

# Transport

# Review: slow start, additive inc

Say MSS = 1 KByte

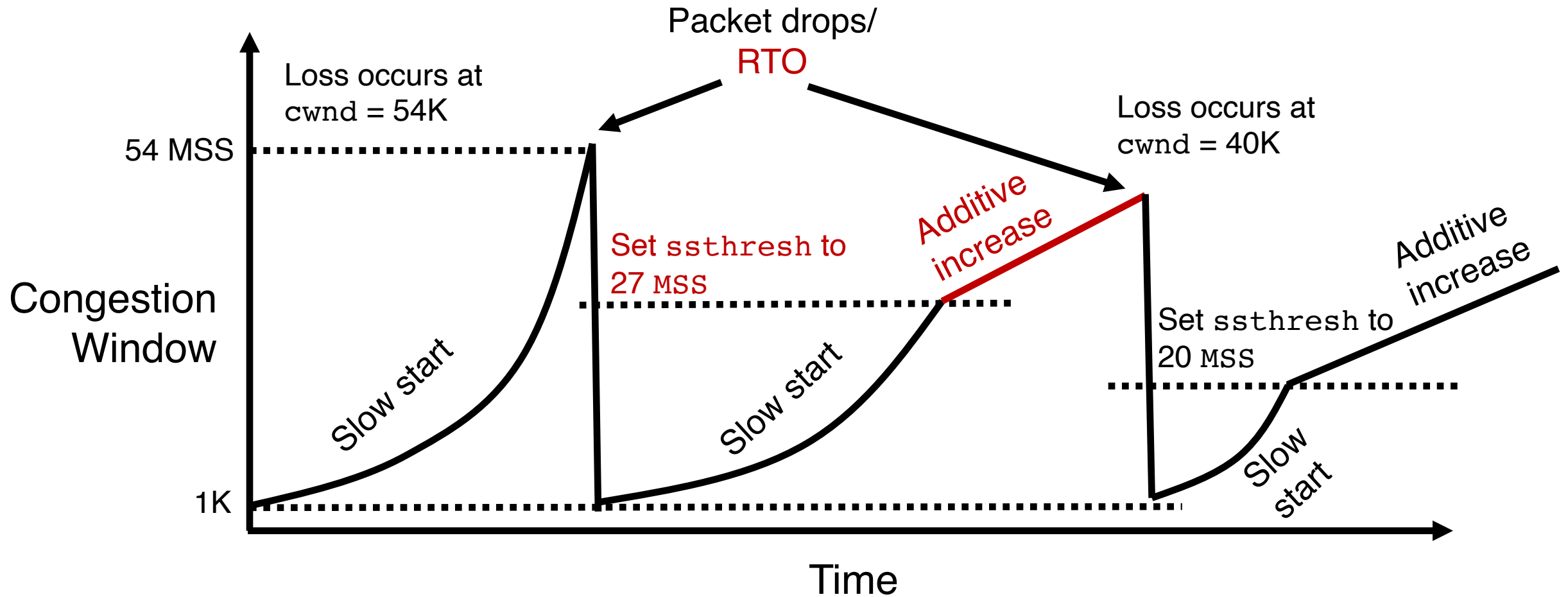
Default ssthresh = 64KB = 64 MSS

AI is slow.

