

Congestion

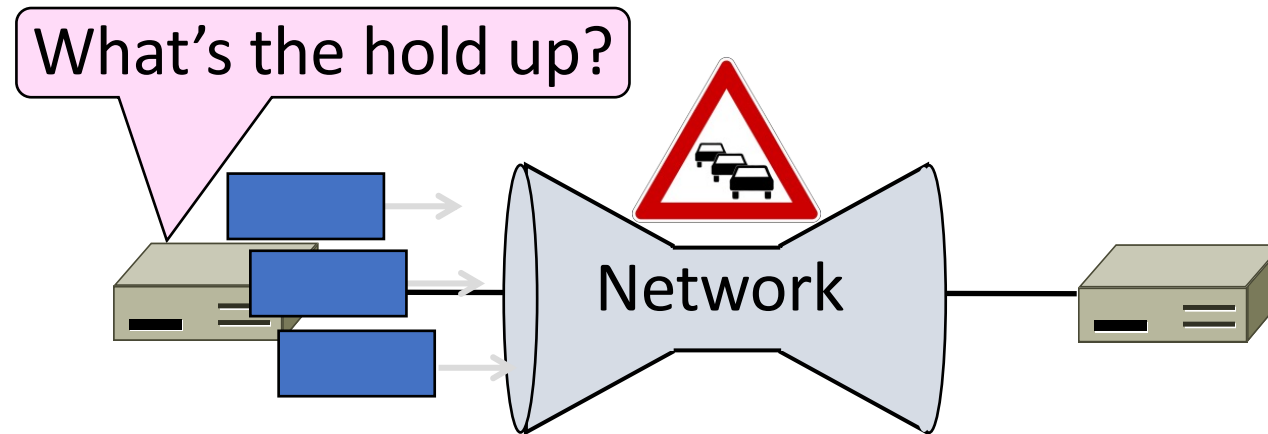
TCP to date:

- We can set up and tear connections
 - Connection establishment and release handshakes
- Keep the sending and receiving buffers from overflowing (flow control)

What's missing?

Network Congestion

- A “traffic jam” in the network
 - Later we will learn how to control it

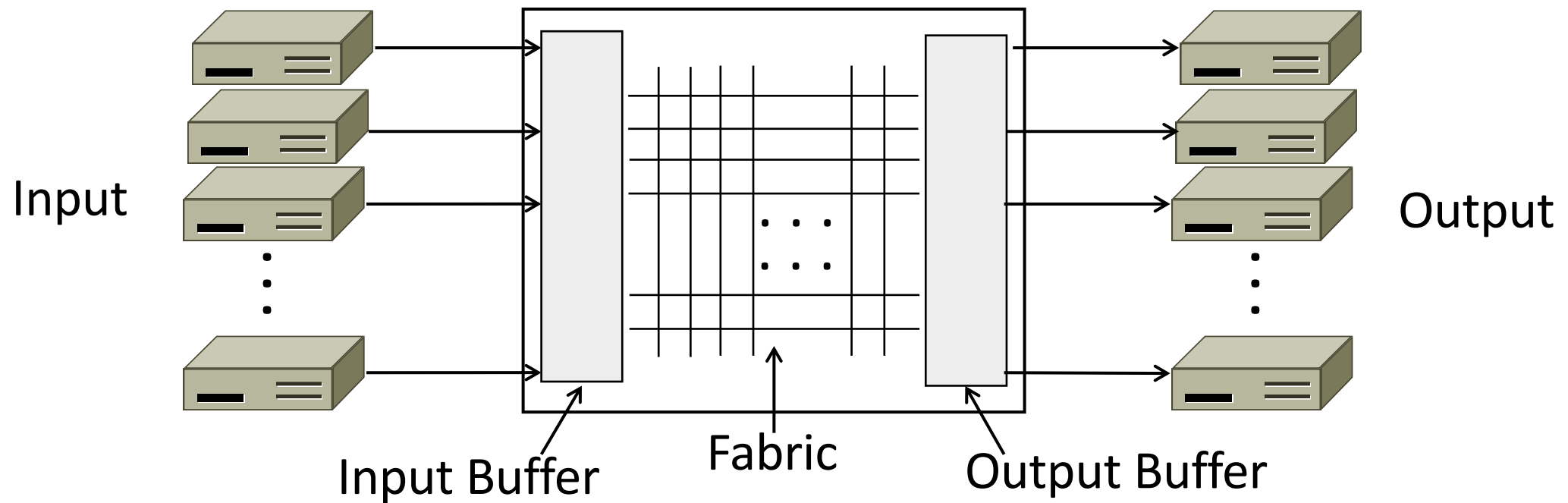


Congestion Collapse in the 1980s

- Early TCP used fixed size window (e.g., 8 packets)
 - Initially fine for reliability
- But something happened as the network grew
 - Links stayed busy but transfer rates fell by orders of magnitude!

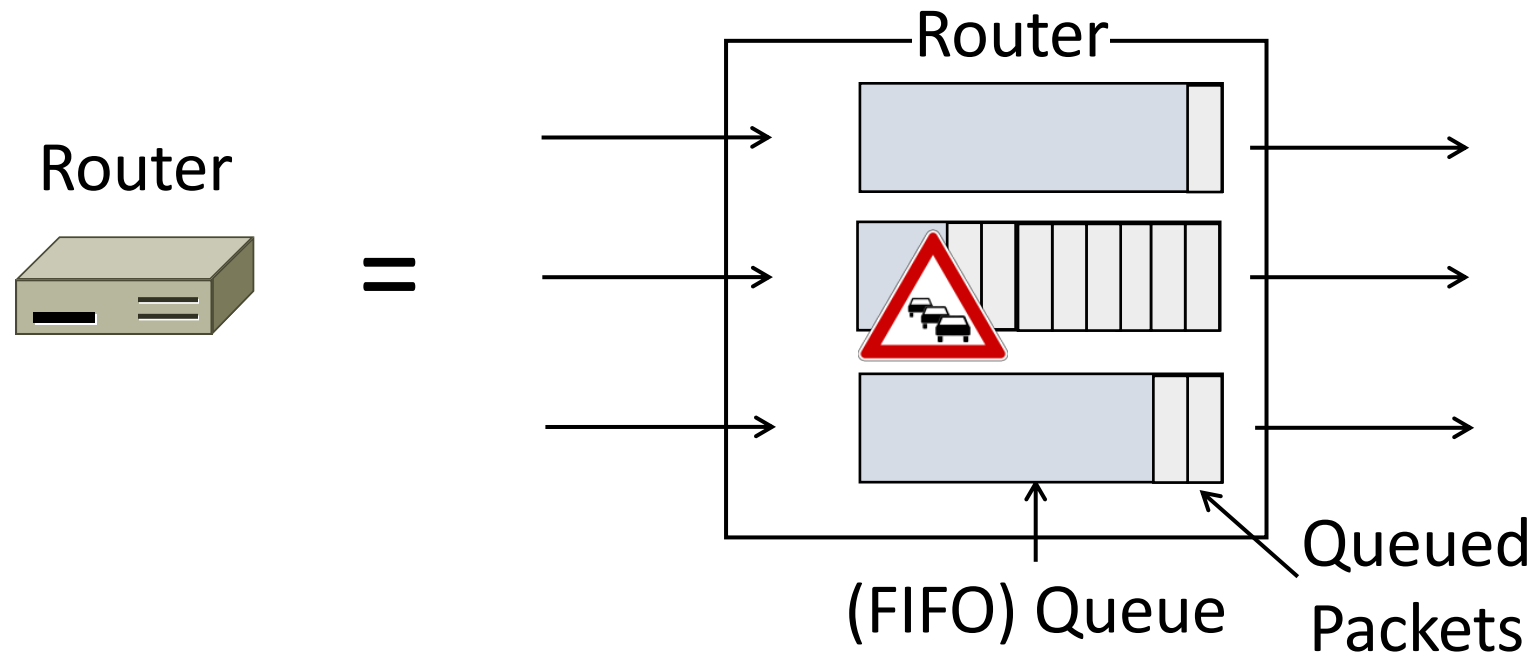
Nature of Congestion

- Routers/switches have internal buffering



Nature of Congestion (2)

- Simplified view of per port output queues
 - Typically FIFO (First In First Out), discard when full



Nature of Congestion (3)

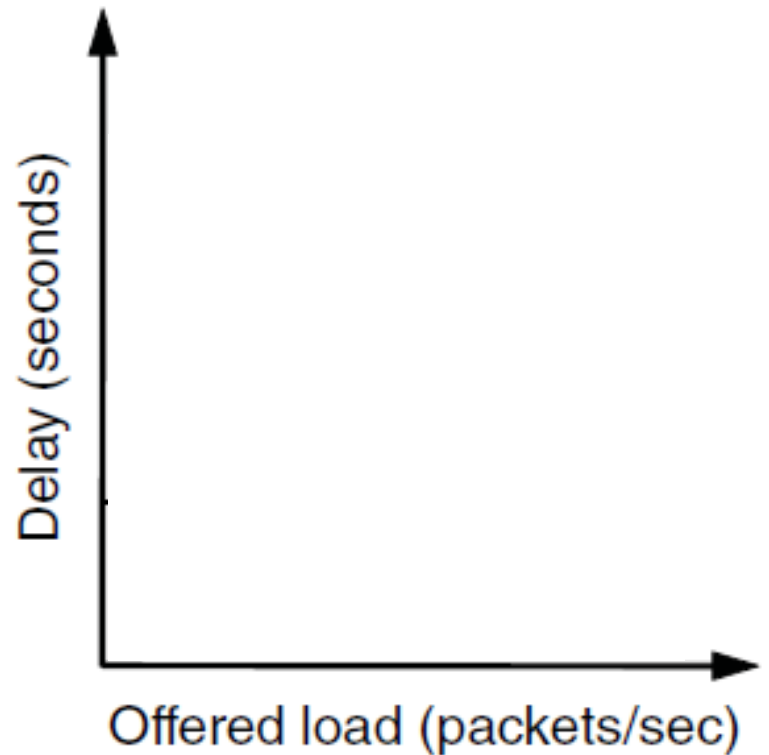
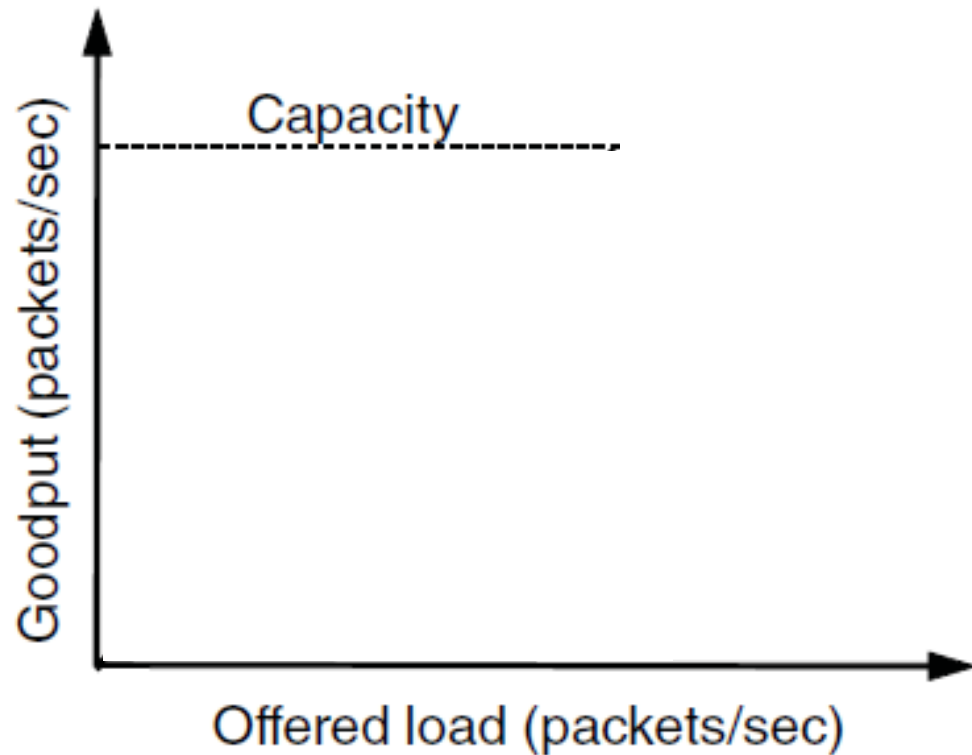
Queues absorb bursts when $\text{input} > \text{output rate}$

