# Network Transport Layer: AIMD; TCP/Reno

Qiao Xiang, Congming Gao, Qiang Su

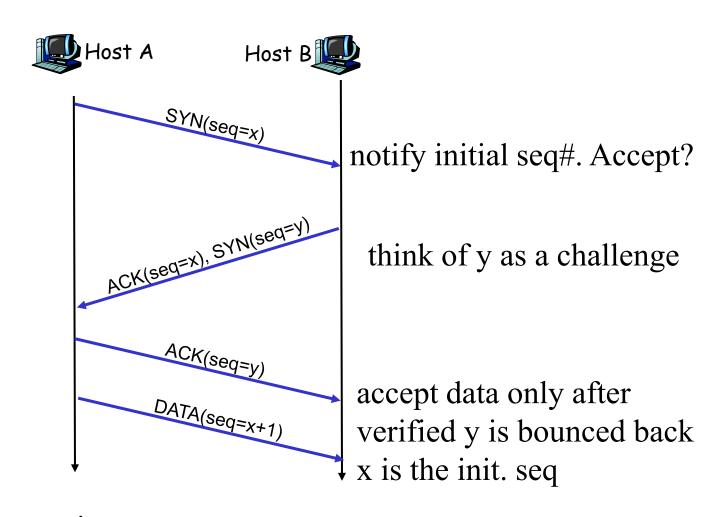
https://sngroup.org.cn/courses/cnnsxmuf25/index.shtml

11/06/2025

### Recap: Transport Design

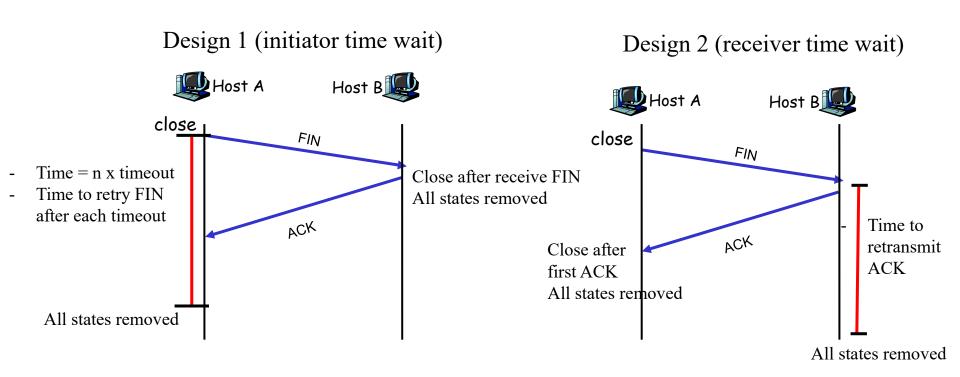
- Basic structure/reliability: sliding window protocols
- □ Determine the "right" parameters
  - Timeout
    - o mean + variation
  - Sliding window size
    - Related w/ congestion control or more generally resource allocation
      - Bad congestion control can lead to congestion collapse (e.g., zombie packets)
    - Goals: distributed algorithm to achieve fairness and efficiency

