# Advanced Computer Networks

## TCP and Fair Queuing

# Introduction to TCP

- Communication abstraction:
  - Reliable
  - Ordered
  - Point-to-point
  - Byte-stream
  - Full duplex
  - Flow and congestion controlled

- Protocol implemented entirely at the ends
  - Fate sharing

- Sliding window with cumulative acks
  - Ack field contains last in-order packet received
  - Duplicate acks sent when out-of-order packet received

# Key Things You Should Know Already

- Port numbers
- TCP/UDP checksum
- Sliding window flow control
  - Sequence numbers
- TCP connection setup
- TCP reliability
  - Timeout
  - Data-driven

# TCP Congestion Control

- Changes to TCP motivated by ARPANET congestion collapse

- Basic principles
  - Additive increase/multiplicative decrease (AIMD)
  - Packet conservation
  - Reaching steady state quickly
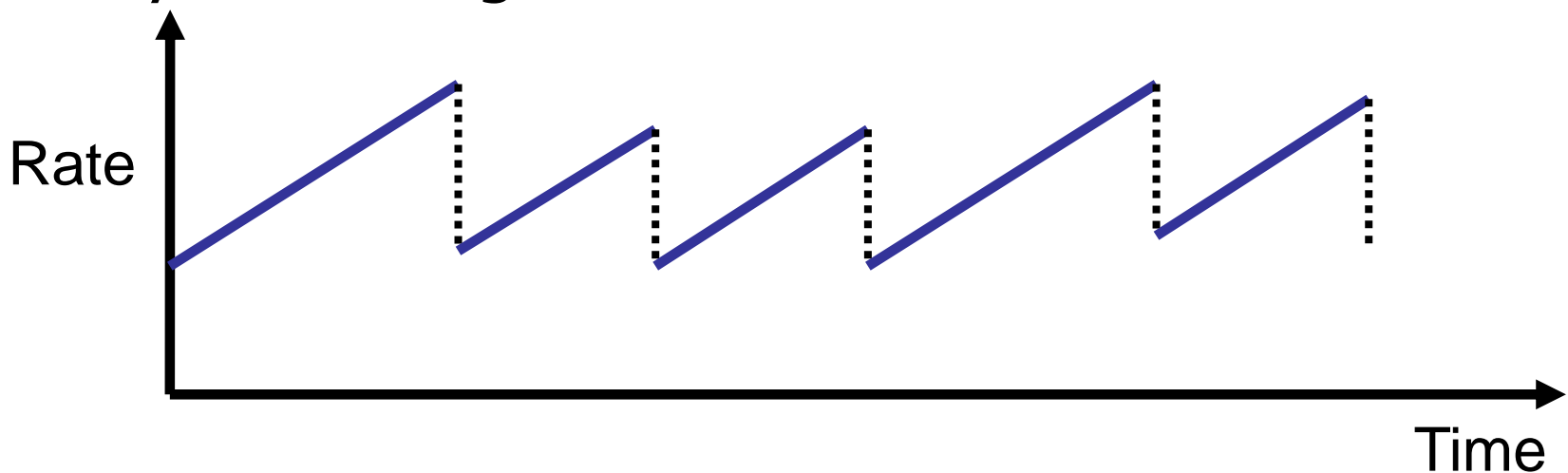  - ACK clocking

# TCP Congestion Control - Solutions

- Reaching equilibrium
  - Slow start
- Eliminates spurious retransmissions
  - Accurate RTO estimation
  - Fast retransmit
- Adapting to resource availability
  - Congestion avoidance