But if $\text{input} > \text{output rate}$ persistently, queue will overflow \rightarrow congestion

Congestion is a function of the traffic patterns – can occur even if every link has the same capacity

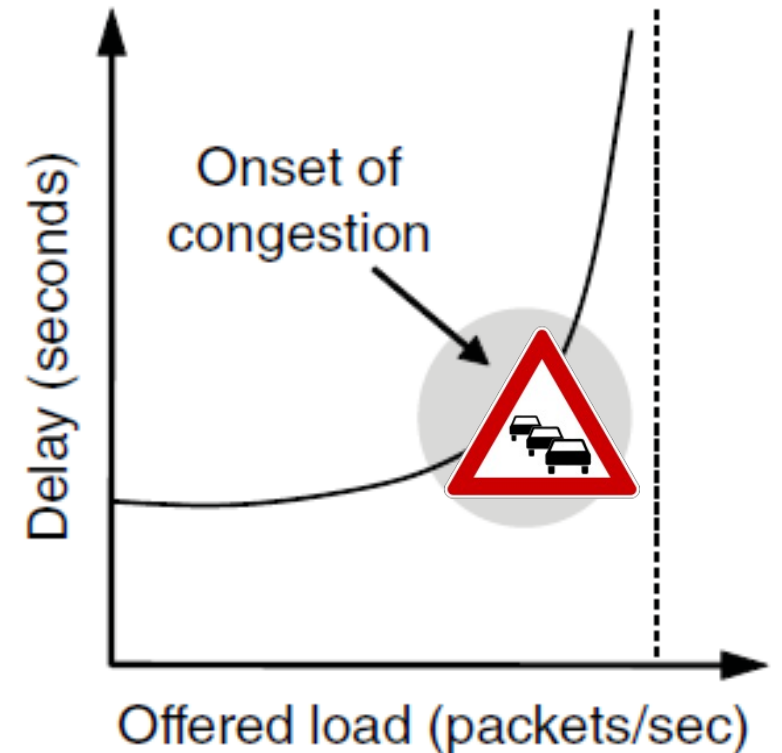
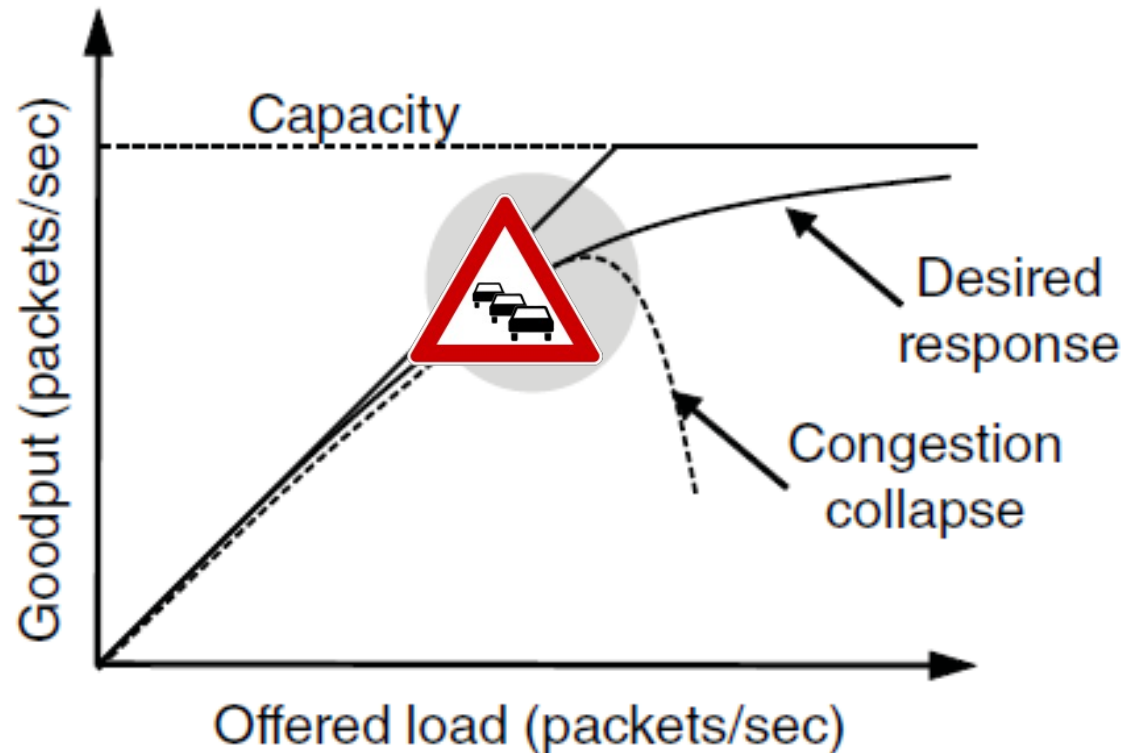
Effects of Congestion

- What happens to performance as we increase load?



Effects of Congestion (2)

- What happens to performance as we increase load?



Effects of Congestion (3)

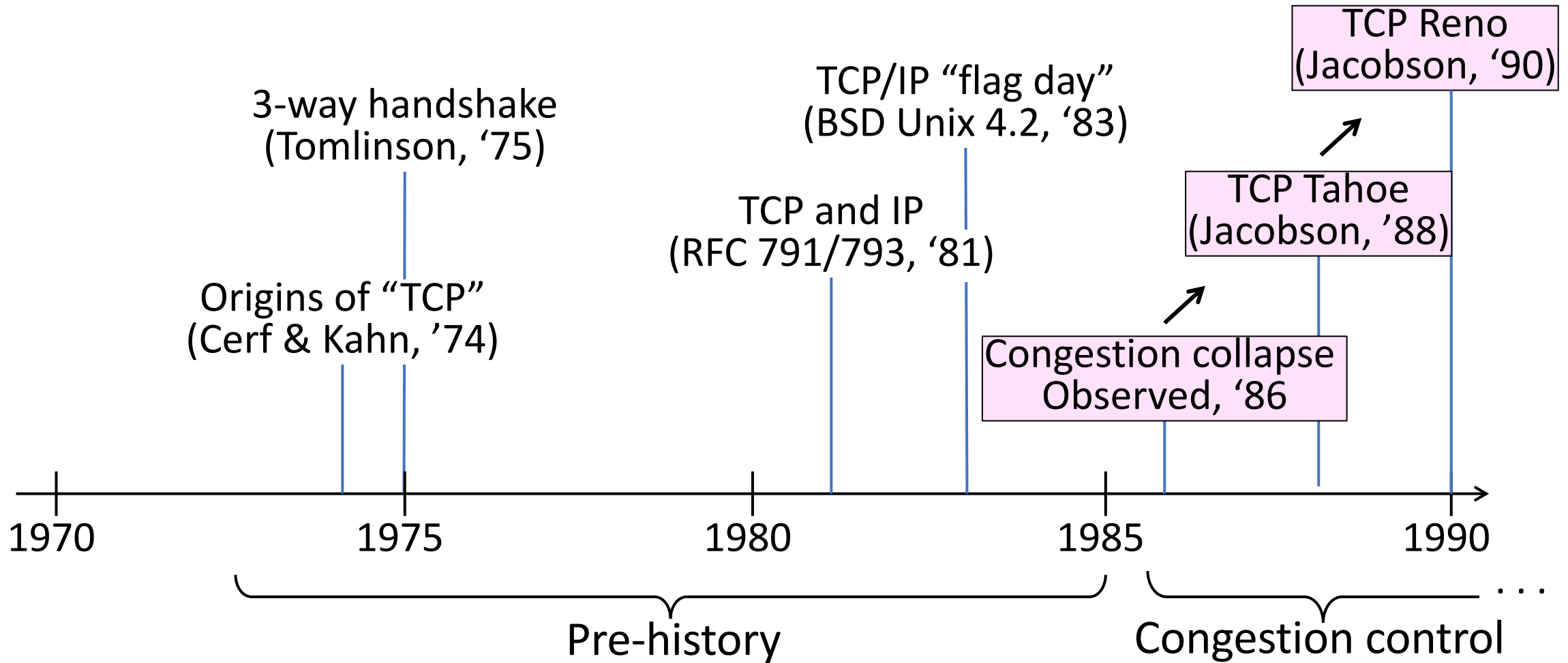
- As offered load rises, congestion occurs as queues begin to fill:
 - Delay and loss rise sharply with load
 - Throughput $<$ load (due to loss)
 - Goodput \ll throughput (due to spurious retransmissions)
- None of the above is good!
 - Want network performance just before congestion



