

COMP 3331/9331: **Computer Networks and** **Applications**

Week 6
Congestion Control (Transport Layer)

Reading Guide: Chapter 3, Sections: 3.6, 3.7

Transport Layer: Outline

3.1 transport-layer services

3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP

3.4 principles of reliable data transfer

3.5 connection-oriented transport: TCP

- segment structure
- reliable data transfer
- flow control
- connection management

3.6 principles of congestion control

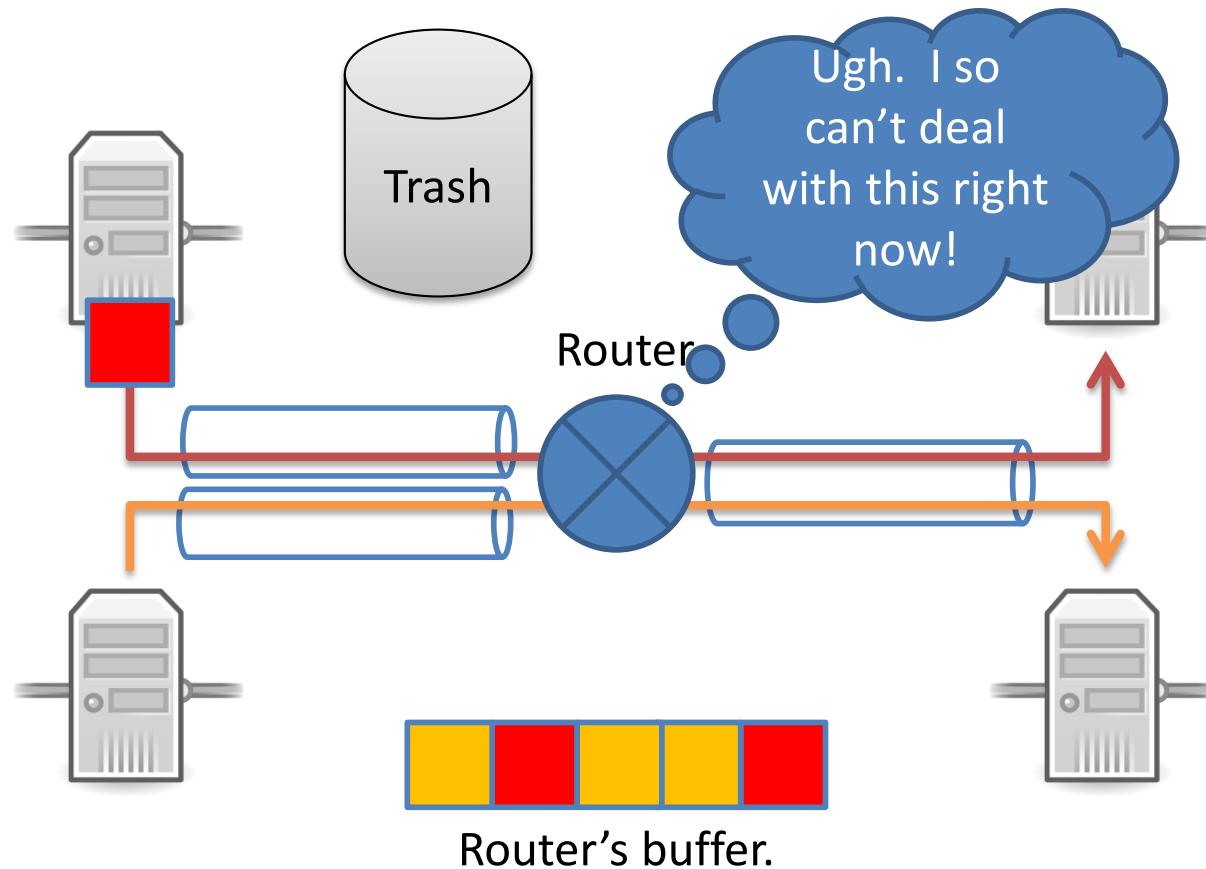
3.7 TCP congestion control

Principles of congestion control

congestion:

- ❖ informally: “too many sources sending too much data too fast for *network* to handle”
- ❖ different from flow control!
- ❖ manifestations:
 - lost packets (buffer overflow at routers)
 - long delays (queueing in router buffers)
- ❖ a top-10 problem!

Congestion



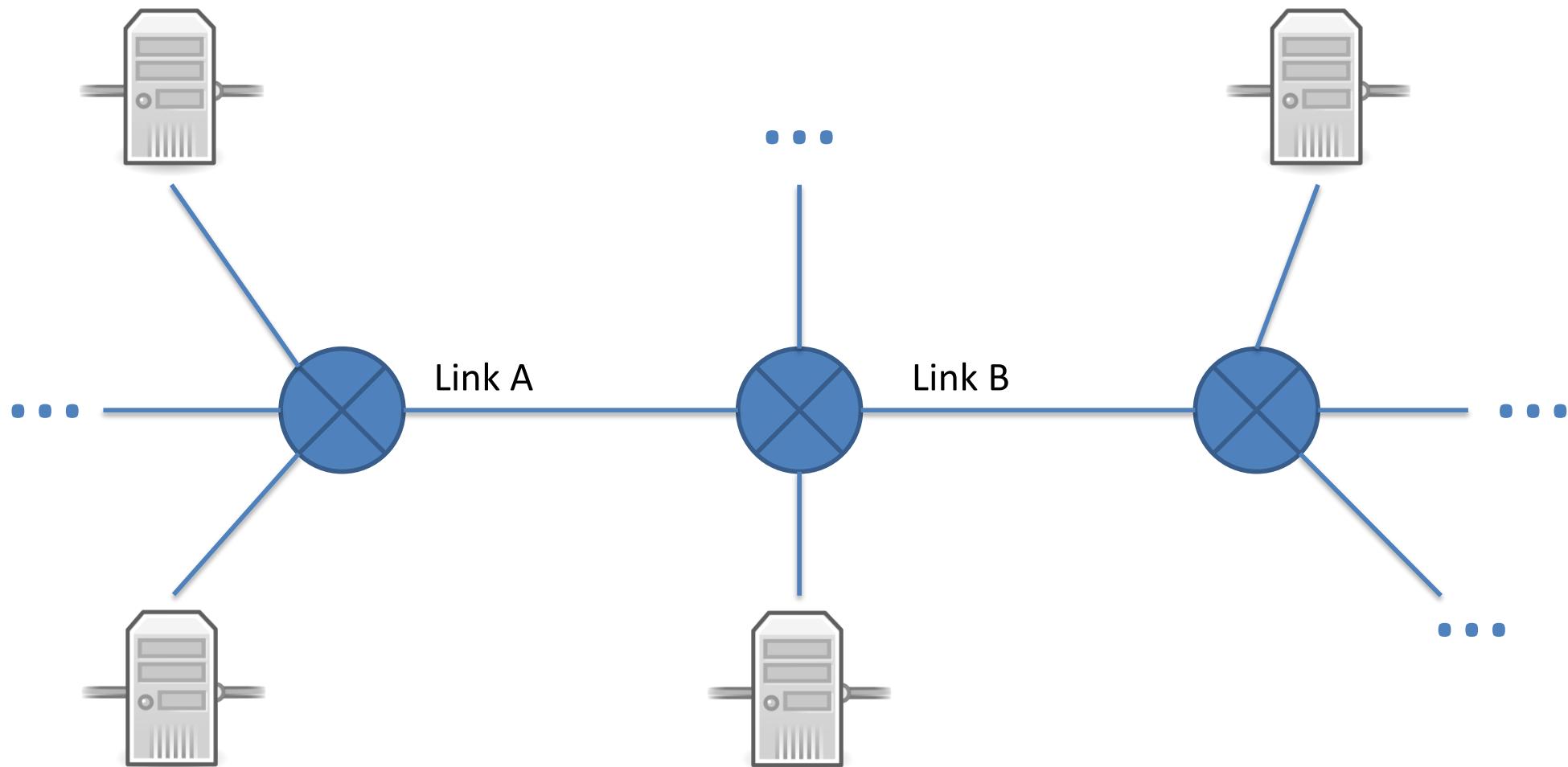
Incoming rate is faster than outgoing link can support.

Quiz: What's the worst that can happen?

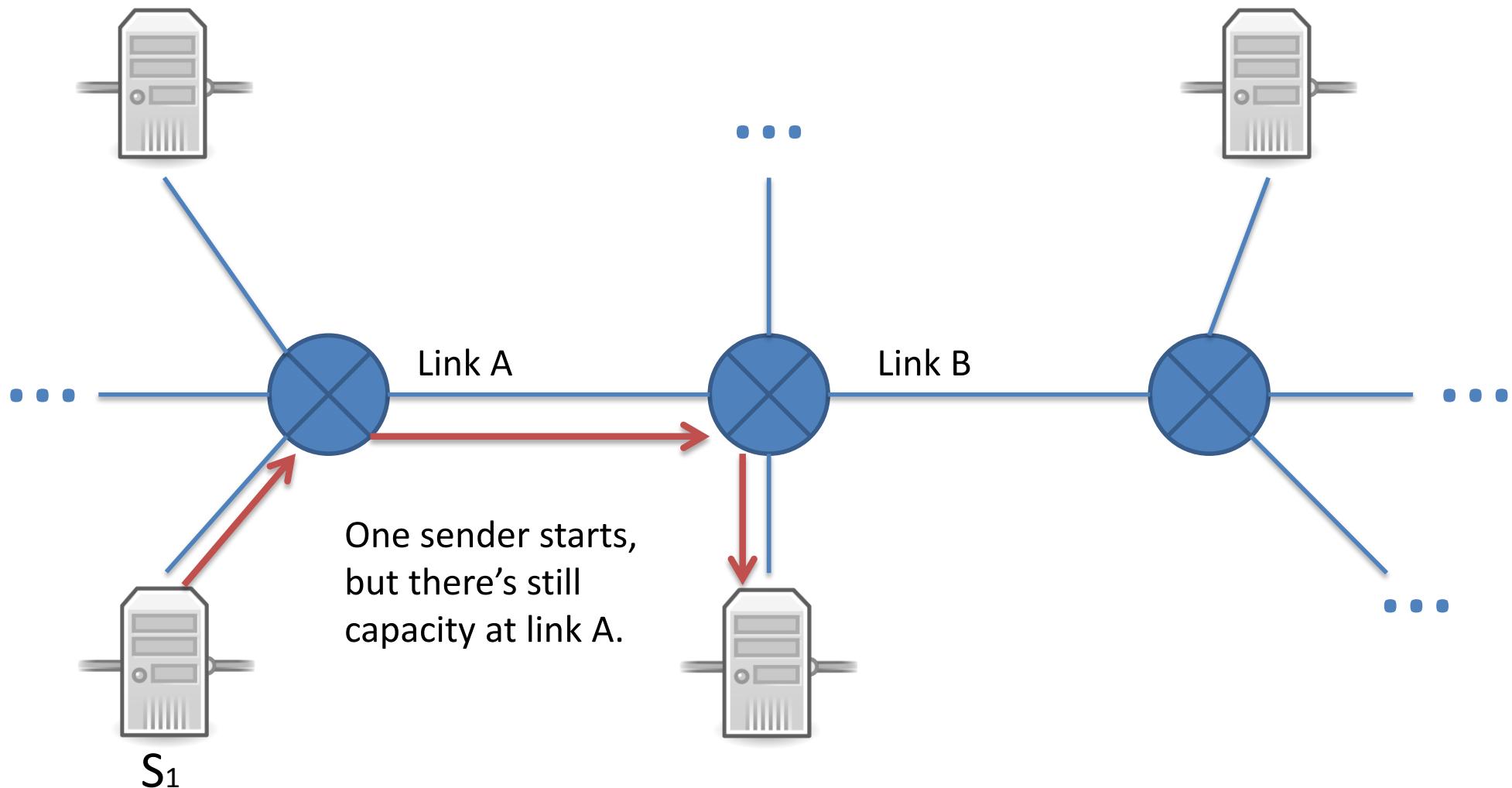


- A: This is no problem. Senders just keep transmitting, and it'll all work out.
- B: There will be retransmissions, but the network will still perform without much trouble.
- C: Retransmissions will become very frequent, causing a serious loss of efficiency
- D: The network will become completely unusable

