



CS120: Computer Networks

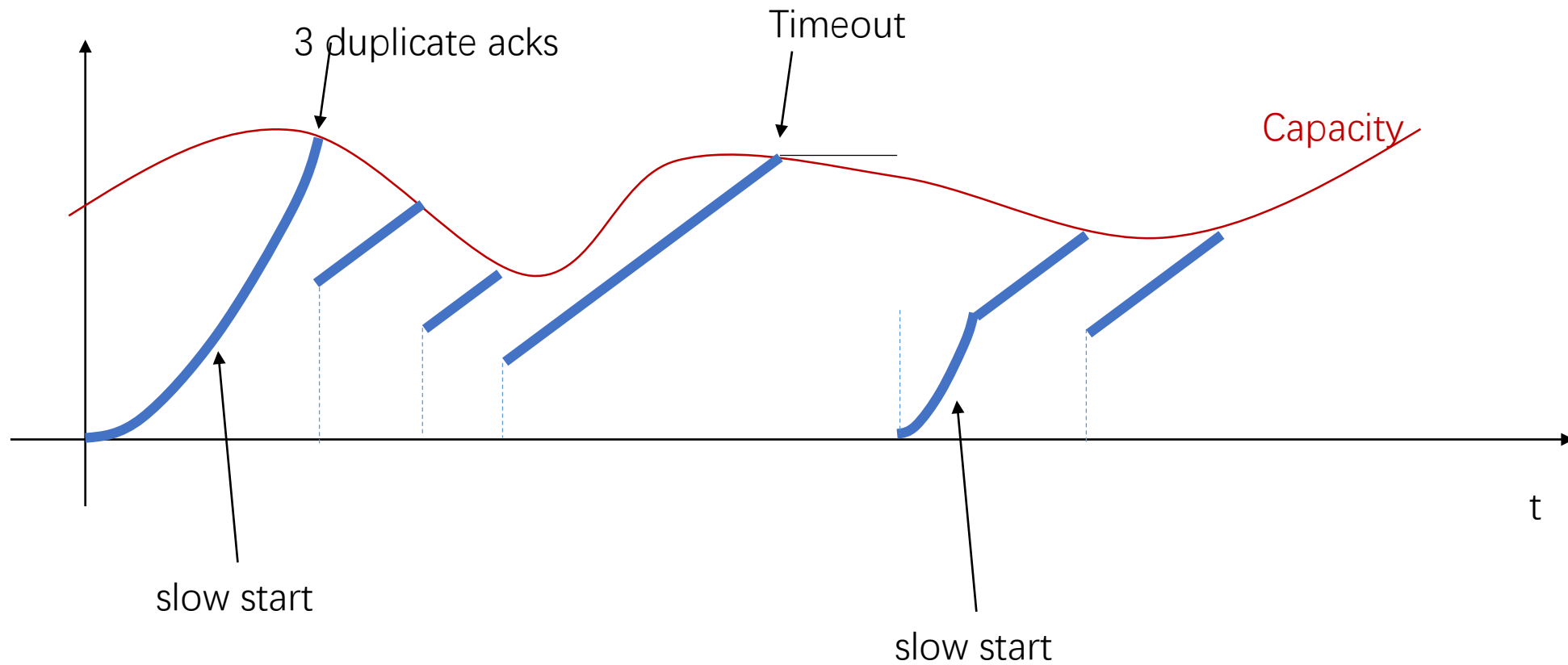
Lecture 19. Congestion Control 3

Zhice Yang

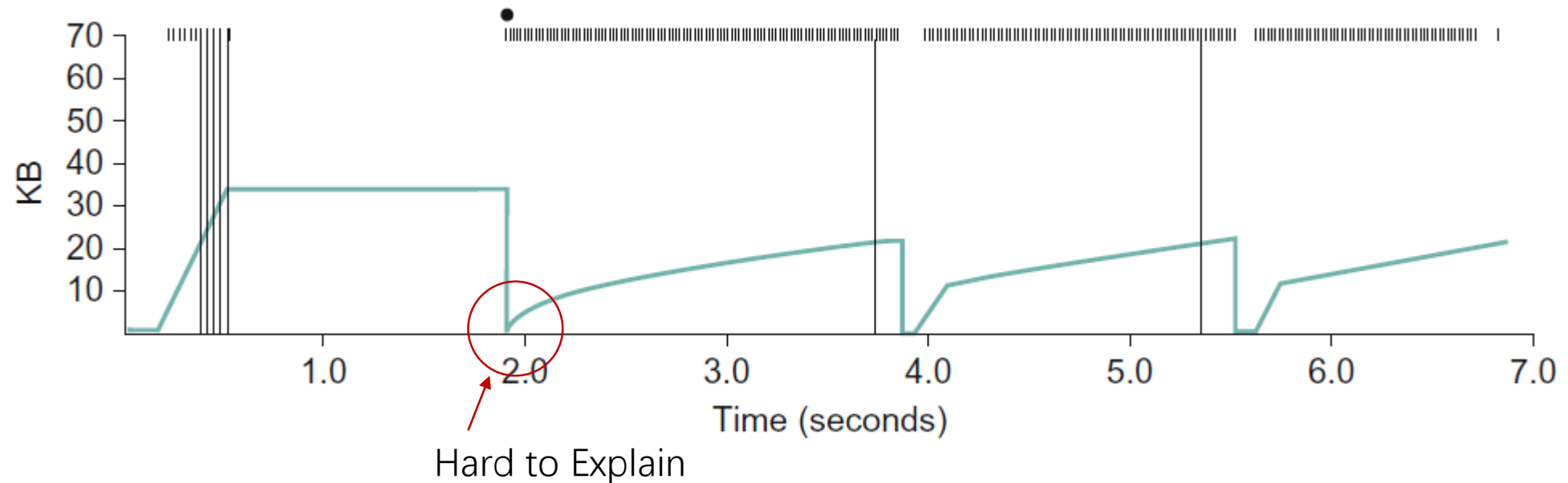
Congestion Control

- Queuing
- Connection Control Methods
 - Congestion Control
 - Congestion Avoidance
- QoS

TCP Reno



Figures in Textbook



No Fast Recovery in the Simulation, i.e., TCP Tahoe

TCP Congestion Control

- Objective: Estimate and adapt to (varying) network capacity
- Approach: Adjust Sliding Window
 - $\text{MaxWindow} = \text{MIN}(\text{CongestionWindow}, \text{AdvertisedWindow})$
 - Decrease **CongestionWindow** upon detecting congestion
 - Increase **CongestionWindow** upon lack of congestion
- Basic Components
 - Additive Increase/Multiplicative Decrease (AIMD)
 - Slow Start
 - Fast Retransmission
 - Fast Recovery
- Other Variations

TCP Congestion Control Algorithms

- ref: https://en.wikipedia.org/wiki/TCP_congestion_control

Variant	Feedback	Required changes	Benefits	Fairness
(New)Reno	Loss	-	-	Delay
Vegas	Delay	Sender	Less loss	Proportional
