# Congestion Control

Lecture 15
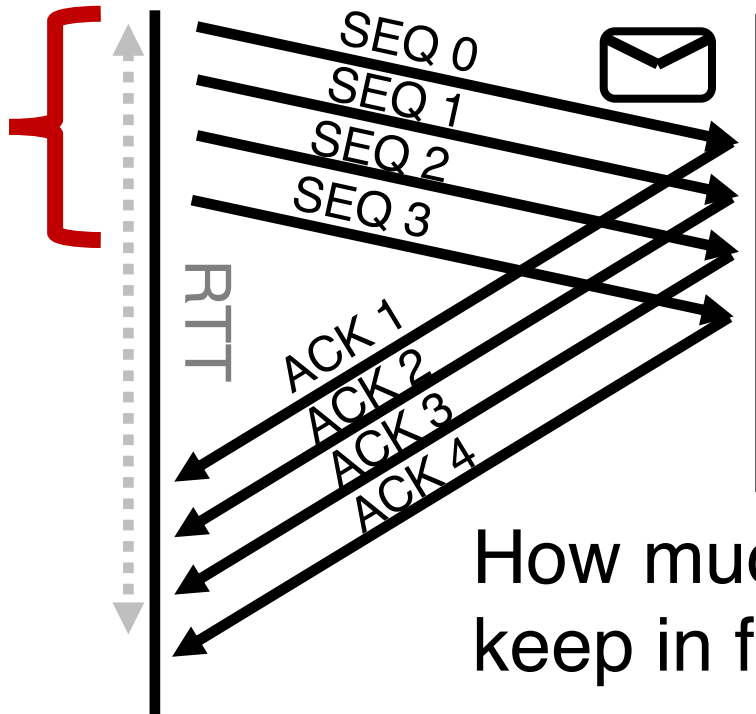Srinivas Narayana

RUTGERS
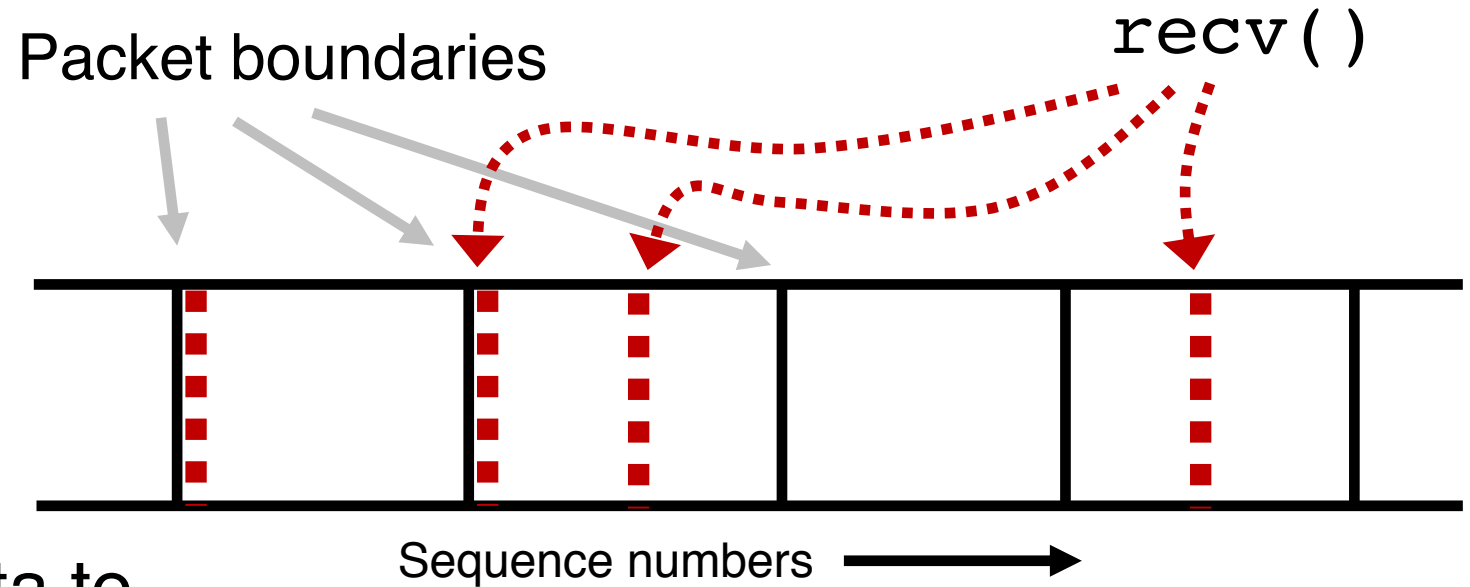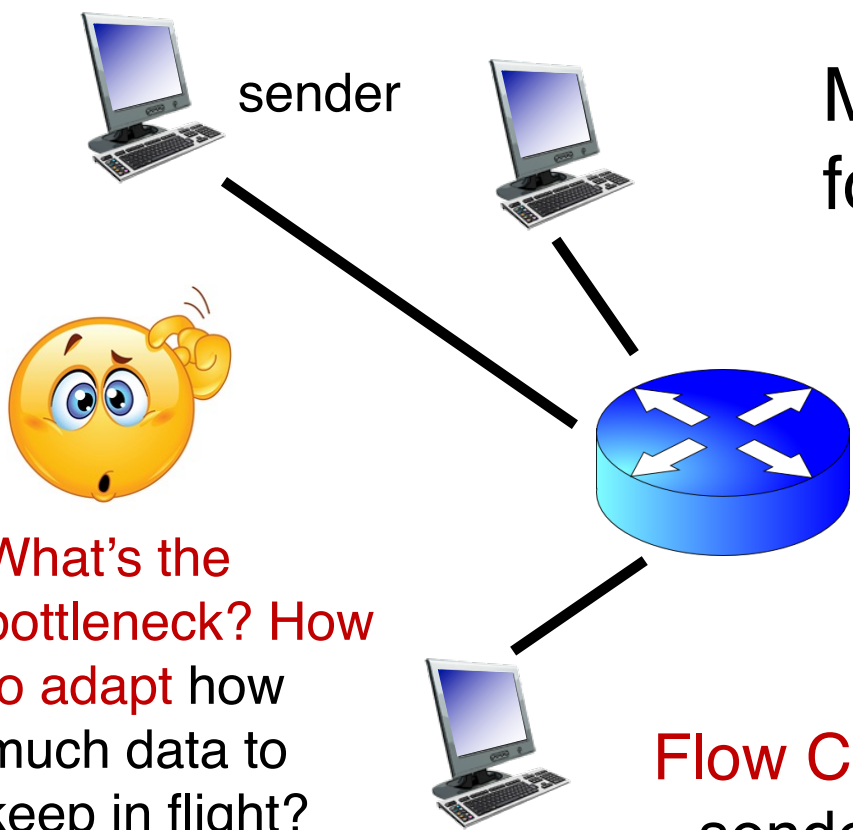UNIVERSITY | NEW BRUNSWICK

