TCP congestion control:

- \rightarrow general data transport
- \rightarrow focus: speedy transport of large files

Issue:

- too fast: significant losses lead to degraded performance
 - \rightarrow e.g., TCP congestion collapse
- too slow: underutilize bandwidth
 - \rightarrow especially pronounced in today's high-speed networks

Example: aggravated in data center networks

 \rightarrow called TCP incast

Data center network topology evolved to accommodate horizontal (or east-west traffic)

- Clos/leaf-spine connectivity
 - \rightarrow facilitates east-west traffic flow of data center apps
- application traffic can generate many-to-one simultaneous data transfer
 - \rightarrow hot spot at spine switches
 - \rightarrow significant packet losses and degraded throughput

After aggressive backoff to prevent TCP throughput collapse, want to recover faster than linear increase to utilize high speed links.

 \rightarrow important application domain of congestion control

Interface between TCP sliding window and congestion control:

where

MaxWindow =
min{ AdvertisedWindow, CongestionWindow }

- \rightarrow how to adjust Congestion Window to increase throughput
- \rightarrow prevent congestion collapse
- \rightarrow based on method B

TCP congestion control components:

- (i) Congestion avoidance
 - → linear increase/exponential decrease
 - → additive increase/exponential decrease (AIMD)

As in method B, increase CongestionWindow linearly, but decrease exponentially

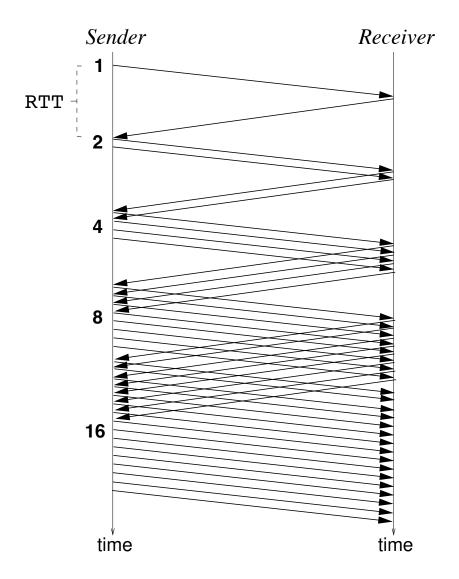
Upon receiving ACK:

 $\label{eq:congestionWindow} \textbf{CongestionWindow} + 1 \\ \textbf{Upon timeout:}$

 $\texttt{CongestionWindow} \leftarrow \texttt{CongestionWindow} / 2$

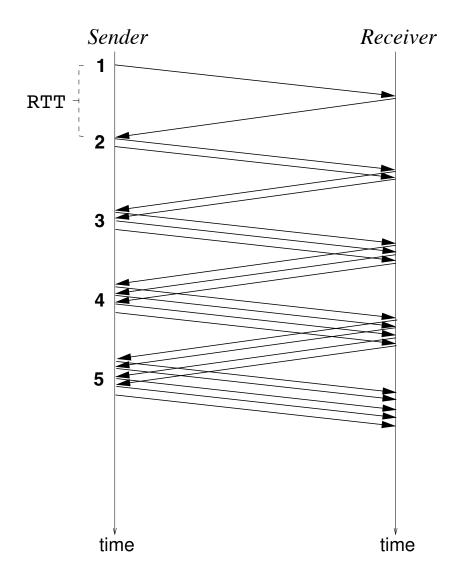
But is it correct...

"Linear increase" time diagram:



 \longrightarrow results in exponential increase

What we want:



 \longrightarrow increase by 1 every window

Thus, linear increase update:

Upon timeout and exponential backoff,

 ${\tt SlowStartThreshold} \, \leftarrow \, {\tt CongestionWindow} \, / \, 2$

(ii) Slow Start

Reset CongestionWindow to 1

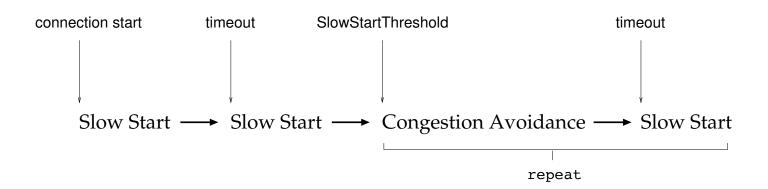
Perform exponential increase

 $\texttt{CongestionWindow} \leftarrow \texttt{CongestionWindow} + 1$

- Until timeout at start of connection
 - \rightarrow rapidly probe for available bandwidth
- Until CongestionWindow hits SlowStartThreshold following Congestion Avoidance
 - \rightarrow rapidly climb to safe level
 - \longrightarrow "slow" is a misnomer
 - \longrightarrow exponential increase is super-fast

Basic dynamics:

- → after connection set-up
- → before connection tear-down

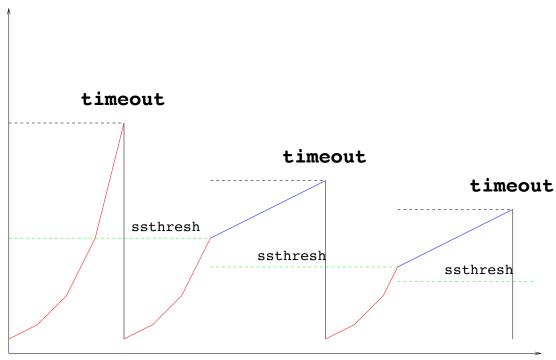


- → many TCP transfers are small
- \longrightarrow small TCP flows don't escape Slow Start

