



# CS120: Computer Networks

## **Lecture 18. Congestion Control 2**

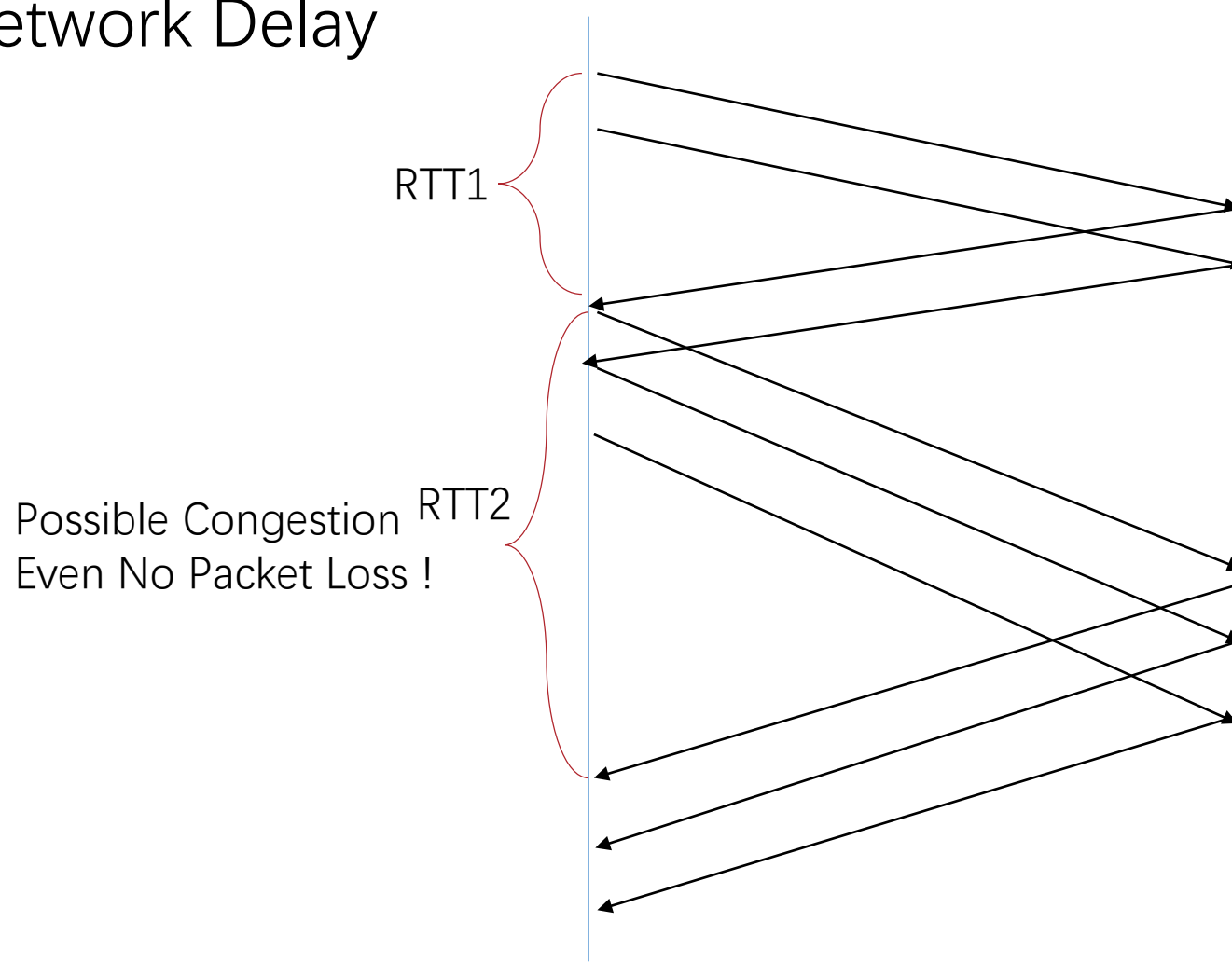
Zhice Yang

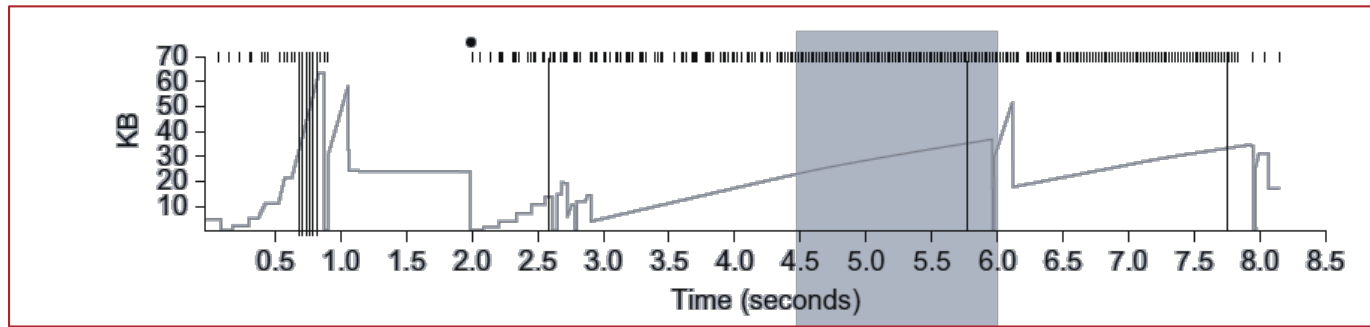
# Congestion Control

- Host-based Congestion Control
  - Packet Loss
    - Delay
- Router-based Congestion Control
  - Queuing

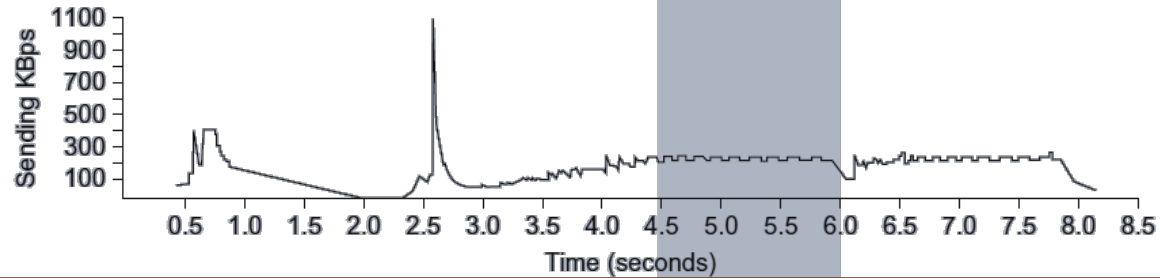
# Avoid Congestion

- Feedback: Network Delay





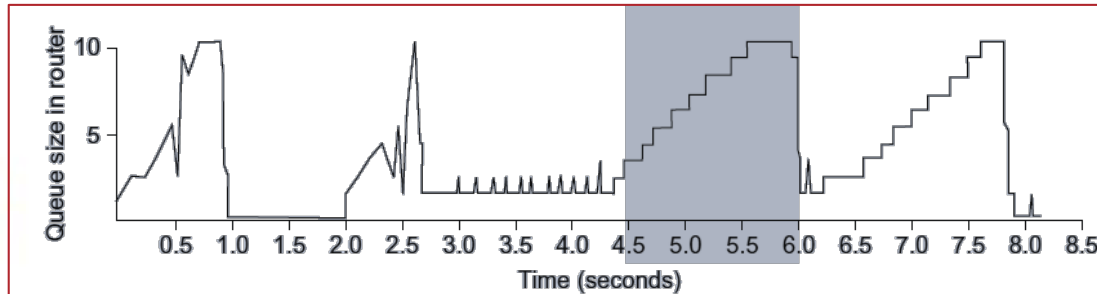
cwnd



Source 1

throughput

throughput



Queue Size

100-Mbps Ethernet

Source 2

Destination

cwnd increases but rate is flat  
when network approach congestion

# TCP Vegas

- Idea: TCP source uses RTT as the sign to **avoid** network congestion
  - Compare RTT with BaseRTT
- Method:
  - Reduce rate when congestion is about to happen
    - If Actual RTT  $\gg$  Base RTT
  - Recover rate soon after bandwidth is available
    - if Actual RTT  $>$  Base RTT
    - Keep certain pressure on network