# Quick recap of concepts


Tp layer

TCP:
Reliability
Ordering

Reordering degrades connection throughput. Apps can't recv(). Packets may even be dropped due to insufficient buffering.
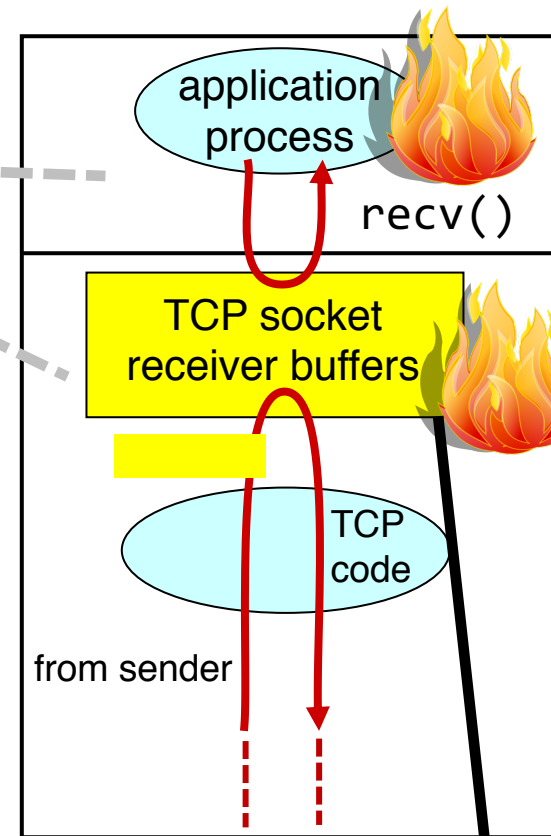

SEQ 0
SEQ 1
SEQ 2
SEQ 3
RTT
ACK 1
ACK 2
ACK 3
ACK 4

How much data to keep in flight?

Packet boundaries

recv()

Sequence numbers

Stream-Oriented Transport

sender

Multiple locations
for bottlenecks

Congestion Control

application
process

recv()

TCP socket
receiver buffers

TCP
code

from sender

Flow Control

What's the
bottleneck? How
to adapt how
much data to
keep in flight?

