

# TCP seq. numbers, ACKs

## sequence numbers:

- byte stream “number” of first byte in segment’s data

## acknowledgements:

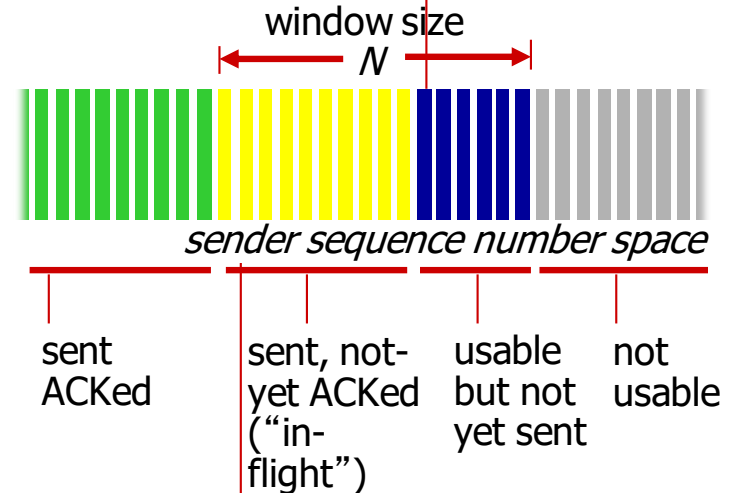
- seq # of next byte expected from other side
- cumulative ACK

**Q:** how receiver handles out-of-order segments

- **A:** TCP spec doesn’t say, - up to implementor

outgoing segment from sender

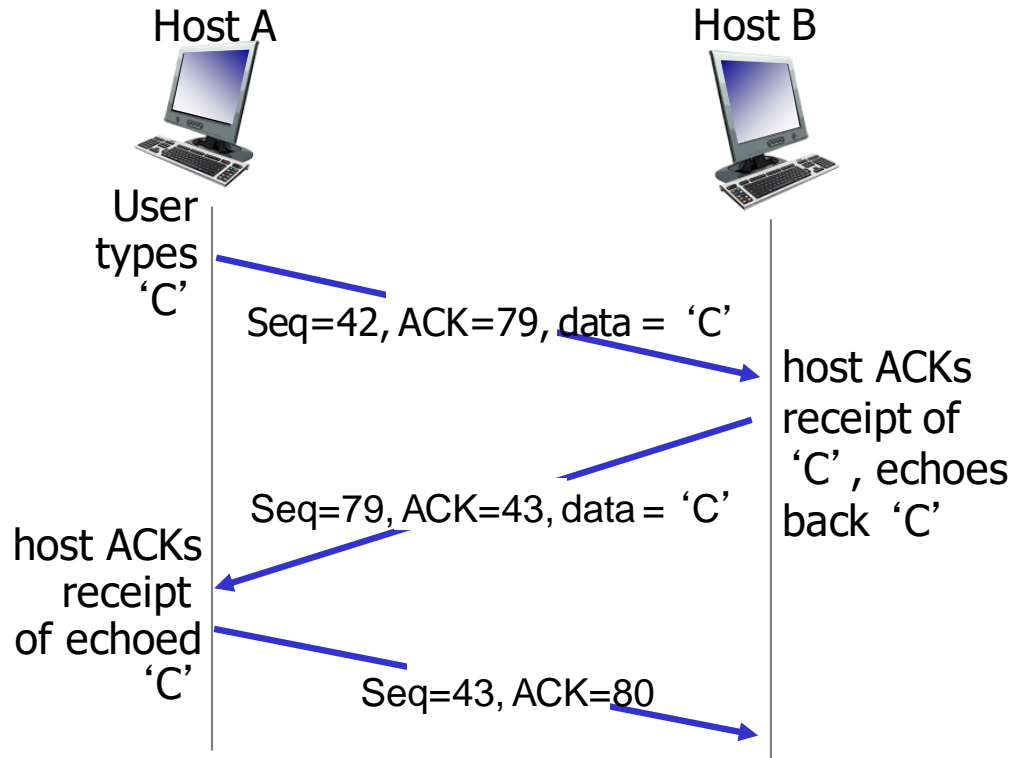
source port #	dest port #
sequence number	
acknowledgement number	
	rwnd
checksum	urg pointer



incoming segment to sender

source port #	dest port #
sequence number	
acknowledgement number	
	A
checksum	urg pointer