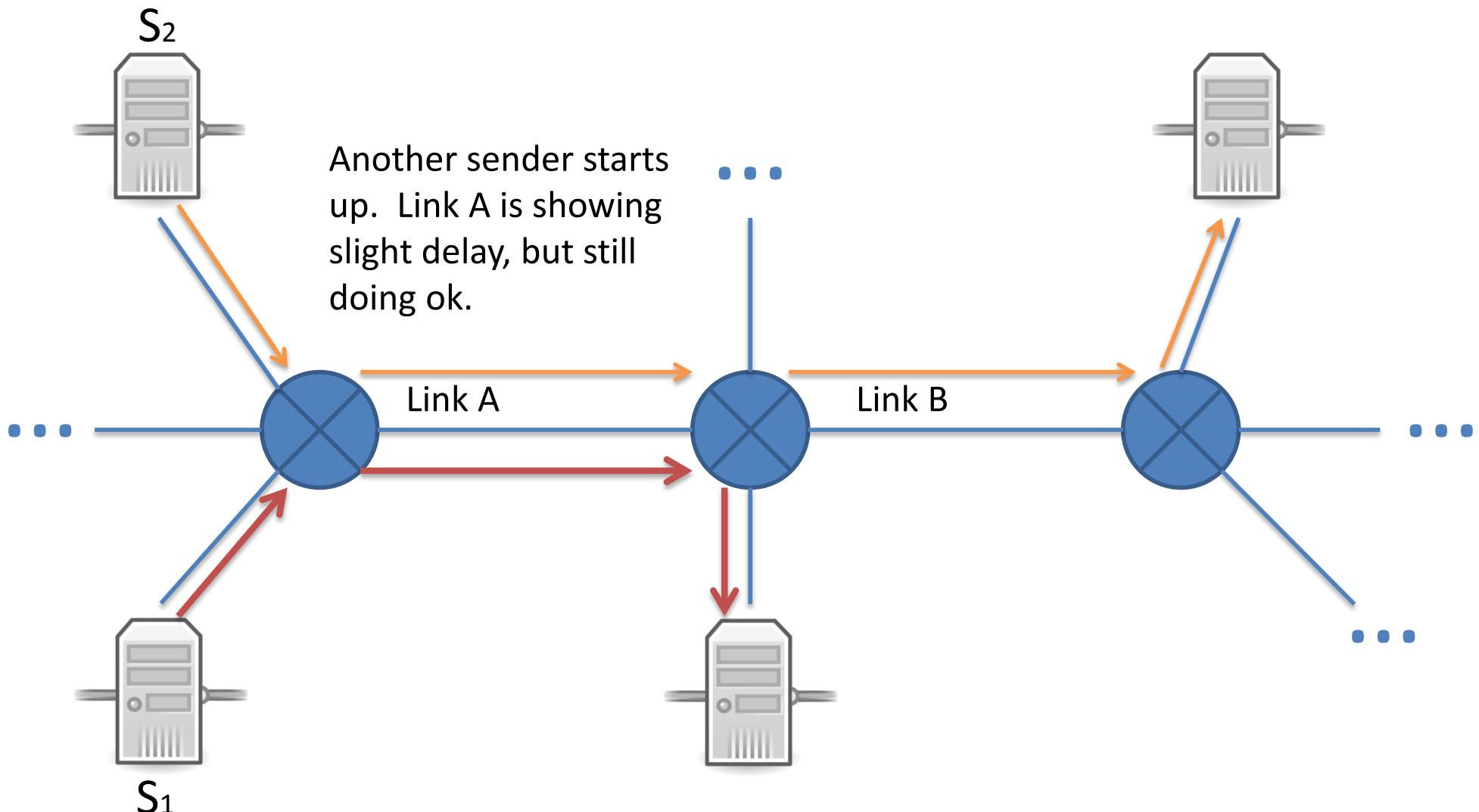
Congestion Collapse



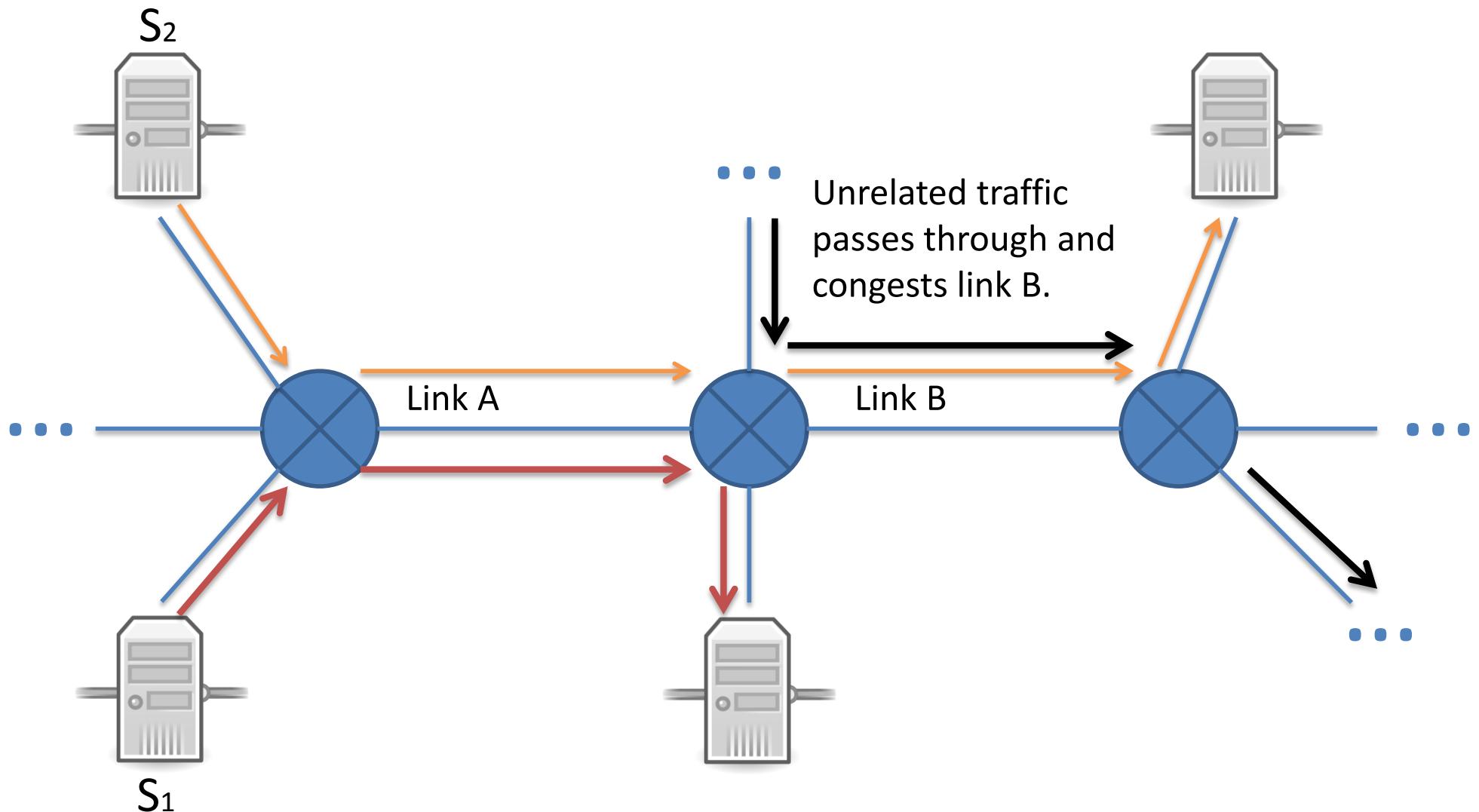
Congestion Collapse



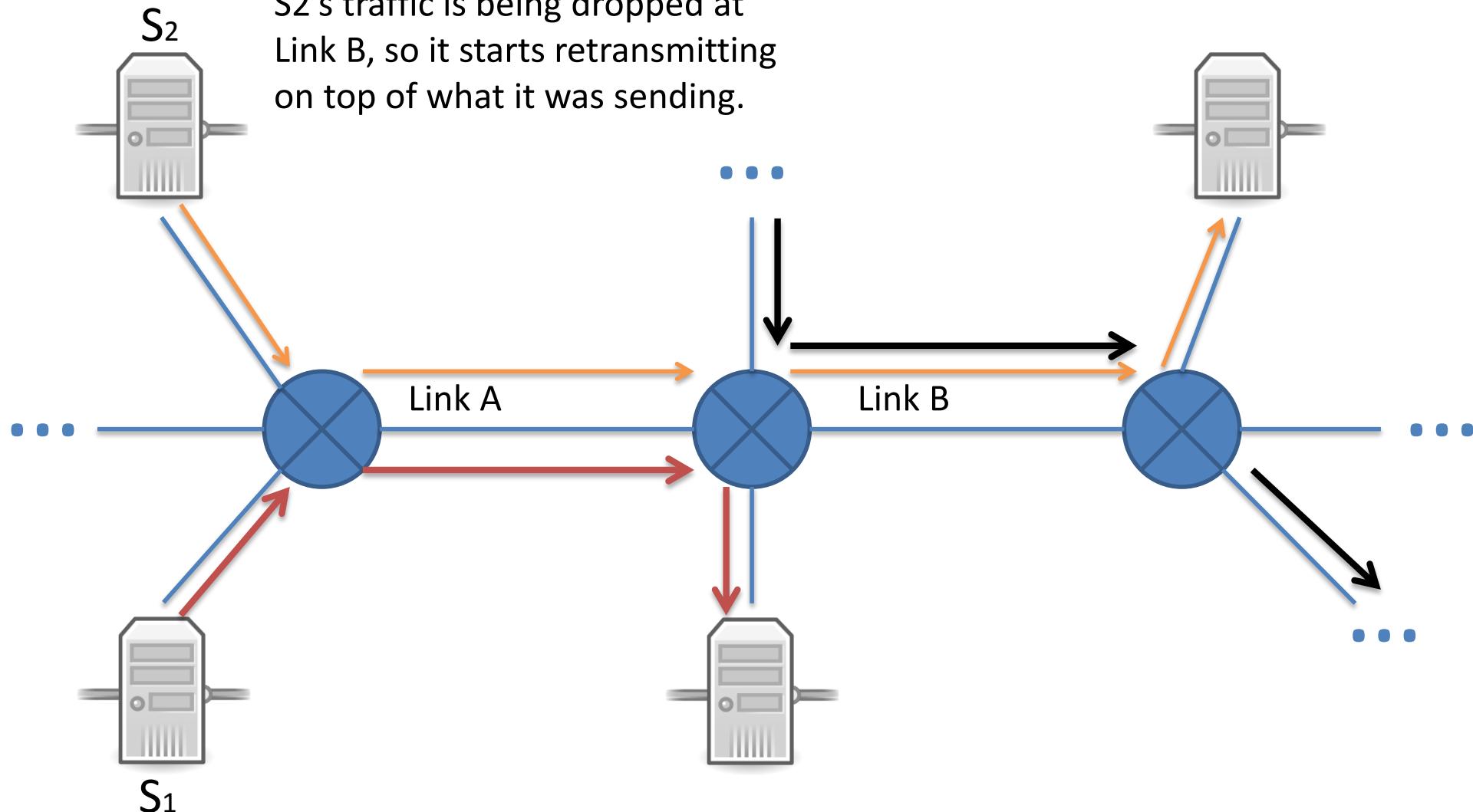
Congestion Collapse



Congestion Collapse

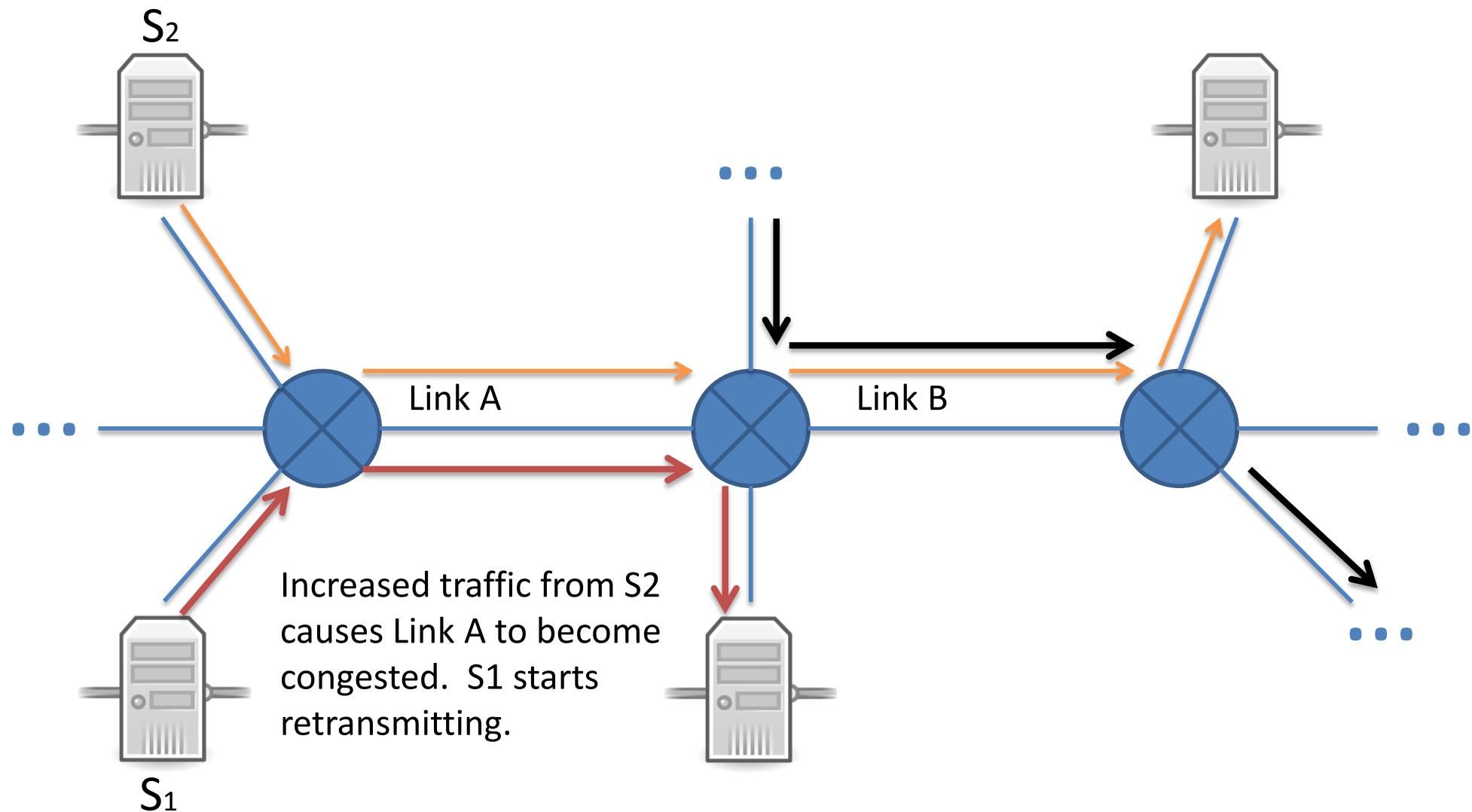


Congestion Collapse

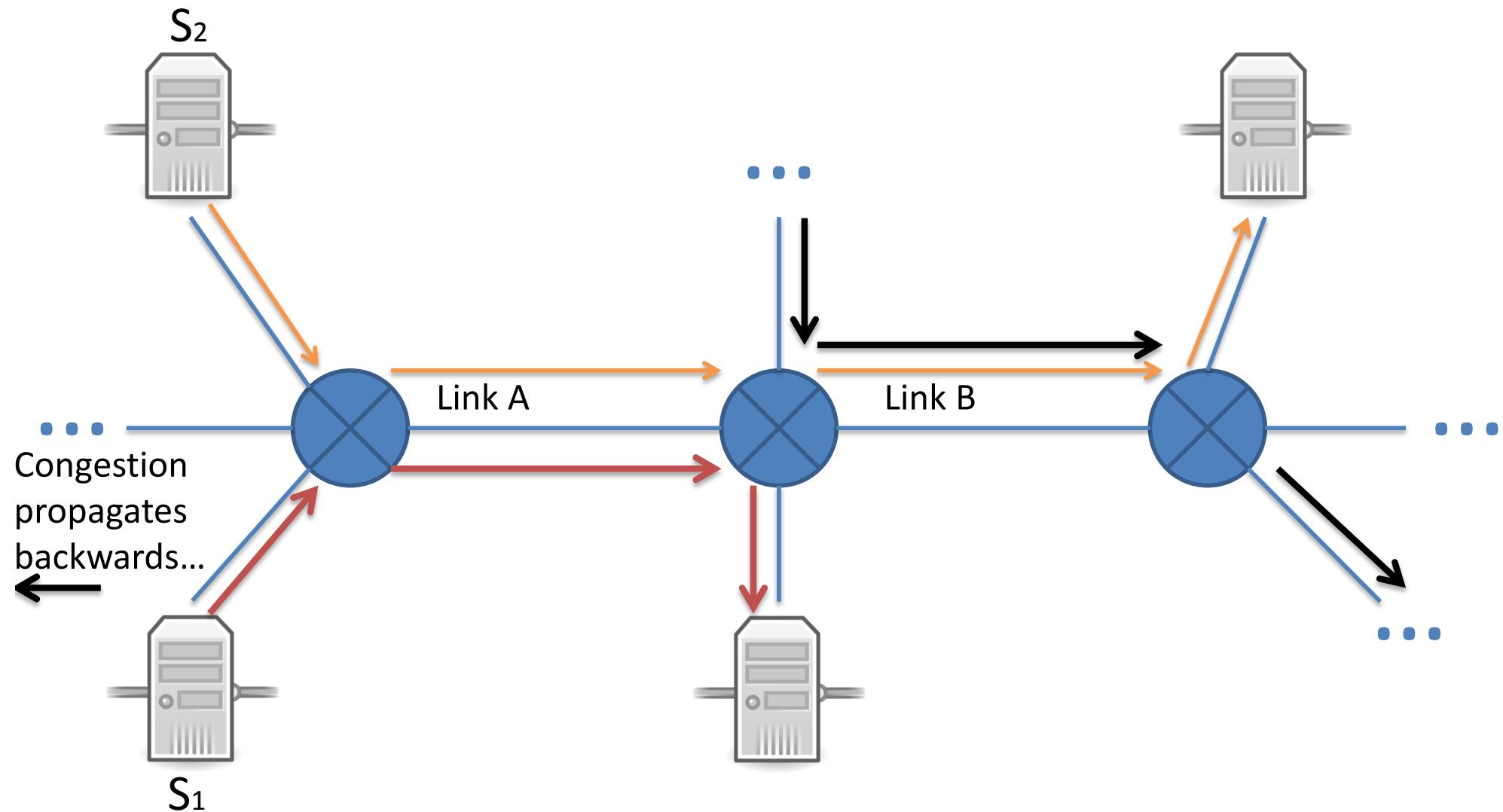


(This is very bad. S2 is now sending lots of traffic over link A that has no hope of crossing link B.) Congestion Control

Congestion Collapse



Congestion Collapse



Without congestion control

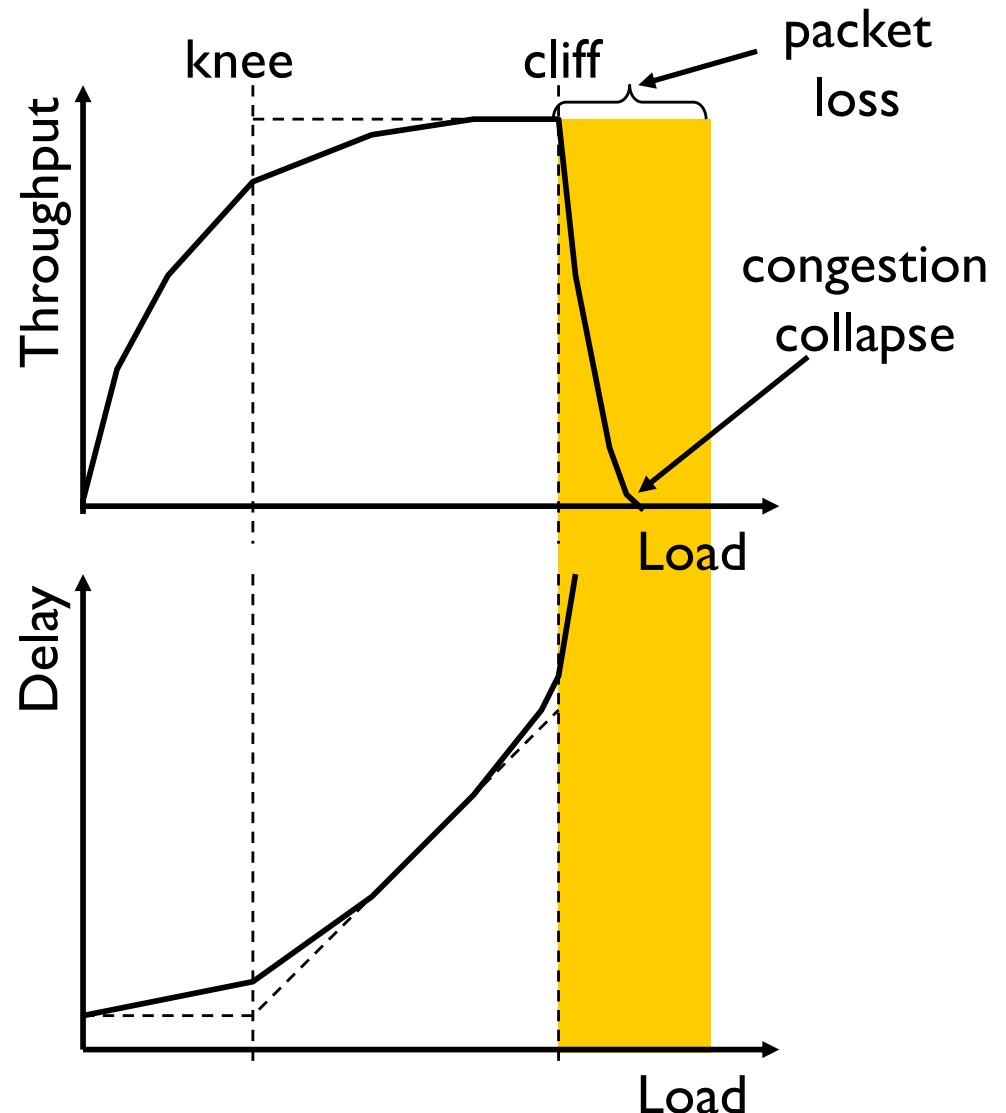
congestion:

- ❖ Increases delivery latency
 - Variable delays
 - If delays > RTO, sender retransmits
- ❖ Increases loss rate
 - Dropped packets also retransmitted
- ❖ Increases retransmissions, many unnecessary
 - Wastes capacity of traffic that is never delivered
 - Increase in load results in decrease in useful work done
- ❖ Increases congestion, cycle continues ...

Cost of Congestion

- ❖ Knee – point after which
 - Throughput increases slowly
 - Delay increases fast

- ❖ Cliff – point after which
 - Throughput starts to drop to zero (congestion collapse)
 - Delay approaches infinity



Congestion Collapse

This happened to the Internet (then NSFnet) in 1986

- ❖ Rate dropped from a *blazing* 32 Kbps to 40bps
- ❖ This happened on and off for two years
- ❖ In 1988, Van Jacobson published “Congestion Avoidance and Control”
- ❖ The fix: senders voluntarily limit sending rate

Approaches towards congestion control

two broad approaches towards congestion control:

end-end congestion control:

- ❖ no explicit feedback from network
- ❖ congestion inferred from end-system observed loss, delay
- ❖ approach taken by TCP

network-assisted congestion control:

- ❖ routers provide feedback to end systems
 - single bit indicating congestion (SNA, DECbit, TCP/IP ECN, ATM)
 - explicit rate for sender to send at

Transport Layer: Outline

3.1 transport-layer services

3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP

3.4 principles of reliable data transfer

