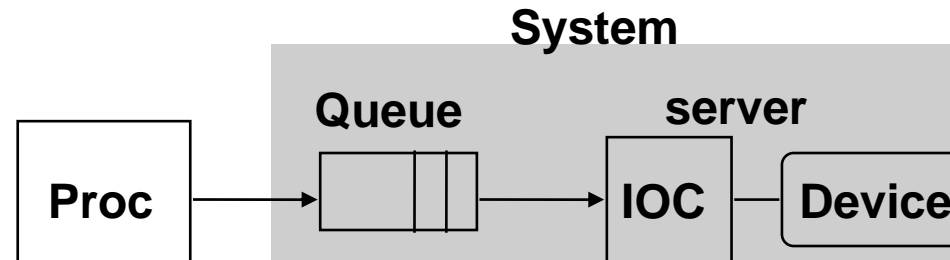


CS252
Graduate Computer Architecture

Lecture 8:
Network 1: Definitions, Metrics, 252
Projects

February 9, 2001
Prof. David A. Patterson
Computer Science 252
Spring 2001

Review: 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: $Length_{system} = rate \times Time_{system}$
(Mean number customers = arrival rate x mean service time)

Review: I/O Benchmarks

- Scaling to track technological change
- TPC: price performance as normalizing configuration feature
- Auditing to ensure no foul play
- Throughput with restricted response time is normal measure
- Benchmarks to measure Availability, Maintainability?

Review: Availability benchmarks

- Availability benchmarks can provide valuable insight into availability behavior of systems
 - reveal undocumented availability policies
 - illustrate impact of specific faults on system behavior
- Methodology is best for *understanding* the availability behavior of a system
 - extensions are needed to distill results for automated system comparison
- A good fault-injection environment is critical
 - need realistic, reproducible, controlled faults
 - system designers should consider building in hooks for fault-injection and availability testing
- Measuring and understanding availability will be crucial in building systems that meet the needs of modern server applications
 - this benchmarking methodology is just 1st step towards goal

Networks

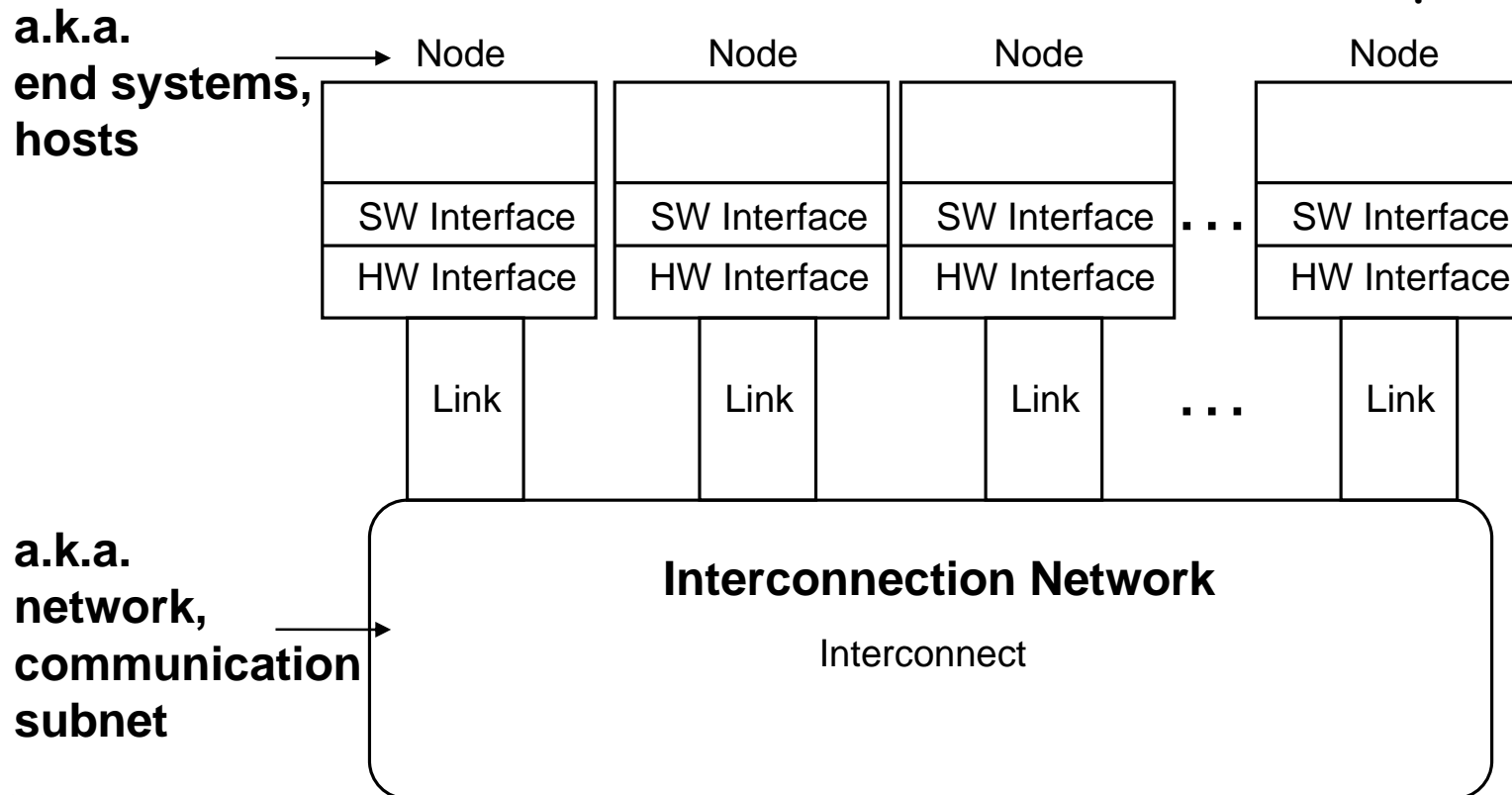
- **Goal:** Communication between computers
- **Eventual Goal:** treat collection of computers as if one big computer, distributed resource sharing
- **Theme:** Different computers must agree on many things
 - Overriding importance of standards and protocols
 - Error tolerance critical as well
- **Warning:** Terminology-rich environment