Flow Control: Receiver informs
sender free buffer over time

receiver

Buffer >= desired W

```
 0                   1                   2                   3
 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|          Source Port          |       Destination Port        |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                        Sequence Number                        |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                     Acknowledgment Number                     |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
| Data  |           |U|A|P|R|S|F|                               |
| Offset| Reserved  |R|C|S|S|Y|I|            Window             |
|       |           |G|K|H|T|N|N|                               |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|           Checksum            |         Urgent Pointer        |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                    Options                    |    Padding     |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                             data                              |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+

                       TCP Header Format

  Note that one tick mark represents one bit position.
```

Last cumulative
ACK'ed seq #

Last transmitted
seq #

Sender's
view:

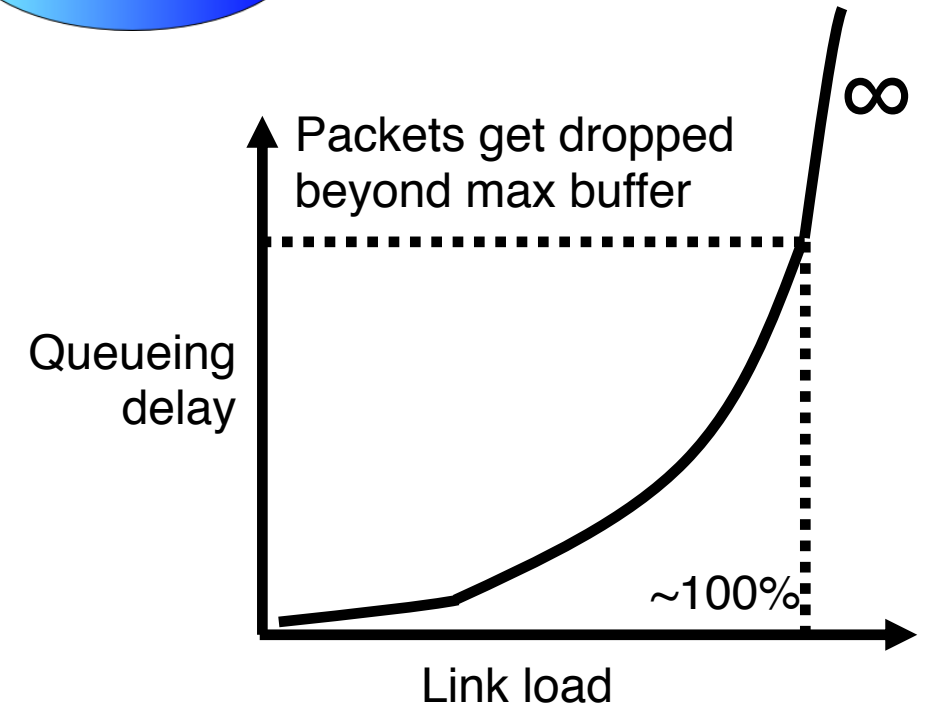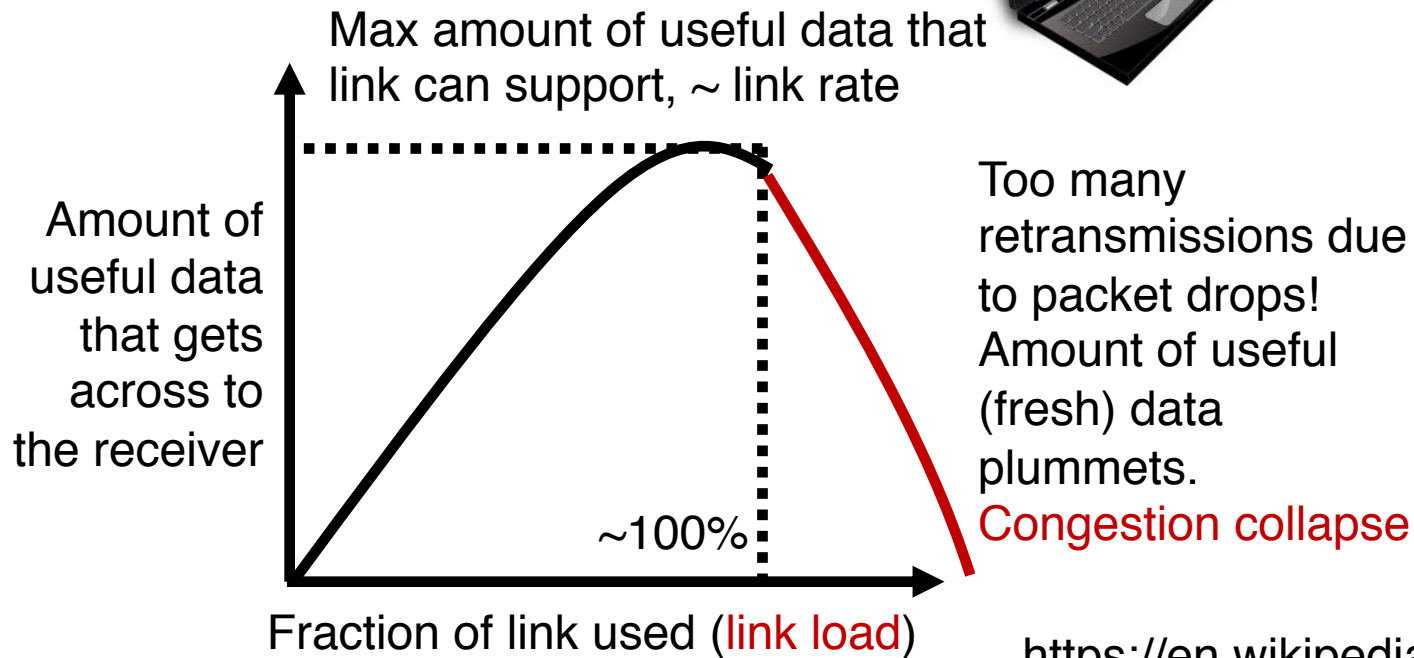| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 0 | 1 |

Window <=   Advertised window

Low socket buffering
==
Poor TCP throughput

# Congestion

Routers have buffers which accommodate queued packets.

