

CS 3700

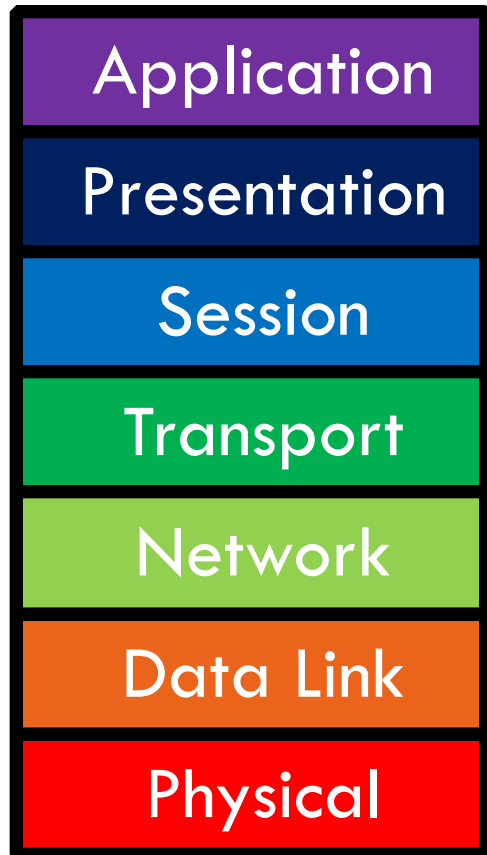
Networks and Distributed Systems

Lecture 10: Congestion Control

Revised 2/9/2014

Transport Layer

2



- Function:
 - ▣ Demultiplexing of data streams
- Optional functions:
 - ▣ Creating long lived connections
 - ▣ Reliable, in-order packet delivery
 - ▣ Error detection
 - ▣ Flow and congestion control
- Key challenges:
 - ▣ Detecting and responding to congestion
 - ▣ Balancing fairness against high utilization

- ❑ Congestion Control
- ❑ Evolution of TCP
- ❑ Problems with TCP

What is Congestion?

4

- Load on the network is higher than capacity

What is Congestion?

4

- Load on the network is higher than capacity
 - ▣ Capacity is not uniform across networks
 - Modem vs. Cellular vs. Cable vs. Fiber Optics
 - ▣ There are multiple flows competing for bandwidth
 - Residential cable modem vs. corporate datacenter
 - ▣ Load is not uniform over time
 - 10pm, Sunday night = Bittorrent Game of Thrones

Why is Congestion Bad?

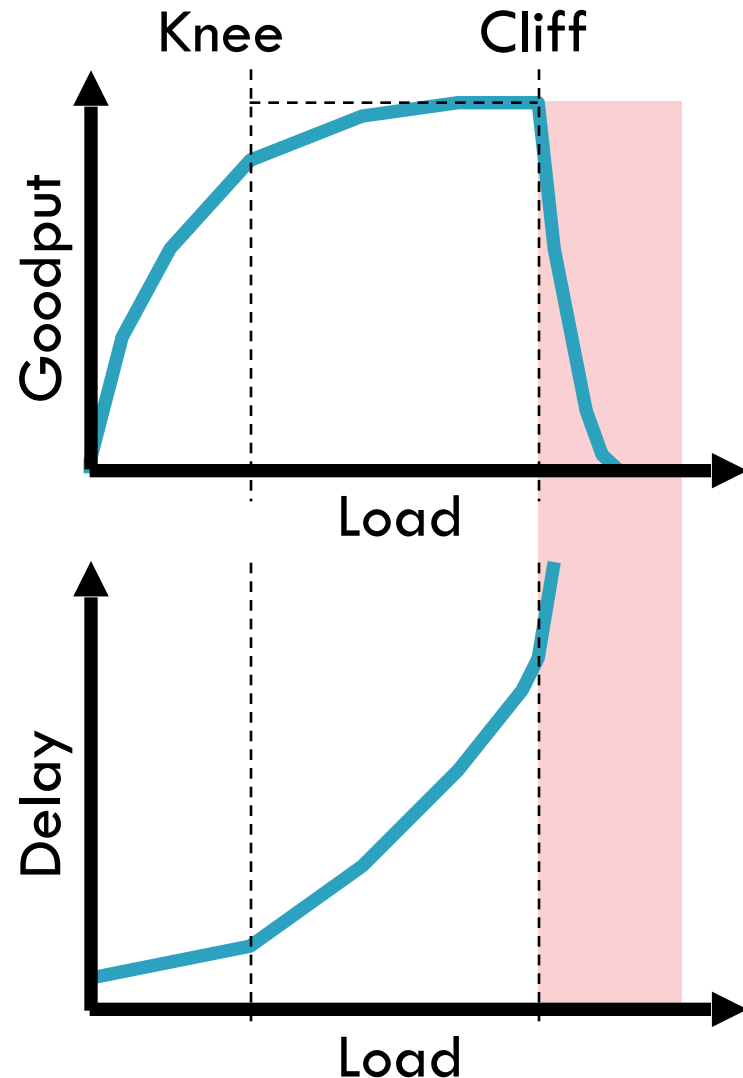
5

- Results in packet **loss**
 - Routers have finite buffers
 - Internet traffic is self similar, no buffer can prevent all drops
 - When routers get overloaded, packets will be dropped
- Practical consequences
 - Router queues build up, **delay** increases
 - Wasted bandwidth from **retransmissions**
 - Low network goodput

The Danger of Increasing Load

6

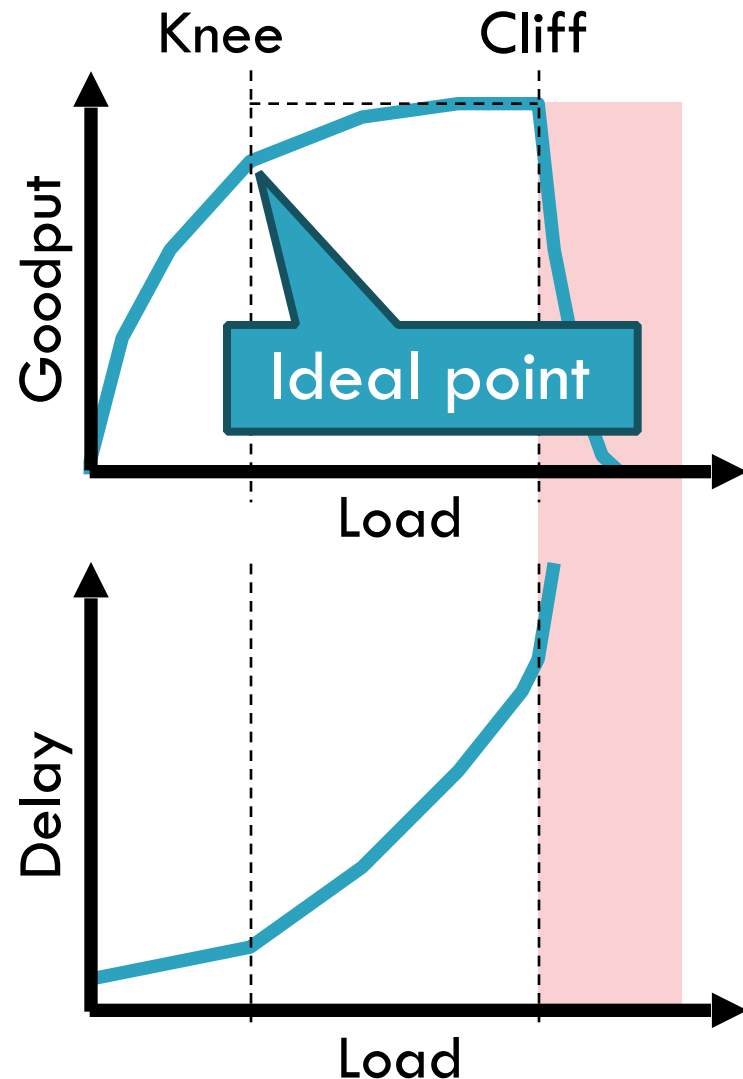
- Knee – point after which
 - ▣ Throughput increases very slow
 - ▣ Delay increases fast
- In an $M/M/1$ queue
 - ▣ Delay = $1 / (1 - \text{utilization})$
- Cliff – point after which
 - ▣ Throughput $\rightarrow 0$
 - ▣ Delay $\rightarrow \infty$



The Danger of Increasing Load

6

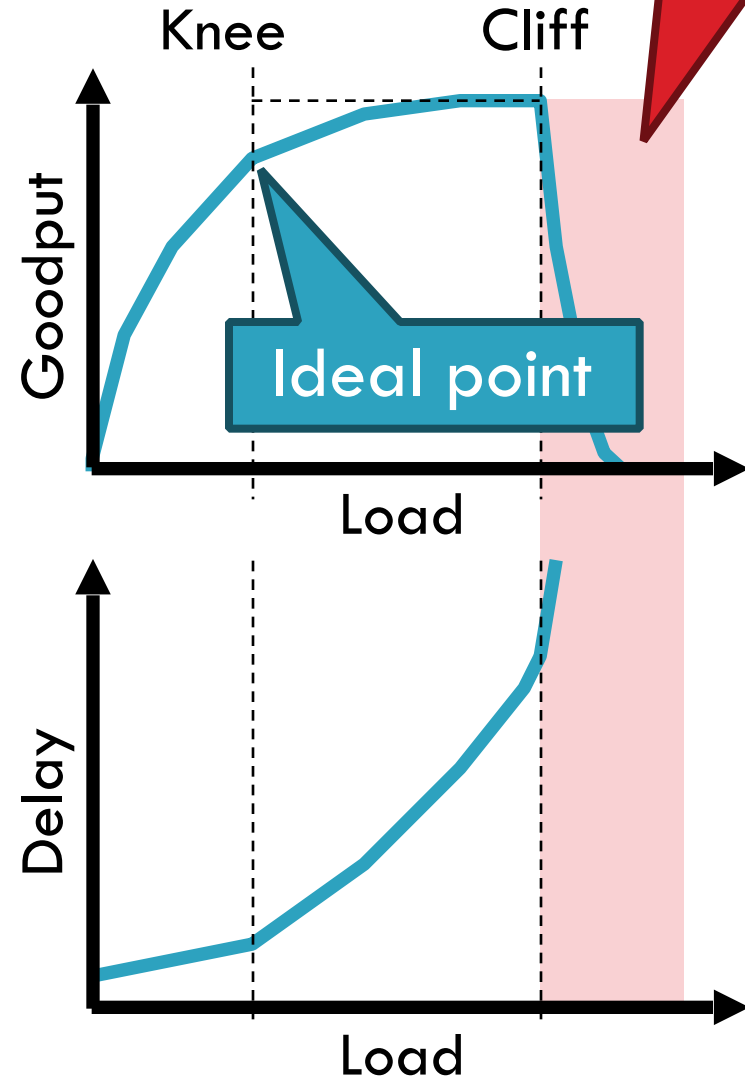
- Knee – point after which
 - ▣ Throughput increases very slow
 - ▣ Delay increases fast
- In an $M/M/1$ queue
 - ▣ Delay = $1 / (1 - \text{utilization})$
- Cliff – point after which
 - ▣ Throughput $\rightarrow 0$
 - ▣ Delay $\rightarrow \infty$



The Danger of Increasing Load

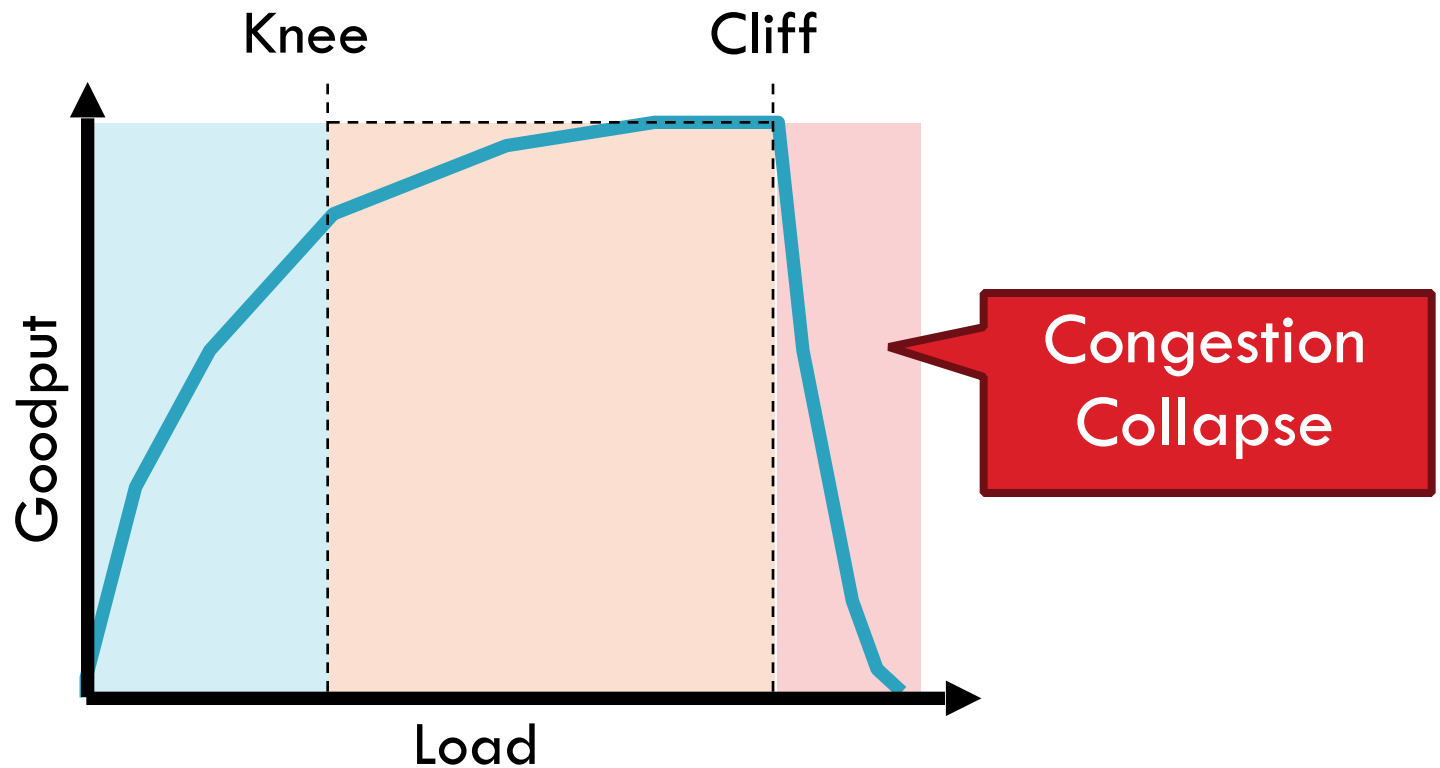
6

- Knee – point after which
 - ▣ Throughput increases very slow
 - ▣ Delay increases fast
- In an $M/M/1$ queue
 - ▣ Delay = $1 / (1 - \text{utilization})$
- Cliff – point after which
 - ▣ Throughput $\rightarrow 0$
 - ▣ Delay $\rightarrow \infty$



Cong. Control vs. Cong. Avoidance

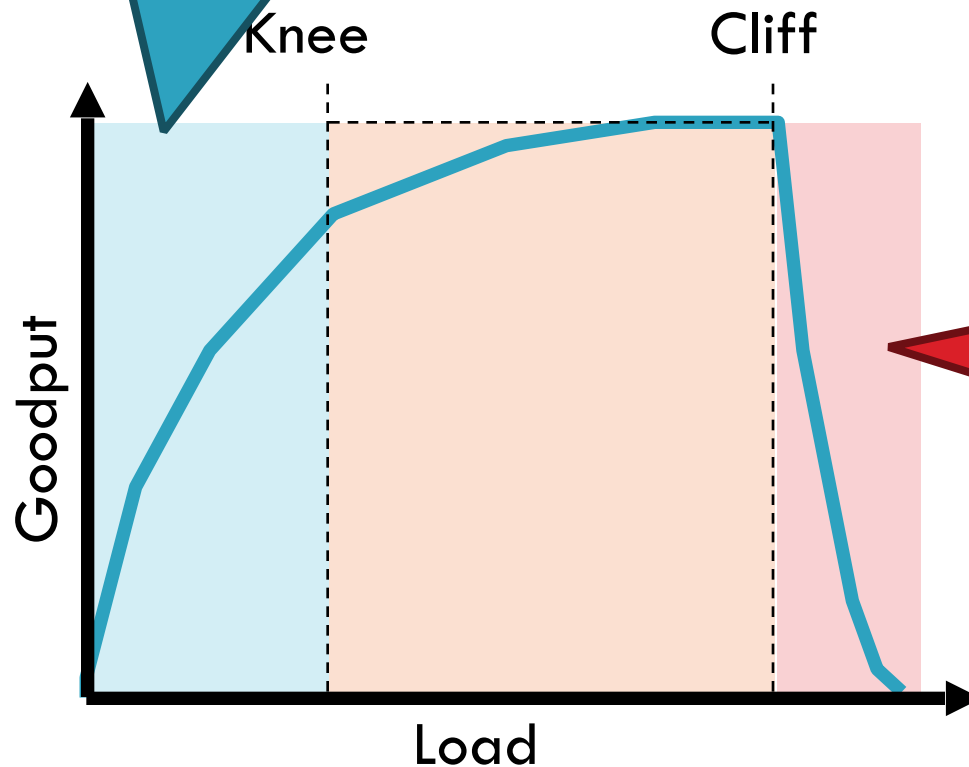
7



Cong. Control vs. Cong. Avoidance

7

Congestion Avoidance:
Stay left of the knee



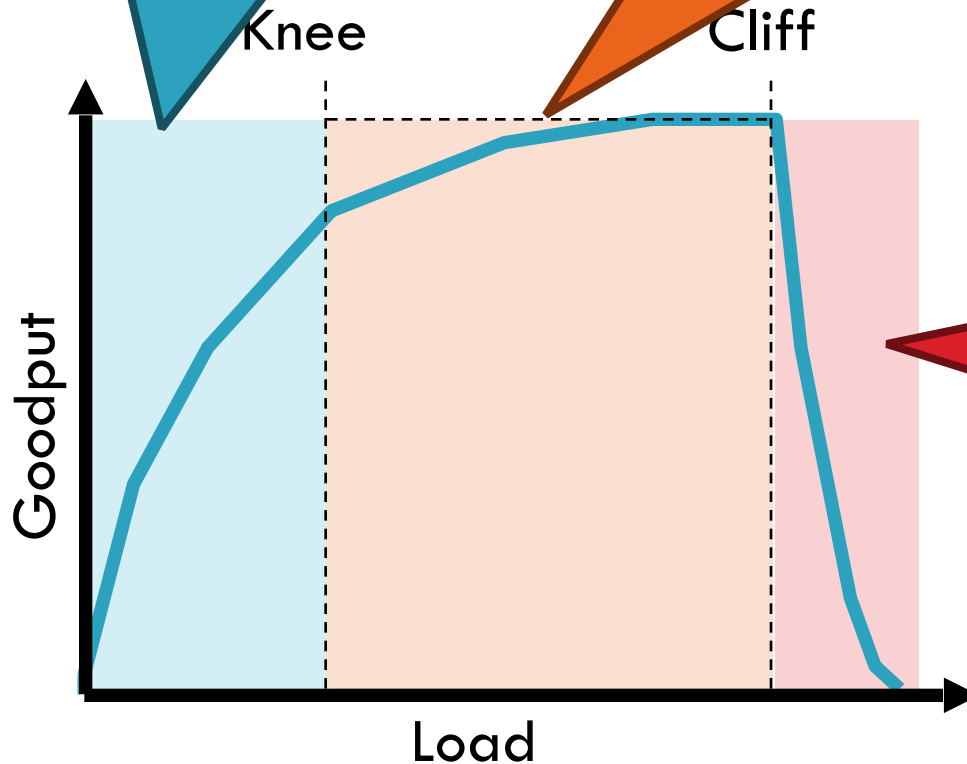
Congestion
Collapse

Cong. Control vs. Cong. Avoidance

7

Congestion Avoidance:
Stay left of the knee

Congestion Control:
Stay left of the cliff



Congestion
Collapse

Advertised Window, Revisited

8

- Does TCP's advertised window solve congestion?

Advertised Window, Revisited

8

- Does TCP's advertised window solve congestion?

NO

- The advertised window only protects the receiver
- A sufficiently fast receiver can max the window
 - ▣ What if the network is slower than the receiver?
 - ▣ What if there are other concurrent flows?

Advertised Window, Revisited

8

- Does TCP's advertised window solve congestion?

NO

- The advertised window only protects the receiver
- A sufficiently fast receiver can max the window
 - ▣ What if the network is slower than the receiver?
 - ▣ What if there are other concurrent flows?
- Key points
 - ▣ Window size determines send rate
 - ▣ Window must be adjusted to prevent congestion collapse

Goals of Congestion Control

Goals of Congestion Control

9

1. Adjusting to the bottleneck bandwidth
2. Adjusting to variations in bandwidth
3. Sharing bandwidth between flows
4. Maximizing throughput

General Approaches

10

- Do nothing, send packets indiscriminately
 - ▣ Many packets will drop, totally unpredictable performance
 - ▣ **May lead to congestion collapse**

General Approaches

10

- Do nothing, send packets indiscriminately
 - ▣ Many packets will drop, totally unpredictable performance
 - ▣ **May lead to congestion collapse**
- Reservations
 - ▣ Pre-arrange bandwidth allocations for flows
 - ▣ Requires negotiation before sending packets
 - ▣ **Must be supported by the network**

General Approaches

10

- Do nothing, send packets indiscriminately
 - ▣ Many packets will drop, totally unpredictable performance
 - ▣ **May lead to congestion collapse**
- Reservations
 - ▣ Pre-arrange bandwidth allocations for flows
 - ▣ Requires negotiation before sending packets
 - ▣ **Must be supported by the network**
- Dynamic adjustment
 - ▣ Use probes to estimate level of congestion
 - ▣ Speed up when congestion is low
 - ▣ Slow down when congestion increases
 - ▣ **Messy dynamics, requires distributed coordination**

General Approaches

10

- Do nothing, send packets indiscriminately
 - ▣ Many packets will drop, totally unpredictable performance
 - ▣ **May lead to congestion collapse**
- Reservations
 - ▣ Pre-arrange bandwidth allocations for flows
 - ▣ Requires negotiation before sending packets
 - ▣ **Must be supported by the network**
- Dynamic adjustment
 - ▣ Use probes to estimate level of congestion
 - ▣ Speed up when congestion is low
 - ▣ Slow down when congestion increases
 - ▣ **Messy dynamics, requires distributed coordination**

TCP Congestion Control

11

- Each TCP connection has a window
 - ▣ Controls the number of unACKed packets
- Sending rate is $\sim \text{window}/\text{RTT}$
- Idea: vary the window size to control the send rate

TCP Congestion Control

11

- Each TCP connection has a window
 - ▣ Controls the number of unACKed packets
- Sending rate is $\sim \text{window}/\text{RTT}$
- Idea: vary the window size to control the send rate
- Introduce a **congestion window** at the sender
 - ▣ Congestion control is sender-side problem

Congestion Window (*cwnd*)

12

- Limits how much data is in transit
 - Denominated in bytes
1. $wnd = \min(cwnd, adv_wnd);$
 2. $effective_wnd = wnd - (last_byte_sent - last_byte_acked);$

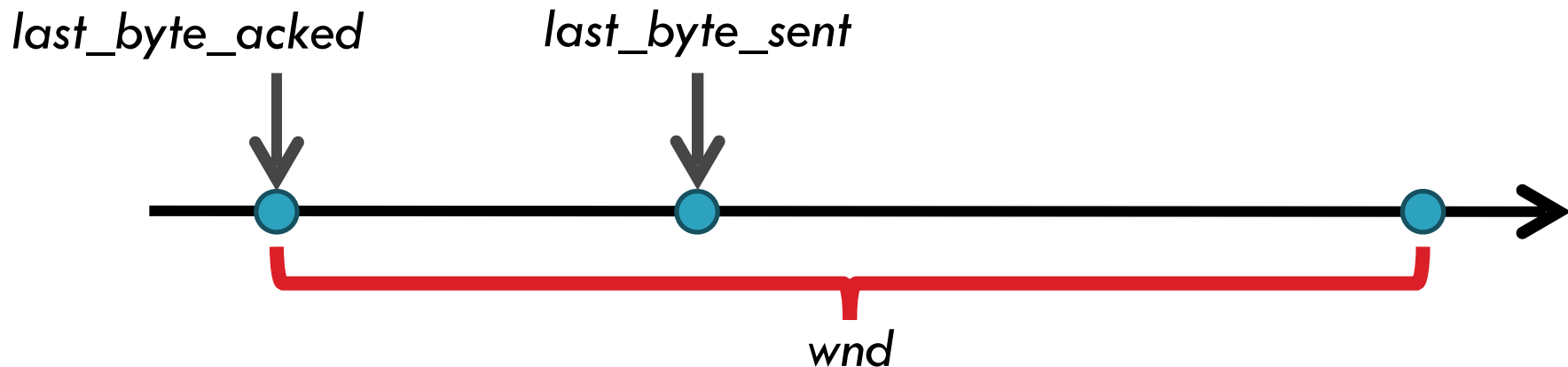
Congestion Window (cwnd)

12

- Limits how much data is in transit
- Denominated in bytes

1. $wnd = \min(cwnd, adv_wnd);$

2. $effective_wnd = wnd - (last_byte_sent - last_byte_acked);$



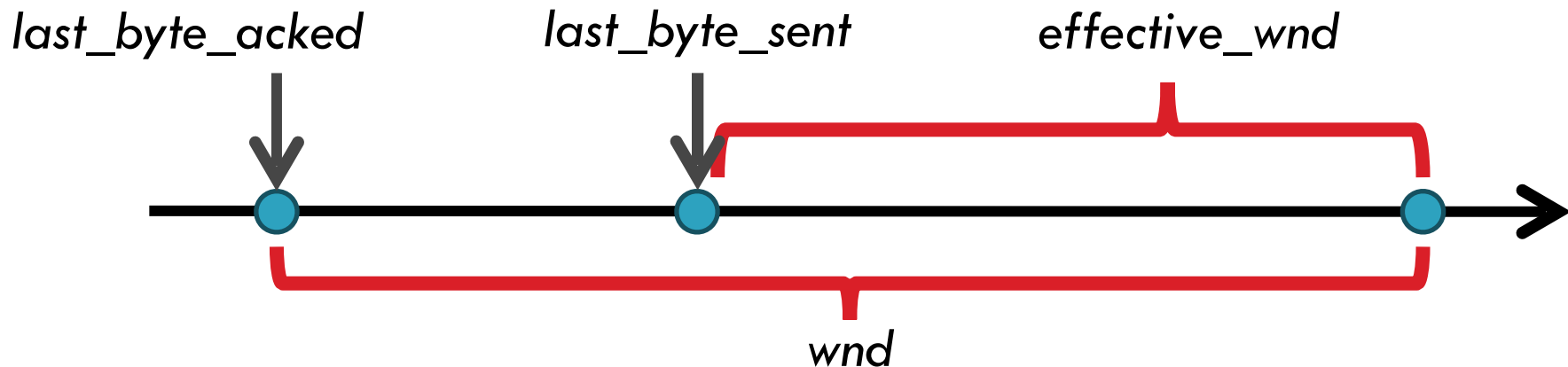
Congestion Window (cwnd)

12

- Limits how much data is in transit
- Denominated in bytes

1. $wnd = \min(cwnd, adv_wnd);$

2. $effective_wnd = wnd - (last_byte_sent - last_byte_acked);$



Two Basic Components

13

1. Detect congestion

Two Basic Components

13

1. Detect congestion

- ▣ Packet dropping is most reliable signal
 - Delay-based methods are hard and risky
- ▣ How do you detect packet drops? ACKs
 - Timeout after not receiving an ACK
 - Several duplicate ACKs in a row (ignore for now)

Two Basic Components

13

1. Detect congestion

- ▣ Packet dropping is most reliable signal
 - Delay-based methods are hard and risky
- ▣ How do you detect packet drops? ACKs
 - Timeout after not receiving an ACK
 - Several duplicate ACKs in a row (ignore for now)

Except on
wireless
networks

Two Basic Components

13

1. Detect congestion

- ❑ Packet dropping is most reliable signal
 - Delay-based methods are hard and risky
- ❑ How do you detect packet drops? ACKs
 - Timeout after not receiving an ACK
 - Several duplicate ACKs in a row (ignore for now)

Except on
wireless
networks

2. Rate adjustment algorithm

- ❑ Modify *cwnd*
- ❑ Probe for bandwidth
- ❑ Responding to congestion

Rate Adjustment

14

- Recall: TCP is ACK clocked
 - Congestion = delay = long wait between ACKs
 - No congestion = low delay = ACKs arrive quickly

Rate Adjustment

14

- Recall: TCP is ACK clocked
 - Congestion = delay = long wait between ACKs
 - No congestion = low delay = ACKs arrive quickly
- Basic algorithm
 - Upon receipt of ACK: increase *cwnd*
 - Data was delivered, perhaps we can send faster
 - *cwnd* growth is proportional to RTT
 - On loss: decrease *cwnd*
 - Data is being lost, there must be congestion

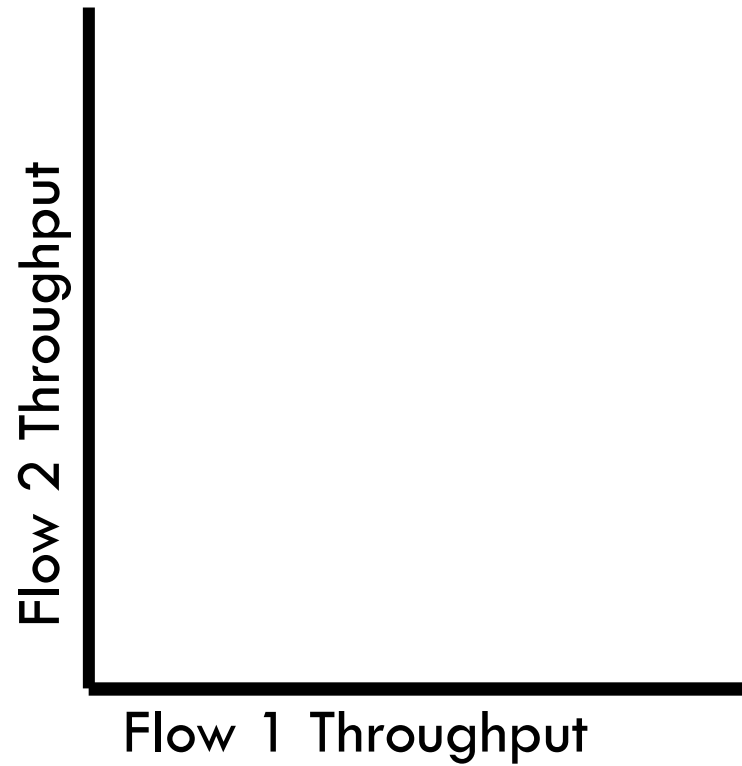
Rate Adjustment

14

- Recall: TCP is ACK clocked
 - ▣ Congestion = delay = long wait between ACKs
 - ▣ No congestion = low delay = ACKs arrive quickly
- Basic algorithm
 - ▣ Upon receipt of ACK: increase *cwnd*
 - Data was delivered, perhaps we can send faster
 - *cwnd* growth is proportional to RTT
 - ▣ On loss: decrease *cwnd*
 - Data is being lost, there must be congestion
- Question: increase/decrease functions to use?

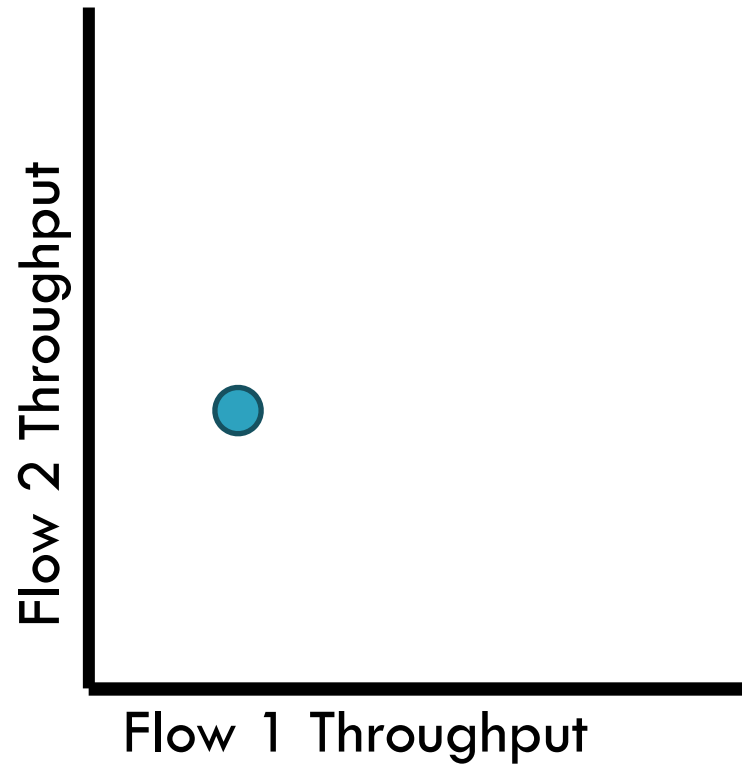
Utilization and Fairness

15



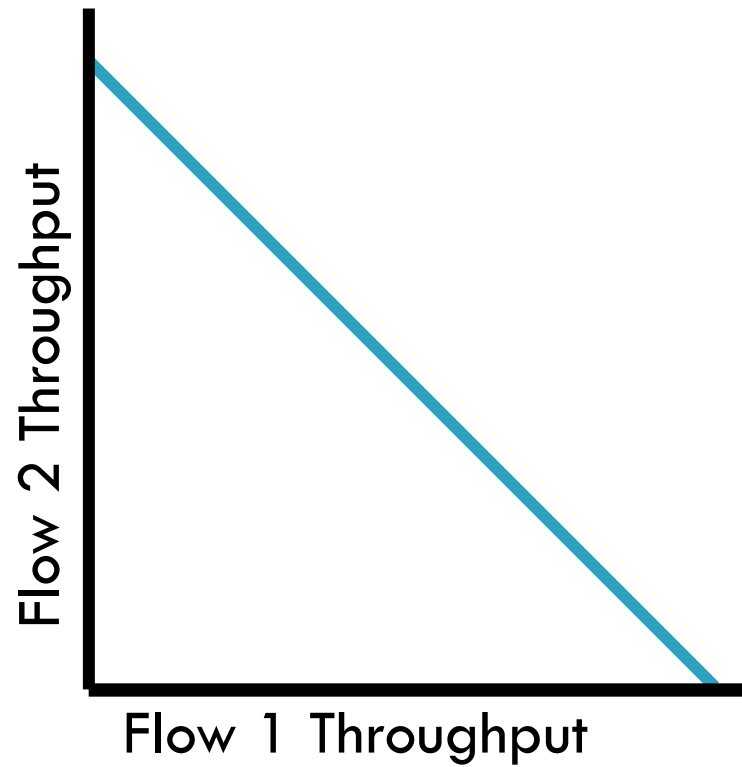
Utilization and Fairness

15



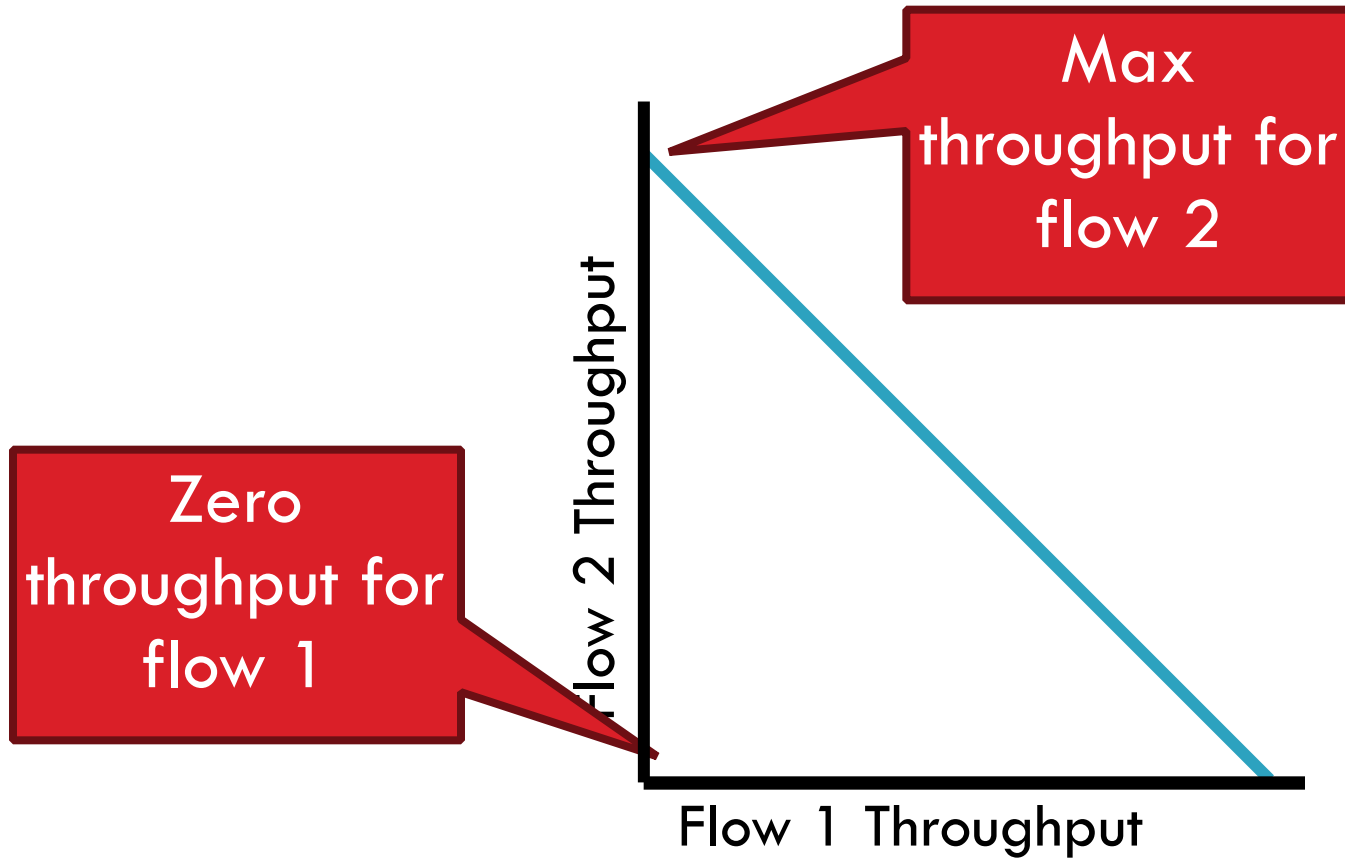
Utilization and Fairness

15



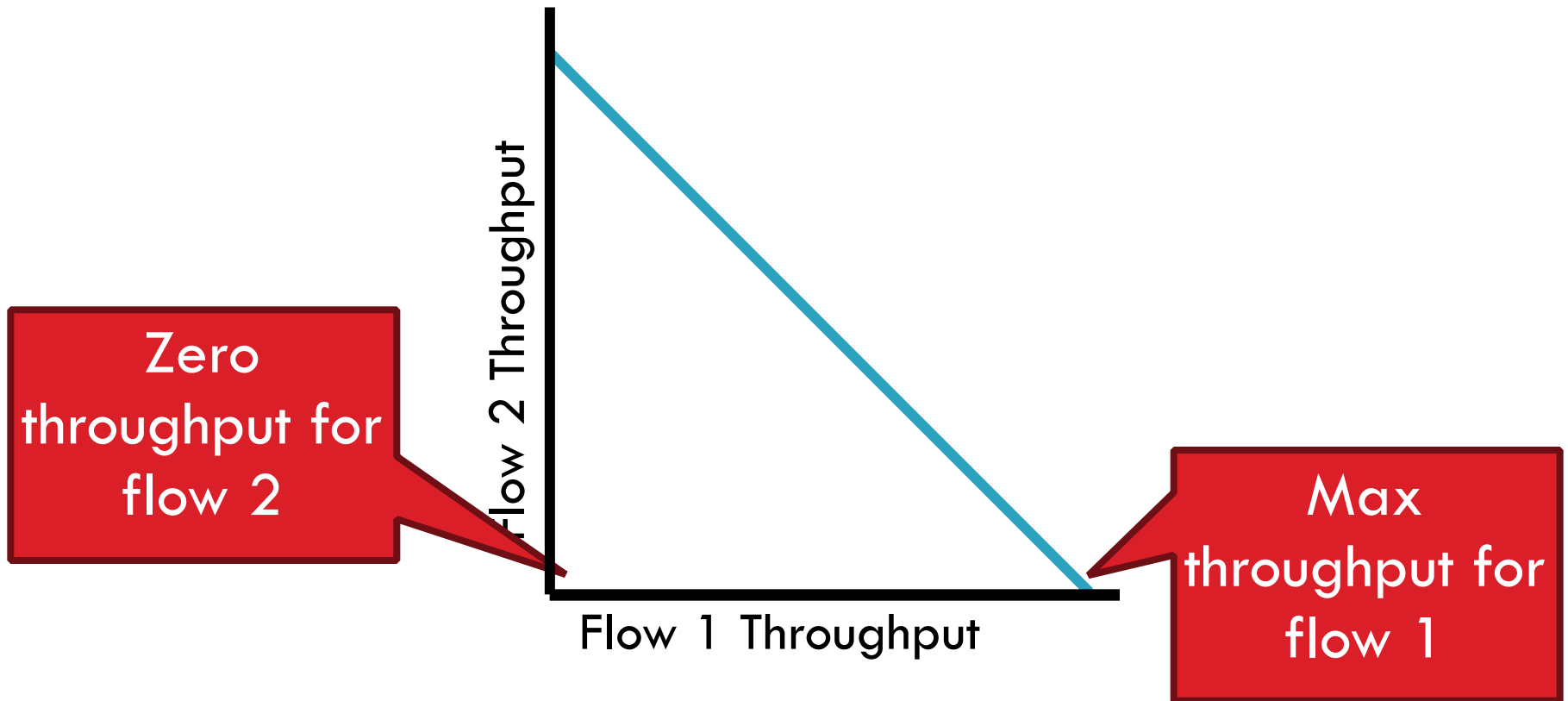
Utilization and Fairness

15



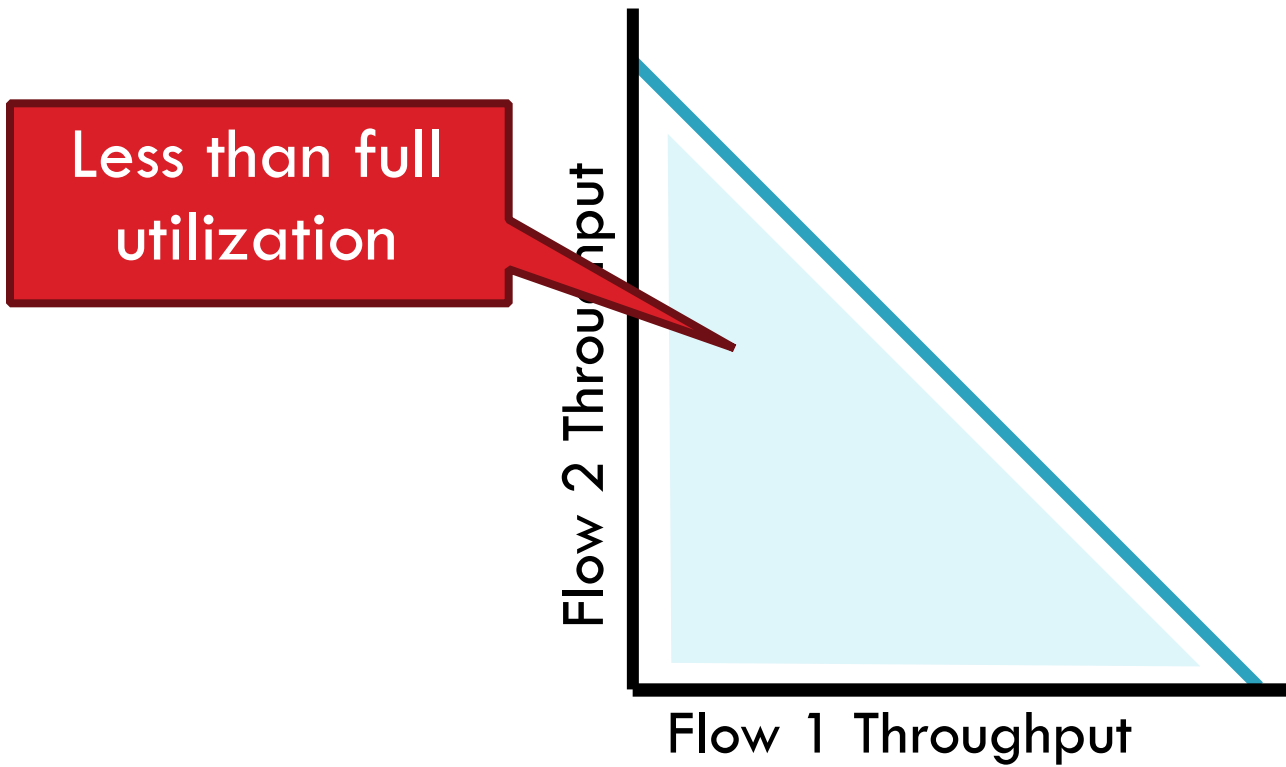
Utilization and Fairness

15



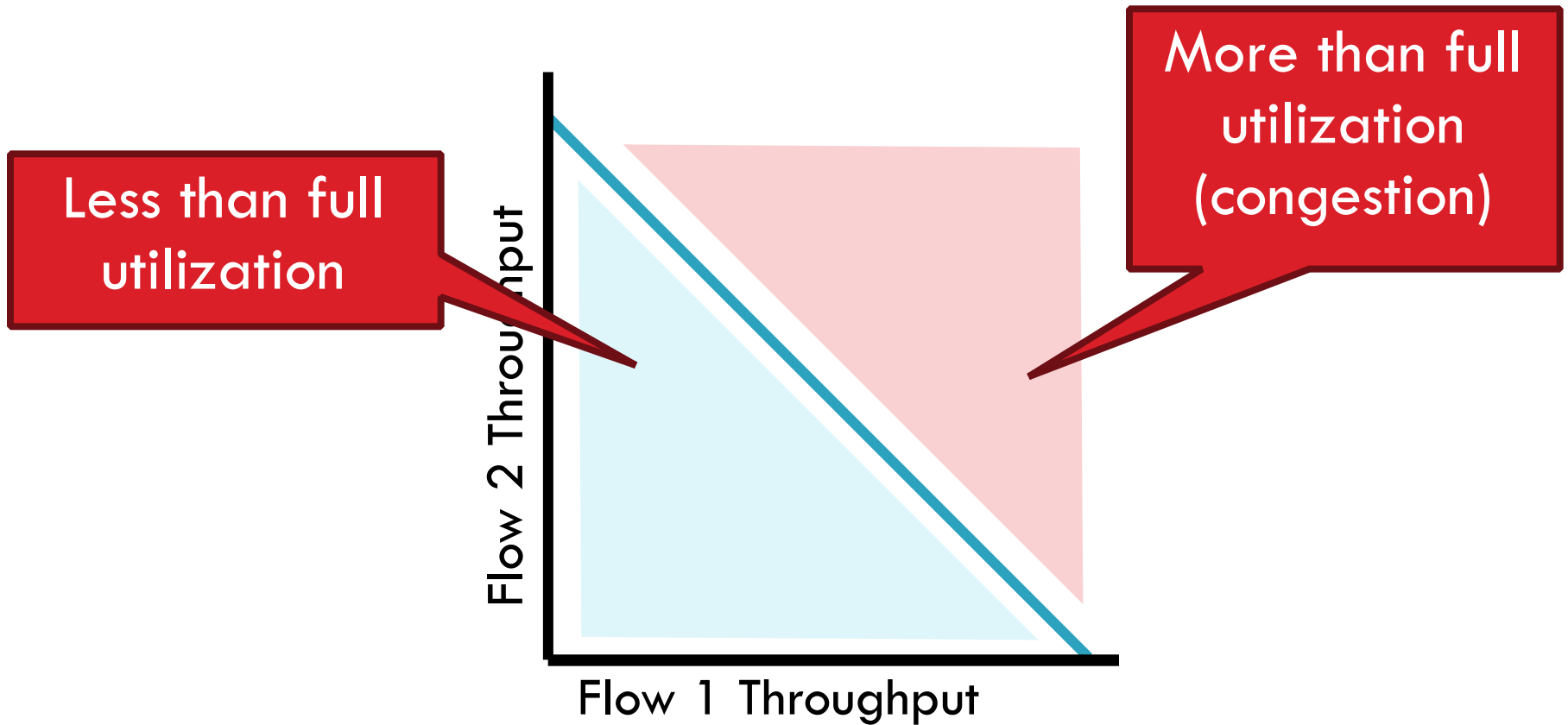
Utilization and Fairness

15



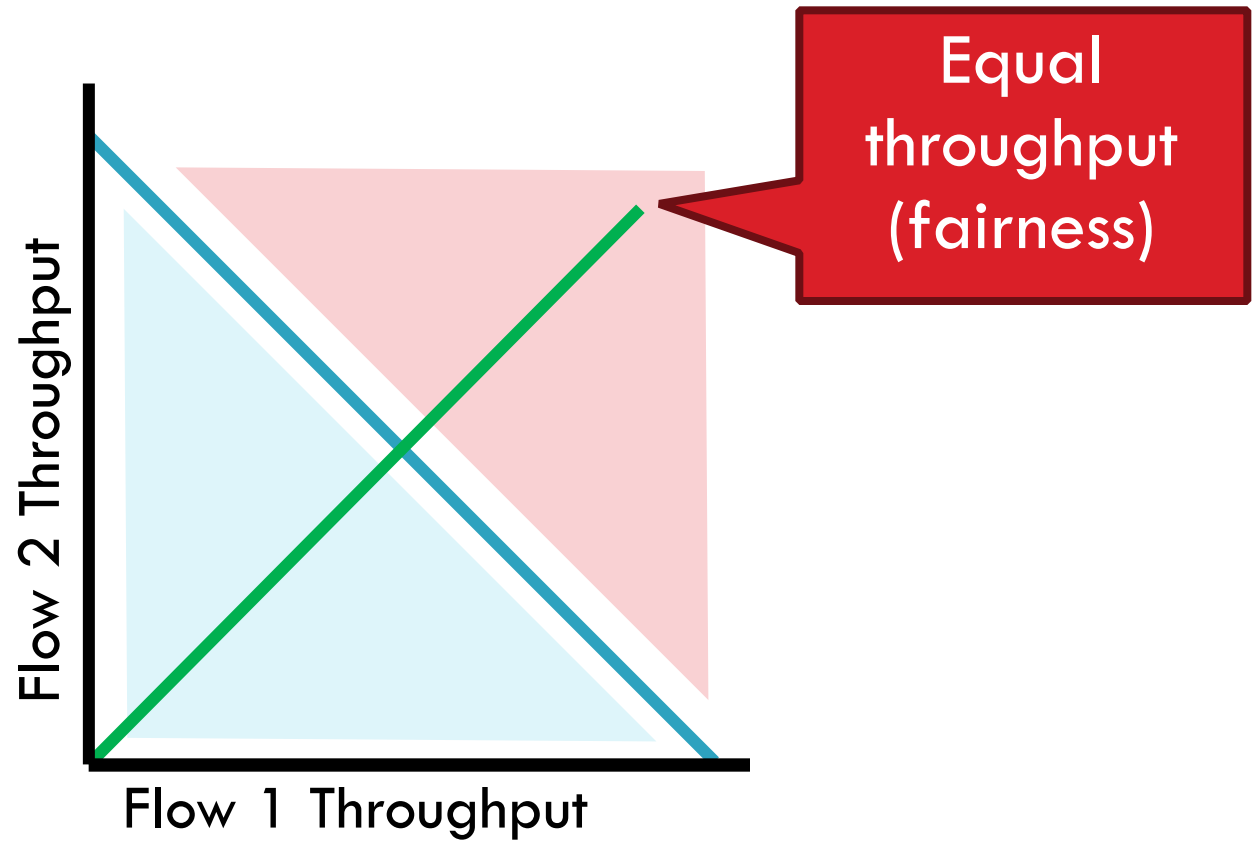
Utilization and Fairness

15



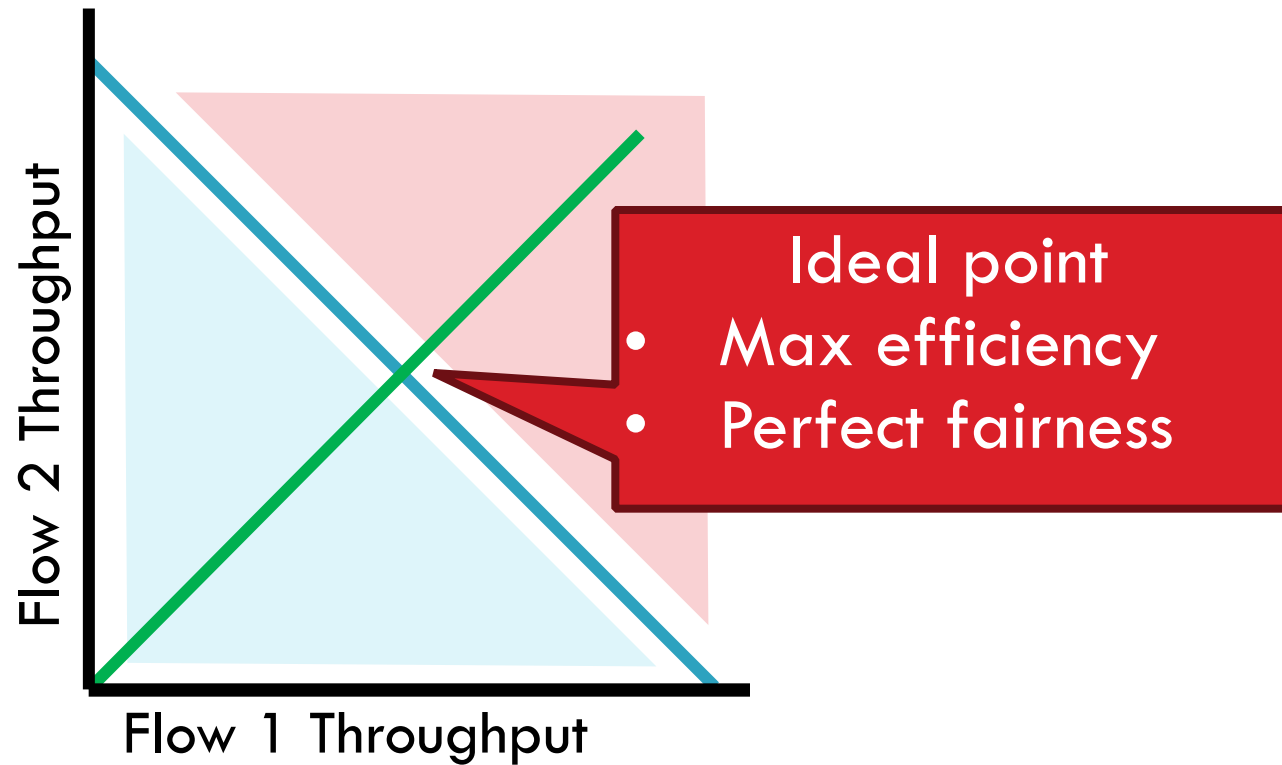
Utilization and Fairness

15



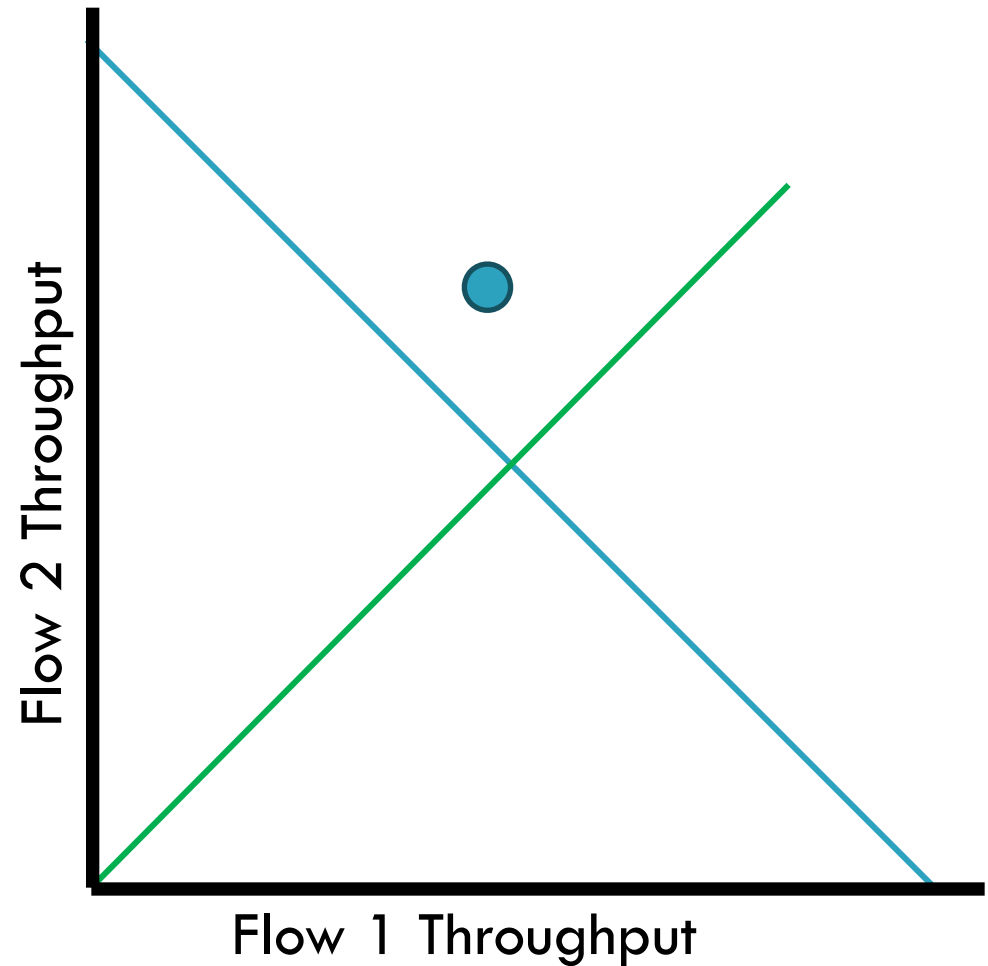
Utilization and Fairness

15



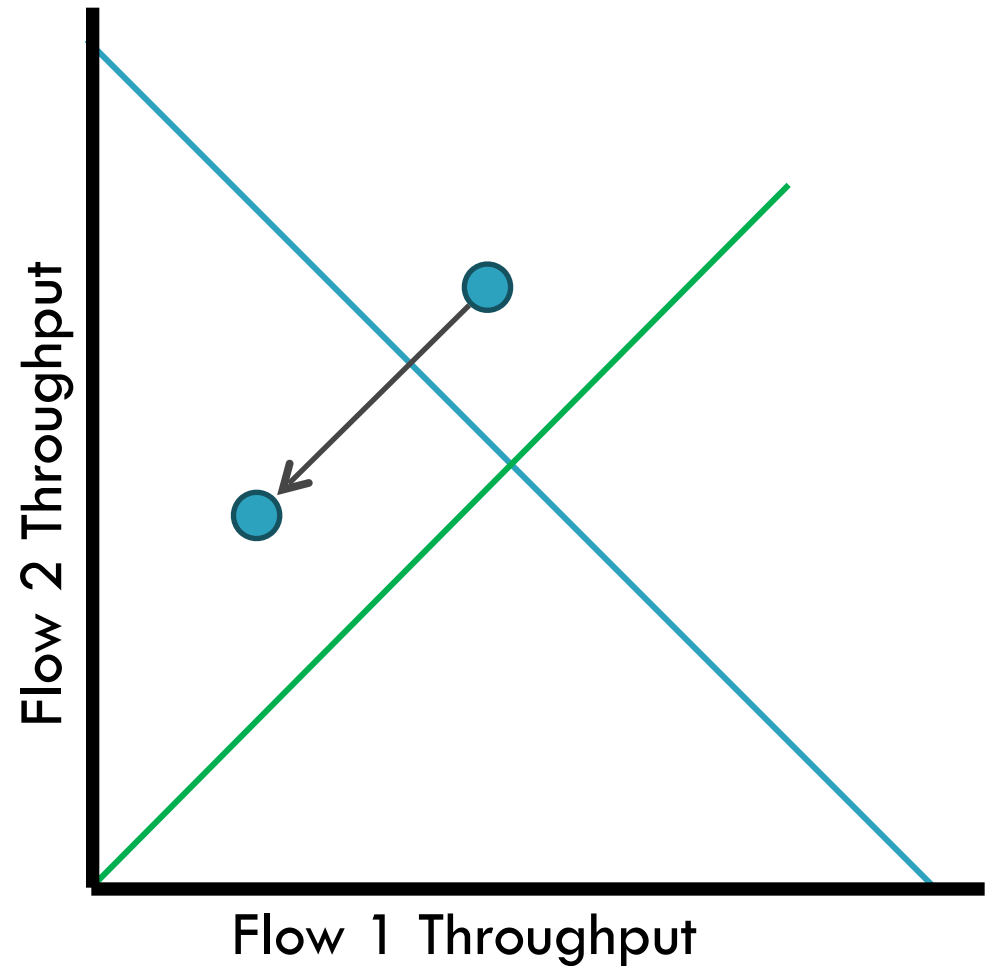
Multiplicative Increase, Additive Decrease

16



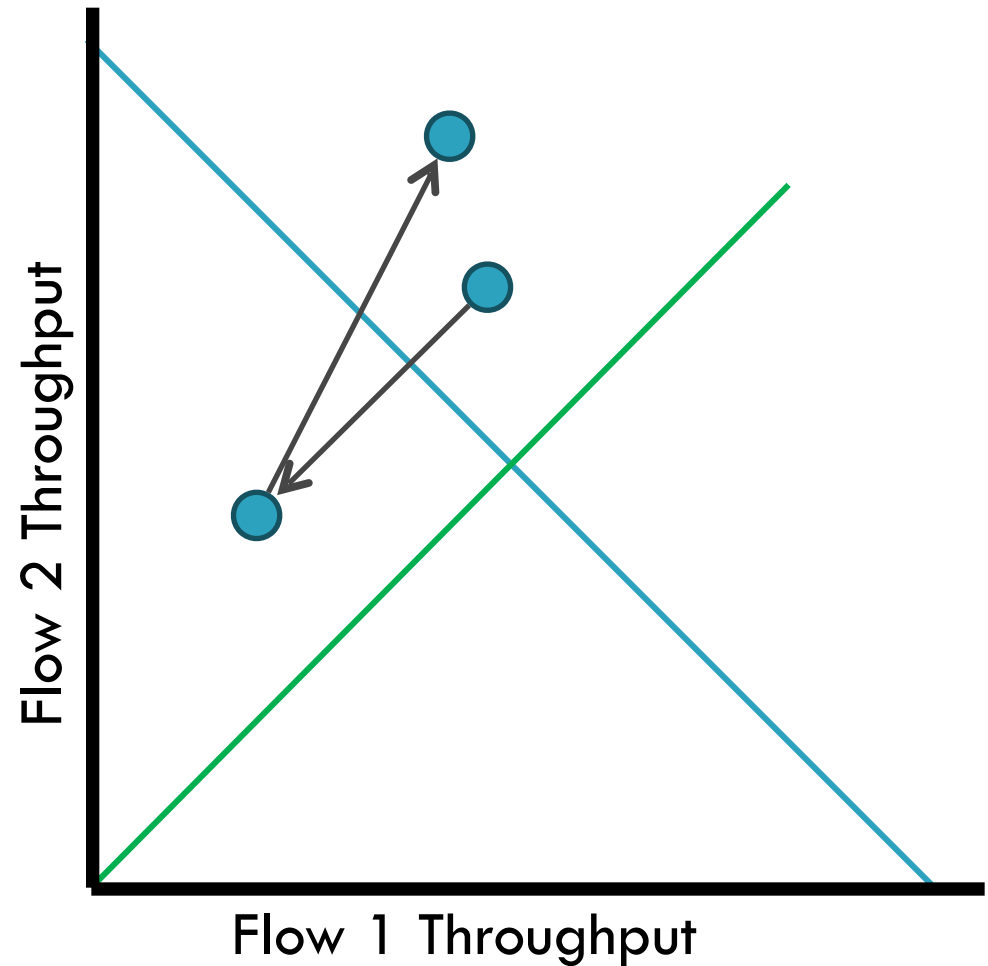
Multiplicative Increase, Additive Decrease

16



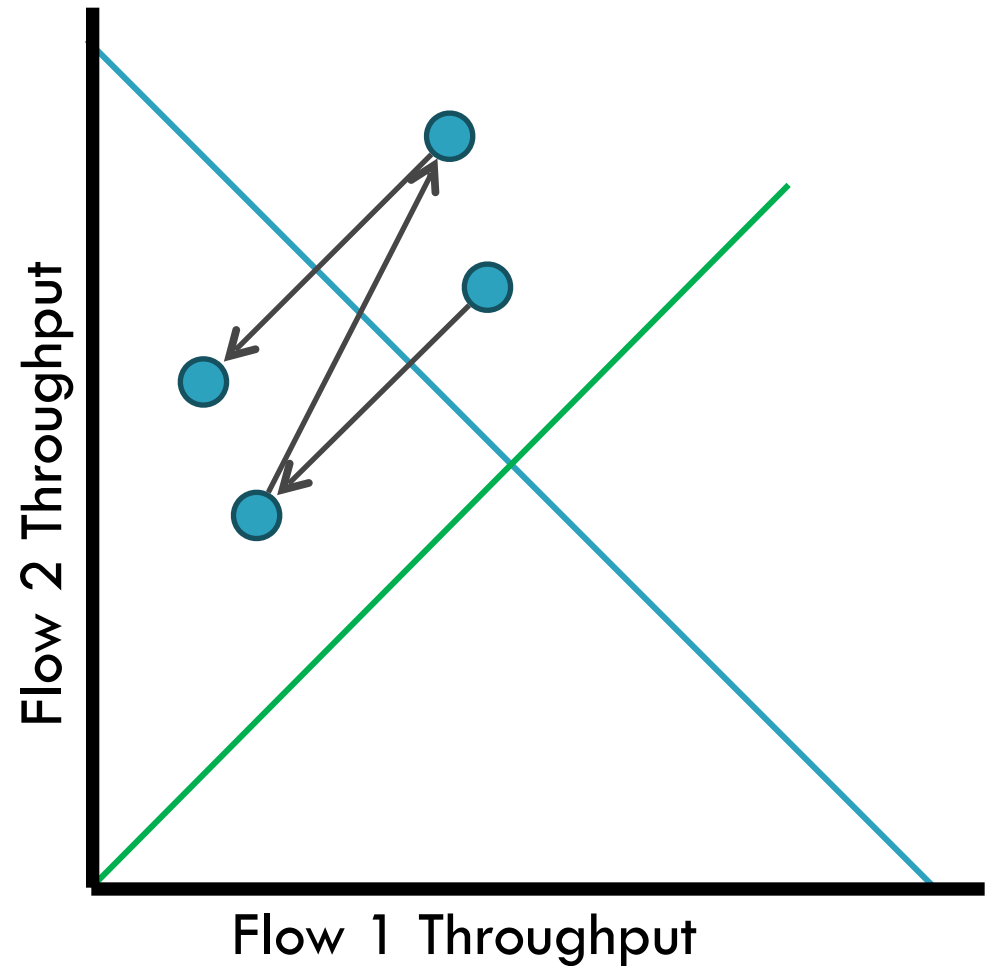
Multiplicative Increase, Additive Decrease

16



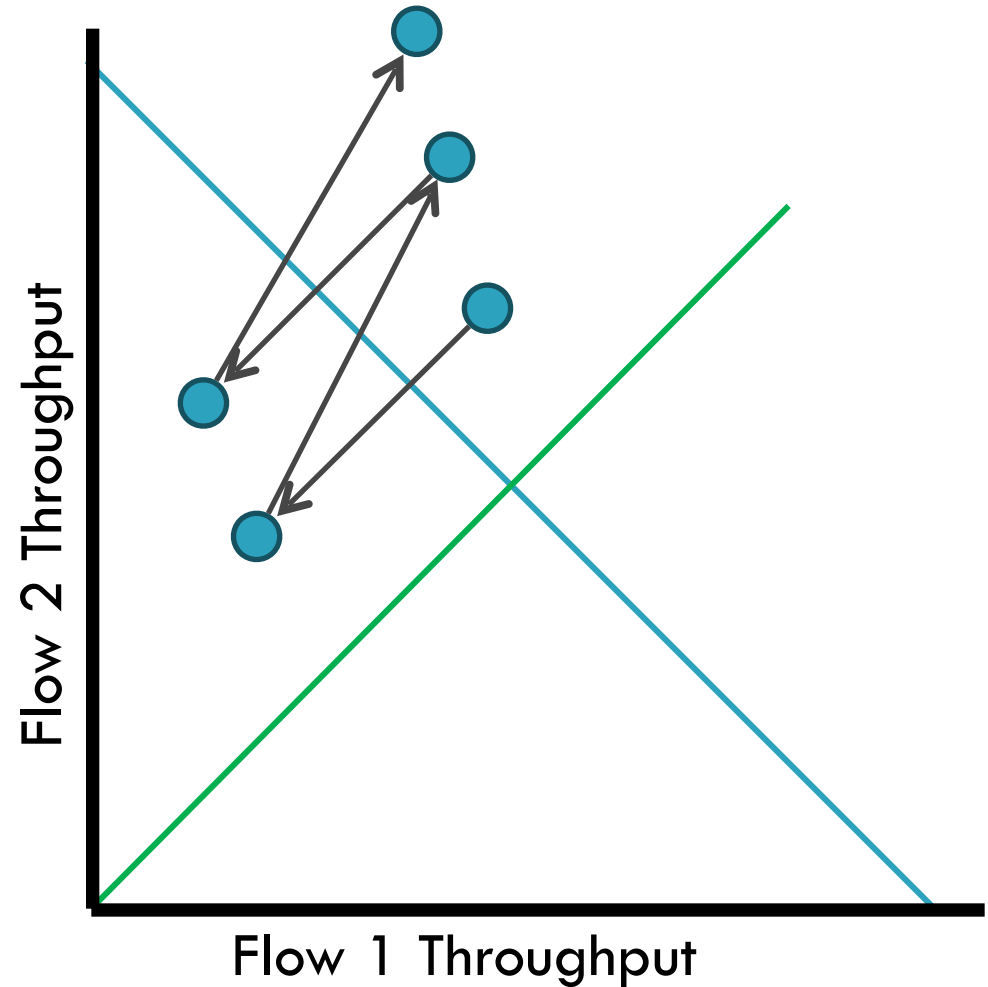
Multiplicative Increase, Additive Decrease

16



Multiplicative Increase, Additive Decrease

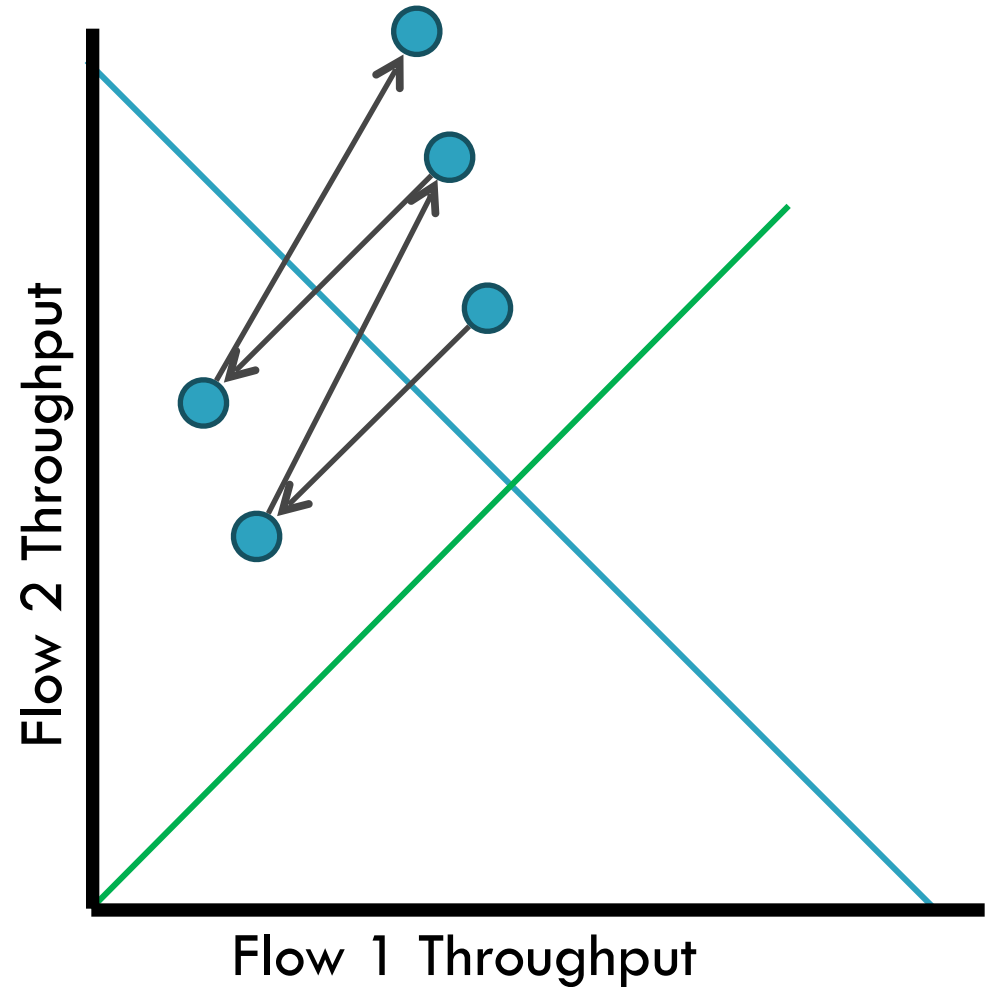
16



Multiplicative Increase, Additive Decrease

16

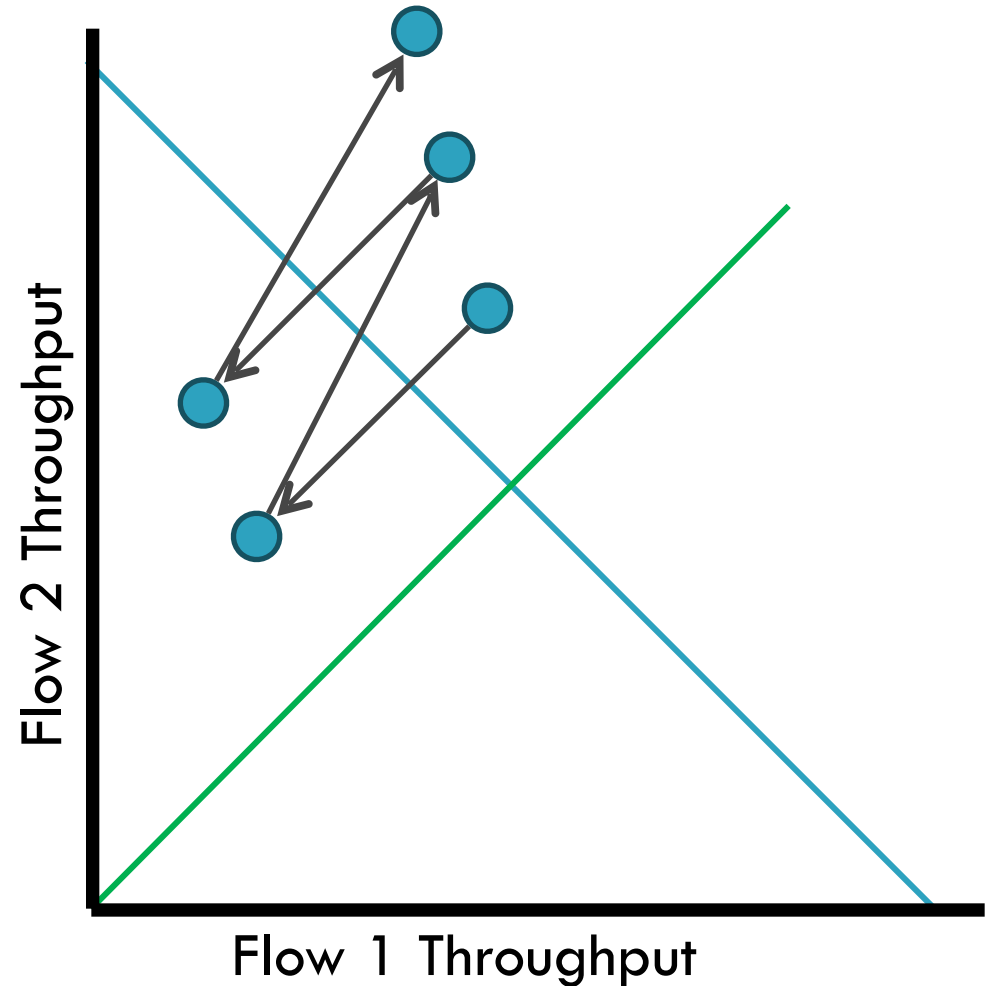
- Not stable!



Multiplicative Increase, Additive Decrease

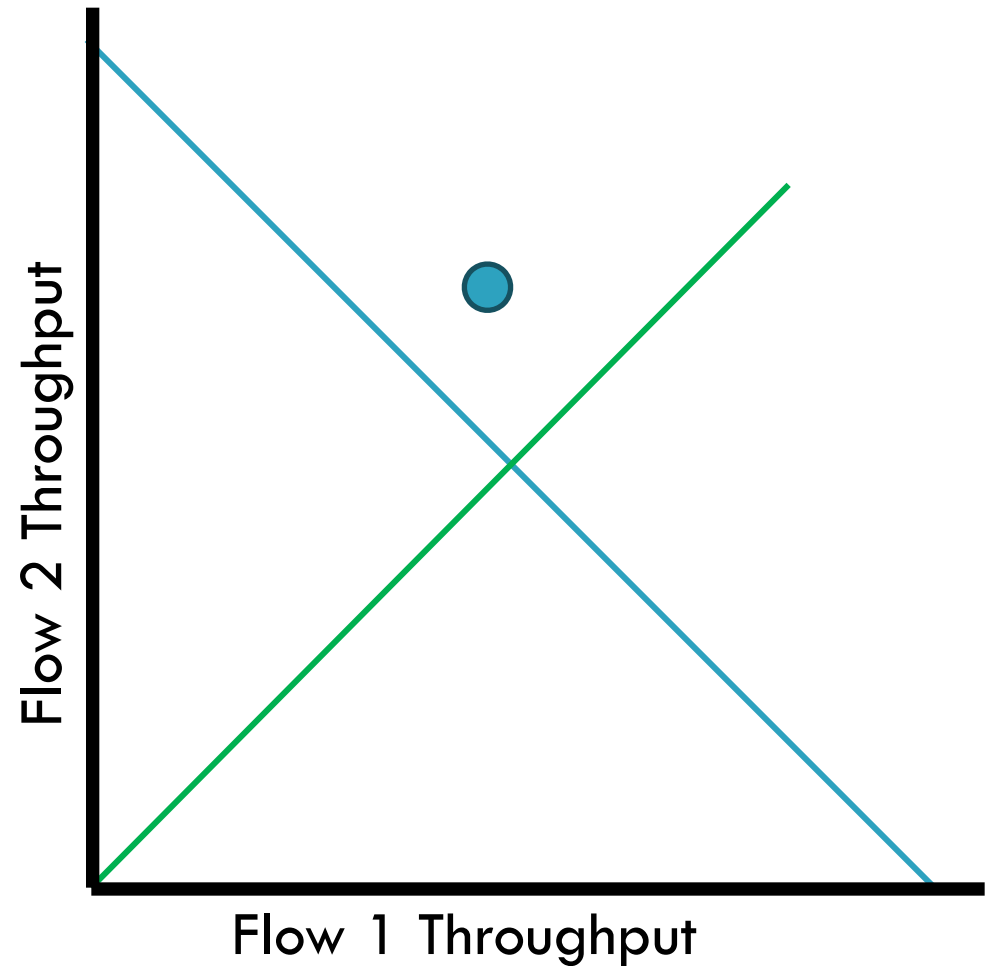
16

- Not stable!
- Veers away from fairness



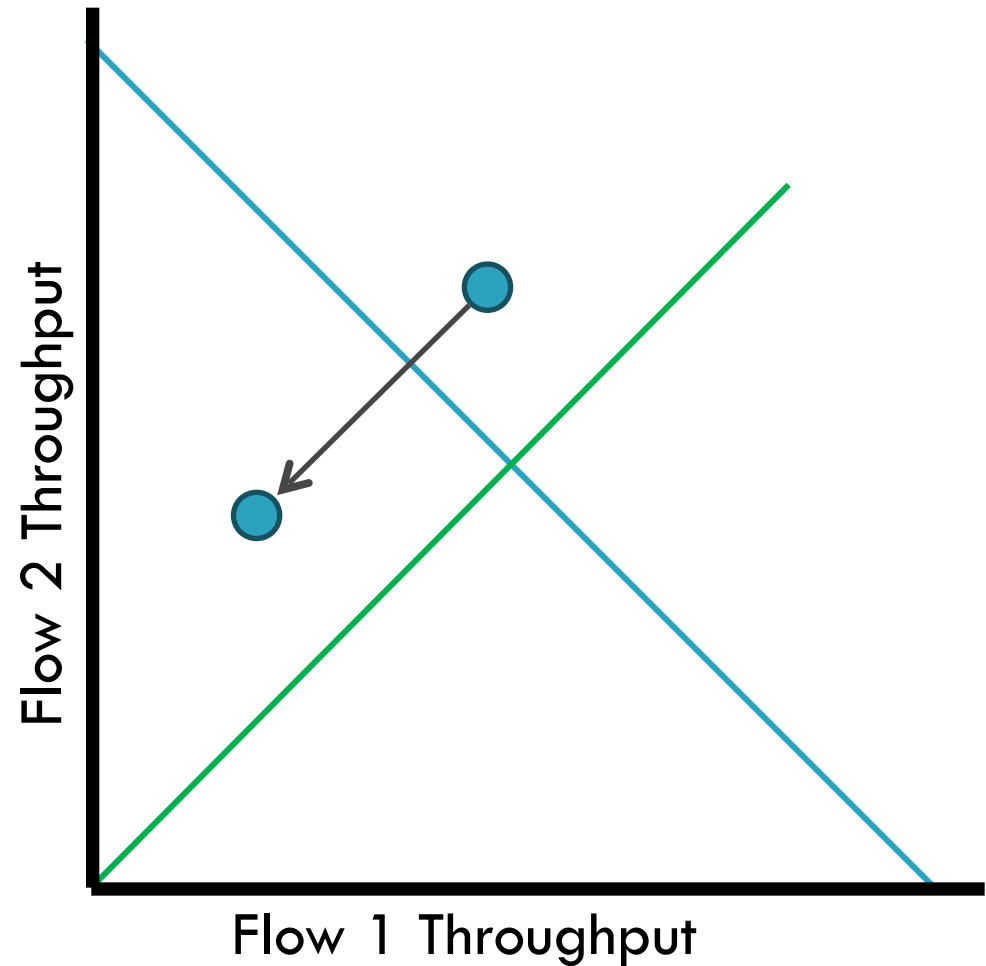
Additive Increase, Additive Decrease

17



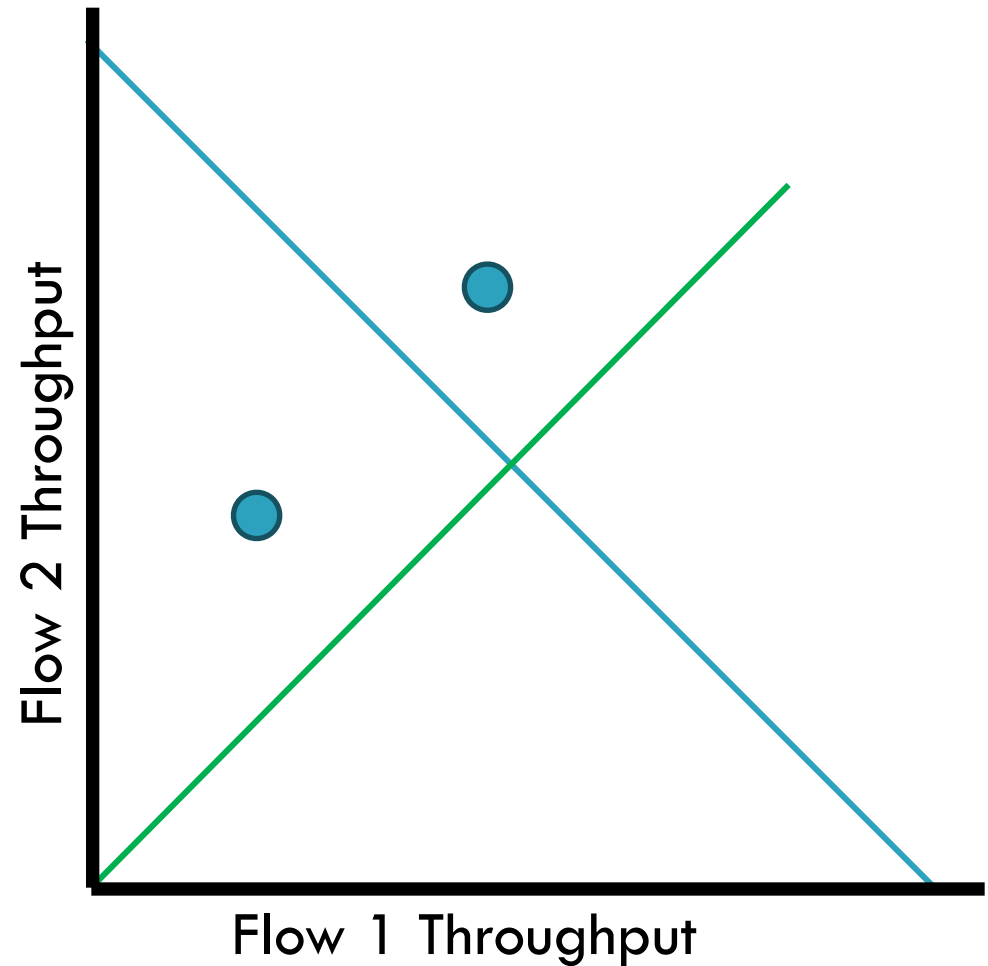
Additive Increase, Additive Decrease

17



Additive Increase, Additive Decrease

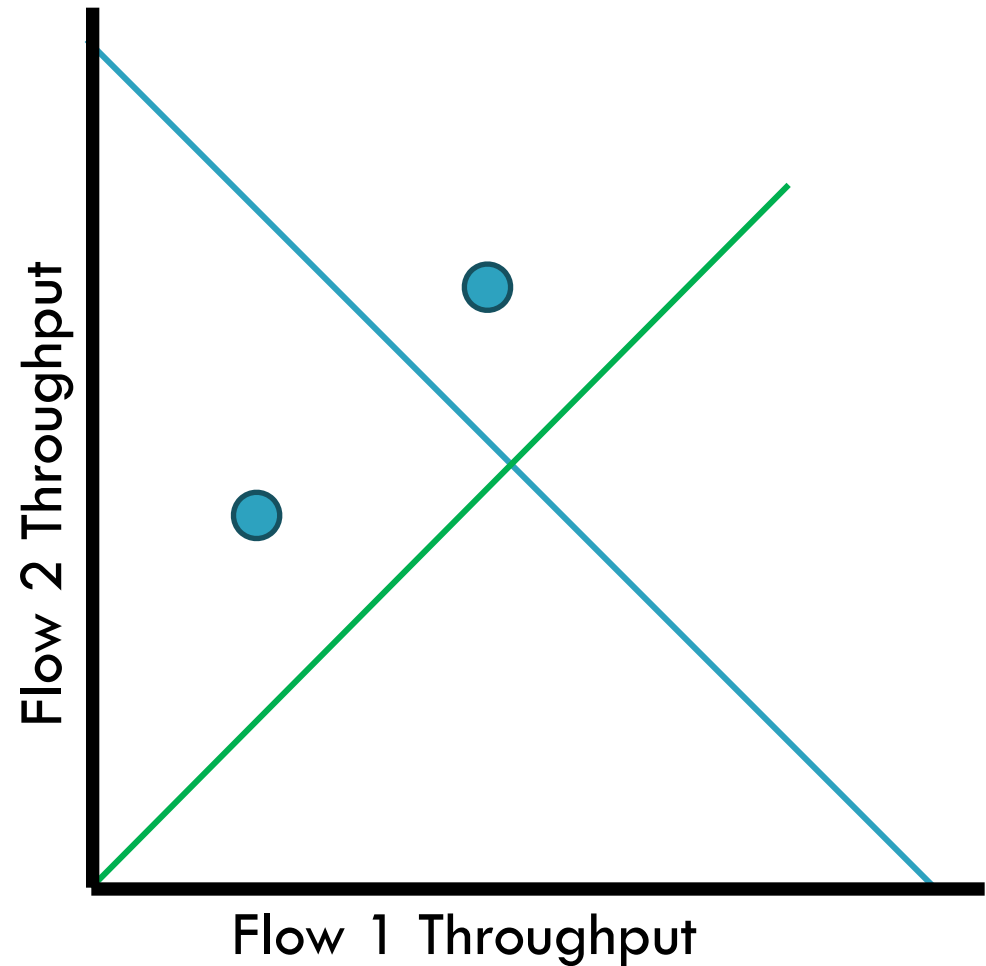
17



Additive Increase, Additive Decrease

17

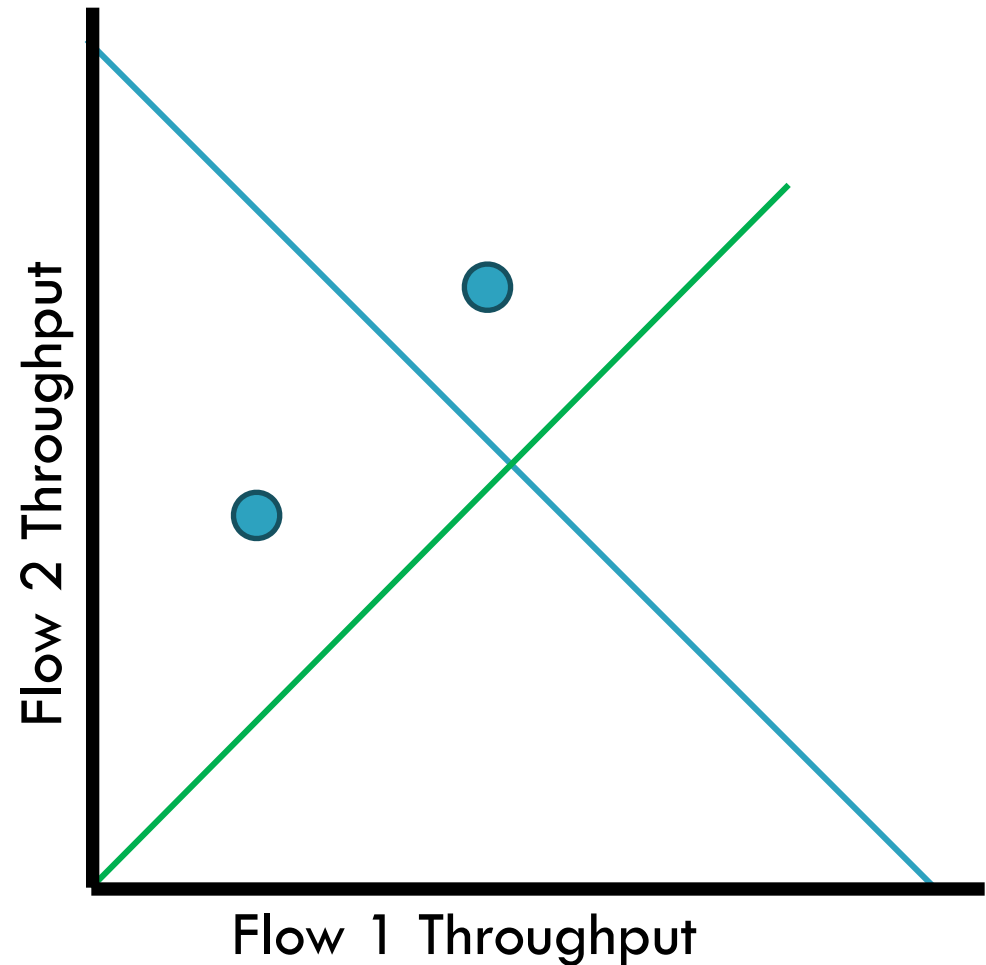
- Stable



Additive Increase, Additive Decrease

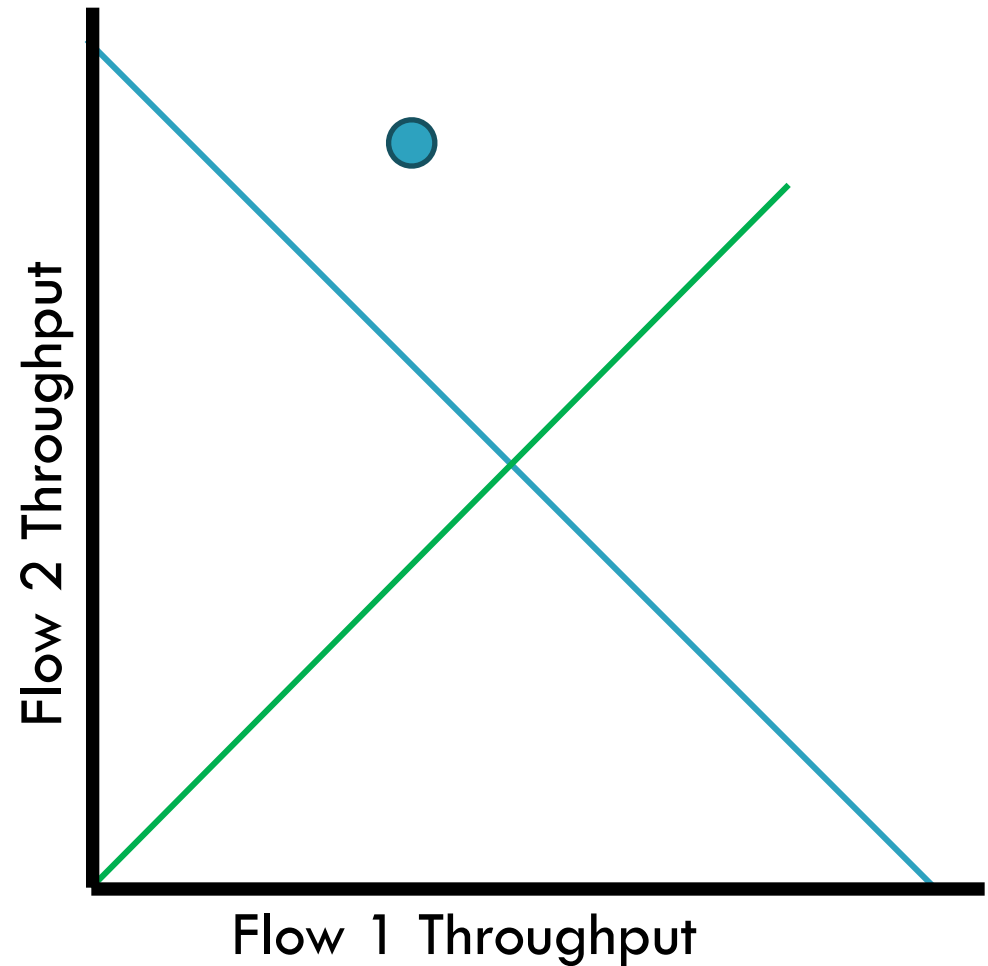
17

- Stable
- But does not converge to fairness



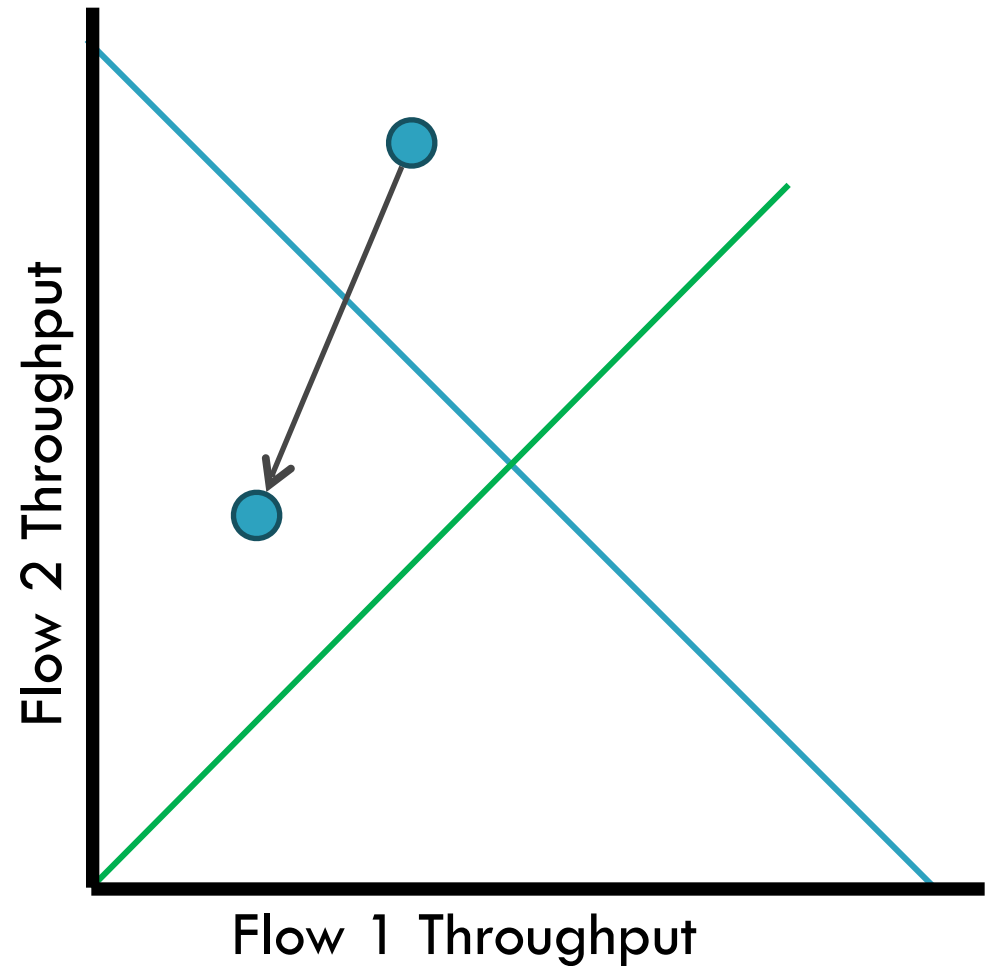
Multiplicative Increase, Multiplicative Decrease

18



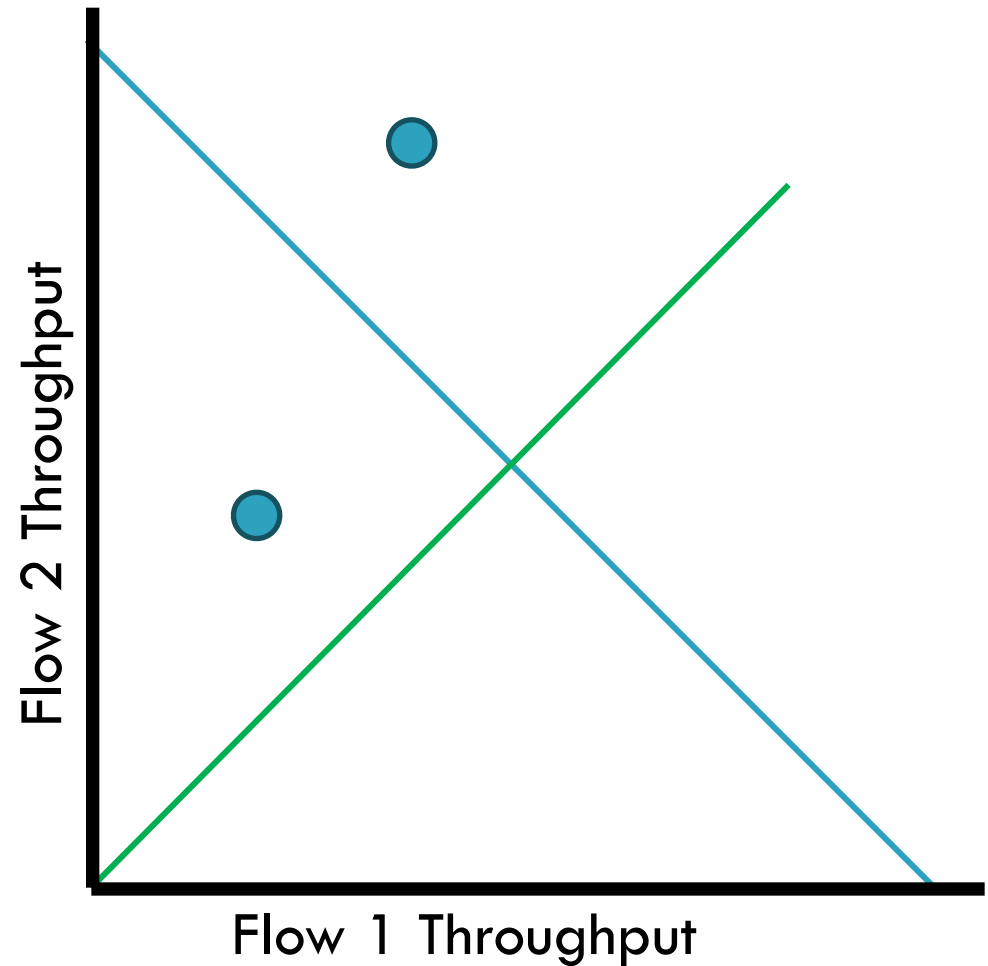
Multiplicative Increase, Multiplicative Decrease

18



Multiplicative Increase, Multiplicative Decrease

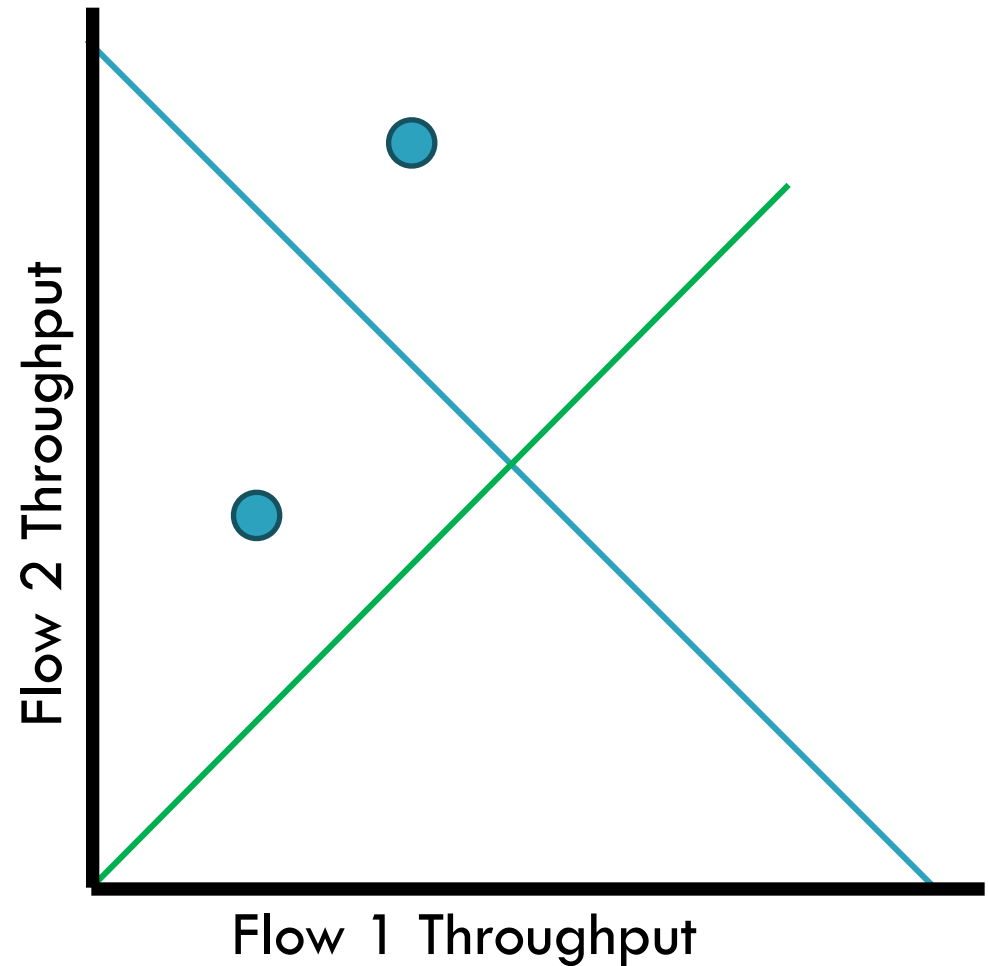
18



Multiplicative Increase, Multiplicative Decrease

18

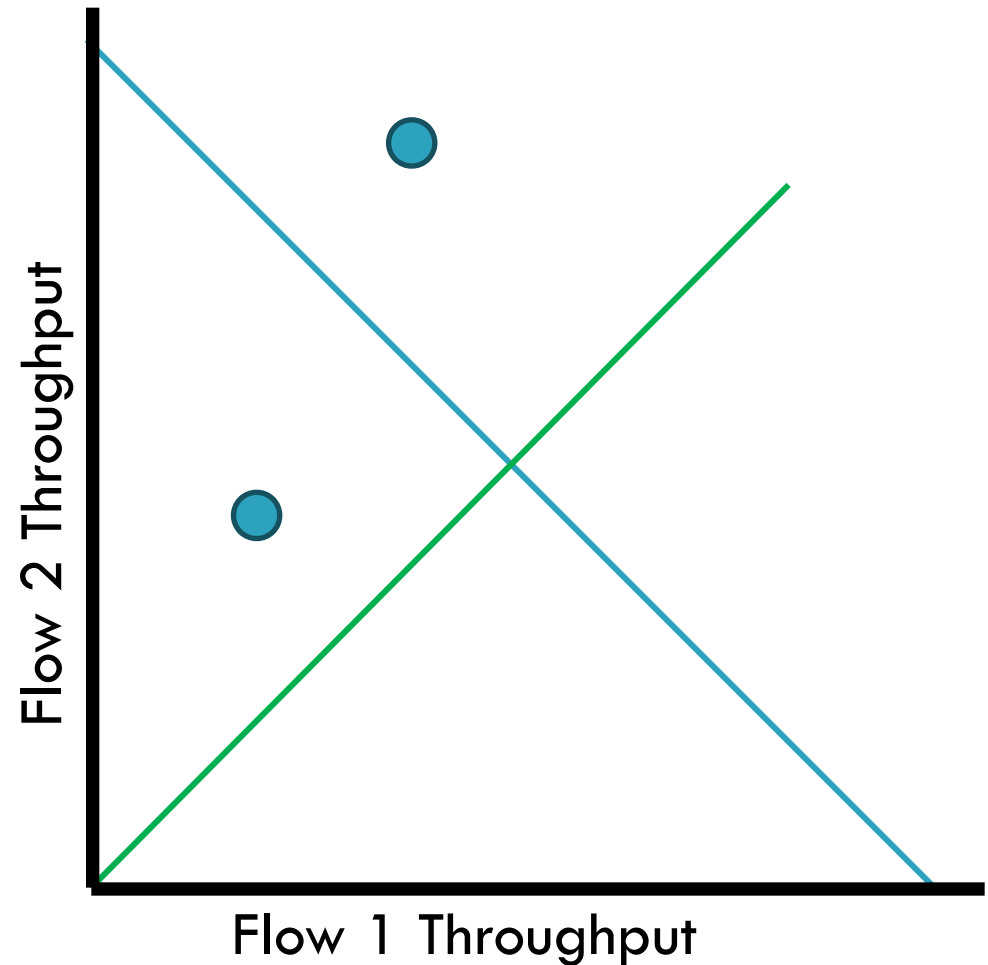
- Stable



Multiplicative Increase, Multiplicative Decrease

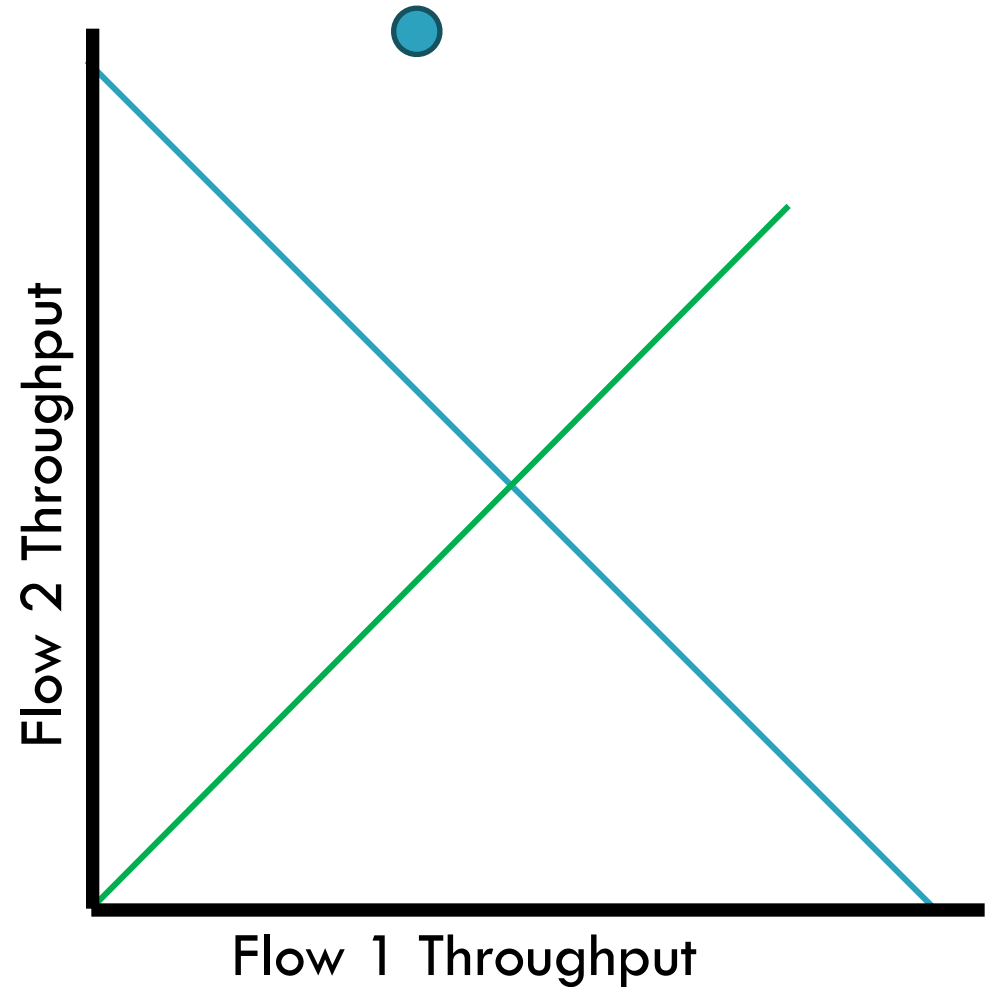
18

- Stable
- But does not converge to fairness



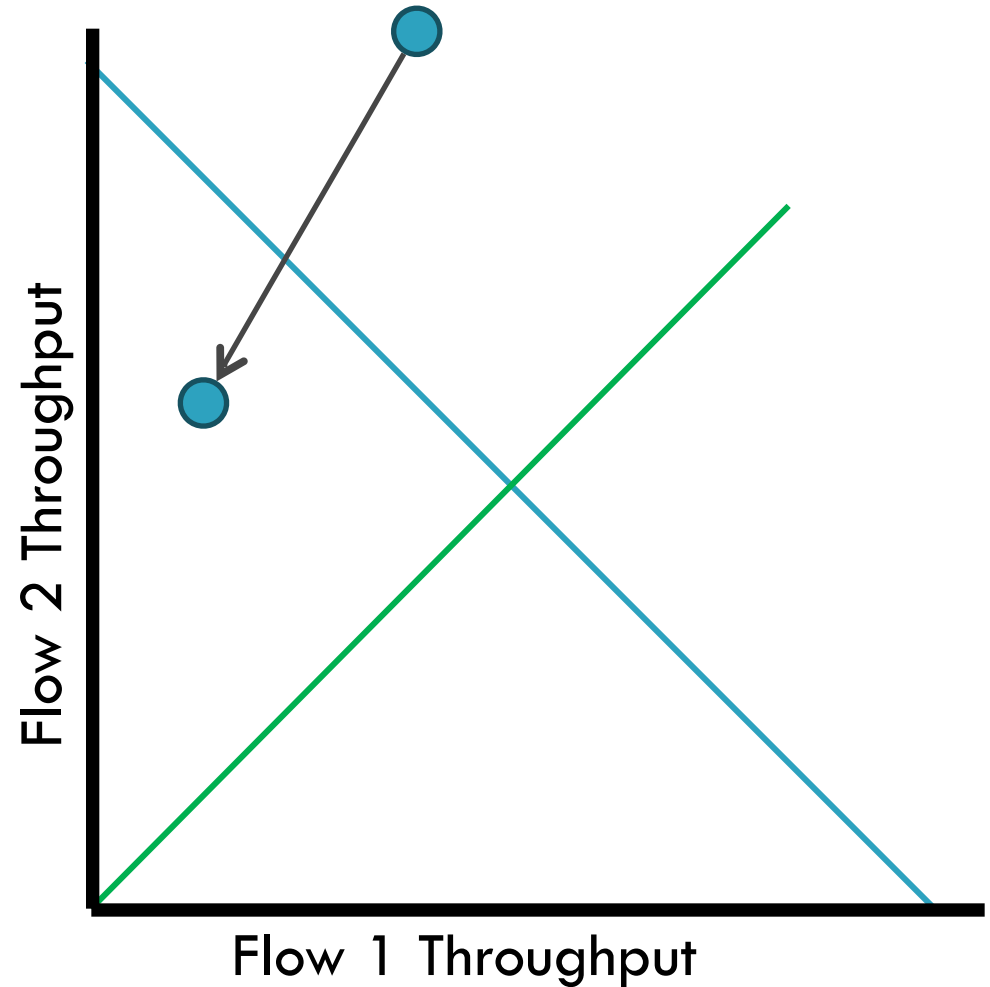
Additive Increase, Multiplicative Decrease

19



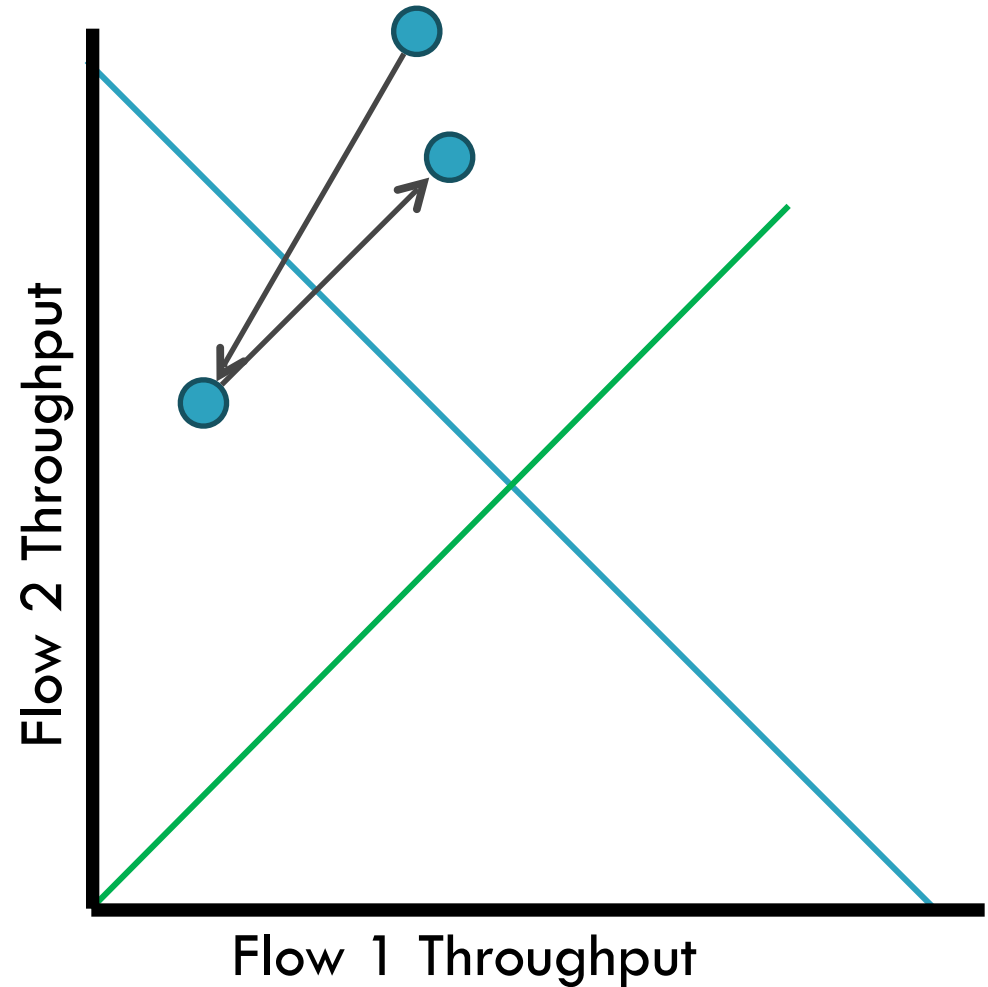
Additive Increase, Multiplicative Decrease

19



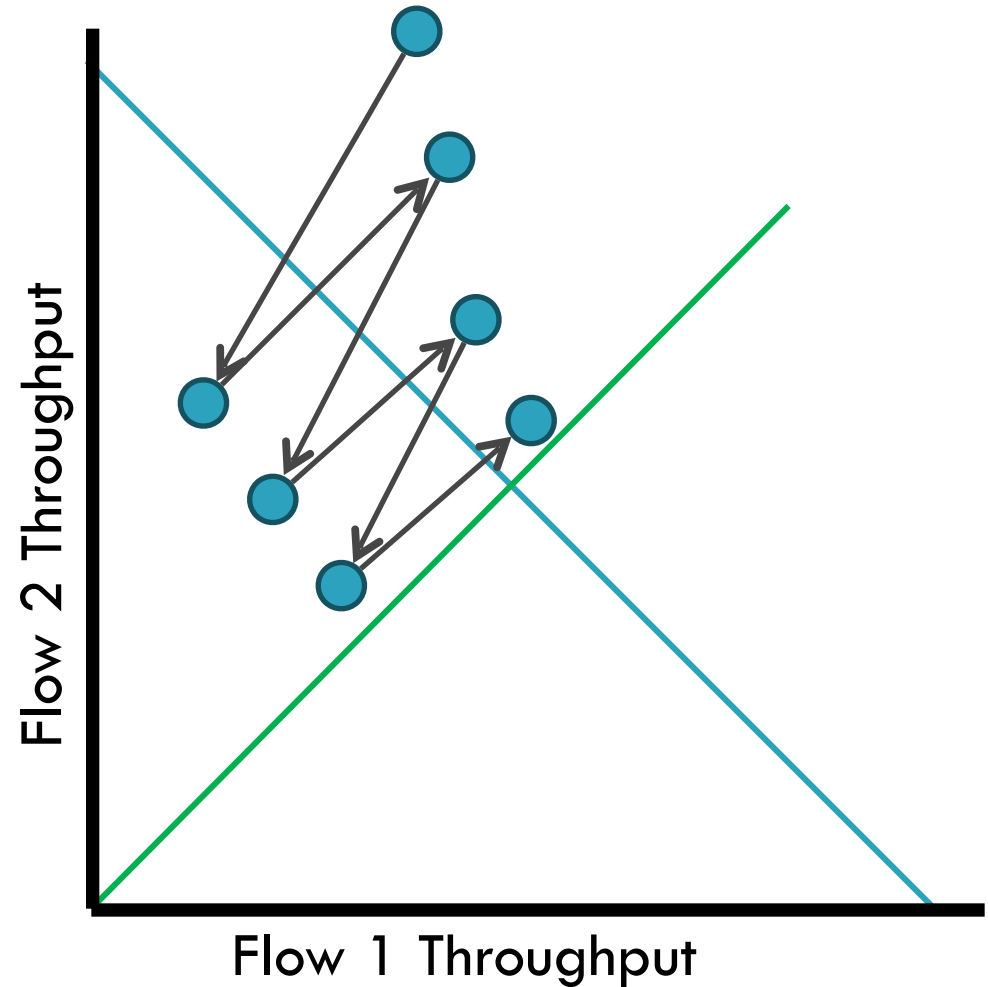
Additive Increase, Multiplicative Decrease

19



Additive Increase, Multiplicative Decrease

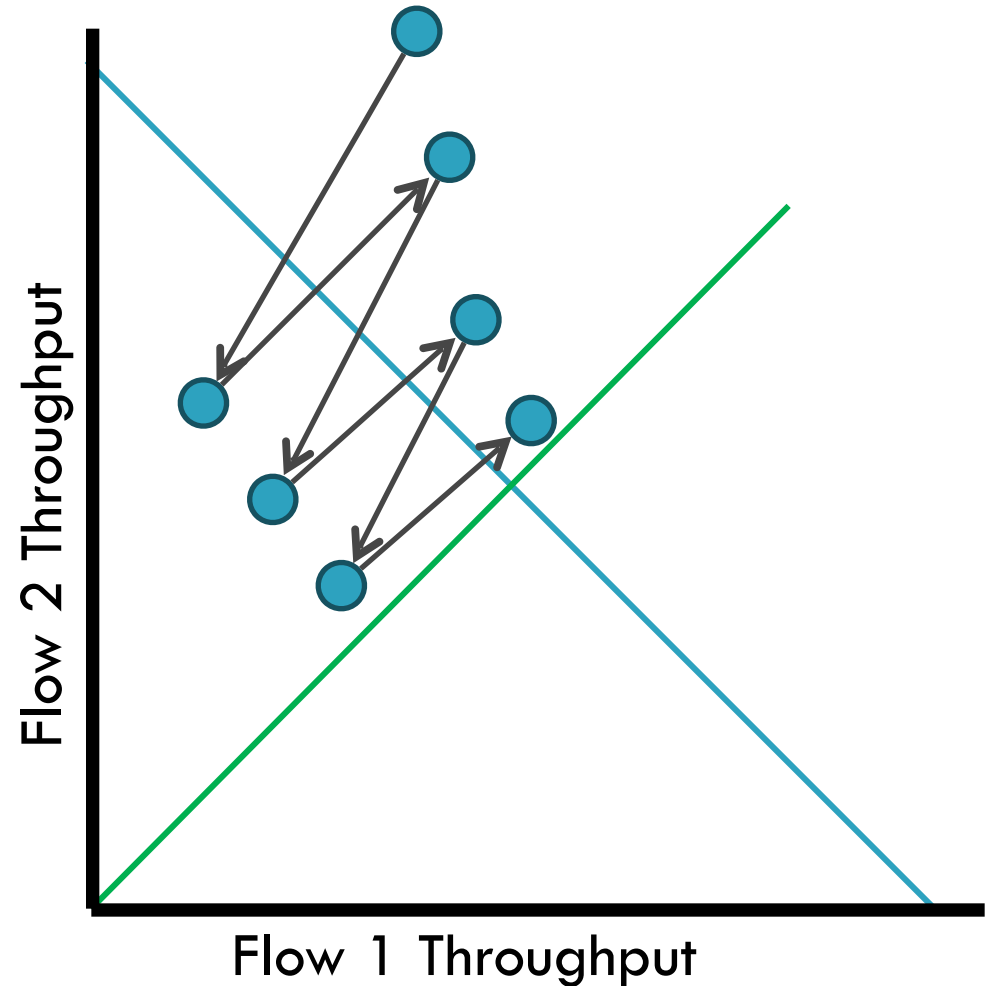
19



Additive Increase, Multiplicative Decrease

19

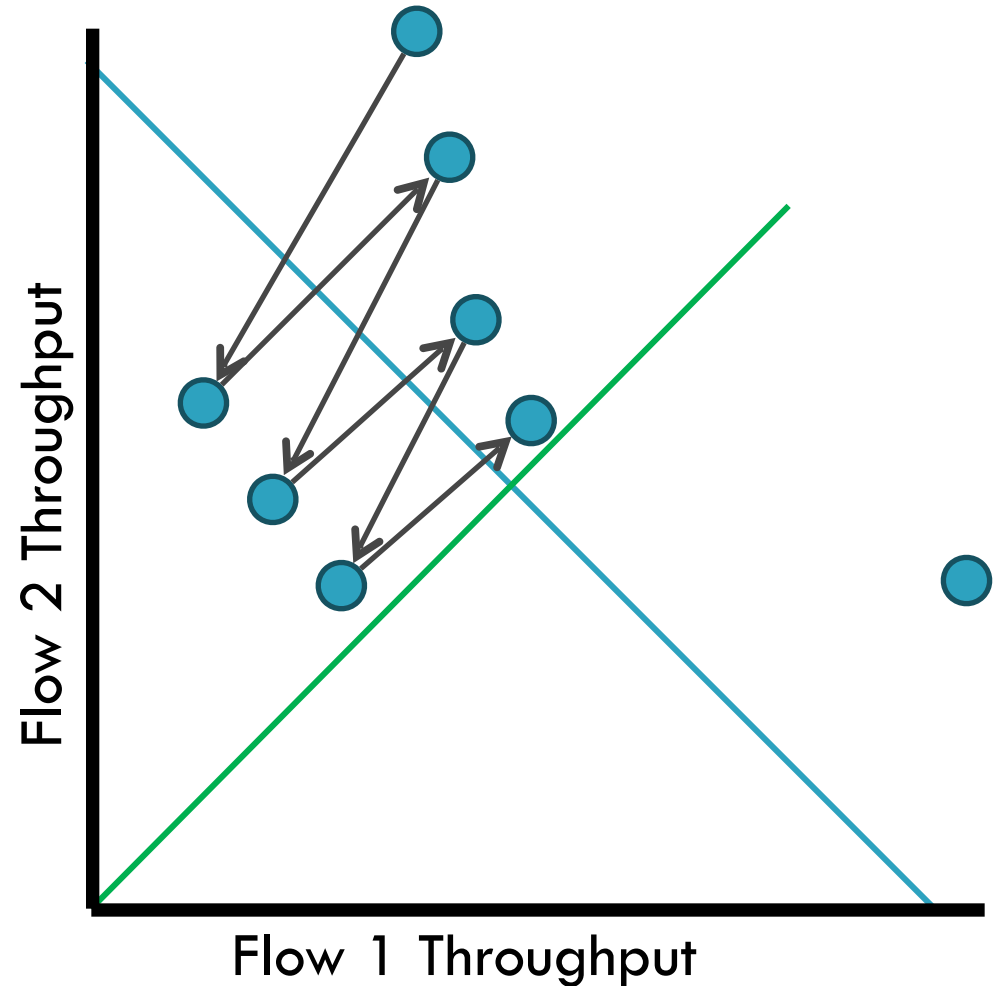
- Converges to stable and fair cycle



Additive Increase, Multiplicative Decrease

19

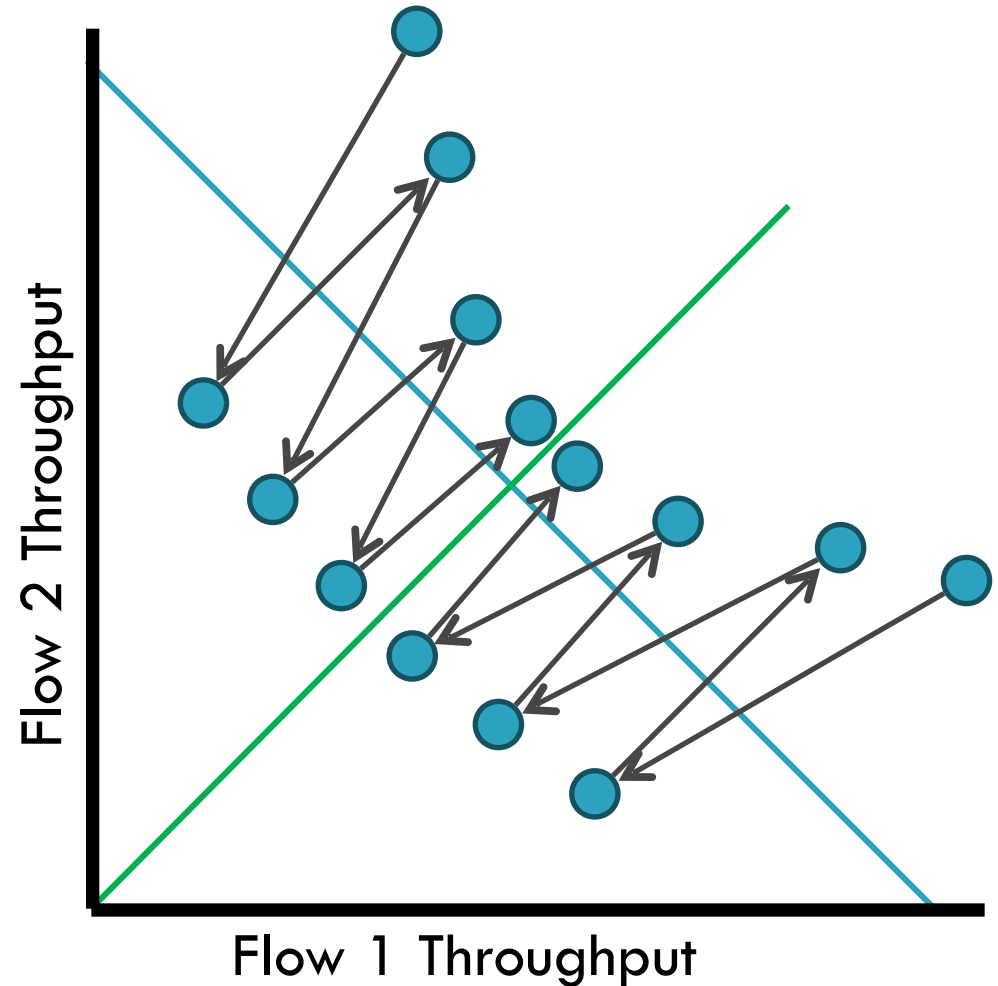
- Converges to stable and fair cycle



Additive Increase, Multiplicative Decrease

19

- Converges to stable and fair cycle
- Symmetric around $y=x$



Implementing Congestion Control

20

- Maintains three variables:
 - *cwnd*: congestion window
 - *adv_wnd*: receiver advertised window
 - *ssthresh*: threshold size (used to update *cwnd*)
- For sending, use: $wnd = \min(cwnd, adv_wnd)$

Implementing Congestion Control

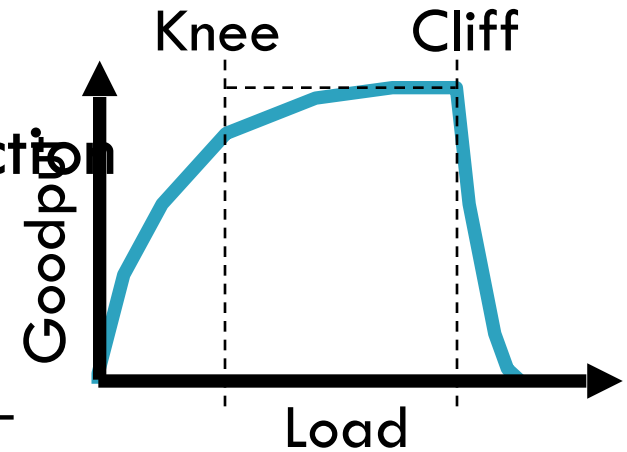
20

- Maintains three variables:
 - *cwnd*: congestion window
 - *adv_wnd*: receiver advertised window
 - *ssthresh*: threshold size (used to update *cwnd*)
- For sending, use: $wnd = \min(cwnd, adv_wnd)$
- Two phases of congestion control
 1. Slow start ($cwnd < ssthresh$)
 - Probe for bottleneck bandwidth
 2. Congestion avoidance ($cwnd \geq ssthresh$)
 - AIMD

Slow Start

21

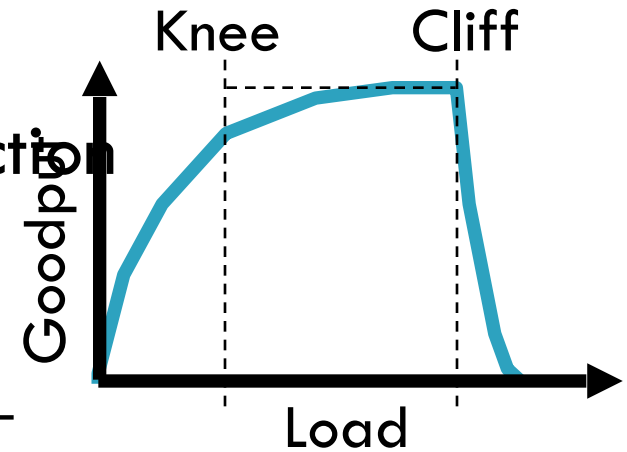
- Goal: reach knee quickly
- Upon starting (or restarting) a connection
 - $cwnd = 1$
 - $ssthresh = adv_wnd$
 - Each time a segment is ACKed, $cwnd++$



Slow Start

21

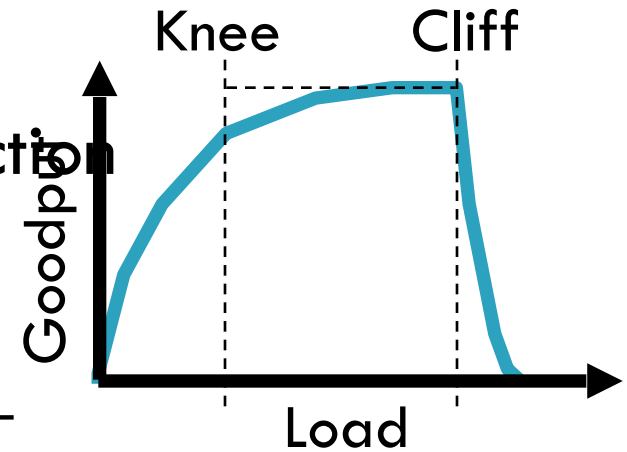
- Goal: reach knee quickly
- Upon starting (or restarting) a connection
 - $cwnd = 1$
 - $ssthresh = adv_wnd$
 - Each time a segment is ACKed, $cwnd++$
- Continues until...
 - $ssthresh$ is reached
 - Or a packet is lost



Slow Start

21

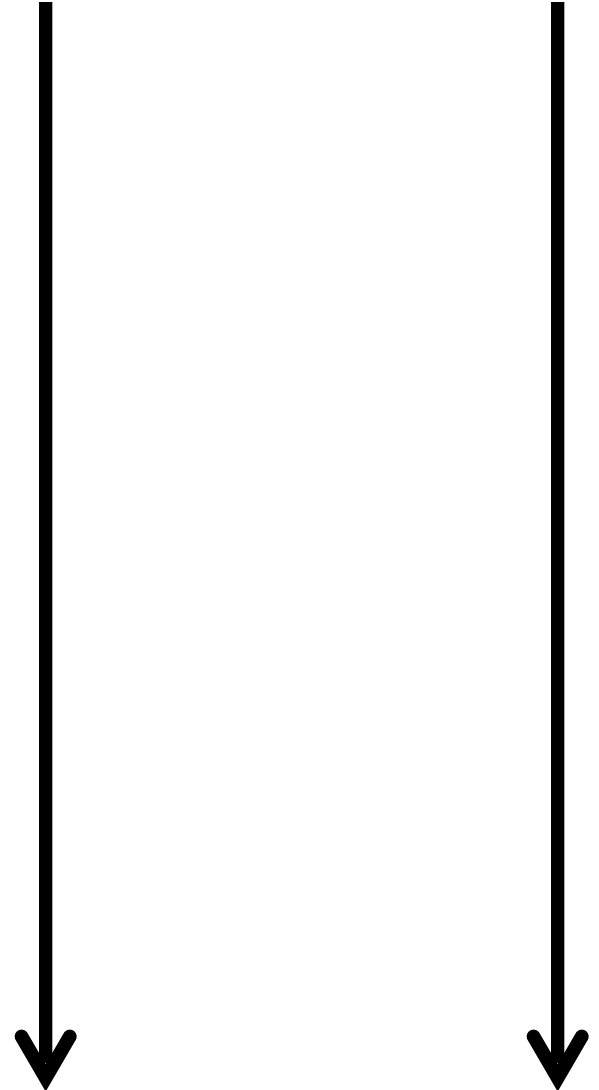
- Goal: reach knee quickly
- Upon starting (or restarting) a connection
 - $cwnd = 1$
 - $ssthresh = adv_wnd$
 - Each time a segment is ACKed, $cwnd++$
- Continues until...
 - $ssthresh$ is reached
 - Or a packet is lost
- Slow Start is not actually slow
 - $cwnd$ increases exponentially



Slow Start Example

22

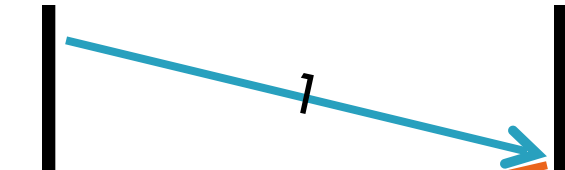
$cwnd = 1$



Slow Start Example

22

$cwnd = 1$



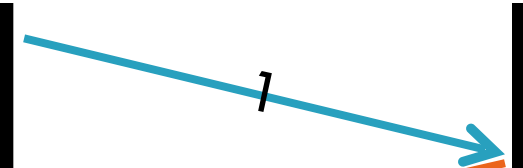
$cwnd = 2$



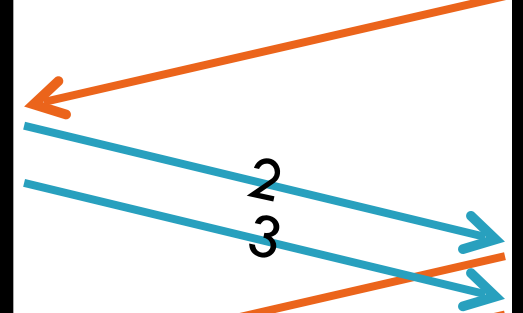
Slow Start Example

22

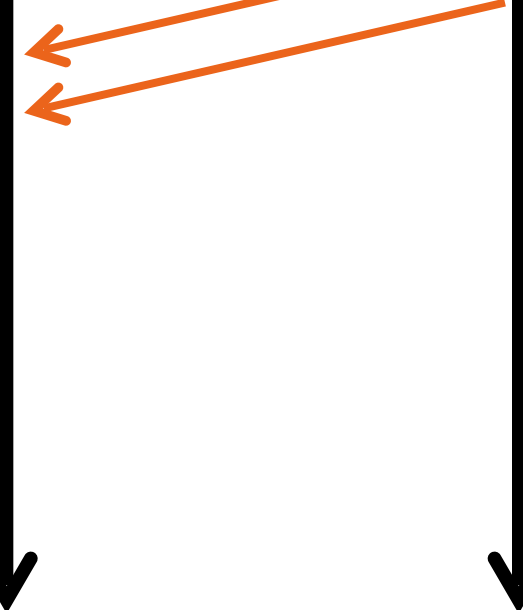
$cwnd = 1$



$cwnd = 2$

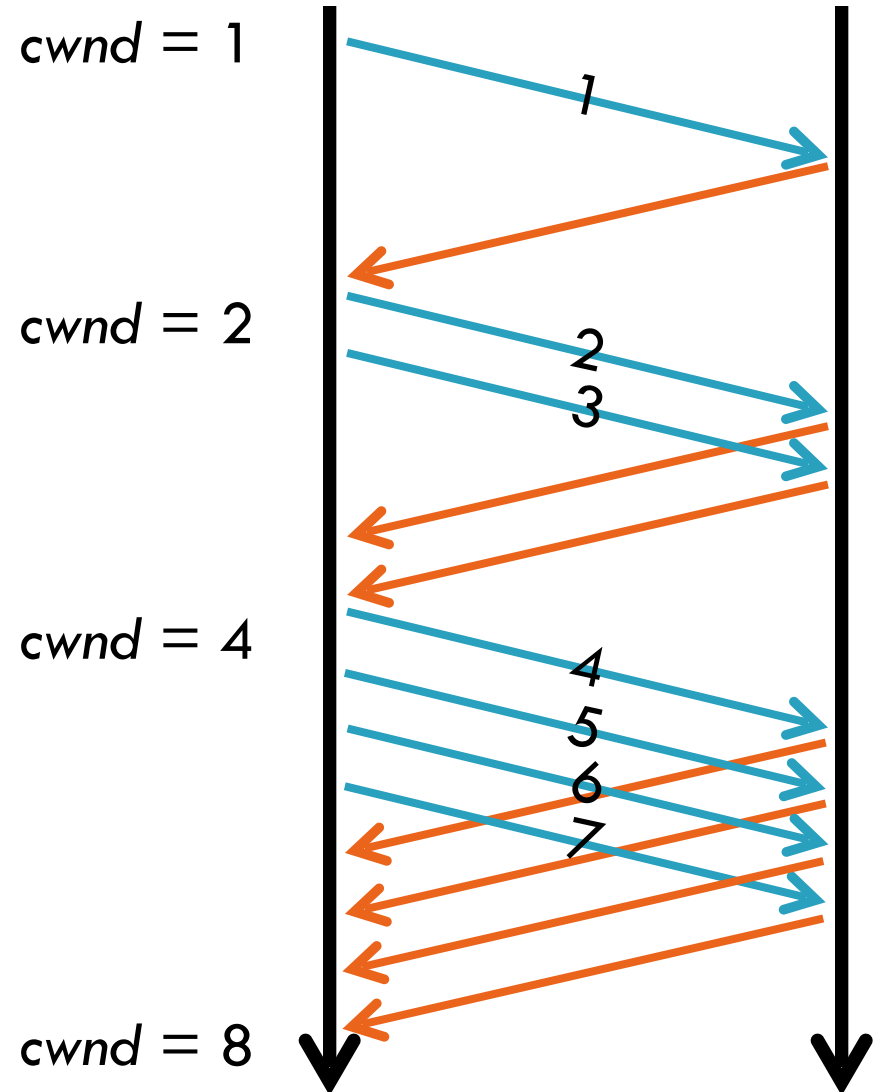


$cwnd = 4$



Slow Start Example

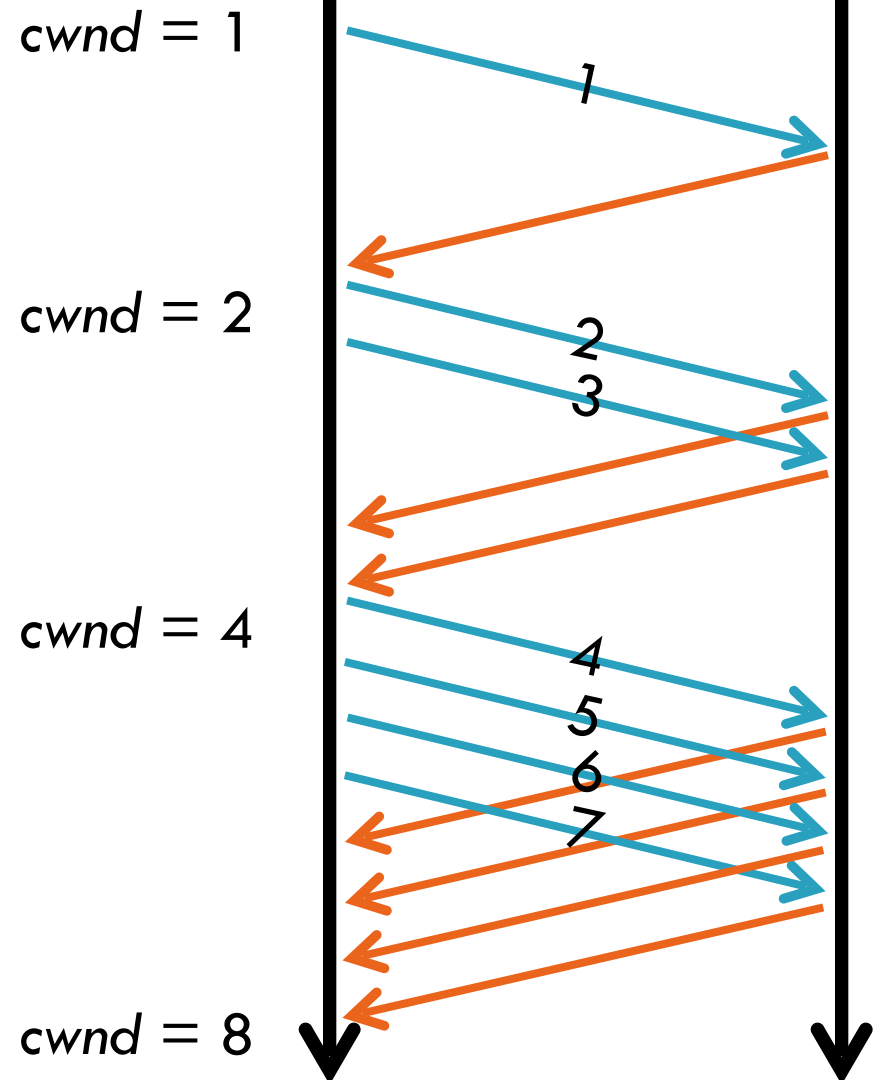
22



Slow Start Example

22

- $cwnd$ grows rapidly
- Slows down when...
 - ▣ $cwnd \geq ssthresh$
 - ▣ Or a packet drops



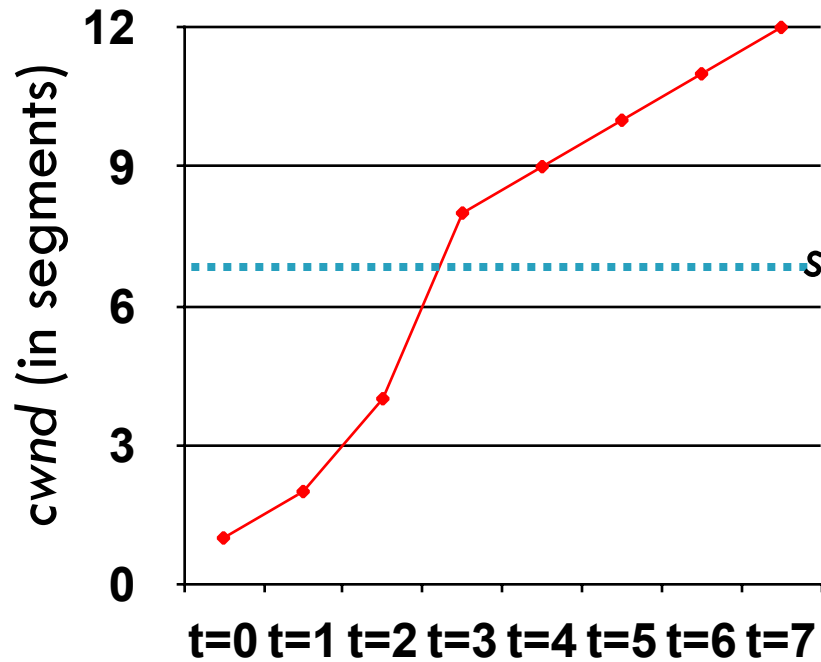
Congestion Avoidance

23

- AIMD mode
- *ssthresh* is lower-bound guess about location of the knee
- **If** $cwnd \geq ssthresh$ **then**
 - each time a segment is ACKed
 - increment $cwnd$ by $1/cwnd$ ($cwnd += 1/cwnd$).
- So $cwnd$ is increased by one only if all segments have been acknowledged

Congestion Avoidance Example

24



$cwnd = 1$

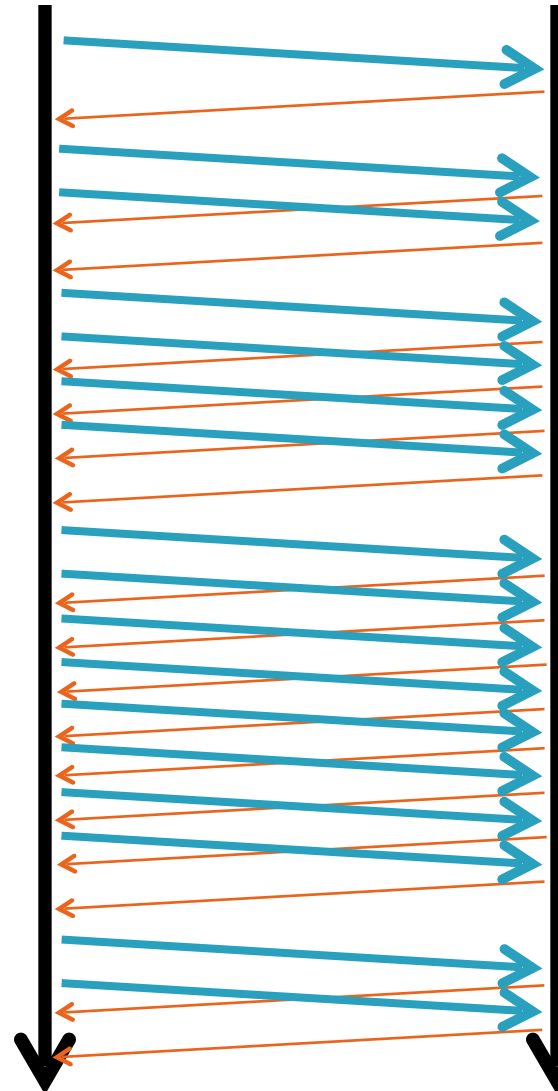
$cwnd = 2$

$cwnd = 4$

$ssthresh = 8$

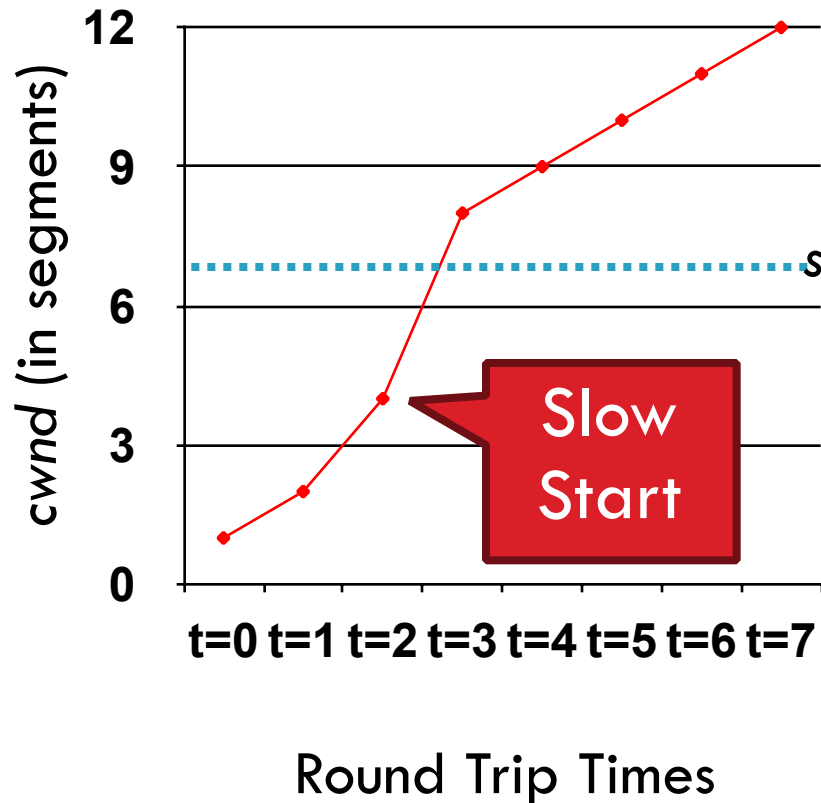
$cwnd = 8$

$cwnd = 9$



Congestion Avoidance Example

24



$cwnd = 1$

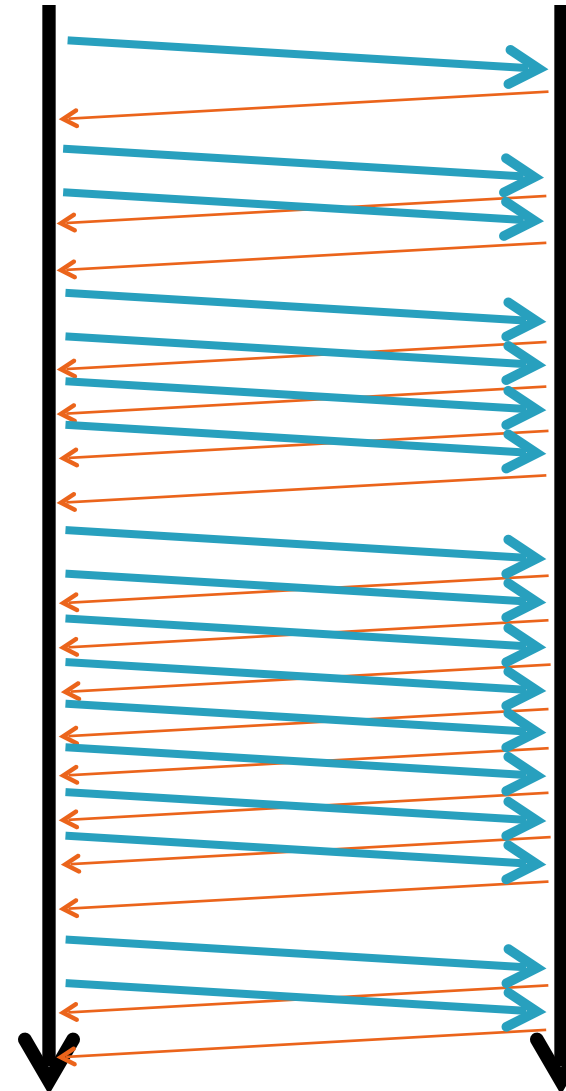
$cwnd = 2$

$cwnd = 4$

$ssthresh = 8$

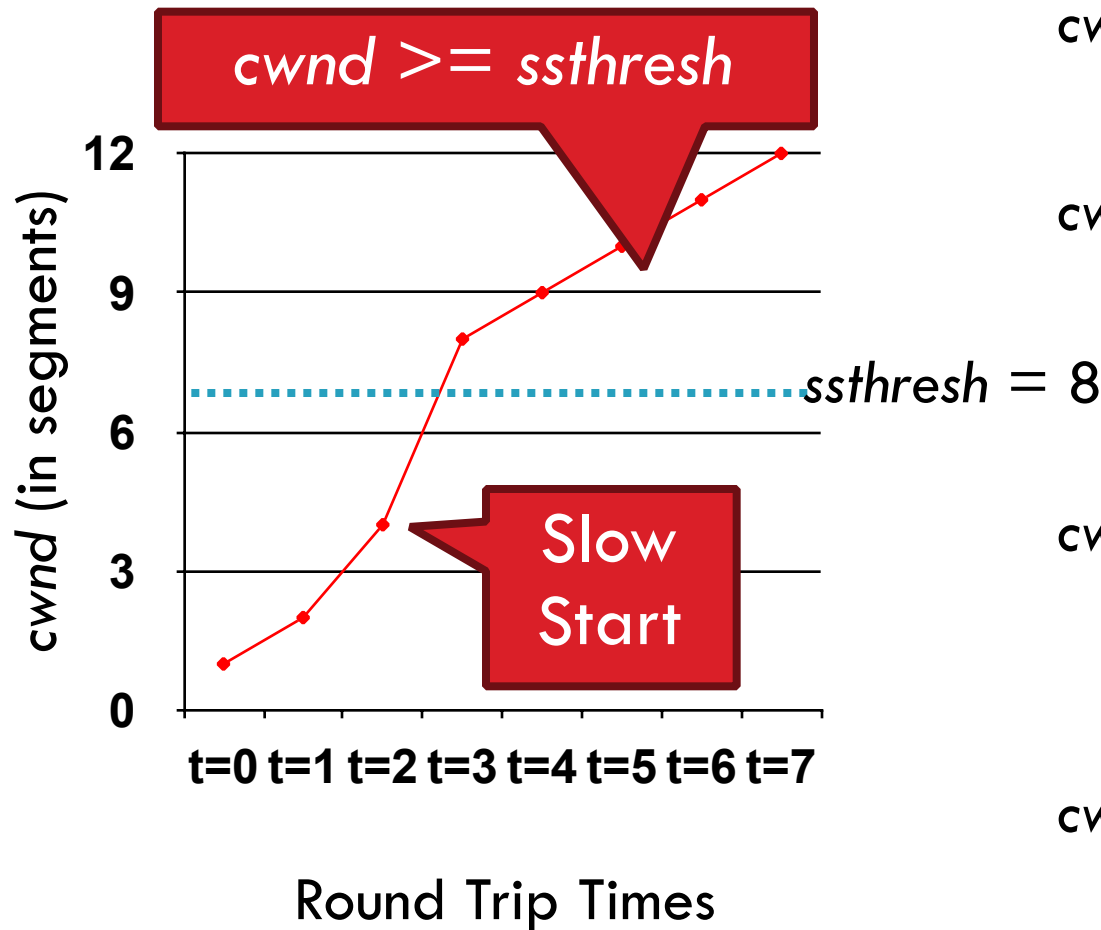
$cwnd = 8$

$cwnd = 9$



Congestion Avoidance Example

24



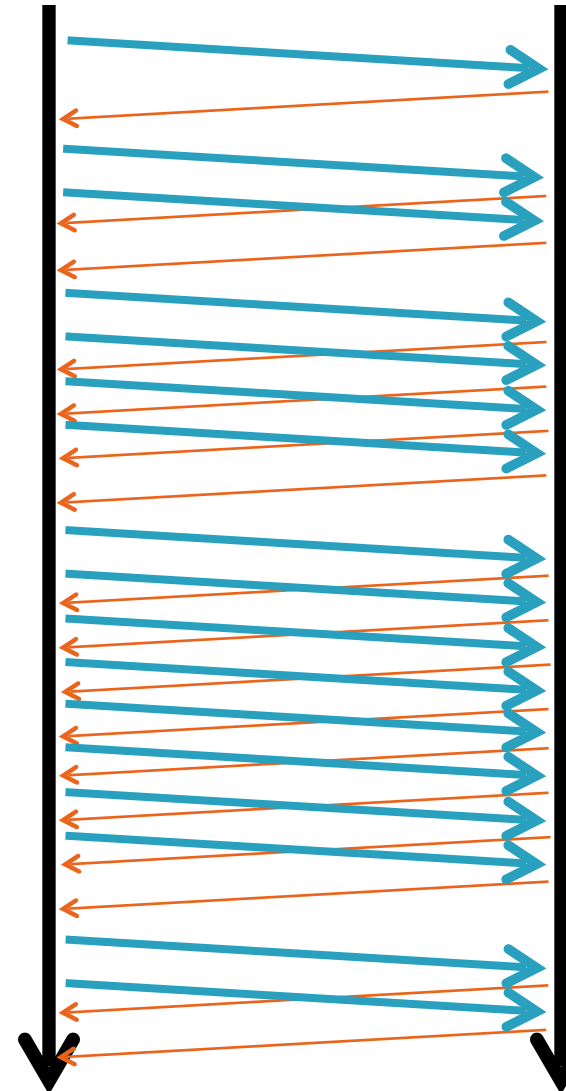
$cwnd = 1$

$cwnd = 2$

$cwnd = 4$

$cwnd = 8$

$cwnd = 9$



TCP Pseudocode

25

Initially:

```
  cwnd = 1;  
  ssthresh = adv_wnd;
```

New ack received:

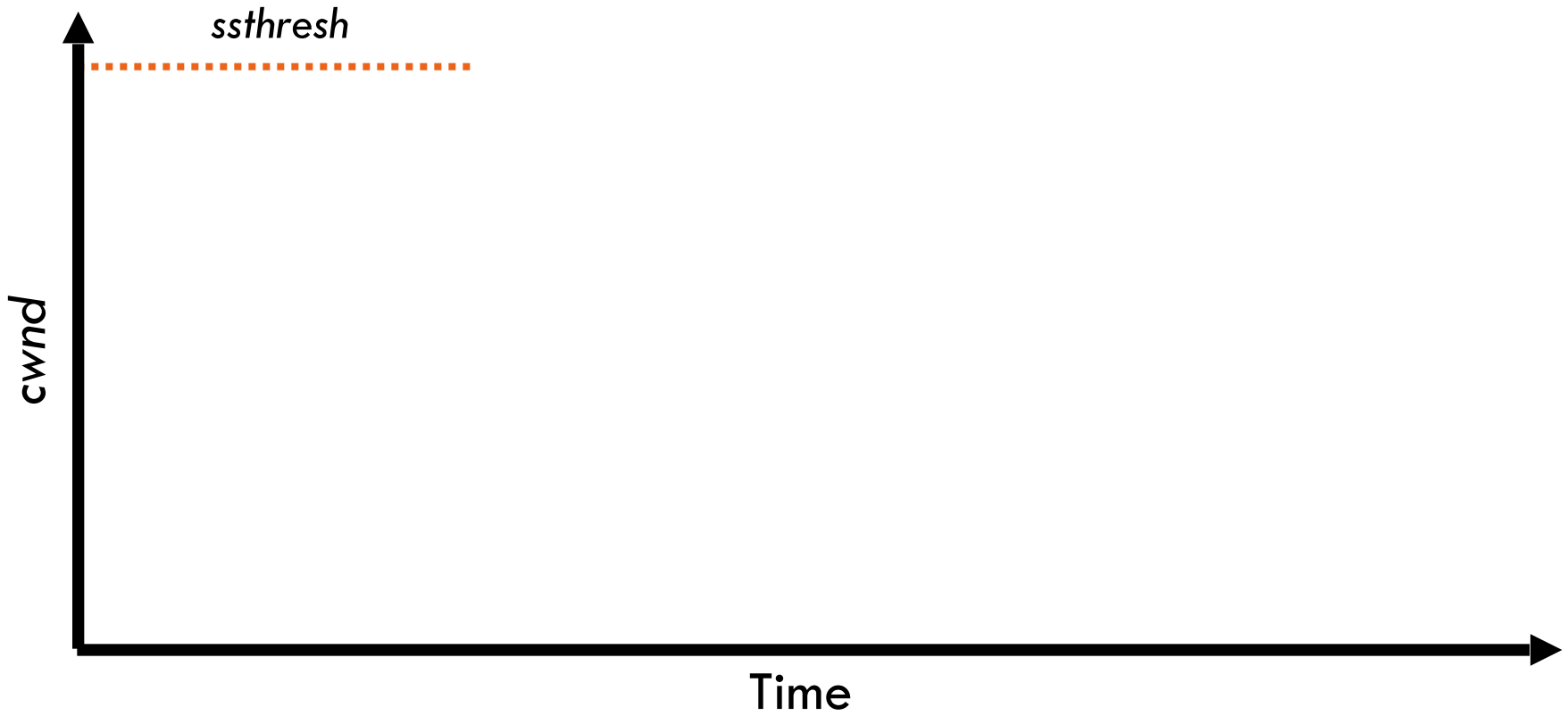
```
  if (cwnd < ssthresh)  
    /* Slow Start */  
    cwnd = cwnd + 1;  
  else  
    /* Congestion Avoidance */  
    cwnd = cwnd + 1 / cwnd;
```

Timeout:

```
  /* Multiplicative decrease */  
  ssthresh = cwnd / 2;  
  cwnd = 1;
```

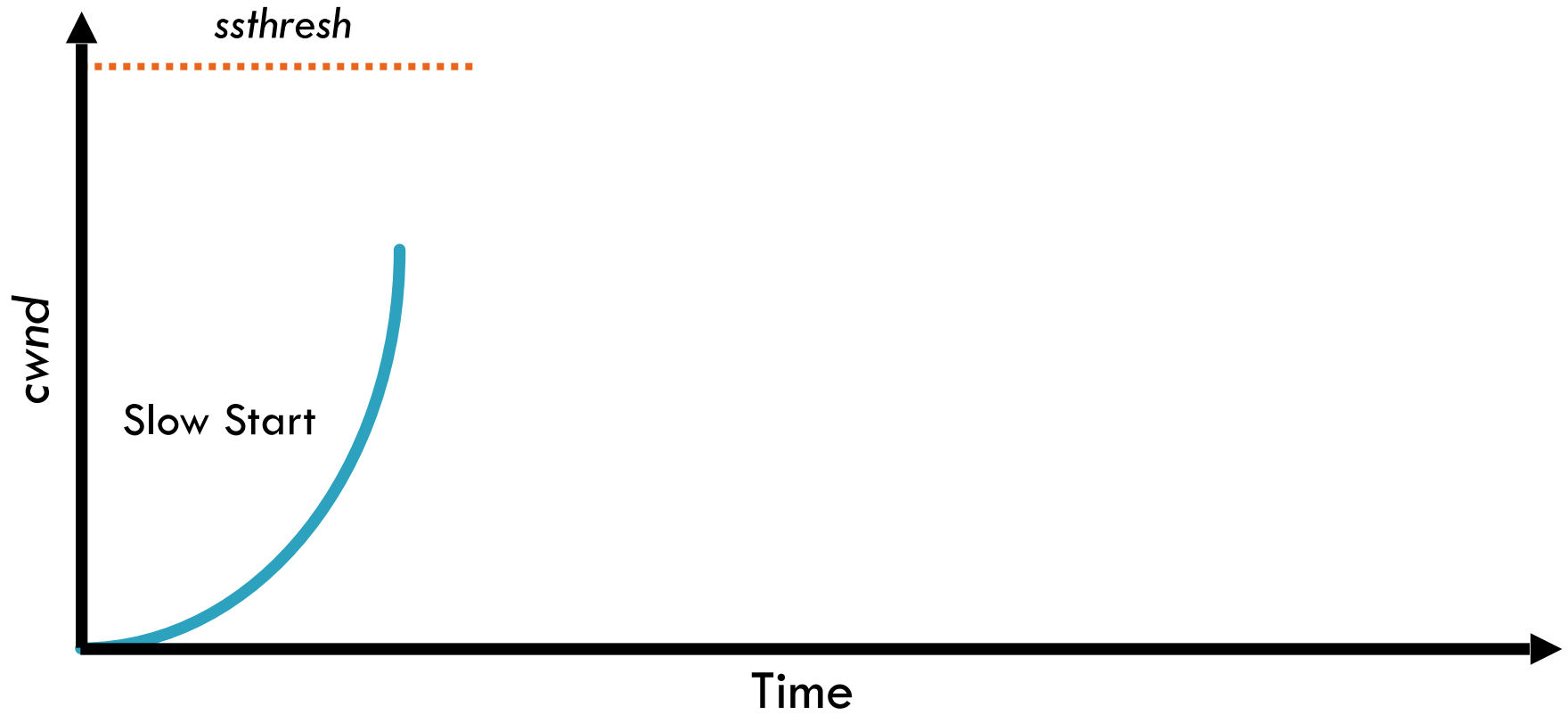
The Big Picture

26



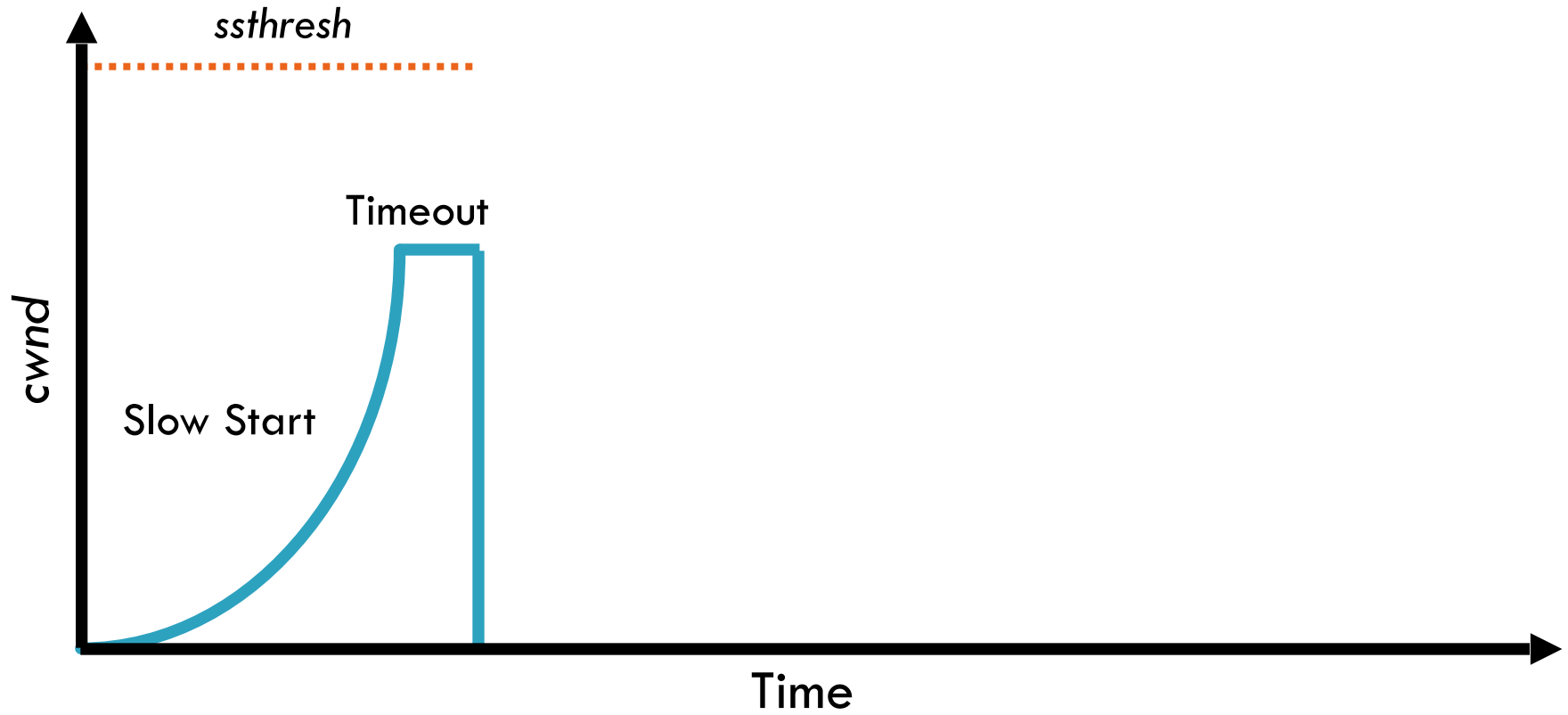
The Big Picture

26



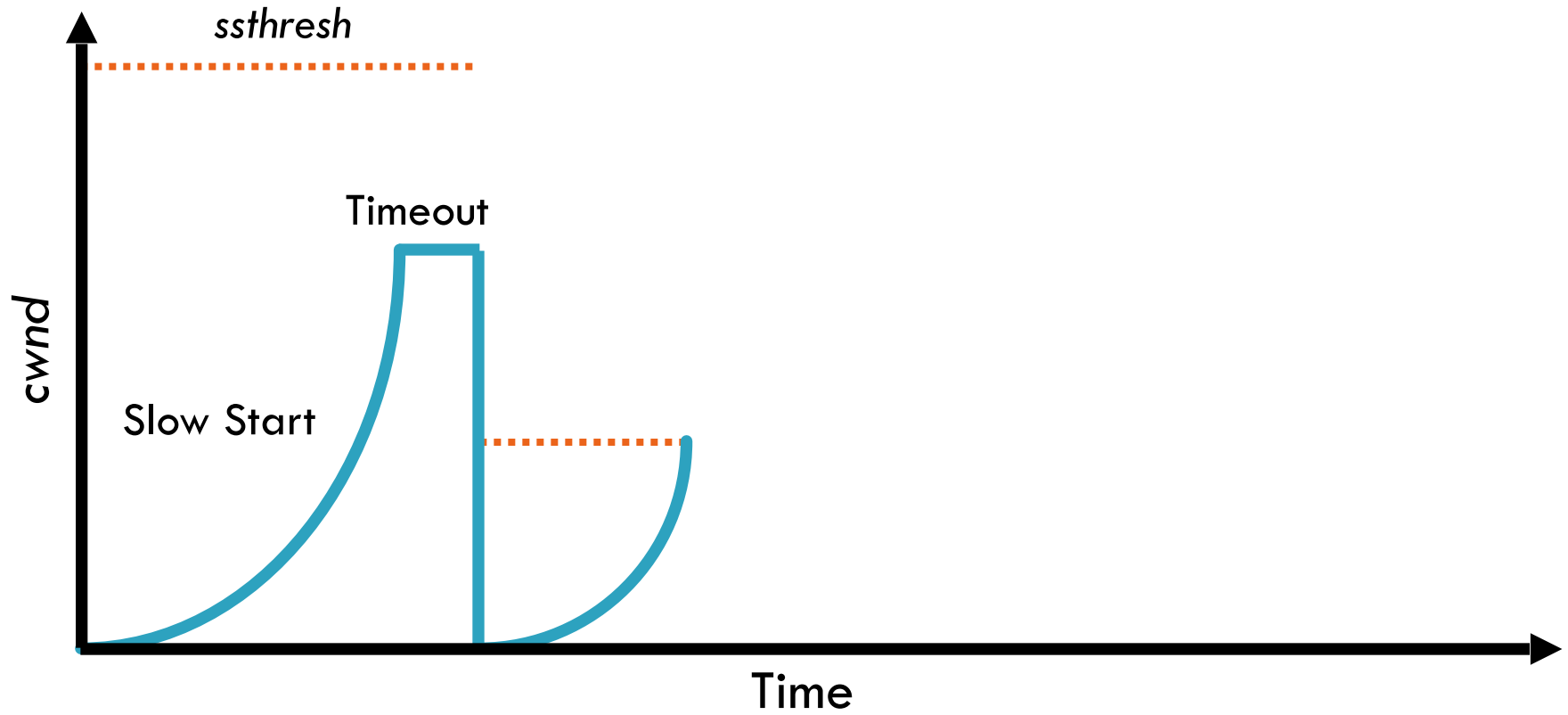
The Big Picture

26



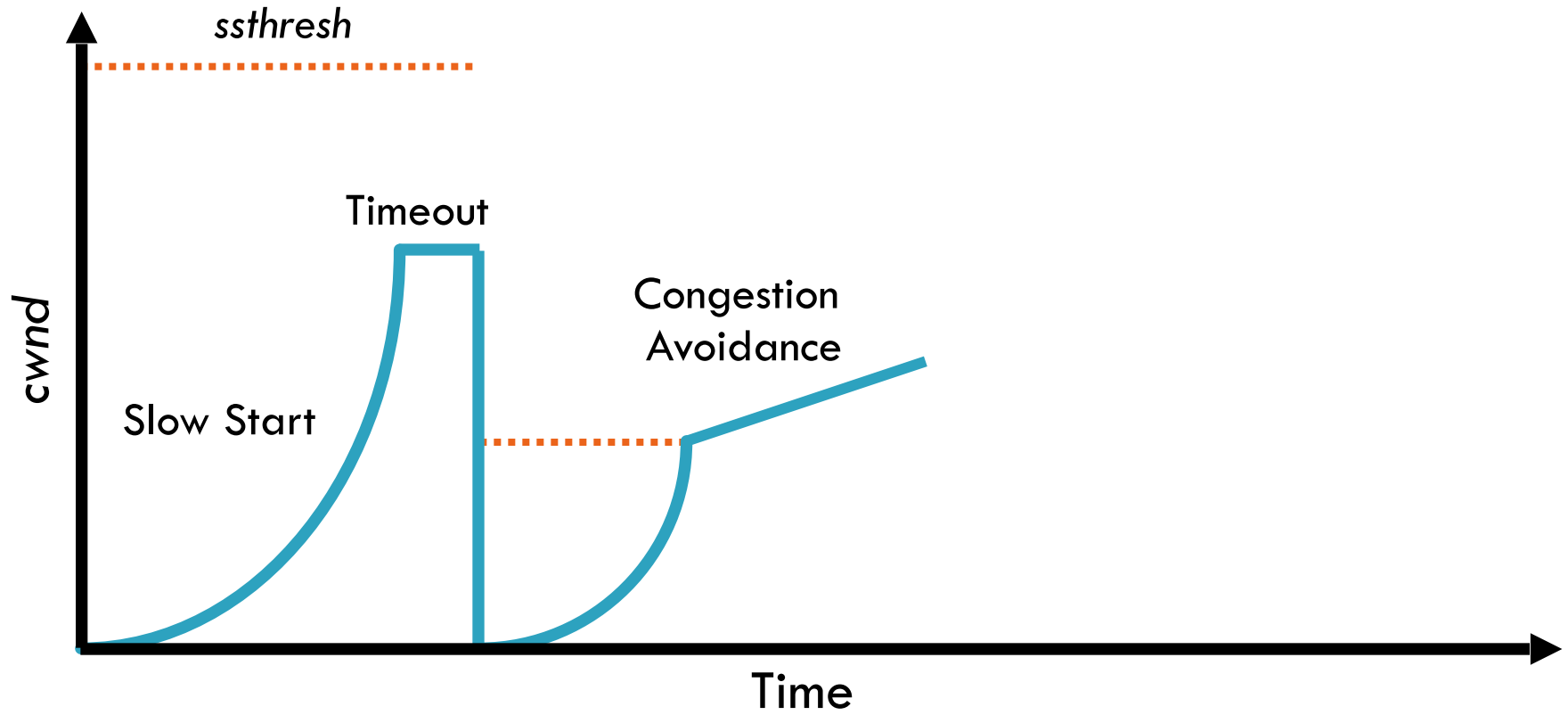
The Big Picture

26



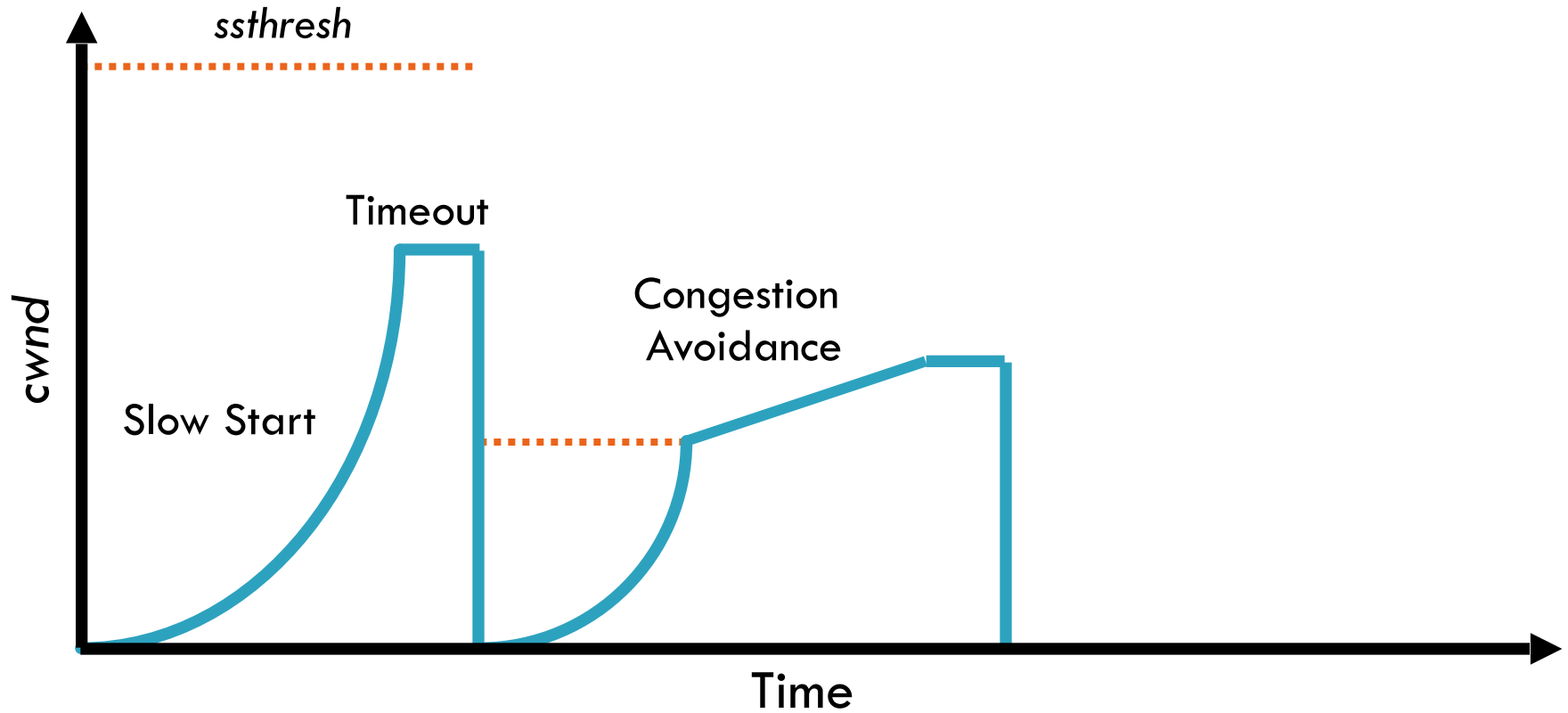
The Big Picture

26



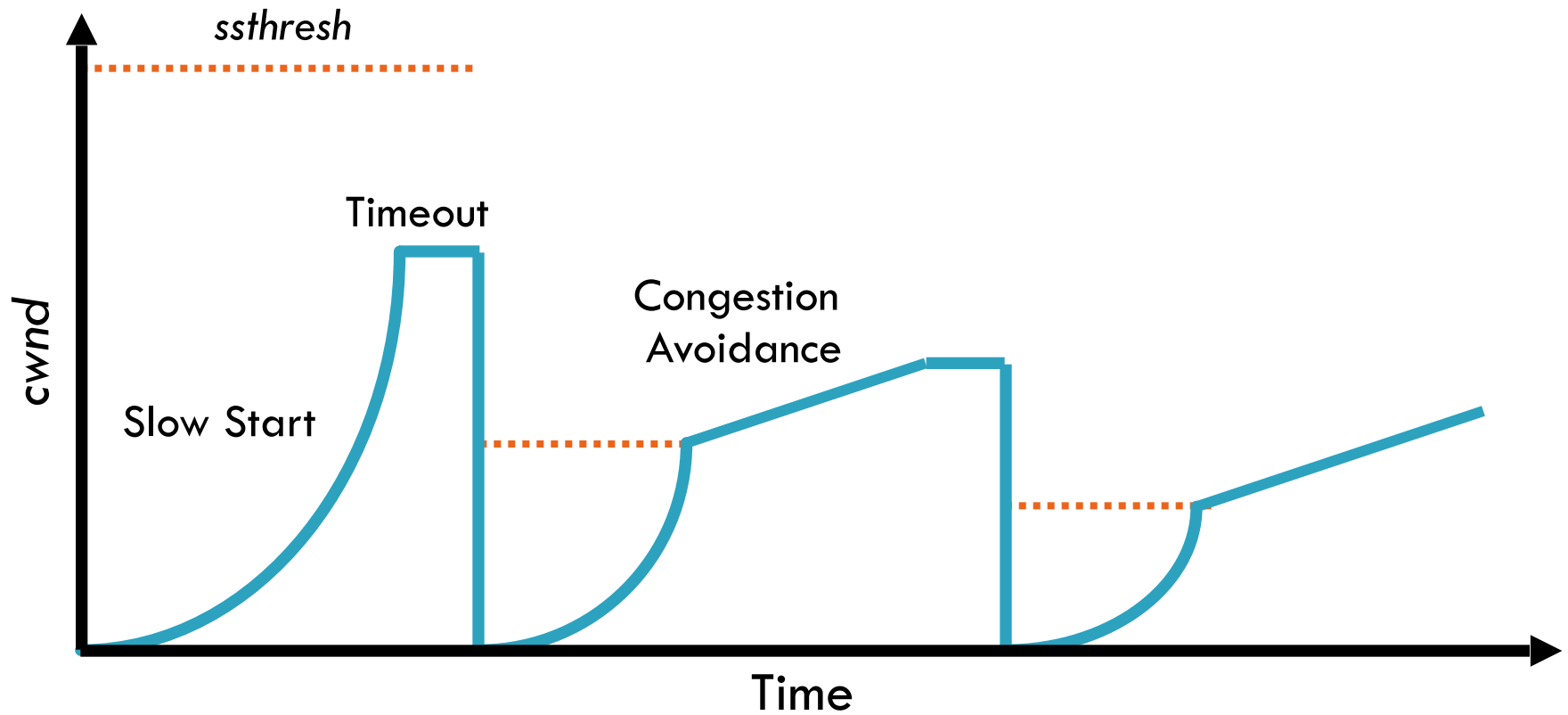
The Big Picture

26



The Big Picture

26



- ❑ Congestion Control
- ❑ Evolution of TCP
- ❑ Problems with TCP

The Evolution of TCP

28

- Thus far, we have discussed TCP Tahoe
 - ▣ Original version of TCP
- However, TCP was invented in 1974!
 - ▣ Today, there are many variants of TCP

The Evolution of TCP

28

- Thus far, we have discussed TCP Tahoe
 - ▣ Original version of TCP
- However, TCP was invented in 1974!
 - ▣ Today, there are many variants of TCP
- Early, popular variant: TCP Reno
 - ▣ Tahoe features, plus...
 - ▣ Fast retransmit
 - ▣ Fast recovery

TCP Reno: Fast Retransmit

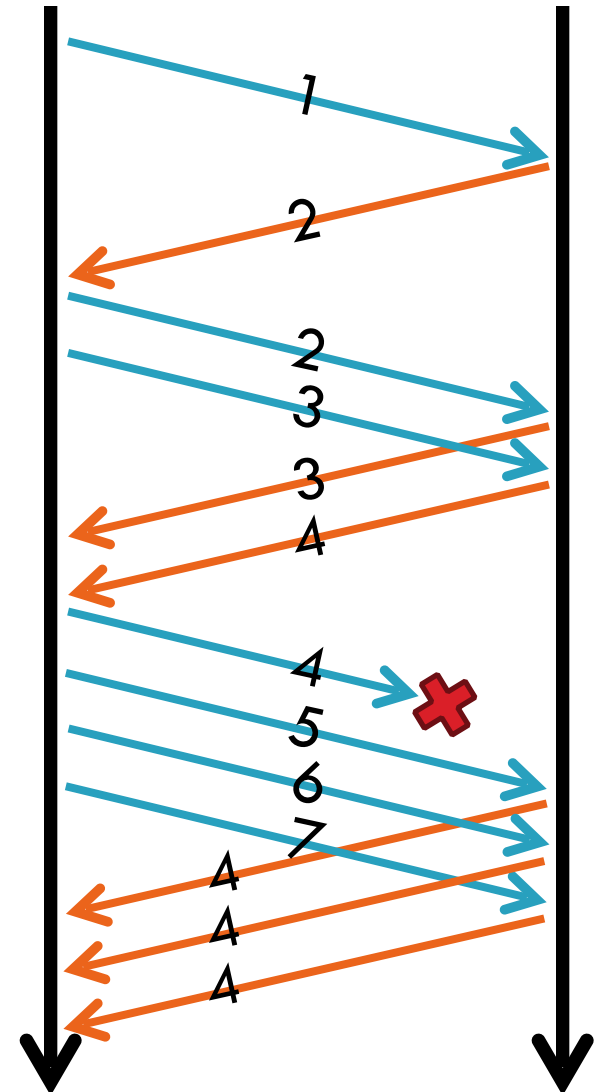
29

- Problem: in Tahoe, if segment is lost, there is a long wait until the RTO
- Reno: retransmit after 3 duplicate ACKs

$cwnd = 1$

$cwnd = 2$

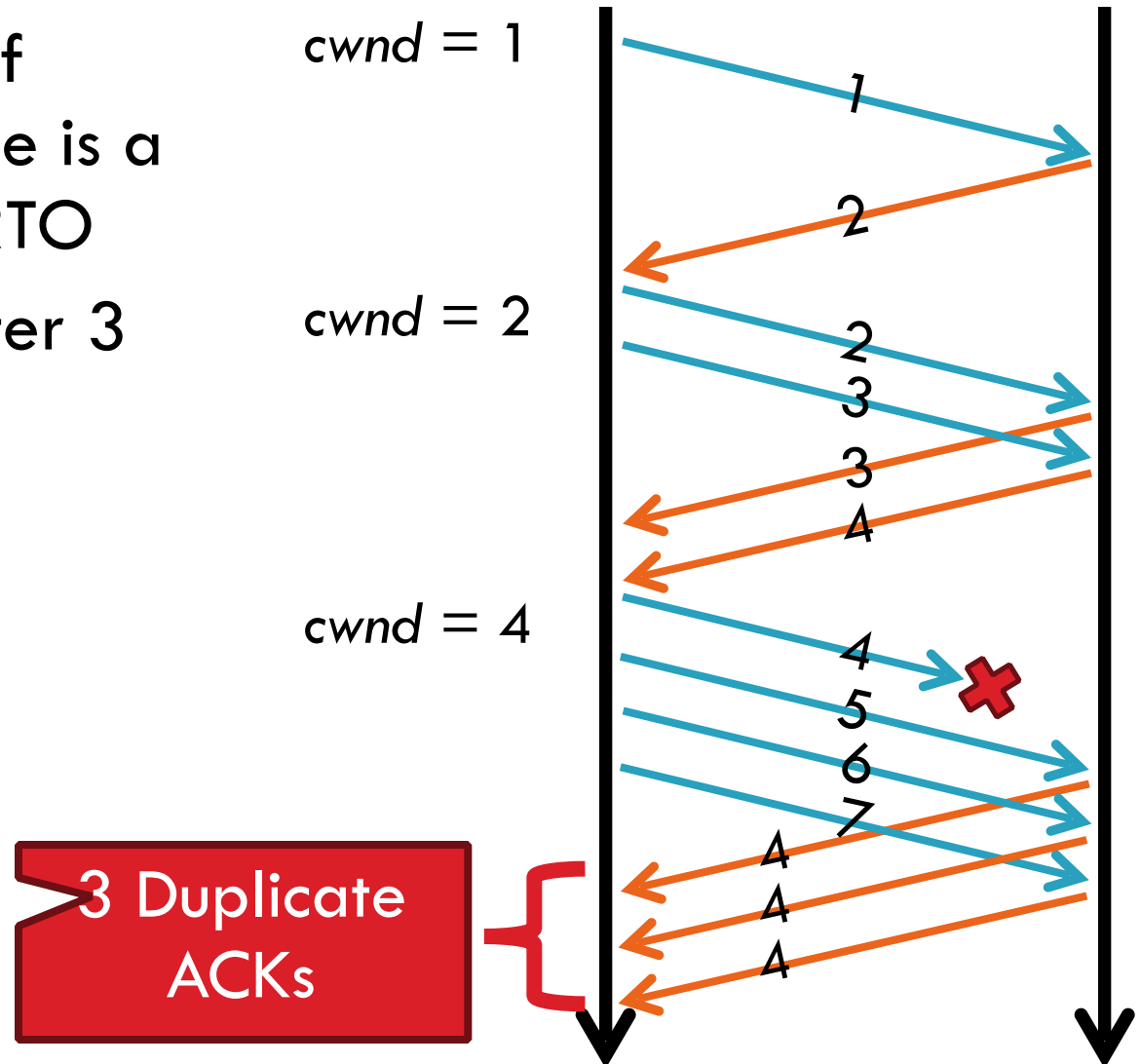
$cwnd = 4$



TCP Reno: Fast Retransmit

29

- Problem: in Tahoe, if segment is lost, there is a long wait until the RTO
- Reno: retransmit after 3 duplicate ACKs



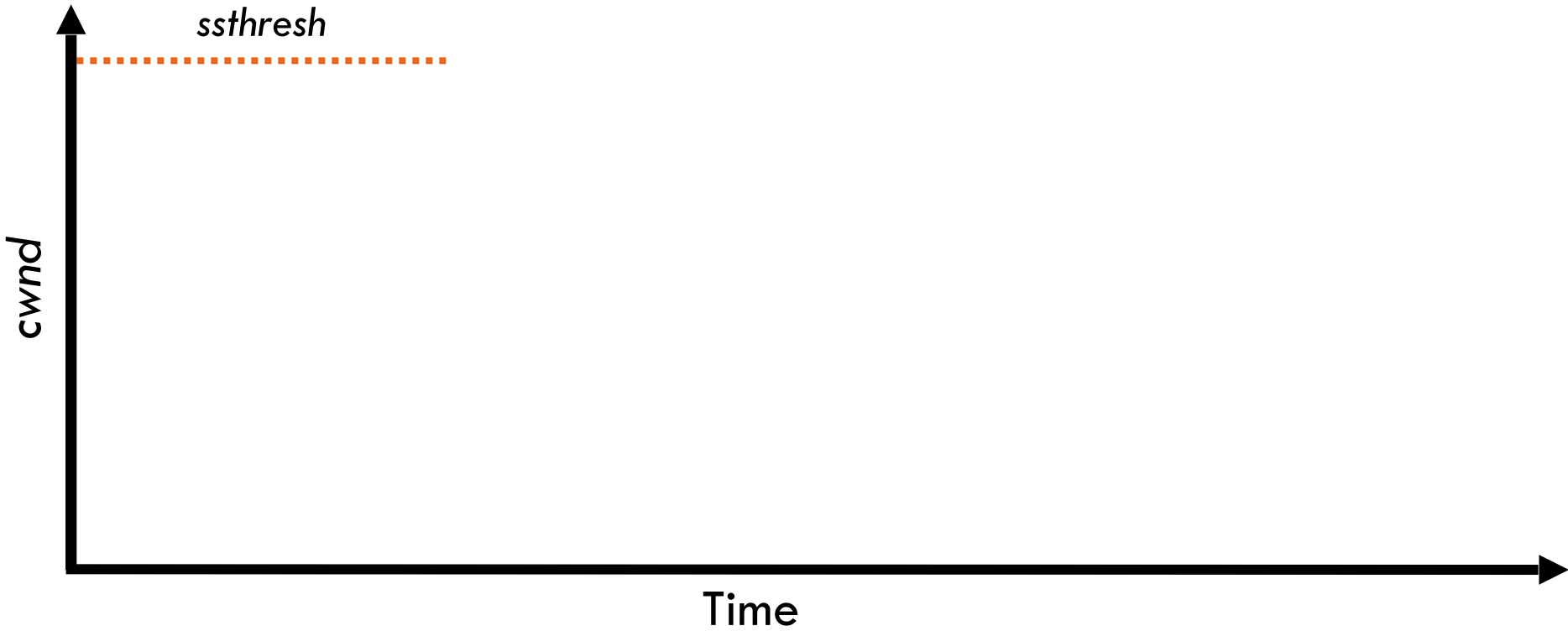
TCP Reno: Fast Recovery

30

- After a fast-retransmit set $cwnd$ to $ssthresh/2$
 - i.e. don't reset $cwnd$ to 1
 - Avoid unnecessary return to slow start
 - Prevents expensive timeouts
- But when RTO expires still do $cwnd = 1$
 - Return to slow start, same as Tahoe
 - Indicates packets aren't being delivered at all
 - i.e. congestion must be really bad

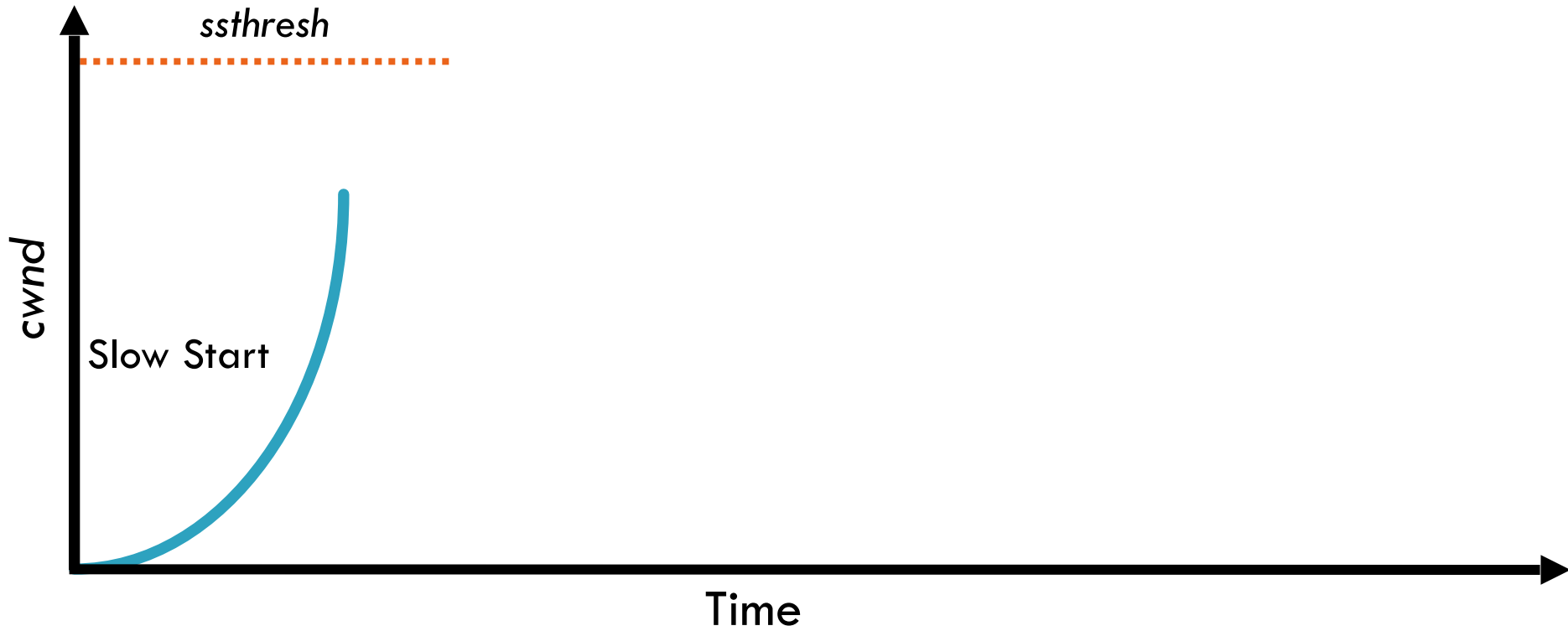
Fast Retransmit and Fast Recovery

31



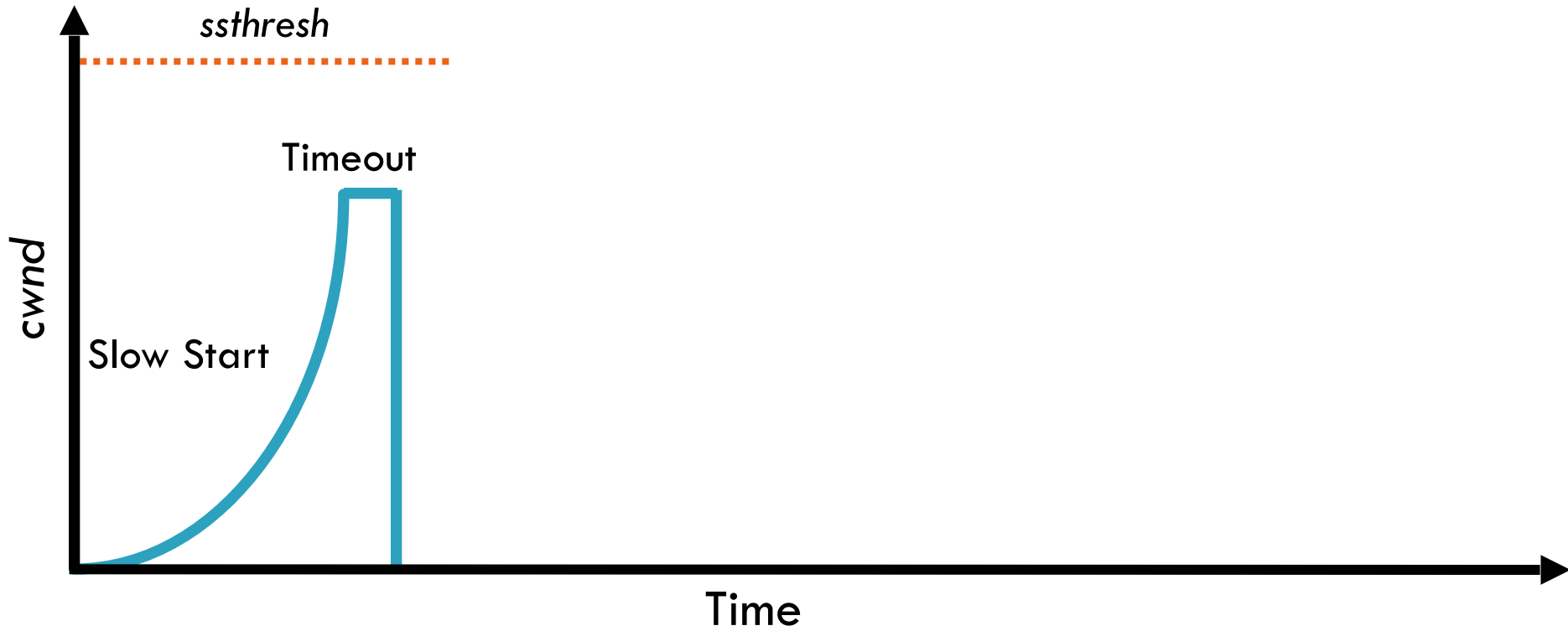
Fast Retransmit and Fast Recovery

31



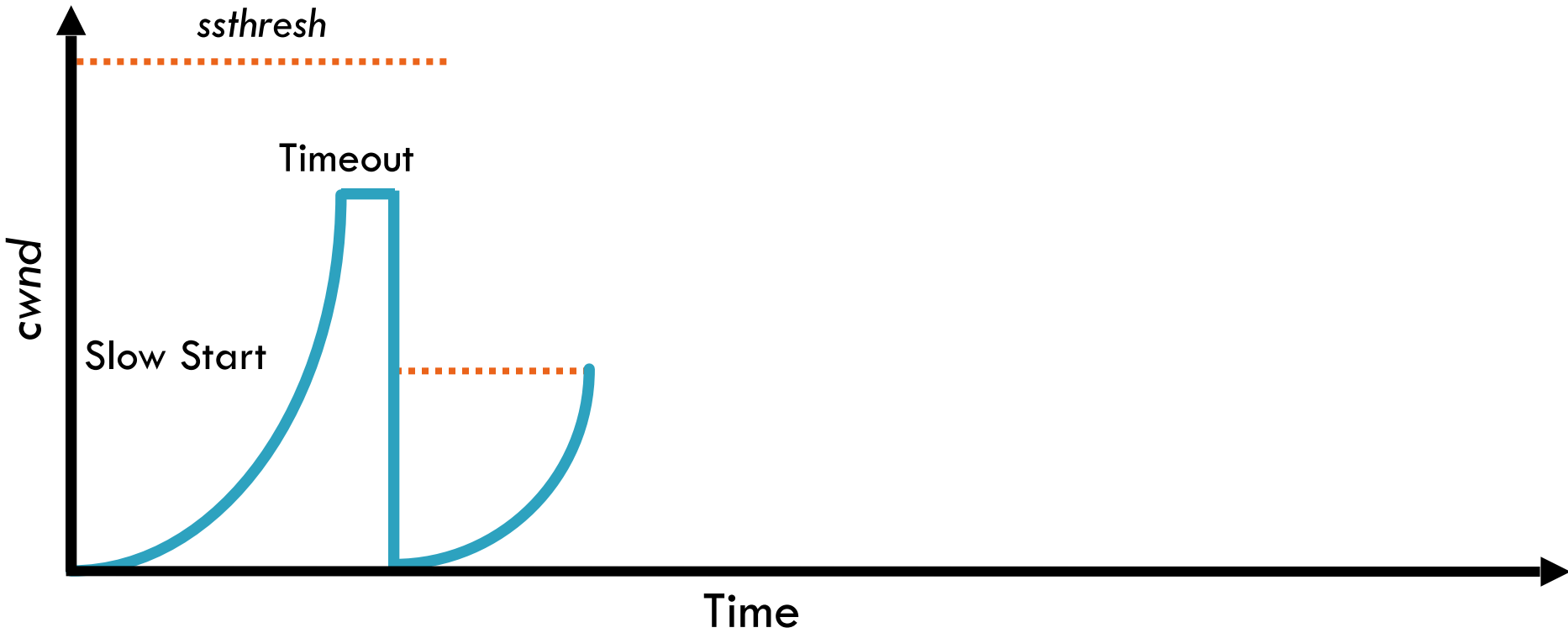
Fast Retransmit and Fast Recovery

31



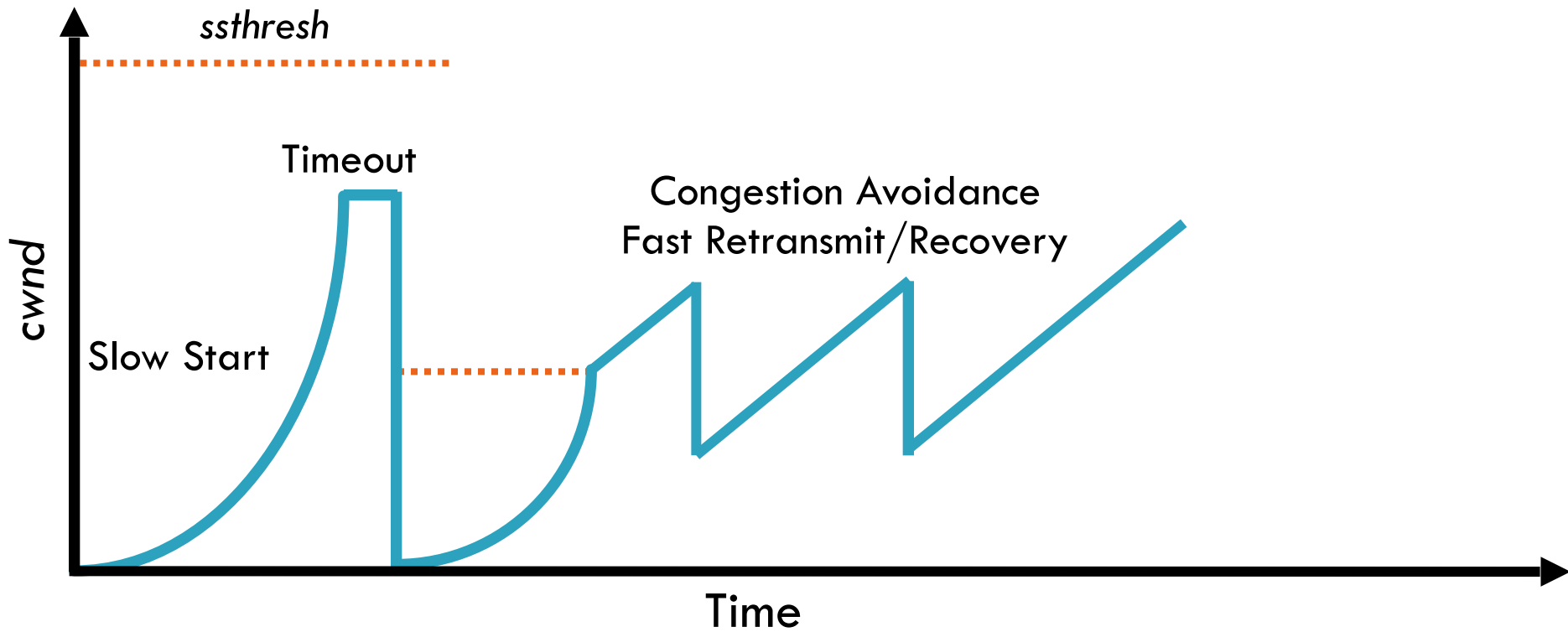
Fast Retransmit and Fast Recovery

31



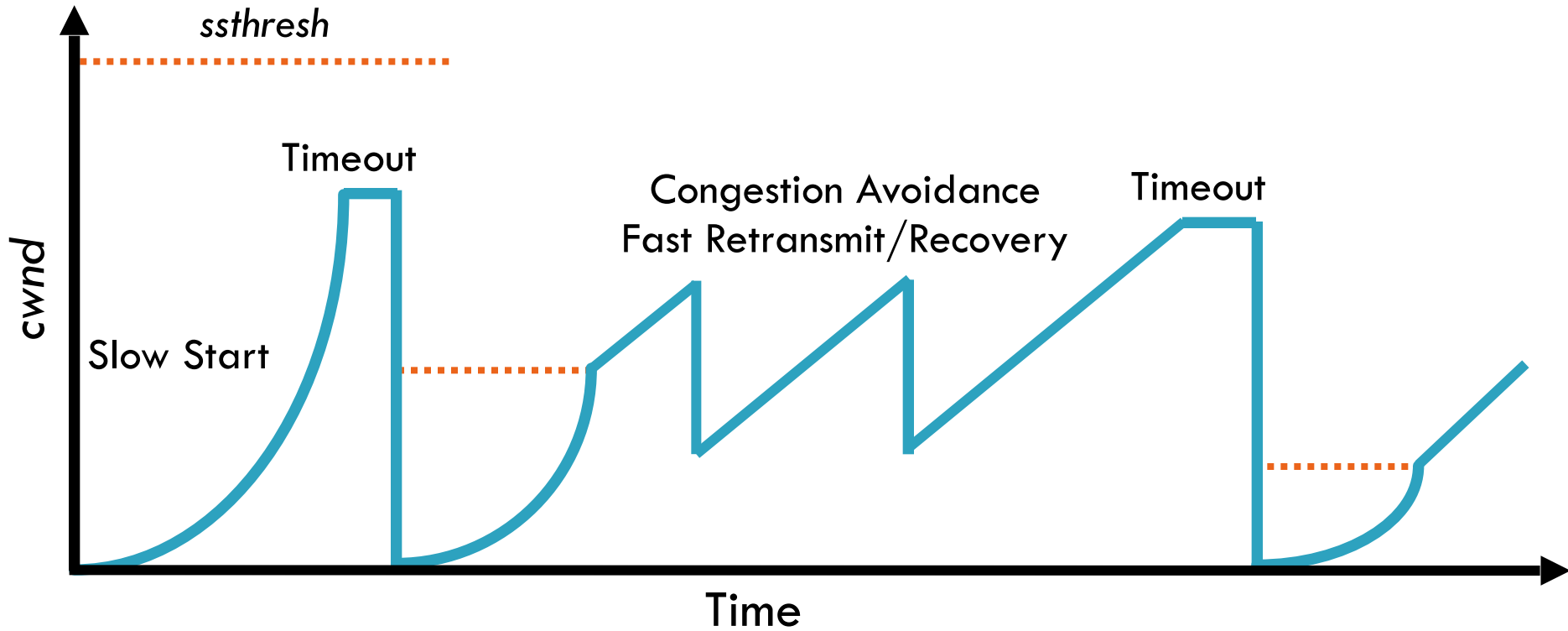
Fast Retransmit and Fast Recovery

31



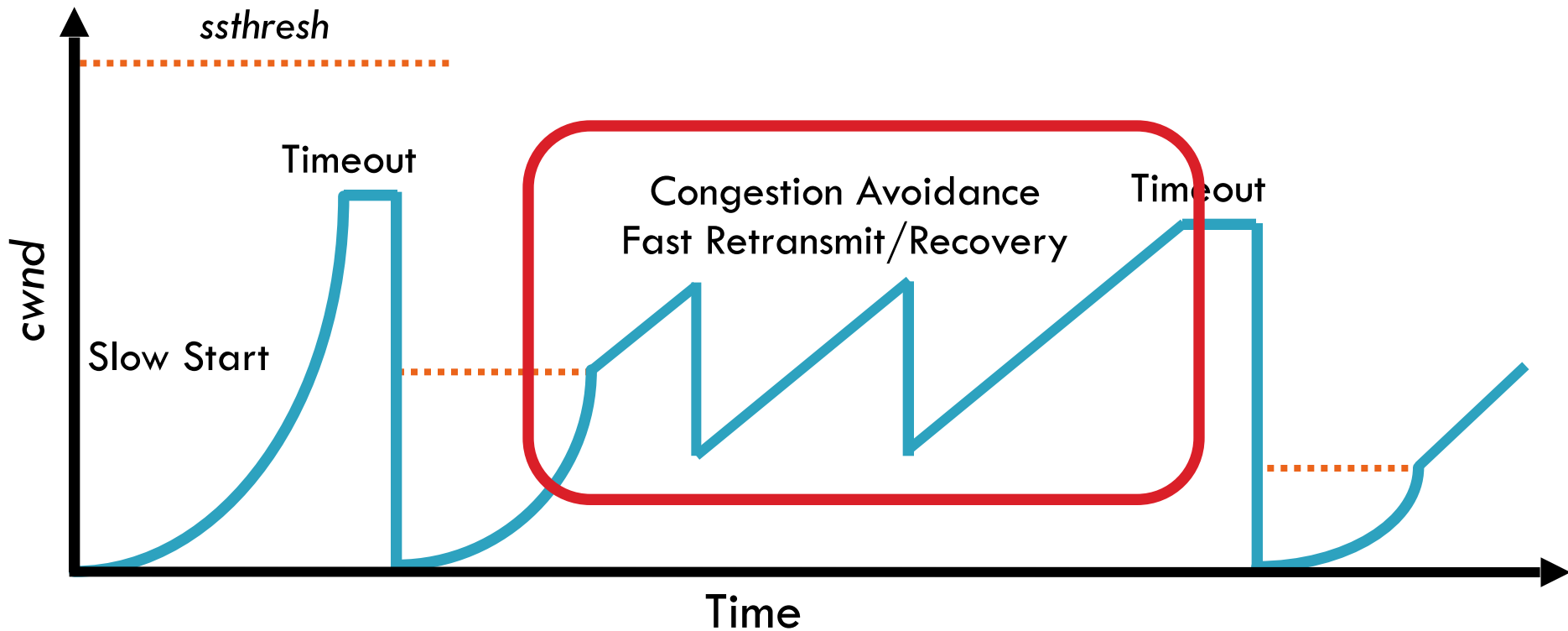
Fast Retransmit and Fast Recovery

31



Fast Retransmit and Fast Recovery

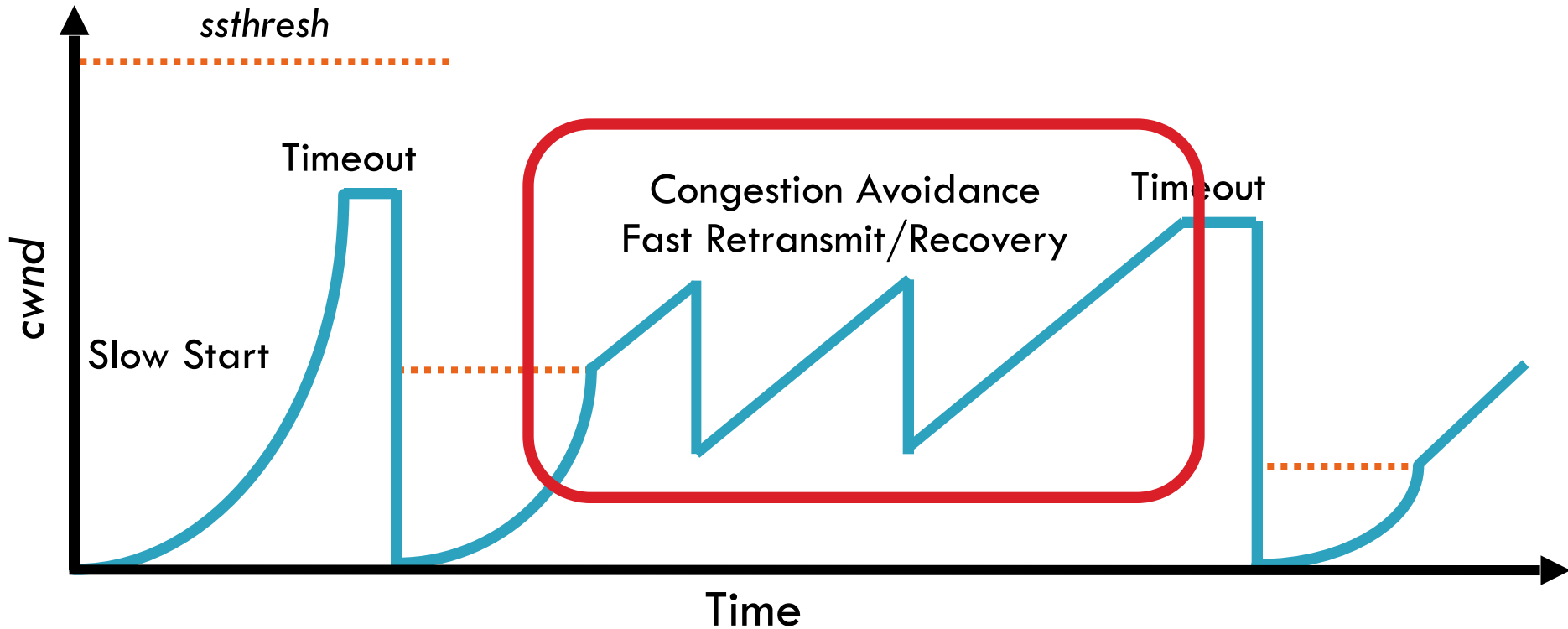
31



- At steady state, $cwnd$ oscillates around the optimal window size

Fast Retransmit and Fast Recovery

31



- At steady state, *cwnd* oscillates around the optimal window size
- TCP always forces packet drops

Many TCP Variants...

32

- Tahoe: the original
 - ▣ Slow start with AIMD
 - ▣ Dynamic RTO based on RTT estimate
- Reno: fast retransmit and fast recovery

Many TCP Variants...

32

- Tahoe: the original
 - ▣ Slow start with AIMD
 - ▣ Dynamic RTO based on RTT estimate
- Reno: fast retransmit and fast recovery
- NewReno: improved fast retransmit
 - ▣ Each duplicate ACK triggers a retransmission
 - ▣ Problem: >3 out-of-order packets causes pathological retransmissions

Many TCP Variants...

32

- Tahoe: the original
 - ▣ Slow start with AIMD
 - ▣ Dynamic RTO based on RTT estimate
- Reno: fast retransmit and fast recovery
- NewReno: improved fast retransmit
 - ▣ Each duplicate ACK triggers a retransmission
 - ▣ Problem: >3 out-of-order packets causes pathological retransmissions
- Vegas: delay-based congestion avoidance

Many TCP Variants...

32

- Tahoe: the original
 - ▣ Slow start with AIMD
 - ▣ Dynamic RTO based on RTT estimate
- Reno: fast retransmit and fast recovery
- NewReno: improved fast retransmit
 - ▣ Each duplicate ACK triggers a retransmission
 - ▣ Problem: >3 out-of-order packets causes pathological retransmissions
- Vegas: delay-based congestion avoidance
- And many, many, many more...

TCP in the Real World

33

- What are the most popular variants today?
 - ▣ Key problem: TCP performs poorly on high bandwidth-delay product networks (like the modern Internet)
 - ▣ Compound TCP (Windows)
 - Based on Reno
 - Uses two congestion windows: delay based and loss based
 - Thus, it uses a *compound* congestion controller
 - ▣ TCP CUBIC (Linux)
 - Enhancement of BIC (Binary Increase Congestion Control)
 - Window size controlled by cubic function
 - Parameterized by the time T since the last dropped packet

High Bandwidth-Delay Product

34

- Key Problem: TCP performs poorly when
 - The capacity of the network (bandwidth) is large
 - The delay (RTT) of the network is large
 - