TCP congestion control:

Recall:

$$\texttt{EffectiveWindow} = \texttt{MaxWindow} -$$
$$(\texttt{LastByteSent} - \texttt{LastByteAcked})$$

where

$$\texttt{MaxWindow} =$$
$$\min\{\,\texttt{AdvertisedWindow},\ \texttt{CongestionWindow}\,\}$$

Key question: how to set `CongestionWindow` which, in turn, affects ARQ's sending rate?

$\longrightarrow$    linear increase/exponential decrease

$\longrightarrow$    AIMD

$\longrightarrow$    method B

TCP congestion control components:

(i) Congestion avoidance

$\longrightarrow$ linear increase/exponential decrease

$\longrightarrow$ additive increase/exponential decrease (AIMD)

As in Method B, increase `CongestionWindow` linearly, but decrease exponentially
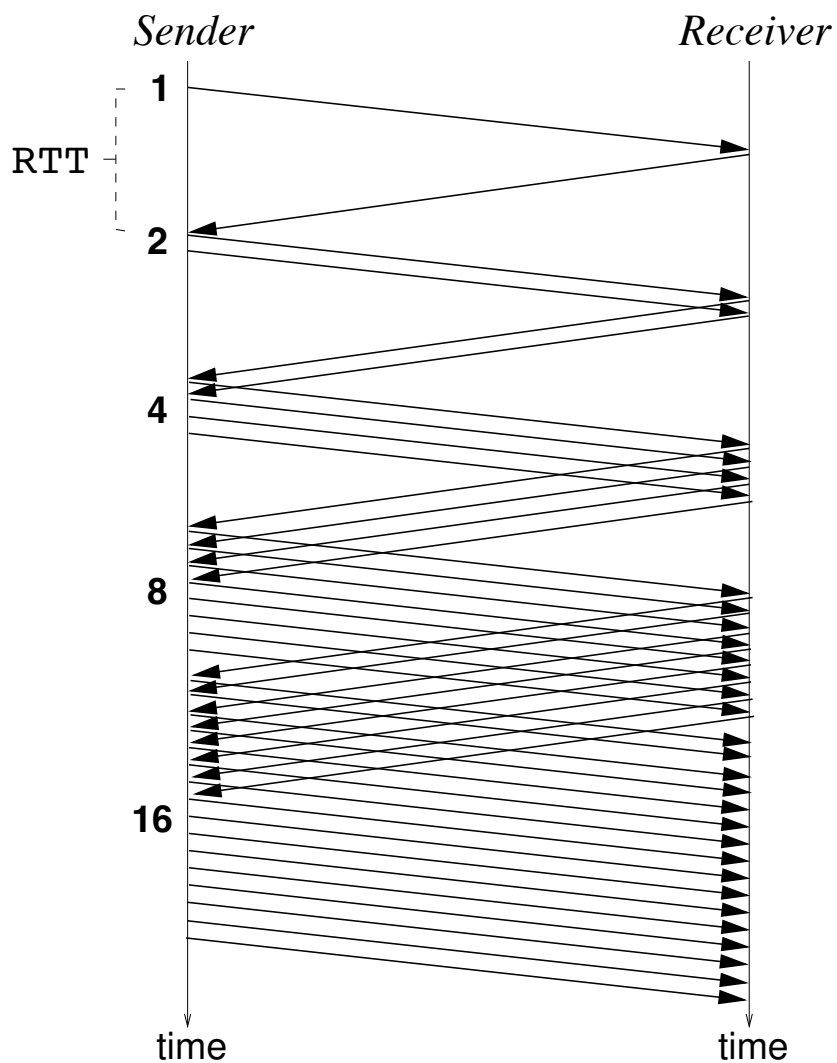
Upon receiving ACK:

`CongestionWindow` $\leftarrow$ `CongestionWindow` $+\,1$

Upon timeout:

`CongestionWindow` $\leftarrow$ `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:

`CongestionWindow` ← `CongestionWindow`
                        + (1 / `CongestionWindow`)

