

Introduction to Computer Networks

Computer Networks: Performance Analysis

<https://pages.cs.wisc.edu/~mgliu/CS640/S26/index.html>

Ming Liu

mgliu@cs.wisc.edu

Outline

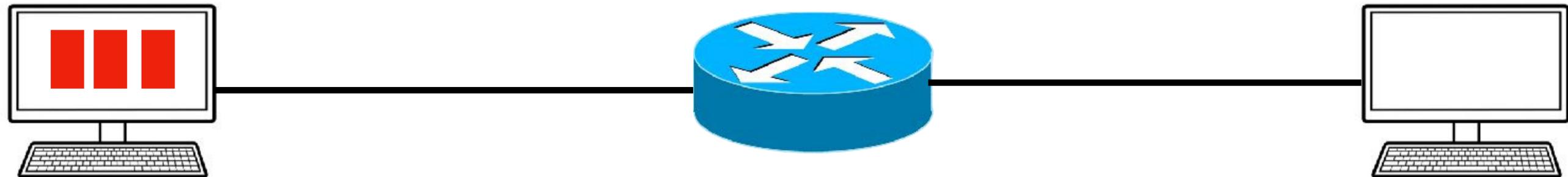
- Last
 - Computer networks: hardware infrastructure
 - Computer networks: software system
- Today
 - Delay
 - Throughput
 - RTT and BDP
- Announcements
 - Lab1 will be released on Jan. 29

Recap

- Key Questions
 - What are communication links?
 - What is an abstract machine model of a network switch/router?
 - What is the key software abstraction of computer networks?
- Terminology
 - Frequency-Division Multiplexing and Time-Division Multiplexing
 - Packet Switching and Circuit Switching
 - Port, Queue, Traffic Manager, Forwarding Table, and Store-and-Forward
 - Layering and Protocol Stack

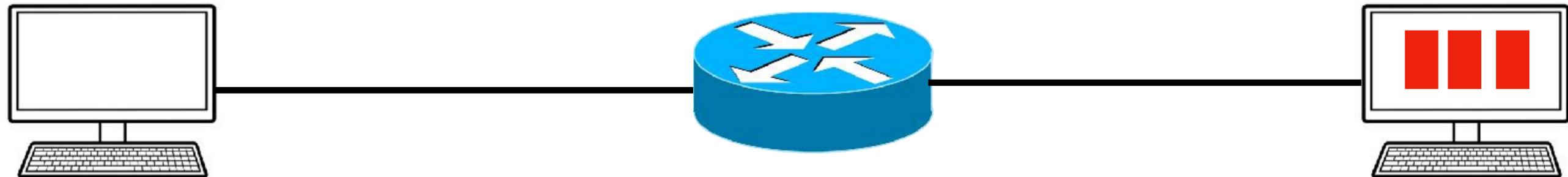
Packet-Switched Networks

- Packet communication path
 - Host 1 → Router A → Router B → ... → Host 2

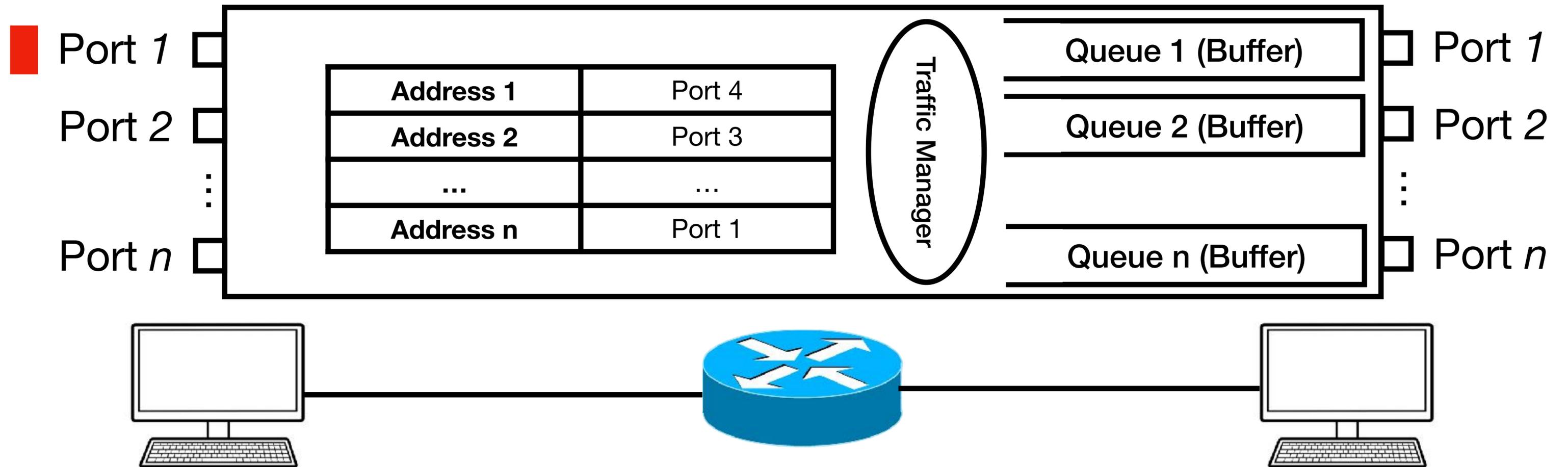


Packet-Switched Networks

- Packet communication path
 - Host 1 → Router A → Router B → ... → Host 2
- Delay
 - The total amount of time transfers N bits across the packet path
 - Four types of delay

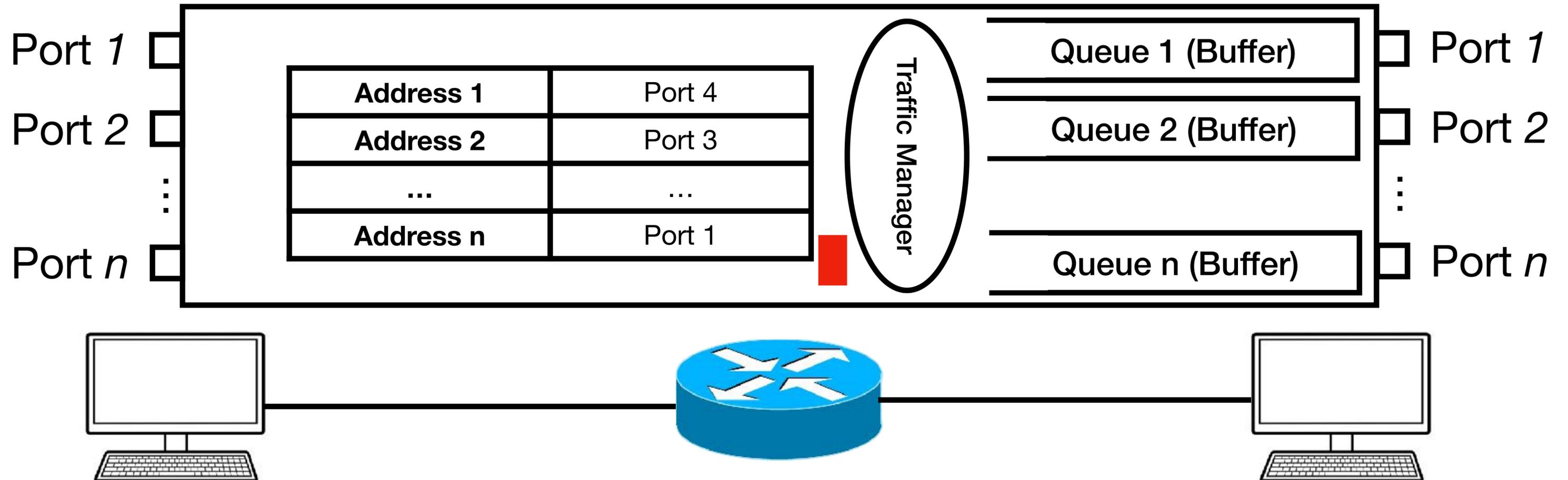


Processing Delay (T_{proc})

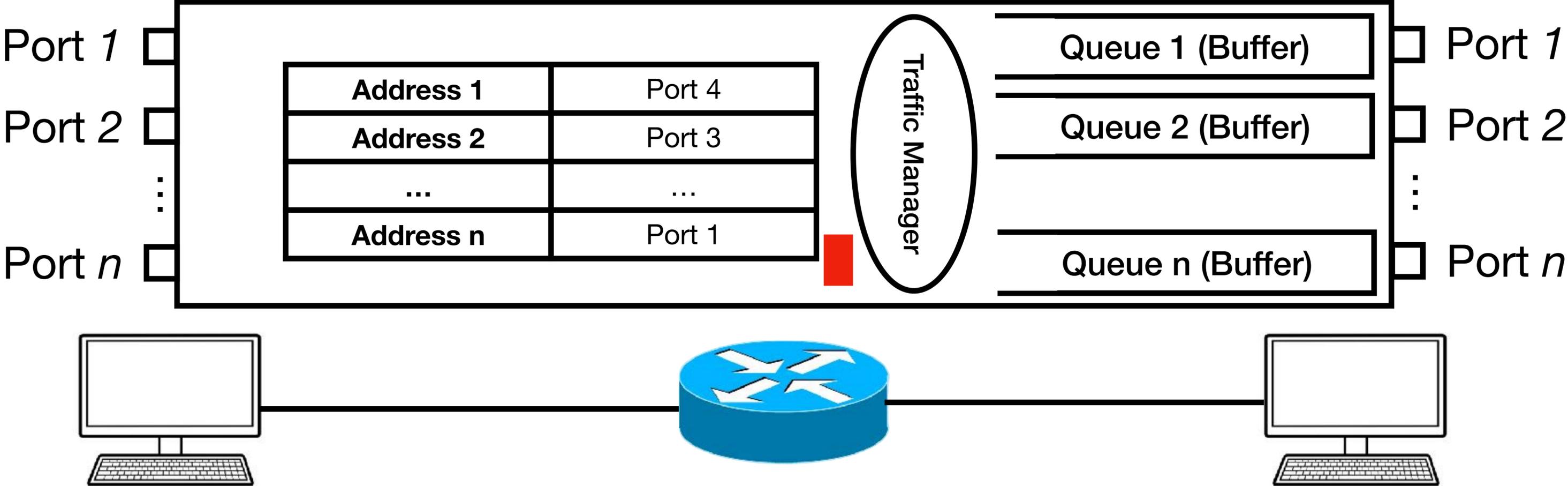


Processing Delay (T_{proc})