Networks

- Facets people talk a lot about:
 - direct (point-to-point) vs. indirect (multi-hop)
 - topology (e.g., bus, ring, DAG)
 - routing algorithms
 - switching (aka multiplexing)
 - wiring (e.g., choice of media, copper, coax, fiber)
- What really matters:
 - latency
 - bandwidth
 - cost
 - reliability

Interconnections (Networks)

- Examples (see Figure 7.19, page 633):
 - **Wide Area Network (ATM)**: 100-1000s nodes; ~ 5,000 kilometers
 - **Local Area Networks (Ethernet)**: 10-1000 nodes; ~ 1-2 kilometers
 - **System/Storage Area Networks (FC-AL)**: 10-100s nodes;
~ 0.025 to 0.1 kilometers per link



SAN: Storage vs. System

- **Storage Area Network (SAN):** A block I/O oriented network between application servers and storage
 - Fibre Channel is an example
- Usually high bandwidth requirements, and less concerned about latency
 - in 2001: 1 Gbit bandwidth and millisecond latency OK
- Commonly a dedicated network (that is, not connected to another network)
- May need to work gracefully when saturated
- Given larger block size, may have higher bit error rate (BER) requirement than LAN

SAN: Storage vs. System

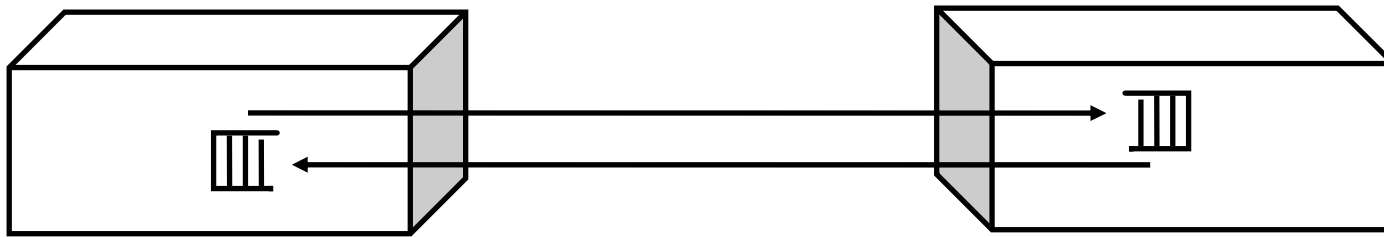
- **System Area Network (SAN):** A network aimed at connecting computers
 - Myrinet is an example
- **Aimed at High Bandwidth AND Low Latency.**
 - in 2001: > 1 Gbit bandwidth and ~ 10 microsecond
- **May offer in order delivery of packets**
- **Given larger block size, may have higher bit error rate (BER) requirement than LAN**

More Network Background

- Connection of 2 or more networks:
Internetworking
- 3 cultures for 3 classes of networks
 - WAN: telecommunications, Internet
 - LAN: PC, workstations, servers cost
 - SAN: Clusters, RAID boxes: latency (System A.N.) or bandwidth (Storage A.N.)
- Try for single terminology
- Motivate the interconnection complexity incrementally

ABCs of Networks

- **Starting Point**: Send bits between 2 computers



- Queue (FIFO) on each end
- Information sent called a “**message**”
- Can send both ways (“**Full Duplex**”)
- Rules for communication? “**protocol**”
 - Inside a computer:
 - » Loads/Stores: Request (Address) & Response (Data)
 - » Need Request & Response signaling

A Simple Example

- What is the format of message?
 - Fixed? Number bytes?

Request/
Response

Address/Data



1 bit

32 bits

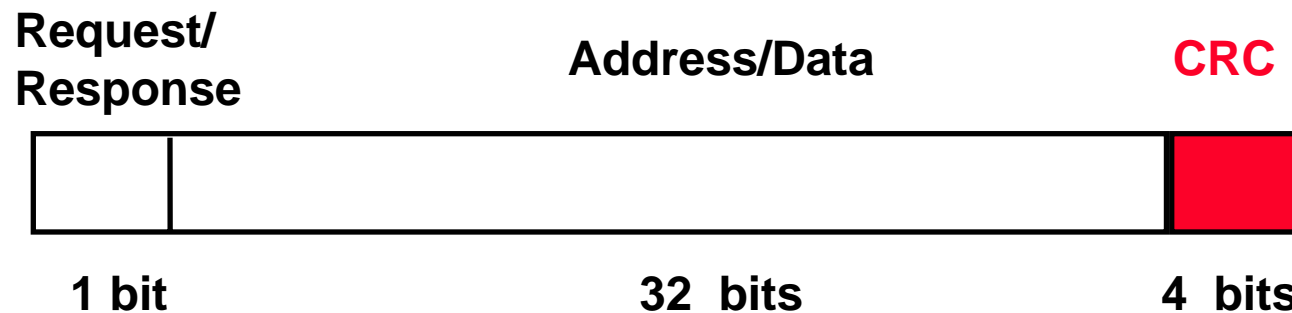
- 0: Please send data from Address
- 1: Packet contains data corresponding to request
- **Header/Trailer**: information to deliver a message
- **Payload**: data in message (1 word above)