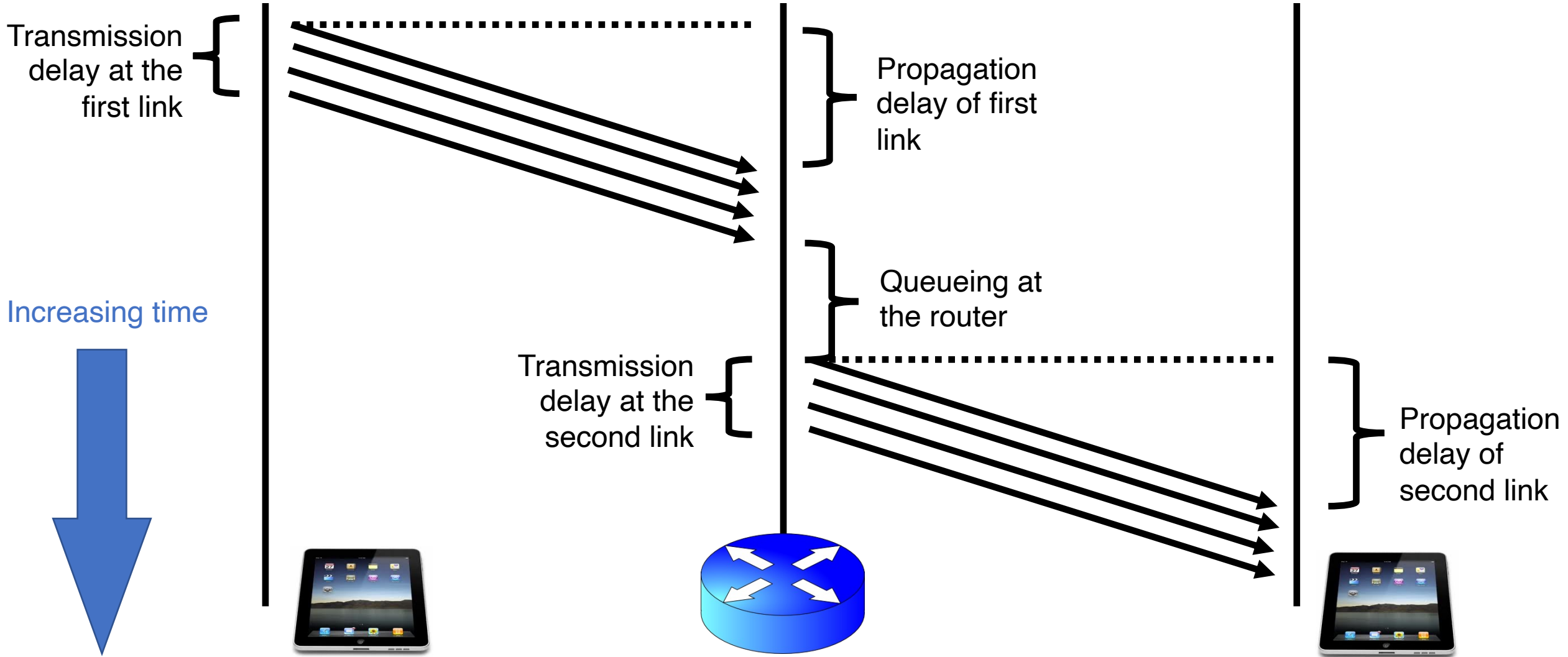
Persistent connections

Large window sizes

Different laws to evolve congestion window



# The components of delay



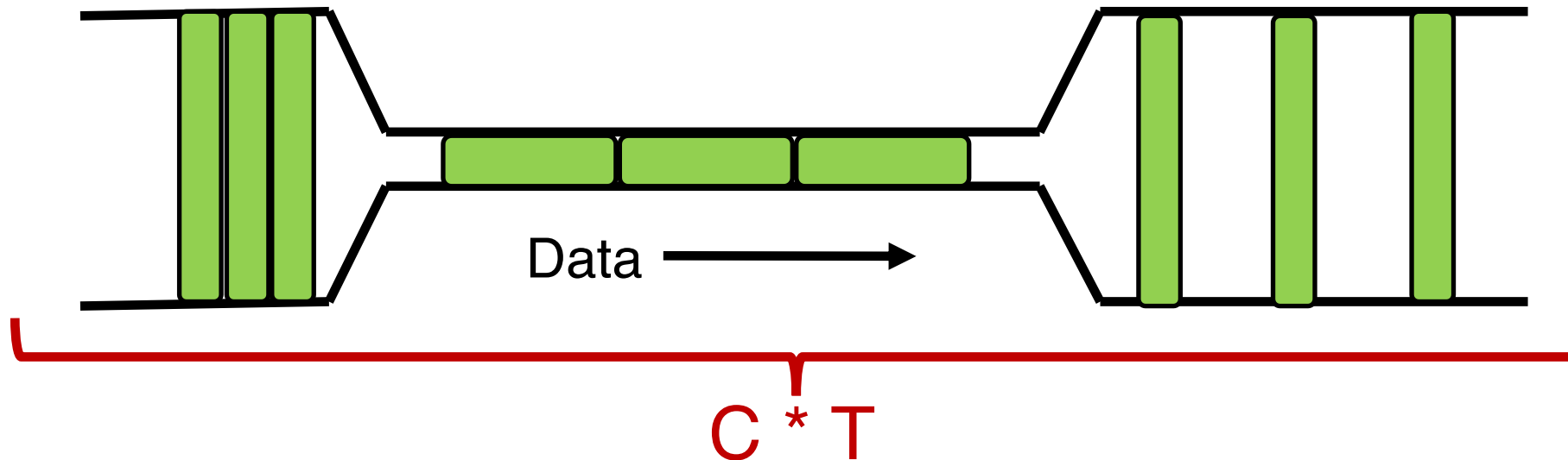
# Bandwidth-Delay Product

# Steady state cwnd for a single flow

- Suppose the bottleneck link has rate  $C$
- Suppose the propagation round-trip delay (propRTT) between sender and receiver is  $T$
- Ignore transmission delays for this example;
- Assume steady state: highest sending rate with no bottleneck congestion
- Q: how much data is in flight over a single RTT?
- $C * T$  data i.e., amount of data unACKed at any point in time
- ACKs take time  $T$  to arrive (without any queueing). In the meantime, sender is transmitting at rate  $C$

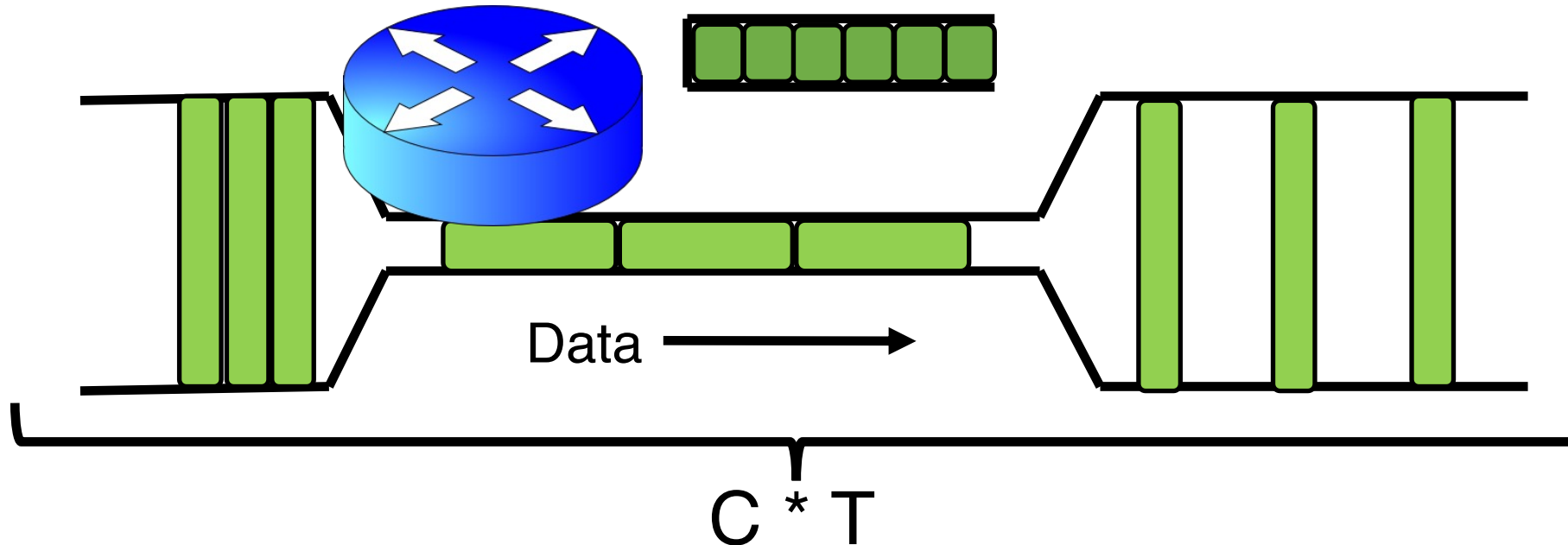
# The Bandwidth-Delay Product

- $C * T$  = **bandwidth-delay product**:
  - The amount of data in flight for a sender transmitting at the ideal rate during the ideal round-trip delay of a packet
- Note: this is just the amount of data “on the pipe”



