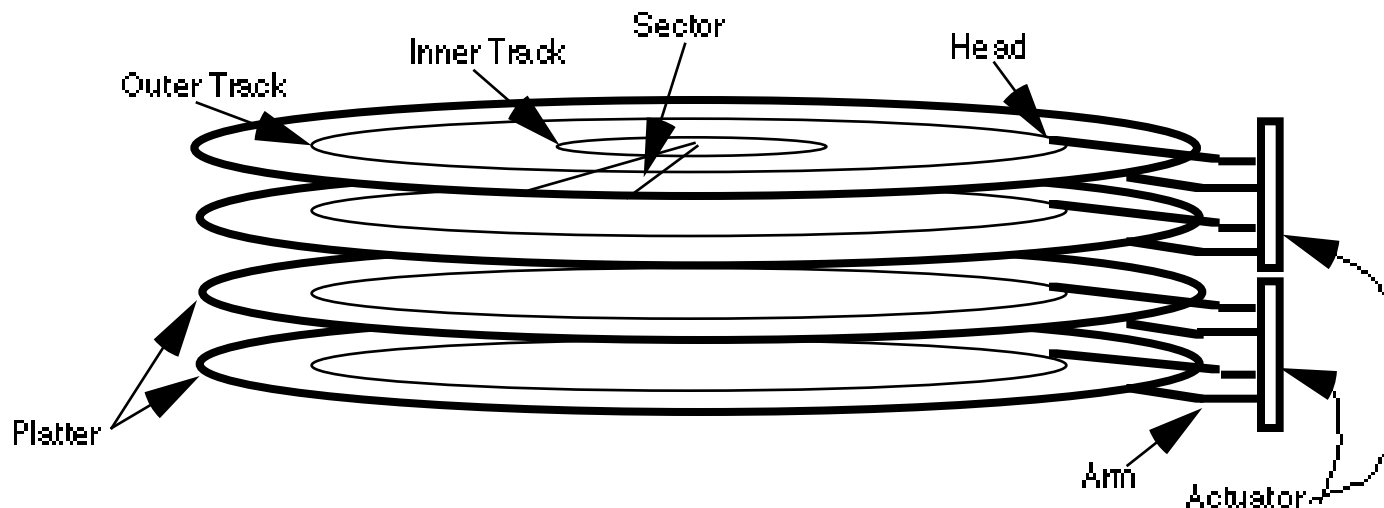


# **Lecture 13:** **I/O: A Little Queuing Theory, RAID**

**Professor David A. Patterson**  
**Computer Science 252**  
**Spring 1998**

# Review: Disk Device Terminology



**Disk Latency = Queuing Time + Seek Time + Rotation Time + Xfer Time**

***Order of magnitude times for 4K byte transfers:***

**Seek: 12 ms or less**

**Rotate: 4.2 ms @ 7200 rpm (8.3 ms @ 3600 rpm )**

**Xfer: 1 ms @ 7200 rpm (2 ms @ 3600 rpm)**

# Review

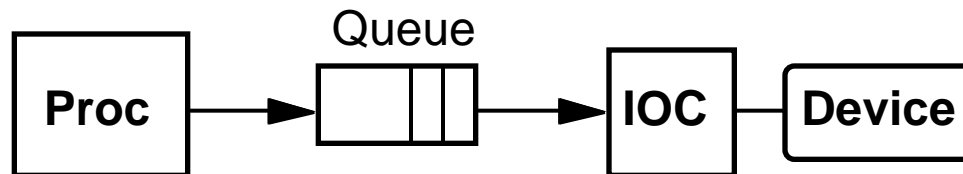
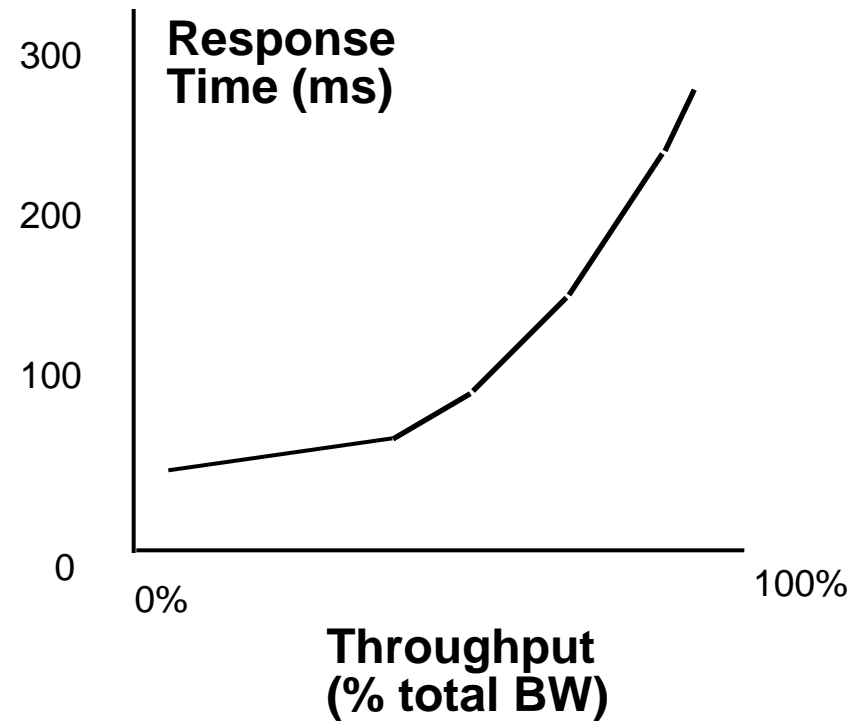
- **Disk industry growing rapidly, improves:**
  - bandwidth 40%/yr ,
  - areal density 60%/year, \$/MB faster?
- **queue + controller + seek + rotate + transfer**
- **Advertised average seek time benchmark much greater than average seek time in practice**
- **Response time vs. Bandwidth tradeoffs**
- **Value of faster response time:**
  - 0.7sec off response saves 4.9 sec and 2.0 sec (70%) total time per transaction => greater productivity
  - everyone gets more done with faster response, but novice with fast response = expert with slow
- **Processor Interface: today peripheral processors, DMA, I/O bus, interrupts**

# Review: Storage System Issues

- *Historical Context of Storage I/O*
- *Secondary and Tertiary Storage Devices*
- *Storage I/O Performance Measures*
- *Processor Interface Issues*
- **A Little Queuing Theory**
- **Redundant Arrays of Inexpensive Disks (RAID)**
- **I/O Buses**
- **ABCs of UNIX File Systems**
- **I/O Benchmarks**
- **Comparing UNIX File System Performance**

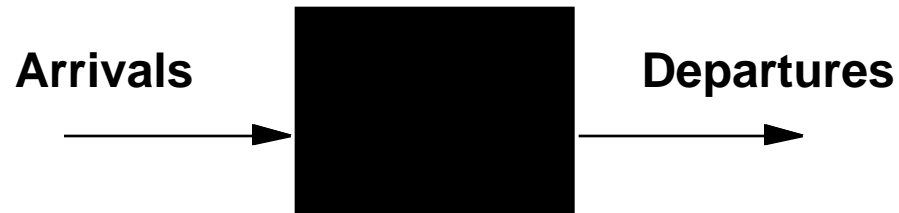
# Review: Disk I/O Performance

**Metrics:**  
**Response Time**  
**Throughput**



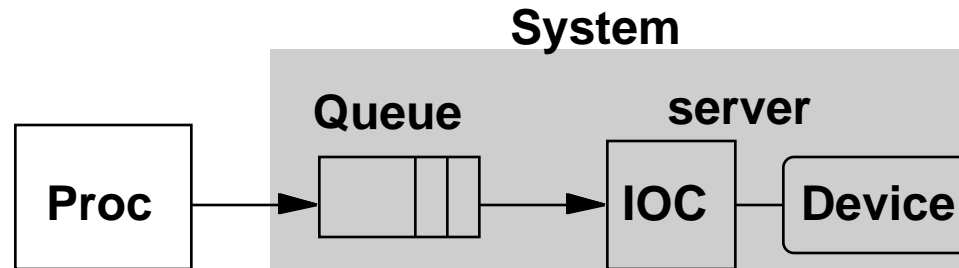
**Response time = Queue + Device Service time**

# Introduction to Queueing Theory



- More interested in long term, steady state than in startup => Arrivals = Departures
- Little's Law:  
Mean number tasks in system = arrival rate x mean response time
  - Observed by many, Little was first to prove
- Applies to any system in equilibrium, as long as nothing in black box is creating or destroying tasks

# A Little Queuing Theory: Notation



- Queuing models assume state of equilibrium: input rate = output rate
- Notation:

$r$  average number of arriving customers/second

$T_{ser}$  average time to service a customer (traditionally  $\mu = 1/ T_{ser}$ )

$u$  server utilization (0..1):  $u = r \times T_{ser}$  (or  $u = r / T_{ser}$ )

$T_q$  average time/customer in queue

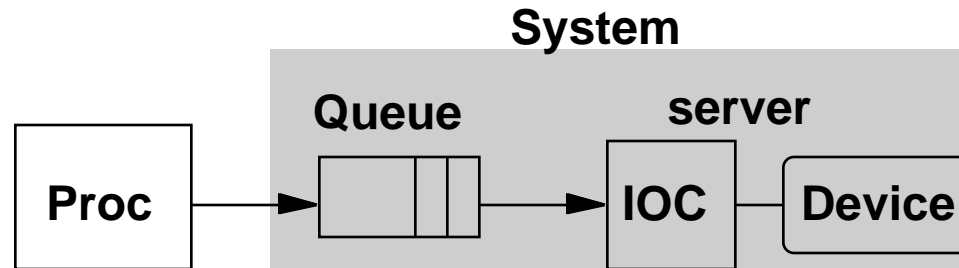
$T_{sys}$  average time/customer in system:  $T_{sys} = T_q + T_{ser}$

$L_q$  average length of queue:  $L_q = r \times T_q$

$L_{sys}$  average length of system:  $L_{sys} = r \times T_{sys}$

- Little's Law: **Length<sub>system</sub> = rate x Time<sub>system</sub>**  
(Mean number customers = arrival rate x mean service time)

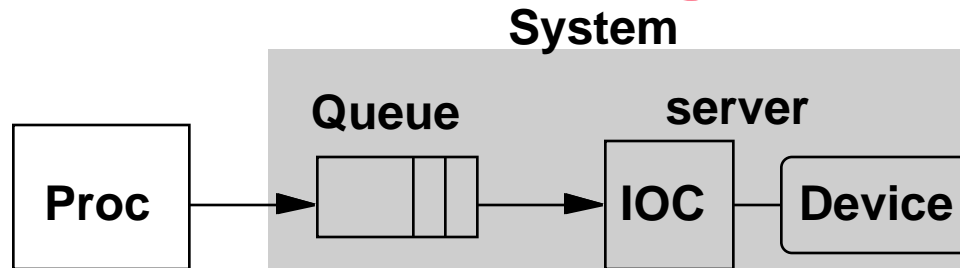
# A Little Queuing Theory



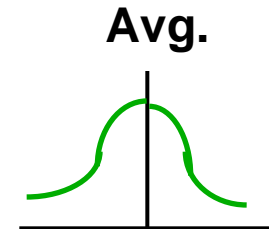
- **Service time completions vs. waiting time for a busy server:** randomly arriving event joins a queue of arbitrary length when server is busy, otherwise serviced immediately
  - Unlimited length queues key simplification
- A **single server queue**: combination of a servicing facility that accomodates 1 customer at a time (**server**) + waiting area (**queue**): together called a **system**
- Server spends a variable amount of time with customers; **how do you characterize variability?**
  - Distribution of a random variable: histogram? curve?