# TCP seq. numbers, ACKs



simple telnet scenario

# TCP round trip time, timeout

- ❖ **timeout interval:** **EstimatedRTT** plus “safety margin”
  - large variation in **EstimatedRTT** -> larger safety margin
- ❖ estimate **SampleRTT** deviation from **EstimatedRTT**:

$$\text{DevRTT} = (1-\beta) * \text{DevRTT} + \beta * |\text{SampleRTT} - \text{EstimatedRTT}|$$

(typically,  $\beta = 0.25$ )

$$\text{TimeoutInterval} = \text{EstimatedRTT} + 4 * \text{DevRTT}$$

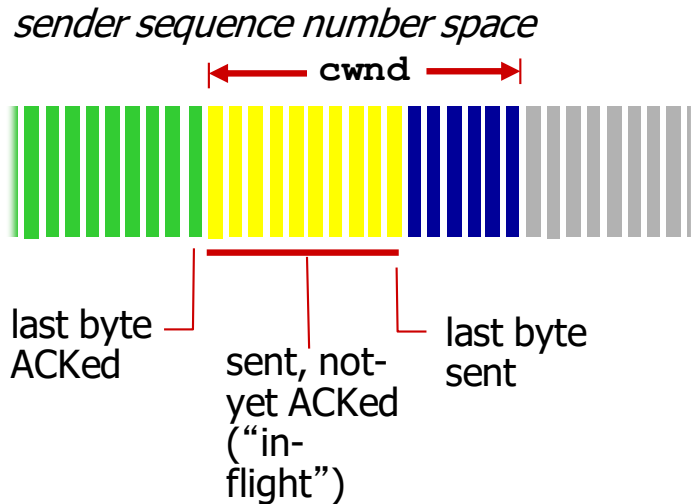


↑  
estimated RTT

↑  
“safety margin”

Retransmissions excluded from TimeoutInterval computation

# TCP Congestion Control: details



- ❖ sender limits transmission:

$$\text{LastByteSent} - \text{LastByteAked} \leq \text{cwnd}$$

- ❖ **cwnd** is dynamic, function of perceived network congestion

TCP sending rate:

- ❖ roughly: send cwnd bytes, wait RTT for ACKS, then send more bytes

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

# Additive Increase/Multiplicative Decrease

- ❖ Objective: adjust to changes in the available capacity
- ❖ New state variable per connection: **CongestionWindow**
  - limits how much data source has in transit

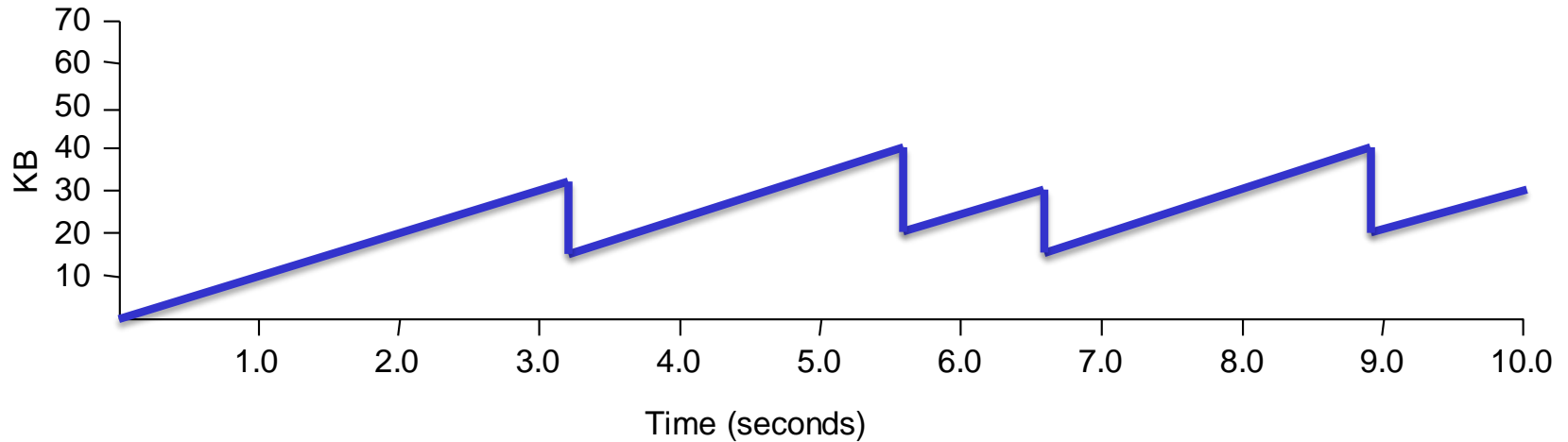
**MaxWin = MIN(CongestionWindow, AdvertisedWindow)**