High Speed	Loss	Sender	High bandwidth	
BIC	Loss	Sender	High bandwidth	
CUBIC	Loss	Sender	High bandwidth	
H-TCP	Loss	Sender	High bandwidth	
FAST	Delay	Sender	High bandwidth	Proportional
Compound TCP	Loss/Delay	Sender	High bandwidth	Proportional
Westwood	Loss/Delay	Sender	L	
Jersey	Loss/Delay	Sender	L	
BBR ^[11]	Delay	Sender	BLVC, Bufferbloat	
CLAMP	Multi-bit signal	Receiver, Routers	V	Max-min
TFRC	Loss	Sender, Receiver	No Retransmission	Minimum delay
XCP	Multi-bit signal	Sender, Receiver, Router	BLFC	Max-min
VCP	2-bit signal	Sender, Receiver, Router	BLF	Proportional
MaxNet	Multi-bit signal	Sender, Receiver, Router	BLFSC	Max-min
JetMax	Multi-bit signal	Sender, Receiver, Router	High bandwidth	Max-min
RFD	Loss	Router	Smaller delay	

Demo

- Sliding Window code in TCP
/net/ipv4/
<https://elixir.bootlin.com/linux/latest/source/net/ipv4>
- Change Sliding Window
 - Show current congestion control scheme
`cat /proc/sys/net/ipv4/tcp_congestion_control`
 - Show/change available congestion control scheme
`sysctl net.ipv4.tcp_available_congestion_control[=XX]`