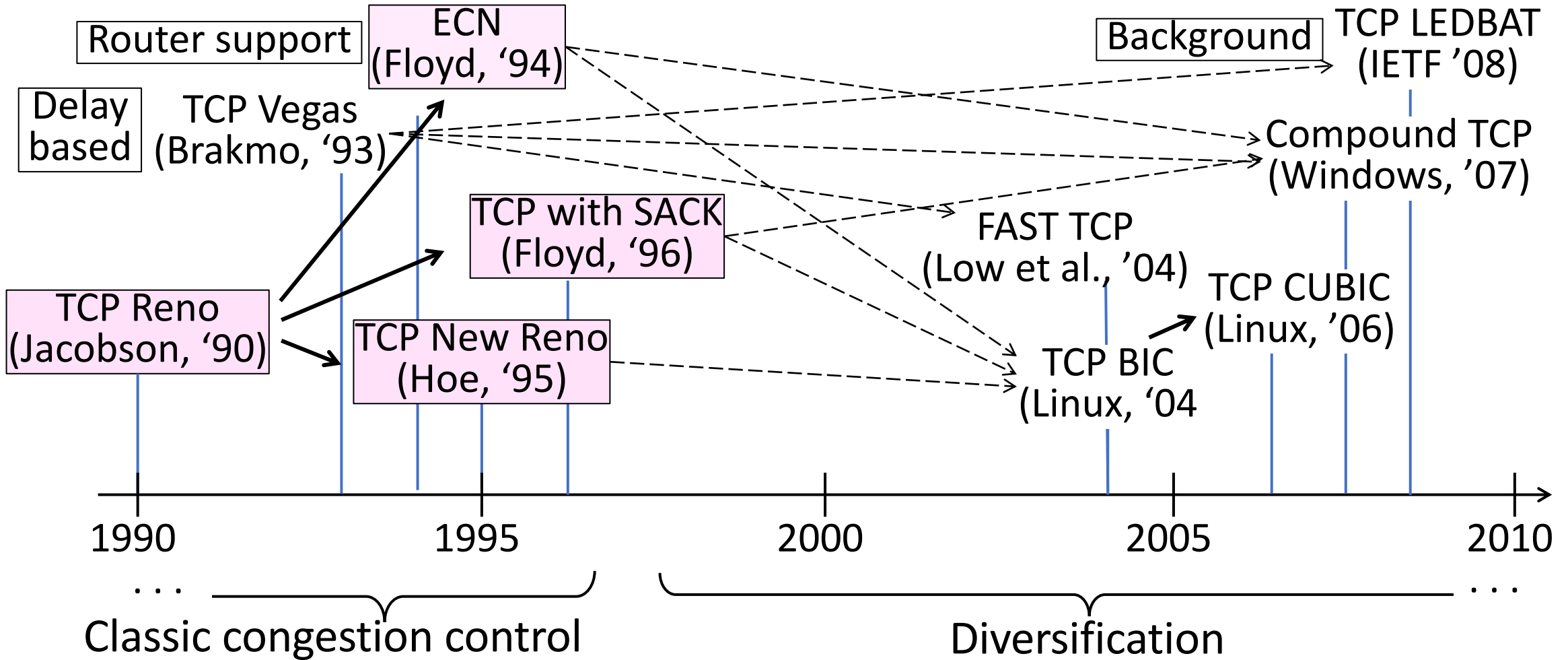
TCP Tahoe/Reno

- TCP extensions and features we will study:
 - AIMD
 - Fair Queuing
 - Slow-start
 - Fast Retransmission
 - Fast Recovery

TCP Timeline



TCP Timeline (2)

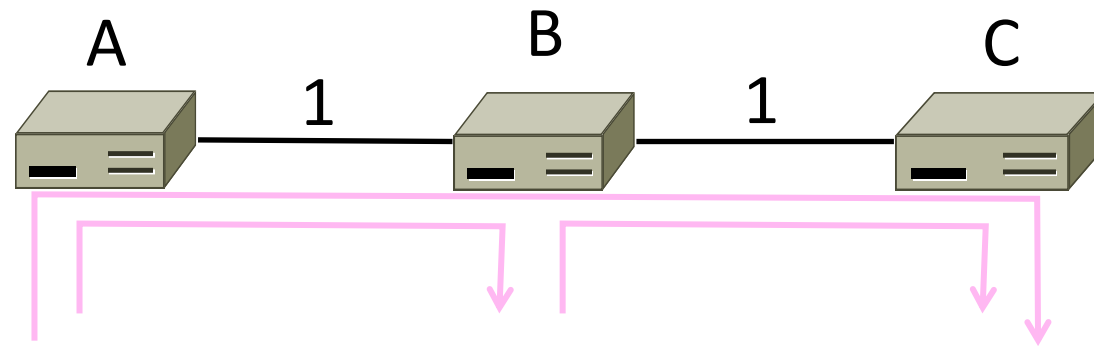


Bandwidth Allocation

- Important task for network is to allocate its capacity to senders
 - Good allocation is both efficient and fair
- Efficient: most capacity is used but there is no congestion
- Fair: every sender gets a reasonable share of the network

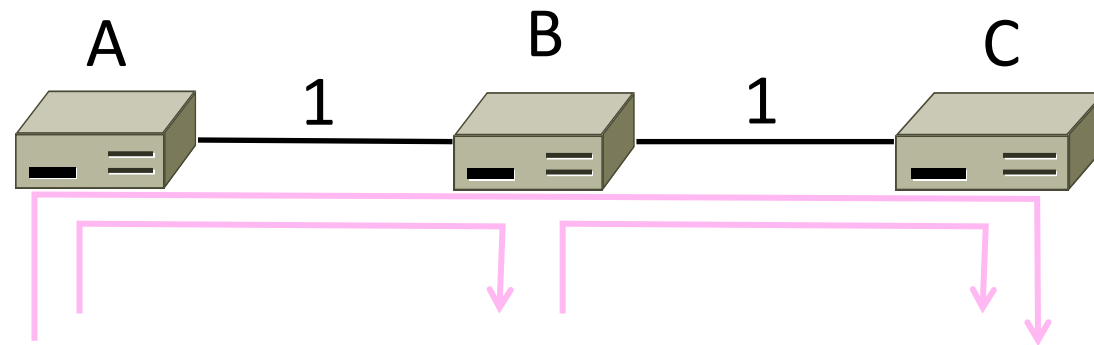
Efficiency vs. Fairness

- Cannot always have both!
 - Example network with traffic:
 - $A \rightarrow B$, $B \rightarrow C$ and $A \rightarrow C$
 - How much traffic can we carry?



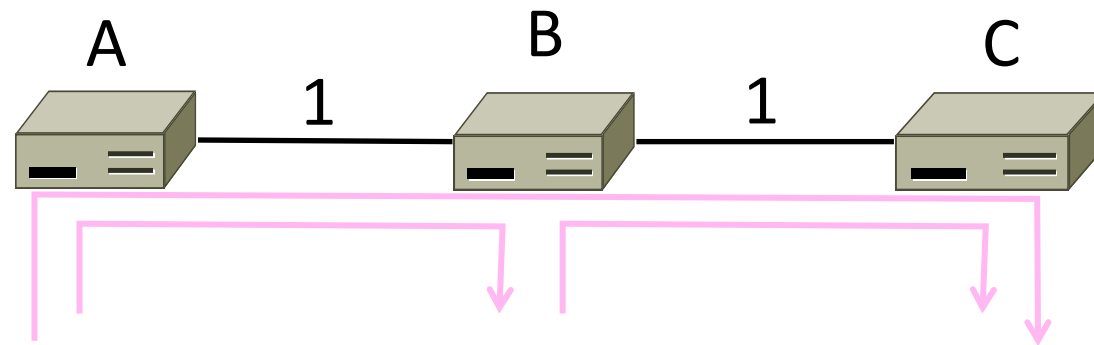
Efficiency vs. Fairness (2)

- If we care about fairness:
 - Give equal bandwidth to each flow
 - $A \rightarrow B$: $\frac{1}{2}$ unit, $B \rightarrow C$: $\frac{1}{2}$, and $A \rightarrow C$, $\frac{1}{2}$
 - Total traffic carried is $1 \frac{1}{2}$ units



Efficiency vs. Fairness (3)

- If we care about efficiency:
 - Maximize total traffic in network
 - $A \rightarrow B$: 1 unit, $B \rightarrow C$: 1, and $A \rightarrow C$, 0
 - Total traffic rises to 2 units!

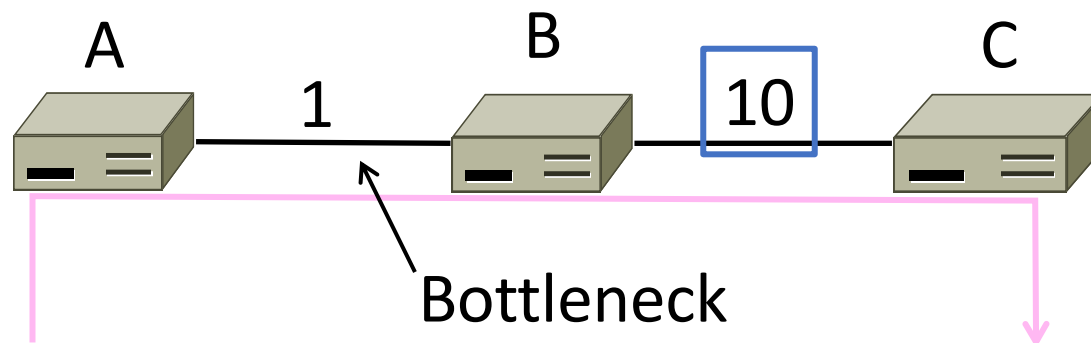


The Slippery Notion of Fairness

- Why is “equal per flow” fair anyway?
 - $A \rightarrow C$ uses more network resources than $A \rightarrow B$ or $B \rightarrow C$
 - Host A sends two flows, B sends one
- Not productive to seek exact fairness
 - More important to avoid starvation
 - A node that cannot use any bandwidth
 - “Equal per flow” is good enough

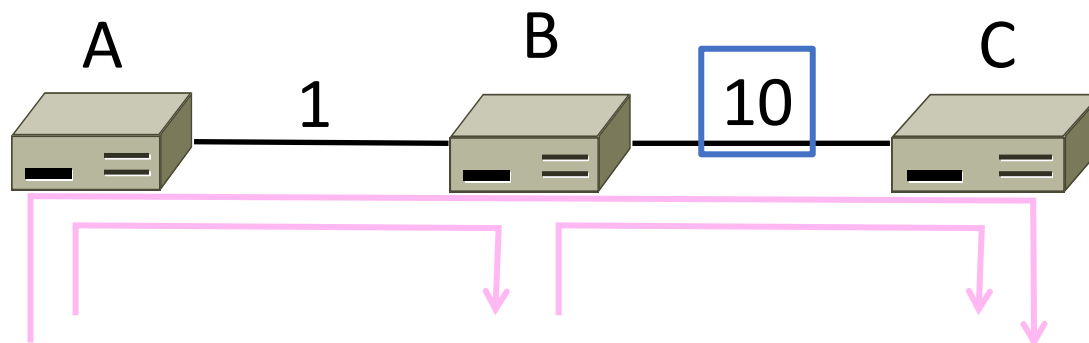
Generalizing “Equal per Flow”

- Bottleneck for a flow of traffic is the link that limits its bandwidth
 - Where congestion occurs for the flow
 - For $A \rightarrow C$, link A–B is the bottleneck



Generalizing “Equal per Flow” (2)

- Flows may have different bottlenecks
 - For $A \rightarrow C$, link $A-B$ is the bottleneck
 - For $B \rightarrow C$, link $B-C$ is the bottleneck
 - Can no longer divide links equally ...