**EffWin = MaxWin - (LastByteSent - LastByteAcked)**

- ❖ Idea:
  - increase **CongestionWindow** when congestion goes down
  - decrease **CongestionWindow** when congestion goes up

# AIMD (cont)

❖ Trace: **sawtooth shape** behavior





# SSTHRESH and CWND

- ❖ SSTHRESH typically very large on connection setup
- ❖ Set to one half of `CongestionWindow` on packet loss
  - So, SSTHRESH goes through multiplicative decrease for each packet loss
  - If loss is indicated by timeout, set `CongestionWindow = 1`
    - SSTHRESH and `CongestionWindow` always  $\geq 1$  MSS
- ❖ After loss, when new data is ACKed, increase CWND
  - Manner depends on whether we're in slow start or congestion avoidance



# Congestion Control Functionality

Until  $\text{cwnd} \leq \text{slow\_start\_threshold}$

## ❖ Slow Start phase (exponential growth)

- Each returning ACK, a new pckt is transmitted
  - $\text{cwnd} \rightarrow \text{cwnd} + 1$
- Every RTT
  - $\text{cwnd} \rightarrow 2 \text{cwnd}$

When  $\text{cwnd} > \text{slow\_start\_threshold}$

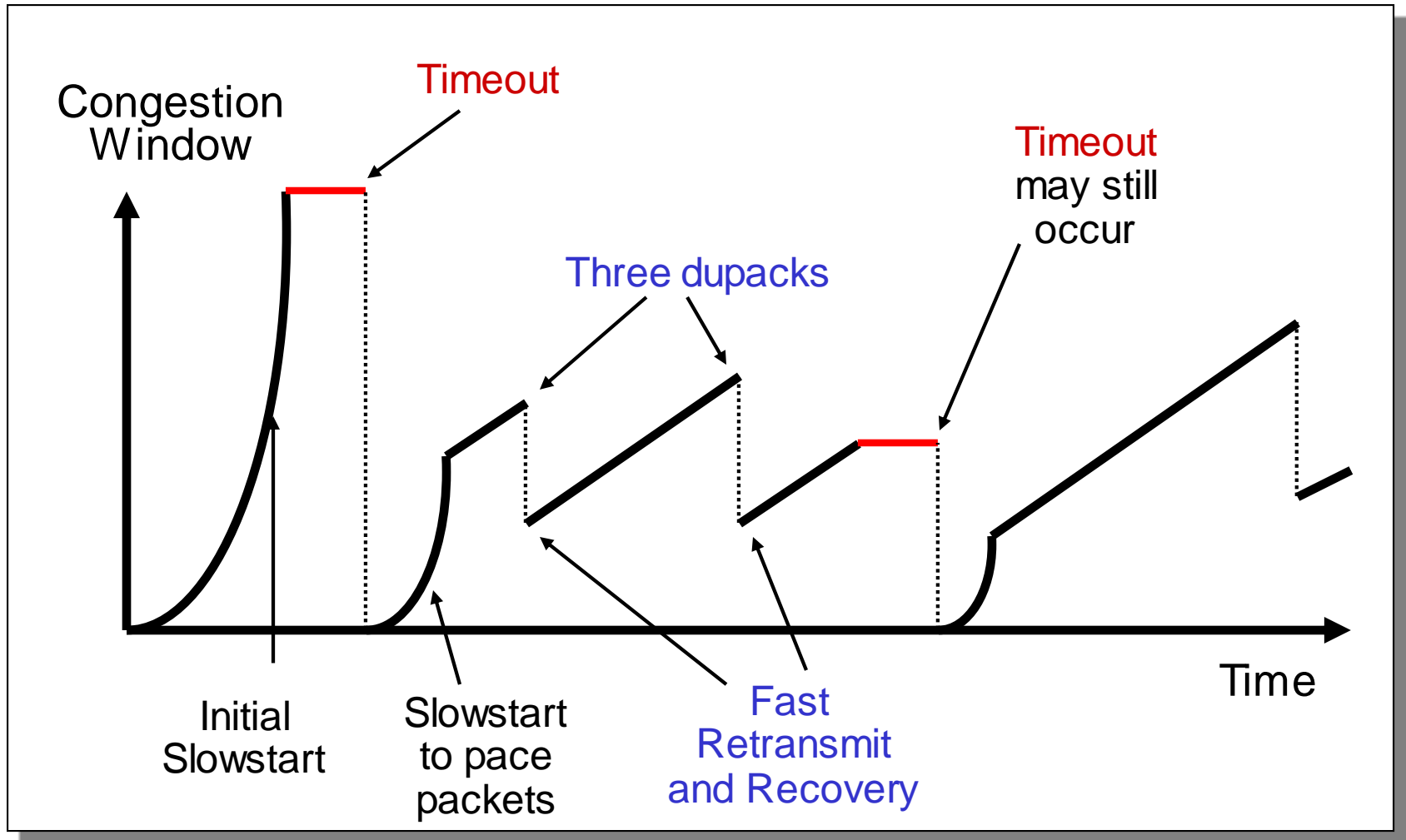
## ❖ Congestion avoidance phase (linear growth)