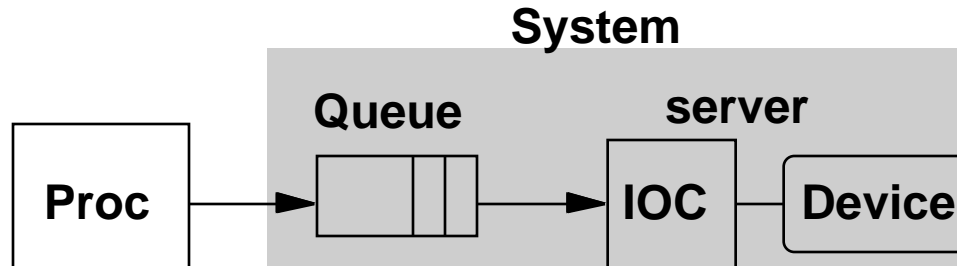
# A Little Queuing Theory



- **Server spends a variable amount of time with customers**
  - Weighted mean  $m1 = (f1 \times T1 + f2 \times T2 + \dots + fn \times Tn)/F$  ( $F=f1 + f2\dots$ )
  - **variance**  $= (f1 \times T1^2 + f2 \times T2^2 + \dots + fn \times Tn^2)/F - m1^2$ 
    - » Must keep track of unit of measure (100 ms<sup>2</sup> vs. 0.1 s<sup>2</sup>)
  - **Squared coefficient of variance:  $C = variance/m1^2$** 
    - » Unitless measure (100 ms<sup>2</sup> vs. 0.1 s<sup>2</sup>)
- **Exponential distribution  $C = 1$**  : most short relative to average, few others long; 90% < 2.3 x average, 63% < average
- **Hypoexponential distribution  $C < 1$**  : most close to average,  $C=0.5 \Rightarrow$  90% < 2.0 x average, only 57% < average
- **Hyperexponential distribution  $C > 1$**  : further from average  $C=2.0 \Rightarrow$  90% < 2.8 x average, 69% < average



# A Little Queuing Theory: Variable Service Time



- **Server spends a variable amount of time with customers**
  - Weighted mean  $m1 = (f1 \times T1 + f2 \times T2 + \dots + fn \times Tn) / F$  ( $F = f1 + f2 + \dots$ )
  - Squared coefficient of variance  $C$
- **Disk response times  $C \approx 1.5$  (majority seeks < average)**
- **Yet usually pick  $C = 1.0$  for simplicity**
- **Another useful value is average time must wait for server to complete task:  $m1(z)$** 
  - Not just  $1/2 \times m1$  because doesn't capture variance
  - Can derive  $m1(z) = 1/2 \times m1 \times (1 + C)$
  - **No variance  $\Rightarrow C = 0 \Rightarrow m1(z) = 1/2 \times m1$**

# A Little Queuing Theory: Average Wait Time

- Calculating average wait time in queue  $T_q$ 
  - If something at server, it takes to complete on average  $m1(z)$
  - Chance server is busy =  $u$ ; average delay is  $u \times m1(z)$
  - All customers in line must complete; each avg  $T_{ser}$

$$T_q = u \times m1(z) + L_q \times T_{ser} = 1/2 \times u \times T_{ser} \times (1 + C) + L_q \times T_{ser}$$

$$T_q = 1/2 \times u \times T_{ser} \times (1 + C) + r \times T_q \times T_{ser}$$

$$T_q = 1/2 \times u \times T_{ser} \times (1 + C) + \frac{u \times T_q}{1 - u}$$

$$T_q \times (1 - u) = T_{ser} \times u \times (1 + C) / 2$$

$$T_q = T_{ser} \times u \times (1 + C) / (2 \times (1 - u))$$

- Notation:

$r$  average number of arriving customers/second

$T_{ser}$  average time to service a customer

$u$  server utilization (0..1):  $u = r \times T_{ser}$

$T_q$  average time/customer in queue

$L_q$  average length of queue:  $L_q = r \times T_q$

# A Little Queuing Theory: M/G/1 and M/M/1

- Assumptions so far:
  - System in equilibrium
  - Time between two successive arrivals in line are random
  - Server can start on next customer immediately after prior finishes
  - No limit to the queue: works First-In-First-Out
  - Afterward, all customers in line must complete; each avg  $T_{ser}$
- Described “memoryless” or Markovian request arrival (M for C=1 exponentially random), General service distribution (no restrictions), 1 server: **M/G/1 queue**
- When Service times have C = 1, **M/M/1 queue**

$$T_q = T_{ser} \times u \times (1 + C) / (2 \times (1 - u)) = T_{ser} \times u / (1 - u)$$

$T_{ser}$  average time to service a customer

$u$  server utilization (0..1):  $u = r \times T_{ser}$

$T_q$  average time/customer in queue

# A Little Queuing Theory: An Example

- processor sends 10 x 8KB disk I/Os per second, requests & service exponentially distrib., avg. disk service = 20 ms
- On average, how utilized is the disk?
  - What is the number of requests in the queue?
  - What is the average time spent in the queue?
  - What is the average response time for a disk request?
- Notation:

$r$  average number of arriving customers/second = 10

$T_{ser}$  average time to service a customer = 20 ms (0.02s)

$u$  server utilization (0..1):  $u = r \times T_{ser} = 10/s \times .02s = 0.2$

$T_q$  average time/customer in queue =  $T_{ser} \times u / (1 - u)$   
=  $20 \times 0.2 / (1 - 0.2) = 20 \times 0.25 = 5 \text{ ms}$  (0.005s)

$T_{sys}$  average time/customer in system:  $T_{sys} = T_q + T_{ser} = 25 \text{ ms}$

$L_q$  average length of queue:  $L_q = r \times T_q$   
=  $10/s \times .005s = 0.05 \text{ requests in queue}$

$L_{sys}$  average # tasks in system:  $L_{sys} = r \times T_{sys} = 10/s \times .025s = 0.25$

# CS 252 Administrivia

- Email URL of initial project home page to TA?
  - Upcoming events in CS 252
- 18-Mar Wed I/O 3: Tertiary Storage & Network Intro
- 19-Mar Thu Send in Project Survey #2
- 20-Mar Fri Networks 2: Interface, Switches, Routing
- 23-Mar to 27-Mar Spring Break

# A Little Queuing Theory: Another Example

- processor sends 20 x 8KB disk I/Os per sec, requests & service exponentially distrib., avg. disk service = 12 ms
- On average, how utilized is the disk?
  - What is the number of requests in the queue?
  - What is the average time a spent in the queue?
  - What is the average response time for a disk request?

- Notation:

$r$  average number of arriving customers/second= 20

$T_{ser}$  average time to service a customer= 12 ms

$u$  server utilization (0..1):  $u = r \times T_{ser} = \underline{\quad}/s \times \underline{\quad} s = \underline{\quad}$

$T_q$  average time/customer in queue =  $T_{ser} \times u / (1 - u)$   
 $= \underline{\quad} \times \underline{\quad} / (\underline{\quad}) = \underline{\quad} \times \underline{\quad} = \underline{\quad} \text{ ms}$

$T_{sys}$  average time/customer in system:  $T_{sys} = T_q + T_{ser} = 16 \text{ ms}$

$L_q$  average length of queue:  $L_q = r \times T_q$   
 $= \underline{\quad}/s \times \underline{\quad} s = \underline{\quad}$  requests in queue

$L_{sys}$  average # tasks in system :  $L_{sys} = r \times T_{sys} = \underline{\quad}/s \times \underline{\quad} s = \underline{\quad}$

# A Little Queuing Theory: Another Example

- processor sends 20 x 8KB disk I/Os per sec, requests & service exponentially distrib., avg. disk service = 12 ms
- On average, how utilized is the disk?
  - What is the number of requests in the queue?
  - What is the average time a spent in the queue?
  - What is the average response time for a disk request?

- Notation:

$r$  average number of arriving customers/second= 20

$T_{ser}$  average time to service a customer= 12 ms

$u$  server utilization (0..1):  $u = r \times T_{ser} = 20/s \times .012s = 0.24$

$T_q$  average time/customer in queue =  $T_{ser} \times u / (1 - u)$   
 $= 12 \times 0.24 / (1 - 0.24) = 12 \times 0.32 = 3.8 \text{ ms}$

$T_{sys}$  average time/customer in system:  $T_{sys} = T_q + T_{ser} = 15.8 \text{ ms}$

$L_q$  average length of queue:  $L_q = r \times T_q$   
 $= 20/s \times .0038s = 0.076 \text{ requests in queue}$

$L_{sys}$  average # tasks in system :  $L_{sys} = r \times T_{sys} = 20/s \times .016s = 0.32$



# A Little Queuing Theory: Yet Another Example

- Suppose processor sends 10 x 8KB disk I/Os per second, squared coef. var.(C) = 1.5, avg. disk service time = 20 ms
- On average, how utilized is the disk?
  - What is the number of requests in the queue?
  - What is the average time a spent in the queue?
  - What is the average response time for a disk request?

- **Notation:**

- $r$  average number of arriving customers/second= 10
- $T_{ser}$  average time to service a customer= 20 ms
- $u$  server utilization (0..1):  $u = r \times T_{ser} = 10/s \times .02s = 0.2$
- $T_q$  average time/customer in queue =  $T_{ser} \times u \times (1 + C) / (2 \times (1 - u))$   
 $= 20 \times 0.2(2.5)/2(1 - 0.2) = 20 \times 0.32 = 6.25$  ms
- $T_{sys}$  average time/customer in system:  $T_{sys} = T_q + T_{ser} = 26$  ms
- $L_q$  average length of queue:  $L_q = r \times T_q$   
 $= 10/s \times .006s = 0.06$  requests in queue
- $L_{sys}$  average # tasks in system :  $L_{sys} = r \times T_{sys} = 10/s \times .026s = 0.26$

# Pitfall of Not using Queuing Theory

- **1st 32-bit minicomputer (VAX-11/780)**
- **How big should write buffer be?**
  - Stores 10% of instructions, 1 MIPS
- **Buffer = 1**
- **=> Avg. Queue Length = 1  
vs. low response time**

# Review: Storage System Issues

- *Historical Context of Storage I/O*
- *Secondary and Tertiary Storage Devices*
- *Storage I/O Performance Measures*
- *Processor Interface Issues*
- *A Little Queuing Theory*
- **Redundant Arrays of Inexpensive Disks (RAID)**
- **I/O Buses**
- **ABCs of UNIX File Systems**
- **I/O Benchmarks**
- **Comparing UNIX File System Performance**

# Network Attached Storage

## *Decreasing Disk Diameters*

14" » 10" » 8" » 5.25" » 3.5" » 2.5" » 1.8" » 1.3" » ...  
high bandwidth disk systems based on arrays of disks

Network provides well defined physical and logical interfaces:  
*separate CPU and storage system!*

**High Performance  
Storage Service  
on a High Speed  
Network**

*Network File Services*

OS structures supporting remote file access

3 Mb/s » 10Mb/s » 50 Mb/s » 100 Mb/s » 1 Gb/s » 10 Gb/s  
networks capable of sustaining high bandwidth transfers

## *Increasing Network Bandwidth*

# Manufacturing Advantages of Disk Arrays

## Disk Product Families

Conventional:  
4 disk  
designs

3.5"

5.25"

10"

14"

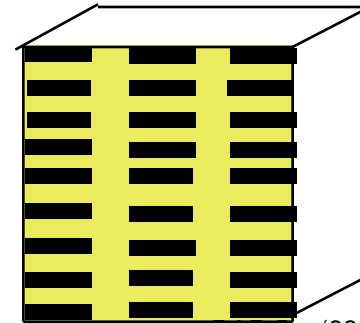
Low End



High End

Disk Array:  
1 disk  
design

3.5"



# Replace Small # of Large Disks with Large # of Small Disks! (1988 Disks)

	IBM 3390 (K)	IBM 3.5" 0061	x70
<b>Data Capacity</b>	20 GBytes	320 MBytes	23 GBytes
<b>Volume</b>	97 cu. ft.	0.1 cu. ft.	11 cu. ft.
<b>Power</b>	3 KW	11 W	1 KW
<b>Data Rate</b>	15 MB/s	1.5 MB/s	120 MB/s
<b>I/O Rate</b>	600 I/Os/s	55 I/Os/s	3900 IOs/s
<b>MTTF</b>	250 KHrs	50 KHrs	??? Hrs
<b>Cost</b>	\$250K	\$2K	\$150K

*Disk Arrays have potential for*

- large data and I/O rates
- high MB per cu. ft., high MB per KW
- reliability?