3.5 connection-oriented transport: TCP

- segment structure
- reliable data transfer
- flow control
- connection management

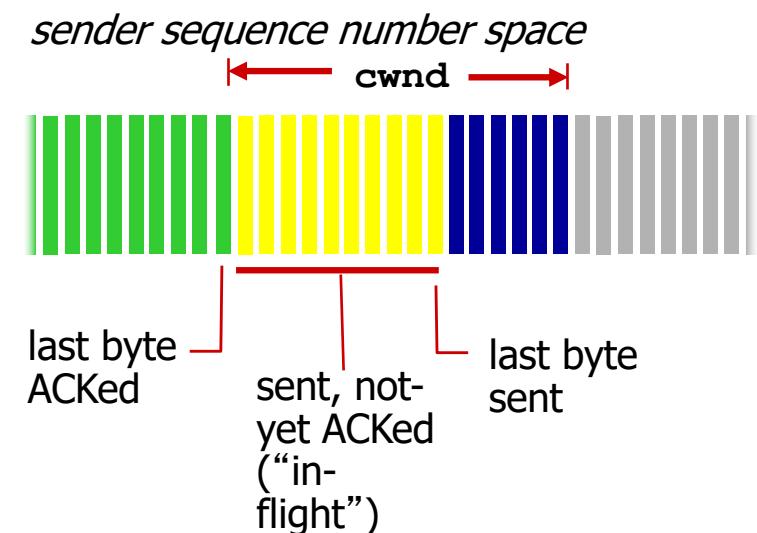
3.6 principles of congestion control

3.7 TCP congestion control

TCP's Approach in a Nutshell

- ❖ TCP connection has window
 - Controls number of packets in flight
- ❖ *TCP sending rate:*
 - roughly: send cwnd bytes, wait RTT for ACKS, then send more bytes

$$\text{rate} \approx \frac{\text{cwnd}}{\text{RTT}} \text{ bytes/sec}$$



- ❖ Vary window size to control sending rate

All These Windows...

- ❖ Congestion Window: **CWND**
 - How many bytes can be sent without overflowing routers
 - Computed by the sender using congestion control algorithm
- ❖ Flow control window: **Advertised / Receive Window (RWND)**
 - How many bytes can be sent without overflowing receiver's buffers
 - Determined by the receiver and reported to the sender
- ❖ Sender-side window = **minimum{CWND, RWND}**
 - Assume for this lecture that RWND >> CWND

CWND

- ❖ This lecture will talk about CWND in units of MSS
 - (Recall MSS: Maximum Segment Size, the amount of payload data in a TCP packet)
 - This is only for pedagogical purposes
- ❖ Keep in mind that real implementations maintain CWND in bytes

Two Basic Questions

- ❖ How does the sender detect congestion?
- ❖ How does the sender adjust its sending rate?

Quiz: What is a “congestion event”



- A: A segment loss (but how can the sender be sure of this?)
- B: Increased delays
- C: Receiving duplicate acknowledgement (s)
- D: A retransmission timeout firing
- E: Some subset of A, B, C & D (what is the subset?)

Quiz: How should we set CWND?

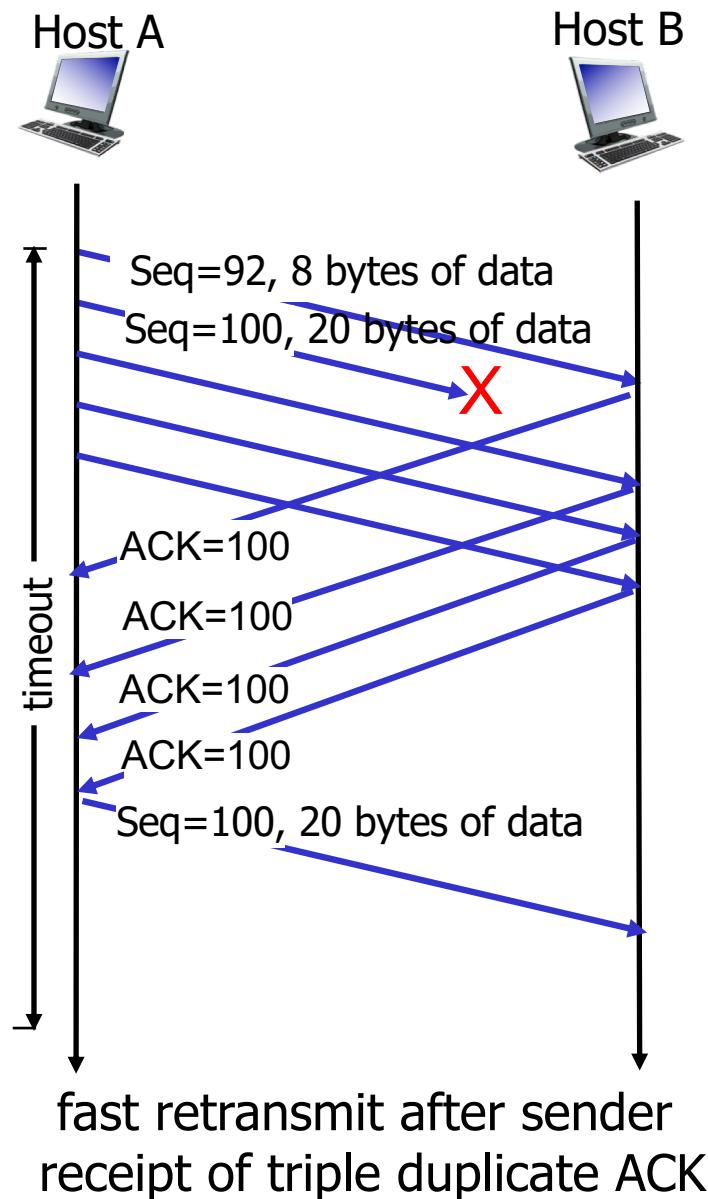


- A: We should keep raising it until a “congestion event” then back off slightly until we notice no more events.
- B: We should raise it until a “congestion event”, then go back to 1 and start raising it again
- C: We should raise it until a “congestion event”, then go back to median value and start raising it again.
- D: We should sent as fast as possible at all times.