- The time required to examine the packet header and determine where to forward
 - Including checking bit-level errors
 - Microsecond

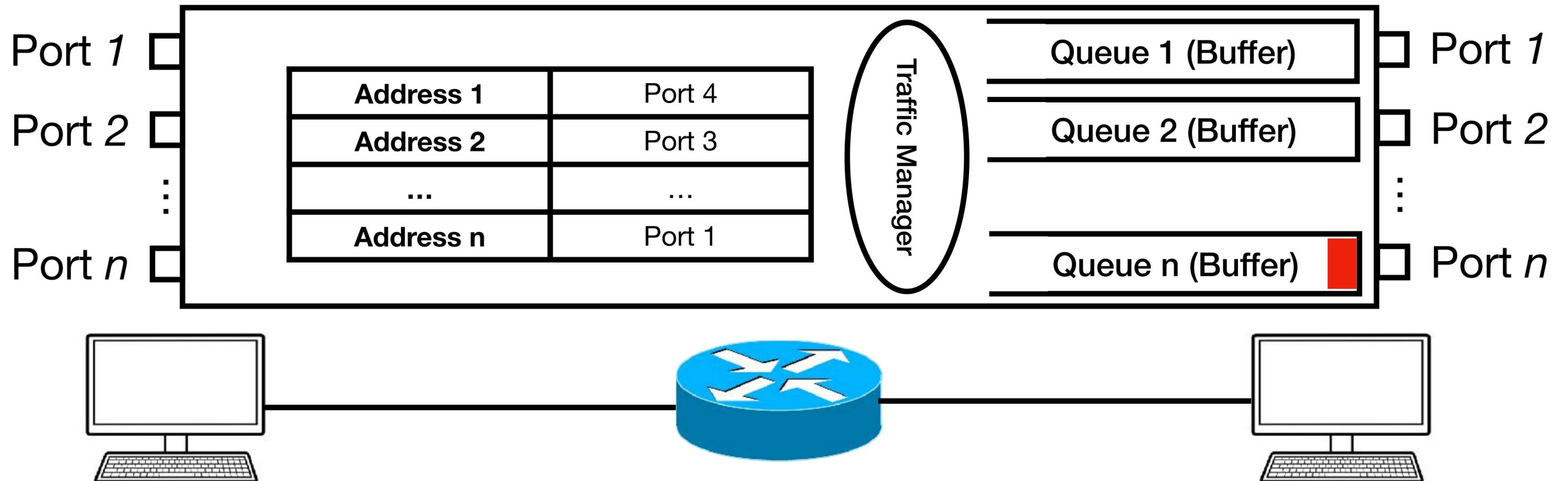


Queueing Delay (T_{queue})



Queueing Delay (T_{queue})

- The time it takes to wait to be transmitted
 - Depend on the number of earlier-arriving packets
 - Microsecond~Millisecond



The Little's Law

- A fundamental theory in mathematical queueing theory
 - John Little (1928-2024): Operational Research
 - Widely applicable to computer networks and systems
- $L = \lambda W$
 - Considering a stationary system
 - L : the long-term average number of customers
 - λ : the long-term average effective arrival rate
 - W : the average time that a customer spends in the system

The Little's Law

- A fundamental theory in mathematical queueing theory
 - John Little (1928-2024): Operational Research
 - Widely applicable to computer networks and systems

If a router processes 10,000 packets per second and a packet takes an average of 5ms to be processed and transmitted, how should we size the buffer?

The Little's Law

- A fundamental theory in mathematical queueing theory
 - John Little (1928-2024): Operational Research
 - Widely applicable to computer networks and systems

If a router processes 10,000 packets per second and a packet takes an average of 5ms to be processed and transmitted, how should we size the buffer?

$$L = \lambda W = 10,000 \text{ pkts/s} \times 0.005\text{s} = 50 \text{ packets}$$

The Little's Law

- A fundamental theory in mathematical queueing theory
 - John Little (1928-2024): Operational Research
 - Widely applicable to computer networks and systems

If a router sees 100 packets on average in the queue and the queue-draining rate is 50 packets/second, what is the average queueing latency?

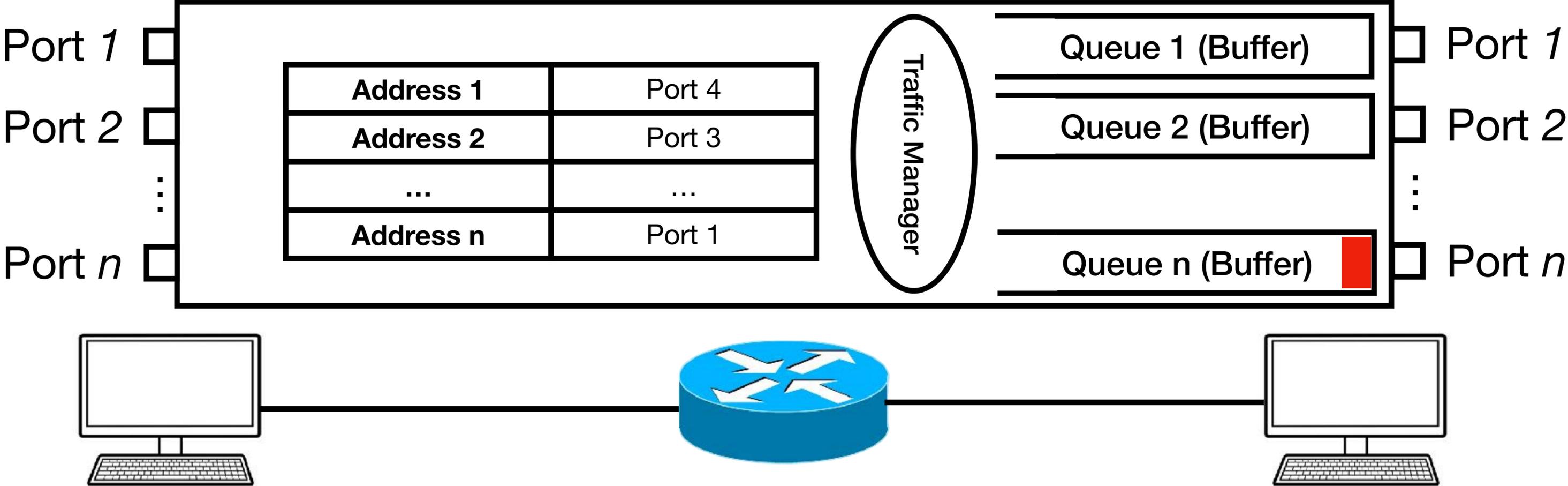
The Little's Law

- A fundamental theory in mathematical queueing theory
 - John Little (1928-2024): Operational Research
 - Widely applicable to computer networks and systems

If a router sees 100 packets on average in the queue and the queue-draining rate is 50 packets/second, what is the average queueing latency?

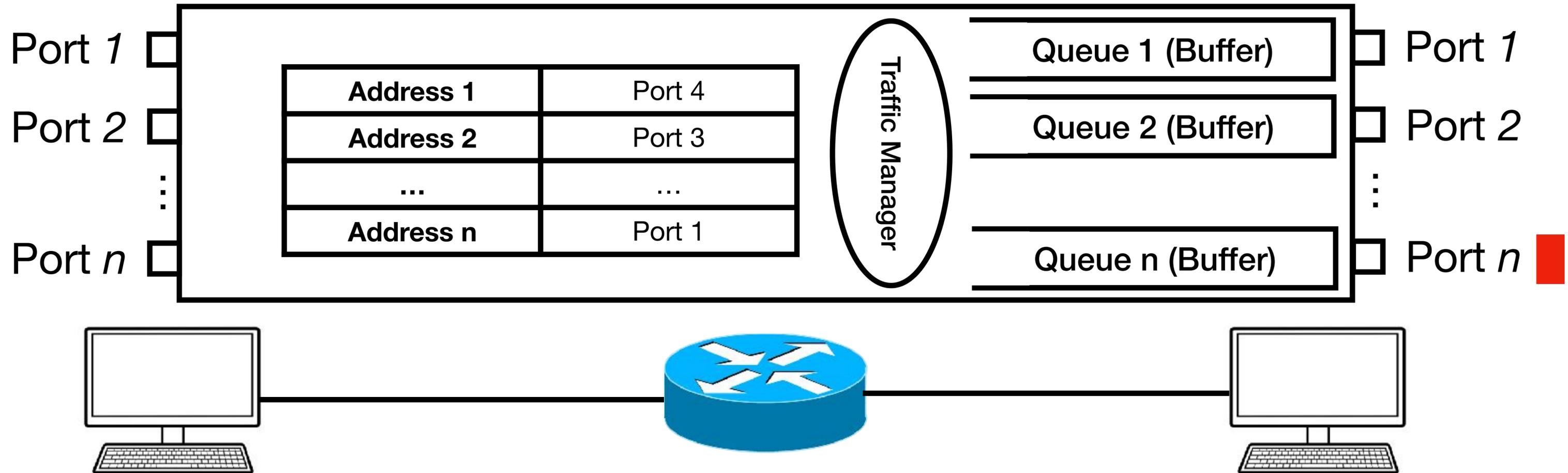
$$W=L/\lambda = 100 \text{ packets} / 50 \text{ packets/s} = 2 \text{ seconds}$$