#### Three Way Handshake (TWH) [Tomlinson 1975]

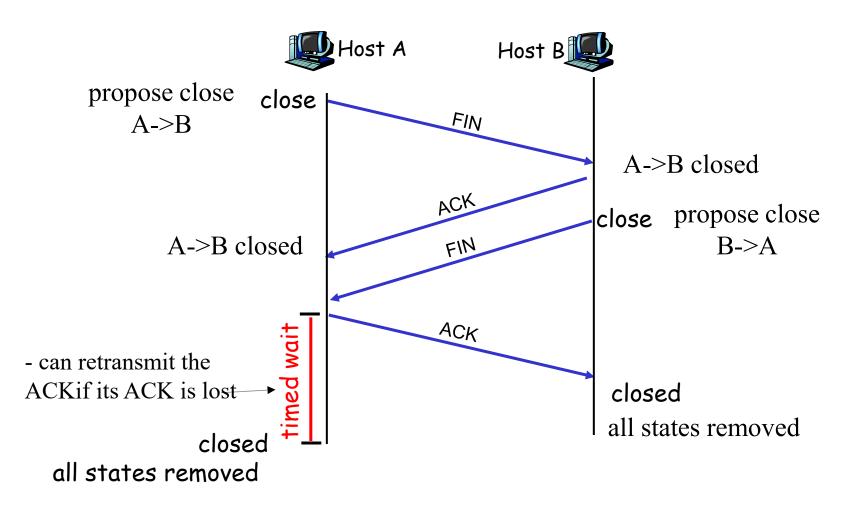


SYN: indicates connection setup

# Time\_Wait Design Options



# TCP Four Way Teardown (For Bi-Directional Transport)



#### Sliding Window Size Function: Rate Control

Transmission rate determined by congestion window size, cwnd, over segments:



cwnd segments, each with MSS bytes sent in one RTT:

Rate = 
$$\frac{\text{cwnd} * MSS}{\text{RTT}}$$
 Bytes/sec

### Some General Questions

#### Big picture question:

□ How to determine a flow's sending rate?

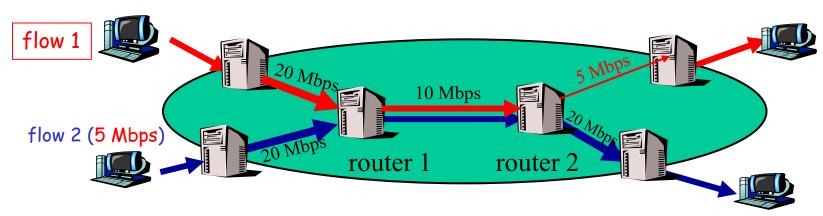
For better understanding, we need to look at a few basic questions:

- What is congestion (cost of congestion)?
- Why are desired properties of congestion control?

# Outline

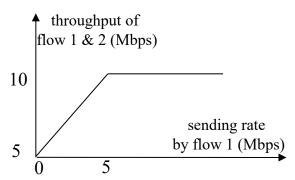
- Admin and recap
- □ TCP Reliability
- Transport congestion control
  - > what is congestion (cost of congestion)

#### Cause/Cost of Congestion: Single Bottleneck

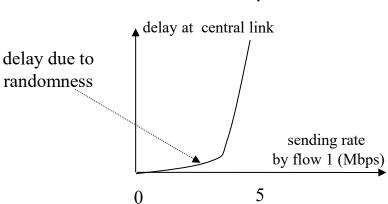


- Flow 2 has a fixed sending rate of 5 Mbps
- We vary the sending rate of flow 1 from 0 to 20 Mbps
- Assume
  - no retransmission; link from router 1 to router 2 has infinite buffer

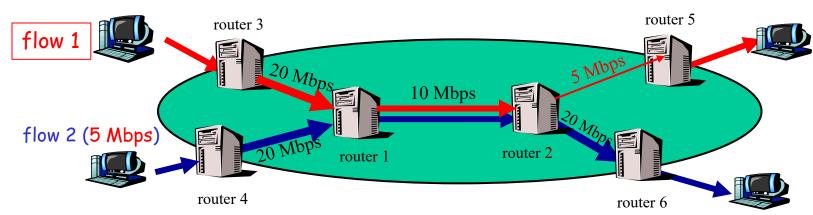
# throughput: e2e packets delivered in unit time





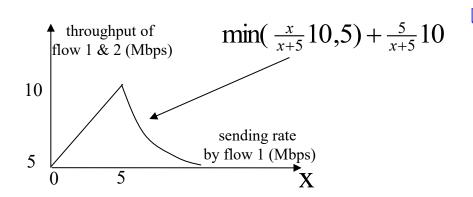


#### Cause/Cost of Congestion: Single Bottleneck



#### □ Assume

- no retransmission
- the link from router 1 to router 2 has finite buffer
- throughput: e2e packets delivered in unit time