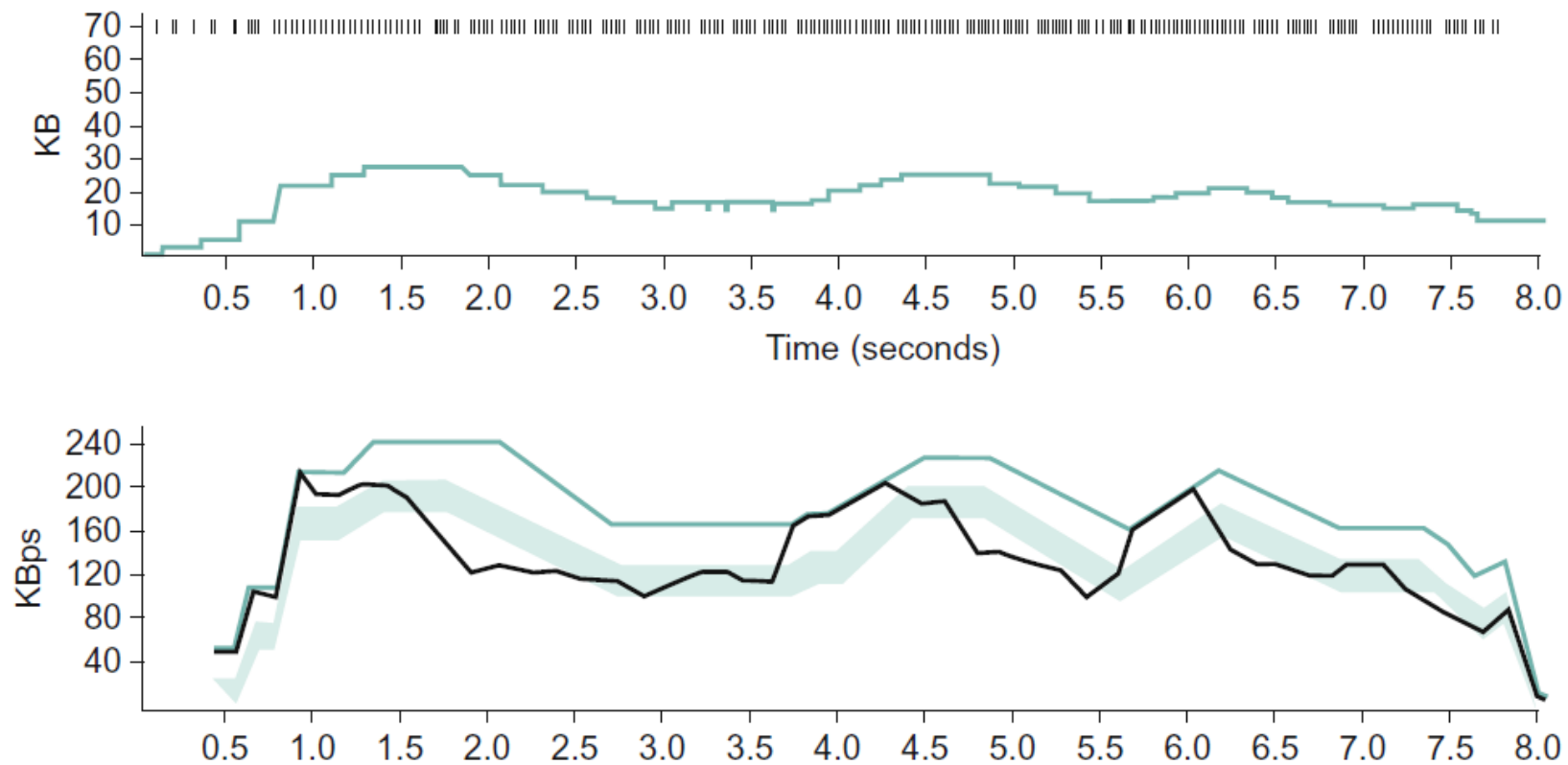
# TCP Vegas Algorithm

- BaseRTT: the reference for increases in RTTs
  - The minimum RTT
    - If  $RTT < BaseRTT$ 
      - $BaseRTT = RTT$
- ExpectRate:
  - $CongestionWindow / BaseRTT$
- ActualRate:
  - ActualRTT: according to timestamps
  - $ActualRate = CongestionWindow / ActualRTT$

# TCP Vegas Algorithm

- Source compares ActualRate with ExpectRate
  - if  $\text{ExpectRate} - \text{ActualRate} < \alpha$ 
    - $\text{cwnd}++$
  - if  $\alpha < \text{ExpectRate} - \text{ActualRate} < \beta$ 
    - $\text{cwnd} = \text{cwnd}$
  - if  $\beta < \text{ExpectRate} - \text{ActualRate}$ 
    - $\text{cwnd}--$