Max amount of useful data that link can support, ~ link rate

Amount of useful data that gets across to the receiver

~100%

Fraction of link used (link load)

Too many retransmissions due to packet drops! Amount of useful (fresh) data plummets.
Congestion collapse

Packets get dropped beyond max buffer

∞

Queueing delay

~100%

Link load

https://en.wikipedia.org/wiki/Network_congestion#Congestive_collapse

# How should multiple endpoints share net?



- It is difficult to know where the bottleneck link is
- It is difficult to know how many other endpoints are using that link
- Endpoints may join and leave at any time
- Network paths may change over time, leading to different bottleneck links (with different link rates) over time

The approach that the Internet takes is to use a distributed algorithm to converge to an efficient and fair outcome.

The approach that the Internet takes is to use a distributed algorithm to converge to an efficient and fair outcome.

No one can centrally view or control all the endpoints and bottlenecks in the Internet.

Every endpoint must try to reach a globally good outcome by itself: i.e., in a distributed fashion.

This also puts a lot of trust in endpoints.

The approach that the Internet takes is to use a distributed algorithm to converge to an <span style="color:red">efficient</span> and fair outcome.

If there is spare capacity in the bottleneck link, the endpoints should use it.

The approach that the Internet takes is to use a distributed algorithm to converge to an efficient and fair outcome.

If there are N endpoints sharing a bottleneck link, they should be able to get equitable shares of the link's capacity.

For example: 1/N'th of the link capacity.

# Flow Control    vs.    Congestion Control

- Avoid overwhelming the <span style="color:red">receiving application</span>

- Sender is managing the <span style="color:red">receiver's socket buffer</span>

- Avoid overwhelming the <span style="color:red">bottleneck network link</span>

- Sender is managing the <span style="color:red">bottleneck link capacity</span> and <span style="color:red">bottleneck router buffers</span>

The approach that the Internet takes is to use a distributed algorithm to converge to an efficient and fair outcome.

How to achieve this?

Approach: sense and react
Example: taking a shower
Use a feedback loop with signals and knobs

H  1  2  C

# Signals and Knobs in Congestion Control

- ## Signals
  - Packets being ACK'ed
  - Packets being dropped (e.g. RTO fires)
  - Packets being delayed (RTT)
  - Rate of incoming ACKs

  Implicit feedback signals measured directly at sender. (There are also explicit signals that the network might provide.)

- ## Knobs
  - What can you change to "probe" the available bottleneck capacity?
  - Suppose receiver buffer is unbounded:
  - Increase window/sending rate: e.g., add x or multiply by a factor of x
  - Decrease window/sending rate: e.g., subtract x or reduce by a factor of x

# Sense and react, sure…but how?

- Where do you want to be?
  - The steady state

- How do you get there?
  - Congestion control algorithms
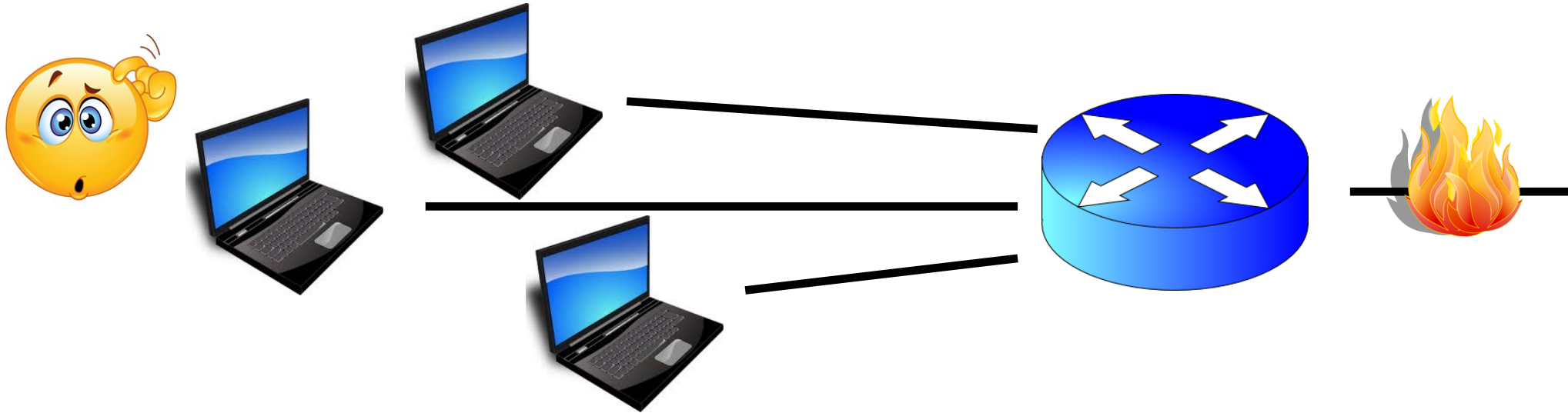
- Sense accurately