CongestionWindow evolution:

 \longrightarrow relevant for larger flows

CongestionWindow



Events (ACK or timeout)

(iii) Exponential timer backoff

 $TimeOut \leftarrow 2 \cdot TimeOut$ if retransmit

(iv) Fast Retransmit

Upon receiving three duplicate ACKs:

- Transmit next expected segment
 - \rightarrow segment indicated by ACK value
- Perform exponential backoff and commence Slow Start
 - → three duplicate ACKs: likely segment is lost
 - → react before timeout occurs

TCP Tahoe: features (i)-(iv)

(v) Fast Recovery

Upon Fast Retransmit: skip Slow Start, continue sending new data before switching to Congestion Avoidance

- \rightarrow assumption: dup ACKs indicate spurious loss
- \rightarrow can afford to be more aggressive
- \rightarrow avoid stalling

Set CongestionWindow to SlowStartThreshold +3

- account for 3 dup ACKs
 - \rightarrow allows transmission of new data
- keep incrementing CongestionWindow for additional dup ACKs
- upon receiving new ACK transition to Congestion Avoidance
 - → TCP Reno

TCP New Reno:

- stay in Fast Recovery longer
 - \rightarrow handle multiple spurious losses within packet train
 - \rightarrow delay transition to Congestion Avoidance
- selective ACK (SACK)
 - \rightarrow receiver informs sender of contigous blocks of data received
 - \rightarrow i.e., sender able to identify holes
- ... drawback: more complex

Recent variants:

 \rightarrow more aggressive than plain linear increase

 \rightarrow targeted at modern high-speed networks

TCP BIC:

- Upon loss set:
 - $\rightarrow W_{\rm max}$ congestion window size at time of loss
 - $\rightarrow W_{\min}$ congestion window size after multiplicative backoff
- Perform binary search between W_{\min} and W_{\max} to rapidly identify sustainable high throughput window size
 - \rightarrow updated congestion window: $(W_{\min} + W_{\max})/2$
 - \rightarrow upper bound window increase

TCP Cubic:

- \rightarrow based on BIC
- \rightarrow window increase follows cubic function
 - from W_{\min} to W_{\max}
 - \rightarrow concave increase: similar to binary search
 - above W_{max}
 - \rightarrow convex increase: rapid probing similar to slow start
 - congestion window update not determined by RTT
 - \rightarrow use time elapsed since last window reduction
 - \rightarrow akin to open-loop window increase
- ... default TCP congestion control in Linux, MacOS, Windows.

Compound TCP: delay-based

- estimate queueing delay from RTT
 - \rightarrow closer to method D
 - \rightarrow other variants of delay based TCP (e.g., Vegas)
- hybrid: uses sum of two windows
 - \rightarrow delay based window that increases rapidly
 - \rightarrow loss based window following linear increase
 - \rightarrow susceptible to congestion collapse

Used in Windows 10 until transition to TCP Cubic.

Data Center TCP (DCTCP):

- \rightarrow invoke router/switch assistance: AQM
- \rightarrow ECN (explicit congestion notification)
- \rightarrow established at TCP connection set-up
- \rightarrow IPv4: 2-bit ECN in TOS field

Congestion control: selfishness and fairness

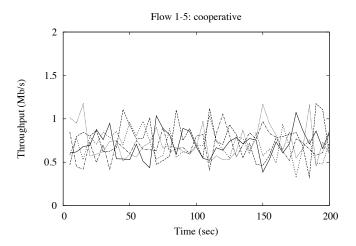
• different versions of TCP co-exist on Internet

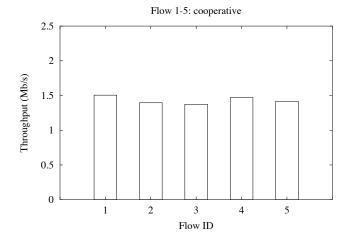
- aggressiveness that increases throughput for one connection may come at the expense of another
- \rightarrow referred to as TCP-friendliness

Example:

- \rightarrow 5 regular (cooperative) TCP flows
- \rightarrow share 11 Mbps WLAN bottleneck link

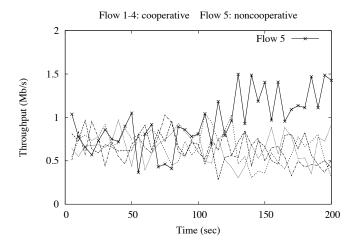
Throughput share of 5 homogenous TCP flows:

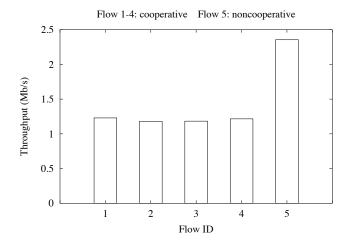




4 regular (cooperative) TCP flows and 1 noncooperative TCP flow:

 \rightarrow starts behaving selfishly at time 100s





Potential for unfairness

- \rightarrow unintentional: TCP friendliness
- \rightarrow intentional: deployment hurdle for noncooperative actor

In general: subject of noncooperative game theory

- \rightarrow e.g., mechanism design
- \rightarrow hurdles to exploiting selfishness
- \rightarrow disincentivize