reads and writes:** eg. Systems handling lots of short transactions like shopping sites, or high performance computing
- A storage system with a certain coding scheme and a certain failure model can perform significantly differently under each of the above usage models
- Unlikely that there is one coding scheme that is suited for all failure and usage models