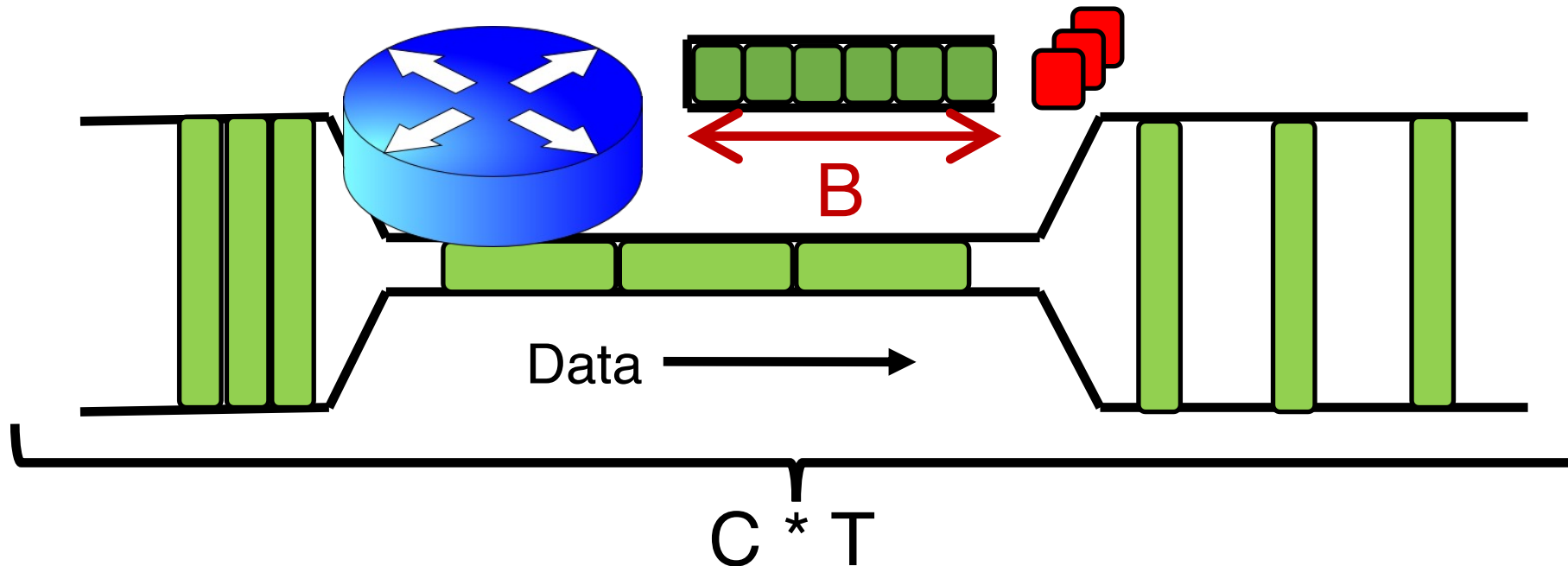
# The Bandwidth-Delay Product

- Q: **What happens if  $cwnd > C * T$ ?**
  - i.e., where are the rest of the in-flight packets?
- A: **Waiting at the bottleneck router queues**



# Router buffers and the max cwnd

- Router buffer memory is finite: queues can only be so long
  - If the router buffer size is  $B$ , there is at most  $B$  data waiting in the queue
- If cwnd increases beyond  $C * T + B$ , data is dropped!





# BDP is a crucial value for a flow

- Bandwidth-Delay Product (BDP) governs the window size of a single flow at steady state
- The bottleneck router buffer size governs how much the cwnd can exceed the BDP before packet drops occur
- BDP is the ideal desired window size to use the full bottleneck link, without any queueing.
  - Accommodating **flow control**, also the min socket buffer size to use the bottleneck link fully

# Demo of the impact of BDP & B

- Utilization
- Congestion window

# Detecting and Reacting Better to Packet Loss

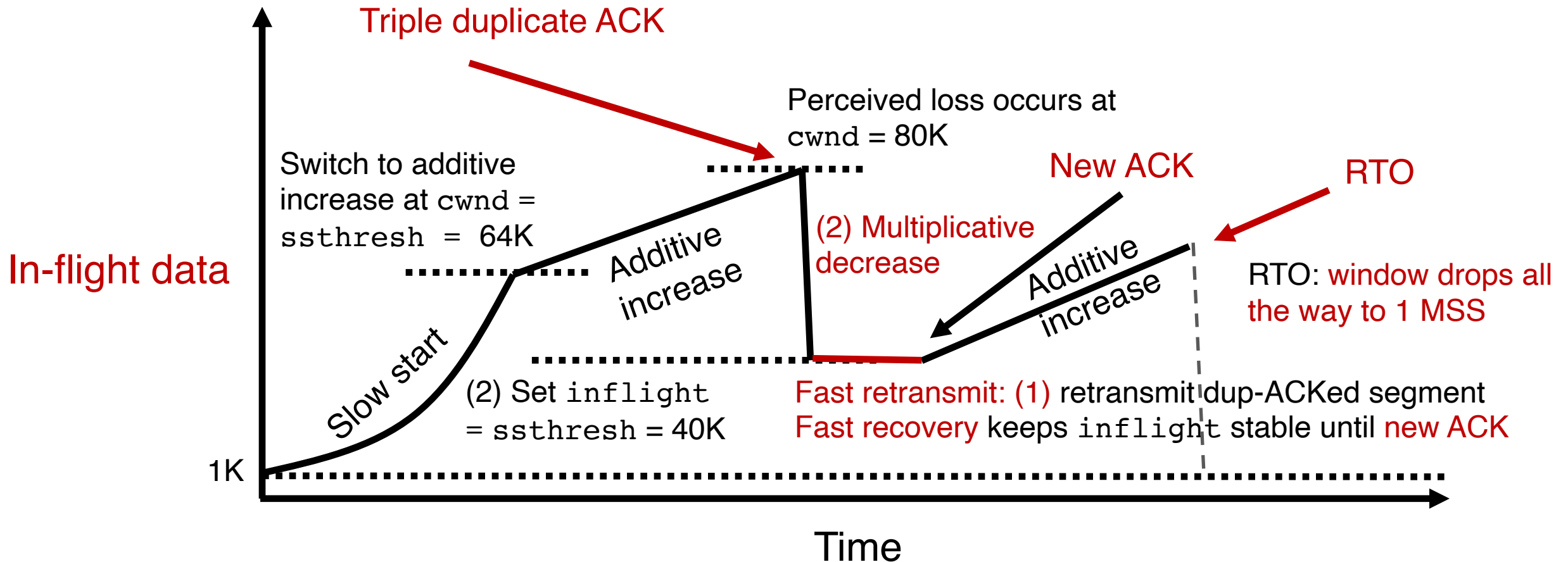
# Can we detect loss earlier than RTO?