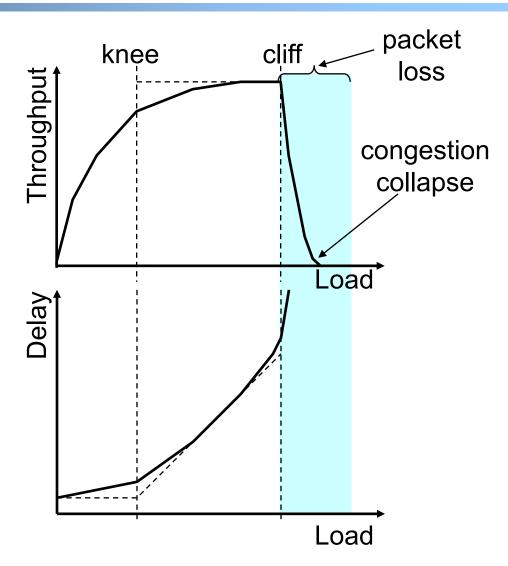
Zombie packet: a packet dropped at the link from router 2 to router 5; the upstream transmission from router 1 to router 2 used for that packet was wasted!

### Summary: The Cost of Congestion

When sources sending rate too high for the network to handle":

- □ Packet loss =>
  - wasted upstream bandwidth when a pkt is discarded at downstream
  - wasted bandwidth due to retransmission (a pkt goes through a link multiple times)





# Outline

- Admin and recap
- □ TCP Reliability
- Transport congestion control
  - what is congestion (cost of congestion)
  - basic congestion control alg.

#### Rate-based vs. Window-based

#### Rate-based:

- Congestion control by explicitly controlling the sending rate of a flow, e.g., set sending rate to 128Kbps
- □ Example: ATM

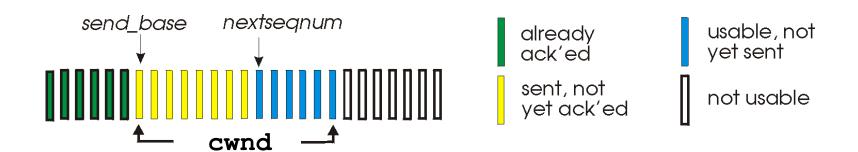
#### Window-based:

- □ Congestion control by controlling the window size of a sliding window, e.g., set window size to 64KBytes
- Example: TCP

Discussion: rate-based vs. window-based

#### Sliding Window Size Function: Rate Control

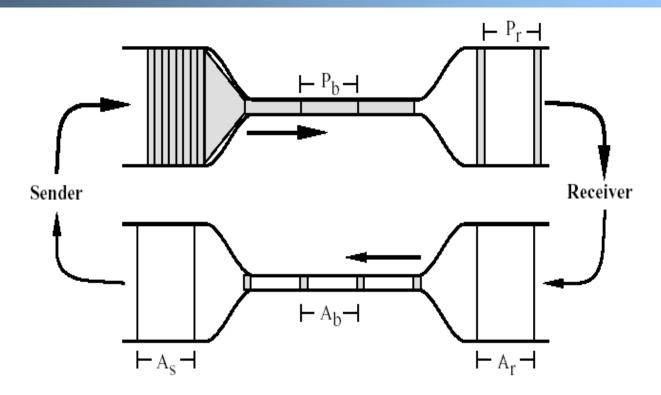
□ Transmission rate determined by congestion window size, cwnd, over segments:



cwnd segments, each with MSS bytes sent in one RTT:

Rate = 
$$\frac{\text{cwnd * MSS}}{\text{RTT}}$$
 Bytes/sec

#### Window-based Congestion Control

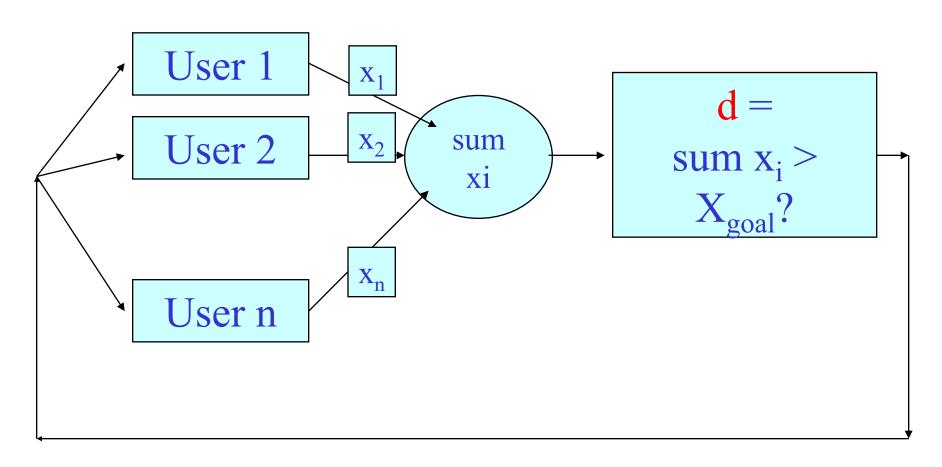


- Window-based congestion control is self-clocking: considers flow conservation, and adjusts to RTT variation automatically.
- □ Hence, for better safety, more designs use window-based design.