Not All Losses the Same

- ❖ Duplicate ACKs: isolated loss
 - dup ACKs indicate network capable of delivering some segments
- ❖ Timeout: much more serious
 - Not enough dup ACKs
 - Must have suffered several losses
- ❖ Will adjust rate differently for each case

RECAP: TCP fast retransmit

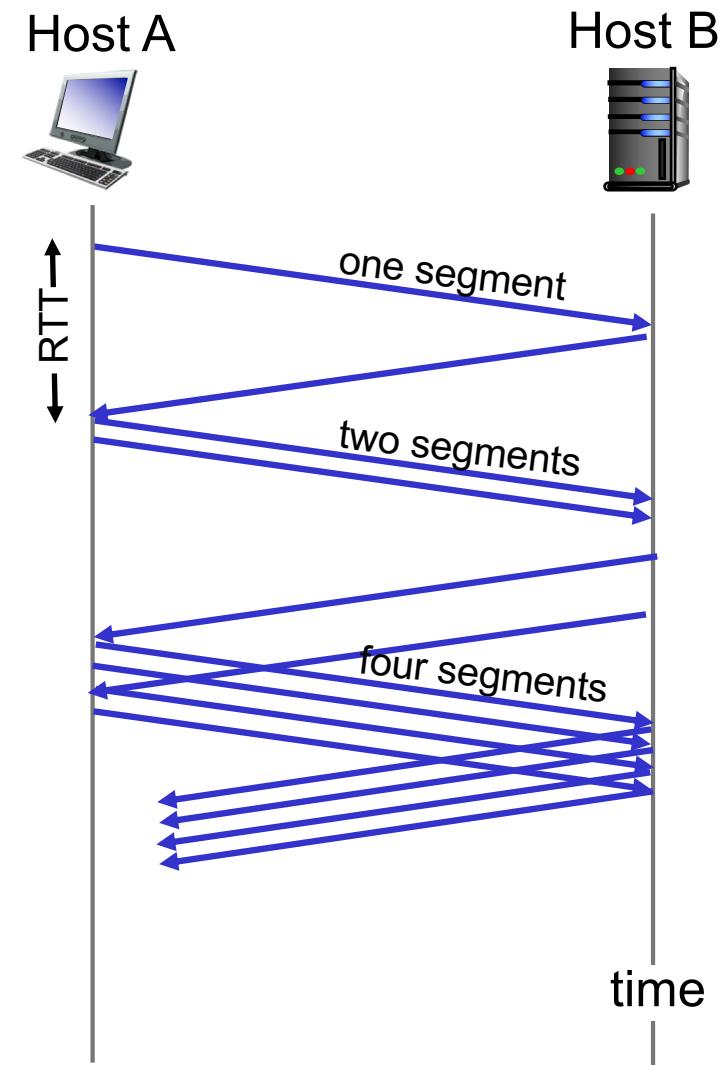


Rate Adjustment

- ❖ Basic structure:
 - Upon receipt of ACK (of new data): increase rate
 - Upon detection of loss: decrease rate
- ❖ How we increase/decrease the rate depends on the phase of congestion control we're in:
 - Discovering available bottleneck bandwidth vs.
 - Adjusting to bandwidth variations

TCP Slow Start (Bandwidth discovery)

- ❖ when connection begins, increase rate **exponentially** until **first loss event**:
 - initially **cwnd** = 1 MSS
 - double **cwnd** every RTT (full ACKs)
 - Simpler implementation achieved by incrementing **cwnd** for every ACK received
 - $cwnd += 1$ for each ACK
- ❖ **summary:** initial rate is slow but ramps up exponentially fast



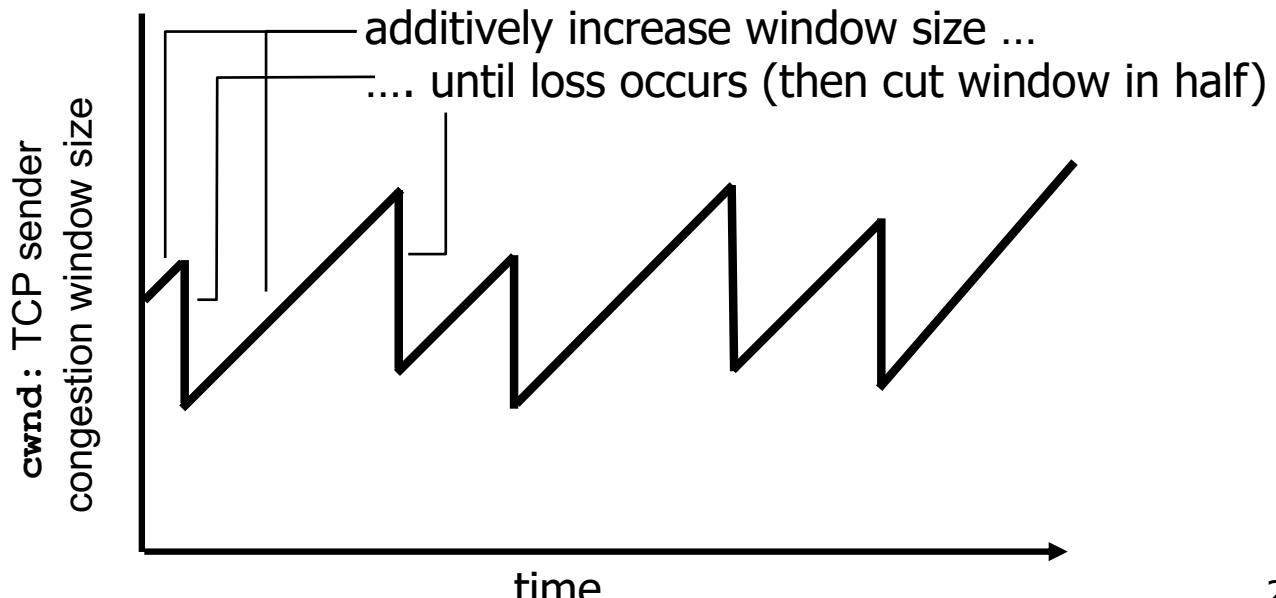
Adjusting to Varying Bandwidth

- ❖ Slow start gave an estimate of available bandwidth
- ❖ Now, want to track variations in this available bandwidth, oscillating around its current value
 - Repeated probing (rate increase) and backoff (rate decrease)
 - Known as Congestion Avoidance (CA)
- ❖ TCP uses: “Additive Increase Multiplicative Decrease” (AIMD)
 - We’ll see why shortly...

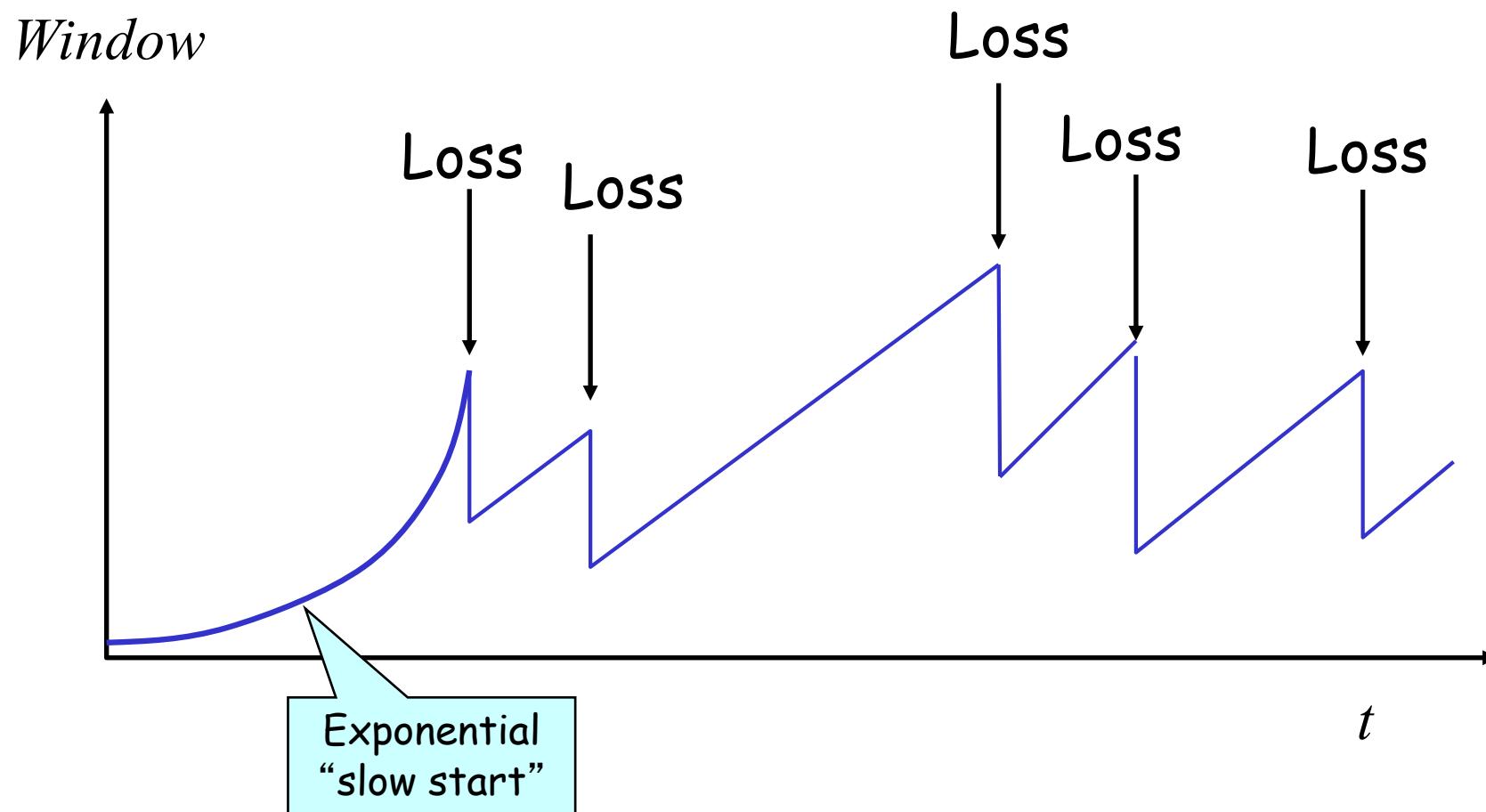
AIMD

- ❖ **approach:** sender increases transmission rate (window size), probing for usable bandwidth, until another congestion event occurs
 - **additive increase:** increase **cwnd** by 1 MSS every RTT until loss detected
 - For each successful RTT (all ACKS), $cwnd = cwnd + 1$
 - Simple implementation: for each ACK, $cwnd = cwnd + 1/cwnd$
 - **multiplicative decrease:** cut **cwnd** in half after loss

AIMD saw tooth behavior: probing for bandwidth



Leads to the TCP “Sawtooth”



Slow-Start vs. AIMD

- ❖ When does a sender stop Slow-Start and start Additive Increase?
- ❖ Introduce a “slow start threshold” (**ssthresh**)
 - Initialized to a large value
- ❖ Convert to AI when $cwnd = ssthresh$, sender switches from slow-start to AIMD-style increase
 - On timeout, $ssthresh = CWND/2$

Implementation

- ❖ State at sender

- CWND (initialized to a small constant)
- ssthresh (initialized to a large constant)
- [Also dupACKcount and timer, as before]

- ❖ Events

- ACK (new data)
- dupACK (duplicate ACK for old data)
- Timeout

Event: ACK (new data)

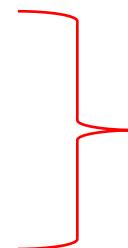
- ❖ If $CWND < ssthresh$

- $CWND += +$

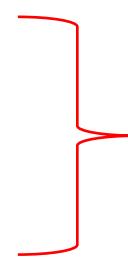
- $2 \times MSS$ packets per ACK
- Hence after one RTT (All ACKs with no drops):
 $CWND = 2 \times CWND$

Event: ACK (new data)