# TCP Vegas

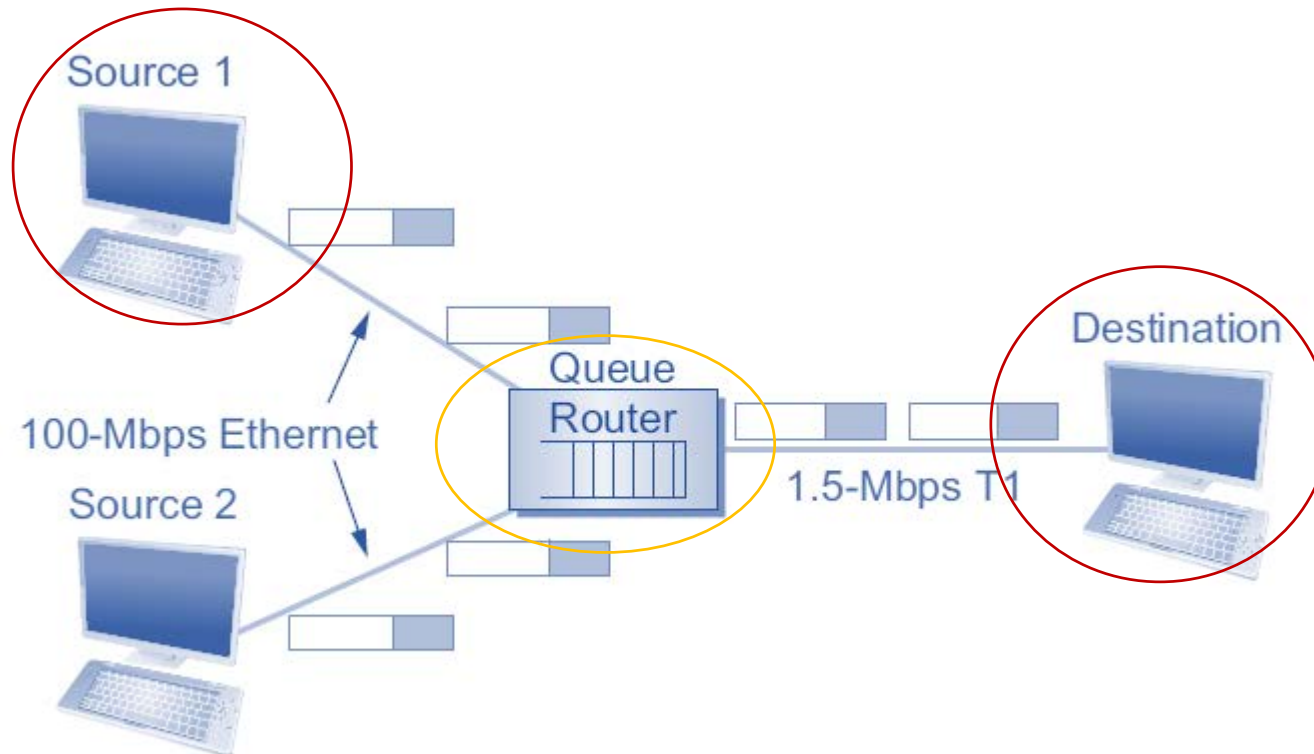


Top, congestion window; bottom, expected (colored line) and actual (black line) throughput. The shaded area is the region between the  $\alpha$  and  $\beta$  threshold



# Two Places to Handle Network Congestion

- End hosts
  - Routers



# Congestion Control

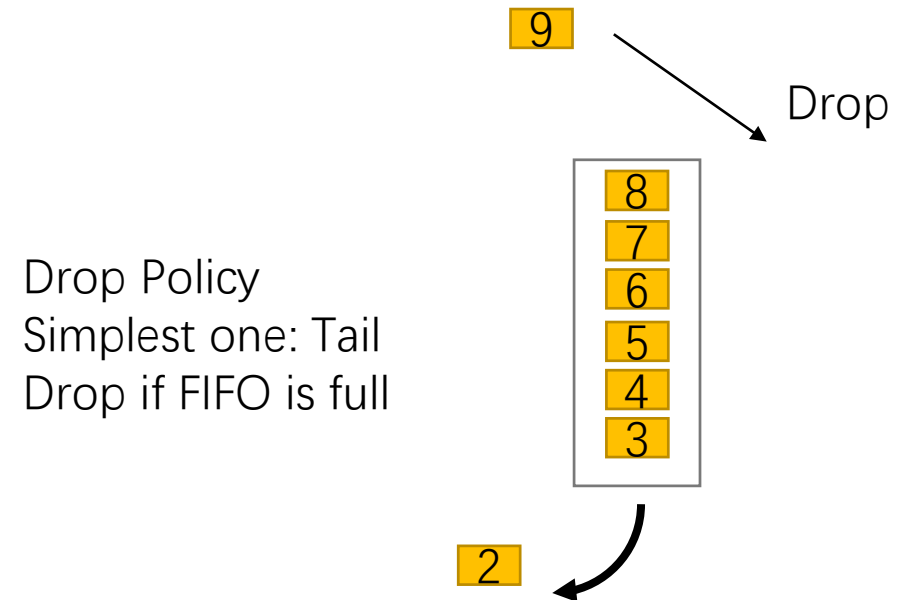
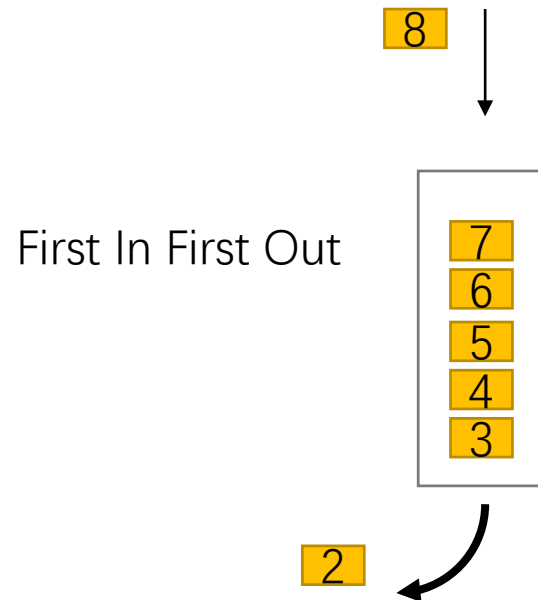
- Host-based Congestion Control
  - Packet Loss
  - Delay
- Router-based Congestion Control
  - Queuing

# Queuing Discipline

- Why ?
  - Queuing discipline in routers determines how packets are transmitted (or dropped)
- Router Resource
  - Bandwidth
    - Which packets get transmitted
  - Queue Buffer
    - Which packets get discarded

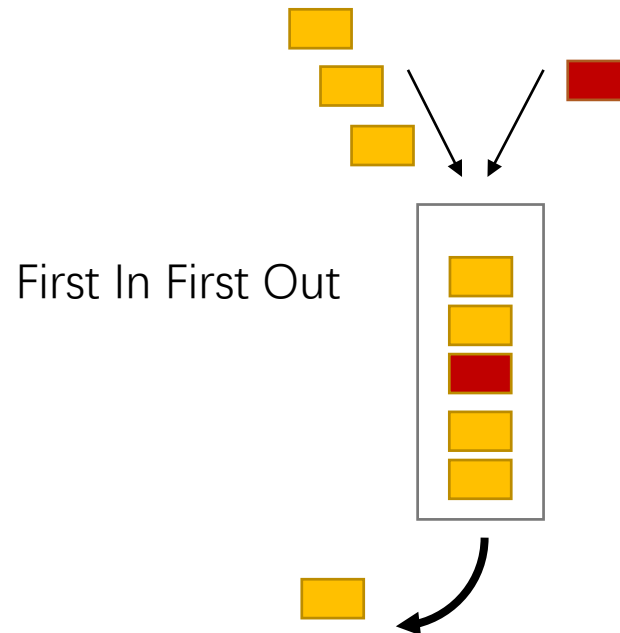
# Queuing Discipline

- First-In-First-Out (FIFO)