- Key idea: use the information in the ACKs. **How?**
- Suppose successive (cumulative) ACKs contain the same ACK#
  - Also called **duplicate ACKs**
  - Occur when network is reordering packets, or one (but not most) packets in the window were lost
  - **Fast retransmit:** (1) Immediately retransmit packet
- Reduce cwnd when you see many duplicate ACKs
  - Consider many dup ACKs a strong indication that packet was lost
  - Default threshold: 3 dup ACKs, i.e., **triple duplicate ACK**
  - Make cwnd reduction gentler than setting cwnd = 1; recover faster
  - Fast retransmit: **(2) reduce window to half of its current value**

# Additive Increase/Multiplicative Decrease

Say MSS = 1 KByte

Default ssthresh = 64KB = 64 MSS



TCP **New Reno** performs additive increase and multiplicative decrease of congestion window.

In short, we often refer to this as **AIMD**.

**Multiplicative decrease** is a part of all TCP algorithms. It is necessary for **fairness** across TCP flows.

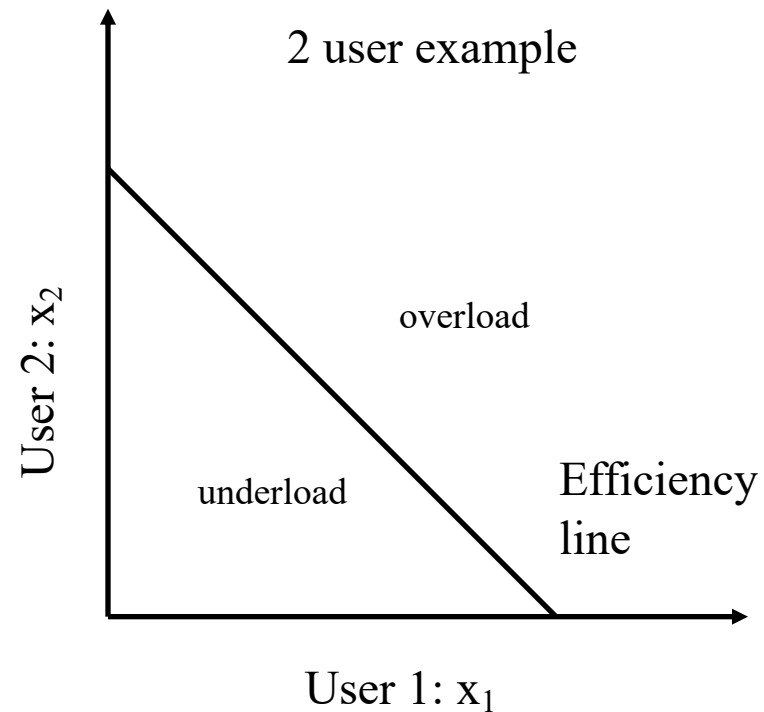
# Why does multiplicative decrease help?

Efficiency and Fairness

Chiu and Jain, “Increase and decrease algorithms for congestion avoidance”

# Efficient allocation

- Don't want sources to transmit either too slow or too fast
  - Slow: Underutilize the network
  - Fast: High delays, lose packets
- Every endpoint is reacting
  - May *all* under/overshoot
  - Large oscillations possible!
- Optimal efficiency:
  - $\sum x_i = X_{goal}$  e.g., link capacity
- Efficiency = 1 - distance from efficiency line



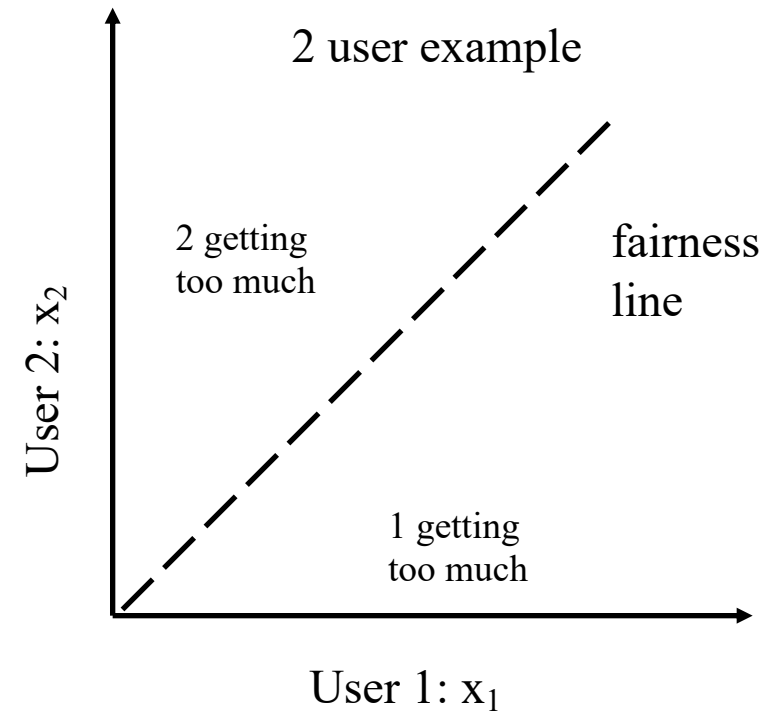


# Fair allocation

- **Max-min fairness**
- Flows which share the same bottleneck get the same amount of bandwidth

$$F(x) = \frac{\left(\sum x_i\right)^2}{n\left(\sum x_i^2\right)}$$

- Fairness = 1 - distance from fairness line



# How should transports react?

- Given efficiency and fairness goals above, how should transports behave?
- Consider  $x(t)$ , window or rate of a source, evolving over time  $t$
- Assume discrete time steps.
- $x(t + 1) =$  function of  $x(t)$ , feedback from the network