- ❖ If $CWND < ssthresh$
 - $CWND += I$
- ❖ Else
 - $CWND = CWND + I/CWND$



Slow start phase



*“Congestion
Avoidance” phase
(additive increase)*

- Hence after one RTT (All ACKs with no drops):
 $CWND = CWND + I$

Event: dupACK

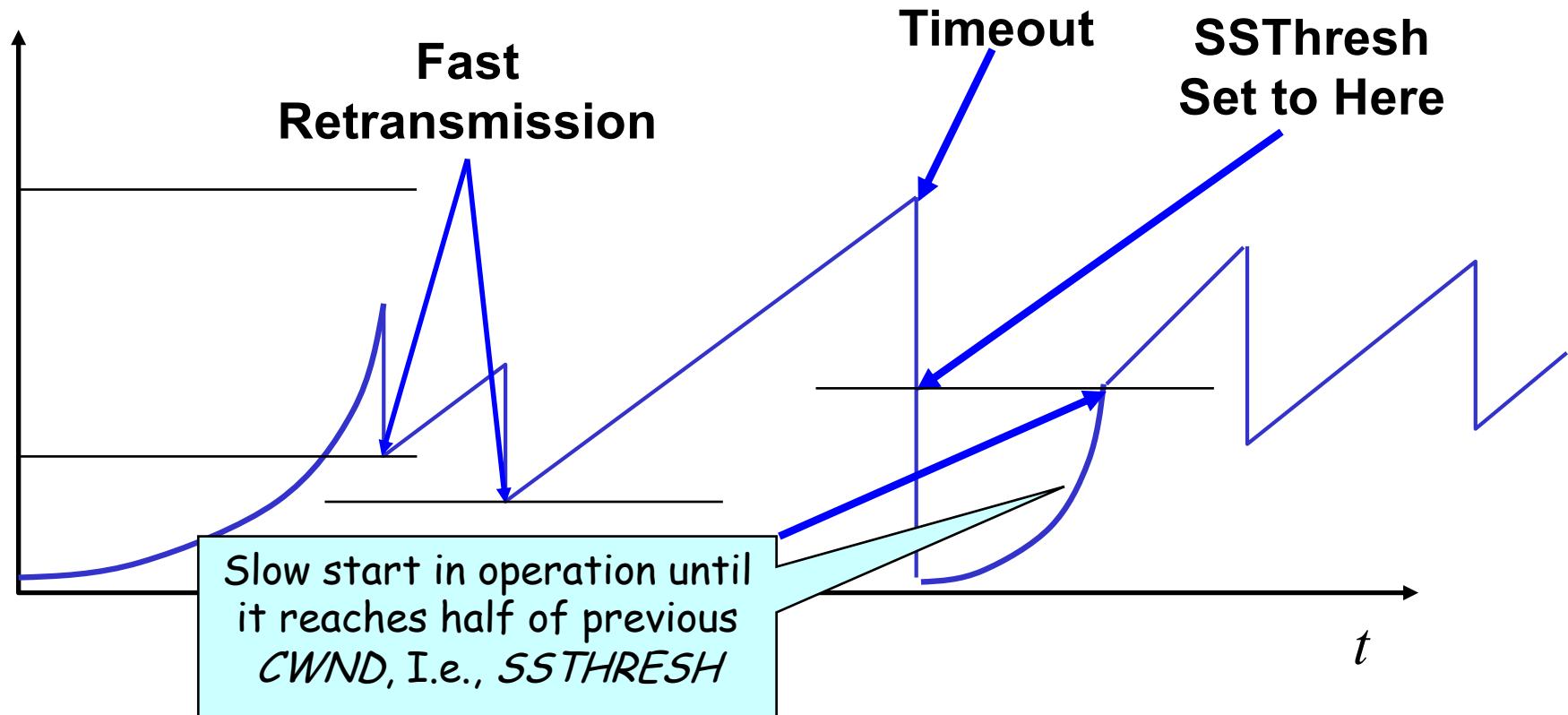
- ❖ dupACKcount ++
- ❖ If dupACKcount = 3 /* fast retransmit */
 - ssthresh = CWND/2
 - **CWND = CWND/2**

Event: TimeOut

- ❖ On Timeout
 - $ssthresh \leftarrow CWND/2$
 - $CWND \leftarrow 1$

Example

Window



Slow-start restart: Go back to $CWND = 1$ MSS, but take advantage of knowing the previous value of $CWND$

One Final Phase: Fast Recovery

- ❖ The problem: Fast retransmit slow in recovering from an isolated loss

Example (window in units of MSS, not bytes)

- ❖ Consider a TCP connection with:
 - CWND=10 packets (of size MSS, which is 100 bytes)
 - Last ACK was for byte # 101
 - i.e., receiver expecting next packet to have seq. no. 101
- ❖ 10 packets [101, 201, 301,..., 1001] are in flight
 - Packet 101 is dropped
 - What ACKs do they generate?
 - And how does the sender respond?

Timeline

- ❖ ACK 101 (due to 201) cwnd=10 dupACK#1 (no xmit)
- ❖ ACK 101 (due to 301) cwnd=10 dupACK#2 (no xmit)
- ❖ ACK 101 (due to 401) cwnd=10 dupACK#3 (no xmit)
- ❖ RETRANSMIT 101 ssthresh=5 cwnd= 5
- ❖ ACK 101 (due to 501) cwnd=5 + 1/5 (no xmit)
- ❖ ACK 101 (due to 601) cwnd=5 + 2/5 (no xmit)
- ❖ ACK 101 (due to 701) cwnd=5 + 3/5 (no xmit)
- ❖ ACK 101 (due to 801) cwnd=5 + 4/5 (no xmit)
- ❖ ACK 101 (due to 901) cwnd=5 + 5/5 (no xmit)
- ❖ ACK 101 (due to 1001) cwnd=6 + 1/5 (no xmit)
- ❖ ACK 1101 (due to 101) ← only now can we transmit new packets
- ❖ Plus no packets in flight so ACK “clocking” (to increase CWND)
stalls for another RTT

Solution: Fast Recovery

Idea: Grant the sender temporary “credit” for each dupACK so as to keep packets in flight

- ❖ If $\text{dupACKcount} = 3$
 - $\text{ssthresh} = \text{cwnd}/2$
 - $\text{cwnd} = \text{ssthresh} + 3$
- ❖ While in fast recovery
 - $\text{cwnd} = \text{cwnd} + 1$ for each additional duplicate ACK
- ❖ Exit fast recovery after receiving new ACK
 - set $\text{cwnd} = \text{ssthresh}$

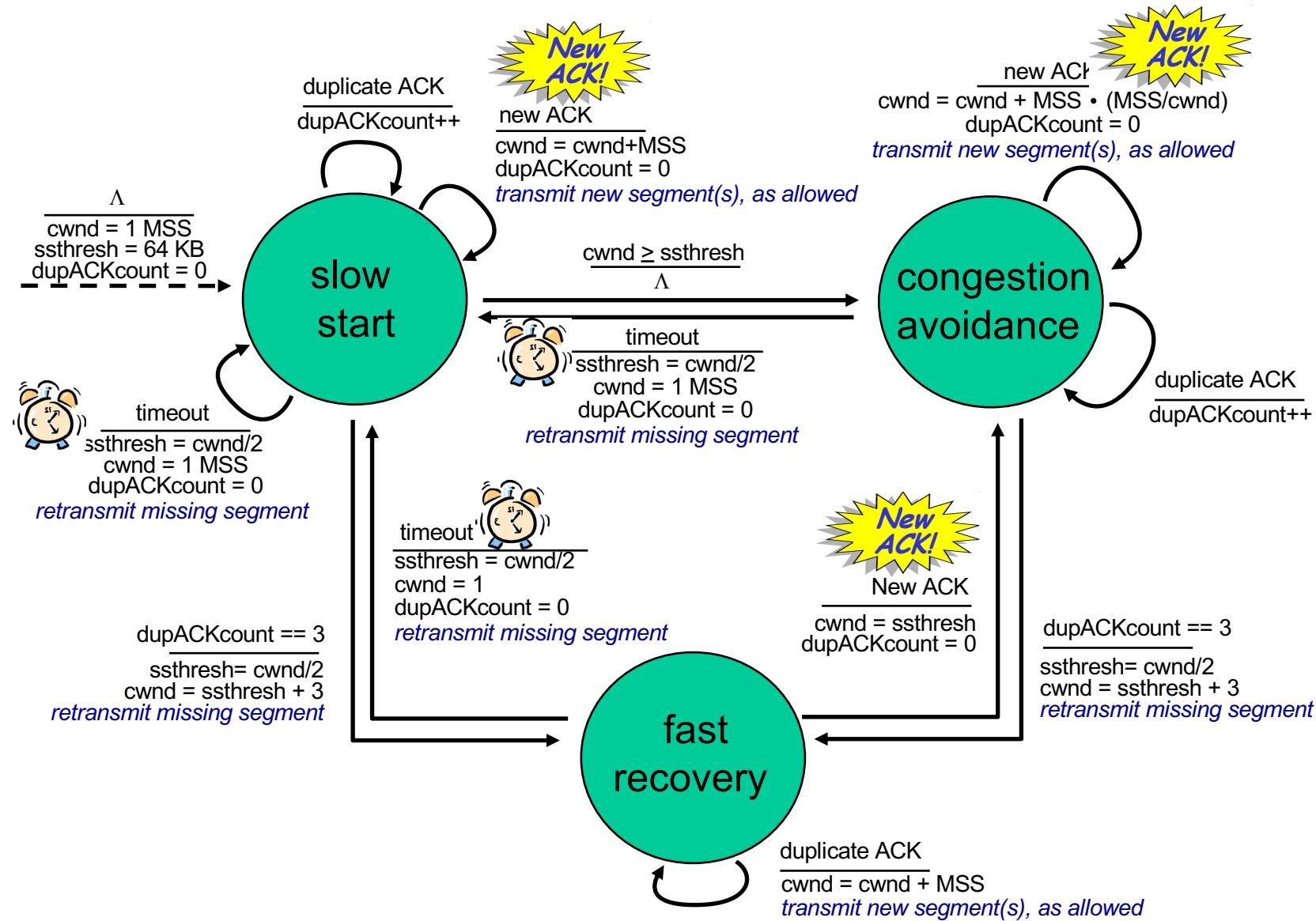
Example

- ❖ Consider a TCP connection with:
 - CWND=10 packets (of size MSS = 100 bytes)
 - Last ACK was for byte # 101
 - i.e., receiver expecting next packet to have seq. no. 101
- ❖ 10 packets [101, 201, 301,..., 1001] are in flight
 - **Packet 101 is dropped**

Timeline

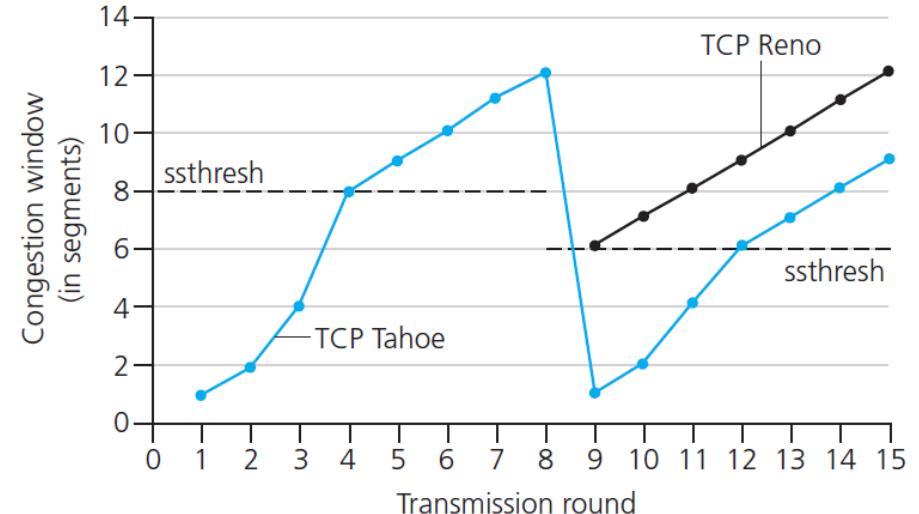
- ❖ ACK 101 (due to 201) cwnd=10 dup#1
- ❖ ACK 101 (due to 301) cwnd=10 dup#2
- ❖ ACK 101 (due to 401) cwnd=10 dup#3
- ❖ REXMIT 101 ssthresh=5 cwnd= 8 (5+3)
- ❖ ACK 101 (due to 501) cwnd= 9 (no xmit)
- ❖ ACK 101 (due to 601) cwnd=10 (no xmit)
- ❖ ACK 101 (due to 701) cwnd=11 (xmit 1101)
- ❖ ACK 101 (due to 801) cwnd=12 (xmit 1201)
- ❖ ACK 101 (due to 901) cwnd=13 (xmit 1301)
- ❖ ACK 101 (due to 1001) cwnd=14 (xmit 1401)
- ❖ ACK 1101 (due to 101) cwnd = 5 (xmit 1501) ← exiting fast recovery
- ❖ Packets 1101-1401 already in flight
- ❖ ACK 1201 (due to 1101) cwnd = 5 + 1/5 ← back in congestion avoidance

Summary: TCP Congestion Control



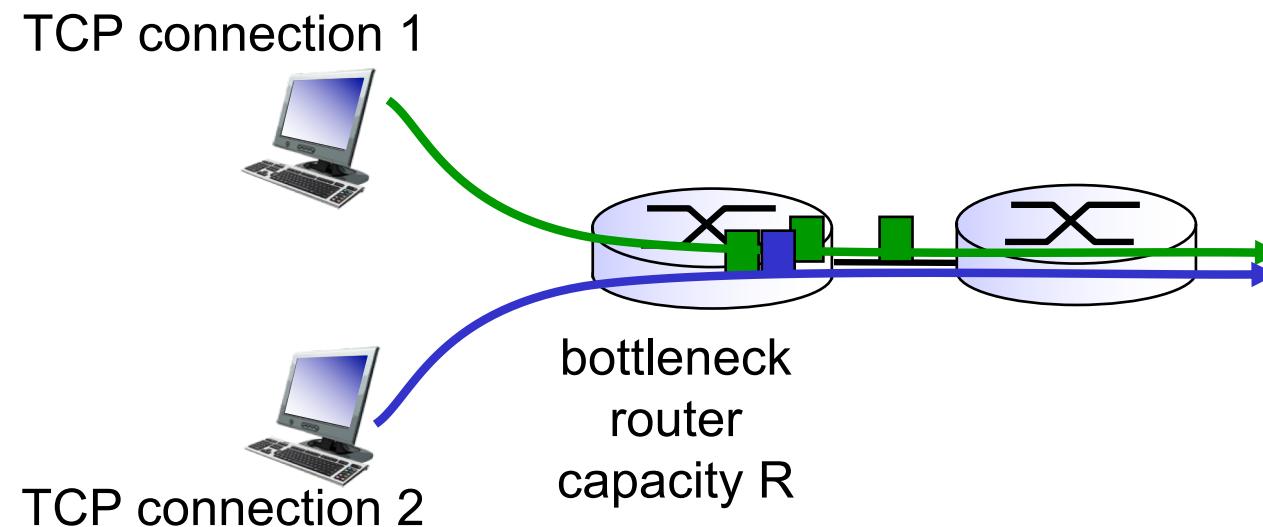
TCP Flavours

- ❖ TCP-Tahoe
 - $cwnd = 1$ on triple dup ACK & timeout
- ❖ TCP-Reno
 - $cwnd = 1$ on timeout
 - $cwnd = cwnd/2$ on triple dup ACK
- ❖ TCP-newReno
 - TCP-Reno + improved fast recovery
- ❖ TCP-SACK (NOT COVERED IN THE COURSE)
 - incorporates selective acknowledgements