- Each returning ACK, a new pckt is transmitted
  - $\text{cwnd} \rightarrow \text{cwnd} + (1/\text{cwnd})$
- Every RTT
  - $\text{cwnd} \rightarrow \text{cwnd} + 1$

# Loss recovery

- ❖ Two ways to detect losses
  - Time outs
  - Three dupacks
  
- ❖ With timeout expiration
  - $ssthresh = cwnd / 2$
  - $cwnd = 1$  (so, restart in slow start phase)
  
- ❖ With three dupacks
  - $ssthresh = cwnd / 2$
  - $cwnd = cwnd / 2$  (so, restart in cong. avoidance phase)

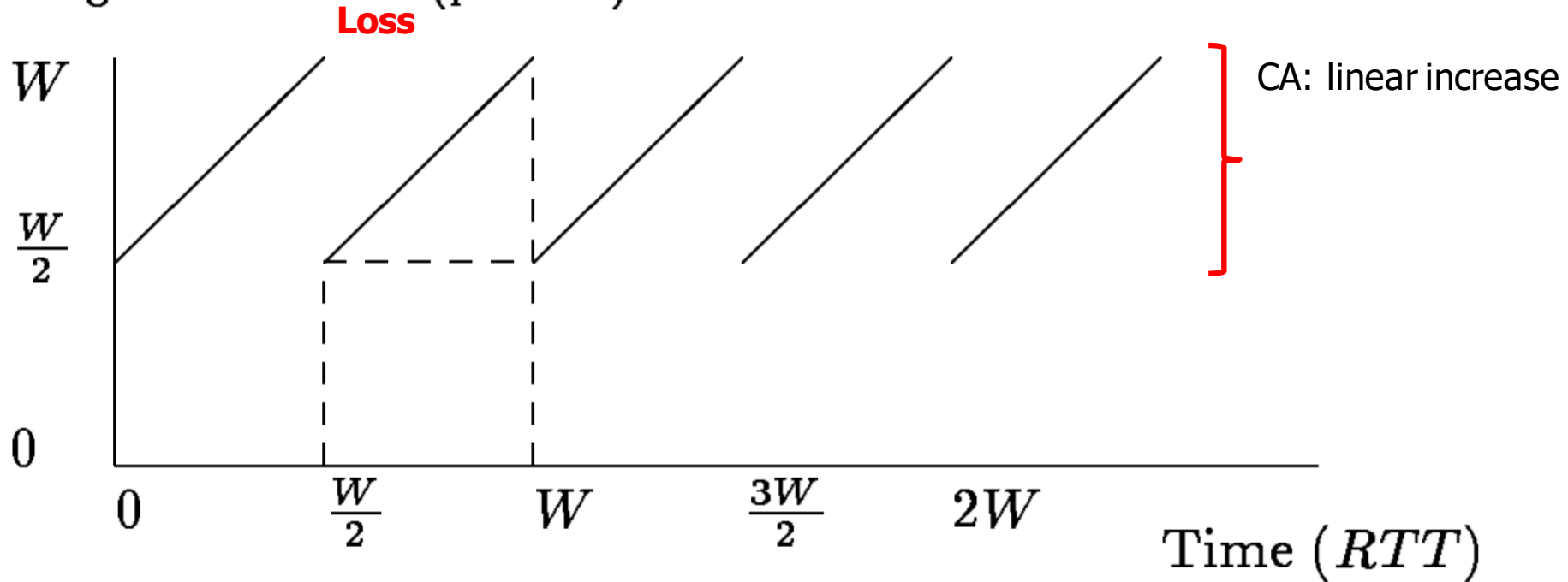
# TCP Saw Tooth



# Mathis et. al, 1997 – Macroscopic TCP Throughput Estimation

$W \rightarrow$  max cwnd  
 $p \rightarrow$  periodic loss prob. at end. cycle  
 $MSS \rightarrow$  Maximum segment size

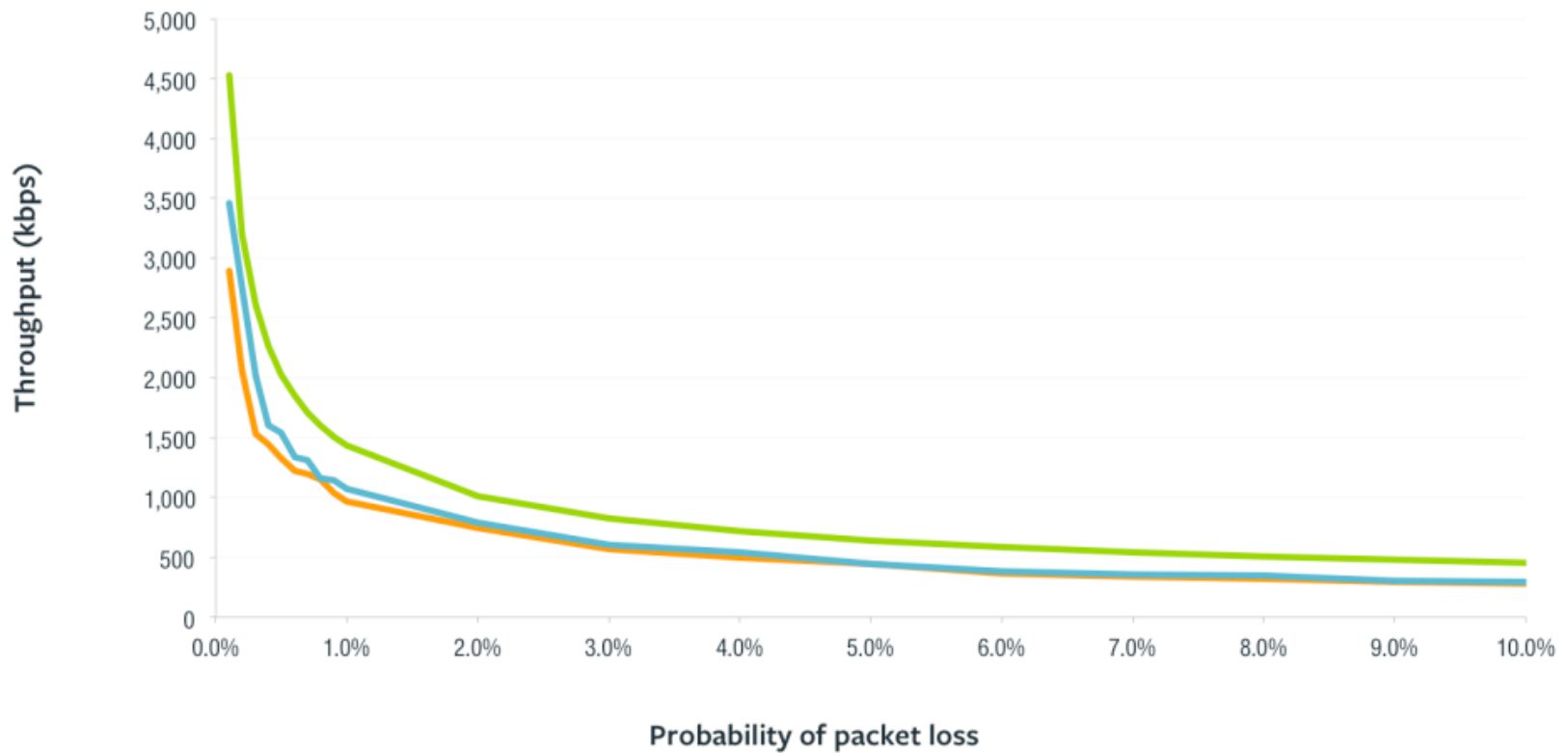
congestion window (packets)