- React proportionately

# The Steady State

Efficiency of a single TCP conversation

# What does efficiency look like?

- Suppose we want to achieve an efficient outcome for one TCP conversation by observing network signals from the endpoint
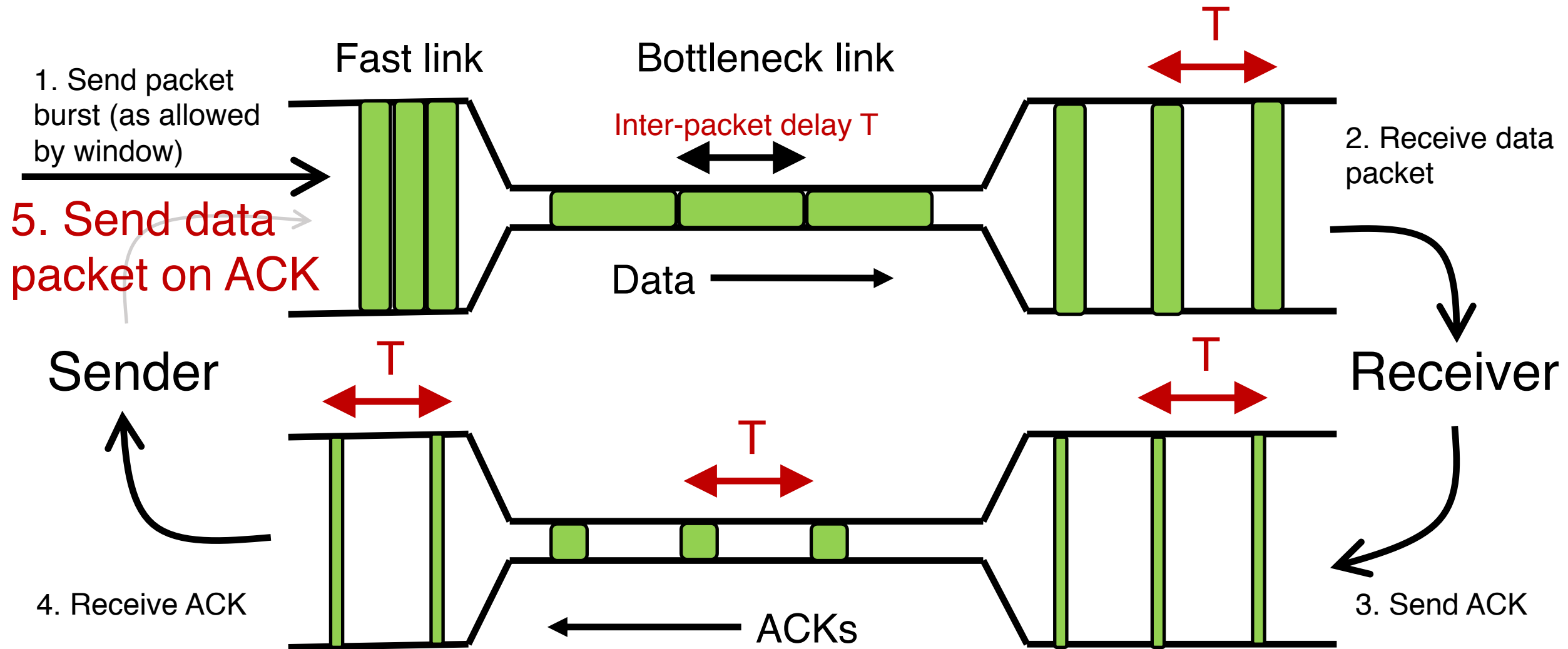


- Q: How should the endpoint behave at steady state?
- Challenge: bottleneck link is remotely located

# Steady state: Ideal goal

- High sending rate: Use the full capacity of the bottleneck link
- Low delay: Minimize the overall delay of packets to get to the receiver
  - Overall delay = propagation + queueing + transmission
  - Assume propagation and transmission components fixed
- "Low delay" reduces to low queueing delay
- i.e., don't push so much data into the network that packets have to wait in queues

- Key question: When to send the next packet?