Congestion Control

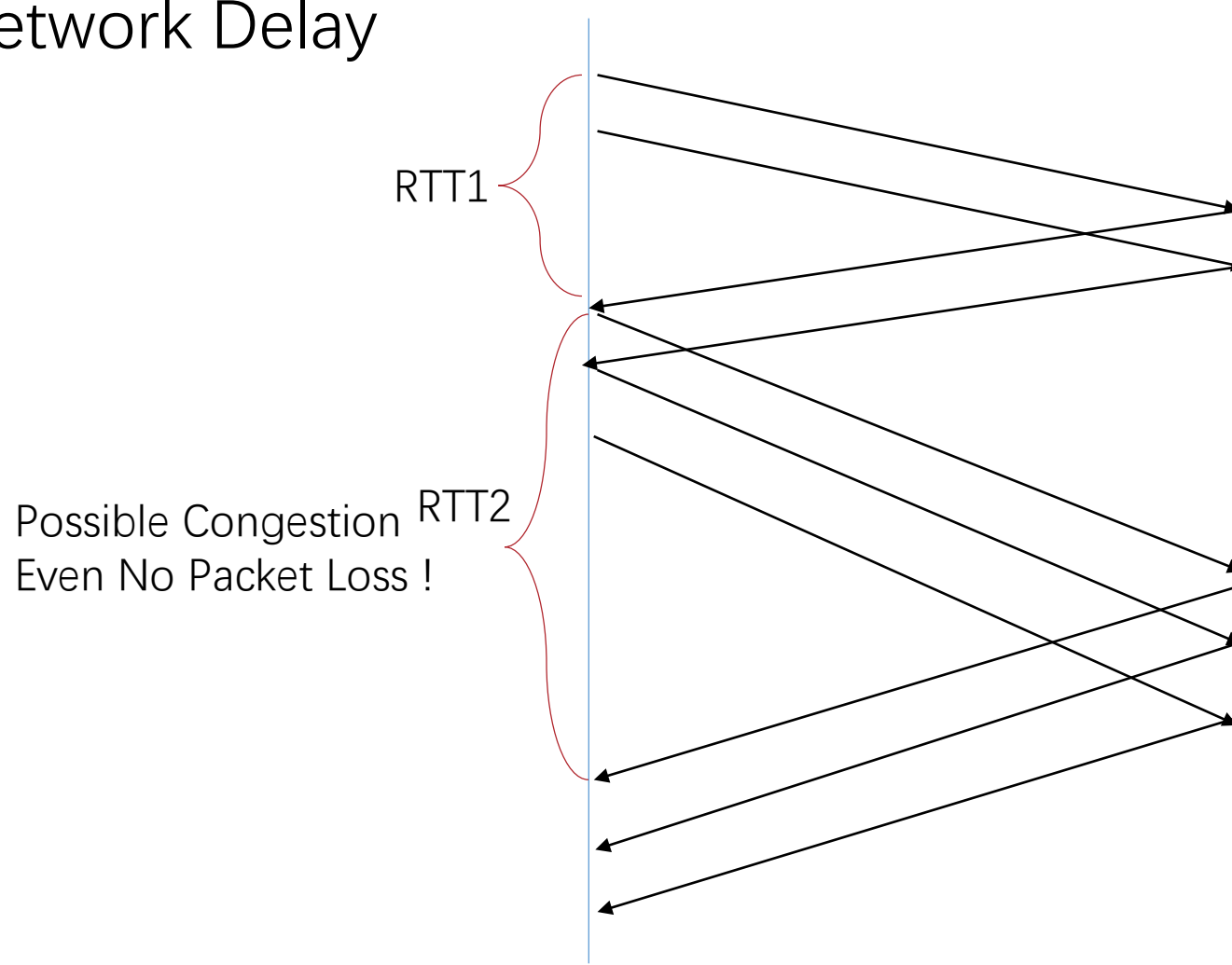
- Queuing
- Connection Control Methods
 - Congestion Control
 - Congestion Avoidance
- QoS