# The Desired Properties of a Congestion Control Scheme

- Efficiency: close to full utilization but low delay
  - fast convergence after disturbance
- □ Fairness (resource sharing)
- Distributedness (no central knowledge for scalability)

# Derive CC: A Simple Model



Flows observe congestion signal d, and locally take actions to adjust rates.

## Linear Control

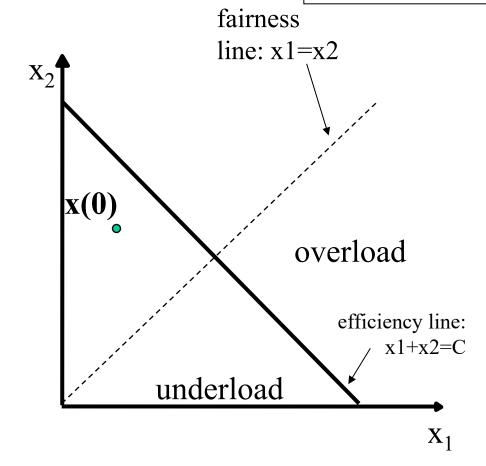
- □ Proposed by Chiu and Jain (1988)
- □ The simplest control strategy

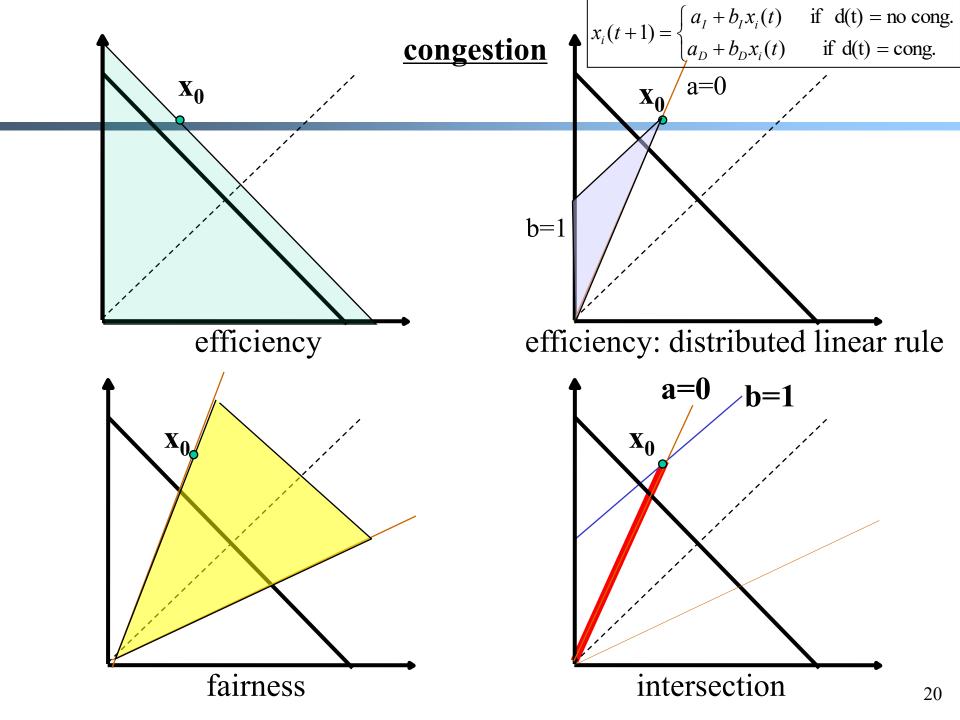
$$x_i(t+1) = \begin{cases} a_I + b_I x_i(t) & \text{if } d(t) = \text{no cong.} \\ a_D + b_D x_i(t) & \text{if } d(t) = \text{cong.} \end{cases}$$

Discussion: values of the parameters?