# When to send the next packet?

# Rationale
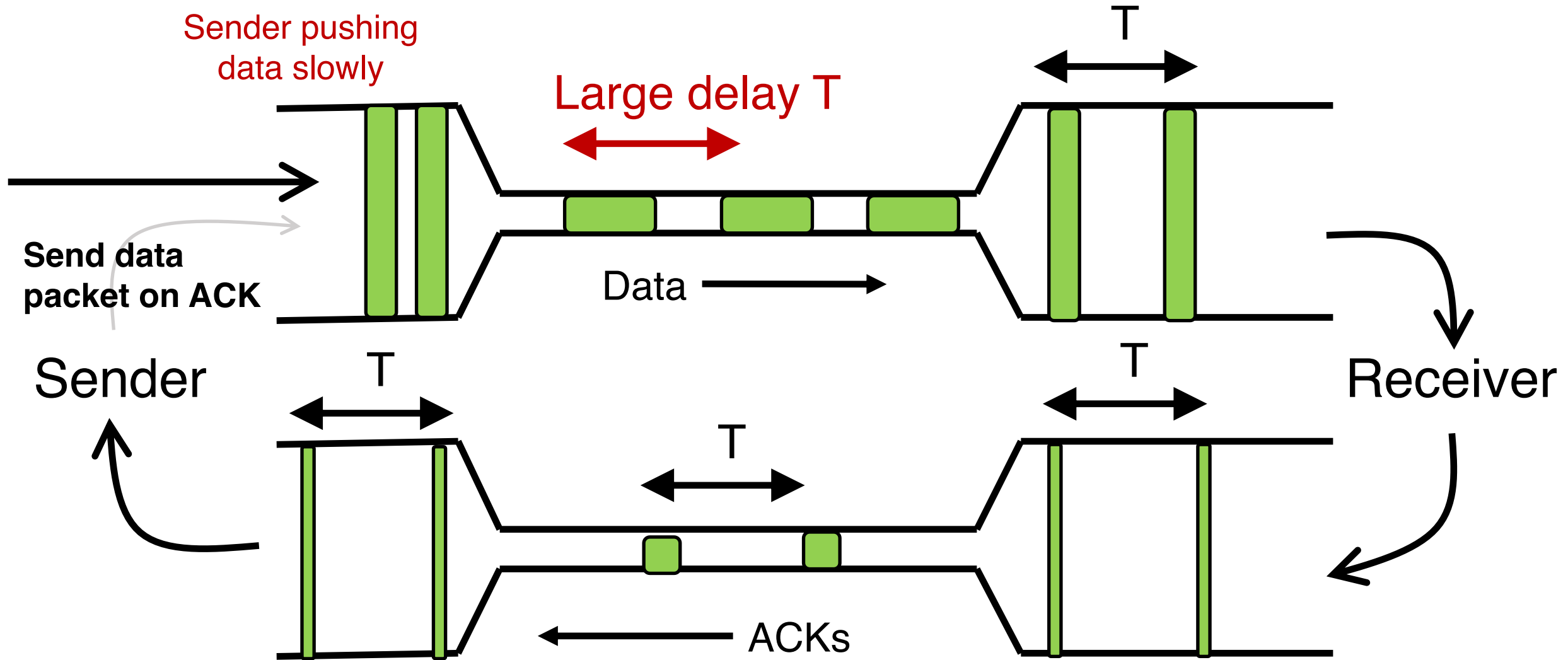
- When the sender receives an ACK, that's a signal that the previous packet has left the bottleneck link (and the rest of the network)

- Hence, <span style="color:red">it must be safe to send another packet without congesting the bottleneck link</span>

- Such transmissions are said to follow <span style="color:red">packet conservation</span>

- <span style="color:red">ACK clocking:</span> "Clock" of ACKs governs packet transmissions
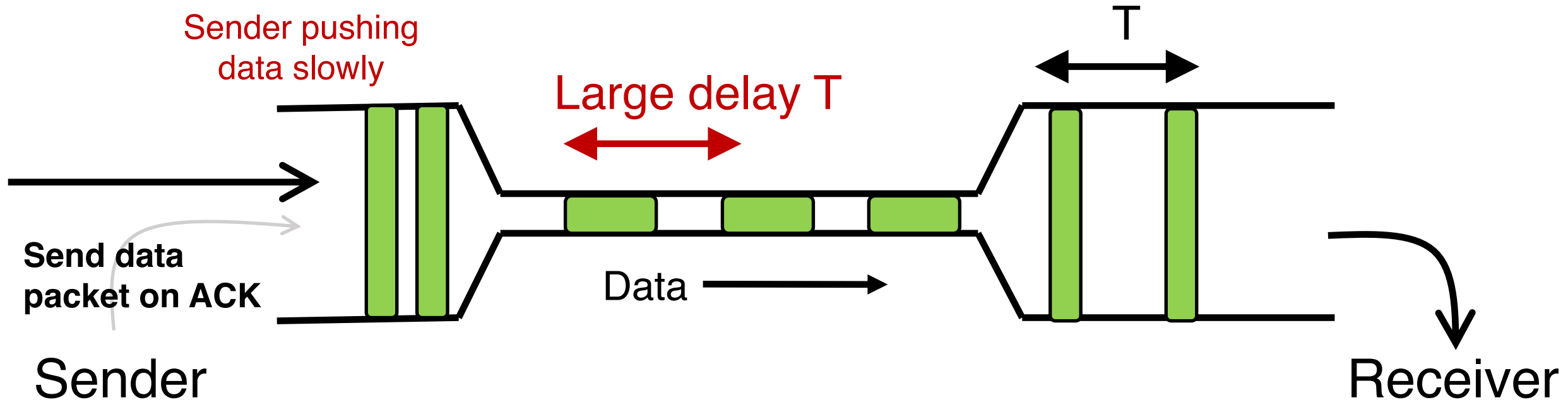
# ACK clocking: analogy

- How to avoid crowding a grocery store?