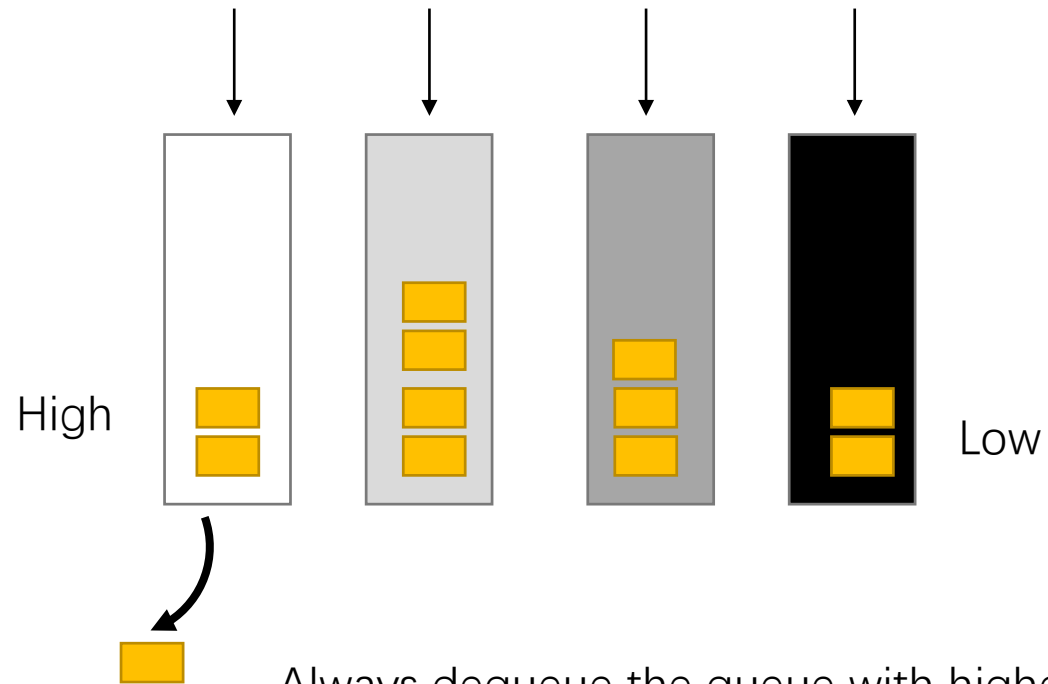
# Queuing Discipline

- Problems in FIFO
  - Too simple to provide resource allocation policies (to avoid congestion)
  - Hard to enforce every network source/flow to follow the same behavior
    - eg. yellow src does not follow congestion control (UDP), can occupy more network source



# Queuing Discipline

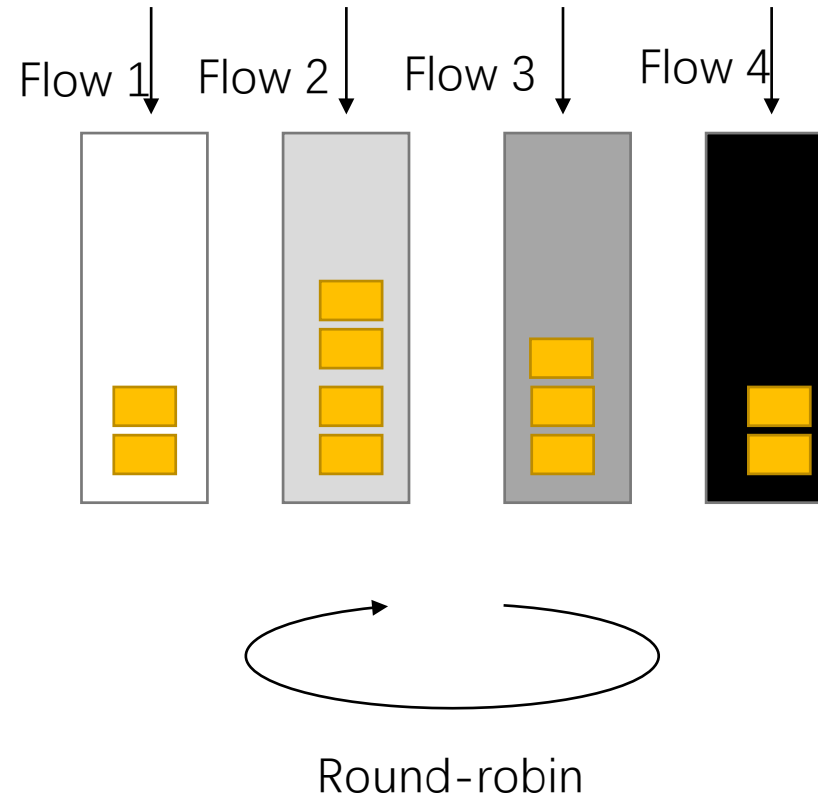
- First-In-First-Out (FIFO) with Priority



Always dequeue the queue with higher priority  
unless it is empty  
Problem: starvation