Upon timeout and exponential backoff,

`SlowStartThreshold` ← `CongestionWindow` / 2

(ii) Slow Start

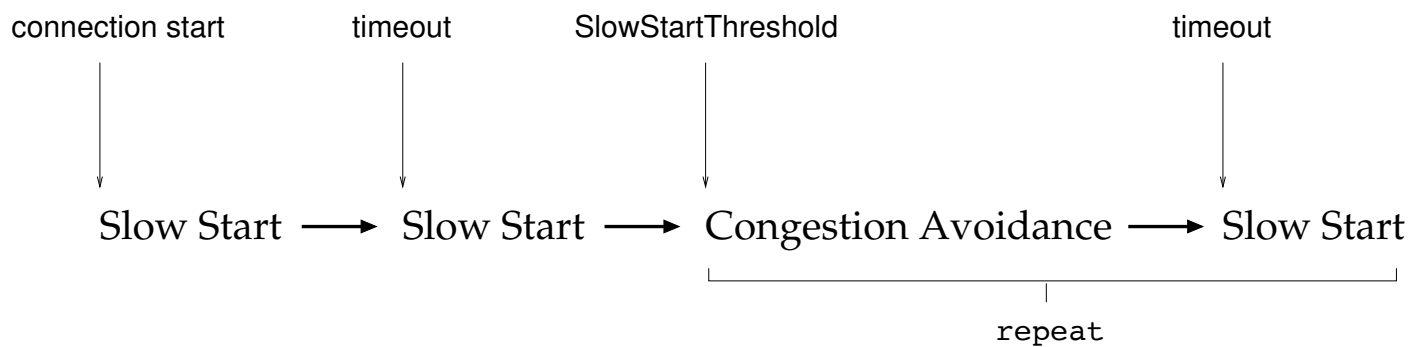Reset `CongestionWindow` to 1

Perform exponential increase

   `CongestionWindow` ← `CongestionWindow` + 1

• Until timeout at start of connection

   → rapidly probe for available bandwidth

• Until `CongestionWindow` hits `SlowStartThreshold`
  following Congestion Avoidance

   → rapidly climb to safe level


   ⟶  "slow" is a misnomer

   ⟶  exponential increase is super-fast

Basic dynamics:

$\longrightarrow$ after connection set-up

$\longrightarrow$ before connection tear-down

| connection start | timeout | SlowStartThreshold | | timeout |
|---|---|---|---|---|

Slow Start $\longrightarrow$ Slow Start $\longrightarrow$ Congestion Avoidance $\longrightarrow$ Slow Start

`repeat`

$\longrightarrow$ many TCP transfers are small

$\longrightarrow$ small TCP flows don't escape `Slow Start`

`CongestionWindow` evolution:

$\longrightarrow$ relevant for larger flows

CongestionWindow



**timeout**

**timeout**

**timeout**

ssthresh

ssthresh

ssthresh

Events (ACK or timeout)

(iii) Exponential timer backoff

$$\texttt{TimeOut} \leftarrow 2 \cdot \texttt{TimeOut} \qquad \text{if retransmit}$$

(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

  $\longrightarrow$ three duplicate ACKs: likely segment is lost

  $\longrightarrow$ react before timeout occurs

TCP Tahoe: features (i)-(iv)

(v) Fast Recovery

Upon Fast Retransmit:

- Skip Slow Start and commence Congestion Avoidance

  $\rightarrow$ dup ACKs: likely spurious loss

- Insert "inflationary" phase just before Congestion Avoidance

Additional changes and recent TCP variants.

Window scaling:

- 16-bit window size field limits receiver buffer size to 64 KB.

- Increase window size by scaling factor.

- During SYN handshake, exchange scaling factor using option field.

- If scaling factor is c, multiply window size by $2^{16+c}$

  $\rightarrow$ shift operation

  $\rightarrow$ $c$ limited to 14

BIC-TCP, TCP CUBIC: loss-based

- Instead of linear increase in Congestion Avoidance, use binary search (BIC)

  $\rightarrow$ concave shape: fast then slow when nearing window size of congestion event ($W_{\max}$)

  $\rightarrow$ convex shape: after $W_{\max}$ switch to probing mode

  $\rightarrow$ TCP CUBIC uses cubic function directly

  $\rightarrow$ Linux

TCP Vegas, Compound TCP: delay-based, hybrid

- Estimate queueing delay from RTT

  $\rightarrow$ use minimum as reference point

- If RTT increases assume queueing at bottleneck link(s)

  $\rightarrow$ slow down linearly

  $\rightarrow$ closer to method D

  $\rightarrow$ susceptible to congestion collapse
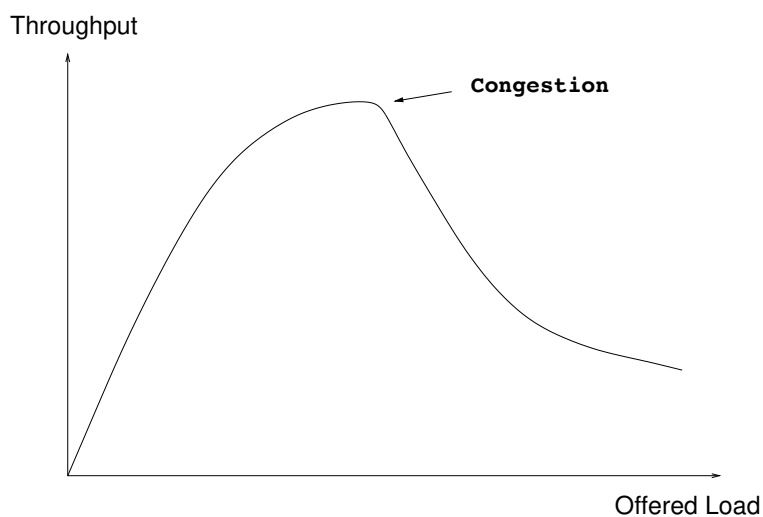
Case for exponential backoff

- For multimedia streaming (e.g., pseudo real-time) with limited prefetch, AIMD (Method B) not suited

  $\rightarrow$ can use Method D, variants

  $\rightarrow$ under long prefetch, can use reliable transport (e.g., TCP)

- For unimodal case—throughput decreases when system load is excessive—instability concern

  $\rightarrow$ asymmetry in control law to curb instability

  $\rightarrow$ worst-case: congestion collapse

Congestion control and selfishness:

$\longrightarrow$  to be or not to be selfish . . .