Transmission Delay (T_{trans})



Transmission Delay (T_{trans})

- The time it takes to transmit through the communication ports
 - First-come-first-served manner
 - Depend on the packet length and transmission rate

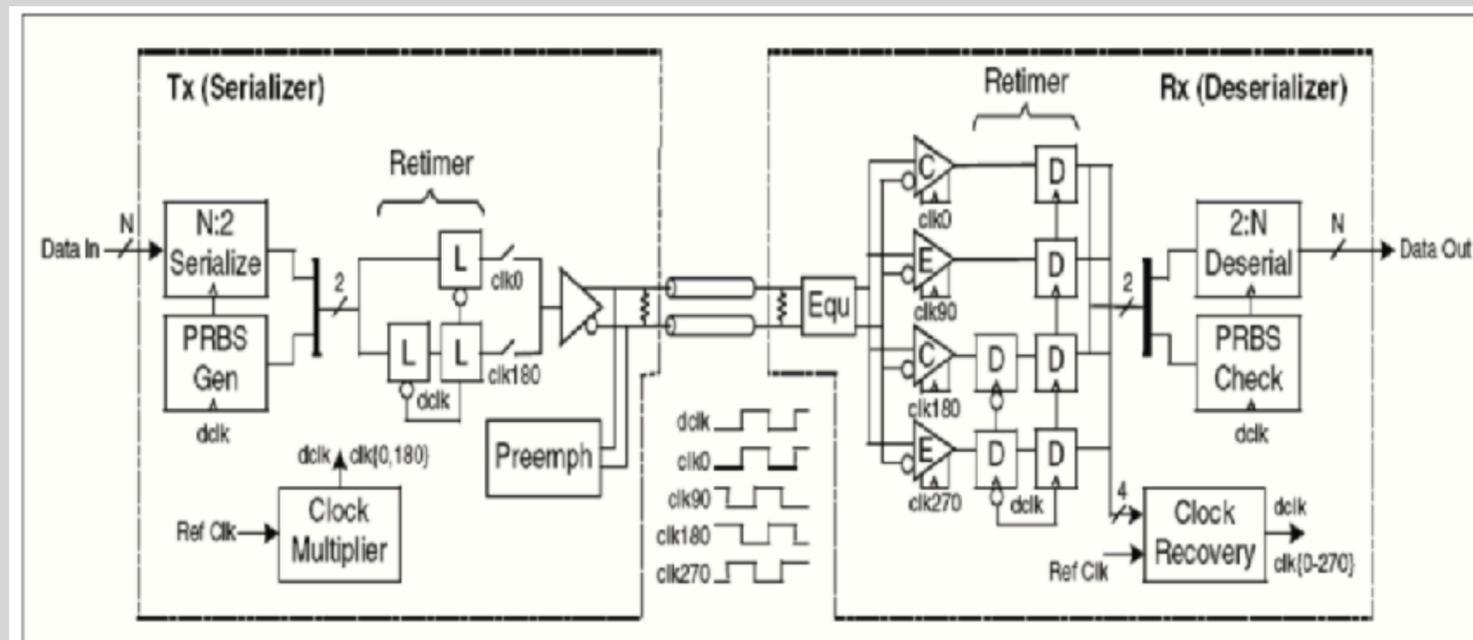


Transmission Delay (T_{trans})

- The time it takes to transmit through the communication ports
 - First-come-first-served manner

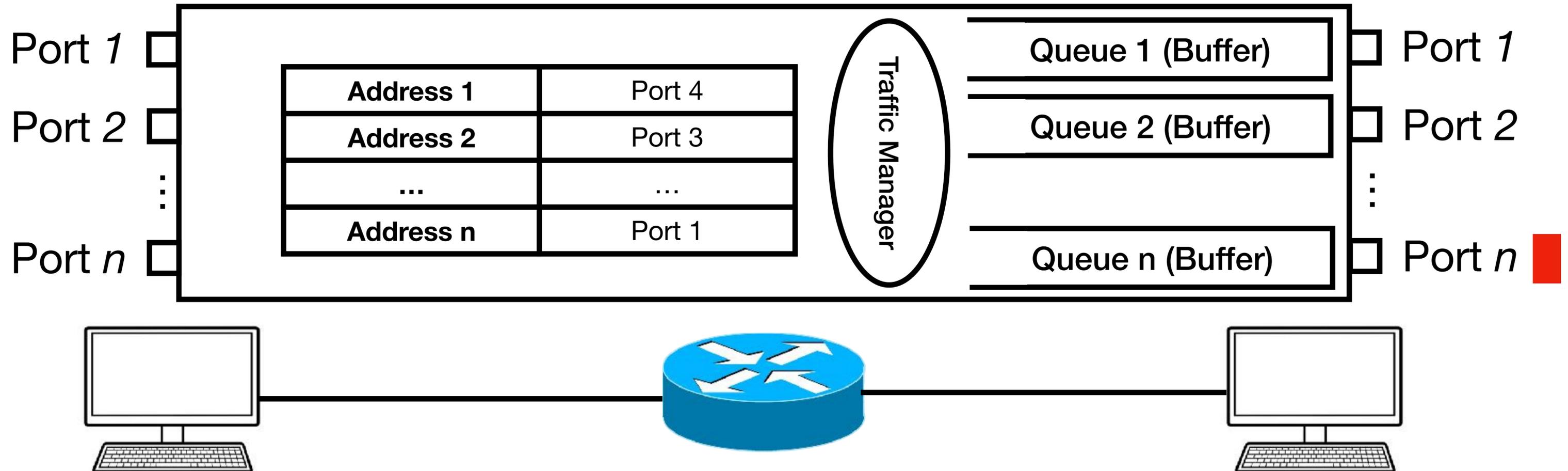
SerDes (Serializer/Deserializer)

- Convert parallel data to serial and vice versa

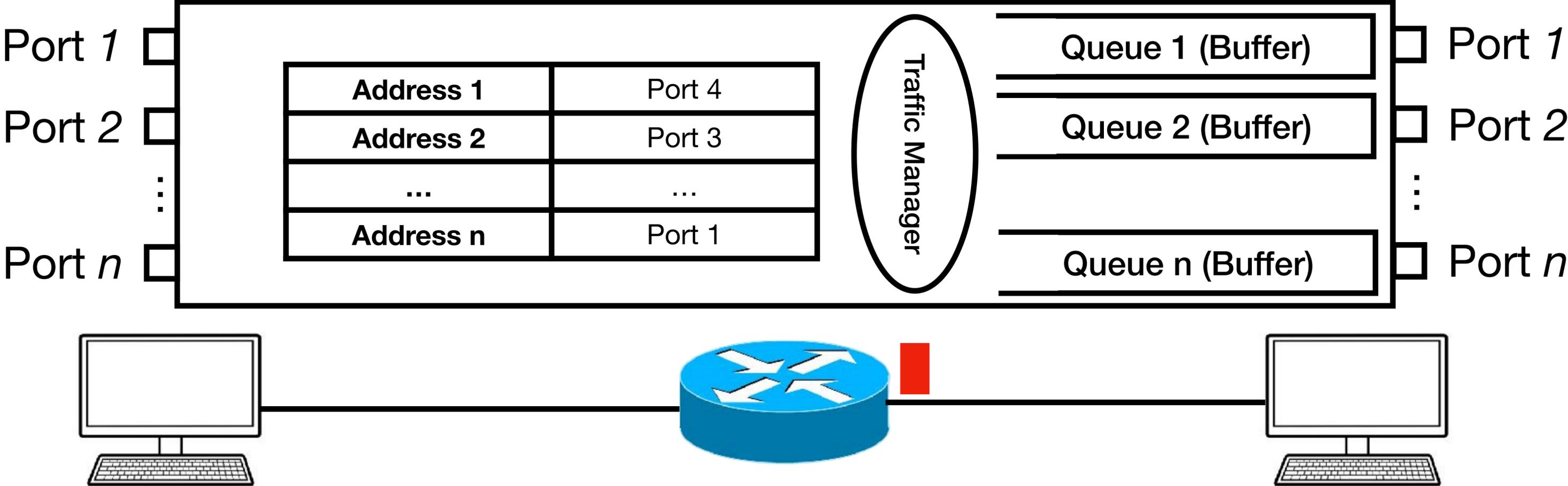


Transmission Delay Calculation

- Packet length (L bits) and Transmission Rate (R bits/sec)
 - Transmission Delay = L/R
 - Suppose $L=1\text{KB}$, $R = 10\text{Mbps}/100\text{Mbps}$, what is the transmission delay?