Packet Loss v.s. Network Delay

- Packet Loss
 - Packet loss is the indication of congestion
 - When packet loss happens
 - Many arriving packets encounter a full queue
 - Many individual flows lose multiple packets
 - Many flows divide sending rate in half
- Network Delay
 - Increase in network delay is an indicator for possible congestion
 - When network delay increases
 - Some packets are queued in routers
 - Slow down one or two flows before congestion happens

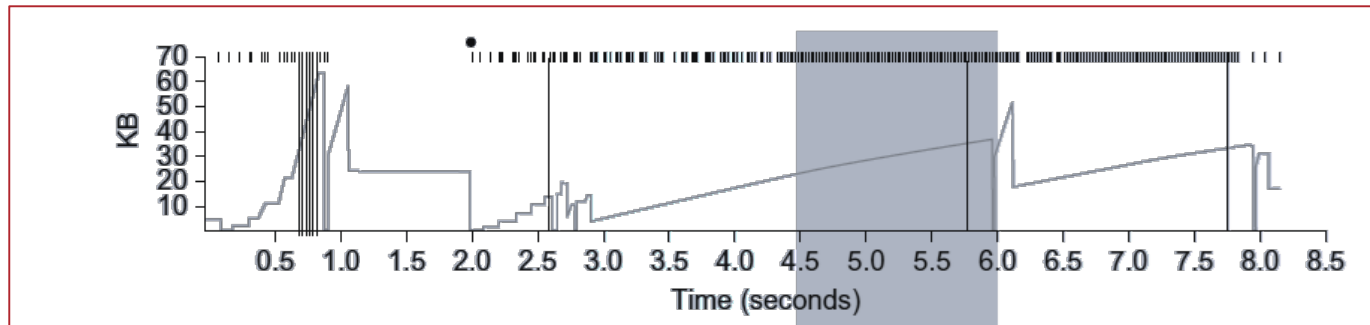
Avoid Congestion

- Feedback: Network Delay

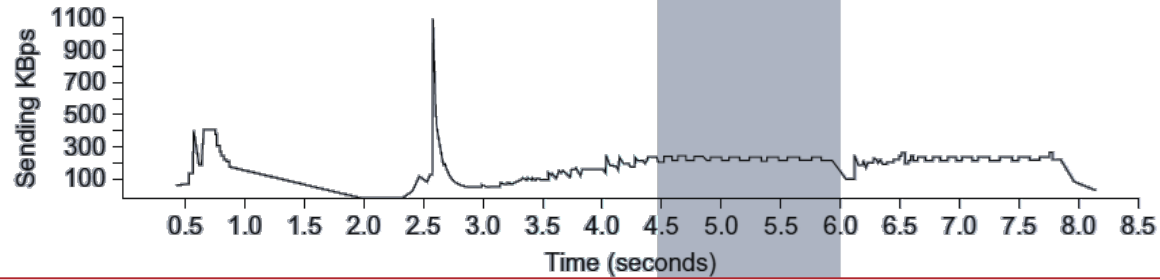


Avoid Congestion

- Feedback: Network Delay
 - Possible Designs
 - Make decision based on current RTT and average RTT
 - Make decision based on current {RTT, Window} and average {RTT, Window}
 - Make decision based rate trend
 - etc...