Max-Min Fairness

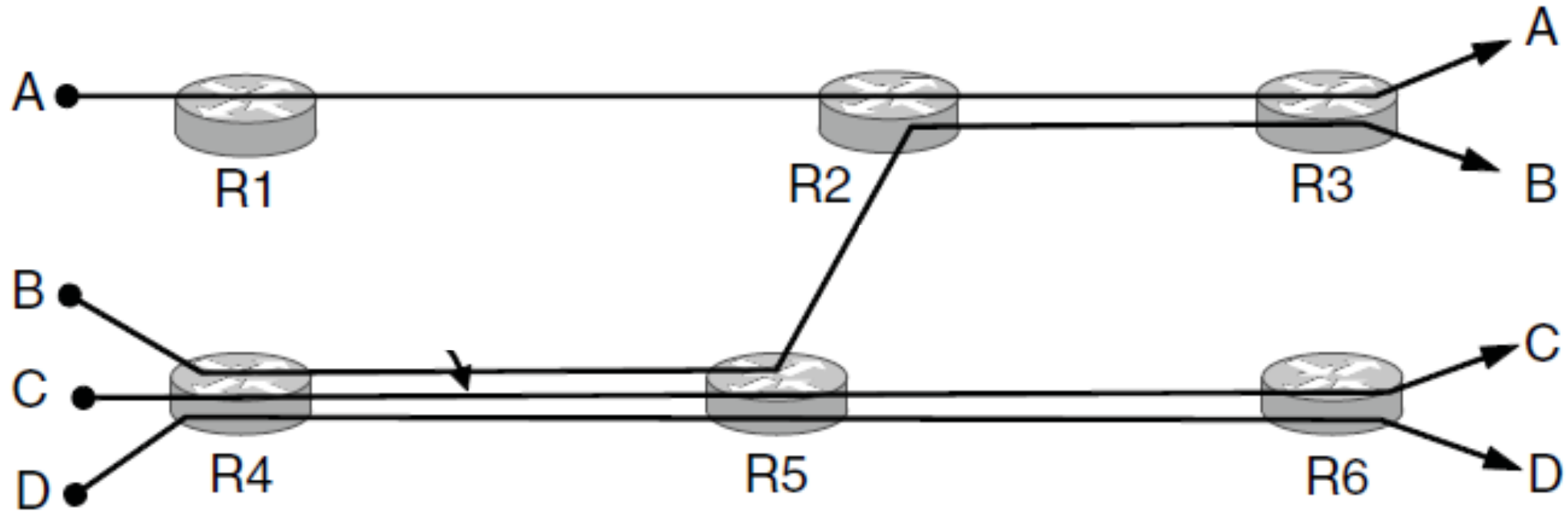
- Intuitively, flows bottlenecked on a link get an equal share of that link
- Max-min fair allocation is one that:
 - Increasing the rate of one flow will decrease the rate of a smaller flow
 - This “maximizes the minimum” flow

Max-Min Fairness (2)

- To find it given a network, imagine “pouring water into the network”
 1. Start with all flows at rate 0
 2. Increase the flows until there is a new bottleneck in the network
 3. Hold fixed the rate of the flows that are bottlenecked
 4. Go to step 2 for any remaining flows

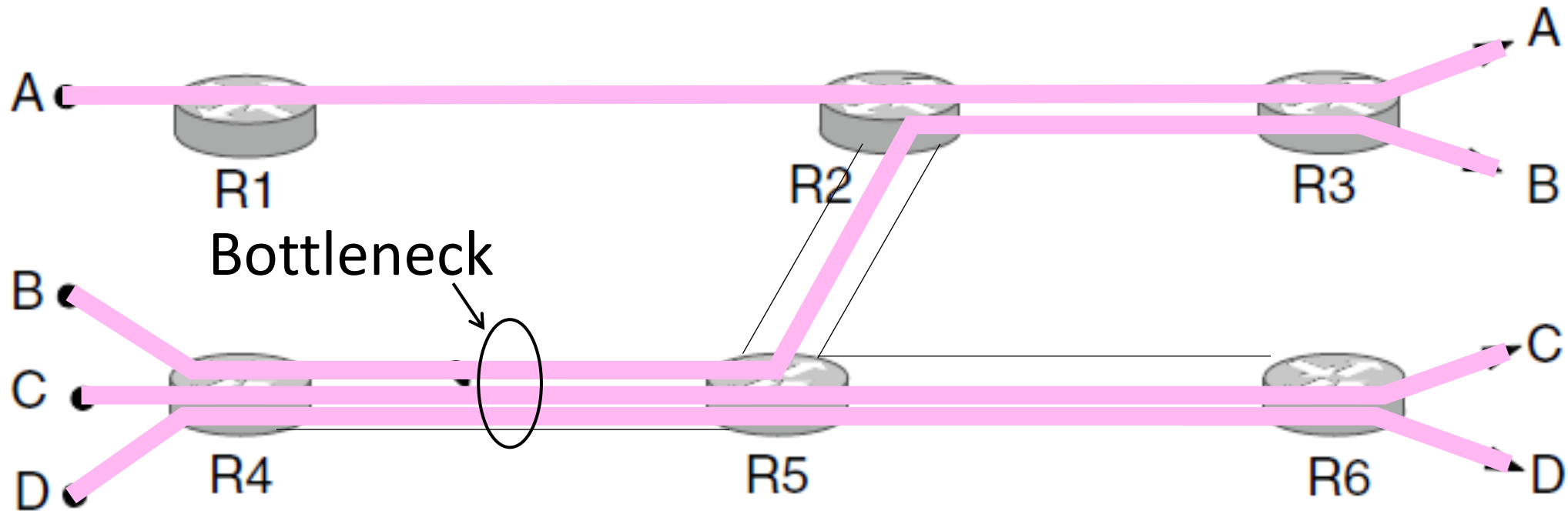
Max-Min Example

- Example: network with 4 flows, link bandwidth = 1
 - What is the max-min fair allocation?



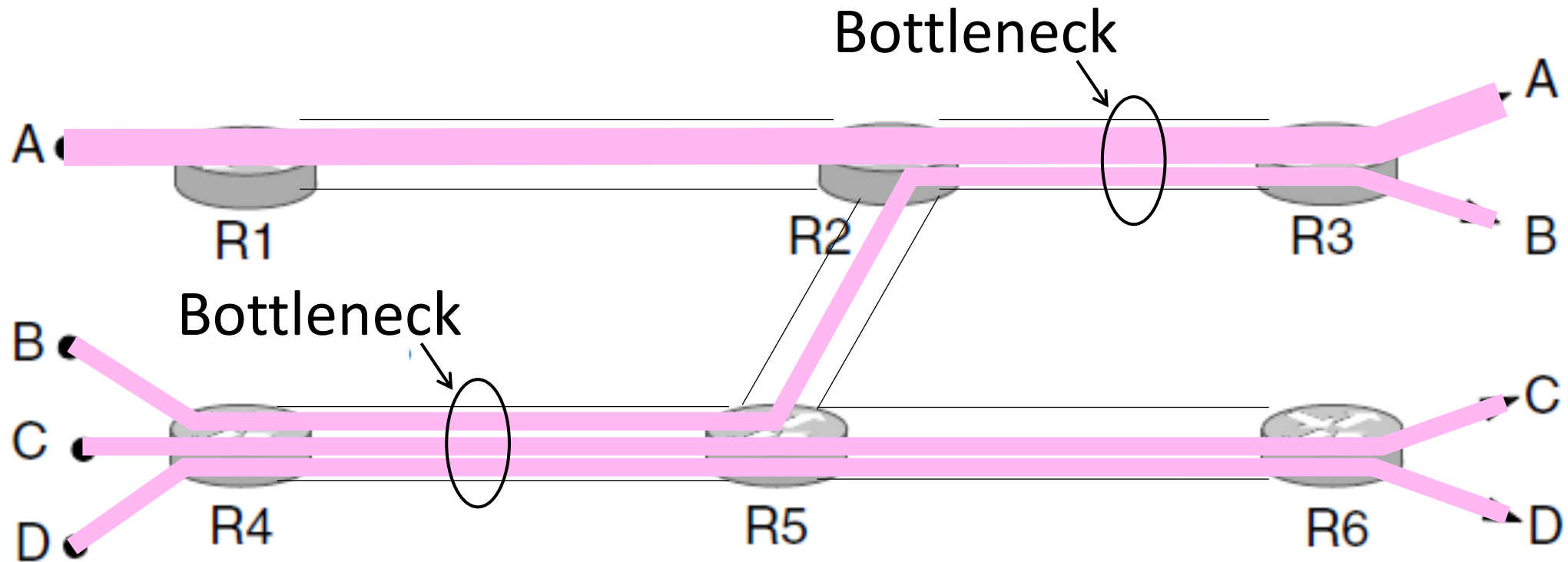
Max-Min Example (2)

- When rate=1/3, flows B, C, and D bottleneck R4—R5
 - Fix B, C, and D, continue to increase A



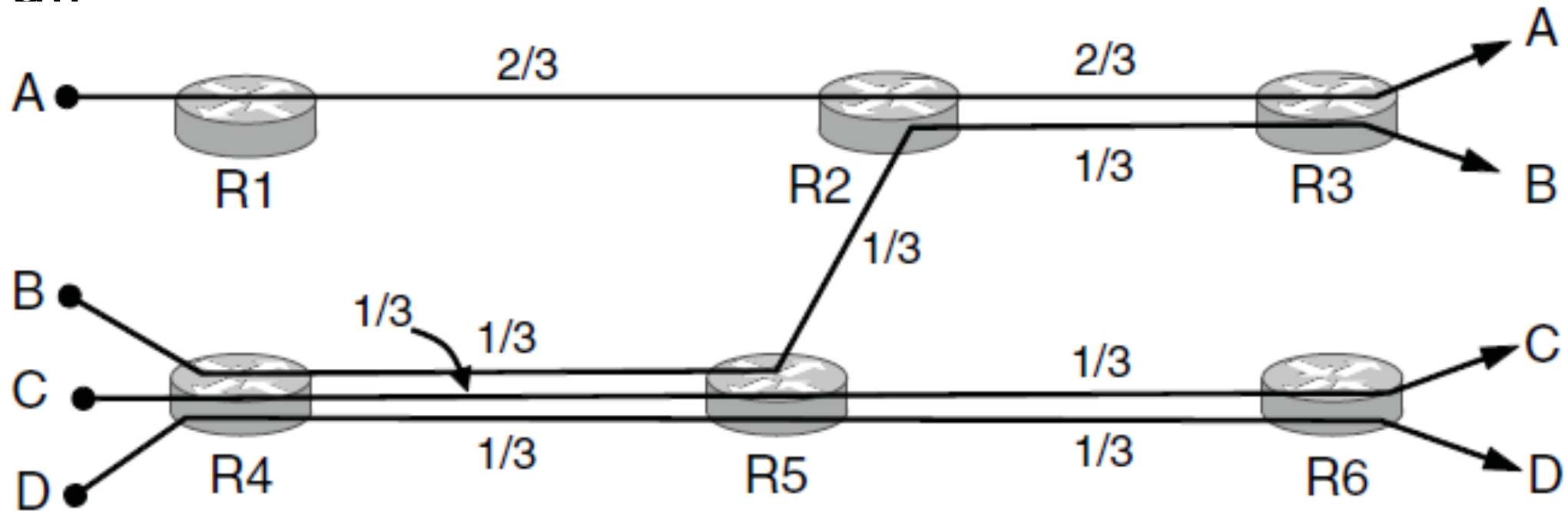
Max-Min Example (3)

- When rate=2/3, flow A bottlenecks R2—R3. Done.