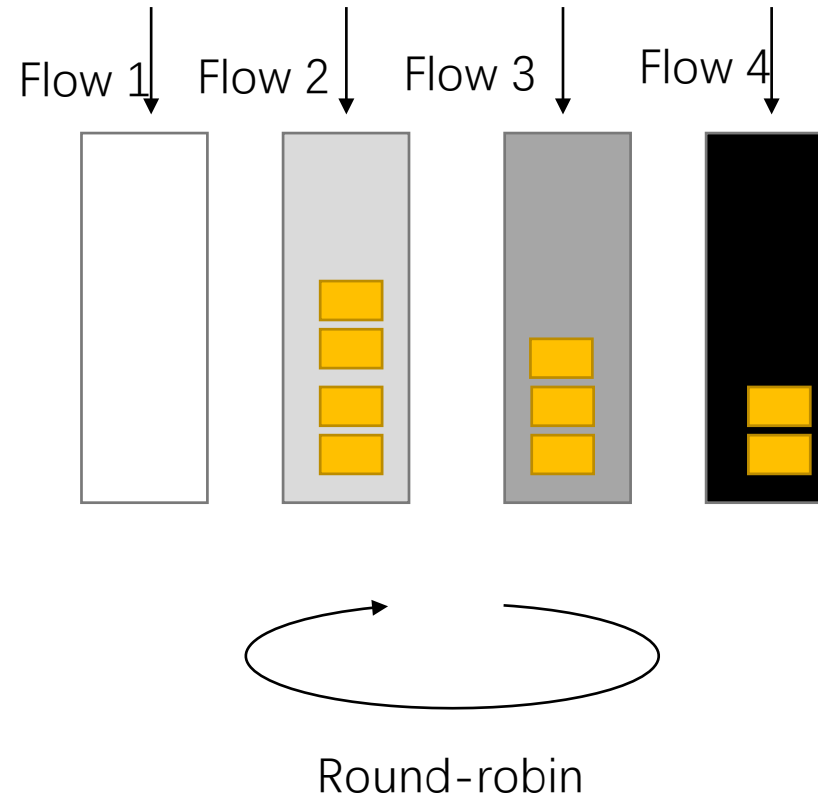
# Queuing Discipline

- Fair Queuing (FQ)
  - Each flow gets 1/4 output bandwidth



# Queuing Discipline

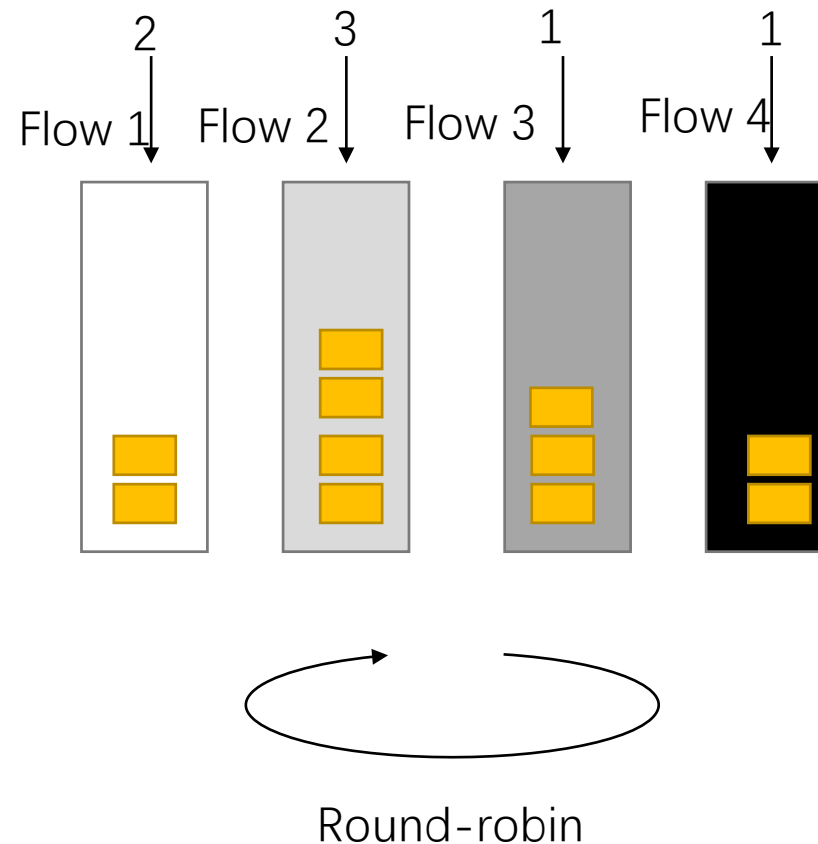
- Fair Queuing (FQ)
  - Each flow gets 1/3 output bandwidth





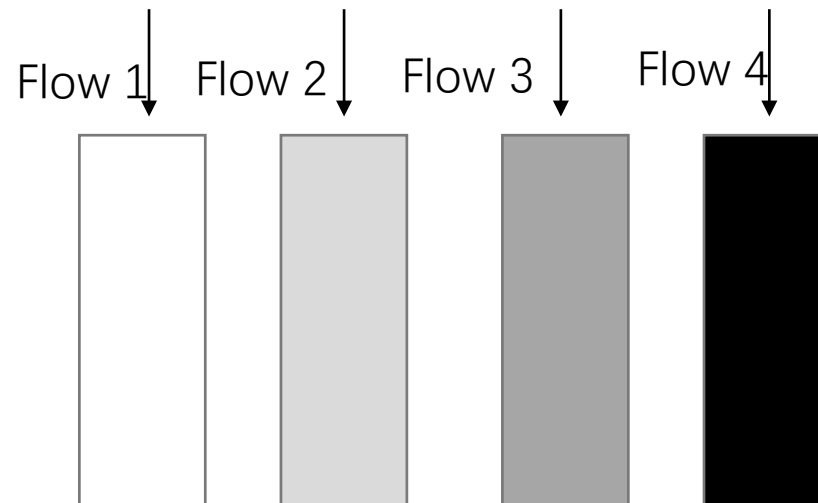
# Queuing Discipline

- Weighted Fair Queuing (FQ)
  - Flows with higher weight get more output bandwidth ( $2/7$ ,  $3/7$ ,  $1/7$ ,  $1/7$ )



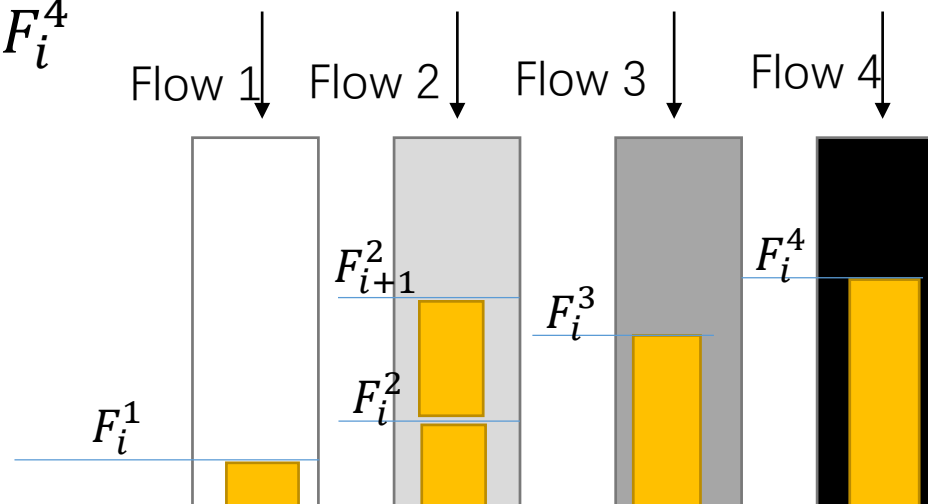
# Queuing Discipline

- Bit-Level Fair Queuing (FQ)
  - Schedule according to earliest finish time
  - Finish time of packet  $i$ :  $F_i = \max(F_{i-1}, A_i) + P_i$
  - $P_i$  is the transmitting duration of packet  $i$ ,  $A_i$  is the arriving time of packet  $i$ ,  $F_{i-1}$  is the finish time of packet  $i-1$  of the same flow



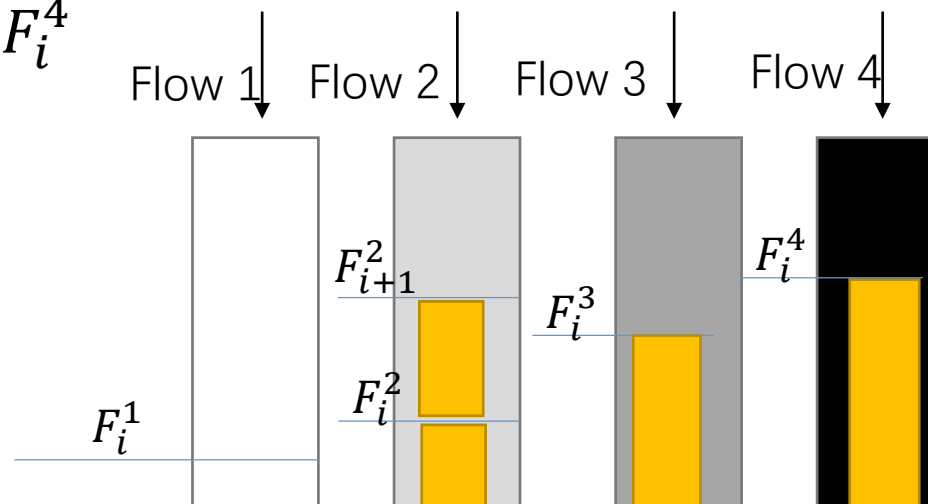
# Queuing Discipline

- Bit-Level Fair Queuing (FQ)
  - Schedule according to earliest finish time
  - Finish time of packet  $i$ :  $F_i = \max(F_{i-1}, A_i) + P_i$
  - $P_i$  is the transmitting duration of packet  $i$ ,  $A_i$  is the arriving time of packet  $i$ ,  $F_{i-1}$  is the finish time of packet  $i-1$  of the same flow
  - Finish time is virtually calculated in bit-level round robin, eg,  $F_i^1 < F_i^2 < F_i^3 < F_{i+1}^2 < F_i^4$