# State Space of Two Flows

$$x_i(t+1) = \begin{cases} a_I + b_I x_i(t) & \text{if } d(t) = \text{no cong.} \\ a_D + b_D x_i(t) & \text{if } d(t) = \text{cong.} \end{cases}$$





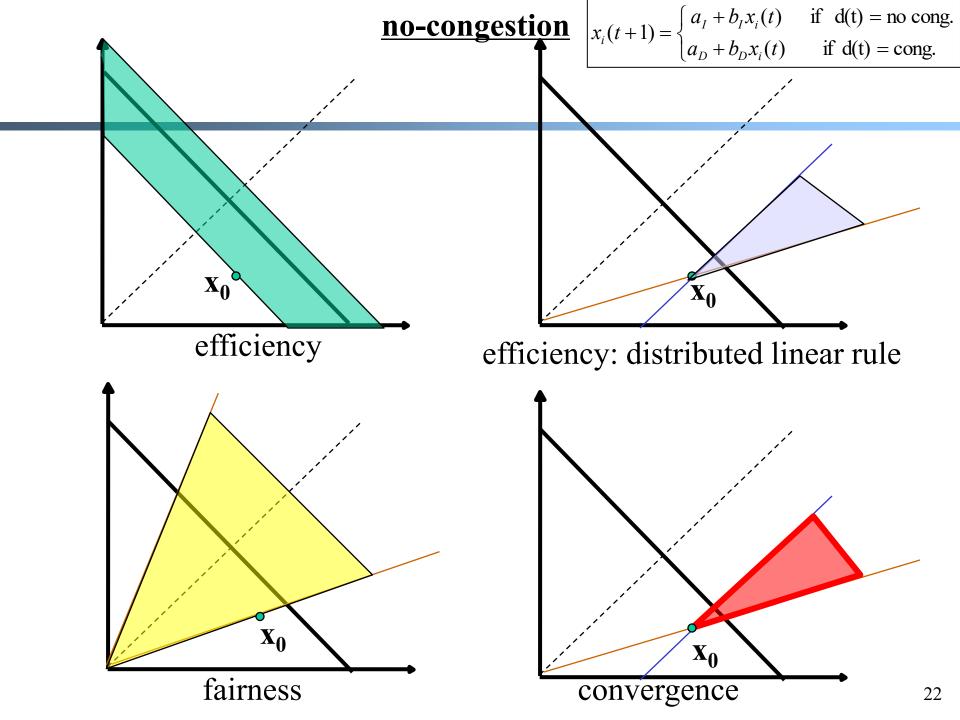
### Implication: Congestion (overload) Case

□ In order to get closer to efficiency and fairness after each update, decreasing of rate must be multiplicative decrease (MD)

$$\circ$$
  $a_D = 0$ 

b<sub>D</sub> < 1</li>

$$x_i(t+1) = \begin{cases} a_I + b_I x_i(t) & \text{if } d(t) = \text{no cong.} \\ b_D x_i(t) & \text{if } d(t) = \text{cong.} \end{cases}$$



# Implication: No Congestion Case

- In order to get closer to efficiency and fairness after each update, additive and multiplicative increasing (AMI), i.e.,
  - $\circ$   $a_{I} > 0$ ,  $b_{I} > 1$

$$x_i(t+1) = \begin{cases} a_I + b_I x_i(t) & \text{if } d(t) = \text{no cong.} \\ b_D x_i(t) & \text{if } d(t) = \text{cong.} \end{cases}$$

- Simply additive increase gives better improvement in fairness (i.e., getting closer to the fairness line)
- □ Multiplicative increase may grow faster