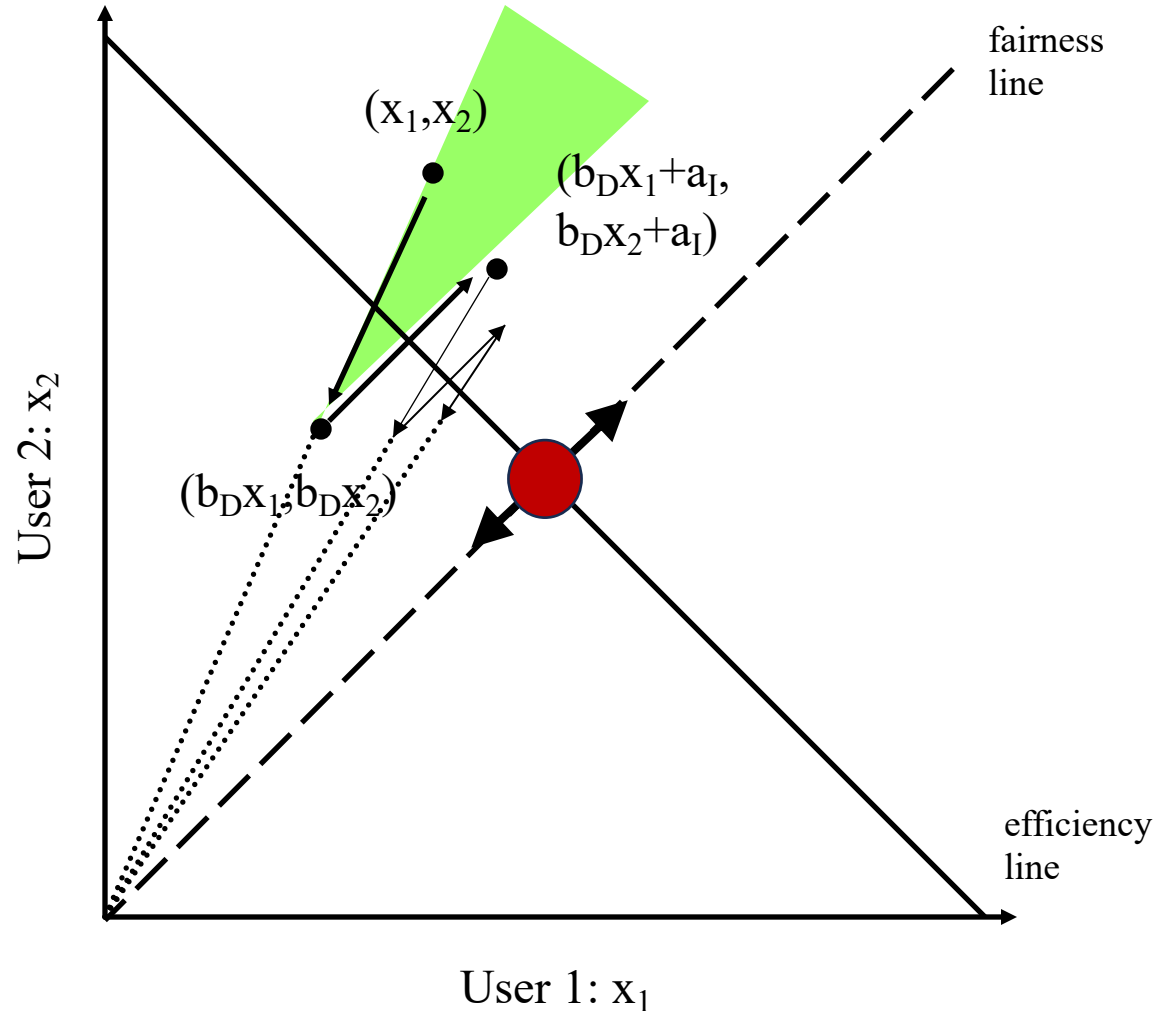
# Linear control rules

$$x_i(t+1) = \begin{cases} a_I + b_I x_i(t) & \text{increase} \\ a_D + b_D x_i(t) & \text{decrease} \end{cases}$$

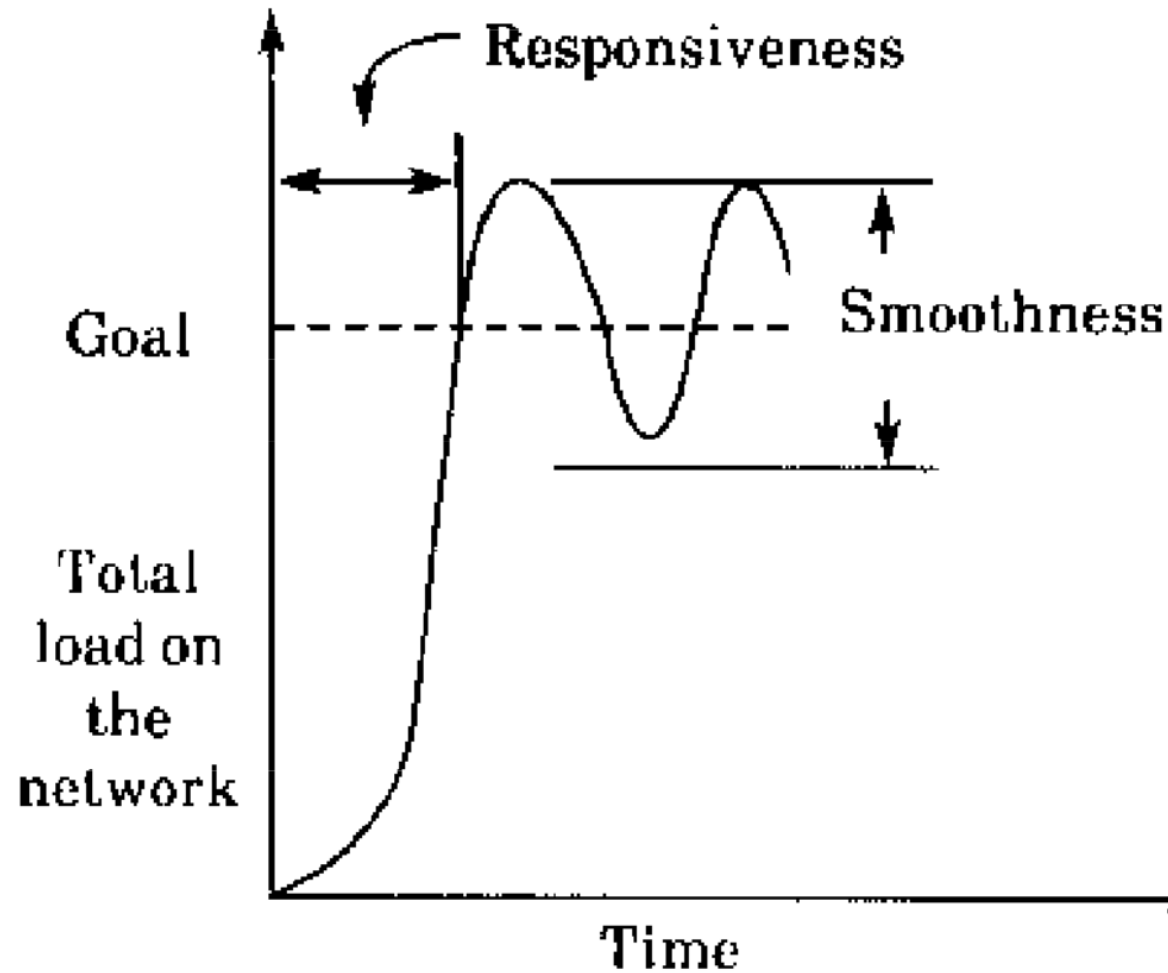
- $x_i(t)$ : window or rate of the  $i^{\text{th}}$  user at time  $t$
- $a_I, a_D, b_I, b_D$ : constant increase/decrease coefficients
- **Assumption: All users receive same network feedback**
  - *Binary* feedback: sense congestion or available capacity
- **Assumption: All users increase or decrease simultaneously**