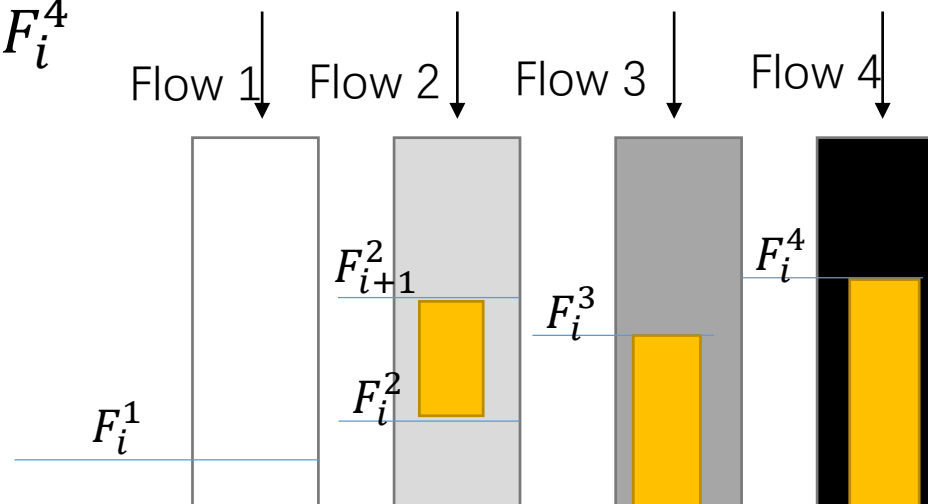
# Queuing Discipline

- Bit-Level Fair Queuing (FQ)
  - Schedule according to earliest finish time
  - Finish time of packet  $i$ :  $F_i = \max(F_{i-1}, A_i) + P_i$
  - $P_i$  is the transmitting duration of packet  $i$ ,  $A_i$  is the arriving time of packet  $i$ ,  $F_{i-1}$  is the finish time of packet  $i-1$  of the same flow
  - Finish time is virtually calculated in bit-level round robin, eg,  $F_i^1 < F_i^2 < F_i^3 < F_{i+1}^2 < F_i^4$



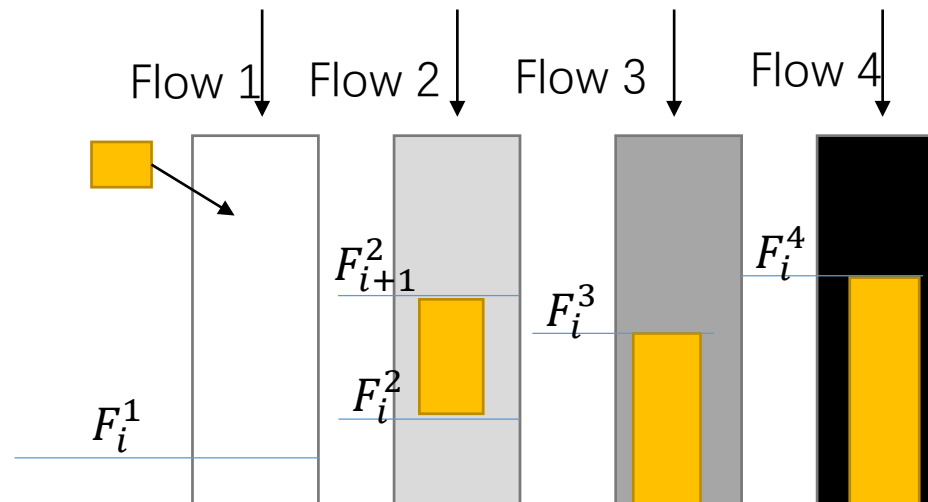
# Queuing Discipline

- Bit-Level Fair Queuing (FQ)
  - Schedule according to earliest finish time
  - Finish time of packet  $i$ :  $F_i = \max(F_{i-1}, A_i) + P_i$
  - $P_i$  is the transmitting duration of packet  $i$ ,  $A_i$  is the arriving time of packet  $i$ ,  $F_{i-1}$  is the finish time of packet  $i-1$  of the same flow
  - Finish time is virtually calculated in bit-level round robin, eg,  $F_i^1 < F_i^2 < F_i^3 < F_{i+1}^2 < F_i^4$



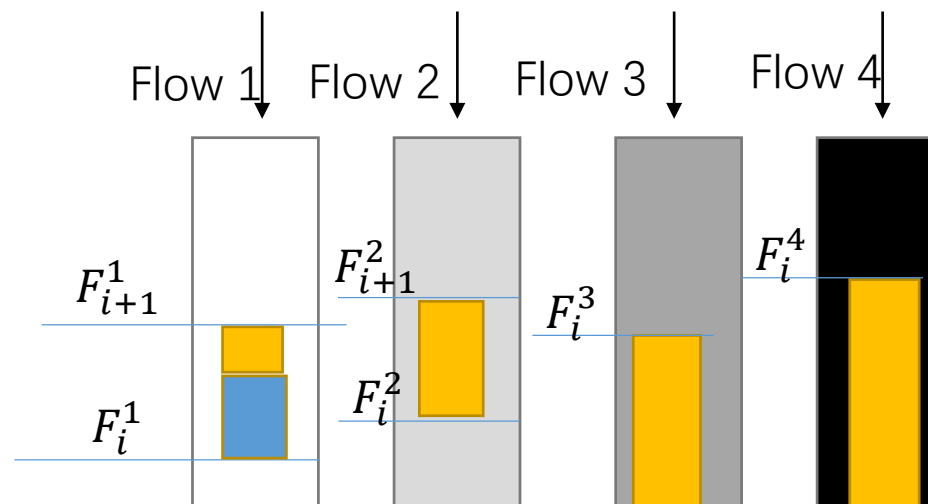
# Queuing Discipline

- Bit-Level Fair Queuing (FQ)
  - Schedule according to earliest finish time
  - Finish time of packet  $i$ :  $F_i = \max(F_{i-1}, A_i) + P_i$
  - $P_i$  is the transmitting duration of packet  $i$ ,  $A_i$  is the arriving time of packet  $i$ ,  $F_{i-1}$  is the finish time of packet  $i-1$  of the same flow
  - Case: when one queue is idle, idle time also count



# Queuing Discipline

- Bit-Level Fair Queuing (FQ)
  - Schedule according to earliest finish time
  - Finish time of packet  $i$ :  $F_i = \max(F_{i-1}, A_i) + P_i$
  - $P_i$  is the transmitting duration of packet  $i$ ,  $A_i$  is the arriving time of packet  $i$ ,  $F_{i-1}$  is the finish time of packet  $i-1$  of the same flow
  - Case: when one queue is idle, idle time also count



# Congestion Control

- Host-based Congestion Control
  - Packet Loss
  - Delay
- Router-based Congestion Control
  - Queuing
    - DECbit
    - Explicit congestion notification (ECN)
    - Random Early Detection (RED)