# AIMD

- Distributed, fair and efficient
- Packet loss is seen as sign of congestion and results in a multiplicative rate decrease
  - Factor of 2
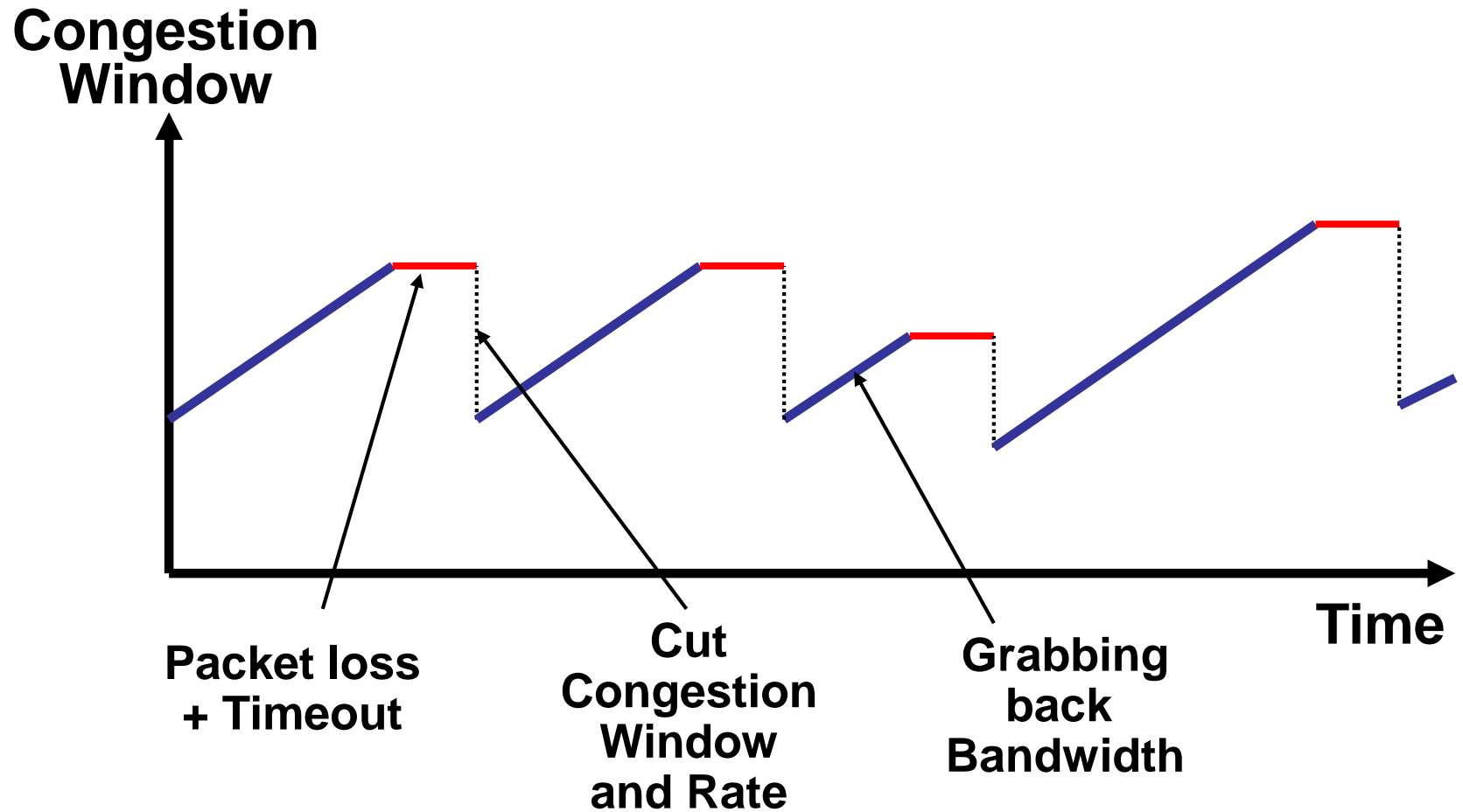- TCP periodically probes for available bandwidth by increasing its rate



Rate

Time

# Congestion Avoidance

- **If loss occurs when cwnd = W**
  - Network can handle 0.5W ~ W segments
  - Set cwnd to 0.5W (multiplicative decrease)
- **Upon receiving ACK**
  - Increase  cwnd by (1 packet)/cwnd
    - What is 1 packet? $\rightarrow$ 1 MSS worth of bytes
    - After cwnd packets have passed by $\rightarrow$ approximately increase of 1 MSS
- **Implements AIMD**
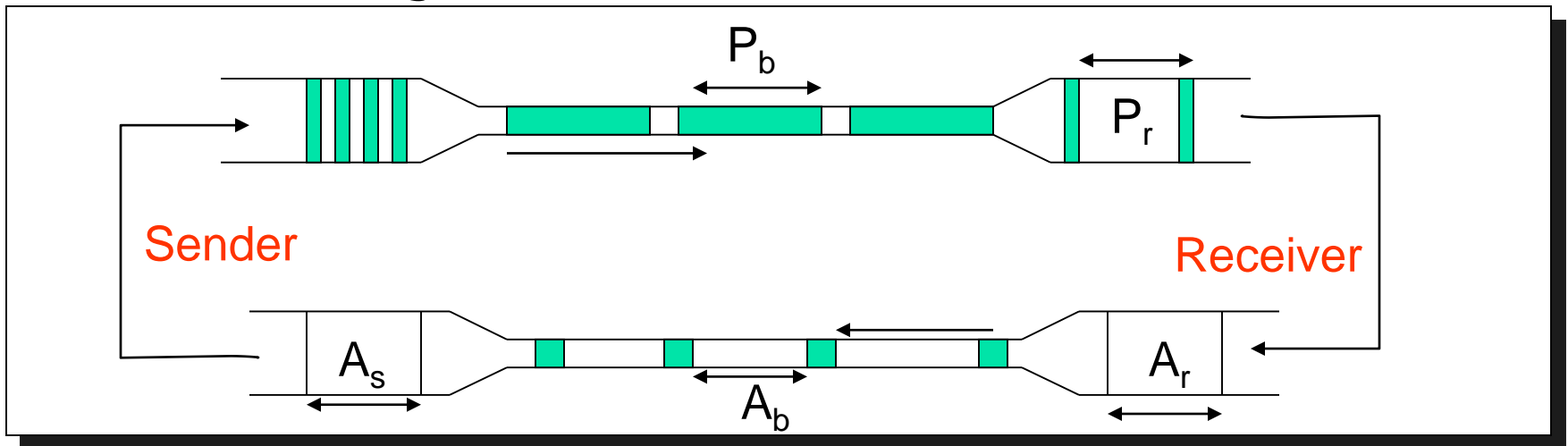
# Congestion Avoidance Behavior