# Model Validation

Throughput vs Loss

Mathis model New Reno CUBIC



# 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

# Case study: ATM ABR congestion control

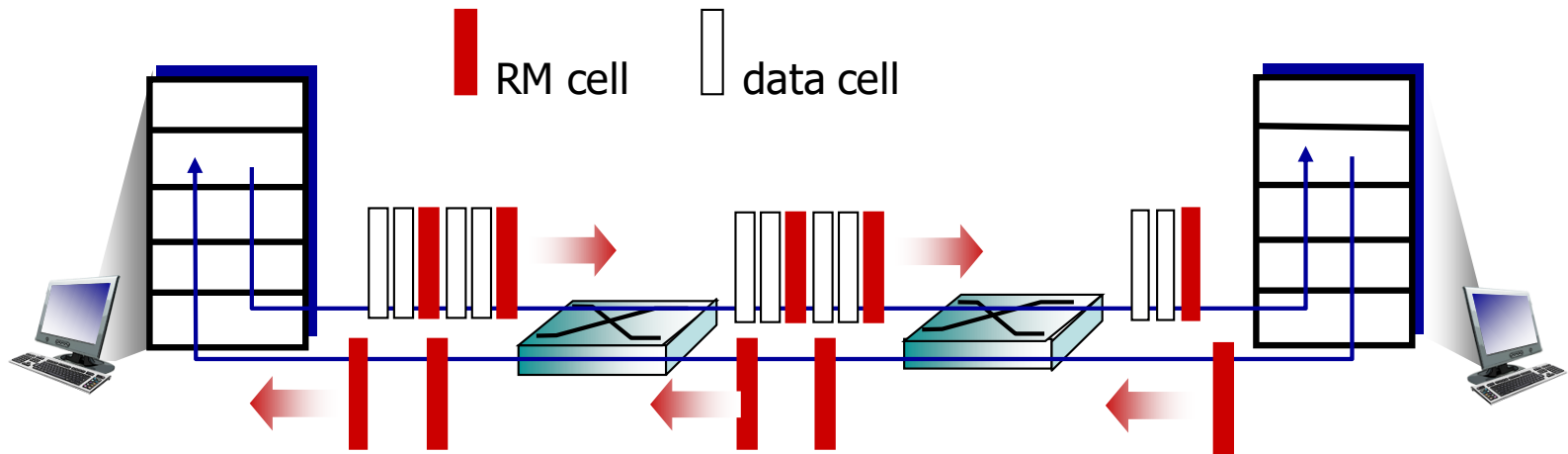
## ABR: available bit rate:

- ❖ “elastic service”
- ❖ if sender’s path “underloaded”:
  - sender should use available bandwidth
- ❖ if sender’s path congested:
  - sender throttled to minimum guaranteed rate

## RM (resource management) cells:

- ❖ sent by sender, interspersed with data cells
- ❖ bits in RM cell set by switches (“*network-assisted*”)
  - *NI bit*: no increase in rate (mild congestion)
  - *CI bit*: congestion indication
- ❖ RM cells returned to sender by receiver, with bits intact

# Case study: ATM ABR congestion control



- ❖ two-byte ER (explicit rate) field in RM cell
  - congested switch may lower ER value in cell
  - senders' send rate thus max supportable rate on path
- ❖ EFCI bit in data cells: set to 1 in congested switch
  - if data cell preceding RM cell has EFCI set, receiver sets CI bit in returned RM cell