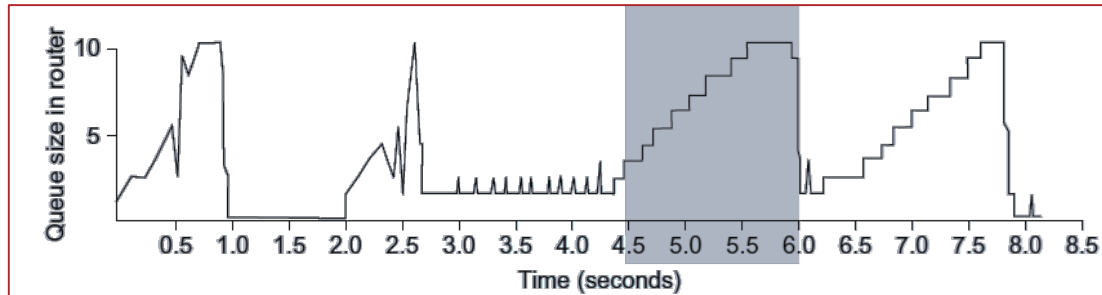
cwnd



Source 1

throughput

throughput



Queue Size

100-Mbps Ethernet

Destination

Source 2

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

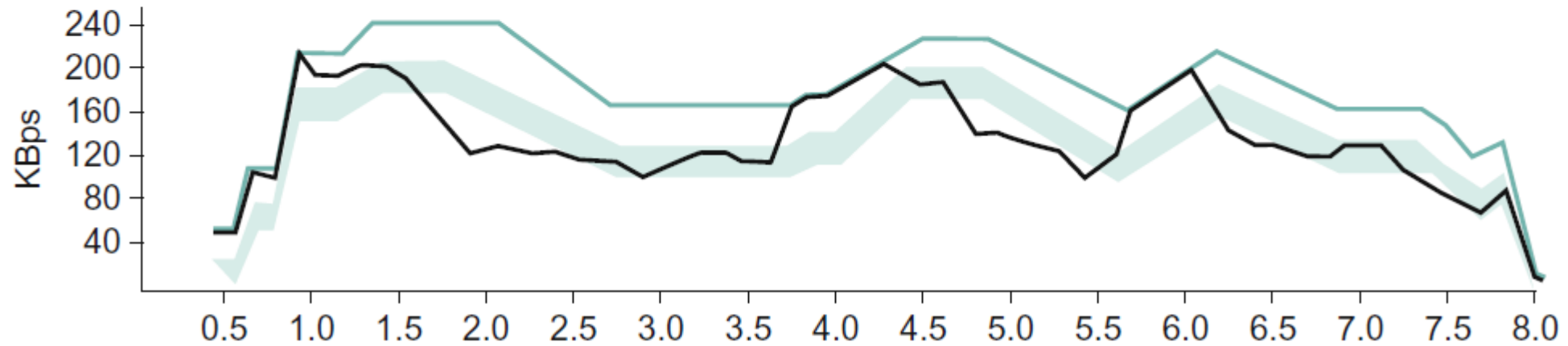
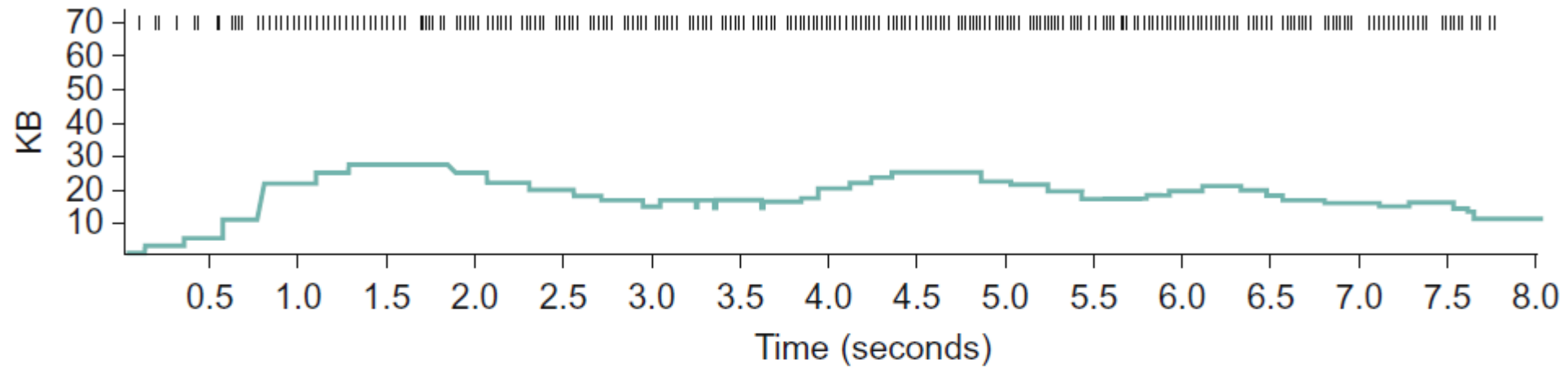
TCP Vegas Algorithm