# Array Reliability

- **Reliability of N disks = Reliability of 1 Disk  $\div$  N**

**50,000 Hours  $\div$  70 disks = 700 hours**

**Disk system MTTF: Drops from 6 years to 1 month!**

- **Arrays (without redundancy) too unreliable to be useful!**

**Hot spares support reconstruction in parallel with access: very high media availability can be achieved**

# ***Redundant Arrays of Disks***

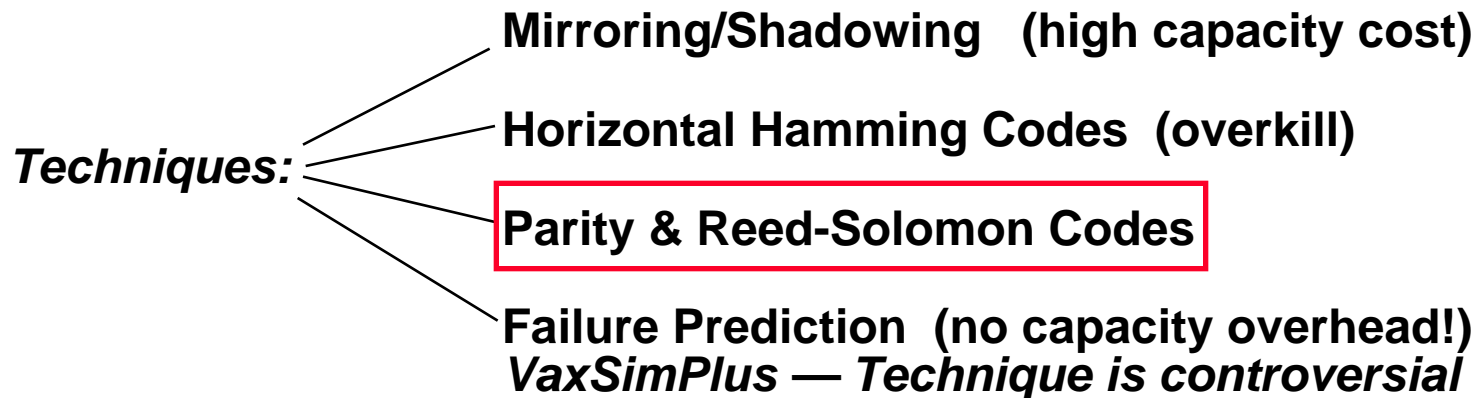
- **Files are "striped" across multiple spindles**
- **Redundancy yields high data availability**

**Disks will fail**

**Contents reconstructed from data redundantly stored in the array**

→ **Capacity penalty to store it**

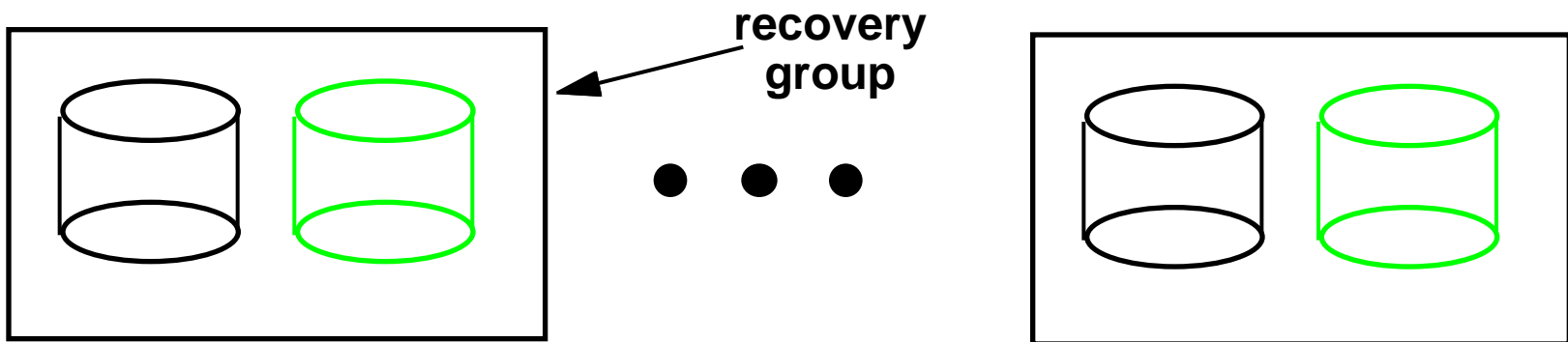
→ **Bandwidth penalty to update**





# Redundant Arrays of Disks

## RAID 1: Disk Mirroring/Shadowing

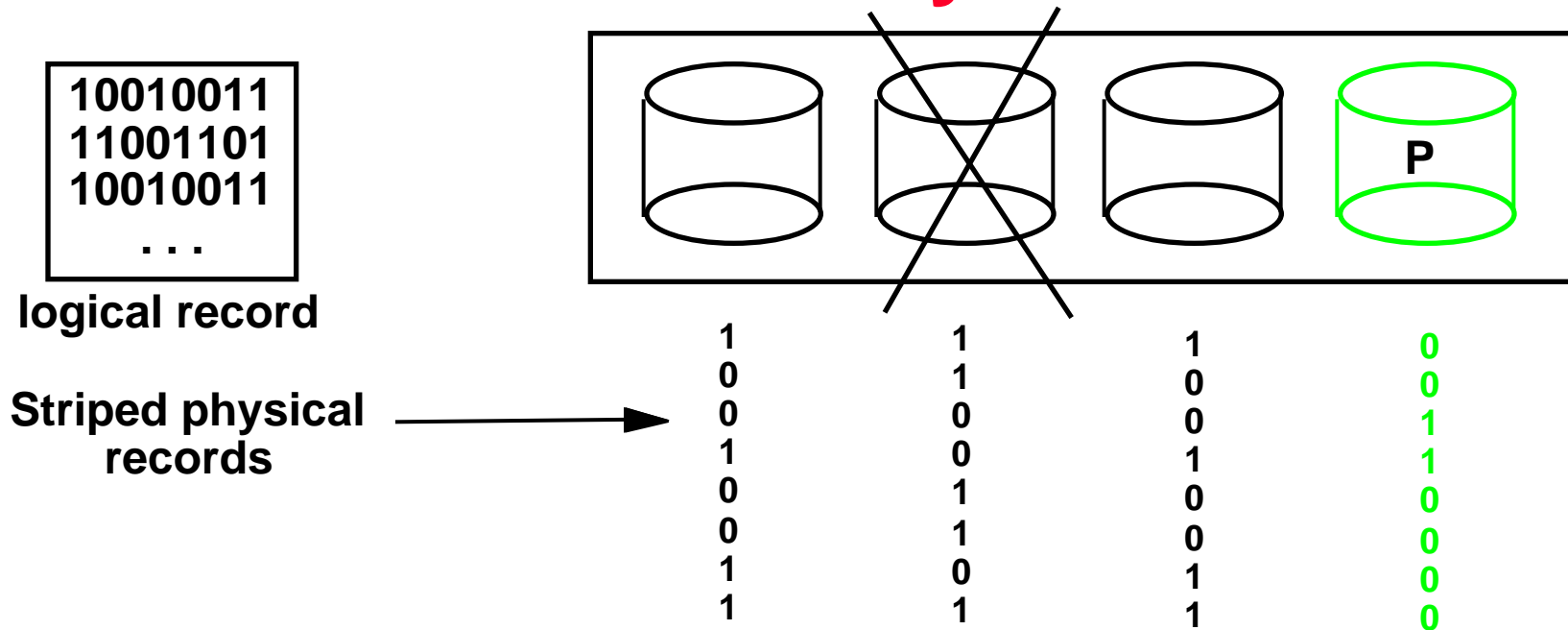


- Each disk is fully duplicated onto its "shadow"  
Very high availability can be achieved
- Bandwidth sacrifice on write:  
Logical write = two physical writes
- Reads may be optimized
- Most expensive solution: 100% capacity overhead

*Targeted for high I/O rate , high availability environments*

# Redundant Arrays of Disks

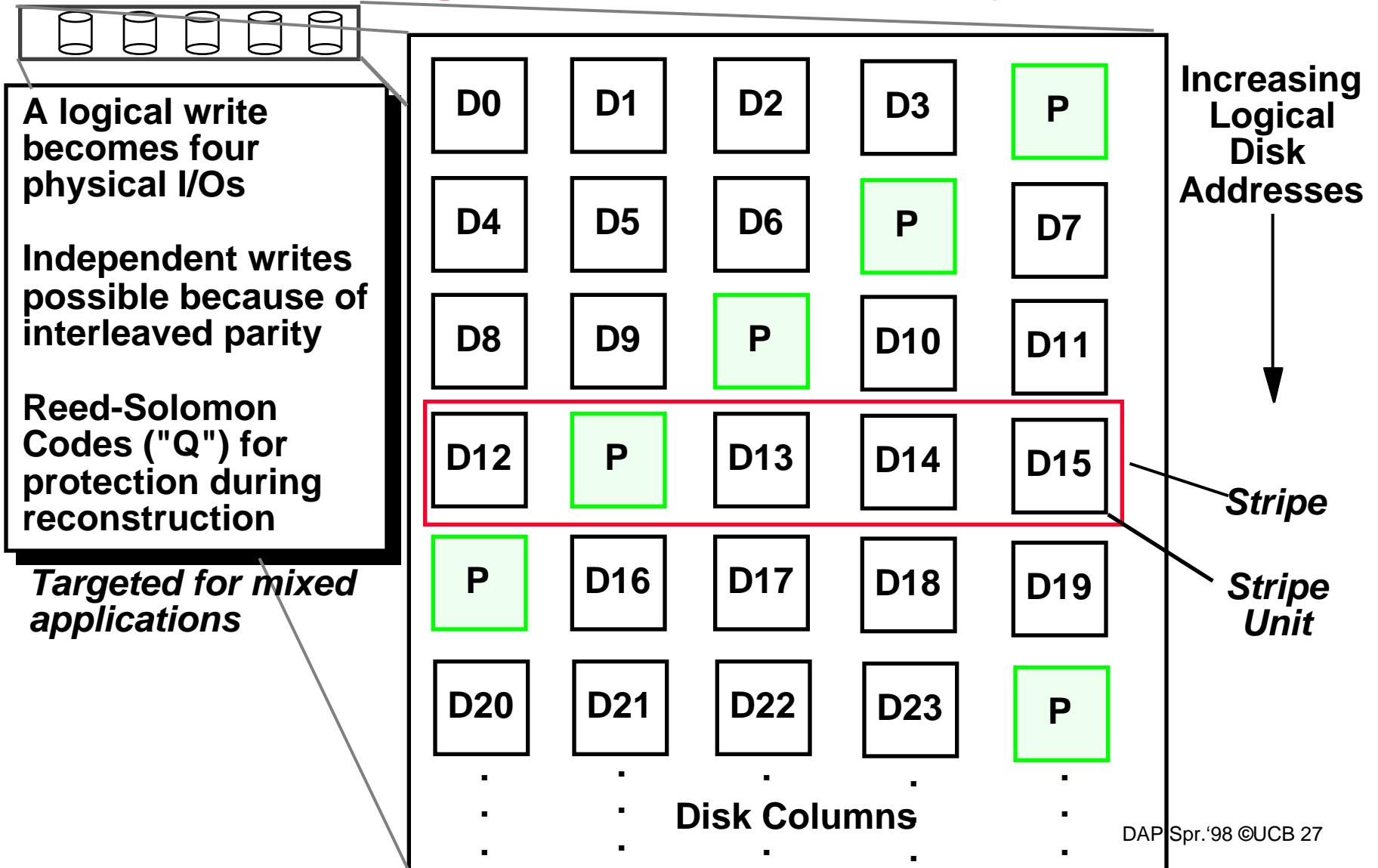
## RAID 3: Parity Disk



- Parity computed across recovery group to protect against hard disk failures
  - 33% capacity cost for parity in this configuration
  - wider arrays reduce capacity costs, decrease expected availability, increase reconstruction time
- Arms logically synchronized, spindles rotationally synchronized
  - logically a single high capacity, high transfer rate disk