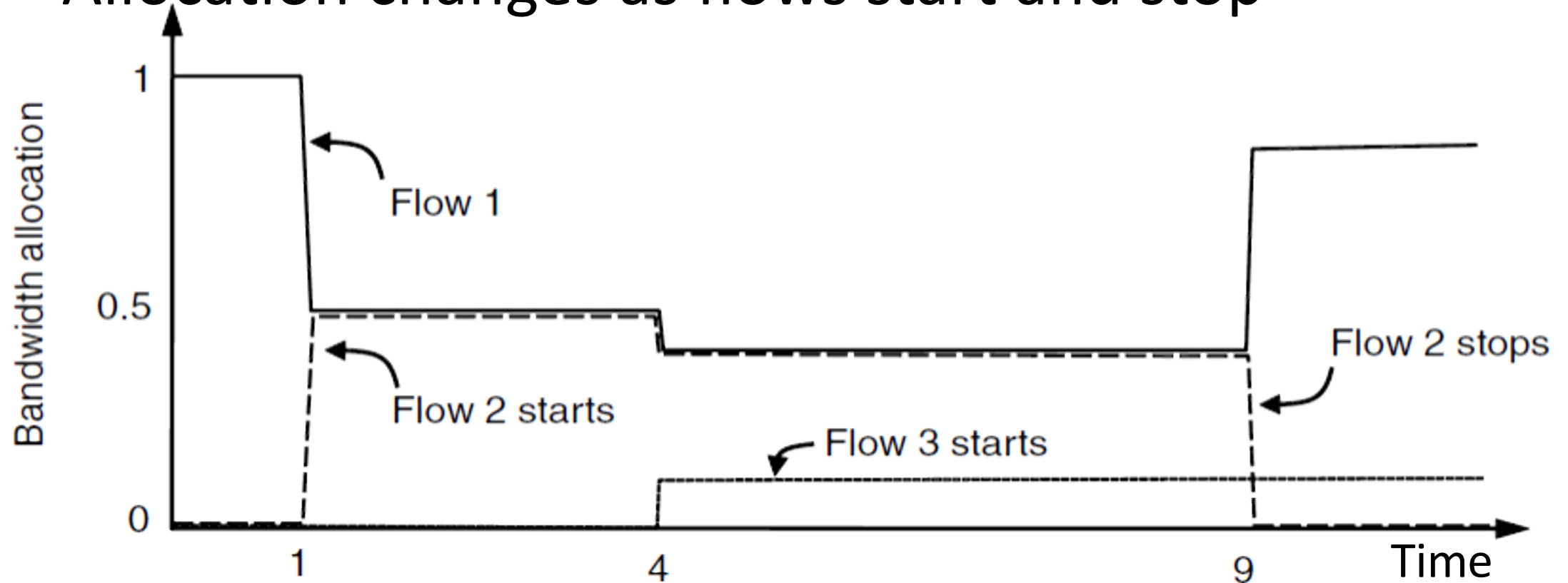
Max-Min Example (4)

- End with $A=2/3$, $B, C, D=1/3$, and $R2-R3$, $R4-R5$ full

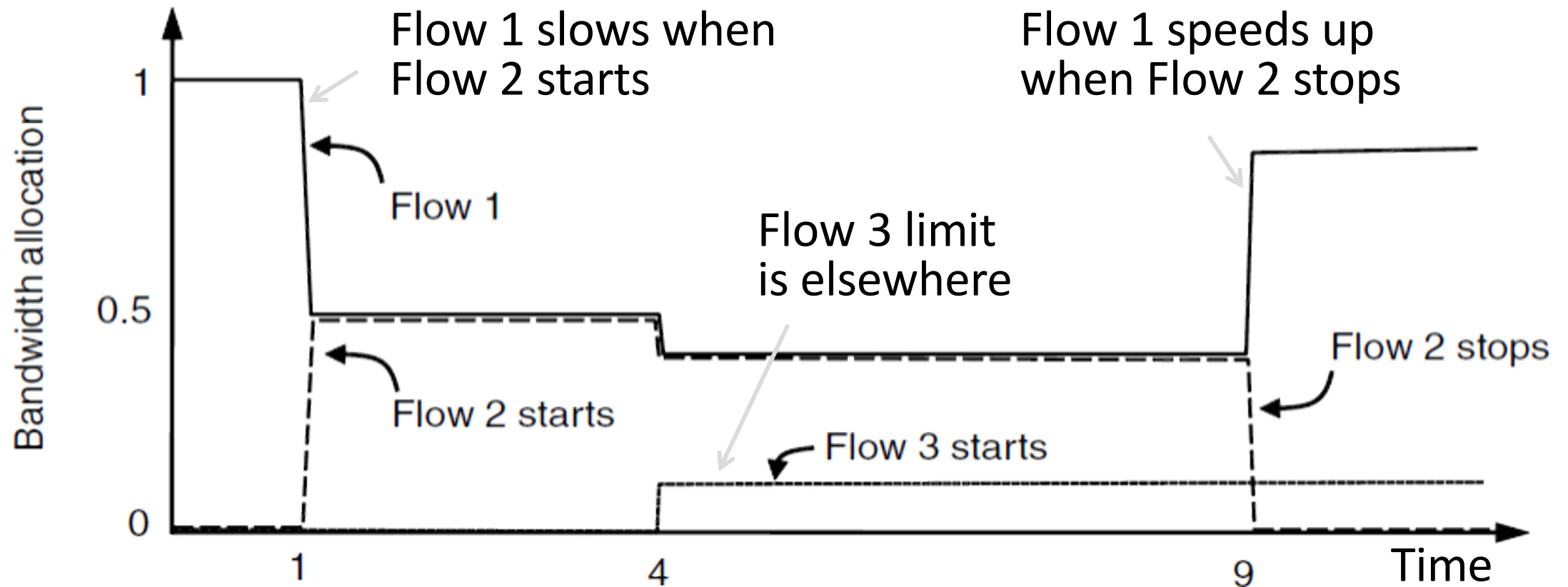


Adapting over Time

- Allocation changes as flows start and stop



Adapting over Time (2)



Why is Bandwidth Allocation hard?

- Number of senders and their offered load changes
- Senders may be limited in other ways
 - Other parts of network or by applications
- Network is distributed; no single party has an overall picture of its state

Bandwidth Allocation Solution Context

In networks without admission control (e.g., Internet)
Transport and Network layers must work together

- Network layer sees congestion
 - Only it can provide direct feedback
- Transport layer causes congestion
 - Only it can reduce load

Bandwidth Allocation Solution Overview

- Senders adapt concurrently based on their own view of the network
- Design this adaptation so the network usage as a whole is efficient and fair
 - In practice, efficiency is more important than fairness
- Adaptation is continuous since offered loads continue to change over time

Bandwidth Allocation Models

- Open loop versus closed loop
 - Open: reserve bandwidth before use
 - Closed: use feedback to adjust rates
- Host versus Network support
 - Who is sets/enforces allocations?
- Window versus Rate based
 - How is allocation expressed?

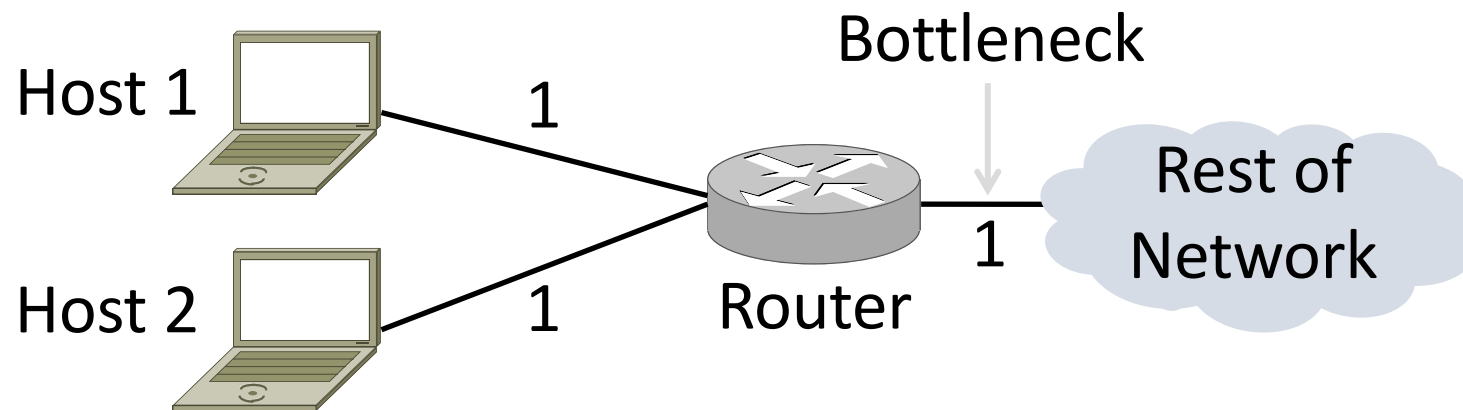
TCP is a closed loop, host-driven, and window-based

Additive Increase Multiplicative Decrease

- AIMD is a control law hosts can use to reach a good allocation
 - Hosts additively increase rate while network not congested
 - Hosts multiplicatively decrease rate when congested
 - Used by TCP
- Let's explore the AIMD game ...