- Strategy: Send the next waiting customer exactly when a customer exits the store

- However, this strategy alone can lead to inefficient use of resources…

# ACK clocking alone can be inefficient

# ACK clocking alone can be inefficient



The sending rate should be high enough to keep the "pipe" full
Analogy: a grocery store with only 1 customer in entire store
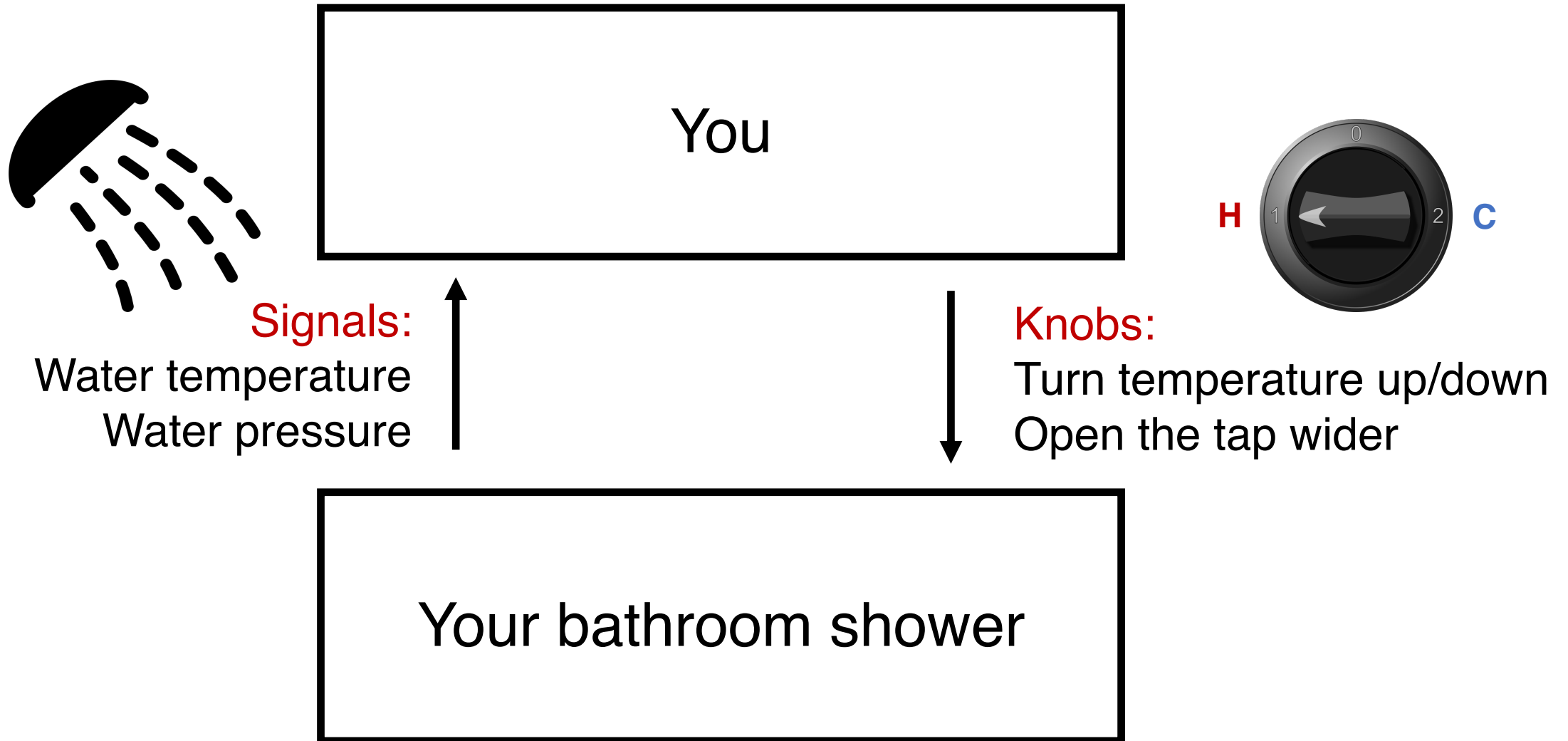If the store isn't "full", you're using store space inefficiently