TCP Fairness

fairness goal: if K TCP sessions share same bottleneck link of bandwidth R , each should have average rate of R/K

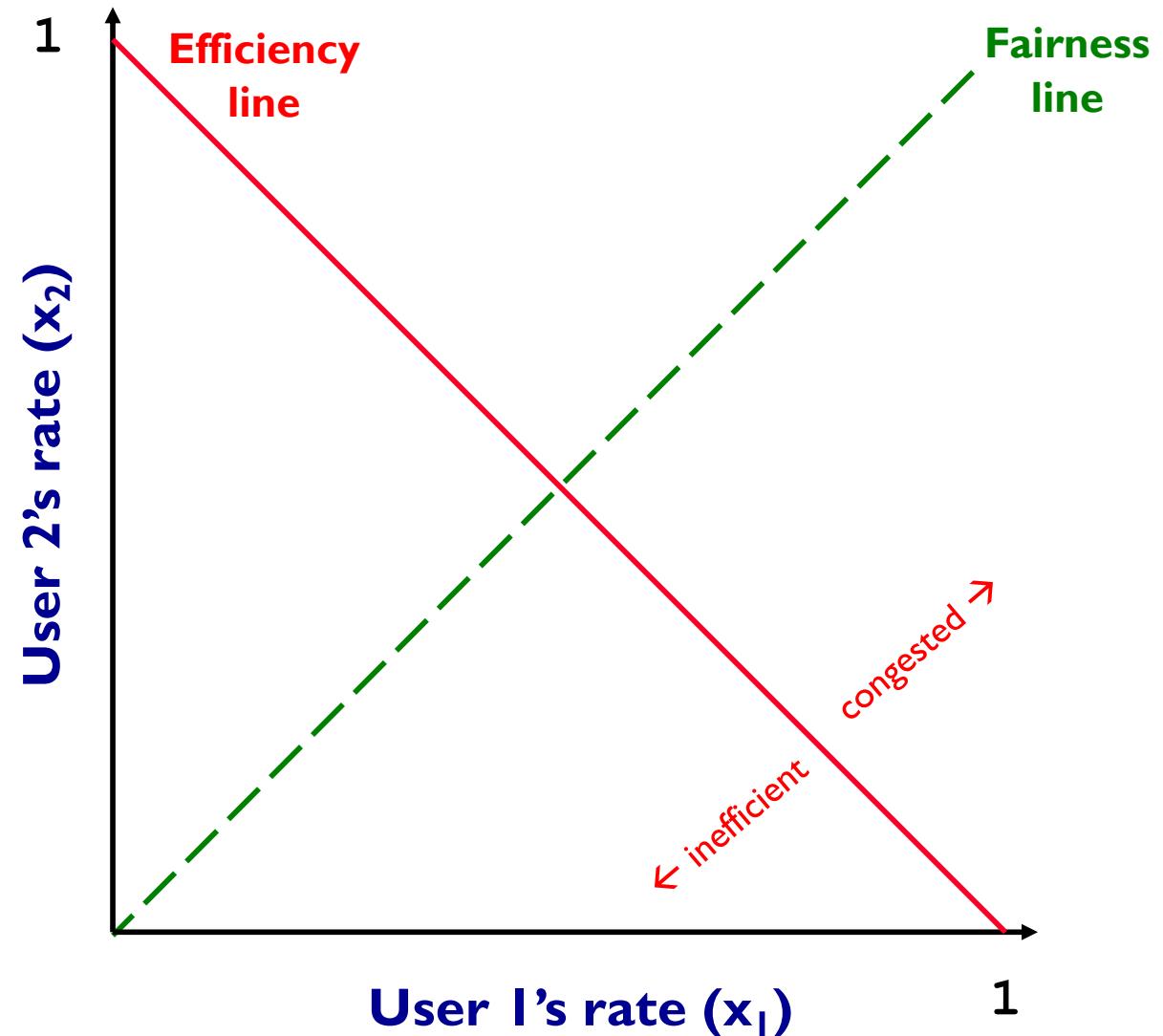


Why AIMD?

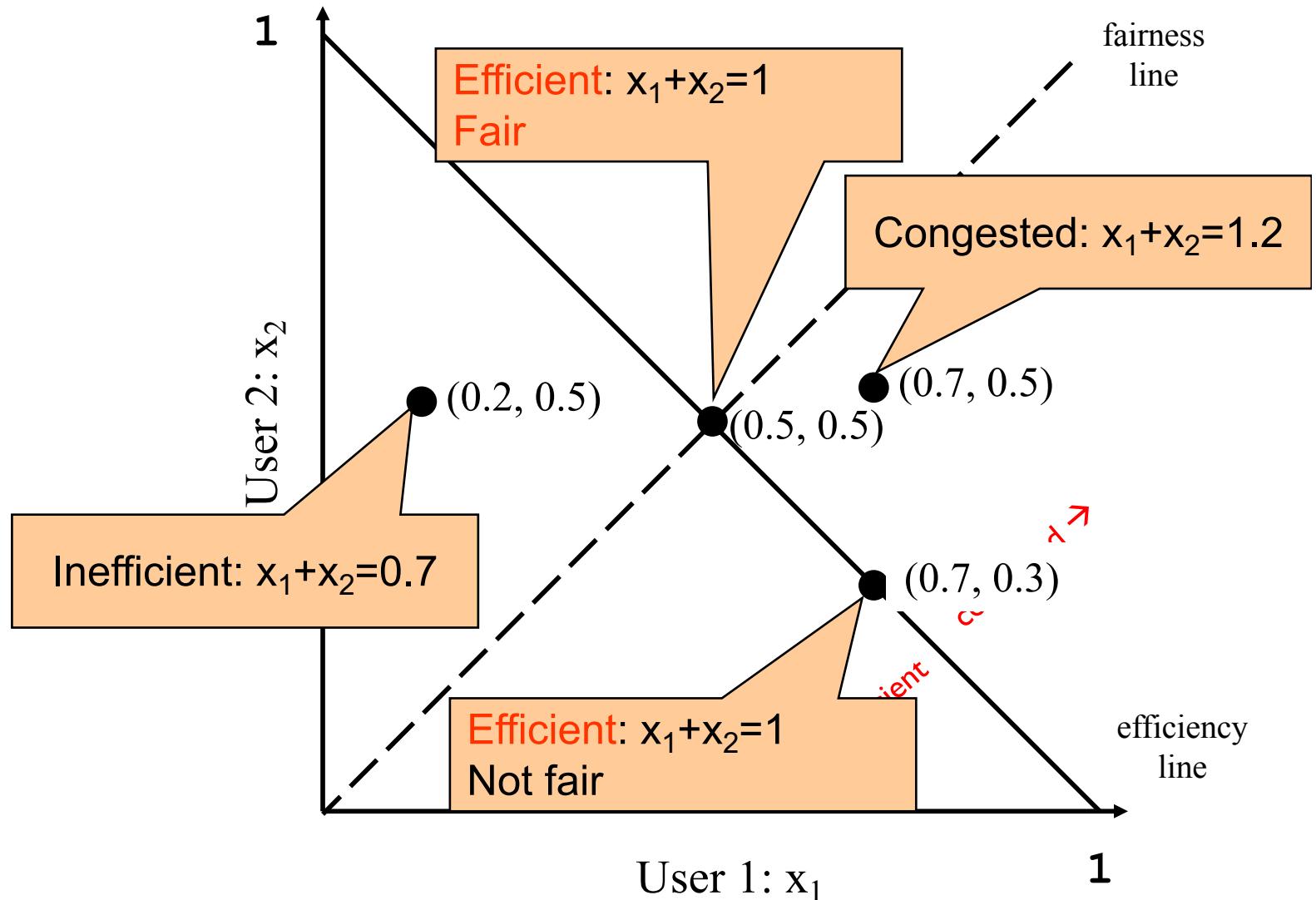
- ❖ Some rate adjustment options: Every RTT, we can
 - Multiplicative increase or decrease: $\text{WND} \rightarrow a * \text{WND}$
 - Additive increase or decrease: $\text{WND} \rightarrow \text{WND} + b$
- ❖ Four alternatives:
 - AIAD: gentle increase, gentle decrease
 - AIMD: gentle increase, drastic decrease
 - MIAD: drastic increase, gentle decrease
 - MIMD: drastic increase and decrease

Simple Model of Congestion Control

- ❖ Two users
 - rates x_1 and x_2
- ❖ Congestion when $x_1+x_2 > 1$
- ❖ Unused capacity when $x_1+x_2 < 1$
- ❖ Fair when $x_1 = x_2$