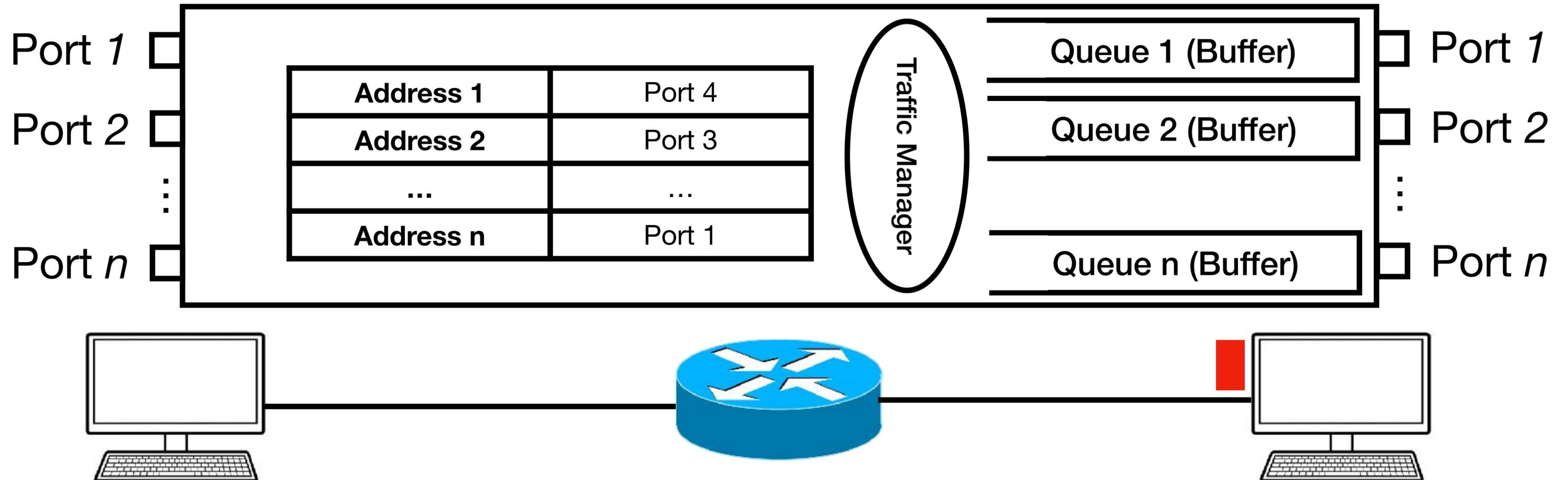


Propagation Delay (T_{prop})



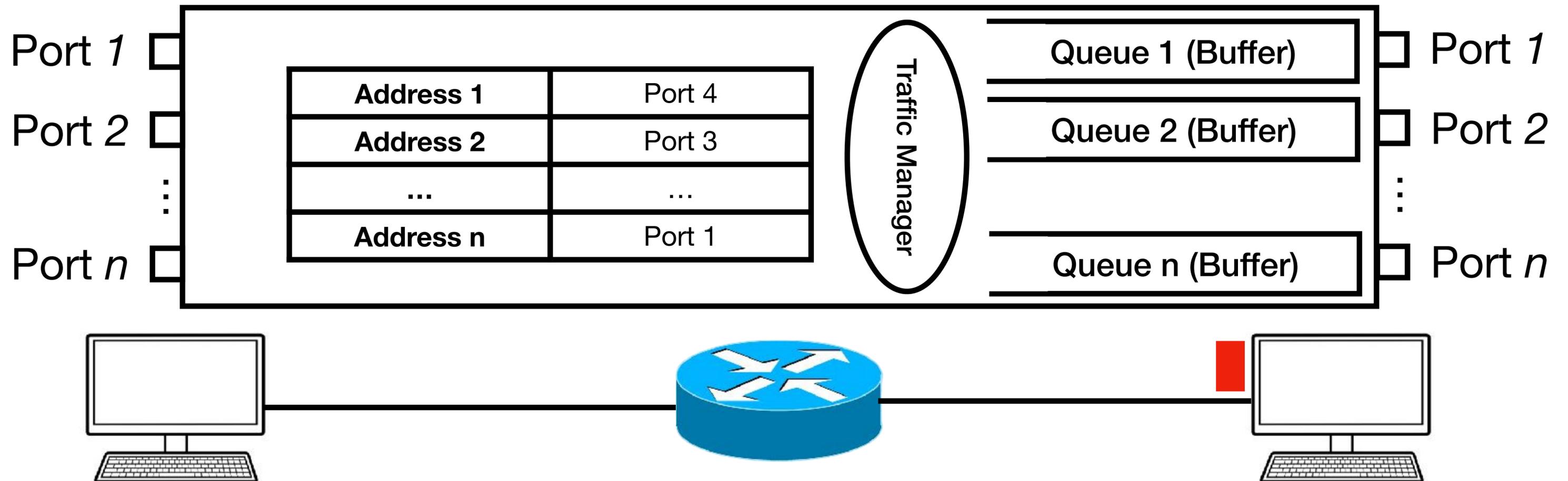
Propagation Delay (T_{prop})

- The time required to propagate over the communication link
 - Depend on the physical media (e.g., fiber optics, copper wire, etc.)
 - $2 \times 10^8 \sim 3 \times 10^8$ meters/second



Propagation Delay Calculation

- Distance (d) and propagation speed (s)
 - Propagation delay = d/s
 - Distance matters!



Transmission Delay v.s. Propagation Delay

- Difference
 - Transmission delay: time taken to push out the packet
 - Propagation delay: time taken to traverse the link

Transmission Delay v.s. Propagation Delay

- Difference
 - Transmission delay: time taken to push out the packet
 - Propagation delay: time taken to traverse the link
- Suppose
 - Packet length (L) = 1.5KB and transmission rate (R) = 1Mbps
 - Distance (d) = 3m and propagation speed (s) = 3^8 m/s
 - What is the transmission delay and propagation delay?

Transmission Delay v.s. Propagation Delay

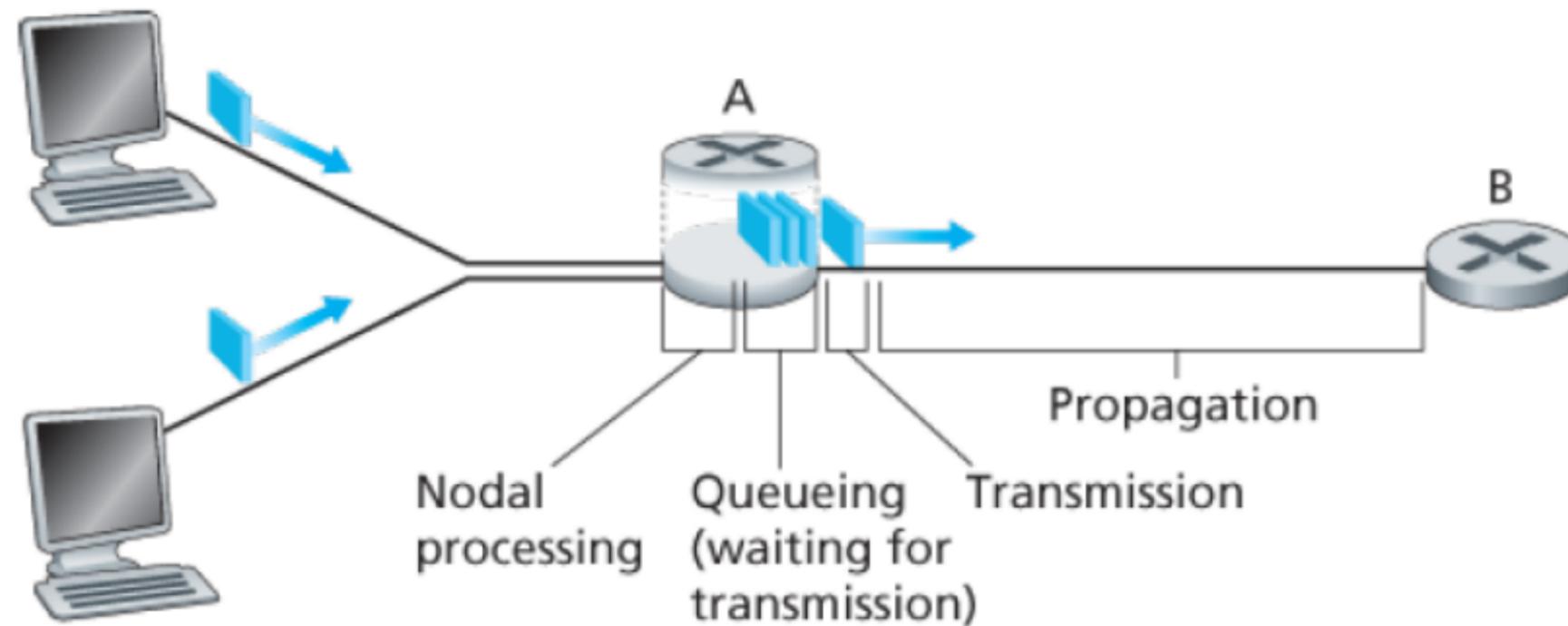
- Difference
 - Transmission delay: time taken to push out the packet
 - Propagation delay: time taken to traverse the link
- Suppose
 - Packet length (L) = 1.5KB and transmission rate (R) = 1Mbps
 - Distance (d) = 3m and propagation speed (s) = 3^8 m/s
 - What is the transmission delay and propagation delay?
 - If we increase the transmission rate to 1Gbps, what happens?

Transmission Delay v.s. Propagation Delay

- Difference
 - Transmission delay: time taken to push out the packet
 - Propagation delay: time taken to traverse the link
- Suppose
 - Packet length (L) = 1.5KB and transmission rate (R) = 1Mbps
 - Distance (d) = 3m and propagation speed (s) = 3^8 m/s
 - What is the transmission delay and propagation delay?
 - If we increase the transmission rate to 1Gbps, what happens?
 - If we increase the distance to 3km, what happens?