# Packet Conservation

- At equilibrium, inject packet into network only when one is removed
  - Sliding window and not rate controlled
  - But still need to avoid sending burst of packets $\rightarrow$ would overflow links
    - Need to carefully pace out packets
    - Helps provide stability
- Need to eliminate spurious retransmissions
  - Accurate RTO estimation
  - Better loss recovery techniques (e.g. fast retransmit)

# TCP Packet Pacing

- Congestion window helps to "pace" the transmission of data packets

- In steady state, a packet is sent when an ack is received

  - Data transmission remains smooth, once it is smooth
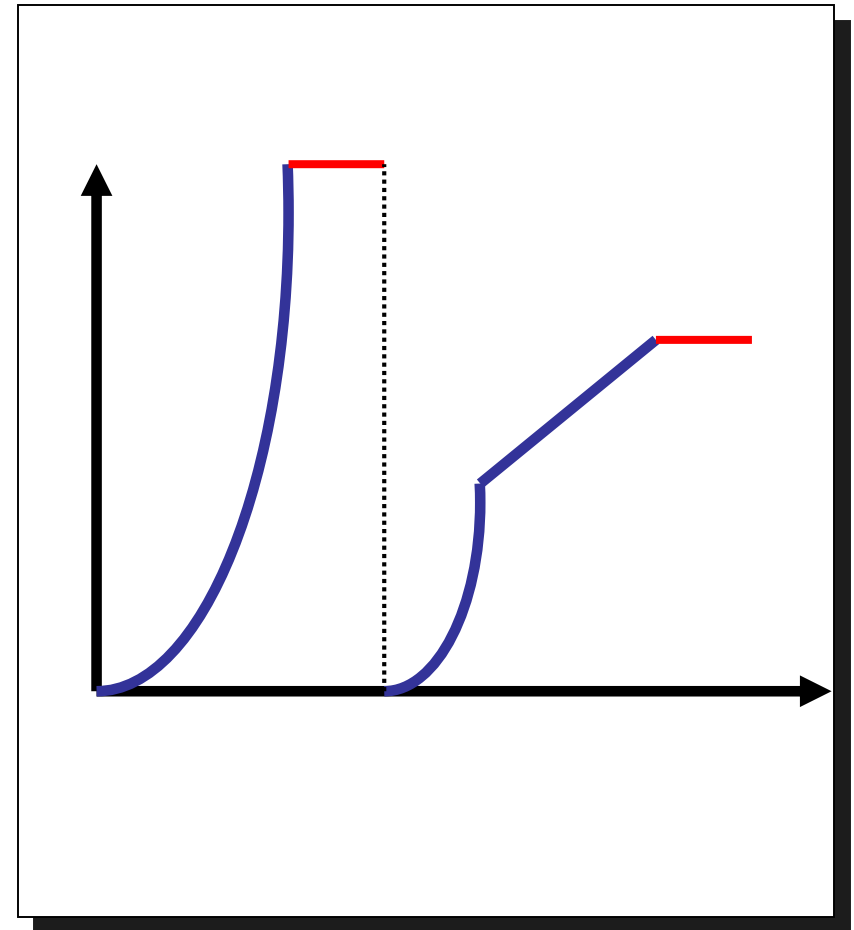  - Self-clocking behavior
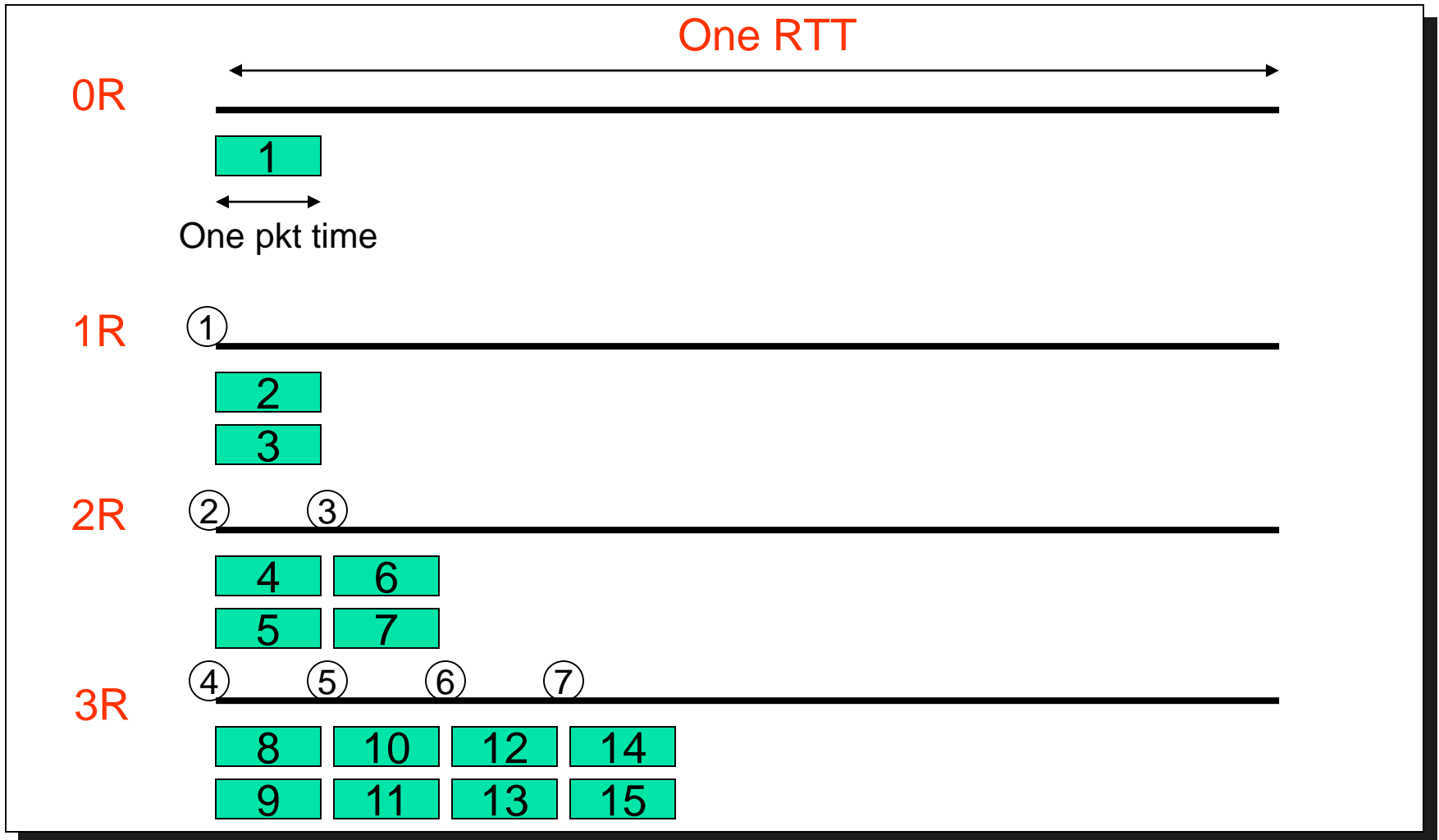
# Reaching Steady State