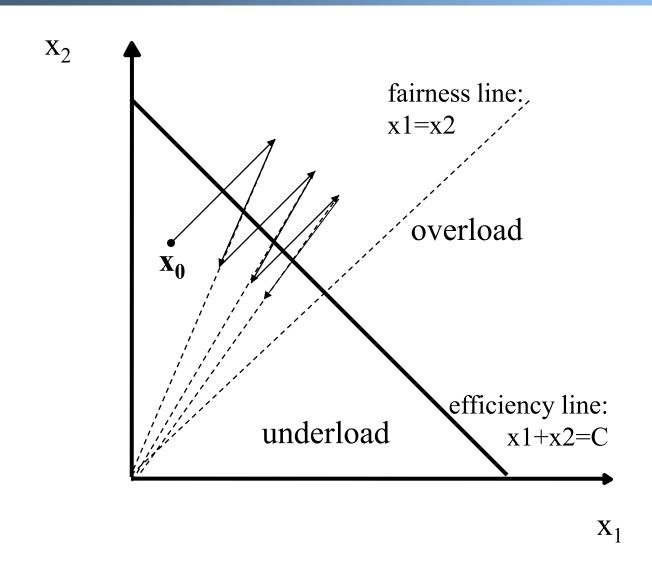
# Intuition: State Trace Analysis of Four Special Cases

	<u>A</u> dditive <u>D</u> ecrease	<u>M</u> ultiplicative <u>D</u> ecrease
<u>A</u> dditive <u>I</u> ncrease	AIAD (b <sub>I</sub> =b <sub>D</sub> =1)	AIMD ( $b_I=1$ , $a_D=0$ )
<u>M</u> ultiplicative <u>I</u> ncrease	MIAD (a <sub>I</sub> =0, b <sub>I</sub> >1, b <sub>D</sub> =1)	$\begin{array}{c} MIMD \\ (a_I=a_D=0) \end{array}$

$$x_i(t+1) = \begin{cases} a_I + b_I x_i(t) & \text{if } d(t) = \text{no cong.} \\ a_D + b_D x_i(t) & \text{if } d(t) = \text{cong.} \end{cases}$$

Discussion: state transition trace.

# AIMD: State Transition Trace



### Intuition: Another Look

- Consider the difference or ratio of the rates of two flows
  - AIAD
    - o difference does not change
  - MIMD
    - o ratio does not change
  - MIAD
    - o difference becomes bigger
  - AIMD
    - o difference does not change

# Outline

- Admin and recap
- □ TCP Reliability
- Transport congestion control
  - what is congestion (cost of congestion)
  - basic congestion control alg.
  - > TCP/reno congestion control

# TCP Congestion Control

- Closed-loop, end-to-end, window-based congestion control
- Designed by Van Jacobson in late 1980s, based on the AIMD alg. of Dah-Ming Chu and Raj Jain
- Worked in a large range of bandwidth values: the bandwidth of the Internet has increased by more than 200,000 times
- Many versions
  - TCP/Tahoe: this is a less optimized version
  - TCP/Reno: many OSs today implement Reno type congestion control
  - TCP/Vegas: not currently used

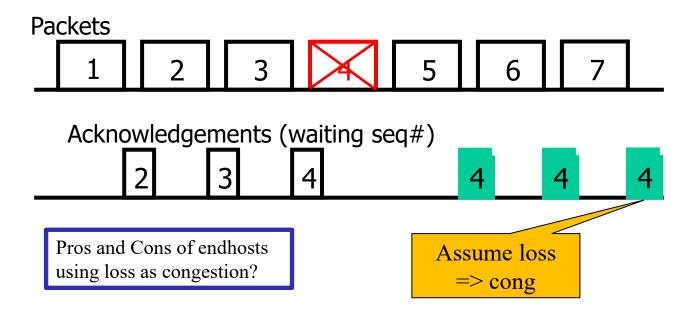
For more details: see TCP/IP illustrated; or read http://lxr.linux.no/source/net/ipv4/tcp\_input.c for linux implementation