- 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

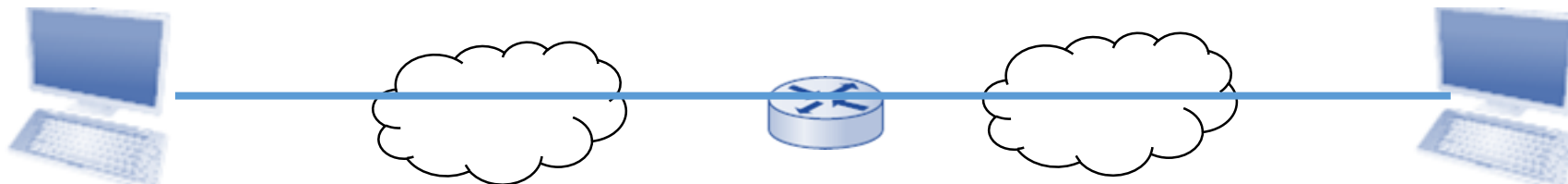


Congestion Control

- Queuing
- Connection Control Methods
 - Congestion Control
 - Congestion Avoidance
- QoS

Congestion Control and Resource Allocation

- Two Sides of the Same Coin
 - Control congestion if (and when) it occurs: reactive
 - Pre-allocate resources so as to avoid congestion: proactive
- Resources in Network
 - Bandwidth
 - Router Queue Buffer
- Two Places of Implementation
 - Hosts at the edges of the network (transport protocol)
 - Routers inside the network (queuing discipline)



Congestion Avoidance with Routers

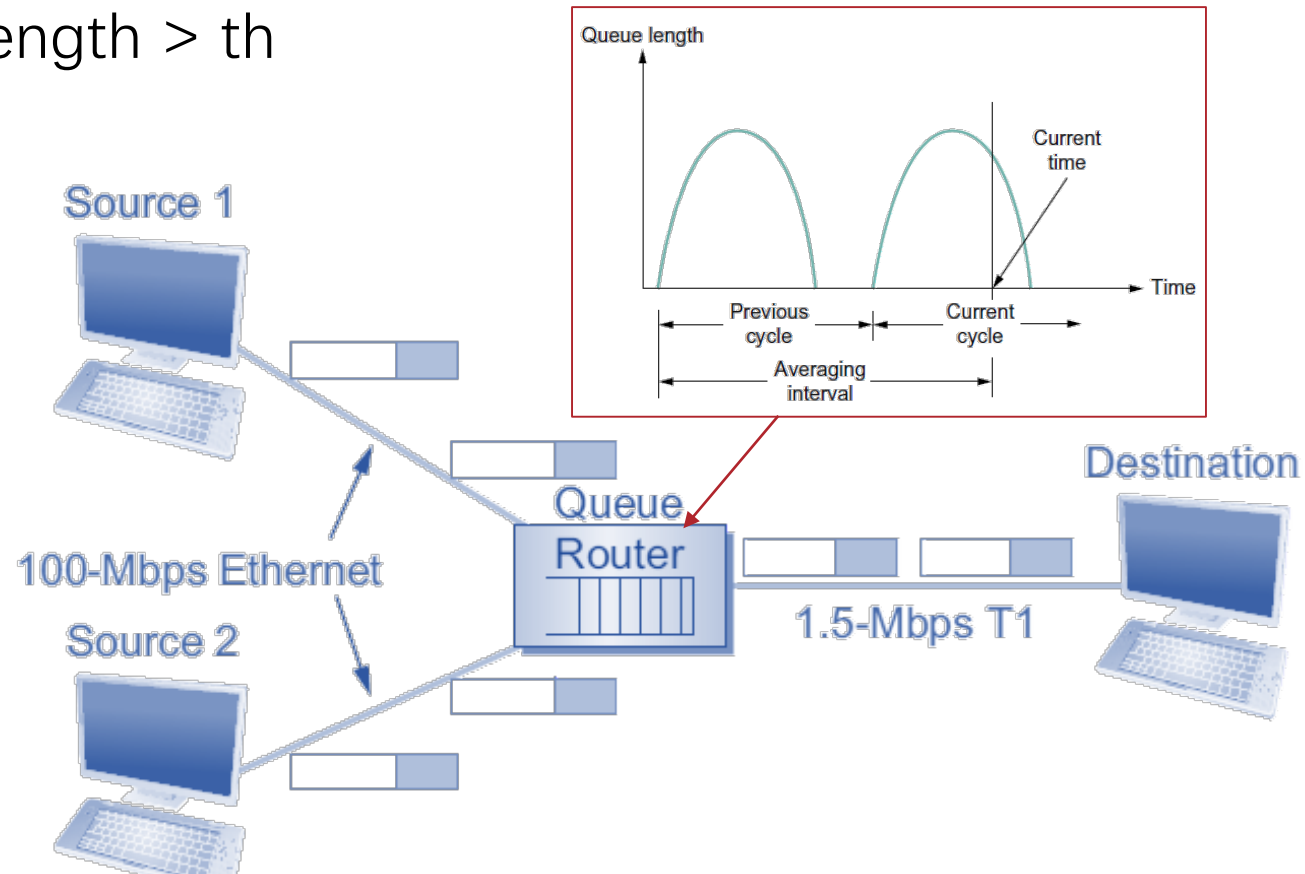
- DECbit
- Random Early Detection (RED)