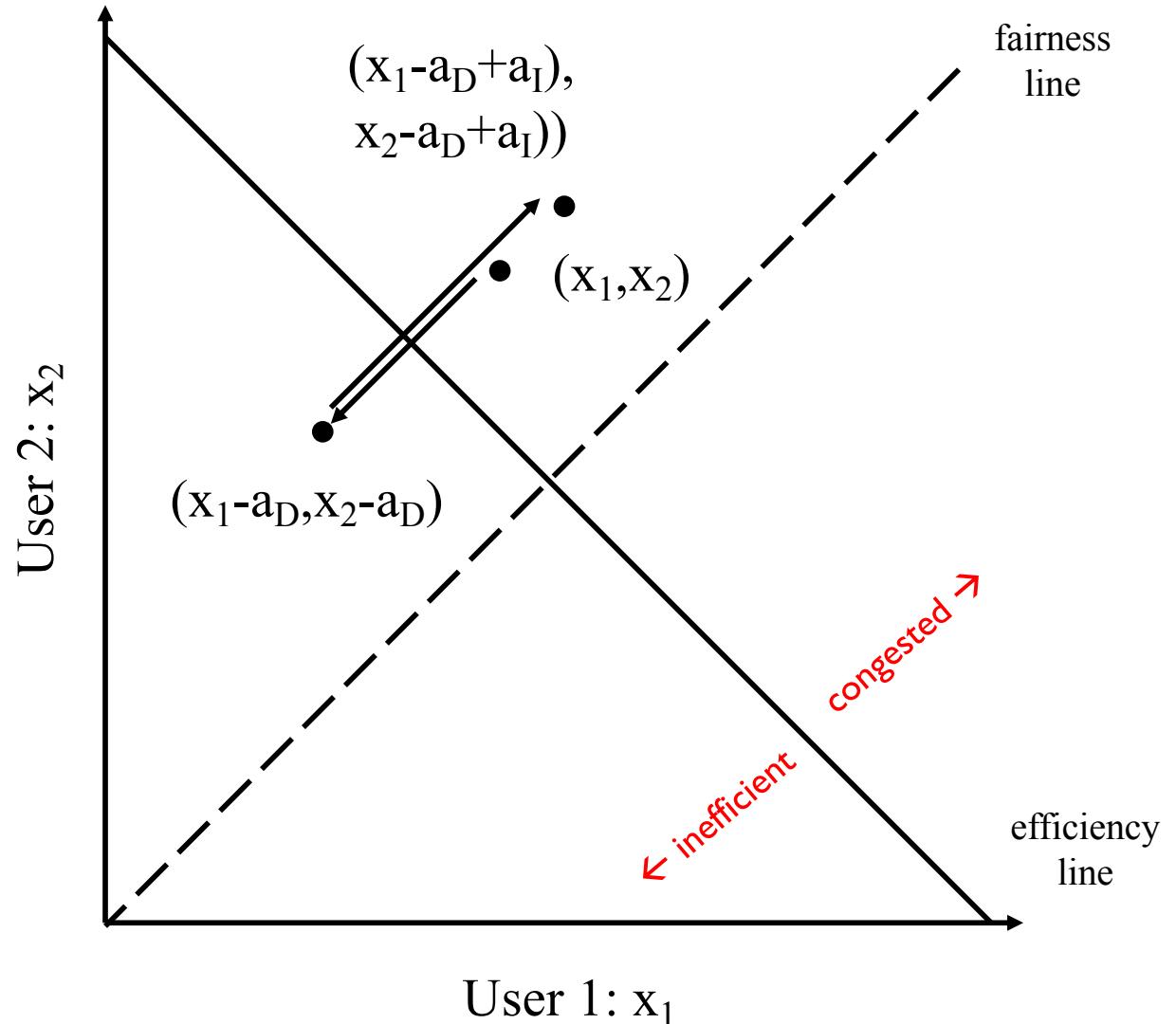


Example

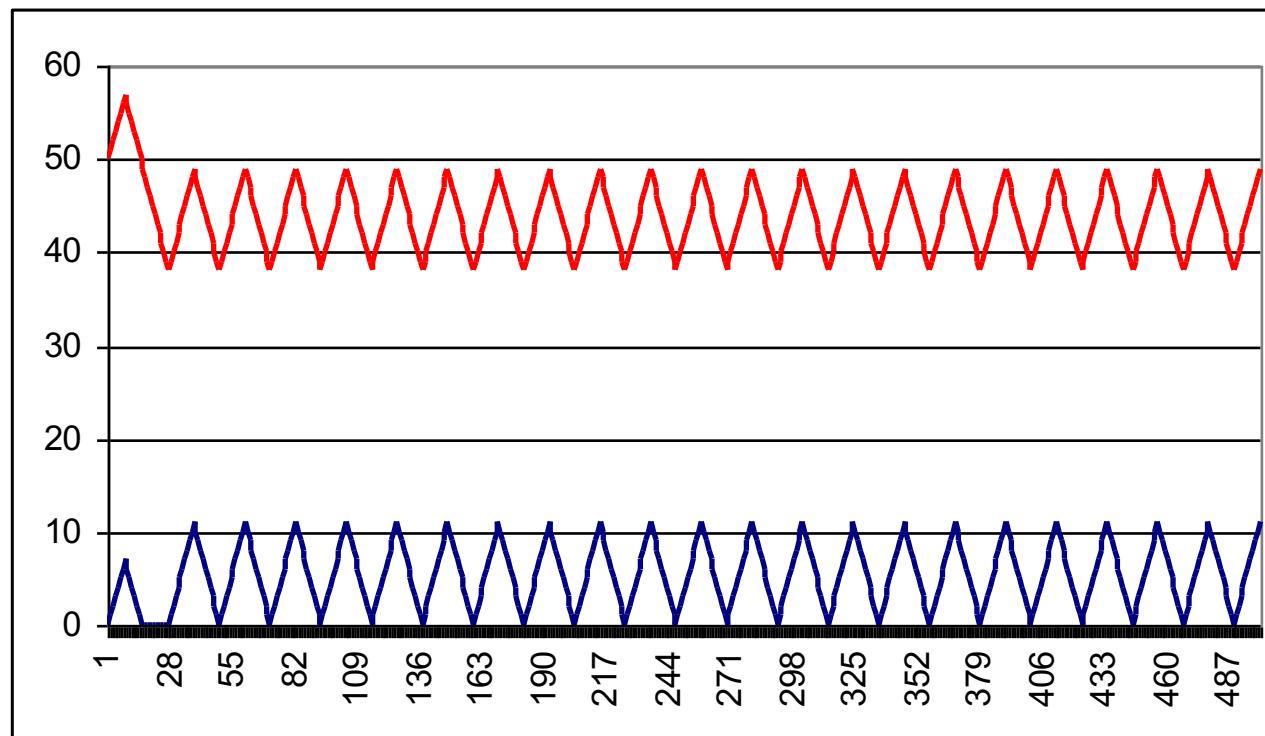
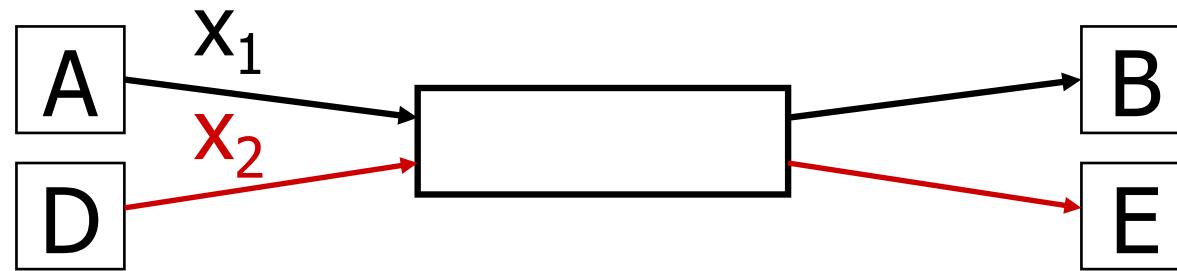


AIAD

- ❖ Increase: $x + a_I$
- ❖ Decrease: $x - a_D$
- ❖ Does not converge to fairness

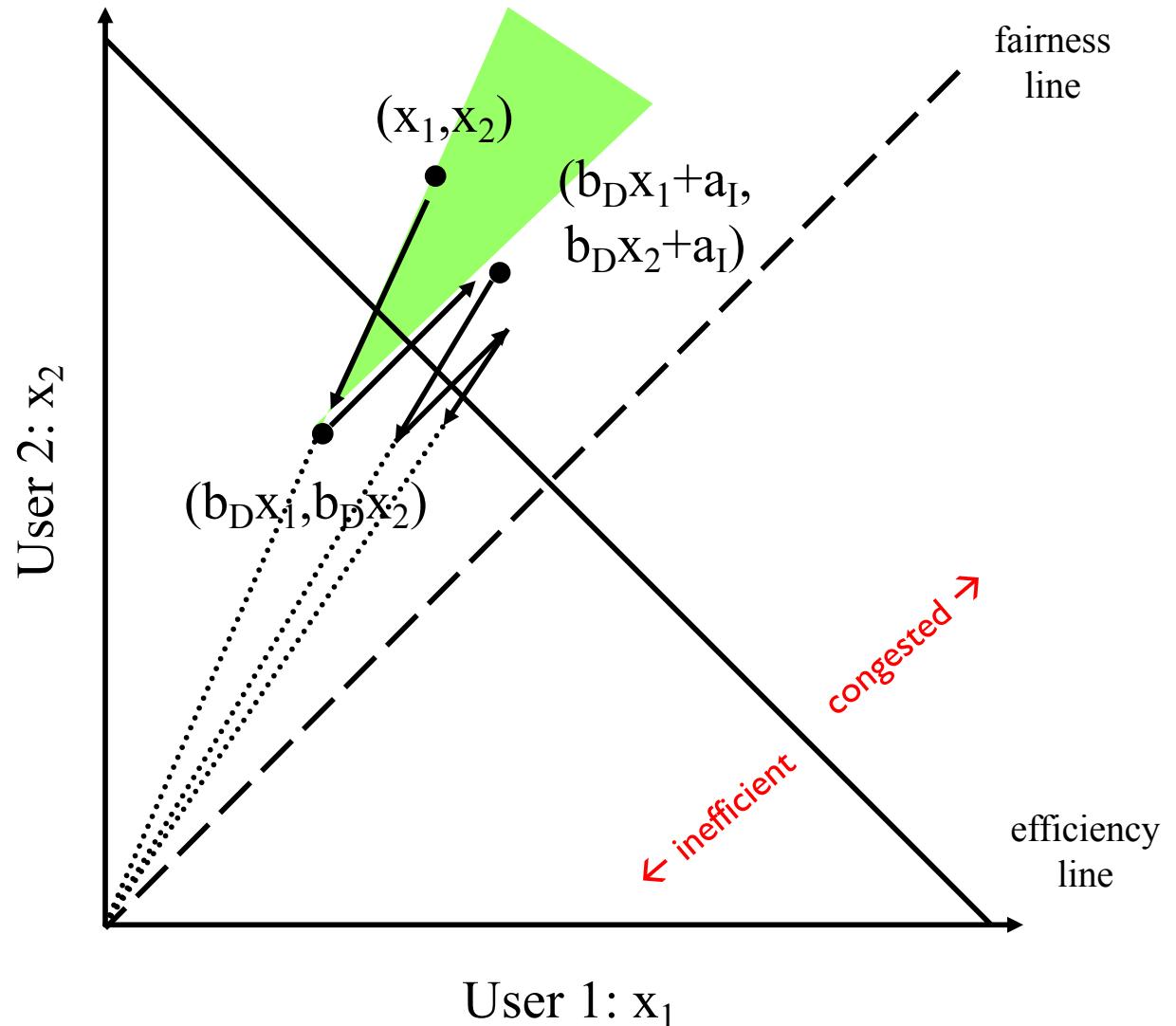


AIAD Sharing Dynamics

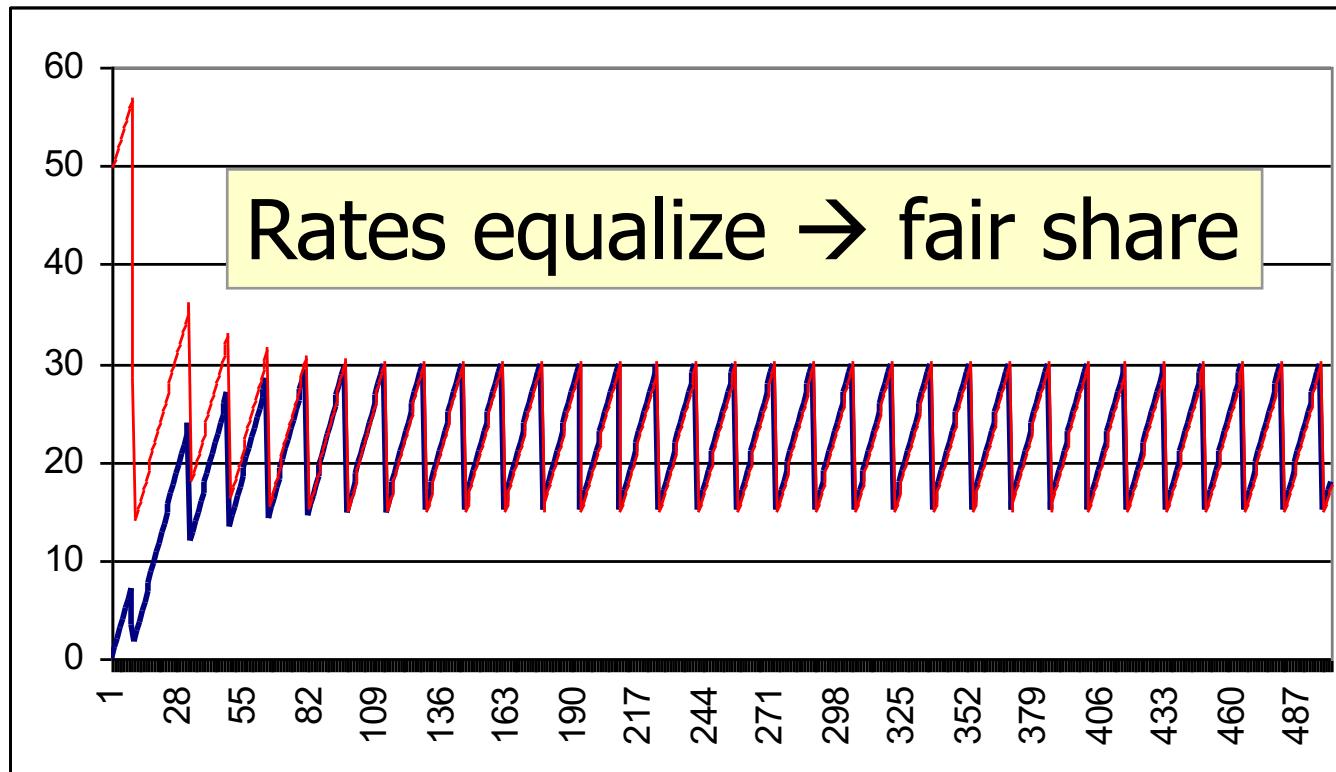


AIMD

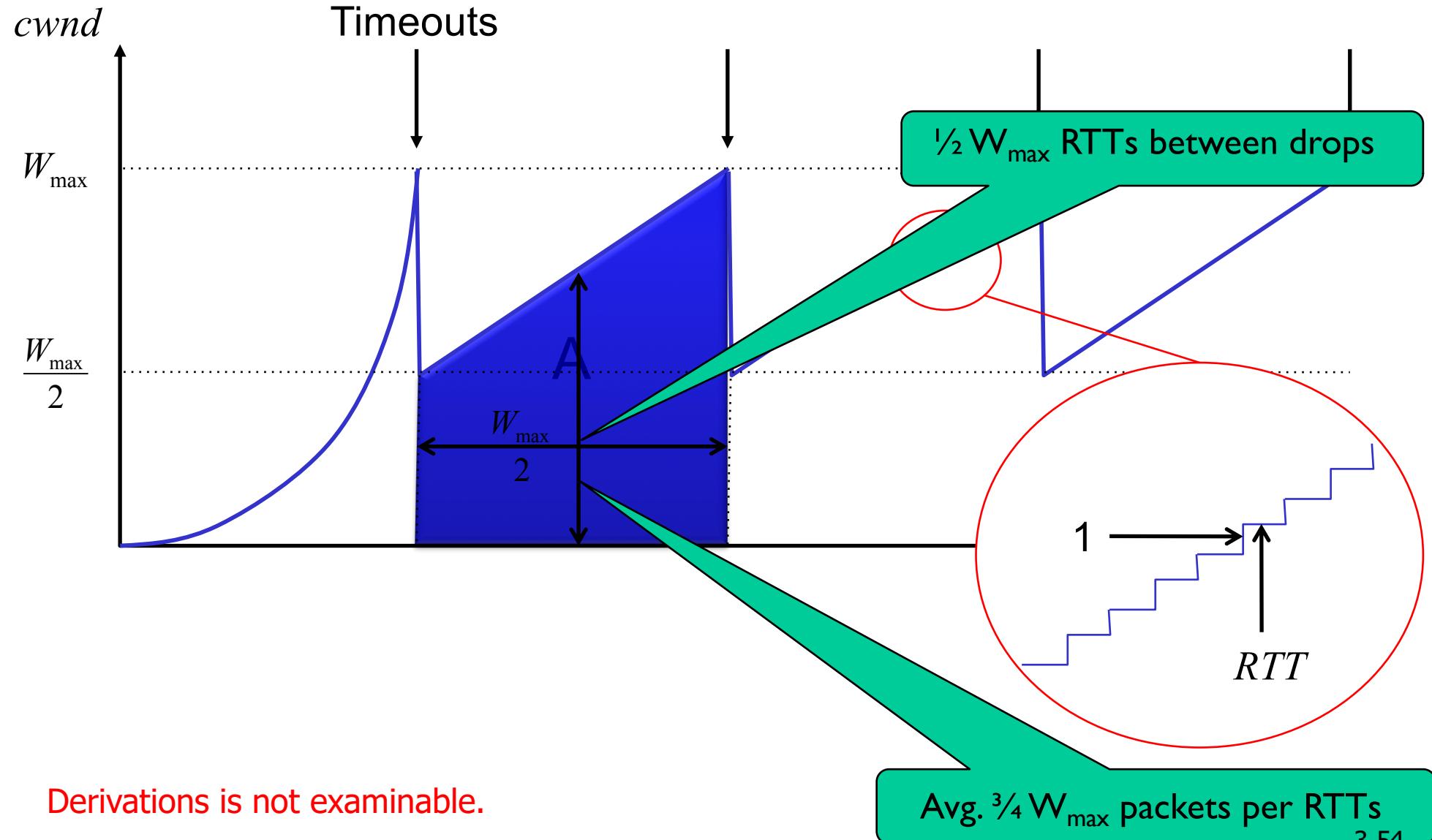
- ❖ Increase: $x+a_I$
- ❖ Decrease: $x*b_D$
- ❖ Converges to fairness