*Targeted for high bandwidth applications: Scientific, Image Processing*

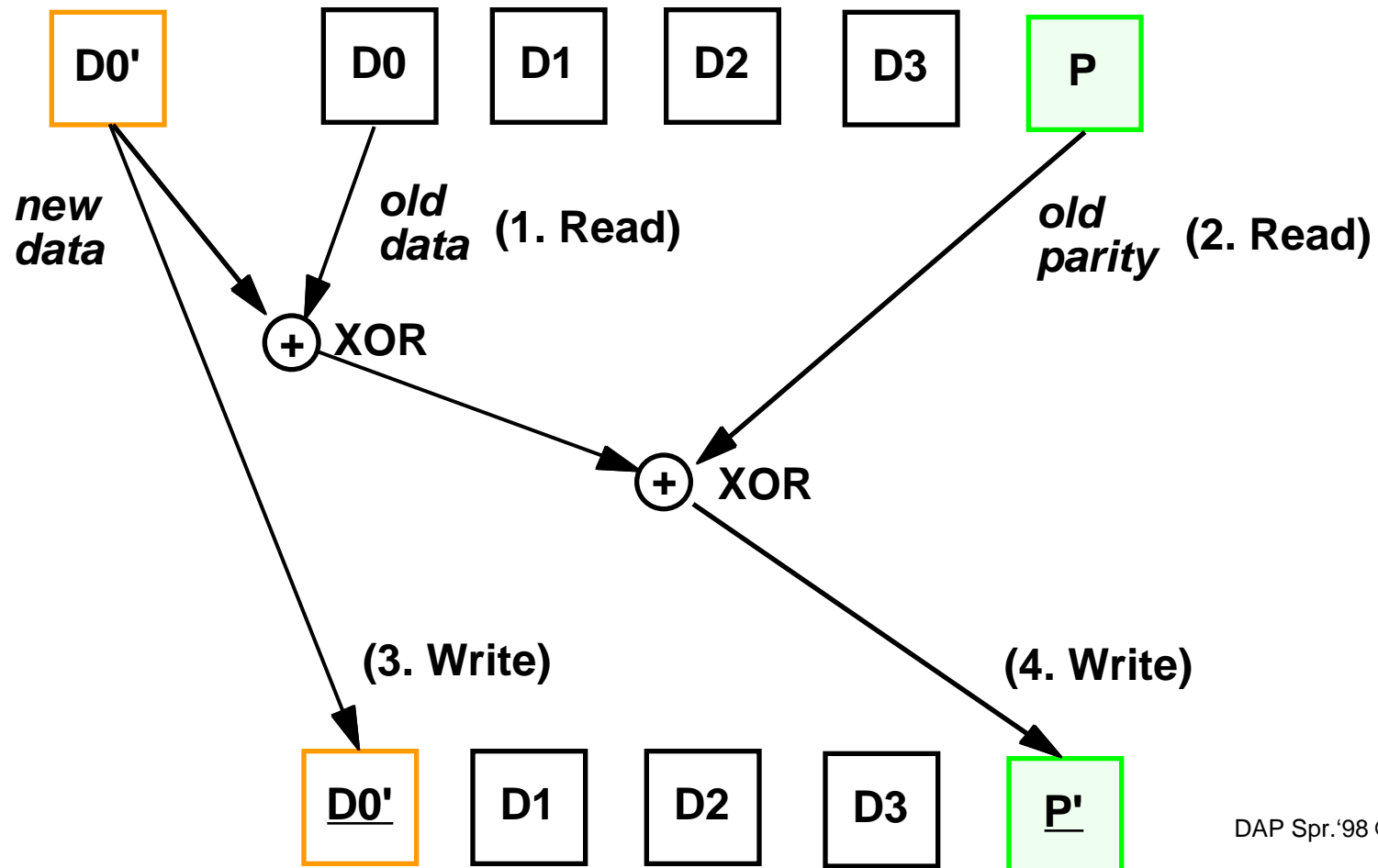
# Redundant Arrays of Disks RAID 5+: High I/O Rate Parity



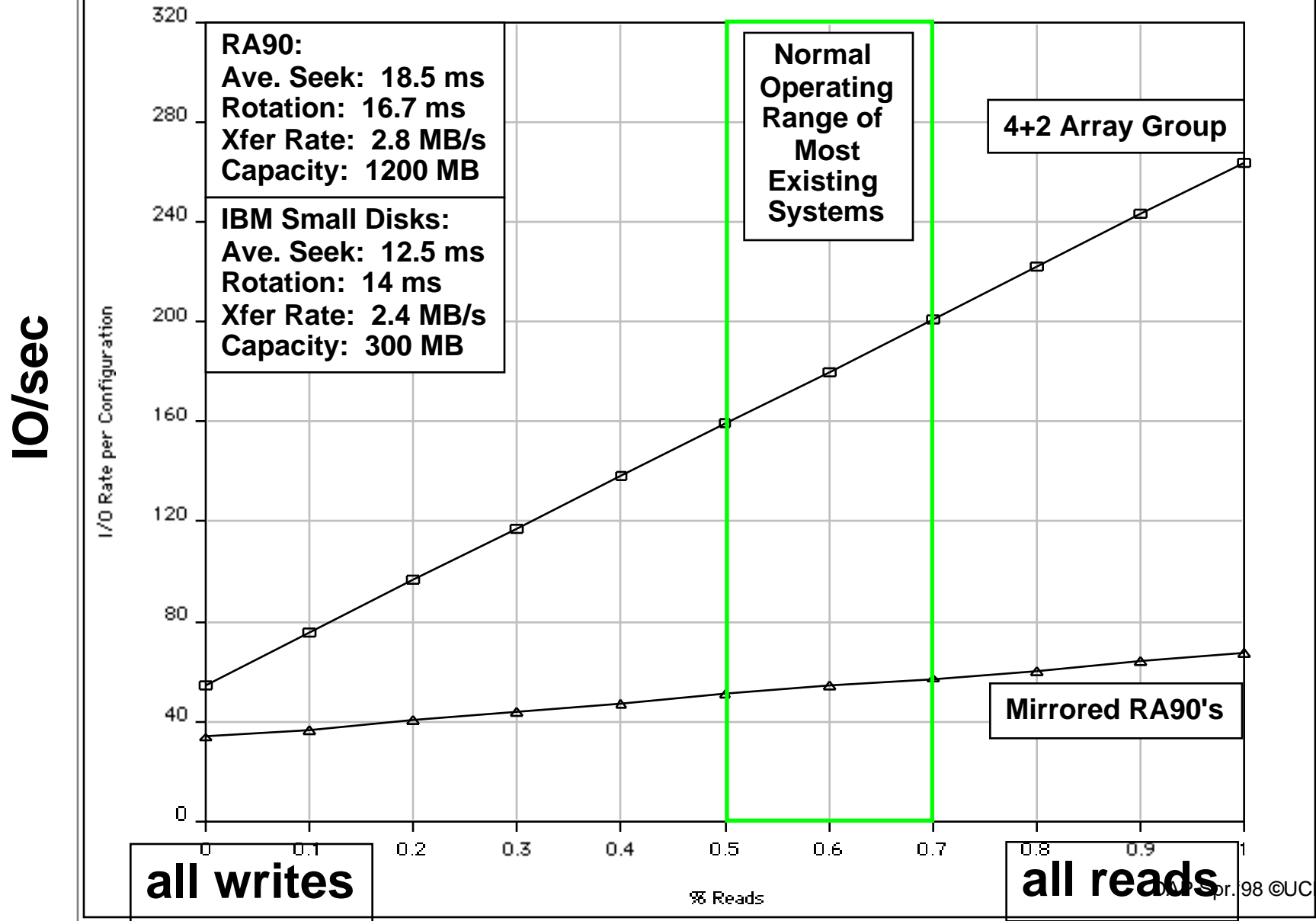
# Problems of Disk Arrays: Small Writes

## RAID-5: Small Write Algorithm

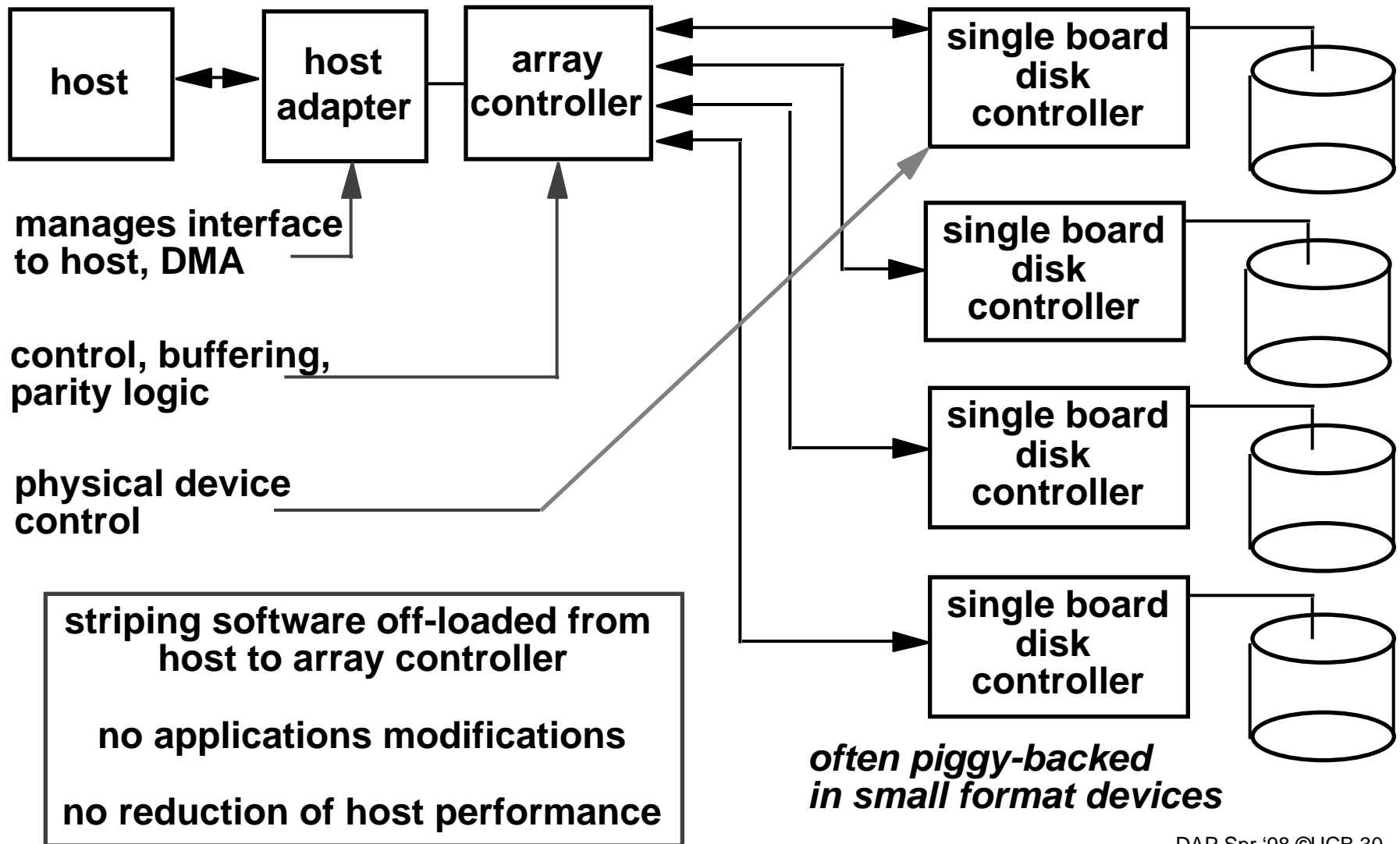
1 Logical Write = 2 Physical Reads + 2 Physical Writes