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

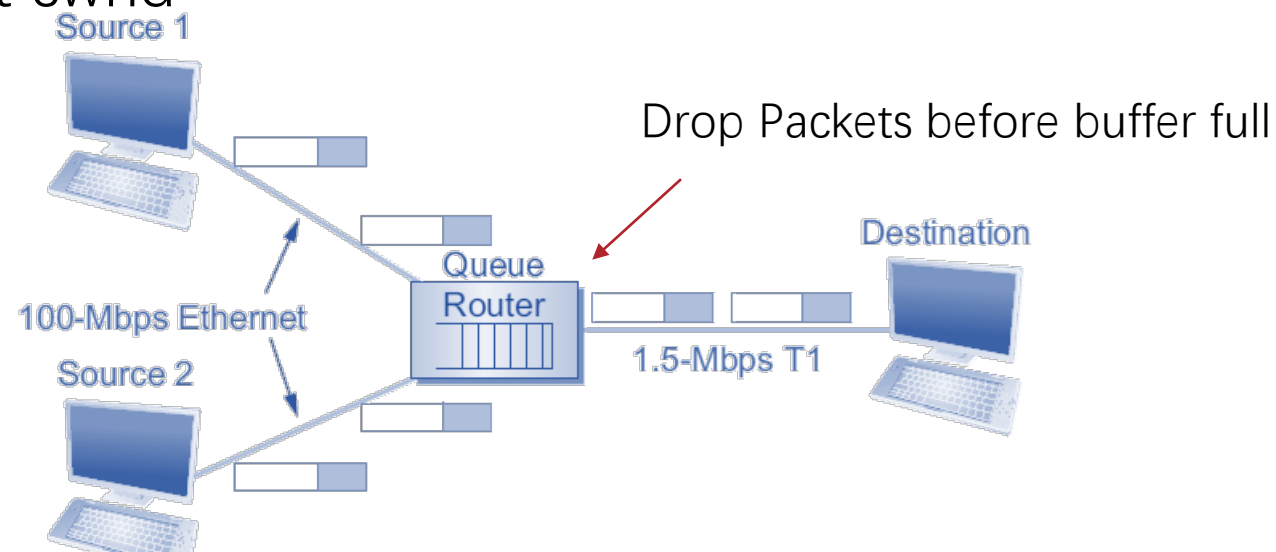


DECbit

- 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
 - $\text{cwnd} \times 0.875$