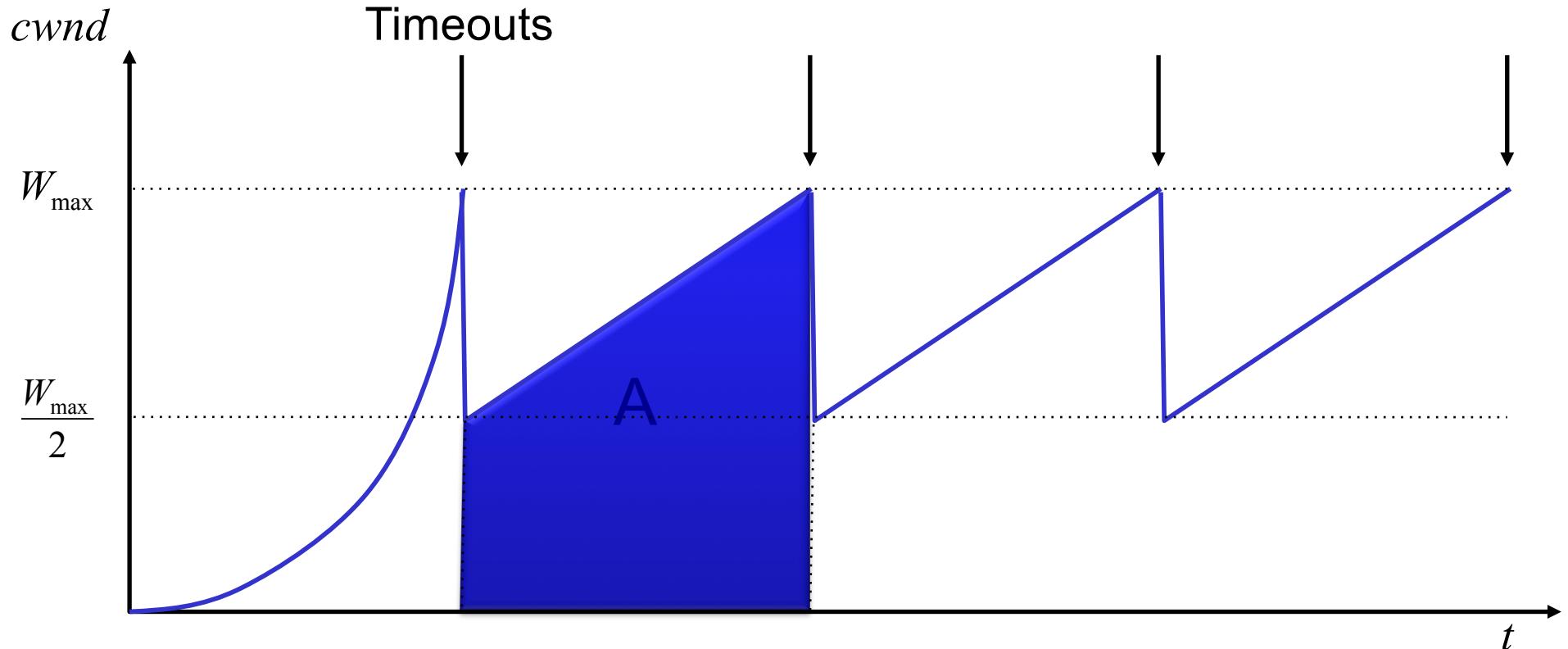
AIMD Sharing Dynamics



A Simple Model for TCP Throughput



A Simple Model for TCP Throughput



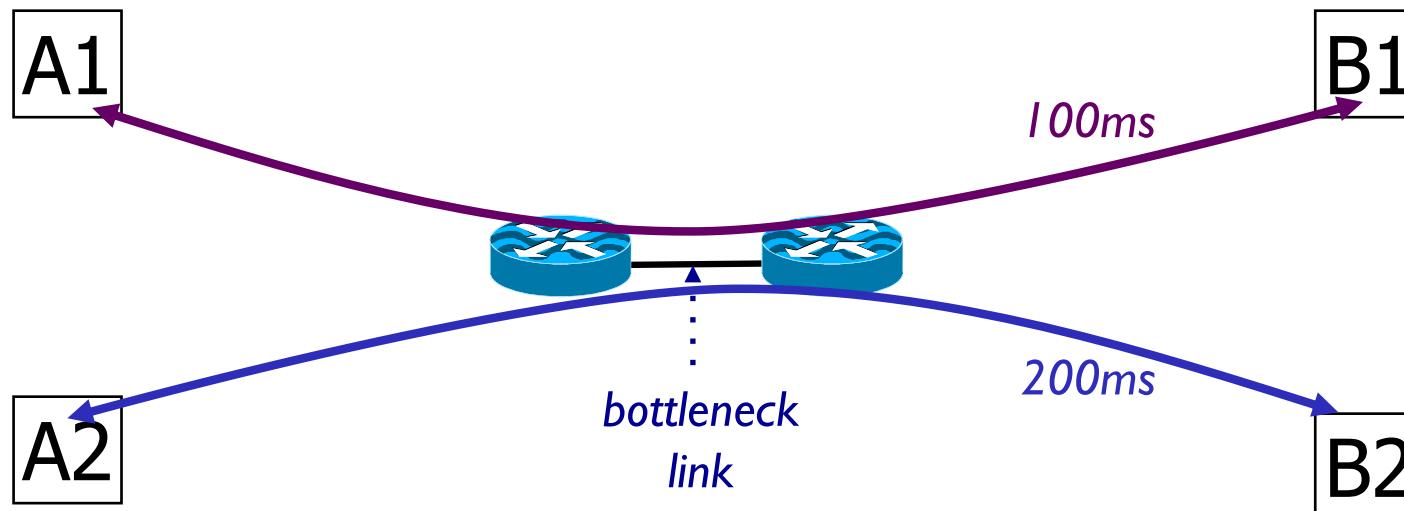
Packet drop rate, $p = 1/A$, where $A = \frac{3}{8}W_{\max}^2$

$$\text{Throughput, } B = \frac{A}{\left(\frac{W_{\max}}{2}\right)RTT} = \sqrt{\frac{3}{2}} \frac{1}{RTT \sqrt{p}}$$

Implications (I): Different RTTs

$$\text{Throughput} = \sqrt{\frac{3}{2}} \frac{1}{RTT\sqrt{p}}$$

- ❖ Flows get throughput inversely proportional to RTT
- ❖ **TCP unfair in the face of heterogeneous RTTs!**



Implications (2): Loss not due to congestion?

- ❖ TCP will confuse corruption with congestion
- ❖ Flow will cut its rate
 - Throughput $\sim 1/\sqrt{p}$ where p is loss prob.
 - Applies even for non-congestion losses!

Implications: (3) How do short flows fare?

- ❖ 50% of flows have < 1500B to send; 80% < 100KB
- ❖ Implication (1): short flows never leave slow start!
 - short flows never attain their fair share
- ❖ Implication (2): too few packets to trigger dupACKs
 - Isolated loss may lead to timeouts
 - At typical timeout values of ~500ms, might severely impact latency

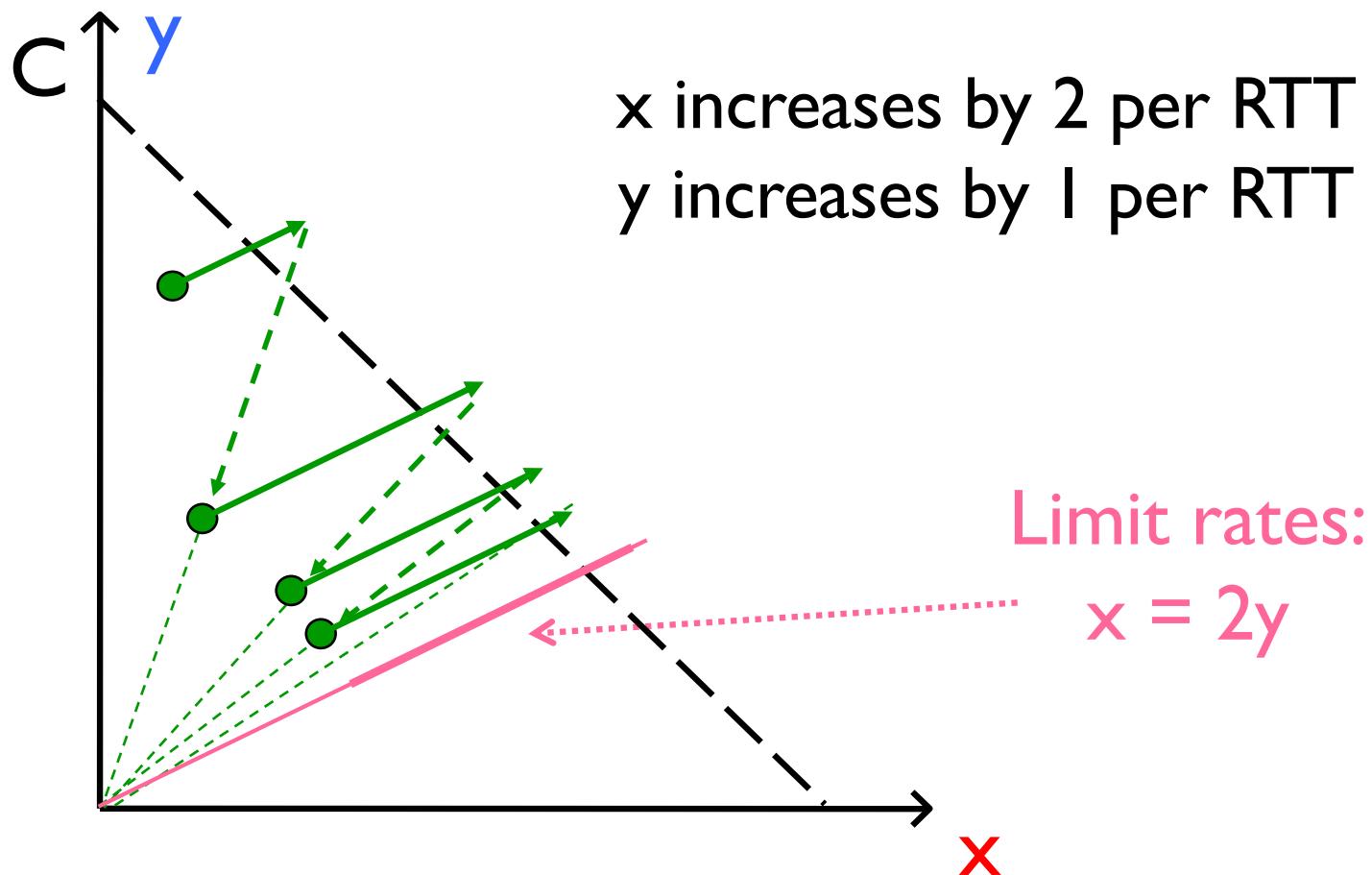
Implications: (4) TCP fills up queues → long delays

- ❖ A flow deliberately overshoots capacity, until it experiences a drop
- ❖ Means that delays are large for everyone
 - Consider a flow transferring a 10GB file sharing a bottleneck link with 10 flows transferring 100B

Implications: (5) Cheating

- ❖ Three easy ways to cheat
 - Increasing CWND faster than +1 MSS per RTT

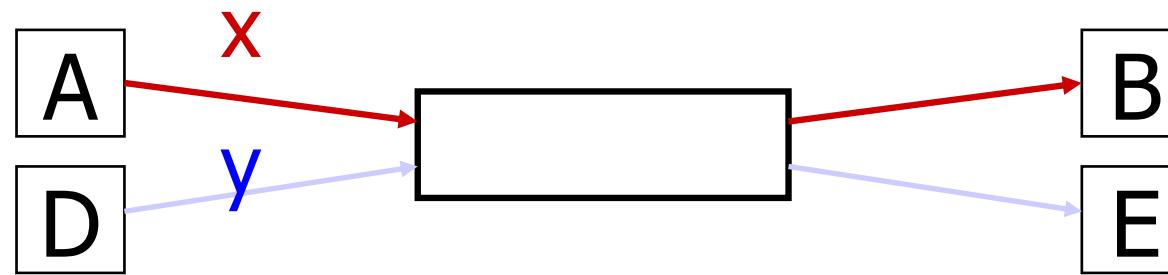
Increasing CWND Faster



Implications: (5) Cheating

- ❖ Three easy ways to cheat
 - Increasing CWND faster than +1 MSS per RTT
 - Opening many connections

Open Many Connections



Assume

- A starts 10 connections to B
- D starts 1 connection to E
- Each connection gets about the same throughput

Then A gets 10 times more throughput than D

Implications: (5) Cheating

- ❖ Three easy ways to cheat
 - Increasing CWND faster than +1 MSS per RTT
 - Opening many connections
 - Using large initial CWND

Transport Layer: Summary

- ❖ principles behind transport layer services:
 - multiplexing, demultiplexing
 - reliable data transfer
 - flow control
 - congestion control
- ❖ instantiation, implementation in the Internet
 - UDP
 - TCP

next:

- ❖ leaving the network “edge” (application, transport layers)
- ❖ into the network “core”