# Additive increase, multiplicative decrease

- $b_I = 1, a_D = 0$
- Multiplicative decrease enables converging to fairness
- Oscillates around the most efficient point



# Convergent doesn't mean static



# Simple models are useful

- Chiu and Jain's model isn't indicative of all TCP/AIMD behavior
  - But it's "realistic" enough
  - Stands the test of time: many more sources, much higher bandwidth, ...
- Models should be simple
  - For us to work with
  - For others to understand
  - But they don't have to mimic the "real" thing in every way
- A real, complex model is likely useless, but a realistic, simple model might teach us something

# Modeling TCP throughput

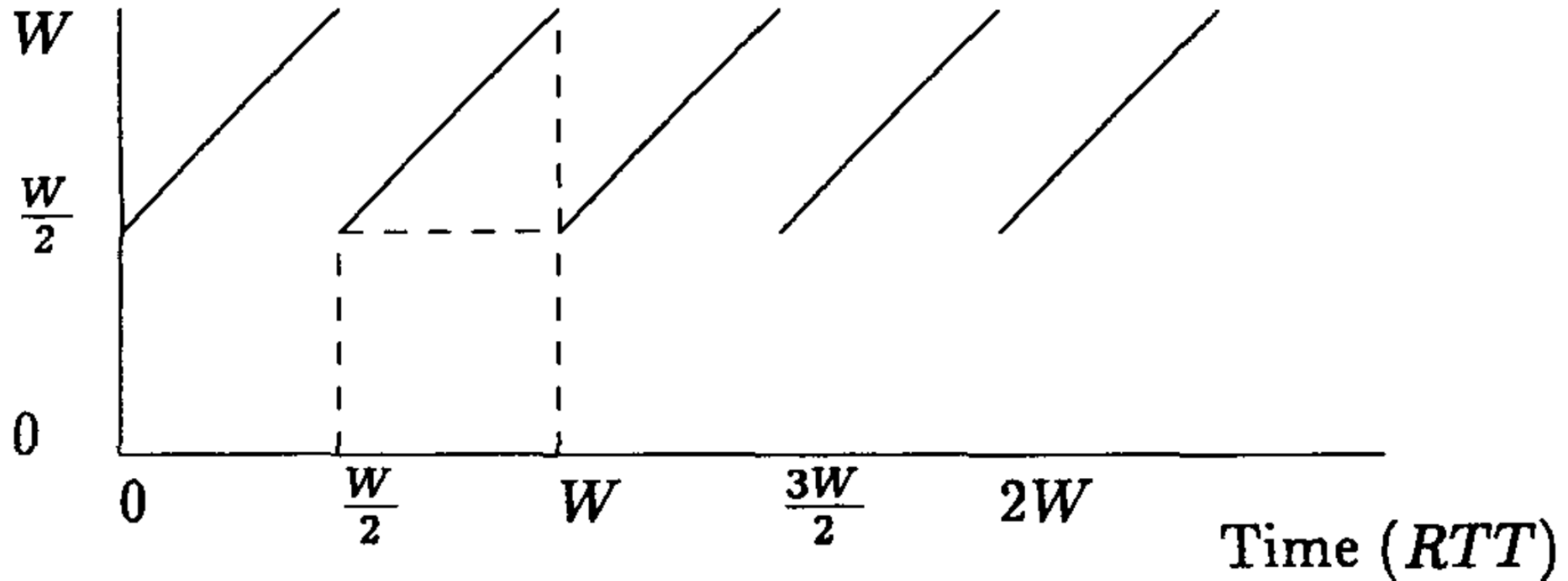
Given network characteristics, how quickly can TCP (New Reno) send data?

Mathis et al., “Macroscopic behavior of TCP congestion avoidance”

# Steady AIMD

$$BW < \left( \frac{MSS}{RTT} \right) \frac{1}{\sqrt{p}}$$

congestion window (packets)





# Implications

- Throughput has a  $1/\sqrt{p}$  dependence on packet loss rate
- RTT unfairness
  - Flows with a smaller RTT get better throughput (ramp up faster)
- Engineering implications:
  - Split TCP (CDNs, data center frontends, ...)
  - Special considerations for long-distance connections