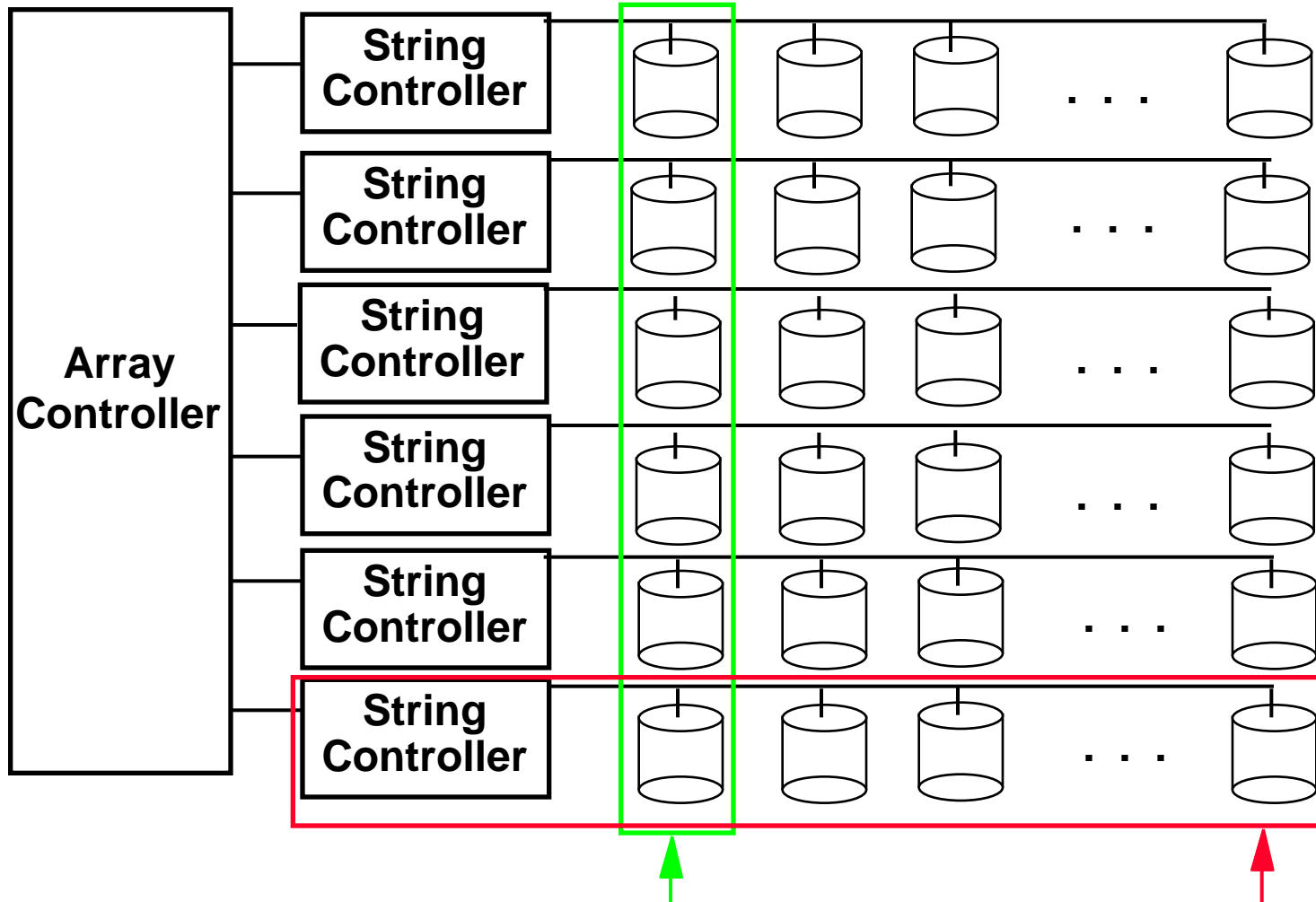
Comparison of I/O Rates for 1 Mirrored Pair of RA90's vs. a 4+2 P&Q Group



# Subsystem Organization



# System Availability: Orthogonal RAIDs

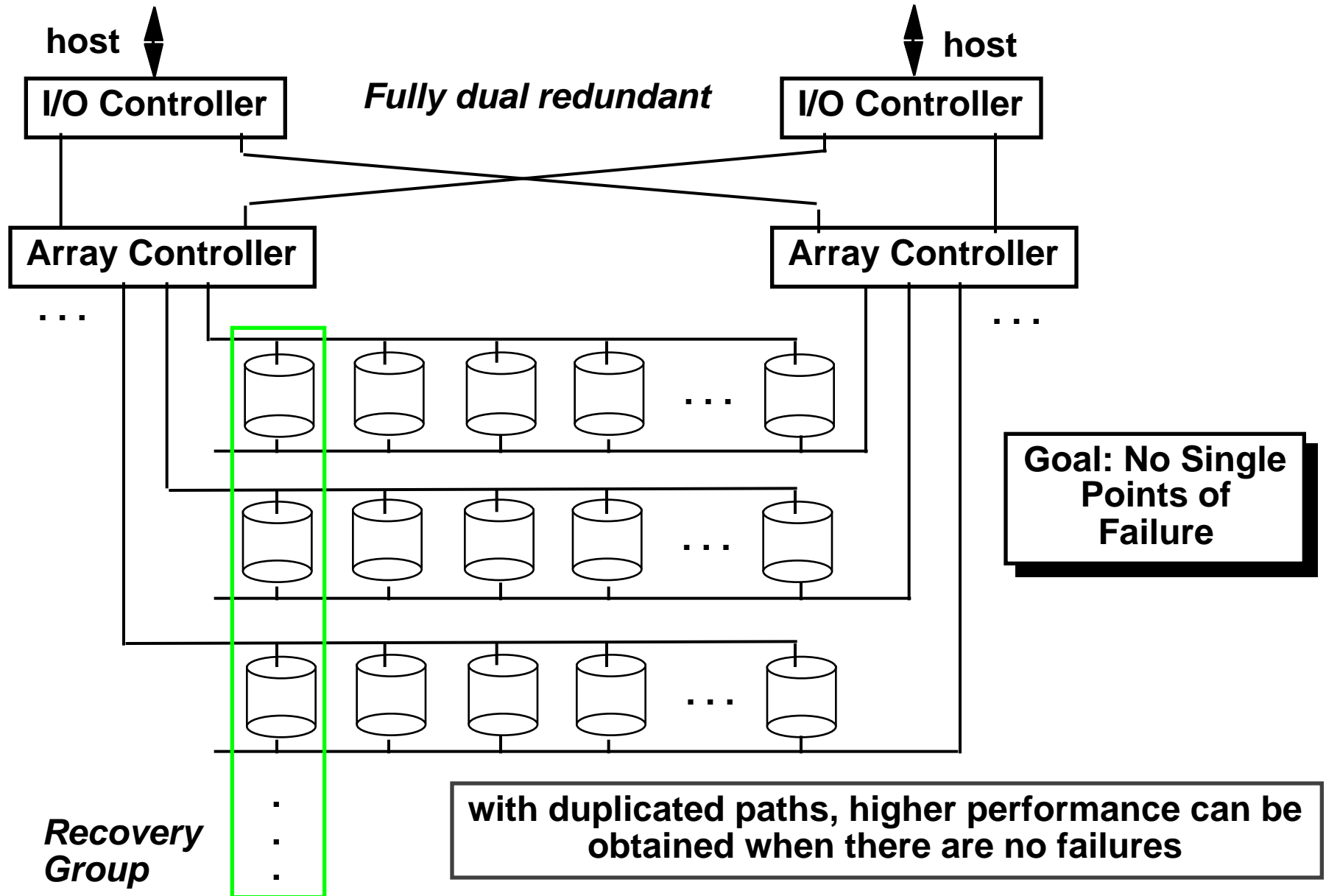


**Data Recovery Group:** unit of data redundancy

**Redundant Support Components:** fans, power supplies, controller, cables

**End to End Data Integrity:** internal parity protected data paths

# System-Level Availability





# Review: Storage System Issues

- *Historical Context of Storage I/O*
- *Secondary and Tertiary Storage Devices*
- *Storage I/O Performance Measures*
- *Processor Interface Issues*
- *A Little Queuing Theory*
- *Redundant Arrays of Inexpensive Disks (RAID)*
- **I/O Buses**
- **ABCs of UNIX File Systems**
- **I/O Benchmarks**
- **Comparing UNIX File System Performance**

# Interconnect Trends

- Interconnect = glue that interfaces computer system components
- High speed hardware interfaces + logical protocols
- Networks, channels, backplanes

	Network	Channel	Backplane
<b>Distance</b>	>1000 m	10 - 100 m	1 m
<b>Bandwidth</b>	10 - 100 Mb/s	40 - 1000 Mb/s	320 - 1000+ Mb/s
<b>Latency</b>	high (>ms)	medium	low (<μs)
<b>Reliability</b>	low Extensive CRC	medium Byte Parity	high Byte Parity

**message-based  
narrow pathways  
distributed arb**



**memory-mapped  
wide pathways  
centralized arb**

# Backplane Architectures

Metric	VME	FutureBus	MultiBus II	SCSI-I
<i>Bus Width (signals)</i>	128	96	96	25
<i>Address/Data Multiplexed?</i>	No	Yes	Yes	na
<i>Data Width</i>	16 - 32	32	32	8
<i>Xfer Size</i>	Single/Multiple	Single/Multiple	Single/Multiple	Single/Multiple
<i># of Bus Masters</i>	Multiple	Multiple	Multiple	Multiple
<i>Split Transactions</i>	No	Optional	Optional	Optional
<i>Clocking</i>	Async	Async	Sync	Either
<i>Bandwidth, Single Word (0 ns mem)</i>	25	37	20	5, 1.5
<i>Bandwidth, Single Word (150 ns mem)</i>	12.9	15.5	10	5, 1.5
<i>Bandwidth Multiple Word (0 ns mem)</i>	27.9	95.2	40	5, 1.5
<i>Bandwidth Multiple Word (150 ns mem)</i>	13.6	20.8	13.3	5, 1.5
<i>Max # of devices</i>	21	20	21	7
<i>Max Bus Length</i>	.5 m	.5 m	.5 m	25 m
<i>Standard</i>	IEEE 1014	IEEE 896	ANSI/IEEE 1296	ANSI X3.131

## Distinctions begin to blur:

**SCSI channel is like a bus**

**FutureBus is like a channel (disconnect/reconnect)**

**HIPPI forms links in high speed switching fabrics**

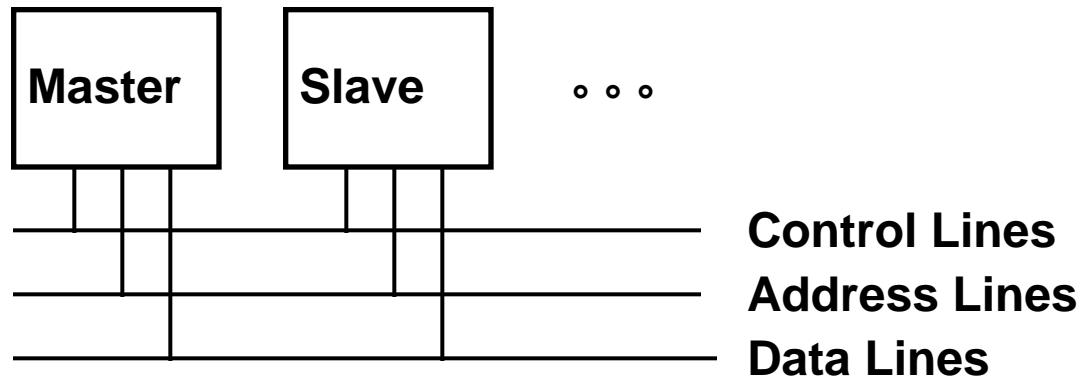
# Bus-Based Interconnect

- **Bus: a shared communication link between subsystems**
  - **Low cost: a single set of wires is shared multiple ways**
  - **Versatility: Easy to add new devices & peripherals may even be ported between computers using common bus**
- **Disadvantage**
  - **A communication bottleneck, possibly limiting the maximum I/O throughput**
- **Bus speed is limited by physical factors**
  - **the bus length**
  - **the number of devices (and, hence, bus loading).**
  - **these physical limits prevent arbitrary bus speedup.**

# Bus-Based Interconnect

- **Two generic types of busses:**
  - I/O busses: lengthy, many types of devices connected, wide range in the data bandwidth), and follow a bus standard (sometimes called a *channel*)
  - CPU–memory buses: high speed, matched to the memory system to maximize memory–CPU bandwidth, single device (sometimes called a *backplane*)
  - To lower costs, low cost (older) systems combine together
- **Bus transaction**
  - Sending address & receiving or sending data