$\longrightarrow$  John von Neumann, John Nash, . . .

Ex.: "tragedy of commons," Garrett Hardin, '68



- if everyone acts selfishly, no one wins

  $\rightarrow$ in fact, everyone loses

- can this be prevented?

Ex.: Prisoner's Dilemma game

   $\longrightarrow$   formalized by Tucker in 1950

   $\longrightarrow$   "cold war"

 • both cooperate (i.e., stay mum): 1 year each

 • both selfish (i.e., rat on the other): 5 years each

 • one cooperative/one selfish: 9 vs. 0 years

<div align="center">

*Bob*

|       | C    | N    |
|-------|------|------|
| **C** | 1, 1 | 9, 0 |
| **N** | 0, 9 | 5, 5 |

*Alice*

</div>

   $\longrightarrow$   payoff matrix

   $\longrightarrow$   what would "rational" prisoners do?

When cast as congestion control game:

$$Bob$$

|  | C | N |
|---|---|---|
| C | $5, 5$ | $0, 9$ |
| N | $9, 0$ | $1, 1$ |

Alice and Bob share network bandwidth

$\rightarrow (a, b)$: throughput (Mbps) achieved by Alice/Bob

$\rightarrow$ large is desirable

Upon congestion: back off or escalate?

$\rightarrow$ equivalent to Prisoner's dilemma

Rational: in the sense of seeking selfish gain

$\rightarrow$ both choose strategy "N"

$\rightarrow$ called Nash equilibrium

$\rightarrow$ steady-state or stable fixed-point

Reason:

$\rightarrow$ whatever choice the other player makes, "N" yields better payoff over "C"
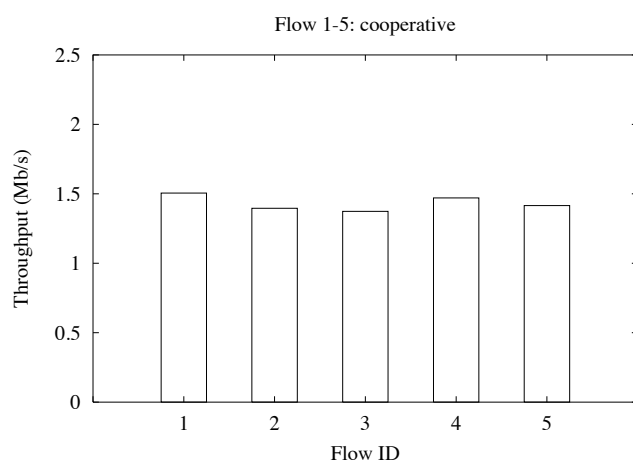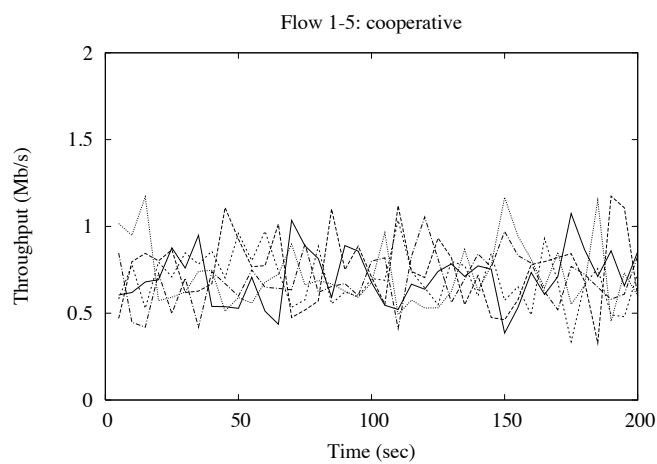
$\rightarrow$ i.e., strategy "N" dominates strategy "C"

In some systems, selfish behavior results in system optimal outcome

$\rightarrow$ theoretical foundation of Adam Smith's "invisible hand"

$\rightarrow$ in general, not the case

$\rightarrow$ cooperation is better but can it be enforced?

Impact in networks:

$\rightarrow$ 5 regular (cooperative) TCP flows

$\rightarrow$ share 11 Mbps WLAN bottleneck link



Flow 1-5: cooperative



Flow 1-5: cooperative

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

$\rightarrow$ starts behaving selfishly at time 100s

Flow 1-4: cooperative     Flow 5: noncooperative



Flow 1-4: cooperative     Flow 5: noncooperative

Potential danger for:

→ unfairness

→ overall system performance

Is it being exploited in today's Internet?

→ no one knows

→ technical implementation issues

→ e.g., interoperability with legacy protocols

→ e.g., shooting oneself in the foot