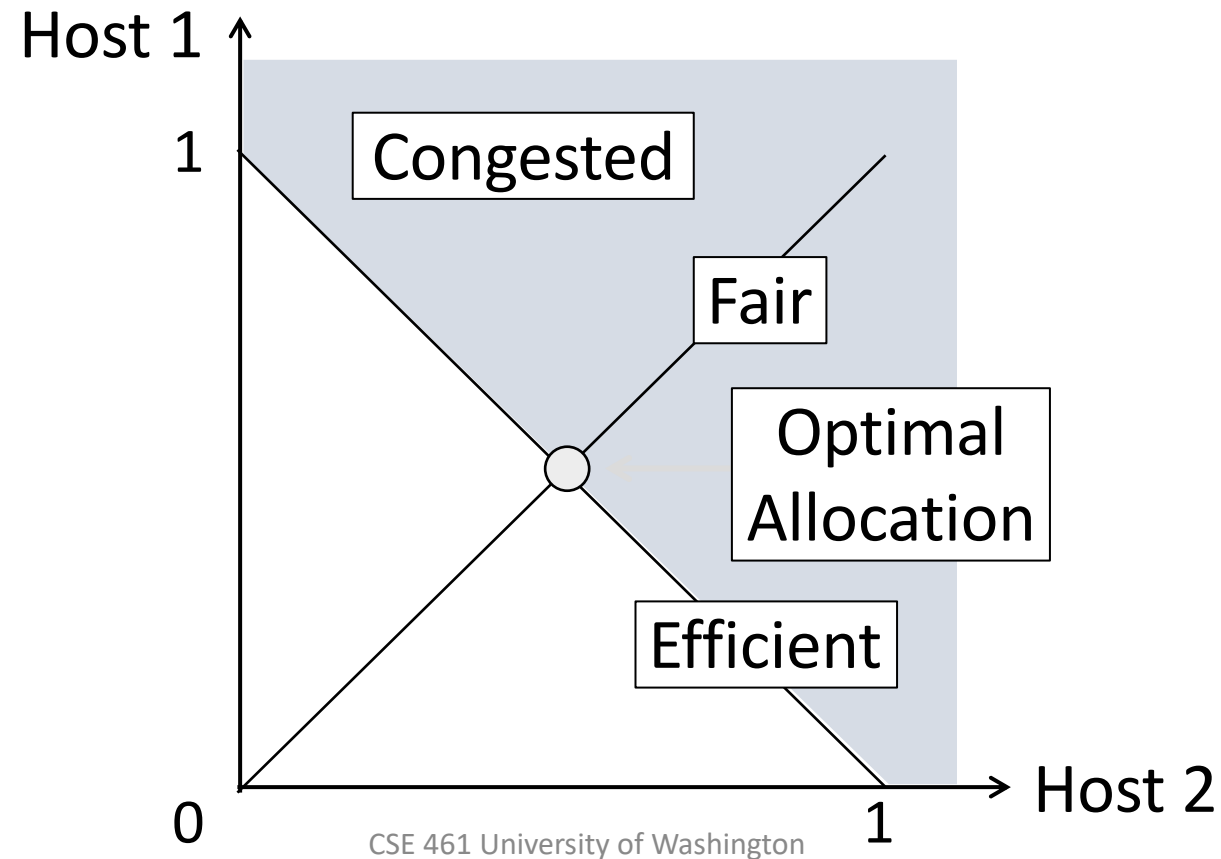
AIMD Game

- Hosts 1 and 2 share a bottleneck
 - But do not talk to each other directly
- Router provides binary feedback
 - Tells hosts if network is congested



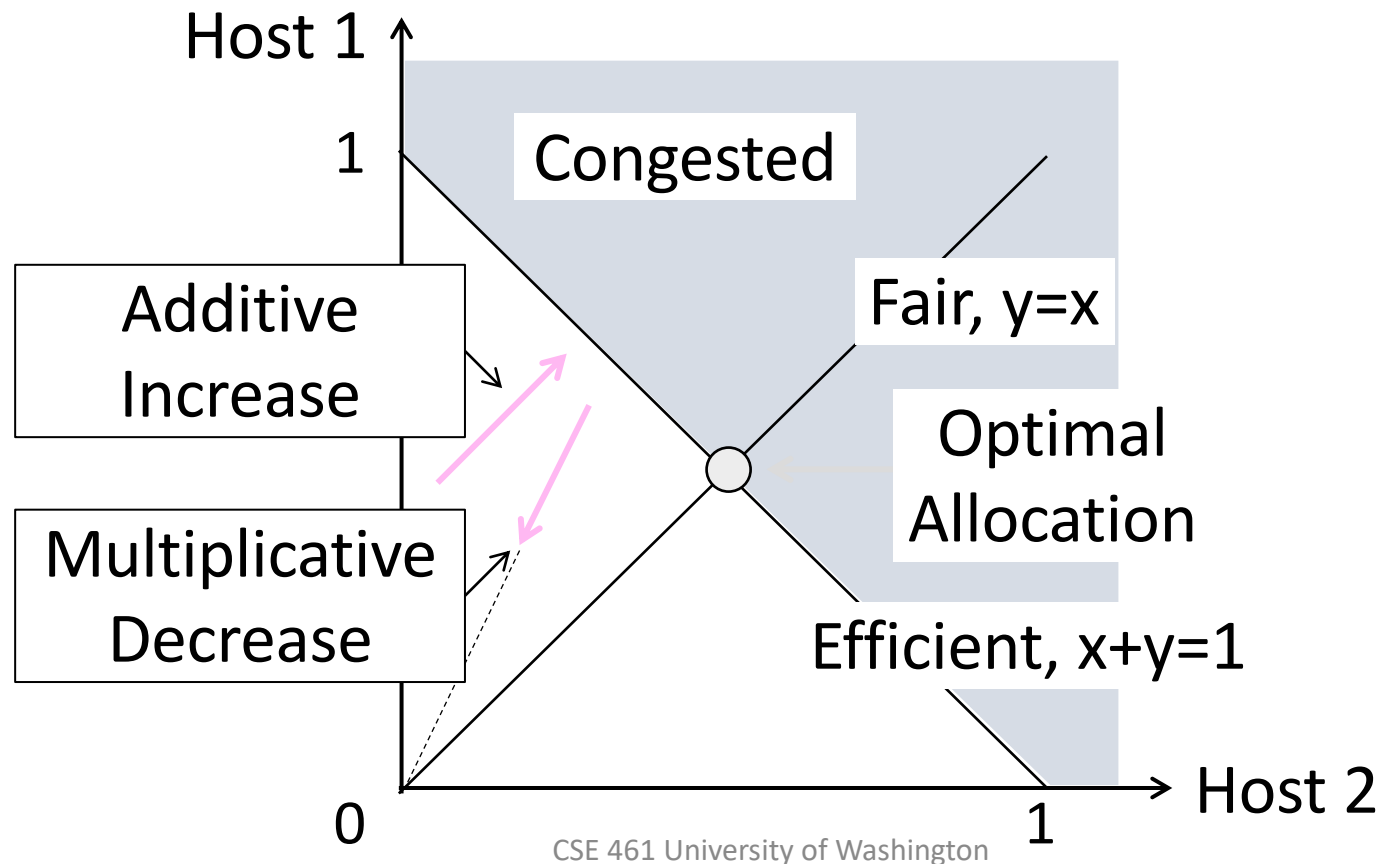
AIMD Game (2)

- Each point is a possible allocation



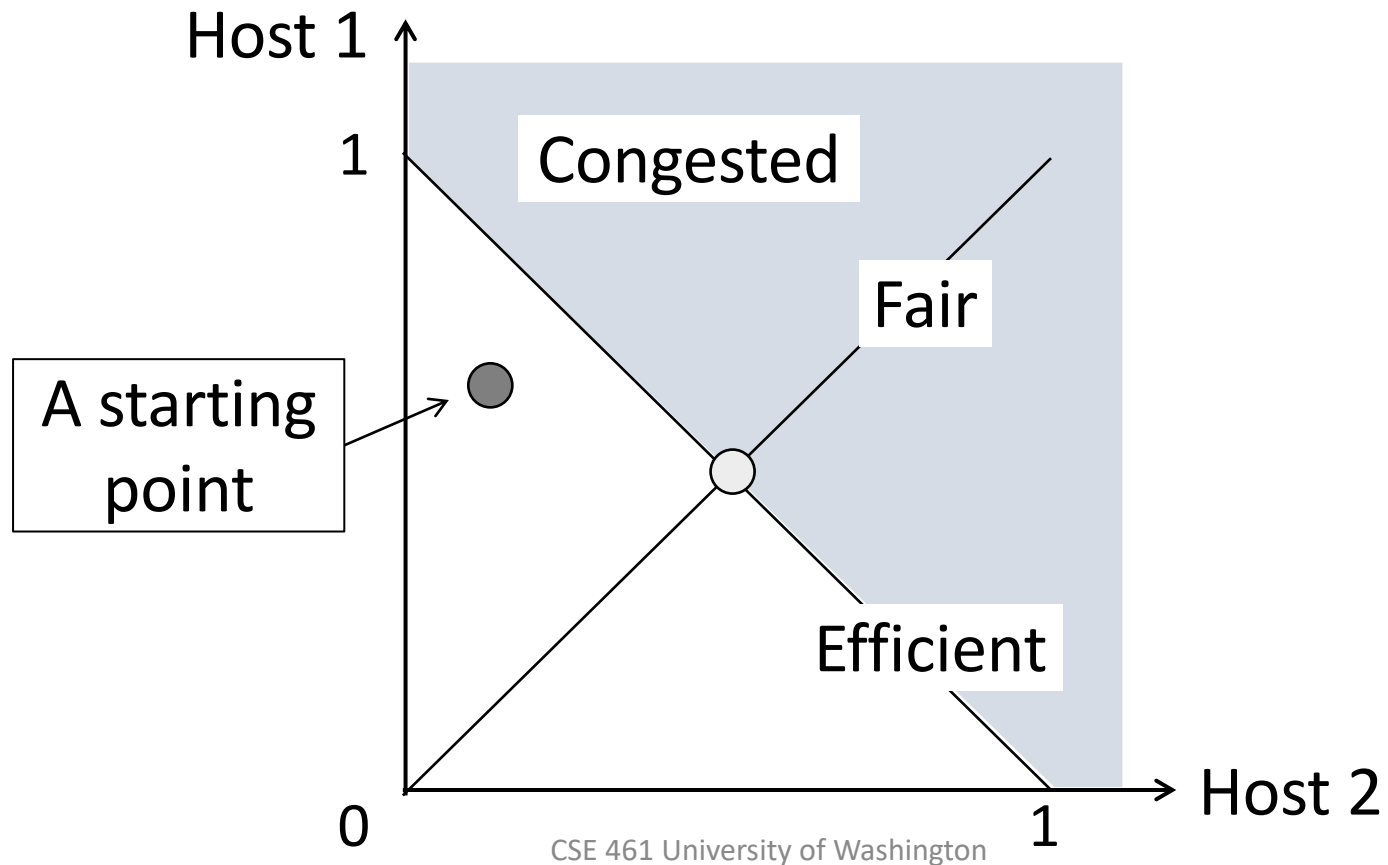
AIMD Game (3)