# Bus Protocols



**Multibus: 20 address, 16 data, 5 control, 50ns Pause**

**Bus Master:** has ability to control the bus, initiates transaction

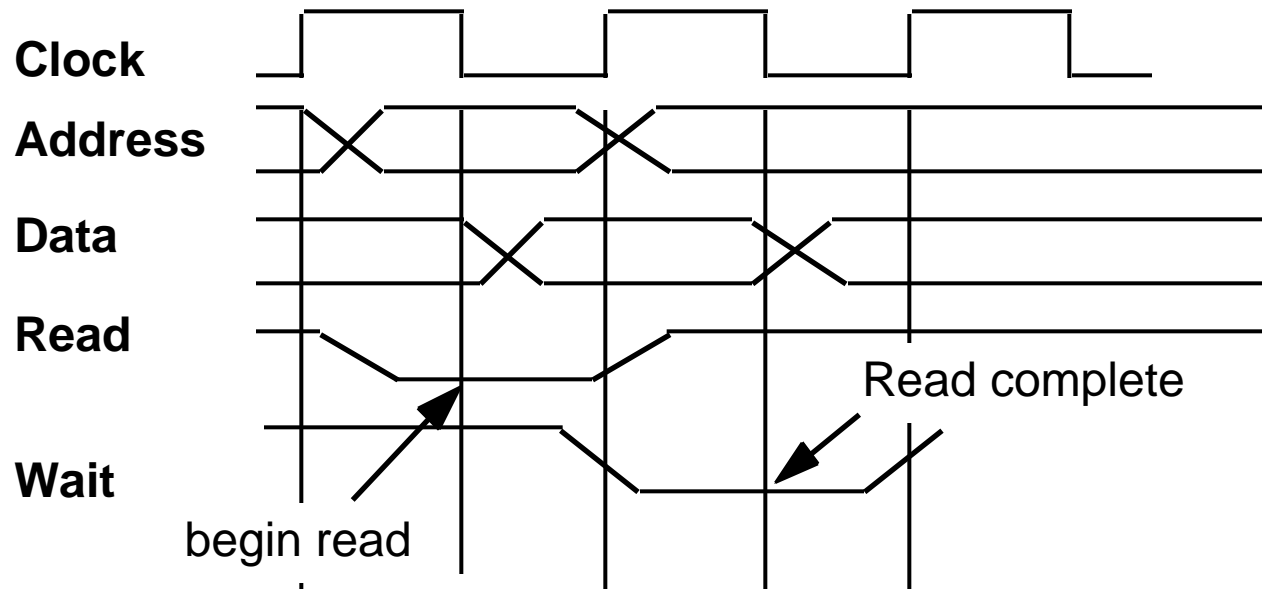
**Bus Slave:** module activated by the transaction

**Bus Communication Protocol:** specification of sequence of events and timing requirements in transferring information.

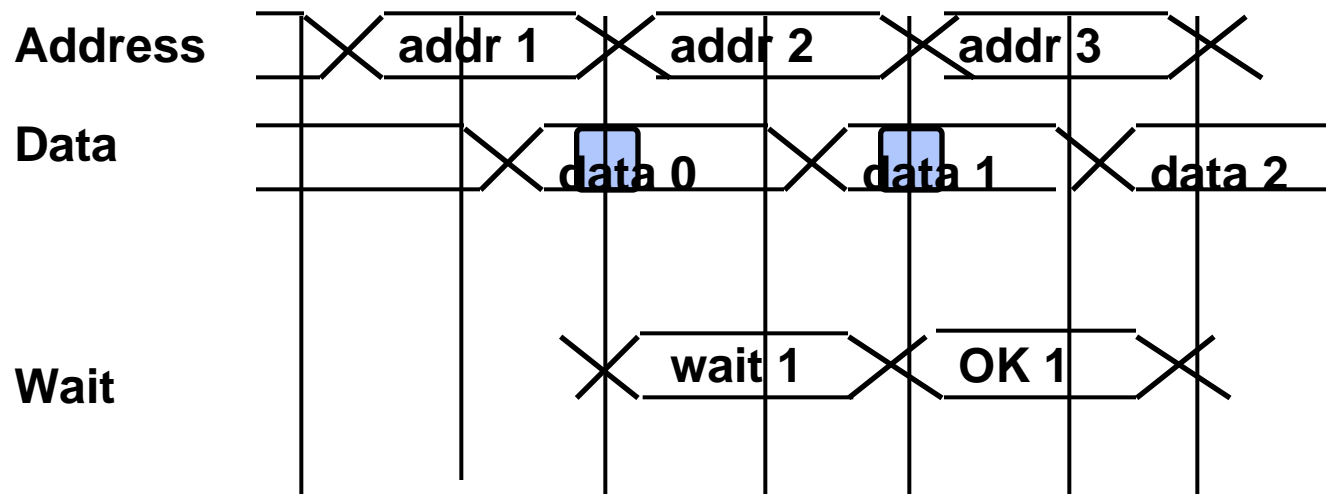
**Asynchronous Bus Transfers:** control lines (req., ack.) serve to orchestrate sequencing

**Synchronous Bus Transfers:** sequence relative to common clock

# Synchronous Bus Protocols

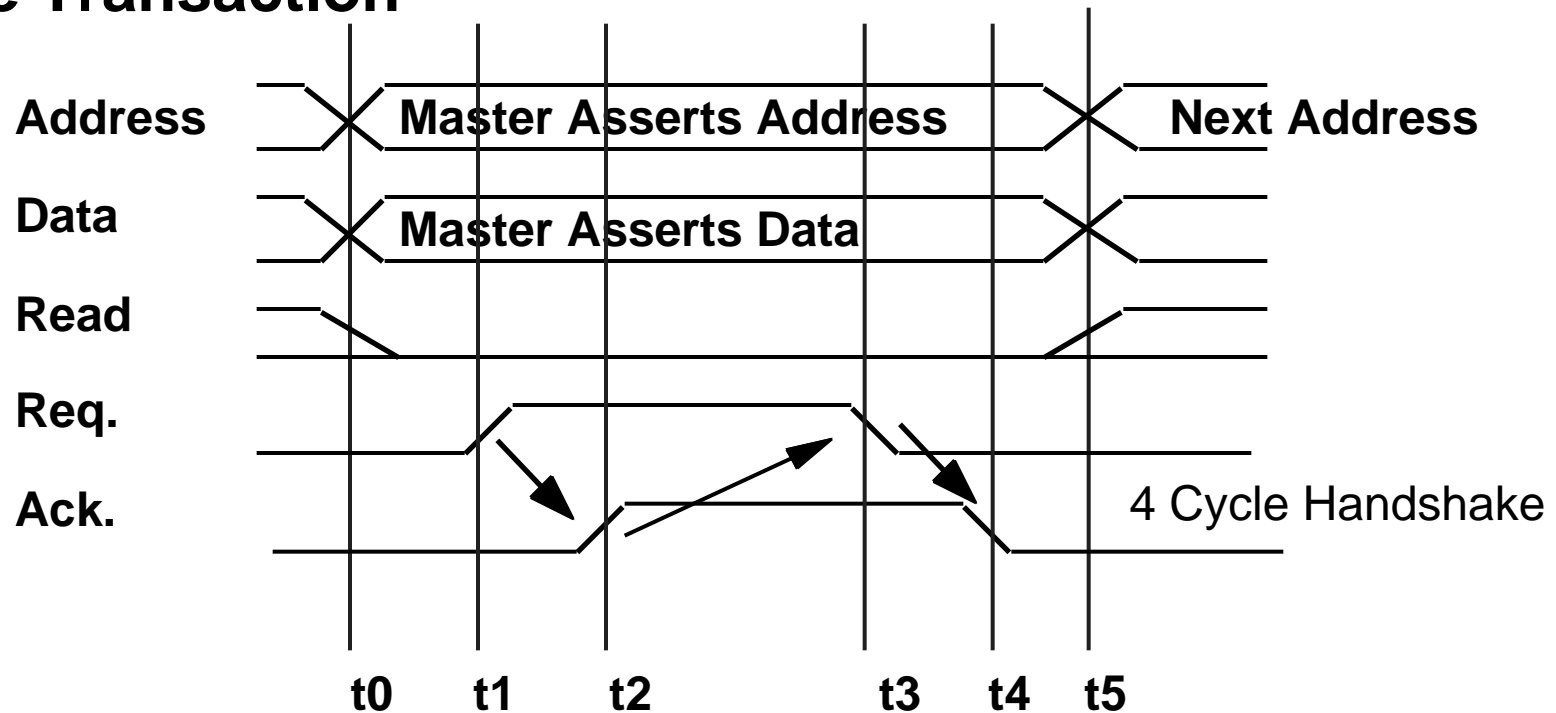


## Pipelined/Split transaction Bus Protocol



# Asynchronous Handshake

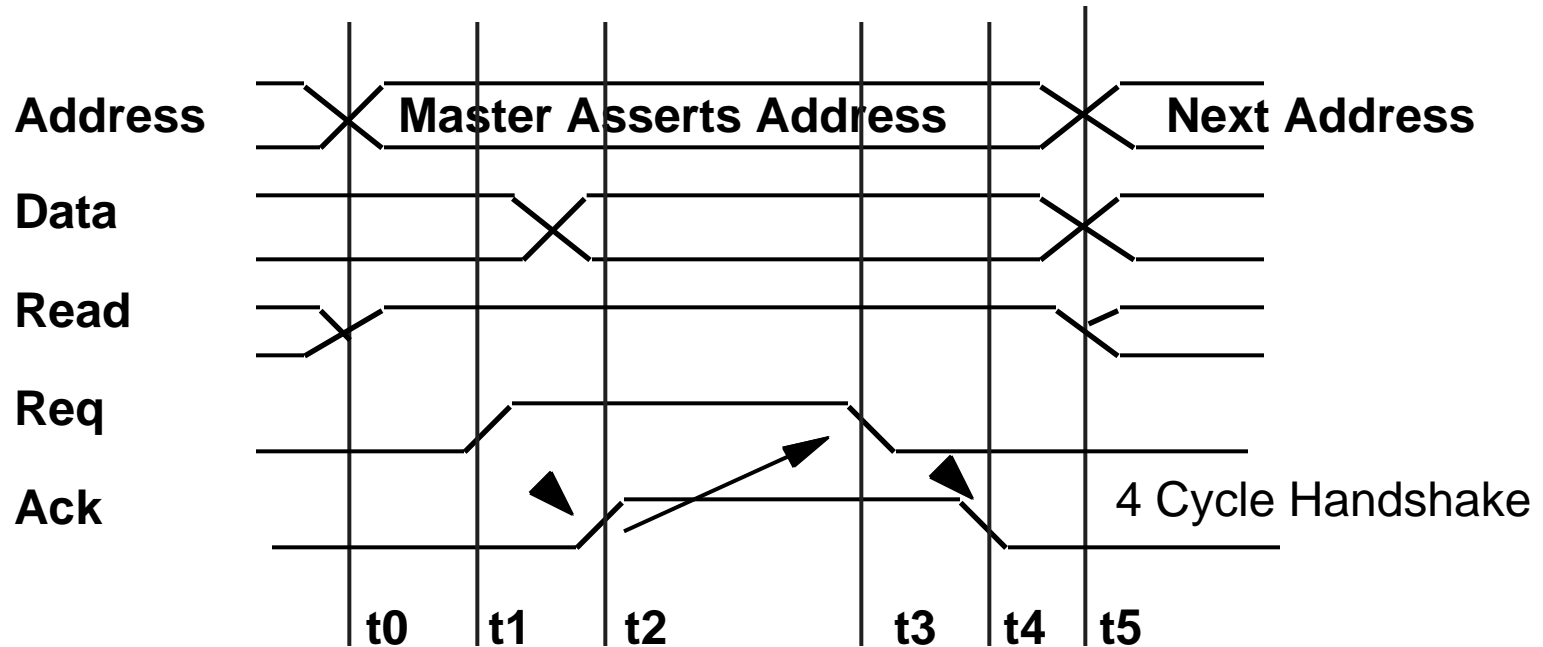
## Write Transaction



- $t_0$  :** Master has obtained control and asserts address, direction, data  
Waits a specified amount of time for slaves to decode target\
- $t_1$  :** Master asserts request line
- $t_2$  :** Slave asserts ack, indicating data received
- $t_3$  :** Master releases req
- $t_4$  :** Slave releases ack



# Read Transaction

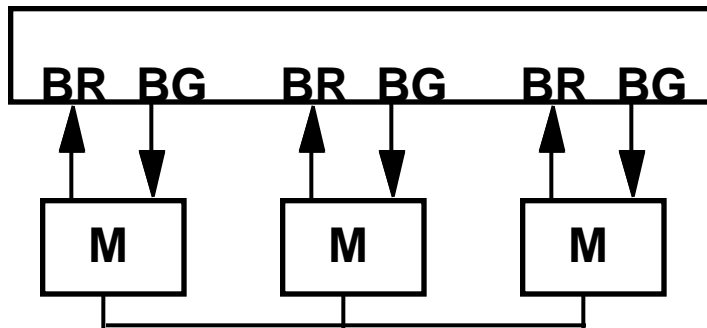


- t0 :** Master has obtained control and asserts address, direction, data  
Waits a specified amount of time for slaves to decode target\
- t1:** Master asserts request line
- t2:** Slave asserts ack, indicating ready to transmit data
- t3:** Master releases req, data received
- t4:** Slave releases ack