Questions About Simple Example

- What if more than 2 computers want to communicate?
 - Need computer “**address field**” (destination) in packet
- What if packet is garbled in transit?
 - Add “**error detection field**” in packet (e.g., Cyclic Redundancy Chk)
- What if packet is lost?
 - More “**elaborate protocols**” to detect loss (e.g., NAK, ARQ, time outs)
- What if multiple processes/machine?
 - Queue per process to provide protection
- Simple questions such as these lead to more complex protocols and packet formats => complexity

A Simple Example Revisted

- What is the format of packet?
 - Fixed? Number bytes?



00: Request—Please send data from Address

01: Reply—Packet contains data corresponding to request

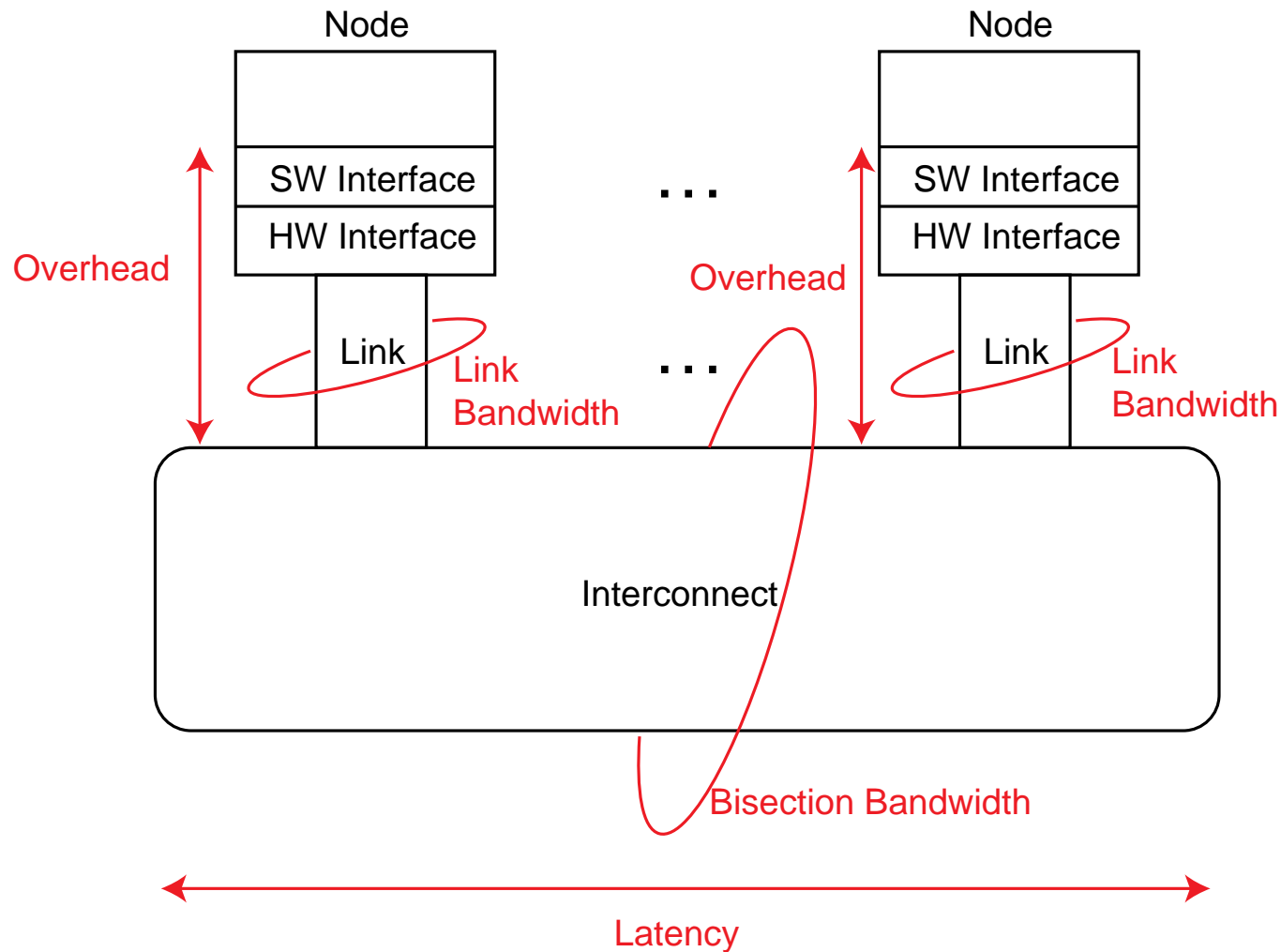
10: Acknowledge request

11: Acknowledge reply

Software to Send and Receive

- **SW Send steps**
 - 1: Application copies data to OS buffer
 - 2: OS calculates checksum, starts timer
 - 3: OS sends data to network interface HW and says start
- **SW Receive steps**
 - 3: OS copies data from network interface HW to OS buffer
 - 2: OS calculates checksum, if matches send ACK; if not, *deletes message* (sender resends when timer expires)
 - 1: If OK, OS copies data to user address space and signals application to continue
- **Sequence of steps for SW: protocol**
 - Example similar to UDP/IP protocol in UNIX