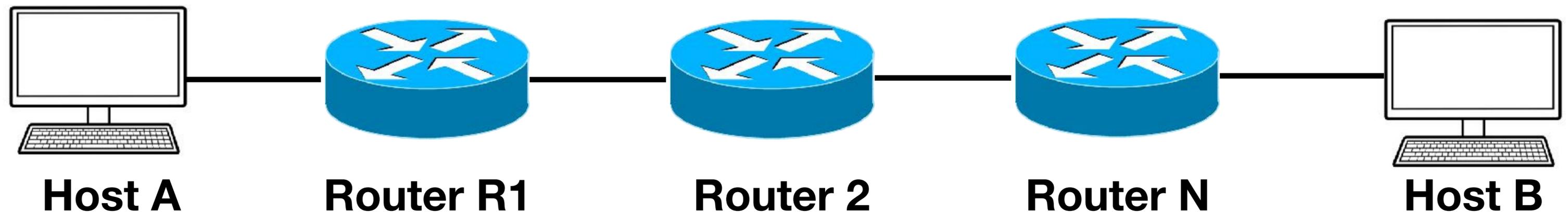
Total (Nodal) Delay

- $T_{\text{total}} = T_{\text{proc}} + T_{\text{queue}} + T_{\text{trans}} + T_{\text{prop}}$
 - Per-node, also called total nodal delay



End-to-End Delay

- Suppose
 - Host A can send a packet infinitely fast
 - Host B can receive a packet infinitely fast
- What is the end-to-end delay for a packet from host A to host B?
 - $T_{\text{end-to-end}} = N (T_{\text{proc}} + T_{\text{queue}} + T_{\text{trans}} + T_{\text{prop}}) + T_{\text{prop}}$



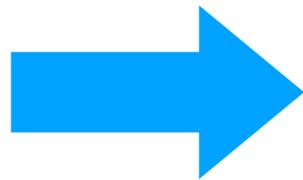
Queueing Delay

- The most complicated and interesting one
 - Varying packet to packet and scenario to scenario
- Statistical metrics
 - Average queueing delay
 - Variance of queueing delay
 - Tail queueing delay
 - The probability that the queueing delay exceeds some value

When the Queue is built up?

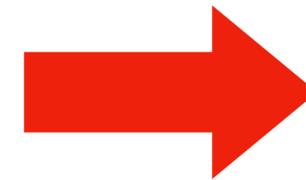
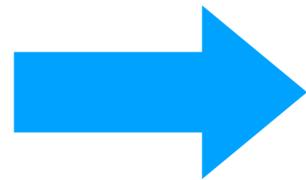


When the Queue is built up?



- Packet arrival rate **a**
- Unit: packets/second

When the Queue is built up?



- Packet arrival rate **a**
- Unit: packets/second

- Transmission rate **R**
- Unit: bits/second

When the Queue is built up?

- Given time T , incoming traffic load $>$ outgoing traffic load
 - Incoming traffic load = $\sum_{i=1}^N pkt_size_i$
 - Out going traffic load = $R * T$



- Packet arrival rate **a**
- Unit: packets/second

- Transmission rate **R**
- Unit: bits/second

Traffic Intensity

- We define the ratio $(\lambda L/R)$ as the traffic intensity
 - Suppose all packets consist of L bits

Traffic Intensity

- We define the ratio $(\lambda a/R)$ as the traffic intensity
 - Suppose all packets consist of L bits
- $\lambda a/R > 1$: queue built up!
 - The average rate at which bits arrive at the queue exceeds the rate at which the bits can be transmitted from the queue

Traffic Intensity

- We define the ratio $(\lambda a/R)$ as the traffic intensity
 - Suppose all packets consist of L bits
- $\lambda a/R > 1$: queue built up!
 - The average rate at which bits arrive at the queue exceeds the rate at which the bits can be transmitted from the queue

Design your system so that the traffic intensity is no greater than 1!

Traffic Intensity

- We define the ratio $(\lambda a/R)$ as the traffic intensity
 - Suppose all packets consist of L bits
- $\lambda a/R > 1$: queue built up!
 - The average rate at which bits arrive at the queue exceeds the rate at which the bits can be transmitted from the queue

Can we see queueing when $\lambda a/R \leq 1$?

Traffic Arrival process

- Suppose one packet arrives every L/R seconds

Traffic Arrival process

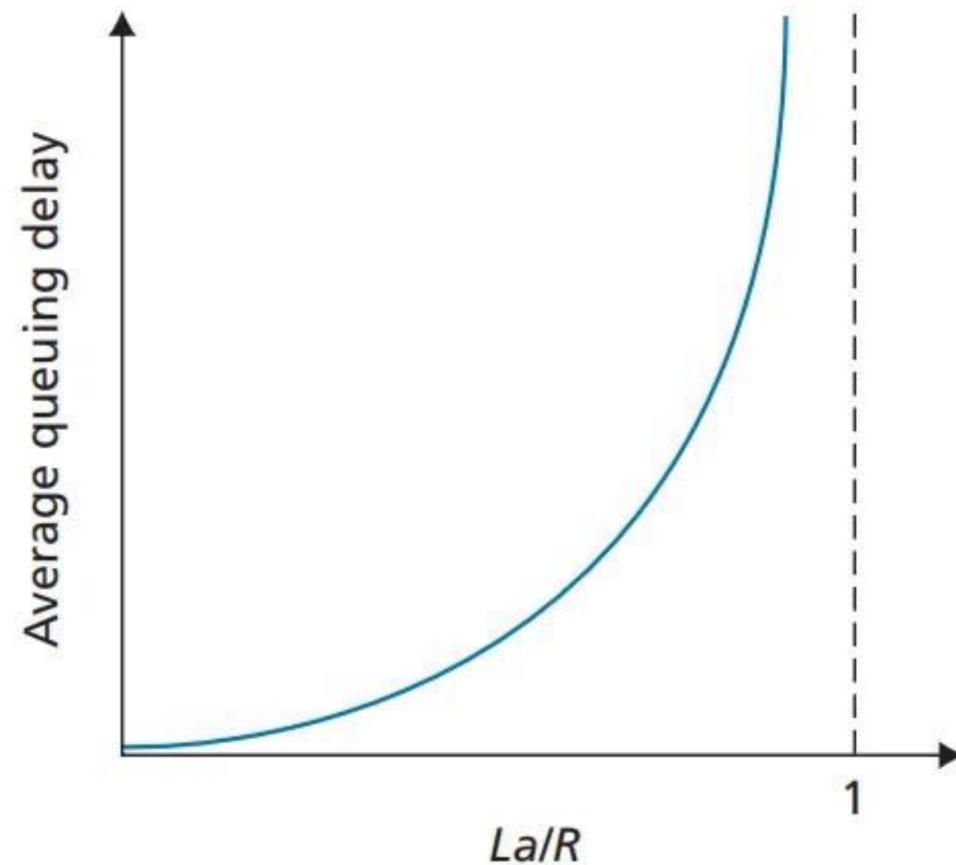
- Suppose one packet arrives every L/R seconds
 - No queueing
- Suppose N packets arrive every $(L/R)N$ seconds

Traffic Arrival process

- Suppose one packet arrives every L/R seconds
 - No queueing
- Suppose N packets arrive every $(L/R)N$ seconds
 - The 1st packet has no queueing delay, $T_{\text{queue}} = 0$
 - The 2nd packet has to wait for the 1st one, $T_{\text{queue}} = L/R \times 1$
 - The 3rd packet has to wait for the 1st and 2nd ones, $T_{\text{queue}} = L/R \times 2$
 - The n th packet has to wait for the $(n-1)$ ones, $T_{\text{queue}} = L/R \times (n-1)$

Traffic Intensity Curve

- Close to 1, the average queueing delay increases
 - A small percentage increase causes a significant delay increase
 - In reality, queue is fixed-sized ==> **Packet Loss**

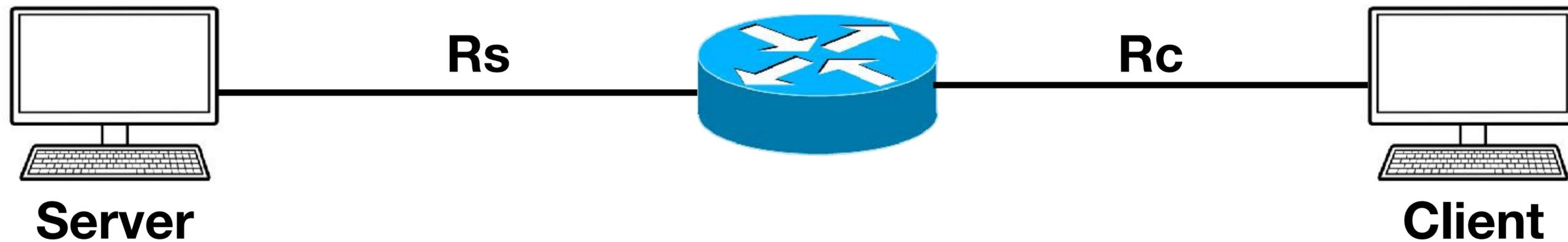


Throughput

- Bandwidth: the number of bits transmitted per second at a communication port and link
 - bps, Kbps, Mbps, Gbps, Tbps, ...
- Throughput: the number of bits transmitted from A to B
 - A and B can be host, switch, etc.
 - bps, Kbps, Mbps, Gbps, Tbps, ...