# Widely Deployed TCPs

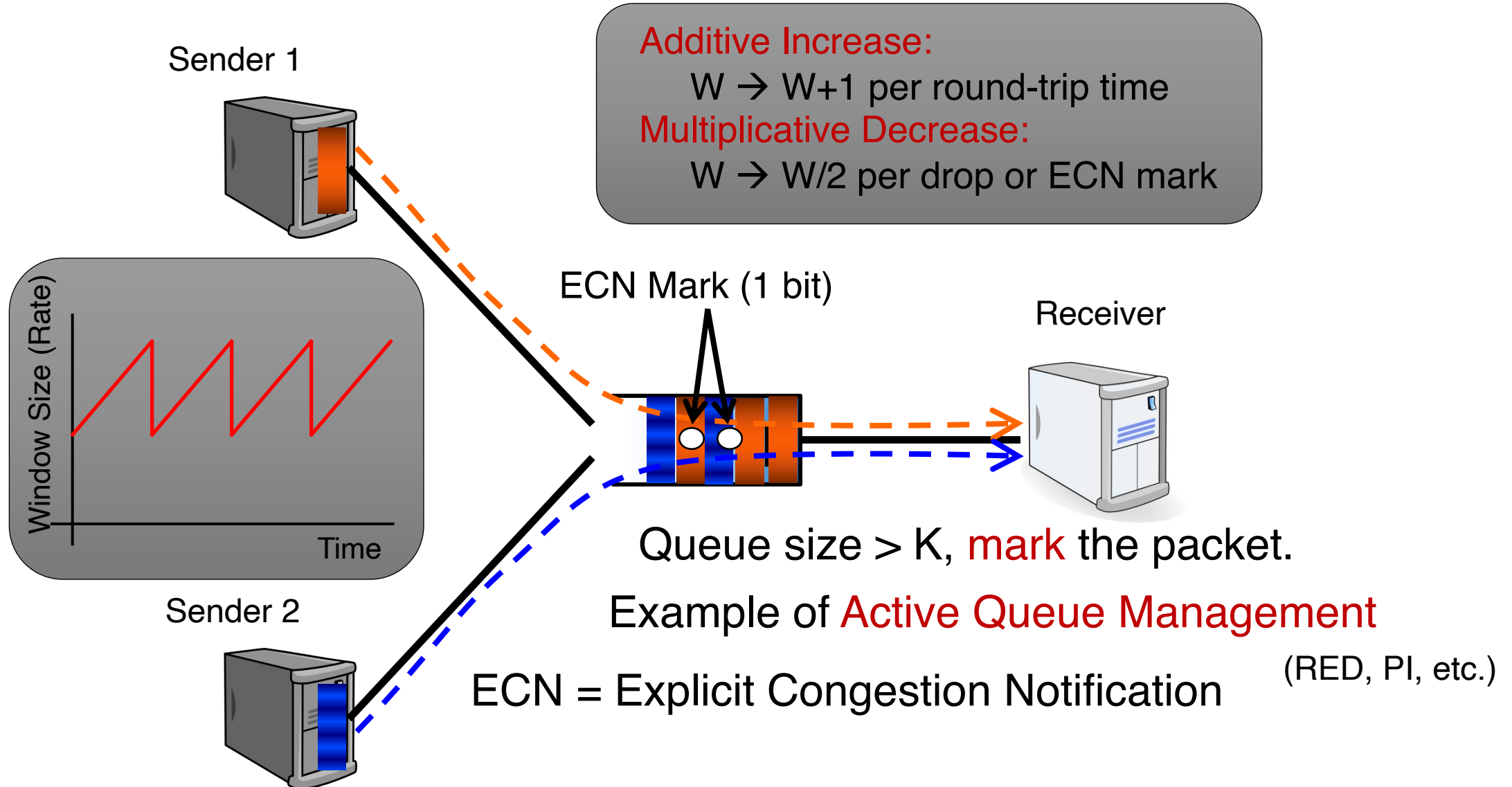
# Data Center TCP

Alizadeh et al.


# Context

- Regular TCP: window evolution with all signals measured end to end
- Data centers: hardware under single administrative control
  - Could network switches do better?
- What if switches provided better feedback than loss?

# Explicit Congestion Notification



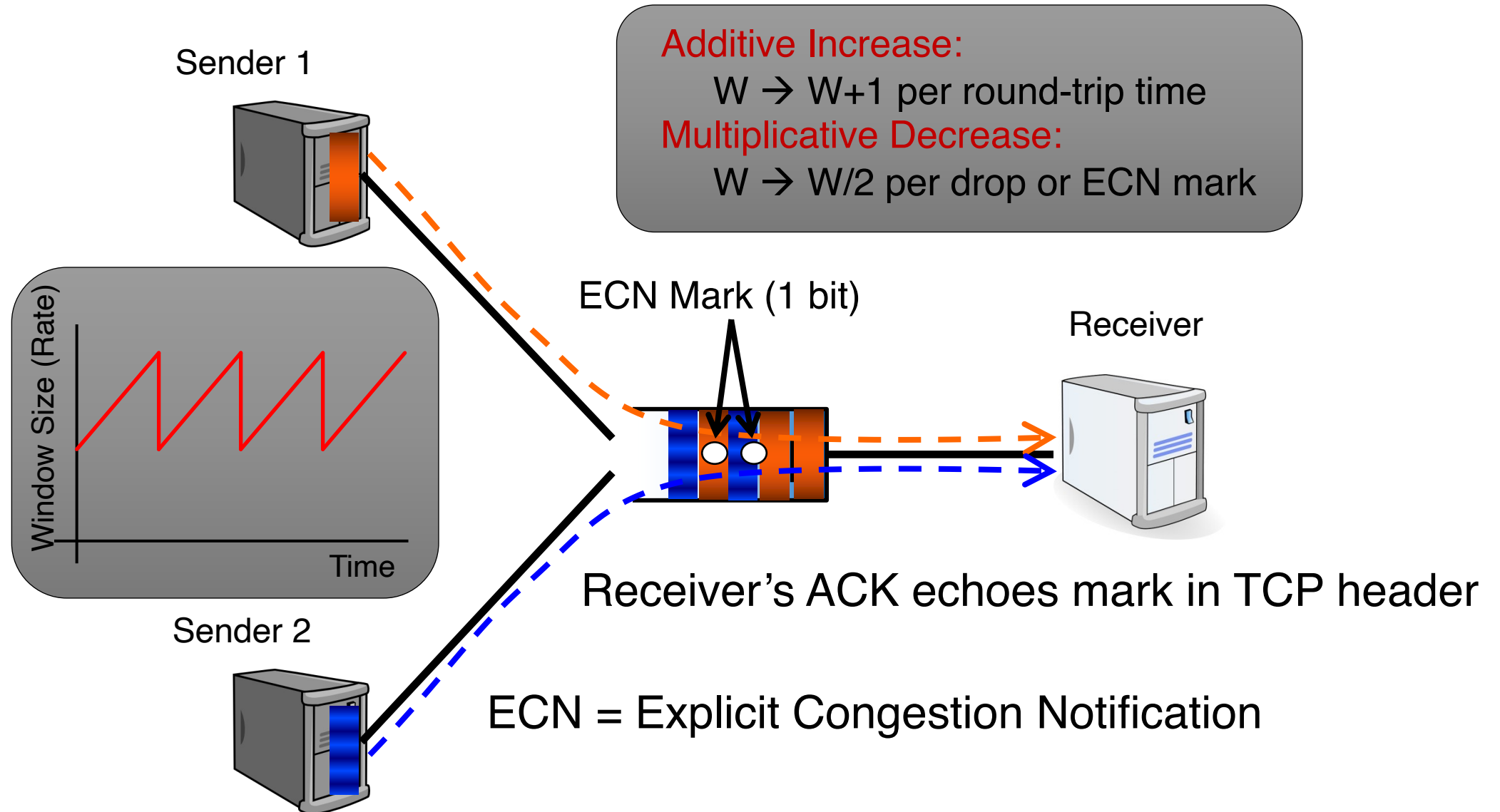
# ECN set on the IP header by routers

- 00 – Not ECN-Capable Transport, Not-ECT
- 01 – ECN Capable Transport(1), ECT(1)
- 10 – ECN Capable Transport(0), ECT(0)
- 11 – Congestion Experienced, CE.  **Dropped** if TCP sender is not ECN enabled

## IPv4 header format

Offsets	Octet	0								1								2								3										
Octet	Bit	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31			
0	0	Version				IHL				DS				ECN				Total Length																		
4	32	Identification																Flags				Fragment Offset														
8	64	Time To Live								Protocol								Header Checksum																		
12	96	Source IP Address																																		
16	128	Destination IP Address																																		
20	160	Options (if IHL > 5)																																		
:	:																																			
56	448																																			