**Time Multiplexed Bus: address and data share lines**

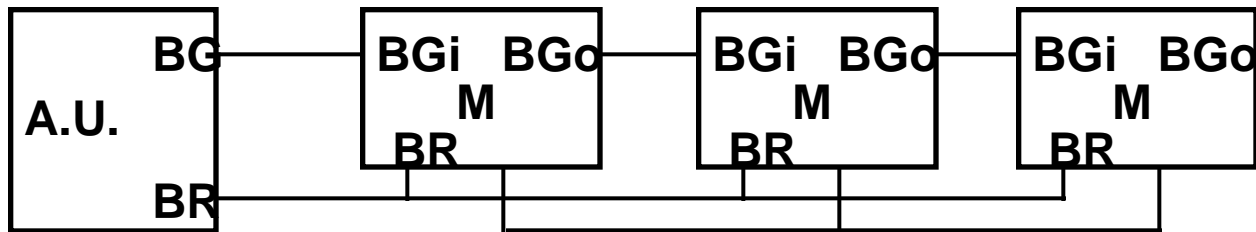
# Bus Arbitration

## Parallel (Centralized) Arbitration

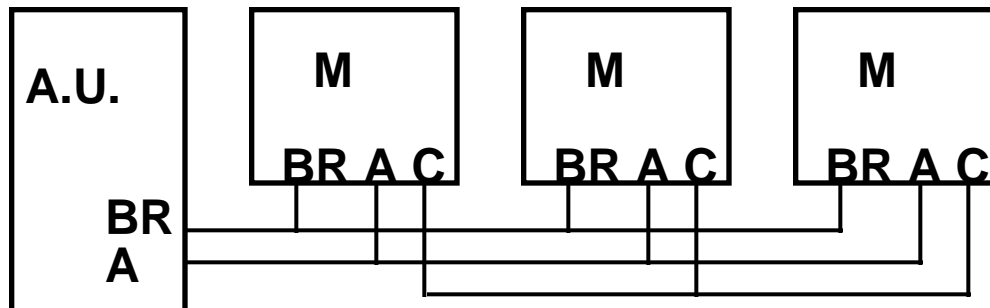


Bus Request  
Bus Grant

## Serial Arbitration (daisy chaining)



## Polling



# Bus Options

<i>Option</i>	<i>High performance</i>	<i>Low cost</i>
<b>Bus width</b>	<b>Separate address &amp; data lines</b>	<b>Multiplex address &amp; data lines</b>
<b>Data width</b>	<b>Wider is faster (e.g., 32 bits)</b>	<b>Narrower is cheaper (e.g., 8 bits)</b>
<b>Transfer size</b>	<b>Multiple words has less bus overhead</b>	<b>Single-word transfer is simpler</b>
<b>Bus masters</b>	<b>Multiple (requires arbitration)</b>	<b>Single master (no arbitration)</b>
<b>Split transaction?</b>	<b>Yes—separate Request and Reply packets gets higher bandwidth (needs multiple masters)</b>	<b>No—continuous connection is cheaper and has lower latency</b>
<b>Clocking</b>	<b>Synchronous</b>	<b>Asynchronous</b>

# 1990 Bus Survey (P&H, 1st Ed)

	VME	FutureBus	Multibus II	IPI	SCSI
Signals	128	96	96	16	8
Addr/Data mux	no	yes	yes	n/a	n/a
Data width	16 - 32	32	32	16	8
Masters	multi	multi	multi	single	multi
Clocking	Async	Async	Sync	Async	either
<b>MB/s (0ns, word)</b>	<b>25</b>	<b>37</b>	<b>20</b>	<b>25</b>	<b>1.5 (asyn)</b> <b>5 (sync)</b>
150ns word	12.9	15.5	10	=	=
<b>0ns block</b>	<b>27.9</b>	<b>95.2</b>	<b>40</b>	<b>=</b>	<b>=</b>
150ns block	13.6	20.8	13.3	=	=
Max devices	21	20	21	8	7
Max meters	0.5	0.5	0.5	50	25
Standard	IEEE 1014	IEEE 896.1	ANSI/IEEE 1296	ANSI X3.129	ANSI X3.131

# VME

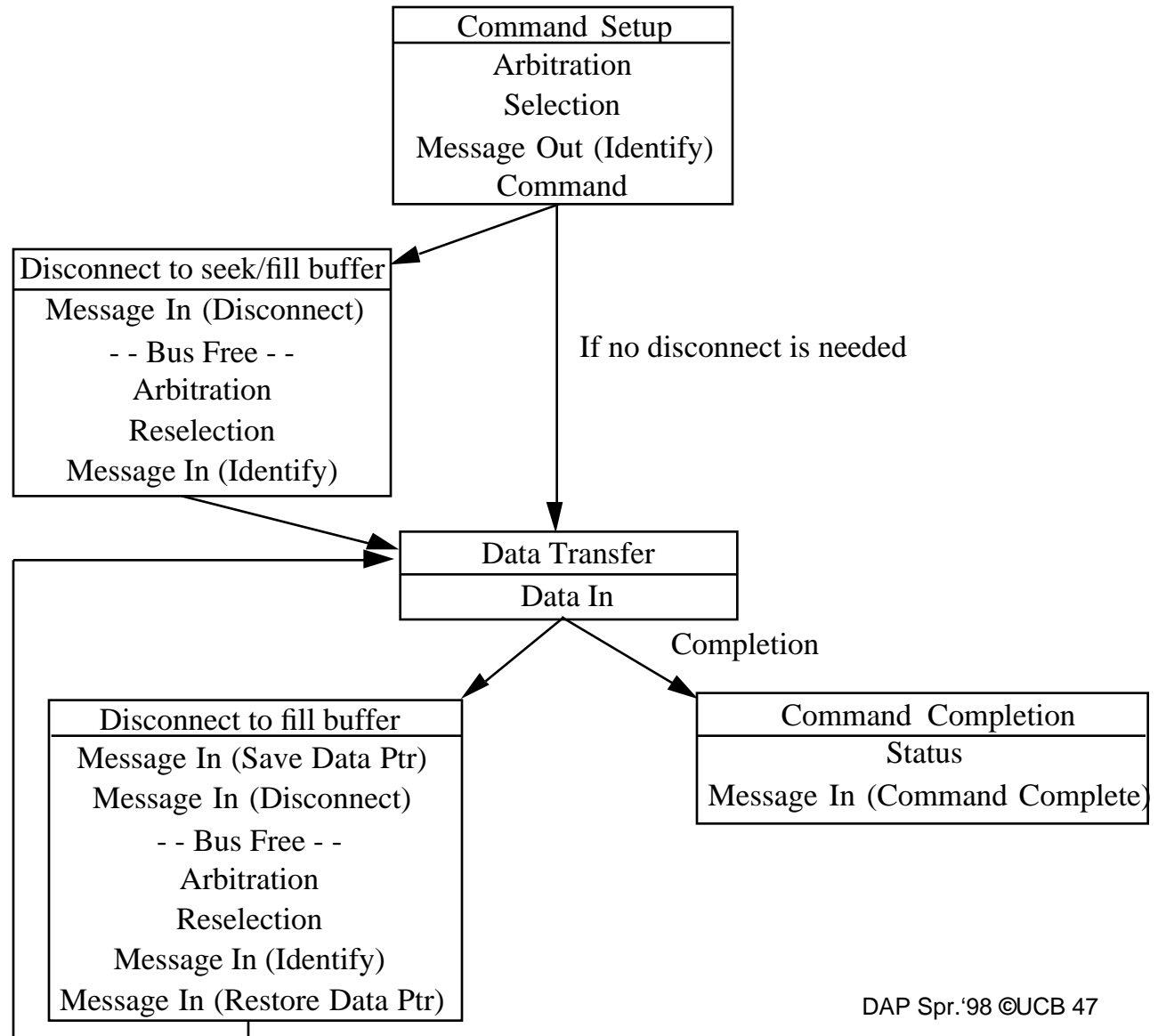
- **3 96-pin connectors**
- **128 defined as standard, rest customer defined**
  - **32 address**
  - **32 data**
  - **64 command & power/ground lines**

# SCSI: Small Computer System Interface

- Clock rate: 5 MHz / 10 MHz (fast) / 20 MHz (ultra)
- Width:  $n = 8$  bits / 16 bits (wide); up to  $n - 1$  devices to communicate on a bus or “string”
- Devices can be slave (“target”) or master (“initiator”)
- SCSI protocol: a series of “phases”, during which specific actions are taken by the controller and the SCSI disks
  - **Bus Free**: No device is currently accessing the bus
  - **Arbitration**: When the SCSI bus goes free, multiple devices may request (arbitrate for) the bus; fixed priority by address
  - **Selection**: informs the target that it will participate (**Reselection** if disconnected)
  - **Command**: the initiator reads the SCSI command bytes from host memory and sends them to the target
  - **Data Transfer**: data in or out, initiator: target
  - **Message Phase**: message in or out, initiator: target (identify, save/restore data pointer, disconnect, command complete)
  - **Status Phase**: target, just before command complete

# SCSI "Bus": Channel Architecture

peer-to-peer protocols  
 initiator/target  
 linear byte streams  
 disconnect/reconnect



# 1993 I/O Bus Survey (P&H, 2nd Ed)

Bus	SBus	TurboChannel	MicroChannel	PCI
Originator	Sun	DEC	IBM	Intel
Clock Rate (MHz)	16-25	12.5-25	async	33
Addressing	Virtual	Physical	Physical	Physical
Data Sizes (bits)	8,16,32	8,16,24,32	8,16,24,32,64	8,16,24,32,64
Master	Multi	Single	Multi	Multi
Arbitration	Central	Central	Central	Central
32 bit read (MB/s)	33	25	20	33
<b>Peak (MB/s)</b>	<b>89</b>	<b>84</b>	<b>75</b>	<b>111 (222)</b>
Max Power (W)	16	26	13	25

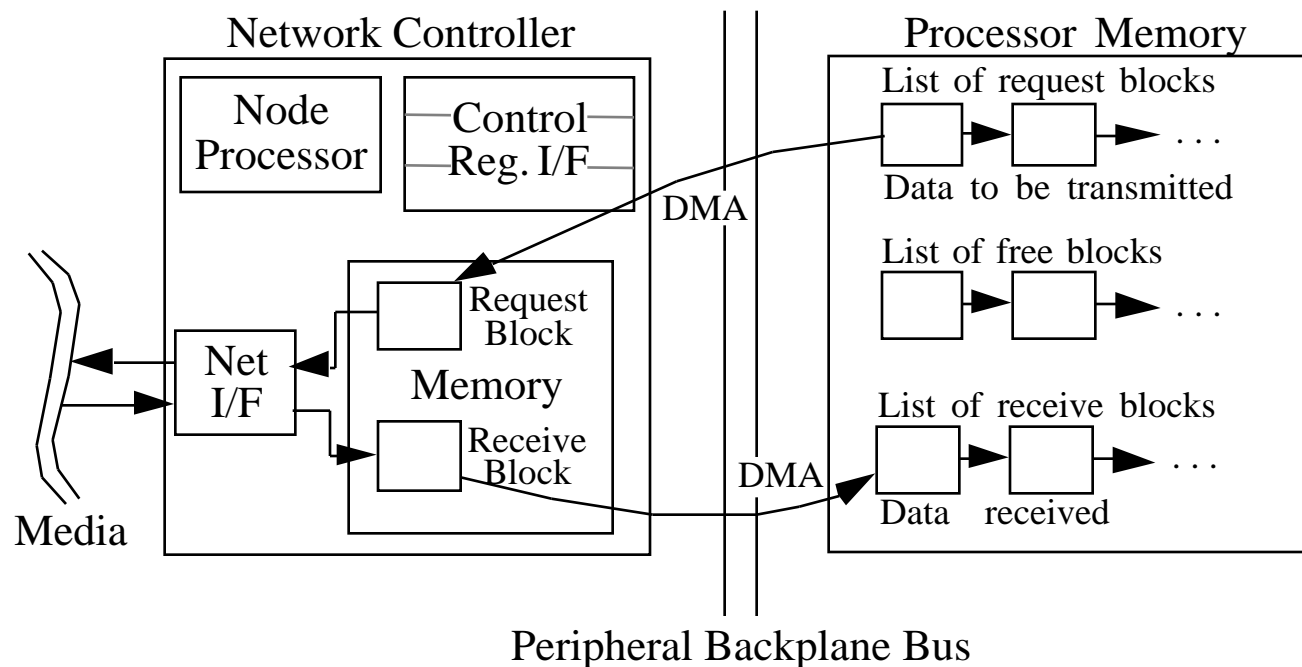


# 1993 MP Server Memory Bus Survey

Bus	Summit	Challenge	XDBus
Originator	HP	SGI	Sun
Clock Rate (MHz)	60	48	66
Split transaction?	Yes	Yes	Yes?
Address lines	48	40	??
Data lines	128	256	144 (parity)
Data Sizes (bits)	512	1024	512
Clocks/transfer	4	5	4?
<b>Peak (MB/s)</b>	<b>960</b>	<b>1200</b>	<b>1056</b>
Master	Multi	Multi	Multi
Arbitration	Central	Central	Central
Addressing	Physical	Physical	Physical
Slots	16	9	10
Busses/system	1	1	2
Length	13 inches	12? inches	17 inches