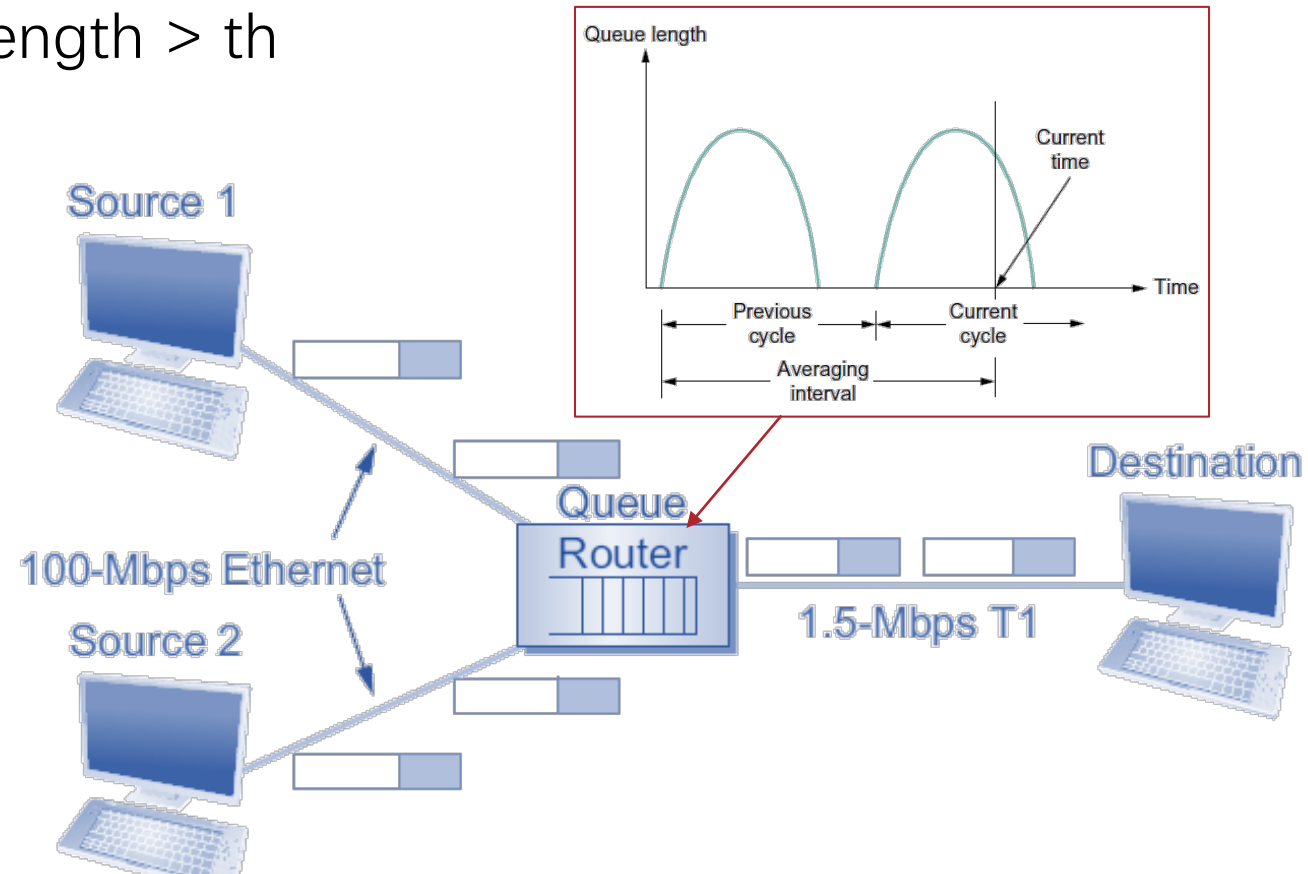


# DECbit

- Developed for the Digital Network Architecture (DNA)
  - Before TCP/IP was “standardized”
- Idea: let routers explicitly indicate congestion
- Approach:
  - Router set congestion bit in passing packets if there is congestion
  - Destination echoes bit back to source through acks
  - Source adjusts cwnd according to congestion bit

# DECbit

- Routers determine congestion
  - Monitor average queue length over last busy+idle cycle
  - If average queue length  $> th$ 
    - set congestion bit

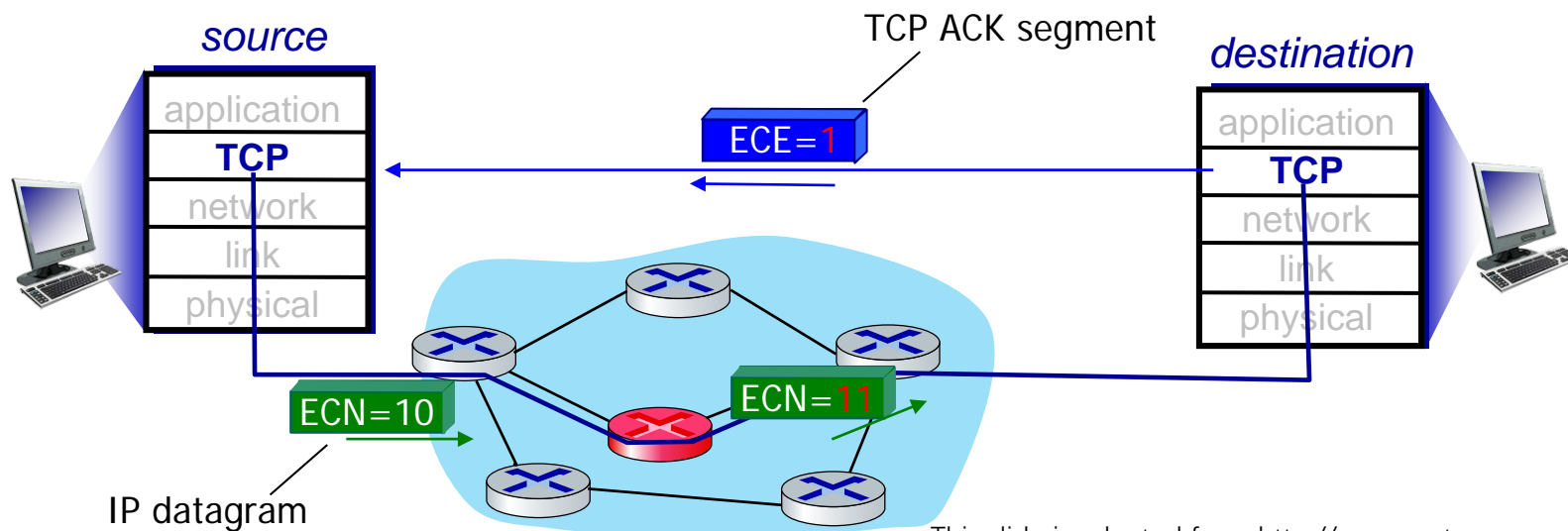


# DECbit

- Source reacts to congestion bit
  - If  $< 50\%$  of last window's packets had bit set
    - $\text{cwnd}++$
  - If  $> 50\%$  of last window's packets had bit set
    - $\text{cwnd} \times 0.875$

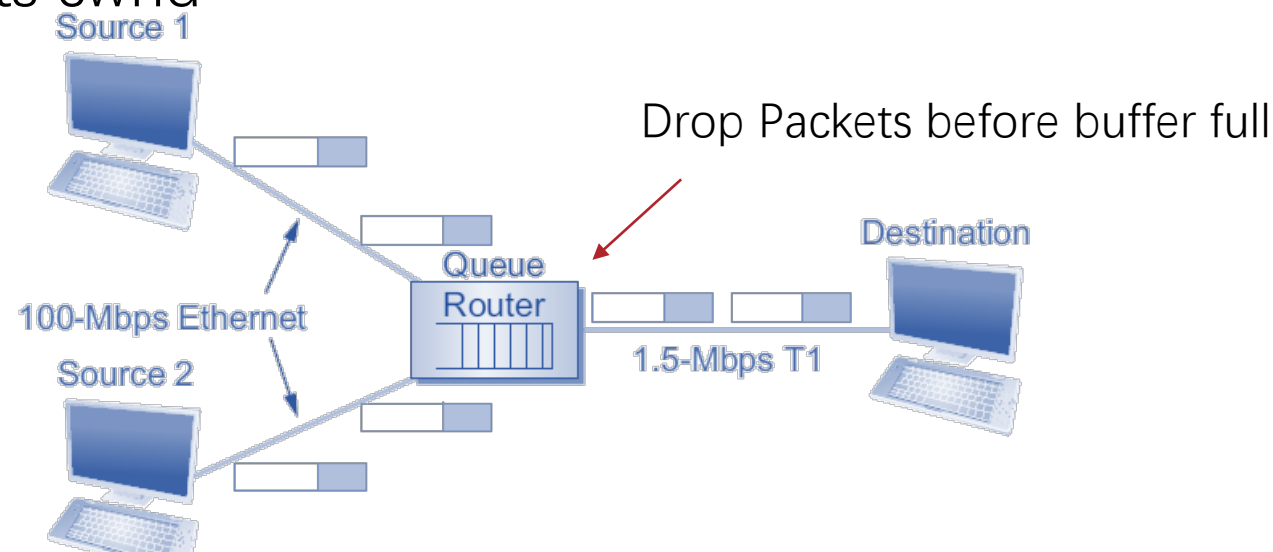
# Explicit congestion notification (ECN)