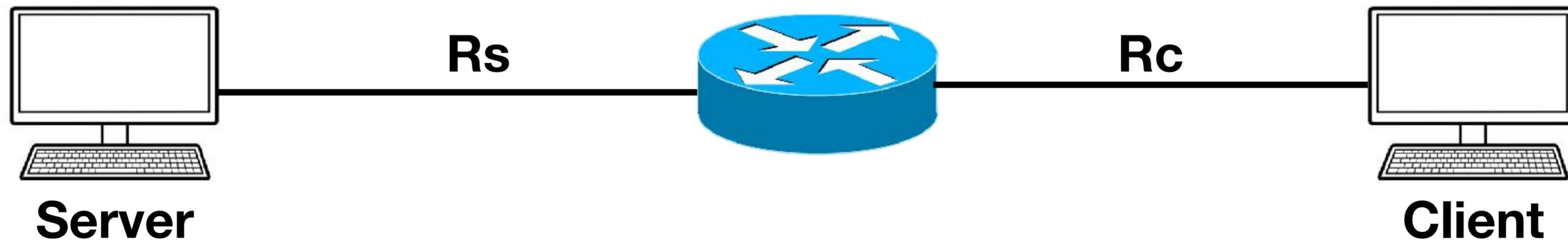
A Simple Throughput Example

- A file transferring from a server to a client



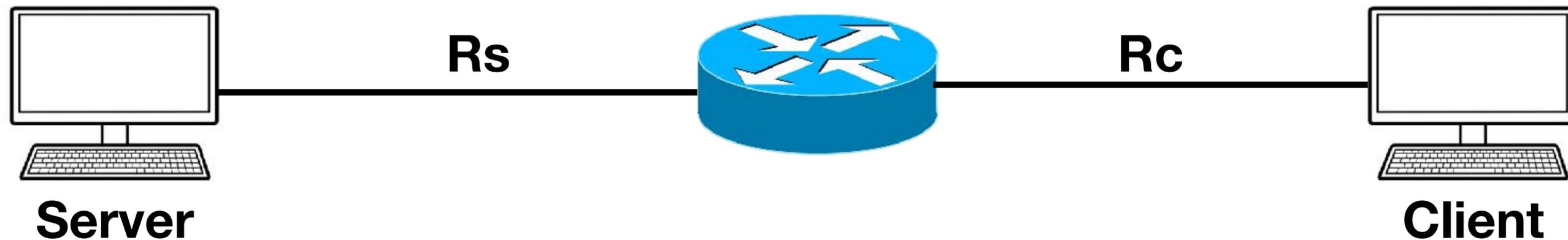
A Simple Throughput Example

- A file transferring from a server to a client
 - If $R_s < R_c \implies$ The client receives the file at R_s
 - If $R_s > R_c \implies$ The client receives the file at R_c , but router is queued



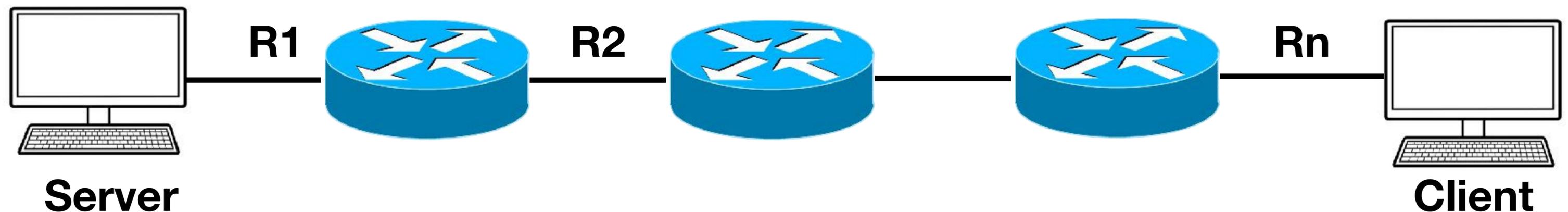
A Simple Throughput Example

- A file transferring from a server to a client
 - If $R_s < R_c \implies$ The client receives the file at R_s
 - If $R_s > R_c \implies$ The client receives the file at R_c , but router is queued
- Throughput = $\min \{R_s, R_c\}$
 - Depend on the bottleneck link
- A log file = 32M bits, $R_s = 2\text{Mbps}$, $R_c = 1\text{Mbps}$
 - Transfer time is 32 seconds



Another Throughput Example

- Communication path: server -> a list of routers -> Client
- Throughput = $\min \{R_1, R_2, R_3, \dots, R_n\}$
 - Depend on the bottleneck link



Video Streaming: a Bandwidth Perspective

- A distributed application that delivers video content over the Internet without downloading the entire file first.
 - Bandwidth intensive
 - Real-time
- Serving 360P video frames at 30 FPS (Frames Per Second)
 - What is the required downloading bandwidth?

Video Streaming: a Bandwidth Perspective

- Serving 360P video frames at 30 FPS (Frames Per Second)

Video Streaming: a Bandwidth Perspective

- Serving 360P video frames at 30 FPS (Frames Per Second)
 - 360P Resolution: 640 X 360 Pixels => 230,400 Pixels

Video Streaming: a Bandwidth Perspective

- Serving 360P video frames at 30 FPS (Frames Per Second)
 - 360P Resolution: 640 X 360 Pixels => 230,400 Pixels
 - Uncompressed frame: 24-bit (3 bytes, RGB) => 691,200 Bytes

Video Streaming: a Bandwidth Perspective

- Serving 360P video frames at 30 FPS (Frames Per Second)
 - 360P Resolution: 640 X 360 Pixels => 230,400 Pixels
 - Uncompressed frame: 24-bit (3 bytes, RGB) => 691,200 Bytes
 - 30 FPS: 691,200 Bytes/Frame X 30FPS => 20,736,000 Bytes

Video Streaming: a Bandwidth Perspective

- Serving 360P video frames at 30 FPS (Frames Per Second)
 - 360P Resolution: 640 X 360 Pixels => 230,400 Pixels
 - Uncompressed frame: 24-bit (3 bytes, RGB) => 691,200 Bytes
 - 30 FPS: 691,200 Bytes/Frame X 30FPS => 20,736,000 Bytes
 - Required downloading bandwidth: **~166 Mbps!**

Video Streaming: a Bandwidth Perspective

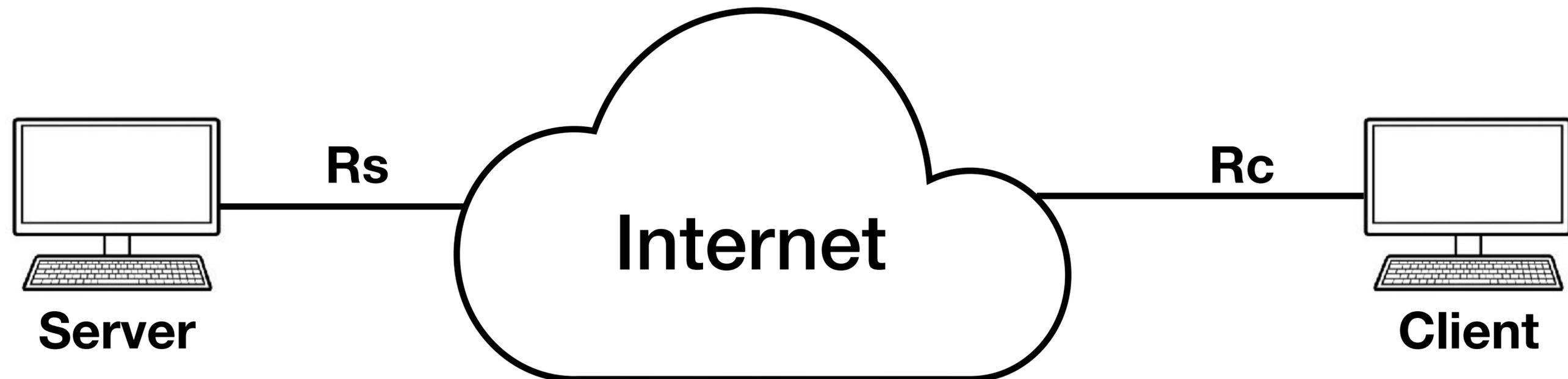
- Serving 360P video frames at 30 FPS (Frames Per Second)
 - 360P Resolution: 640 X 360 Pixels => 230,400 Pixels
 - Uncompressed frame: 24-bit (3 bytes, RGB) => 691,200 Bytes
 - 30 FPS: 691,200 Bytes/Frame X 30FPS => 20,736,000 Bytes
 - Required downloading bandwidth: **~166 Mbps!**
- Let's do H.264 video encoding
 - H.264 compressed 360P frame: ~1.6KB — 5.2KB

Video Streaming: a Bandwidth Perspective

- Serving 360P video frames at 30 FPS (Frames Per Second)
 - 360P Resolution: 640 X 360 Pixels => 230,400 Pixels
 - Uncompressed frame: 24-bit (3 bytes, RGB) => 691,200 Bytes
 - 30 FPS: 691,200 Bytes/Frame X 30FPS => 20,736,000 Bytes
 - Required downloading bandwidth: **~166 Mbps!**
- Let's do H.264 video encoding
 - H.264 compressed 360P frame: ~1.6KB — 5.2KB
 - 30FPS: 48KB/s — 156KB/s
 - Required downloading bandwidth: **384Kbps — 1.248Mbps!**

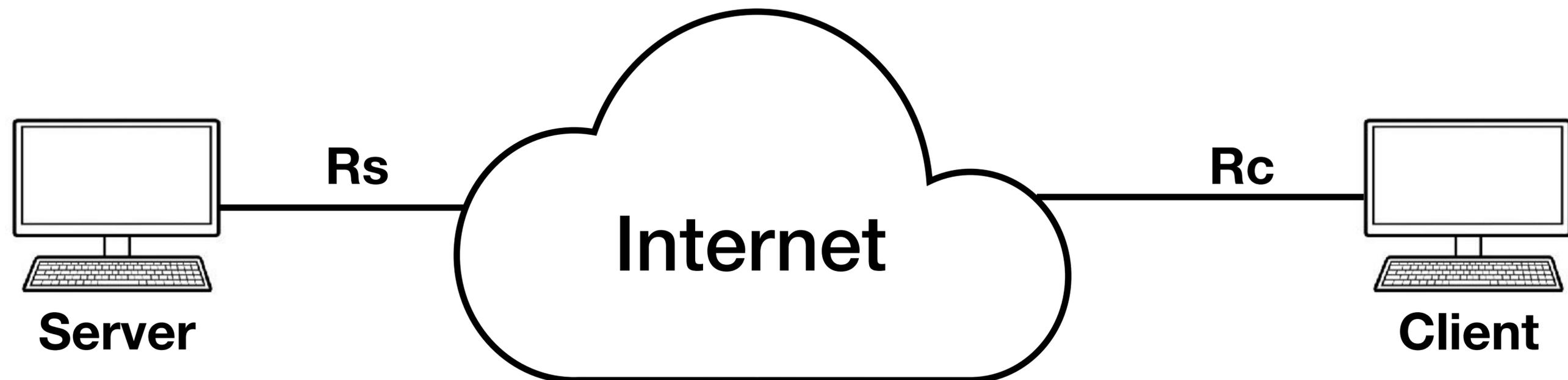
Throughput in a Shared Network

- Impossible to know the communication path details
- What is the throughput to transfer a file from a server to a client?