### Mapping A(M)I-MD to Protocol

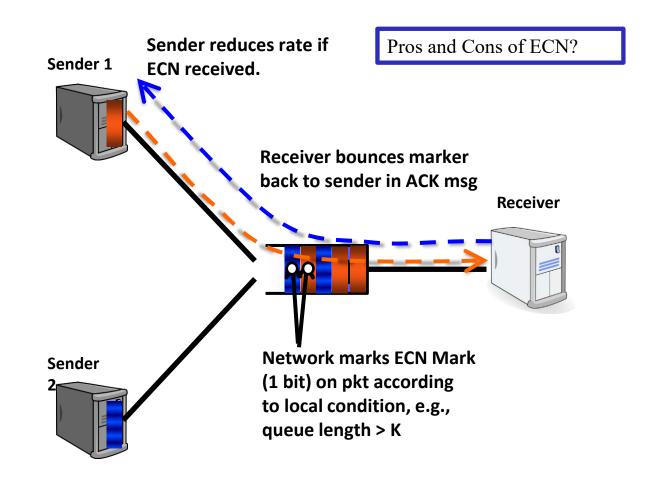
- Basic questions to look at:
  - How to obtain d(t)--the congestion signal?
  - What values do we choose for the formula?
  - How to map formula to code?

$$x_i(t+1) = \begin{cases} a_I + x_i(t) & \text{if } d(t) = \text{no cong.} \\ b_D x_i(t) & \text{if } d(t) = \text{cong.} \end{cases}$$

# Obtain d(t) Approach 1: End Hosts Consider Loss as Congestion



# Obtain d(t) Approach 2: Network Feedback (FCN: Explicit Congestion Notification)



### Mapping A(M)I-MD to Protocol

- Basic questions to look at:
  - o How to obtain d(t)--the congestion signal?
  - What values do we choose for the formula?
  - How to map formula to code?

$$x_i(t+1) = \begin{cases} a_I + x_i(t) & \text{if } d(t) = \text{no cong.} \\ b_D x_i(t) & \text{if } d(t) = \text{cong.} \end{cases}$$

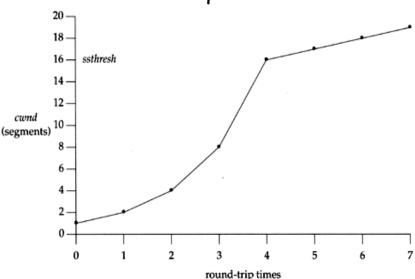
#### TCP/Reno Formulas

- Multiplicative Increase (MI)
  - o double the rate: x(t+1) = 2 x(t)
- Additive Increase (AI)
  - Linear increase the rate: x(t+1) = x(t) + 1
- Multiplicative decrease (MD)
  - half the rate: x(t+1) = 1/2 x(t)

# TCP/Reno Formula Switching (Control Structure)

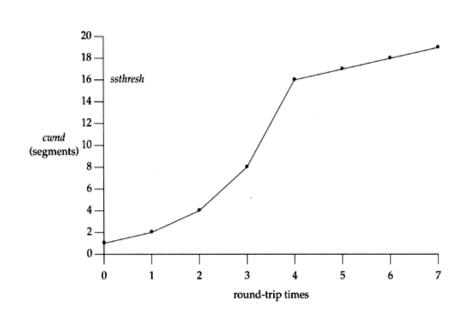
#### ■ Two "phases"

- slow-start
  - Goal: getting to equilibrium gradually but quickly, to get a rough estimate of the optimal of cwnd
  - · Formula: MI
- congestion avoidance
  - Goal: Maintains equilibrium and reacts around equilibrium
  - Formula: AI MD



# TCP/Reno Formula Switching (Control Structure)

- Important variables:
  - cwnd: congestion window size
  - ssthresh: threshold between the slow-start phase and the congestion avoidance phase
- ☐ If cwnd < ssthresh
  - MI
- Else
  - AIMD