Network Performance Measures

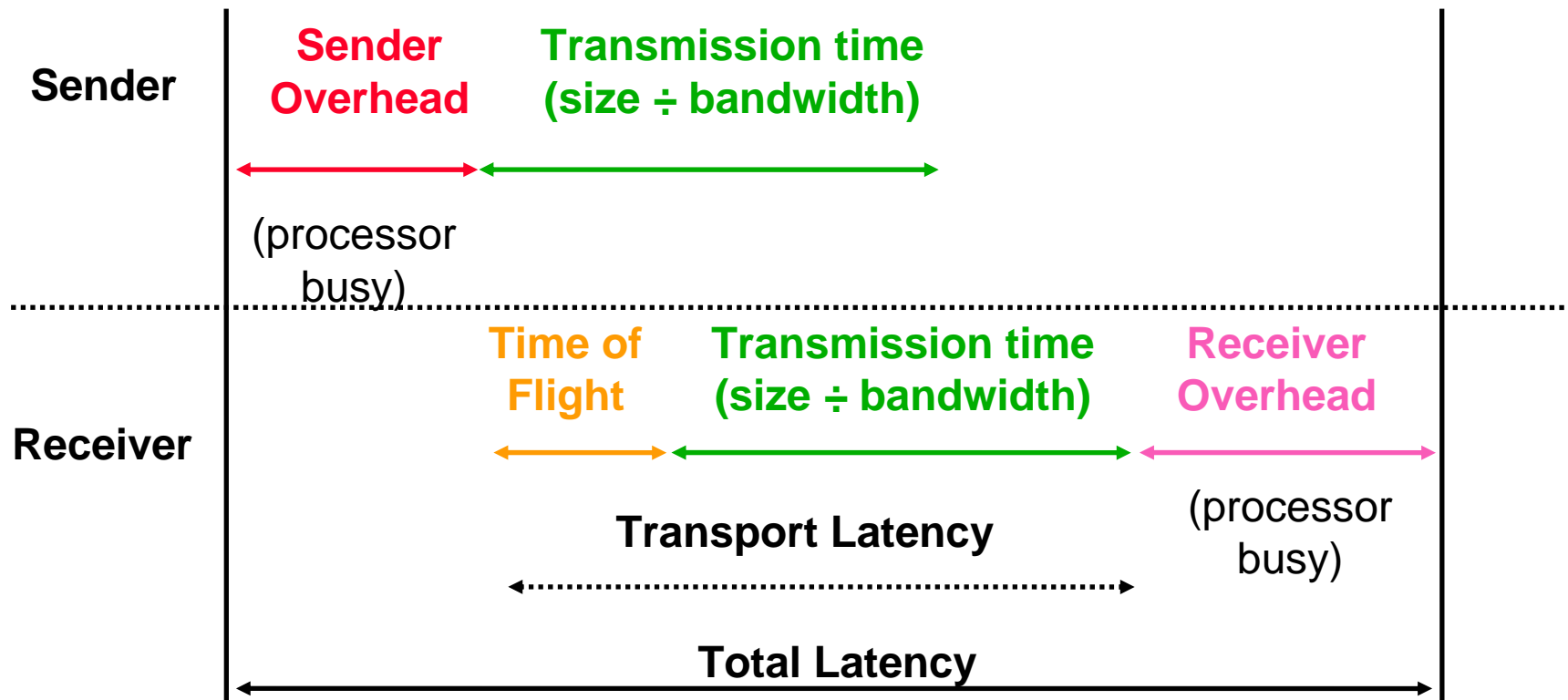


- **Overhead**: latency of interface vs. **Latency**: network

CS 252 Administtrivia

- HW #1 due Saturday: send electronically to TA
- Pick a partner, project by Monday; send electronically to me, TA
- I'll be available Monday afternoon to talk

Universal Performance Metrics



$$\text{Total Latency} = \text{Sender Overhead} + \text{Time of Flight} + \text{Message Size} \div \text{BW} + \text{Receiver Overhead}$$

Includes header/trailer in BW calculation?

Total Latency Example

- 1000 Mbit/sec., sending overhead of 80 μ sec & receiving overhead of 100 μ sec.
- a 10000 byte message (including the header), allows 10000 bytes in a single message
- 3 situations: distance 1000 km v. 0.5 km v. 0.01
- Speed of light \sim 300,000 km/sec (1/2 in media)
- Latency_{0.01km} =
- Latency_{0.01km} =
- Latency_{1000km} =