# Steady State of Congestion Control

- Send at the highest rate possible (to keep the pipe full)
- while being ACK-clocked (to avoid congesting the pipe)
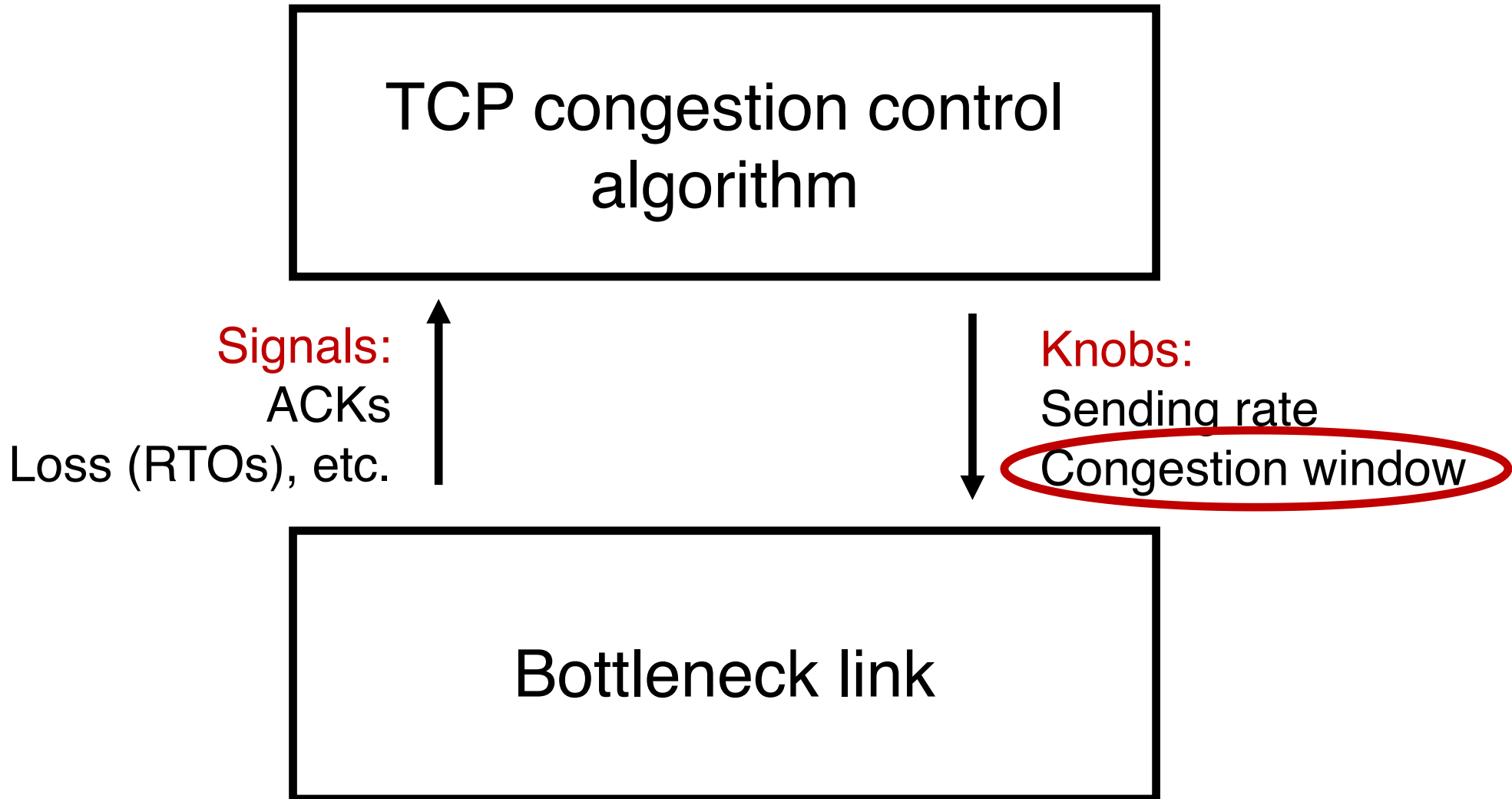
- So, how to get to steady state?

# Getting to Steady State

The Feedback Loop

# An example of a feedback loop



**You**

**Your bathroom shower**

Signals:
Water temperature
Water pressure

Knobs:
Turn temperature up/down
Open the tap wider

# The congestion control feedback loop

TCP congestion control algorithm

Signals:
ACKs
Loss (RTOs), etc.

Knobs:
Sending rate
Congestion window

Bottleneck link

# Congestion window

- The sender maintains an estimate of the amount of in-flight data needed to keep the pipe full without congesting it.

- This estimate is called the congestion window (cwnd)

- Recall: There is a relationship between the sending rate (throughput) and the sender's window:  sender transmits a window's worth of data over an RTT duration
    - Rate = window / RTT

# Interaction b/w flow & congestion control

- Use window = min(congestion window, receiver advertised window)

- Overwhelm neither the receiver nor network links & routers

Window <= Congestion window (congestion control)
Window <= Advertised window (flow control)

Sender's view:

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 0 | 1 |

Last cumulative ACK'ed seq #

Last transmitted seq #