- AI and MD move the allocation



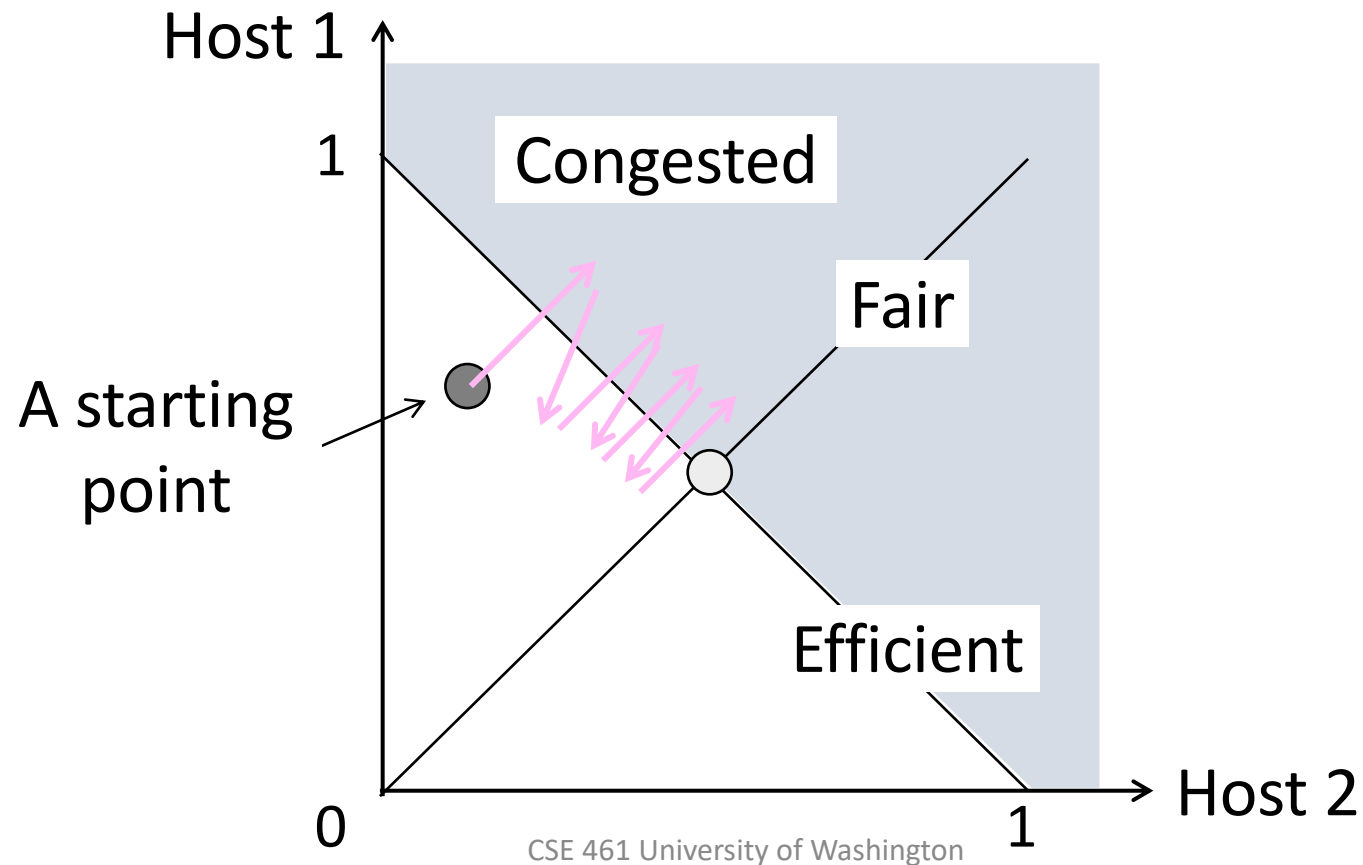
AIMD Game (4)

- Play the game!



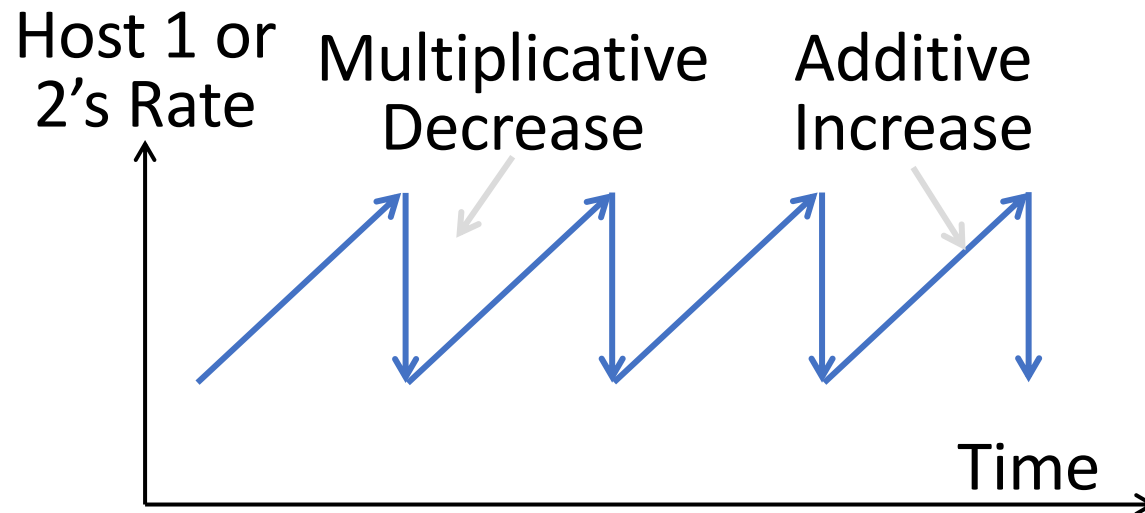
AIMD Game (5)

- Always converge to good allocation!



AIMD Sawtooth

- Produces a “sawtooth” pattern over time for rate of each host
 - This is the TCP sawtooth (later)



AIMD Properties

- Converges to an allocation that is efficient and fair when hosts run it
 - Holds for more general topologies
- Other increase/decrease control laws do not! (Try MIAD, MIMD, MIAD)
- Requires only binary feedback from the network

Feedback Signals

- Several possible signals, with different pros/cons
 - We'll look at classic TCP that uses packet loss as a signal

Signal	Example Protocol	Pros / Cons
Packet loss	TCP NewReno Cubic TCP (Linux)	Hard to get wrong Hear about congestion late Other events can cause loss
Packet delay	BBR (Google)	Hear about congestion early Need to infer congestion
Router indication	TCPs with Explicit Congestion Notification	Hear about congestion early Require router support

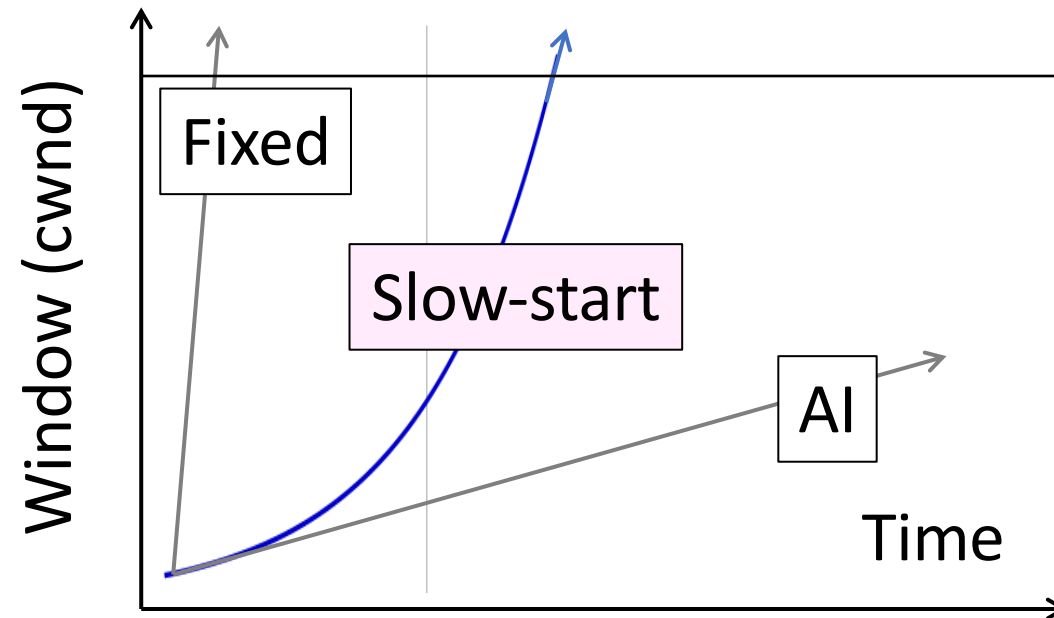
Slow Start (TCP Additive Increase)

TCP “Slow Start” Problem

- We want to quickly reach the right rate, $cwnd_{IDEAL}$, but it varies greatly
 - Fixed sliding window doesn't adapt and is rough on the network (loss!)
 - Additive Increase with small bursts adapts $cwnd$ gently, but might take a long time to become efficient

Slow-Start Solution

- Start by doubling cwnd every RTT
 - Exponential growth (1, 2, 4, 8, 16, ...)
 - Start slow, quickly reach large values