Total Latency Example

- 1000 Mbit/sec., sending overhead of 80 μ sec & receiving overhead of 100 μ sec.
- a 10000 byte message (including the header), allows 10000 bytes in a single message
- 2 situations: distance 100 m vs. 1000 km
- Speed of light \sim 300,000 km/sec
- $\text{Latency}_{0.01\text{km}} = 80 + 0.01\text{km} / (50\% \times 300,000) + 10000 \times 8 / 1000 + 100 = 260 \mu\text{sec}$
- $\text{Latency}_{0.5\text{km}} = 80 + 0.5\text{km} / (50\% \times 300,000) + 10000 \times 8 / 1000 + 100 = 263 \mu\text{sec}$
- $\text{Latency}_{1000\text{km}} = 80 + 1000 \text{ km} / (50\% \times 300,000) + 10000 \times 8 / 1000 + 100 = 6931$
- Long time of flight \Rightarrow complex WAN protocol

Universal Metrics

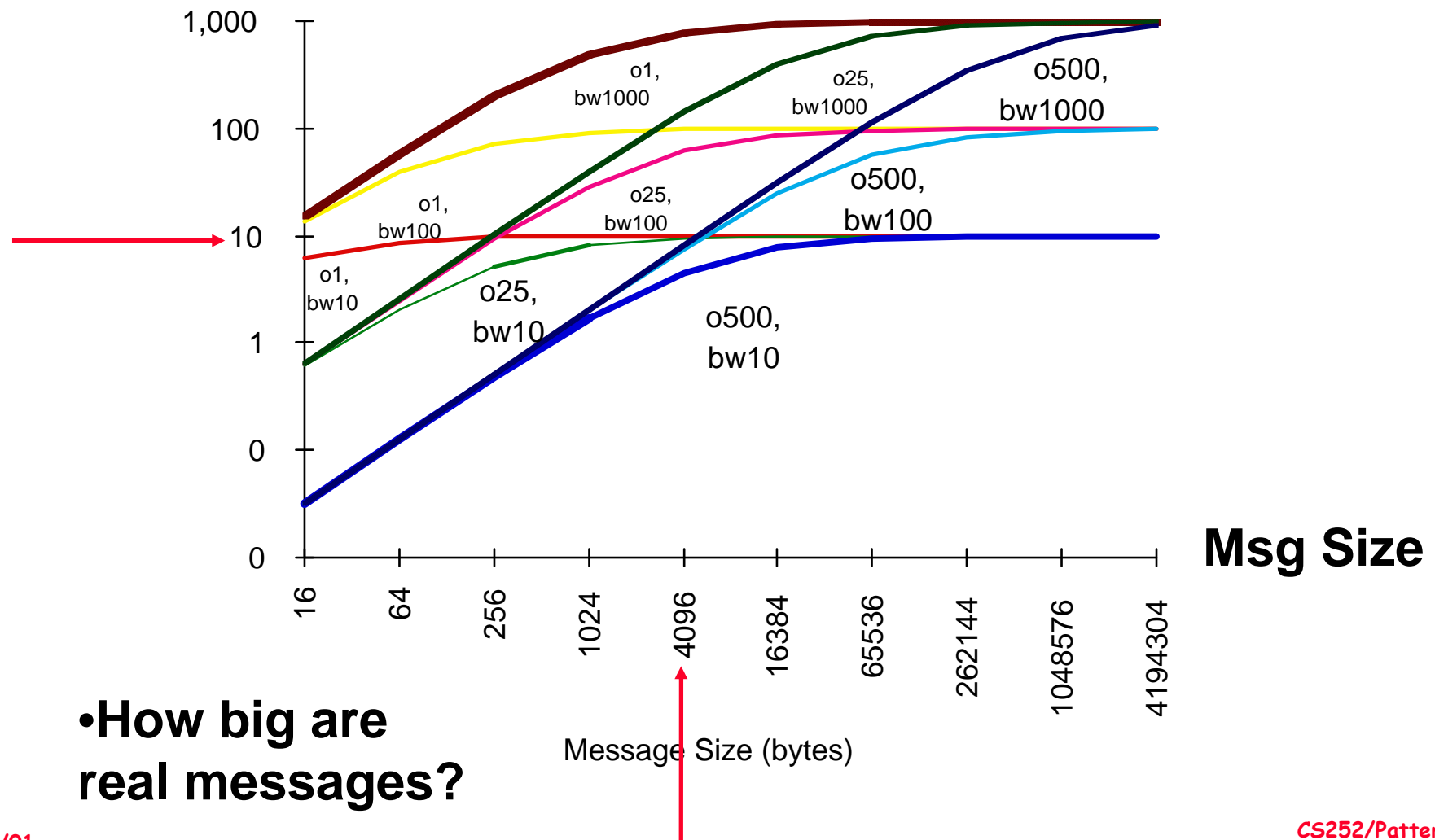
- Apply recursively to all levels of system
- inside a chip, between chips on a board, between computers in a cluster, ...
- Look at WAN v. LAN v. SAN