- Doing AIMD is fine in steady state but slow...

- How does TCP know what is a good initial rate to start with?

  - Should work both for a CDPD (10s of Kbps or less) and for supercomputer links (10 Gbps and growing)

- Quick initial phase to help get up to speed (slow start)

# Slow Start Packet Pacing

- How do we get this clocking behavior to start?
  - Initialize cwnd = 1
  - Upon receipt of every ack, cwnd = cwnd + 1
- Implications
  - Window actually increases to W in RTT * $\log_2(W)$
  - Can overshoot window and cause packet loss

# Slow Start Example

One RTT

**0R**

| 1 |

One pkt time

**1R** ①

| 2 |
| 3 |

**2R** ② ③

| 4 | 6 |
| 5 | 7 |

**3R** ④ ⑤ ⑥ ⑦

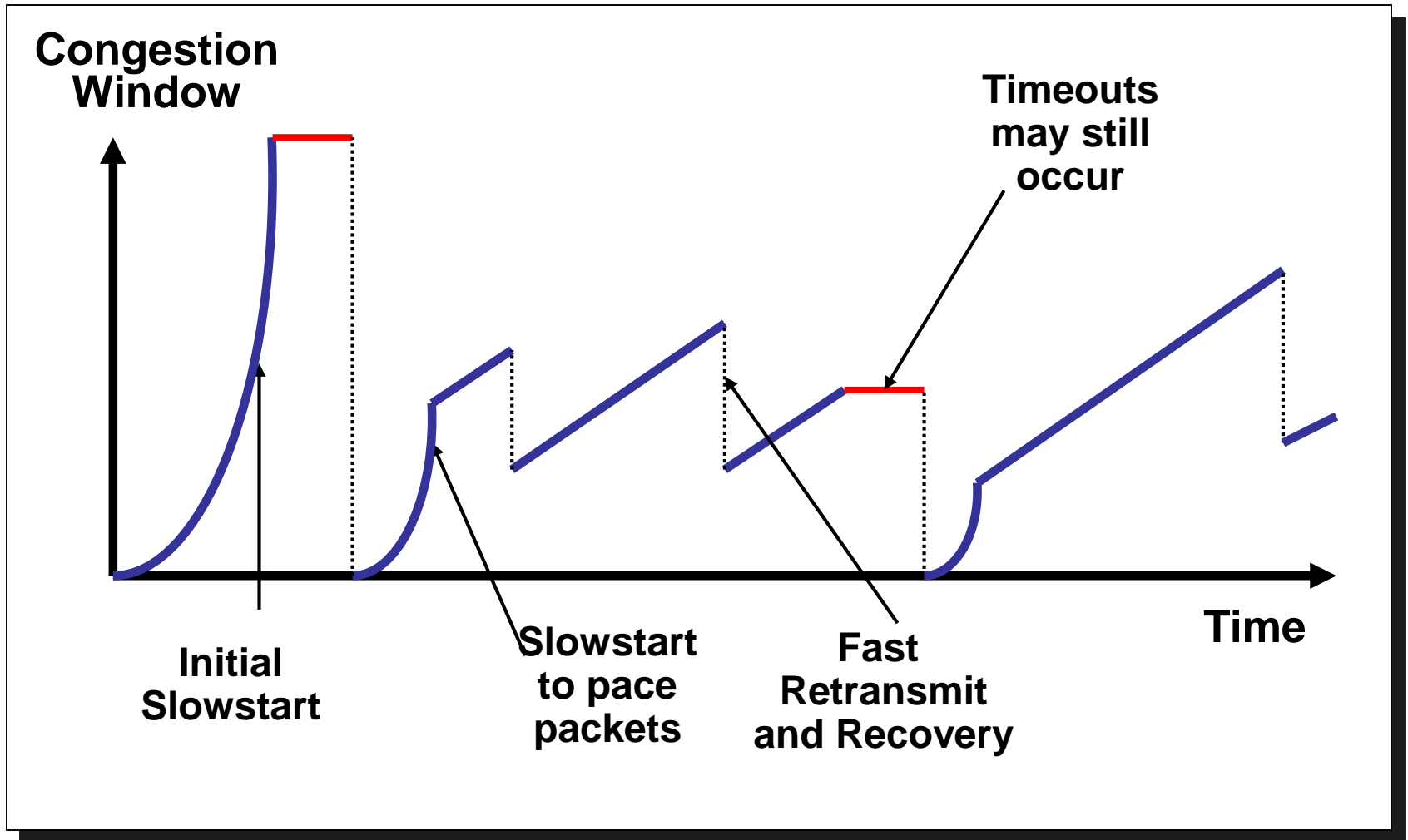| 8 | 10 | 12 | 14 |
| 9 | 11 | 13 | 15 |

# Return to Slow Start

- If packet is lost we lose our self clocking as well
  - Need to implement slow-start and congestion avoidance together
- When timeout occurs set ssthresh to 0.5w
  - If cwnd < ssthresh, use slow start
  - Else use congestion avoidance

# TCP Saw Tooth Behavior

**Congestion Window**

**Timeouts may still occur**

**Time**

**Initial Slowstart**

**Slowstart to pace packets**

**Fast Retransmit and Recovery**

# Changing Workloads

- New applications are changing the way TCP is used
- 1980's Internet
    - Telnet & FTP → long lived flows
    - Well behaved end hosts
    - Homogenous end host capabilities
    - Simple symmetric routing
- 2000's Internet
    - Web & more Web → large number of short xfers
    - Wild west – everyone is playing games to get bandwidth
    - Cell phones and toasters on the Internet
    - Policy routing
- How to accommodate new applications?

# TCP Friendliness

- **What does it mean to be TCP friendly?**
  - TCP is not going away
  - Any new congestion control must compete with TCP flows
    - Should not clobber TCP flows and grab bulk of link
    - Should also be able to hold its own, i.e. grab its fair share, or it will never become popular
- **How is this quantified/shown?**
  - Has evolved into evaluating loss/throughput behavior
  - If it shows 1/sqrt(p) behavior it is ok

# TCP Friendly Rate Control (TFRC)

- **Equation 1 – real TCP response**

$$T = \frac{s}{R\sqrt{\frac{2p}{3}} + t_{RTO}(3\sqrt{\frac{3p}{8}})p(1 + 32p^2)}$$

  - 1st term corresponds to simple derivation
  - 2nd term corresponds to more complicated timeout behavior
    - Is critical in situations with > 5% loss rates → where timeouts occur frequently

- **Key parameters**
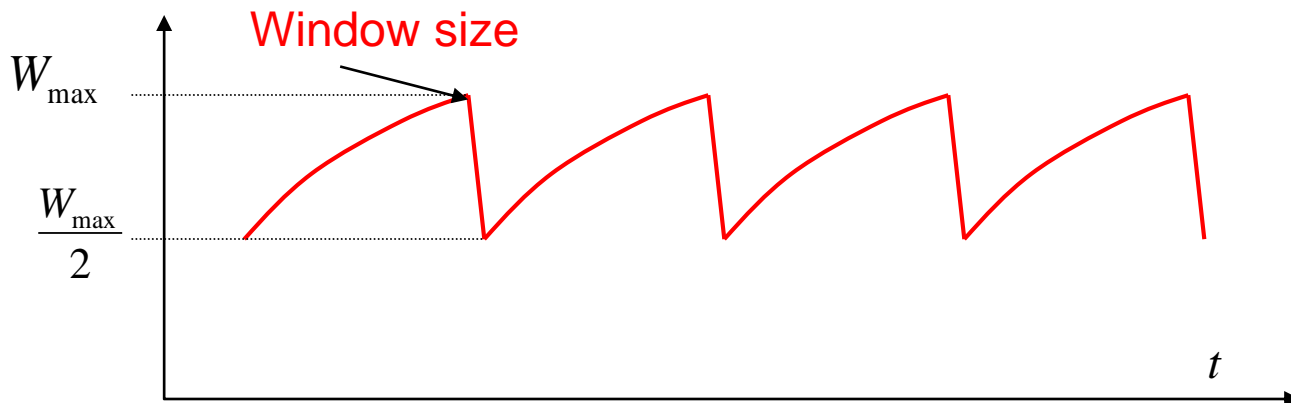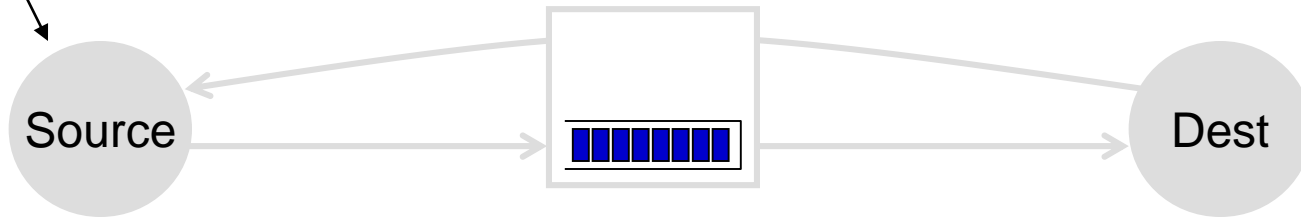  - RTO
  - RTT
  - Loss rate

# TCP Performance

- Can TCP saturate a link?
- Congestion control
  - Increase utilization until… link becomes congested
  - React by decreasing window by 50%
  - Window is proportional to rate * RTT
- Doesn't this mean that the network oscillates between 50 and 100% utilization?
  - Average utilization = 75%??
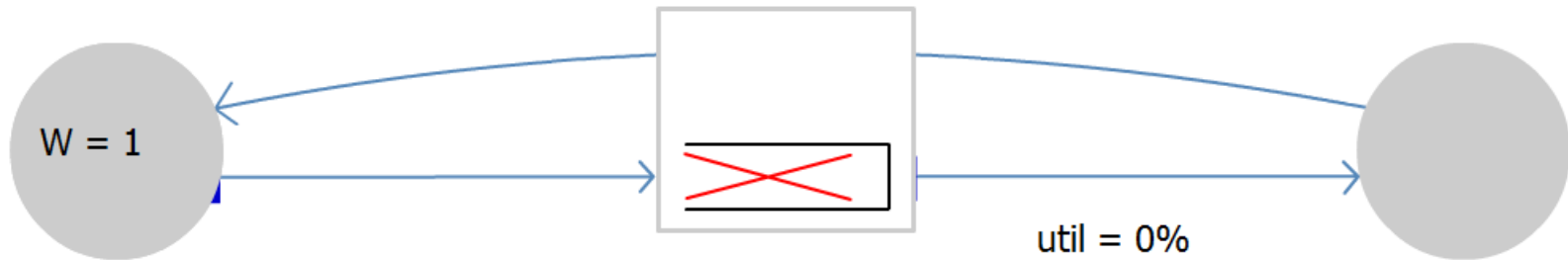  - No…this is *not* right!

# TCP Congestion Control

Only W packets may be outstanding

## Rule for adjusting $W$

- If an ACK is received: $W \leftarrow W + 1/W$
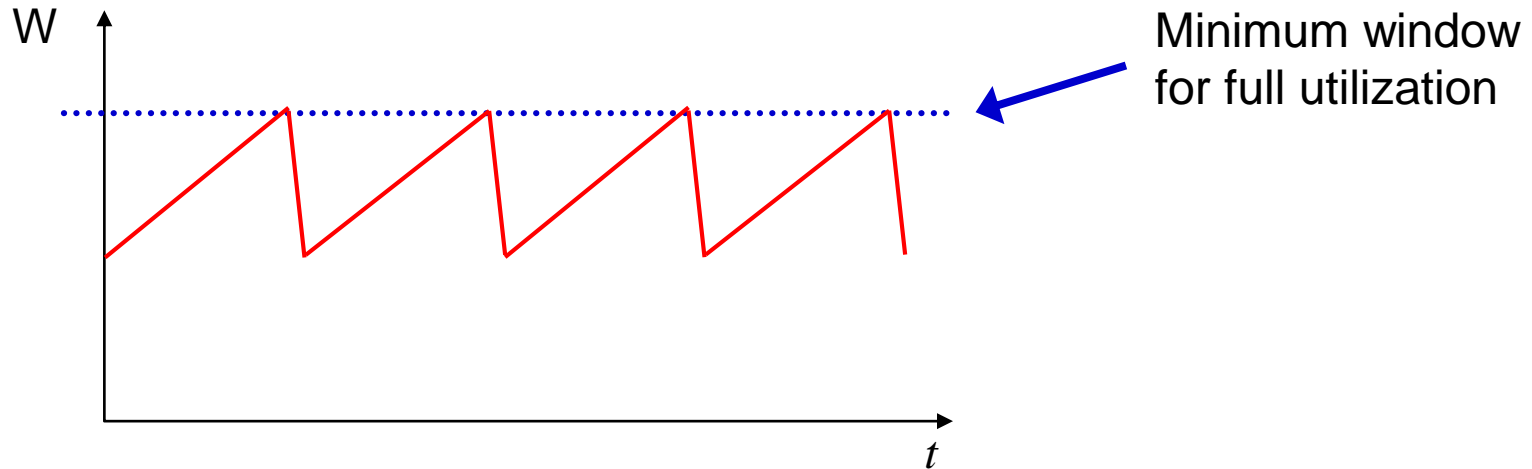- If a packet is lost: $W \leftarrow W/2$

Source

Dest

Window size

$W_{max}$

$\dfrac{W_{max}}{2}$

$t$

# Single TCP Flow
## Router without buffers



W = 1

util = 0%

W

time

# Summary Unbuffered Link

W

Minimum window
for full utilization

*t*

- The router can't fully utilize the link
  - If the window is too small, link is not full
  - If the link is full, next window increase causes drop
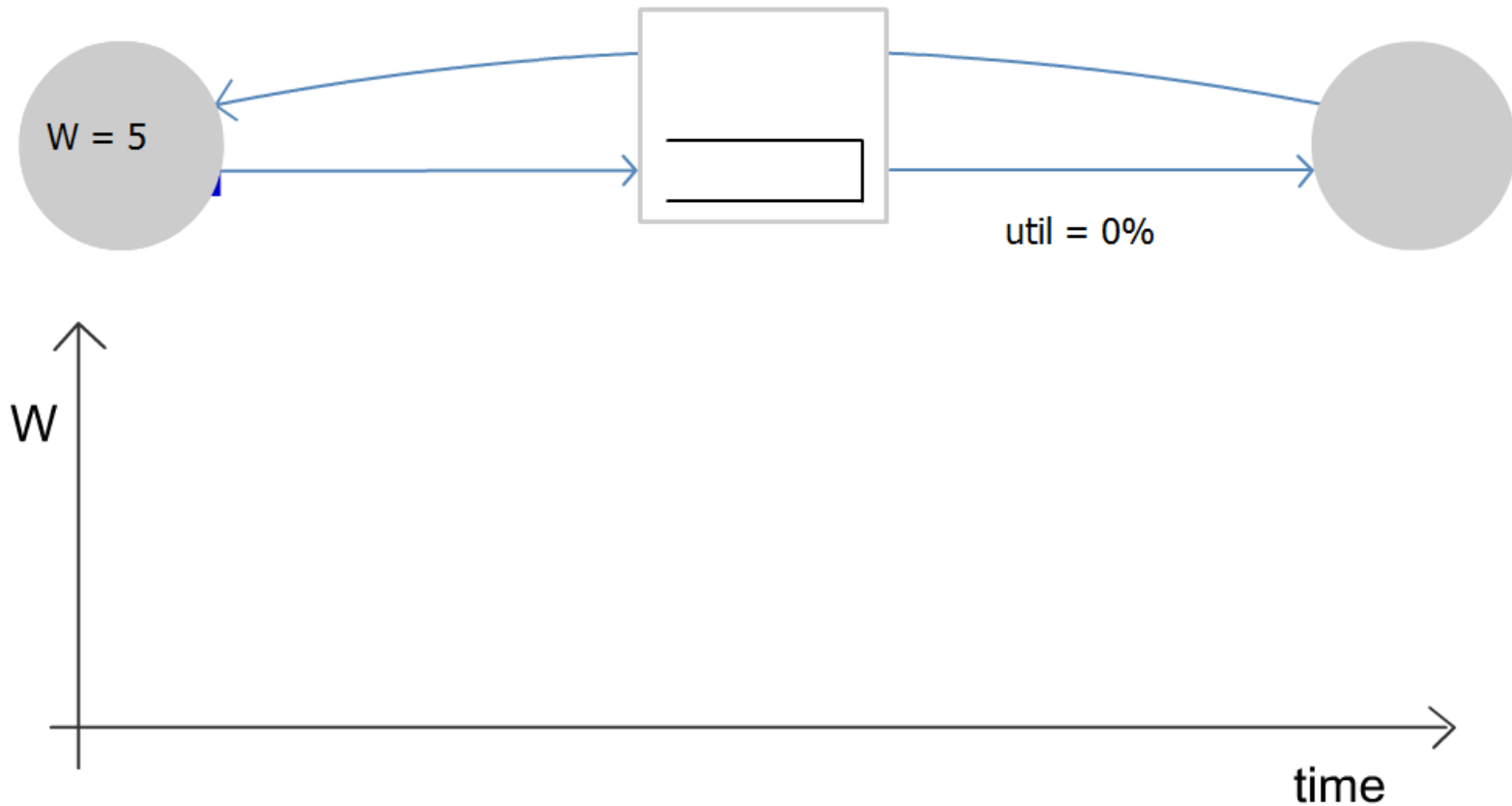  - With no buffer it still achieves 75% utilization

# TCP Performance

- In the real world, router queues play important role
  - Window is proportional to rate * RTT
    - But, RTT changes as well the window
  - Window to fill links = propagation RTT * bottleneck bandwidth
    - If window is larger, packets sit in queue on bottleneck link

# Single TCP Flow
Router with large enough buffers for full link utilization
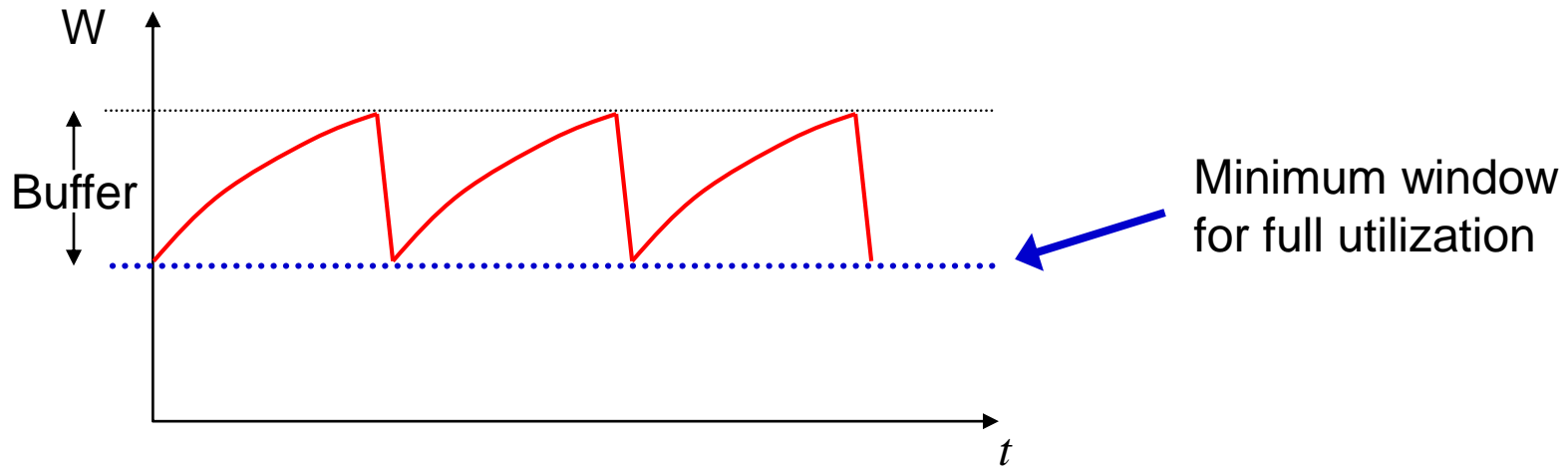


W = 5

util = 0%

W

time

# TCP Performance with Queue

- If we have a large router queue → can get <u>100% utilization</u>

  - But, router queues can cause large delays

- How big does the queue need to be?

  - Windows vary from W → W/2

    - Must make sure that link is always full

    - W/2 > RTT * BW

    - W = RTT * BW + Qsize

    - Therefore, Qsize > RTT * BW

  - <u>Ensures 100% utilization</u>

  - Delay?

    - Varies between RTT and 2 * RTT

# Summary Buffered Link



W

Buffer

Minimum window
for full utilization

t

- With sufficient buffering we achieve full link utilization
  - The window is always above the critical threshold
  - Buffer absorbs changes in window size
    - Buffer Size = Height of TCP Sawtooth
    - Minimum buffer size needed is 2T*C, where T=RTT, C=BW
  - This is the origin of the rule-of-thumb

# Queuing Disciplines

- Each router **must** implement some queuing discipline

- Queuing allocates both bandwidth and buffer space:
  - Bandwidth: which packet to serve (transmit) next
  - Buffer space: which packet to drop next (when required)
- Queuing also affects latency

# Typical Internet Queuing

- FIFO + drop-tail
  - Simplest choice
  - Used widely in the Internet
- FIFO (first-in-first-out)
  - Implies single class of traffic
- Drop-tail
  - Arriving packets get dropped when queue is full regardless of flow or importance
- Important distinction:
  - FIFO: scheduling discipline
  - Drop-tail: drop policy