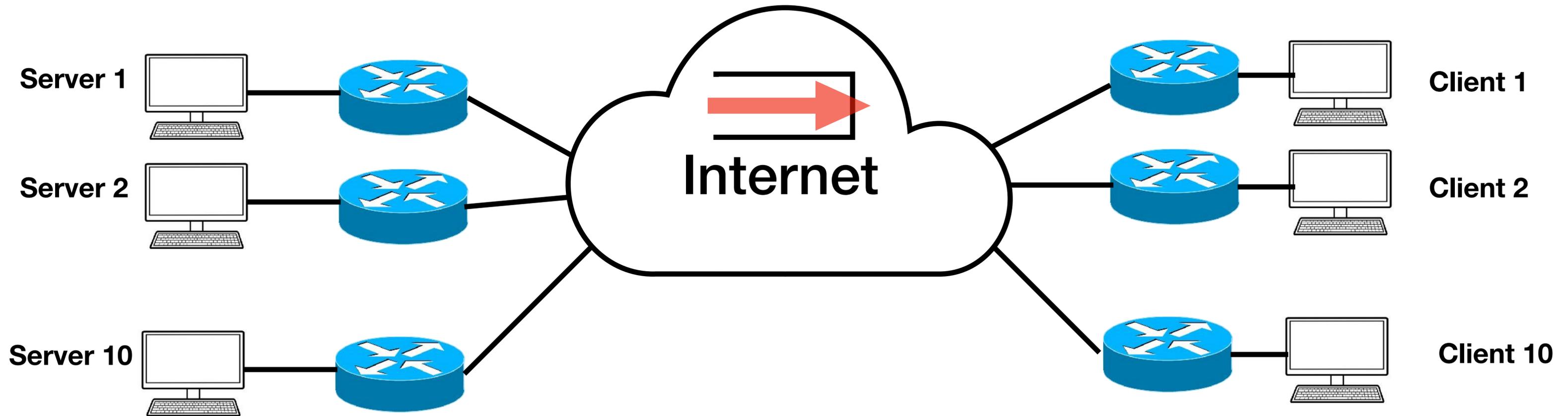
Throughput in a Shared Network

- Impossible to know the communication path details
- What is the throughput to transfer a file from a server to a client?
 - Actual Throughput = File Size / Total Transfer Time
 - Actual Throughput $\leq \min \{R_s, R_c\}$



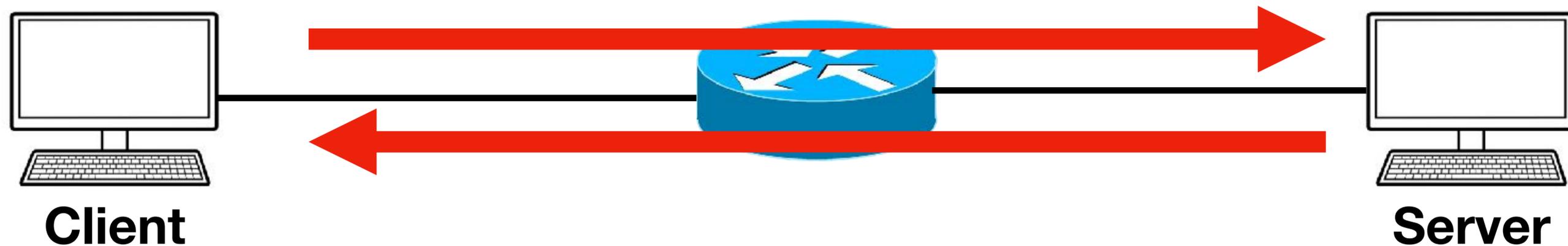
Throughput under Concurrent Transmissions

- Throughput also depends on intervening traffic



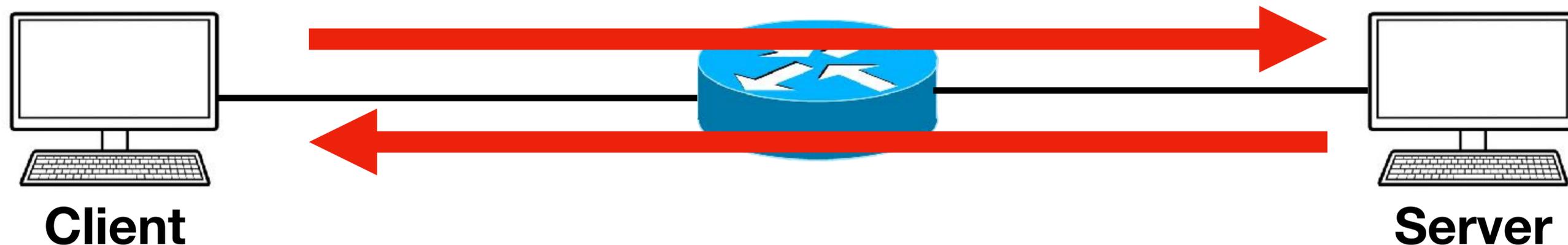
RTT (Round-Trip Time)

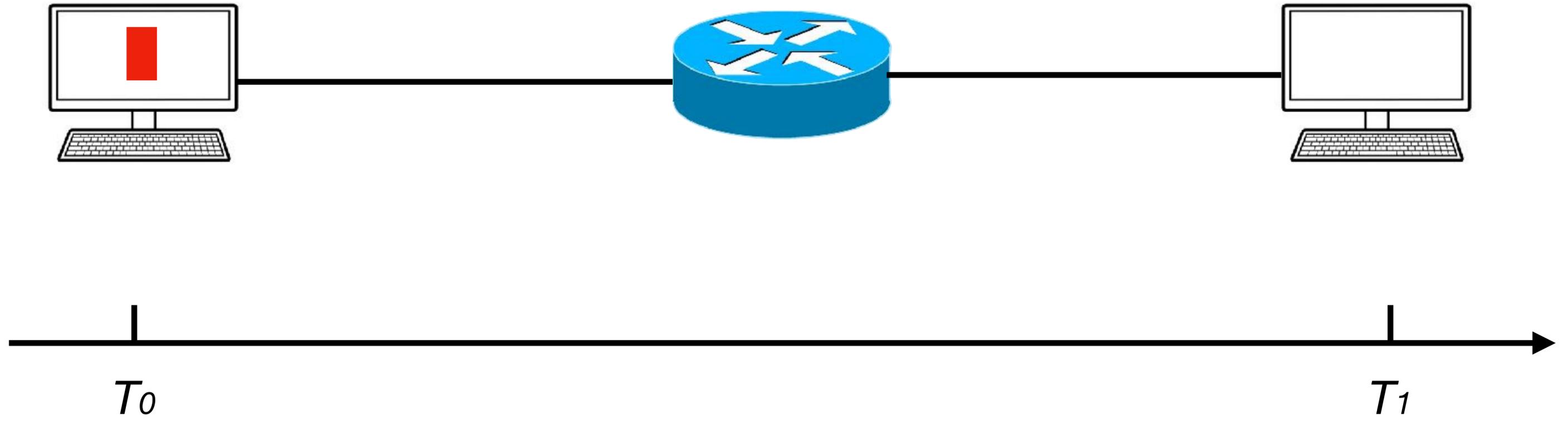
- The time it takes to send a request from a client to a server and receive the response back
 - $RTT = E2E \text{ delay (Client} \rightarrow \text{Server)} + E2E \text{ delay (Server} \rightarrow \text{Client)}$

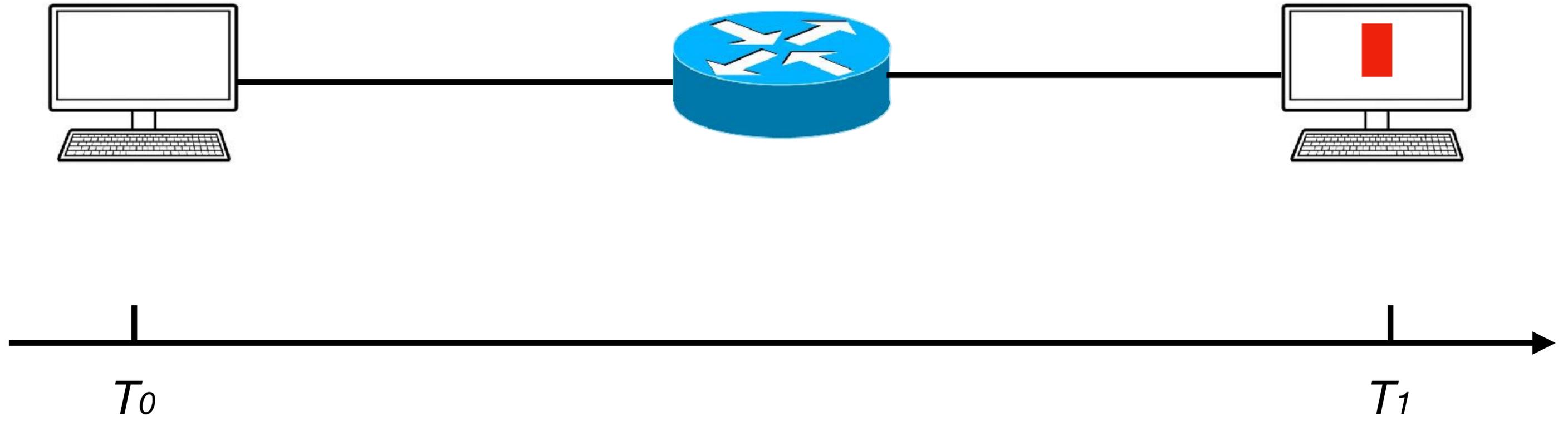


RTT (Round-Trip Time)