Simplified Latency Model

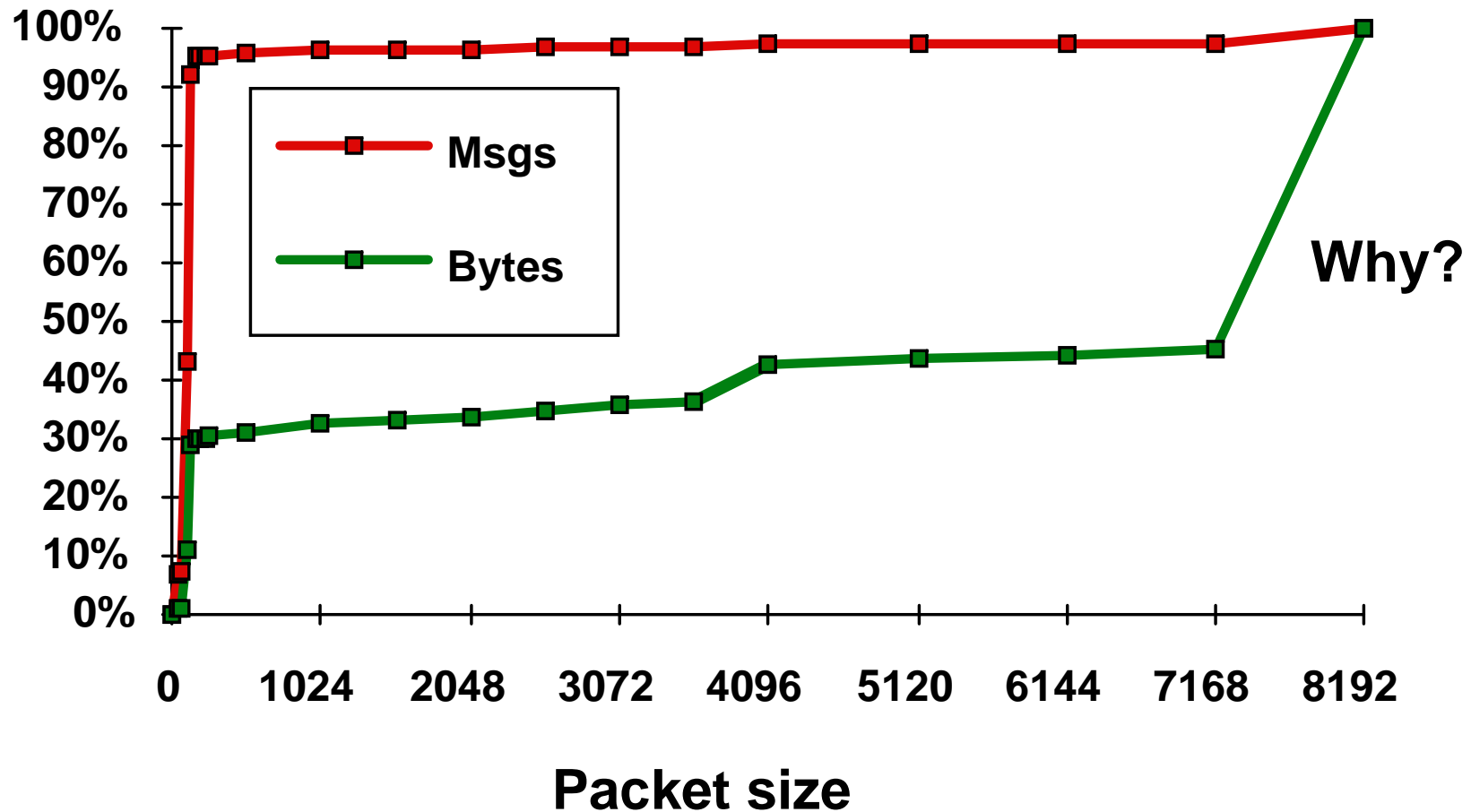
- Total Latency - **Overhead** + Message Size / BW
- **Overhead** = Sender Overhead + Time of Flight + Receiver Overhead
- Example: show what happens as vary
 - Overhead: 1, 25, 500 μ sec
 - BW: 10, 100, 1000 Mbit/sec (factors of 10)
 - Message Size: 16 Bytes to 4 MB (factors of 4)
- If overhead 500 μ sec,
how big a message > 10 Mb/s?