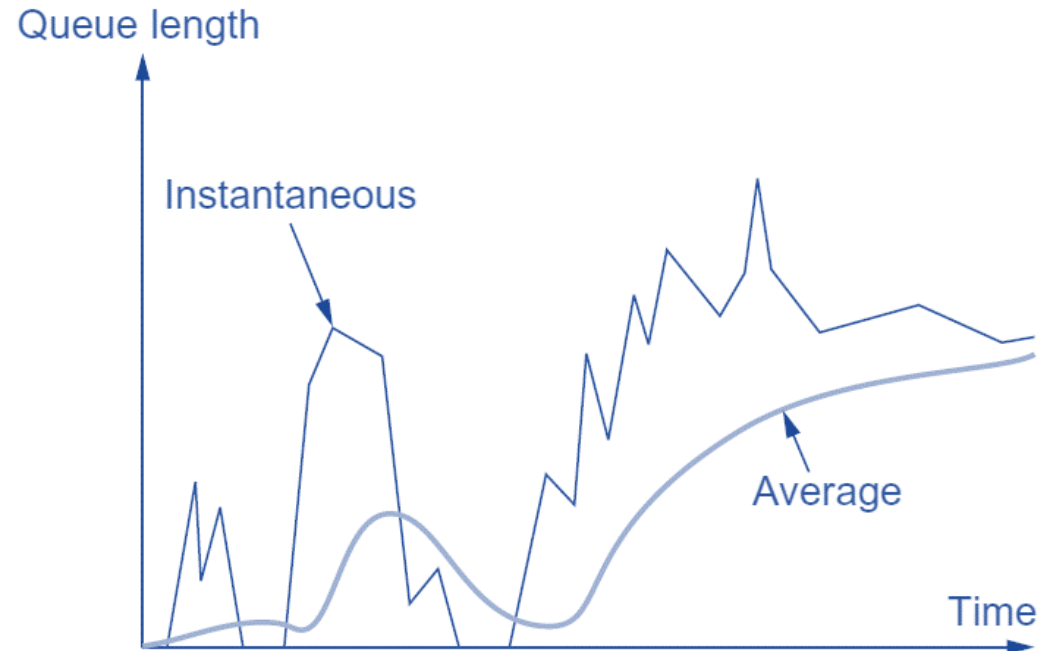
Random Early Detection (RED)

- Developed for TCP/IP, 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 adjust cwnd



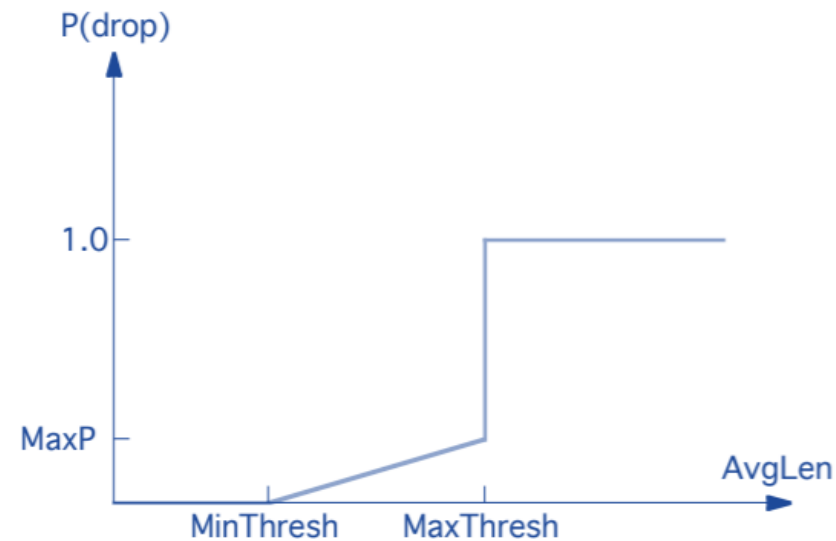
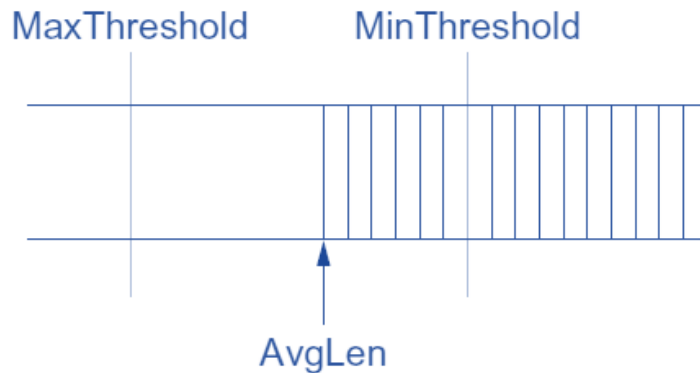
RED Algorithm

- 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