# Communications Networks

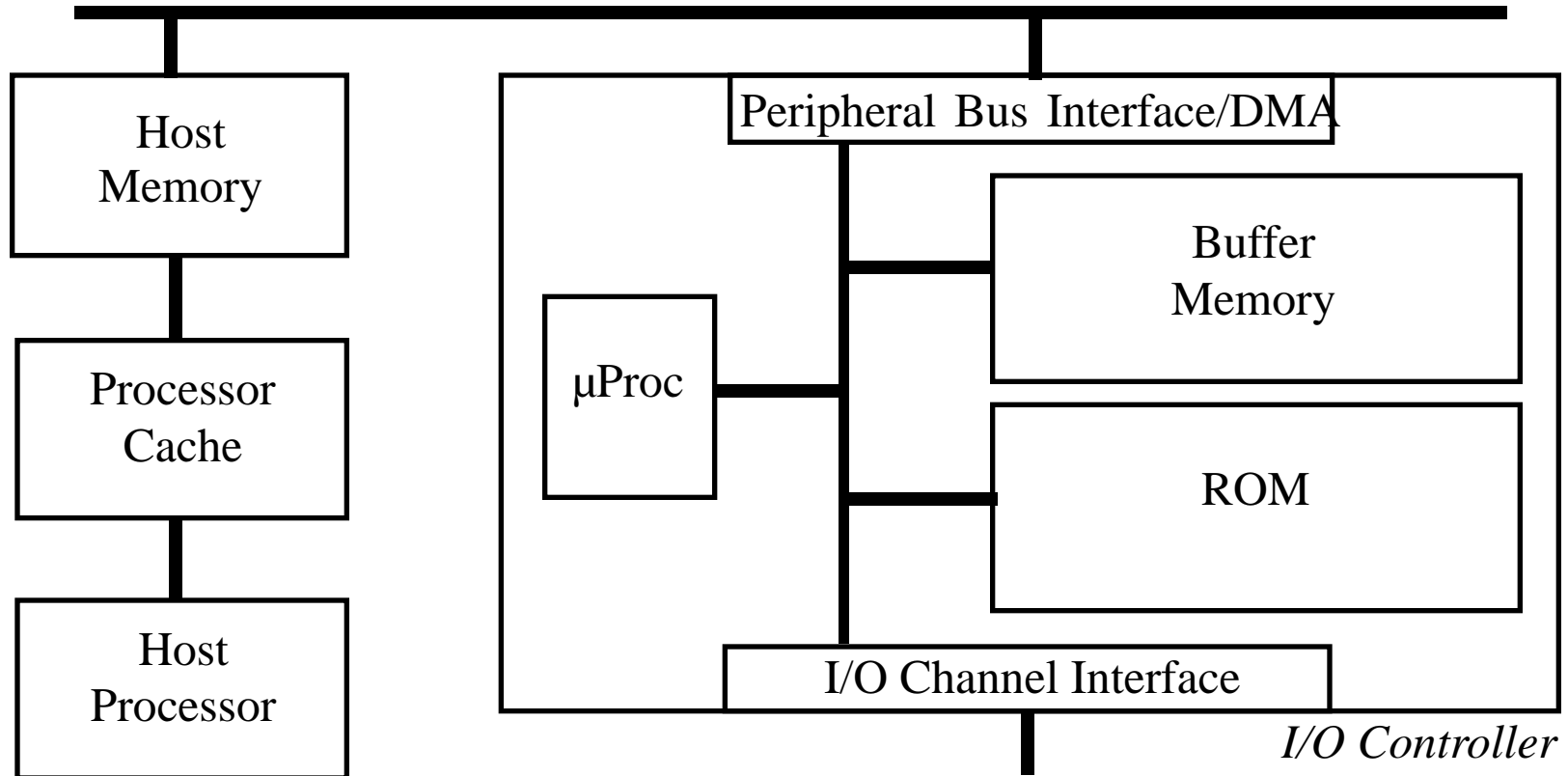
Performance limiter is memory system, OS overhead



- **Send/receive queues in processor memories**
- **Network controller copies back and forth via DMA**
- **No host intervention needed**
- **Interrupt host when message sent or received**

# I/O Controller Architecture

Peripheral Bus (VME, FutureBus, etc.)

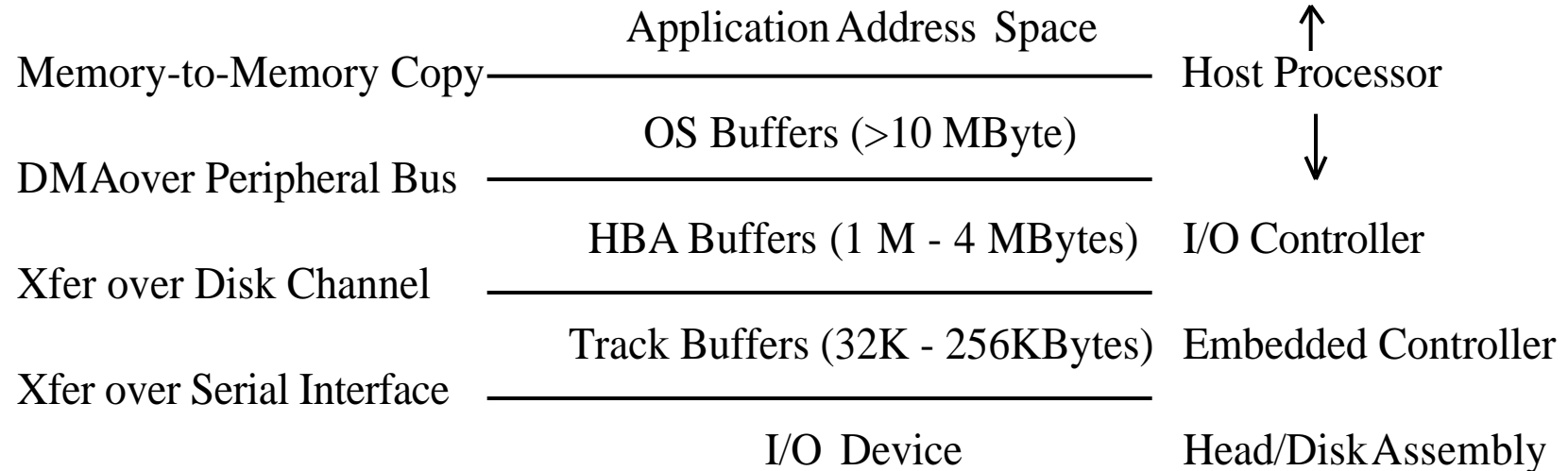


**Request/response block interface**

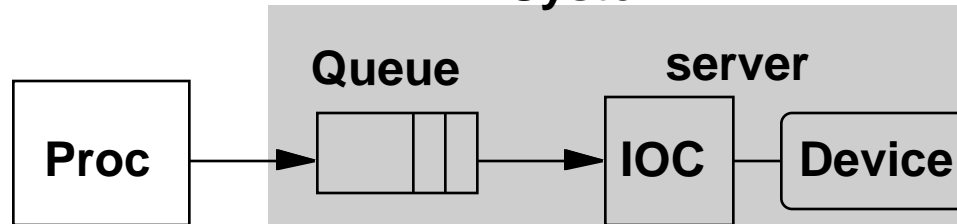
**Backdoor access to host memory**

# I/O Data Flow

**Impediment to high performance: multiple copies, complex hierarchy**



# Summary: A Little Queuing Theory



- Queuing models assume state of equilibrium: input rate = output rate

- Notation:

$r$  average number of arriving customers/second

$T_{ser}$  average time to service a customer (traditionally  $\mu = 1/ T_{ser}$ )

$u$  server utilization (0..1):  $u = r \times T_{ser}$

$T_q$  average time/customer in queue

$T_{sys}$  average time/customer in system:  $T_{sys} = T_q + T_{ser}$

$L_q$  average length of queue:  $L_q = r \times T_q$

$L_{sys}$  average length of system :  $L_{sys} = r \times T_{sys}$

- Little's Law: **Length<sub>system</sub> = rate x Time<sub>system</sub>**  
(Mean number customers = arrival rate x mean service time)

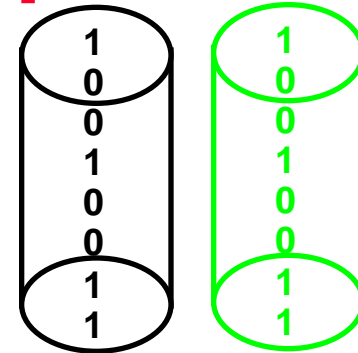
# Summary: Redundant Arrays of Disks (RAID) Techniques

- **Disk Mirroring, Shadowing (RAID 1)**

Each disk is fully duplicated onto its "shadow"

Logical write = two physical writes

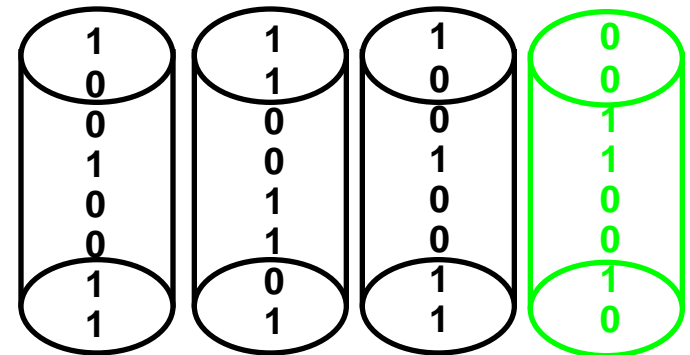
100% capacity overhead



- **Parity Data Bandwidth Array (RAID 3)**

Parity computed horizontally

Logically a single high data bw disk



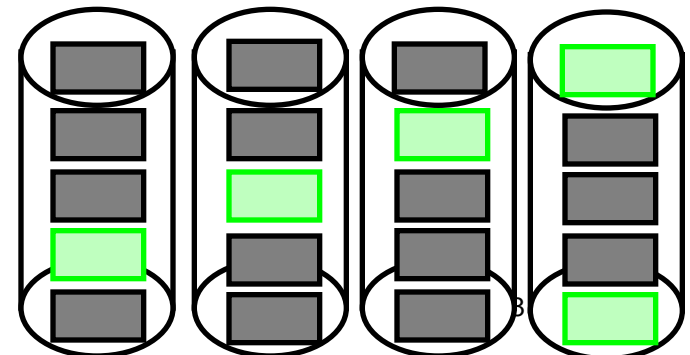
- **High I/O Rate Parity Array (RAID 5)**

Interleaved parity blocks

Independent reads and writes

Logical write = 2 reads + 2 writes

Parity + Reed-Solomon codes



# Review: Storage System Issues

- *Historical Context of Storage I/O*
- *Secondary and Tertiary Storage Devices*
- *Storage I/O Performance Measures*
- *Processor Interface Issues*
- **A Little Queuing Theory**
- **Redundant Arrays of Inexpensive Disks (RAID)**
- **I/O Buses**
- **ABCs of UNIX File Systems**
- **I/O Benchmarks**
- **Comparing UNIX File System Performance**