- The time it takes to send a request from a client to a server and receive the response back
 - $RTT = E2E \text{ delay (Client} \rightarrow \text{Server)} + E2E \text{ delay (Server} \rightarrow \text{Client)}$
- RTT is application-dependent
 - Web browsing: page download time (time to retrieve the first object)
 - Cloud gaming: interactive latency
 - Video conferencing and streaming: Time to First Frame
 - LLM: Time to first token

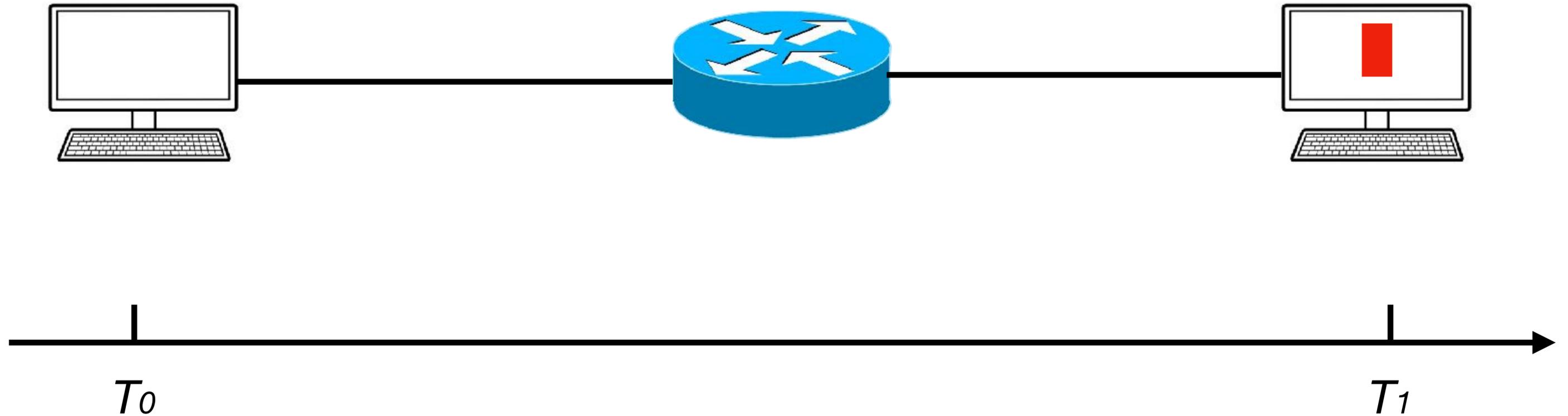






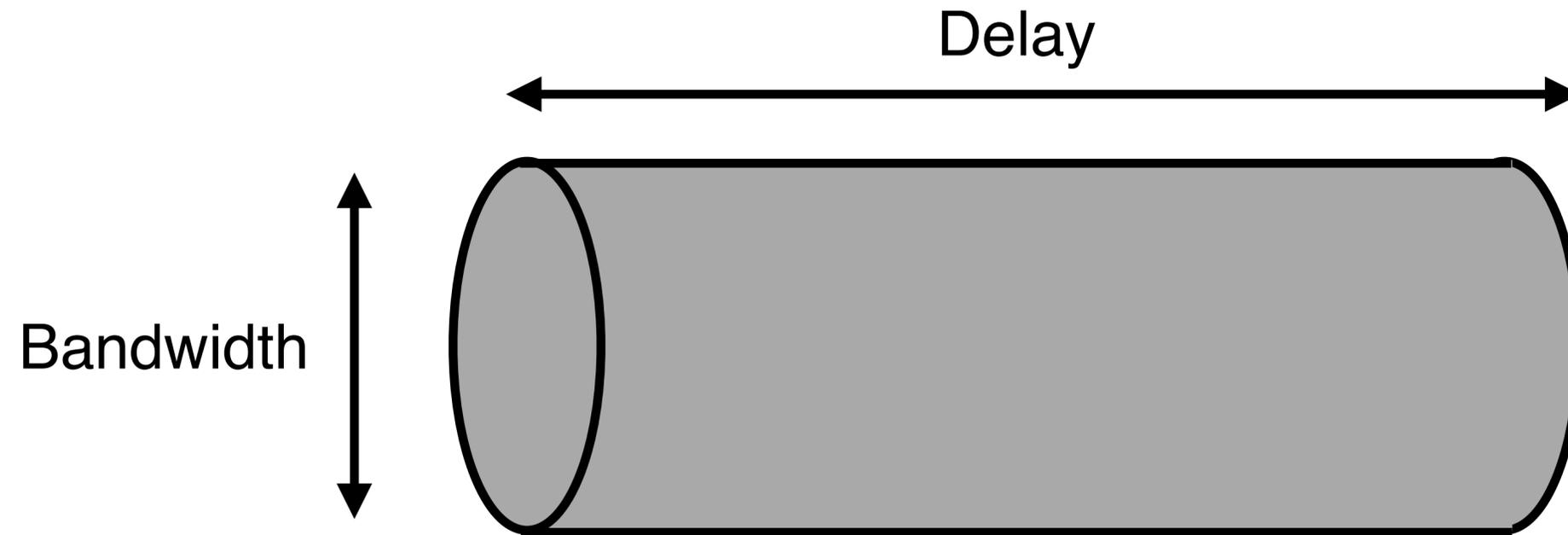
Where is the packet between T_0 and T_1 ?

- Networks become an implicit queue.



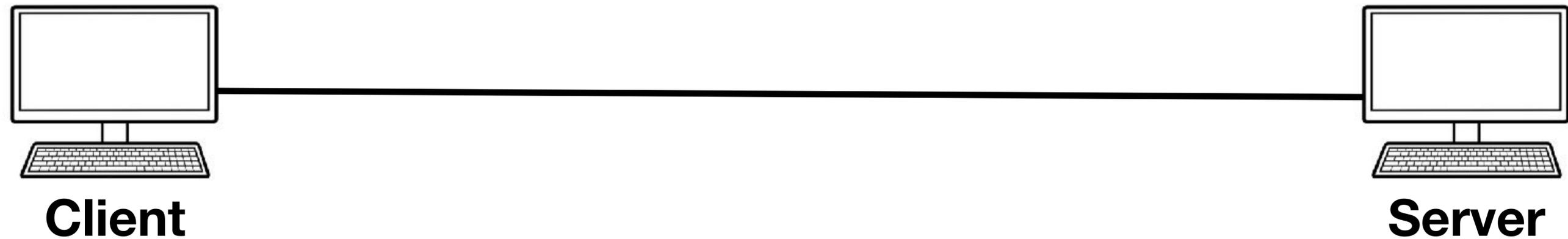
BDP (Bandwidth-Delay Product)

- The volume of a data pipe for one transmission
 - Bandwidth * Delay
 - The number of bits have left the sender and are yet to reach the receiver



BDP of a Communication Link

- Delay=Propagation delay \implies Link BDP
 - The number of bits over the wire



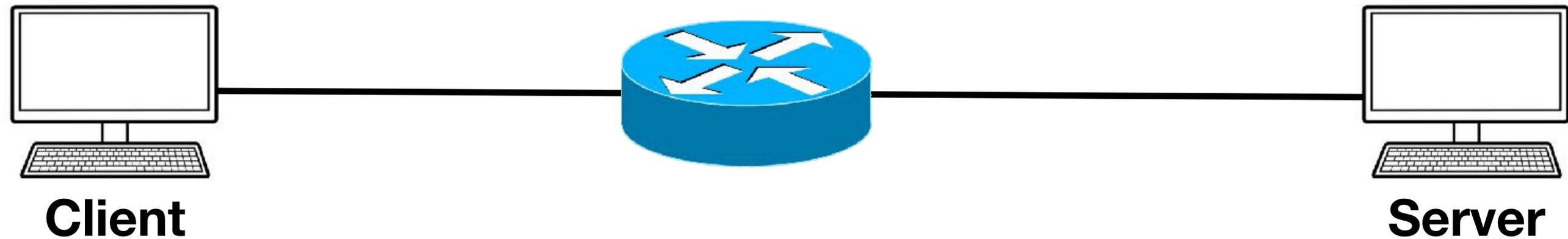
BDP of a Router

- Delay=Propagation delay \implies Link BDP
 - The number of bits over the wire
- Delay=Processing + Queueing + Transmission \implies Router BDP
 - The number of bits a router can hold



BDP of Client \rightarrow Server

- Delay=Propagation delay \Rightarrow Link BDP
 - The number of bits over the wire
- Delay=Processing + Queueing + Transmission \Rightarrow Router BDP
 - The number of bits a router can hold
- Delay=Total delay \Rightarrow End-to-End BDP
 - The number of bits that stay in-flight between two hosts



Summary

- Today
 - Computer networks: performance analysis
- Next lecture
 - Physical Layer: Encoding