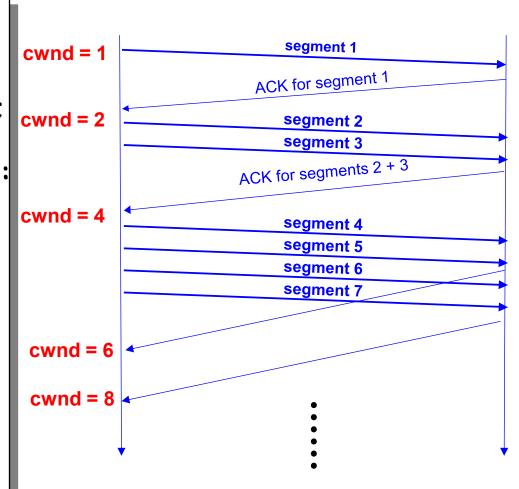


# MI: Slow Start

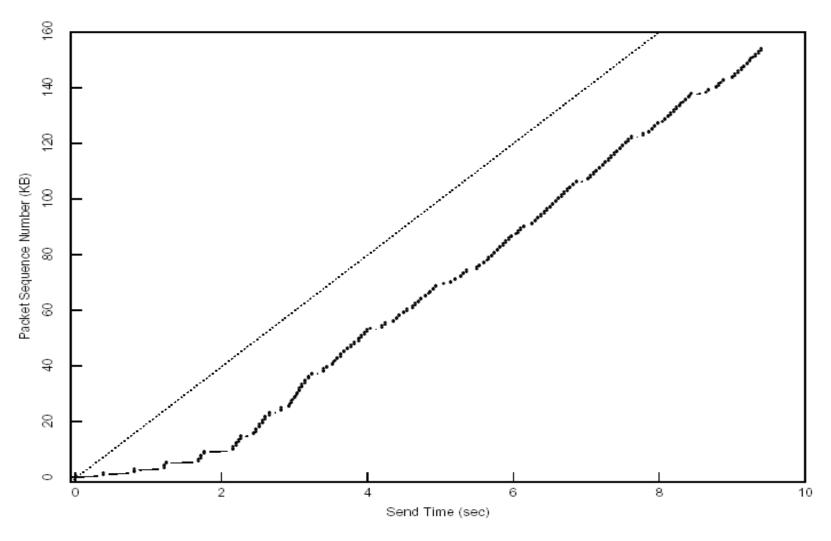
- □ Algorithm: MI
  - o double cwnd every RTT until network congested
- Goal: getting to equilibrium gradually but quickly, to get a rough estimate of the optimal of cwnd

## MI: Slow-start

```
Initially:
   cwnd = 1;
  ssthresh = infinite (e.g., 64K);
For each newly ACKed segment:
   if (cwnd < ssthresh)
     /* MI: slow start*/
     cwnd = cwnd + 1;
```



### Startup Behavior with Slow-start



See [Jac89]

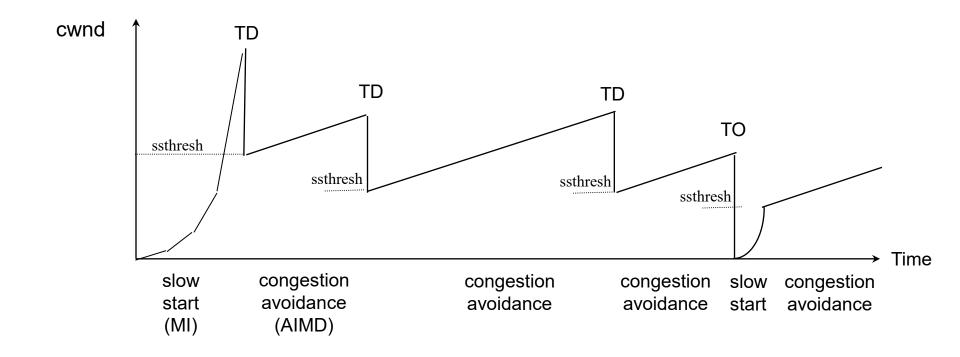
# AIMD: Congestion Avoidance

- □ Algorithm: AIMD
  - increases window by 1 per round-trip time (how?)
  - cuts window size
    - to half when detecting congestion by 3DUP
    - to 1 if timeout
    - if already timeout, doubles timeout
- Goal: Maintains equilibrium and reacts around equilibrium

# TCP/Reno Full Alg

```
Initially:
   cwnd = 1;
   ssthresh = infinite (e.g., 64K);
For each newly ACKed segment:
   if (cwnd < ssthresh) // slow start: MI
     cwnd = cwnd + 1;
   else
                             // congestion avoidance; AI
     cwnd += 1/cwnd:
Triple-duplicate ACKs:
                             //MD
   cwnd = ssthresh = cwnd/2:
Timeout:
  ssthresh = cwnd/2; // reset
   cwnd = 1;
(if already timed out, double timeout value; this is called exponential backoff)
```

# TCP/Reno: Big Picture



TD: Triple duplicate acknowledgements

TO: Timeout