- RFC 3168
  - Two bits in IP header (ECN in IP's ToS field) marked by network router to indicate congestion
  - Mechanism is similar to DECbit
  - Destination sets ECE bit (in TCP Header) on ACK segment to notify sender of congestion



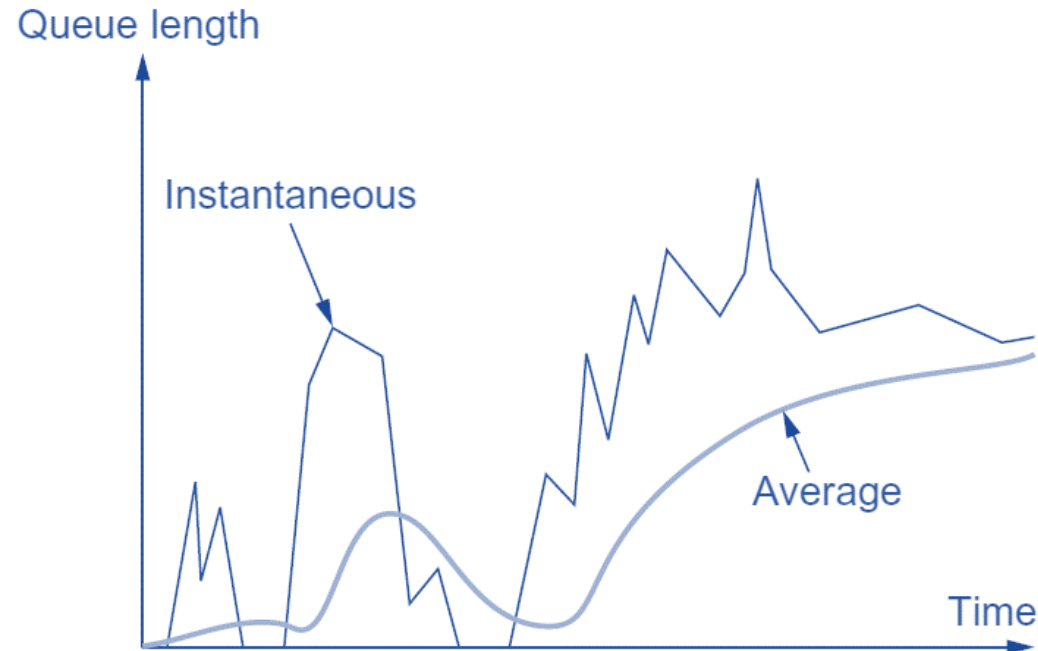
# Random Early Detection (RED)

- Another way to offload a part of congestion control function to routers
- Idea: let routers implicitly indicate congestion
  - Router notices that the queue is getting backlogged
  - Router randomly drops packets to signal congestion
  - Source adjusts cwnd



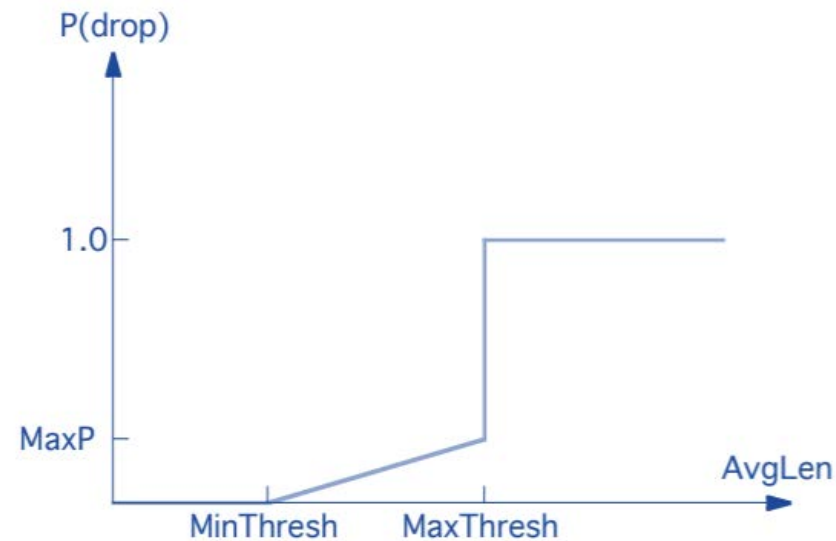
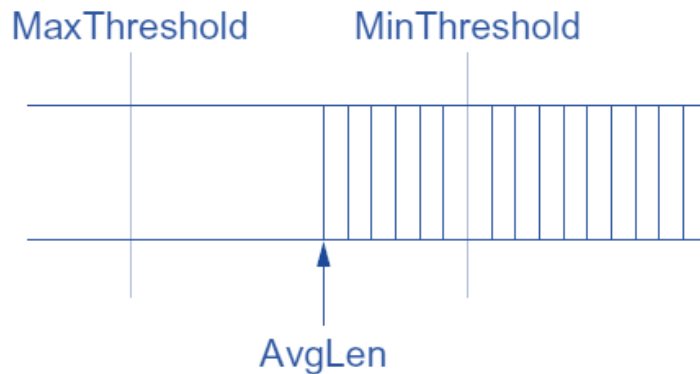
# RED Algorithm

- Compute Average Queue Length
  - Moving average
    - $\text{AvgLen} = (1 - \text{Weight}) * \text{AvgLen} + \text{Weight} * \text{SampleLen}$



# RED Algorithm

- Two queue length thresholds
  - if  $\text{AvgLen} \leq \text{MinThreshold}$ 
    - Enqueue the packet
  - if  $\text{MinThreshold} < \text{AvgLen} < \text{MaxThreshold}$ 
    - Drop arriving packet with probability  $P$
  - if  $\text{MaxThreshold} \leq \text{AvgLen}$  then drop arriving packet



# RED Algorithm

- Computing probability  $P$ 
  - $\text{TempP} = \text{MaxP} * (\text{AvgLen} - \text{MinThreshold}) / (\text{MaxThreshold} - \text{MinThreshold})$
  - $P = \text{TempP} / (1 - \text{count} * \text{TempP})$ 
    - Count: number of continuously queued packets without drop
- Why?
  - Drops are spaced out in time
    - $\text{Count} == 0 \Rightarrow P = \text{TempP}$
    - Count is large  $\Rightarrow P = 1$



# Reference

- Textbook 6.4