# Explicit Congestion Notification



# ECN on the TCP header

## TCP segment header

Offsets	Octet	0								1								2								3							
Octet	Bit	7	6	5	4	3	2	1	0	7	6	5	4	3	2	1	0	7	6	5	4	3	2	1	0	7	6	5	4	3	2	1	0
0	0	Source port																Destination port															
4	32	Sequence number																															
8	64	Acknowledgment number (if ACK set)																															
12	96	Data offset				Reserved 0000				C	E	U	A	P	R	S	F	Window Size															
										W	C	R	C	S	S	Y	I																
										R	E	G	K	H	T	N	N																
16	128	Checksum																Urgent pointer (if URG set)															
20	160	Options (if <i>data offset</i> > 5. Padded at the end with "0" bits if necessary.)																															
:	:																																
56	448																																

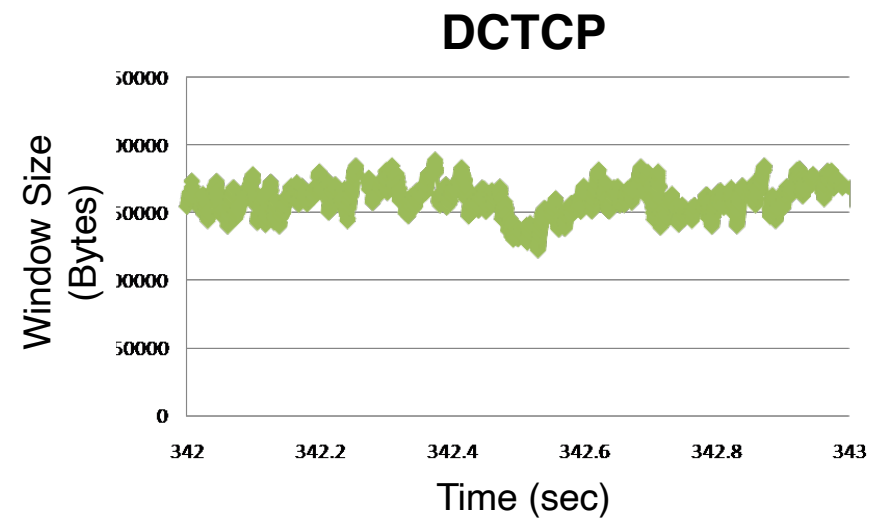
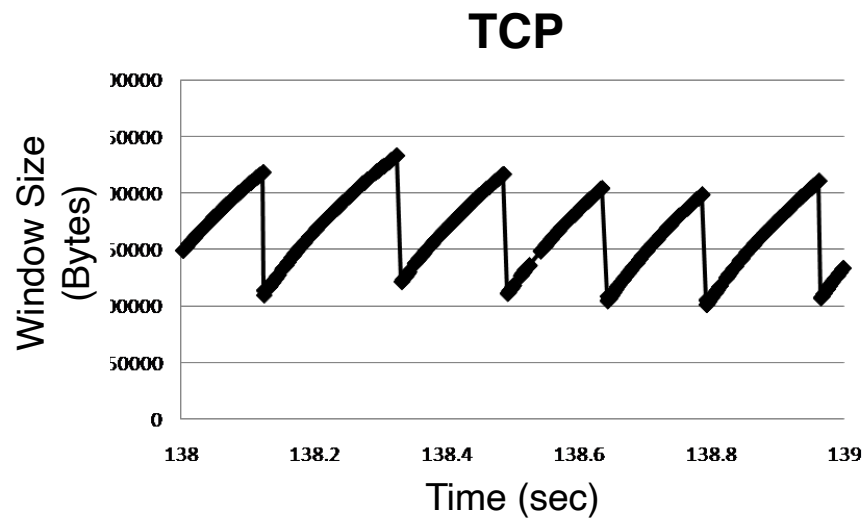


# DCTCP: Main idea

- Extract multi-bit feedback from single-bit stream of ECN marks
  - Reduce window size based on **fraction** of marked packets

# DCTCP: Main idea

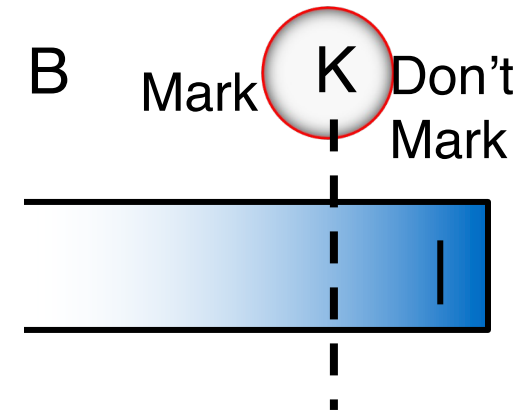
ECN Marks	TCP	DCTCP
1 0 1 1 1 1 0 1 1 1	Cut window by <b>50%</b>	Cut window by <b>40%</b>
0 0 0 0 0 0 0 0 0 1	Cut window by <b>50%</b>	Cut window by <b>5%</b>



# DCTCP algorithm

## Switch side:

- Mark packets when Queue Length  $> K$ .



## Sender side:

- Maintain running average of *fraction* of packets marked ( $\alpha$ ).

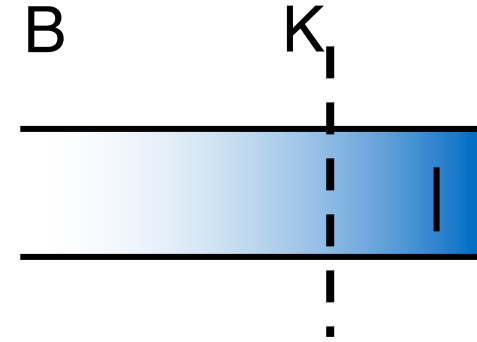
$$\text{each RTT: } F = \frac{\# \text{ of marked ACKs}}{\text{Total \# of ACKs}} \Rightarrow \alpha \leftarrow (1 - g)\alpha + gF$$

- Adaptive window decreases:**  $W \leftarrow (1 - \frac{\alpha}{2})W$ 
  - Note: decrease factor between 1 and 2.

Reacting to and controlling queue size **distribution**, specifically, the region above  $K$ .

# Setting parameters: Marking threshold $K$

- Tradeoff!
- Mark too late: too much queueing & pkt loss
- Mark too early: queues too small, **lose throughput**
- How small can queues be without loss of throughput?



$$K > (1/7) C \times RTT$$



for TCP:  
$$K > C \times RTT$$

# Buffering requirements

- How much buffering does DCTCP need for 100% throughput?
- Need to quantify queue size **oscillations** (stability).