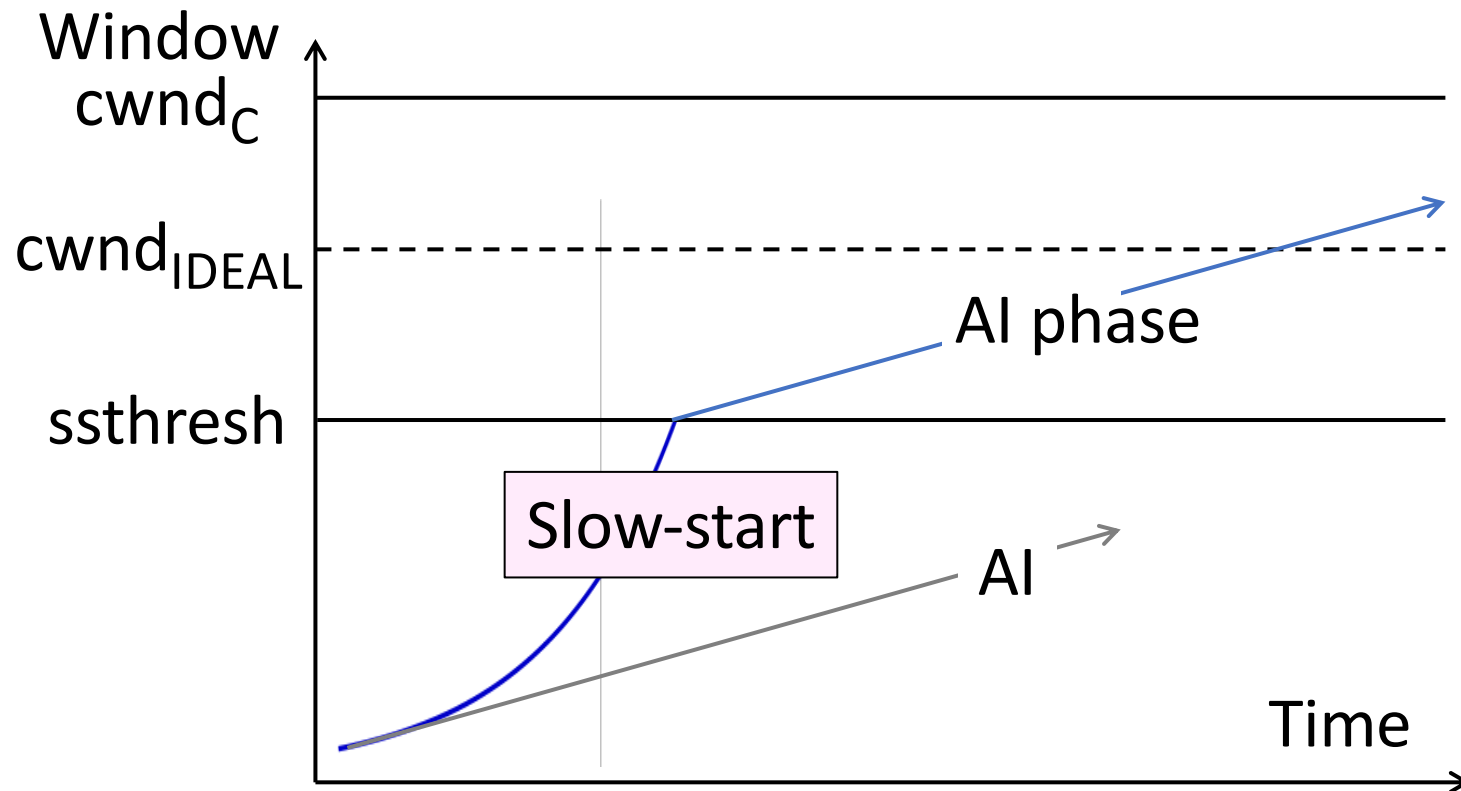


Slow-Start Solution (2)

- Eventually packet loss will occur when the network is congested
 - Loss timeout tells us cwnd is too large
 - Next time, switch to AI beforehand
 - Slowly adapt cwnd near right value
- In terms of cwnd:
 - Expect loss for $\text{cwnd}_c \approx 2BD + \text{queue}$
 - Use $\text{ssthresh} = \text{cwnd}_c / 2$ to switch to AI

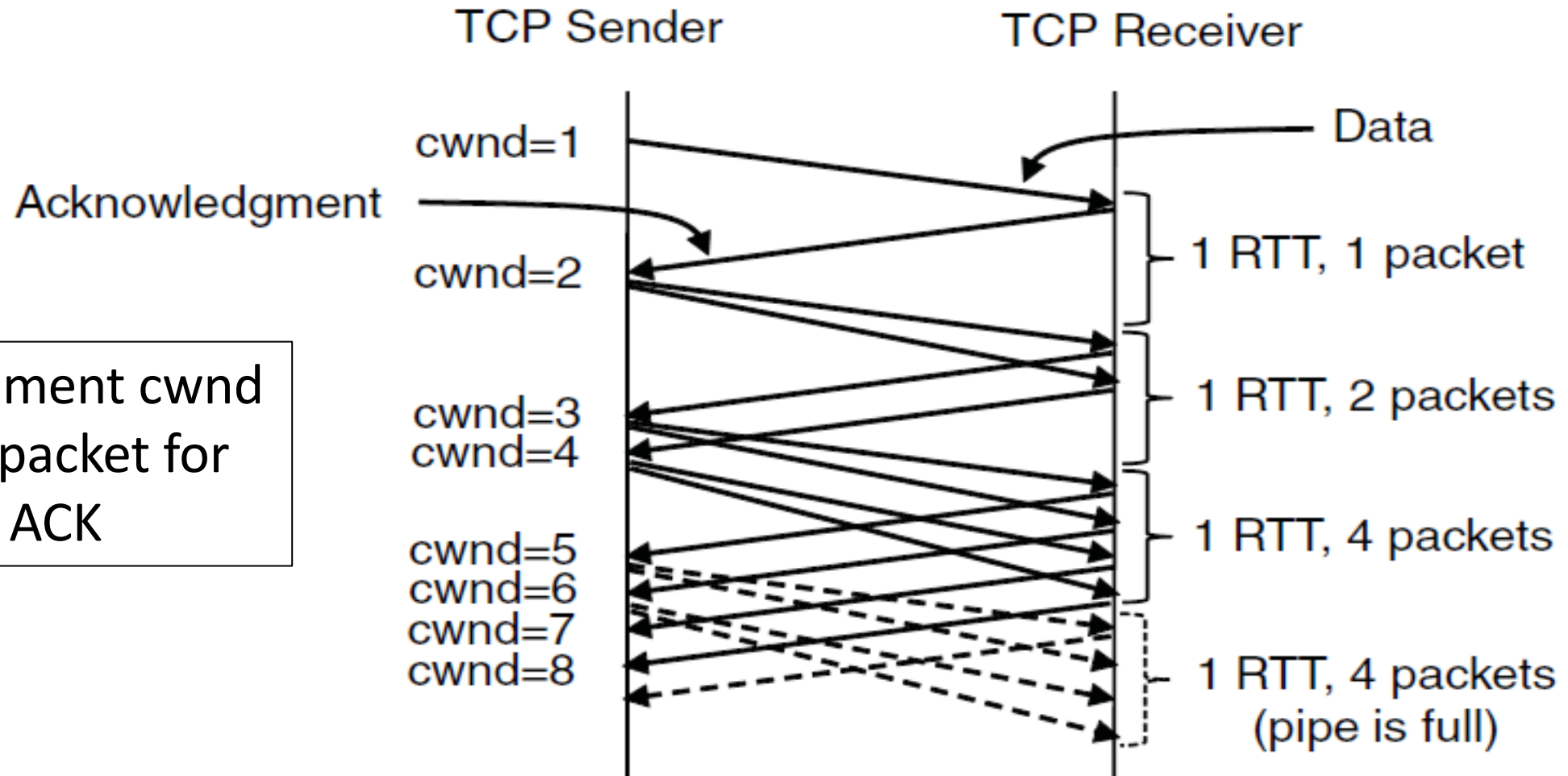
Slow-Start Solution (3)

- Combined behavior, **after first time**
 - Most time spent near right value



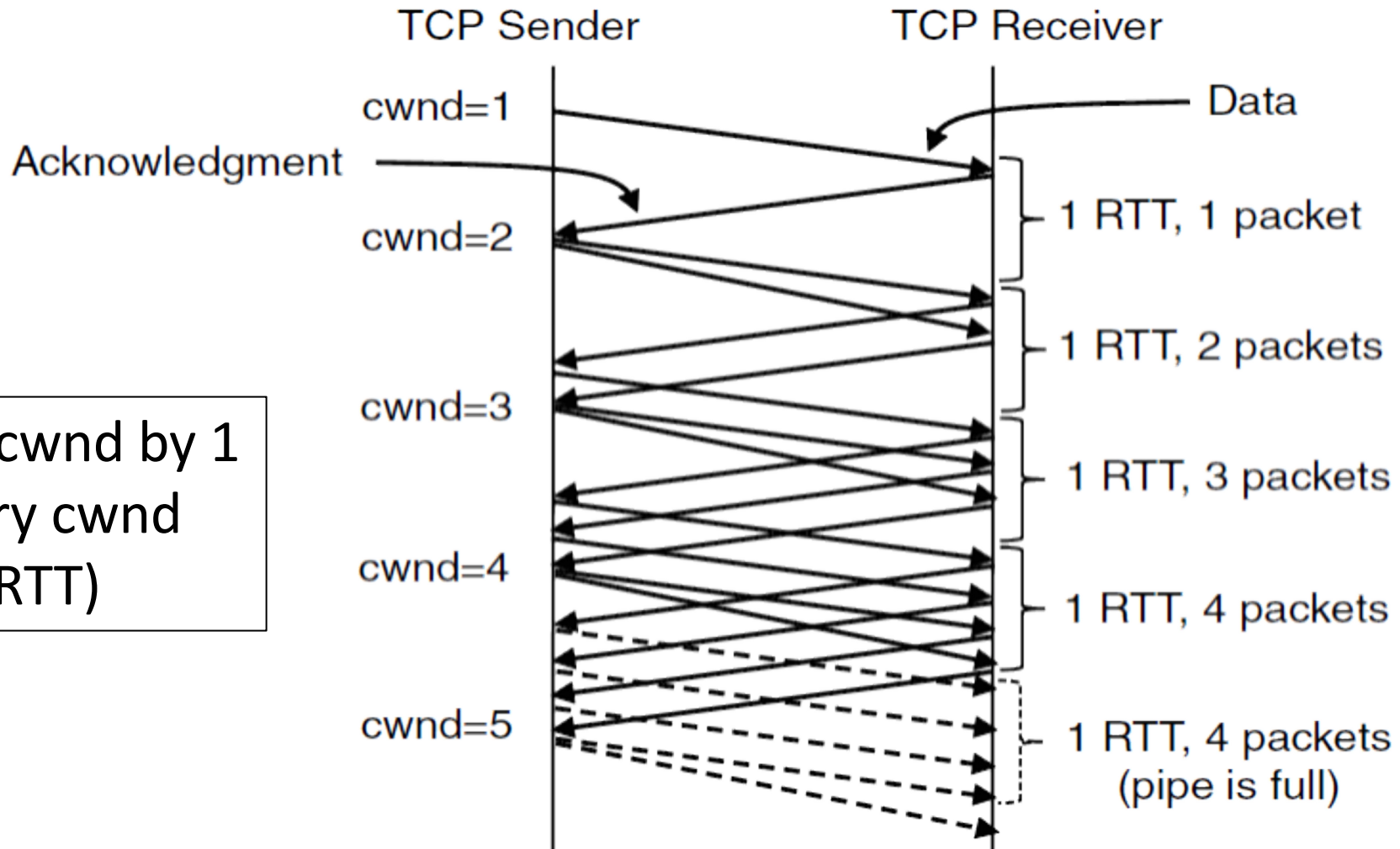
Slow-Start (Doubling) Timeline

Increment cwnd by 1 packet for each ACK



Additive Increase Timeline

Increment cwnd by 1 packet every cwnd ACKs (or 1 RTT)



TCP Tahoe (Implementation)

- Initial slow-start (doubling) phase
 - Start with $cwnd = 1$ (or small value)
 - $cwnd += 1$ packet per ACK
- Later Additive Increase phase
 - $cwnd += 1/cwnd$ packets per ACK
 - Roughly adds 1 packet per RTT
- Switching threshold (initially infinity)
 - Switch to AI when $cwnd > ssthresh$
 - Set $ssthresh = cwnd/2$ after loss
 - Begin with slow-start after timeout