Overhead, BW, Size

Delivered BW

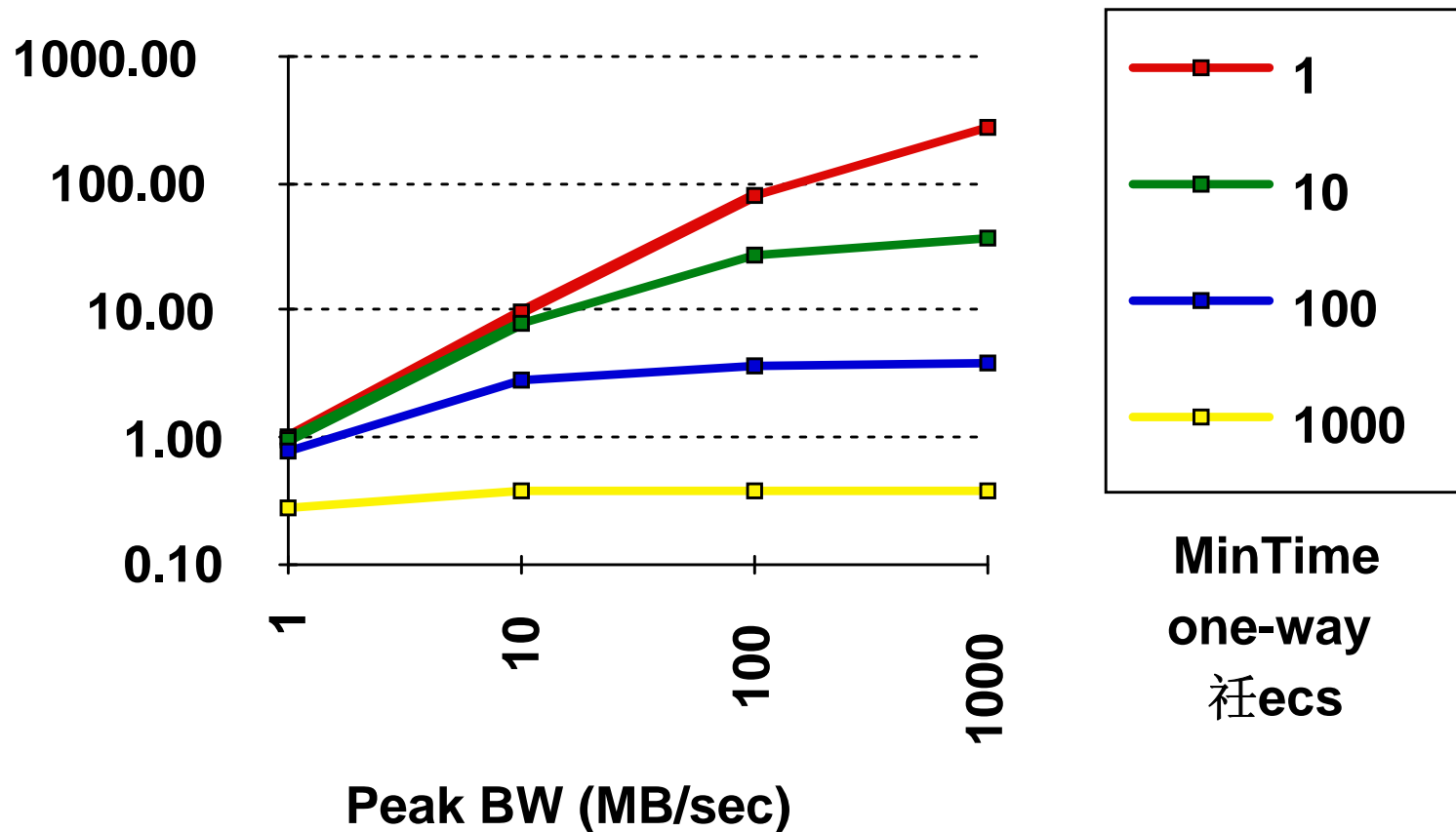


Measurement: Sizes of Message for NFS



- 95% Msgs, 30% bytes for packets ~ 200 bytes
- > 50% data transferred in packets = 8KB

Impact of Overhead on Delivered BW



- BW model: $\text{Time} = \text{overhead} + \text{msg size} / \text{peak BW}$

Interconnect Issues

- Performance Measures
- Network Media

Network Media

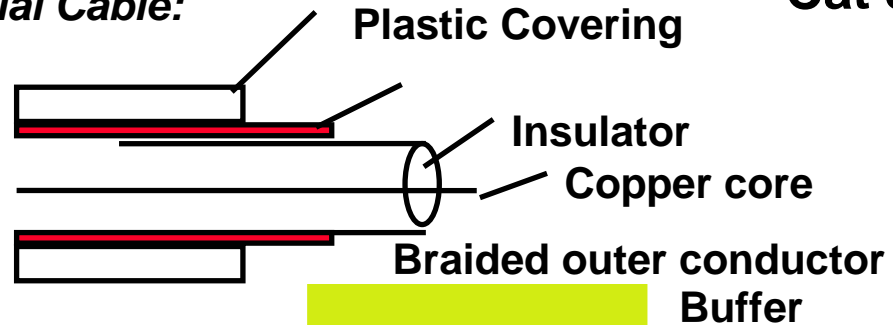
Twisted Pair:



Copper, 1mm thick, twisted to avoid antenna effect (telephone)

"Cat 5" is 4 twisted pairs in bundle

Coaxial Cable:

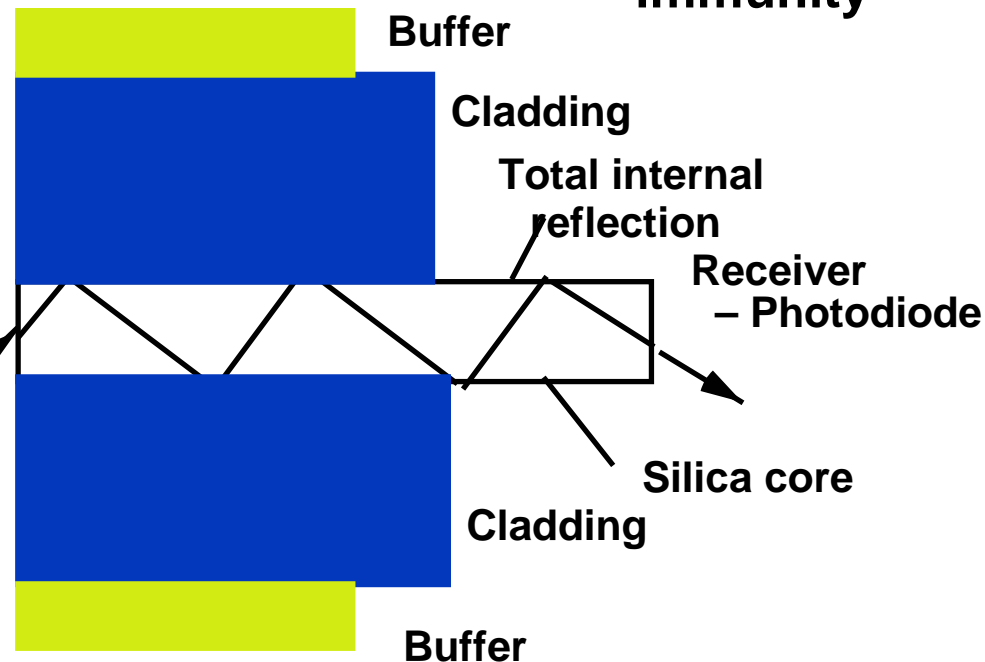


Used by cable companies:
high BW, good noise immunity

Fiber Optics

Transmitter
– L.E.D
– Laser Diode

light source



Light: 3 parts
are cable, light
source, light
detector.

Note fiber is
unidirectional;
need 2 for full
duplex

Fiber

- **Multimode fiber**: ~ 62.5 micron diameter vs. the 1.3 micron wavelength of infrared light. Since wider it has more dispersion problems, limiting its length at 1000 Mbits/s for 0.1 km, and 1-3 km at 100 Mbits/s. Uses LED as light
- **Single mode fiber**: "single wavelength" fiber (8-9 microns) uses laser diodes, 1-5 Gbits/s for 100s kms
 - Less reliable and more expensive, and restrictions on bending
 - Cost, bandwidth, and distance of single-mode fiber affected by power of the light source, the sensitivity of the light detector, and the attenuation rate (loss of optical signal strength as light passes through the fiber) per kilometer of the fiber cable.
 - Typically glass fiber, since has better characteristics than the less expensive plastic fiber

Wave Division Multiplexing Fiber

- Send N independent streams on single fiber!
- Just use different wavelengths to send and demultiplex at receiver
- WDM in 2000: 40 Gbit/s using 8 wavelengths
- Plan to go to 80 wavelengths \Rightarrow 400 Gbit/s!
- A figure of merit: $\text{BW} \times \text{max distance}$ (Gbit-km/sec)
- 10X/4 years, or 1.8X per year

Compare Media

- Assume 40 2.5" disks, each 25 GB, Move 1 km
- Compare Cat 5 (100 Mbit/s), Multimode fiber (1000 Mbit/s), single mode (2500 Mbit/s), and car
- Cat 5: $1000 \times 1024 \times 8 \text{ Mb} / 100 \text{ Mb/s} = 23 \text{ hrs}$
- MM: $1000 \times 1024 \times 8 \text{ Mb} / 1000 \text{ Mb/s} = 2.3 \text{ hrs}$
- SM: $1000 \times 1024 \times 8 \text{ Mb} / 2500 \text{ Mb/s} = 0.9 \text{ hrs}$
- Car: $5 \text{ min} + 1 \text{ km} / 50 \text{ kph} + 10 \text{ min} = 0.25 \text{ hrs}$
- Car of disks = high BW media

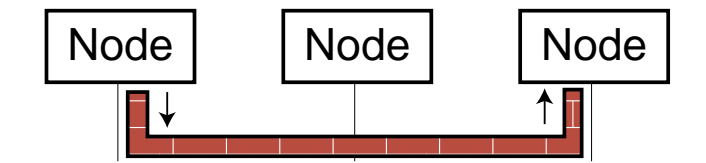
Interconnect Issues

- Performance Measures
- Network Media
- Connecting Multiple Computers

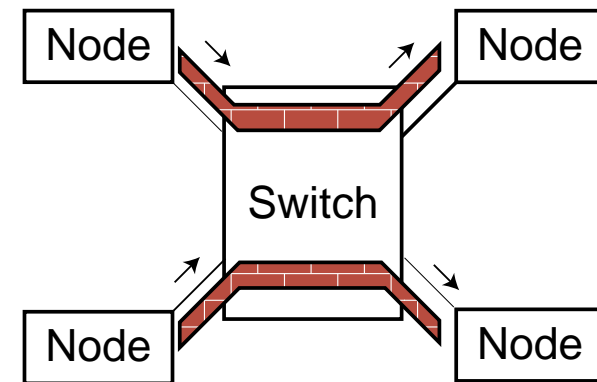
Connecting Multiple Computers

- Shared Media vs. Switched:
pairs communicate at same time:
“**point-to-point**” connections
- Aggregate BW in switched network is many times shared
 - point-to-point faster since no arbitration, simpler interface
- Arbitration in Shared network?
 - Central arbiter for LAN?
 - Listen to check if being used (“**Carrier Sensing**”)
 - Listen to check if collision (“**Collision Detection**”)
 - Random resend to avoid repeated collisions; not fair arbitration;
 - OK if low utilization

Shared Media (Ethernet)



Switched Media (CM-5, ATM)



(A. K. A. data switching interchanges, multistage interconnection networks, interface message processors)

Summary: Interconnections

- Communication between computers
- Packets for standards, protocols to cover normal and abnormal events
- Performance issues: HW & SW overhead, interconnect latency, bisection BW
- Media sets cost, distance
- Shared vs. Switched Media determines BW

Projects

- See www.cs/~pattrsn/252S01/suggestions.html

If time permits

- Discuss Hennessy paper. "The future of systems research." Computer, vol.32, (no.8), IEEE Comput. Soc, Aug. 1999
- Microprocessor Performance via ILP Analogy?
- What is key metric if services via servers is killer app?
- What is new focus for PostPC Era?
- How does he define availability vs. textbook?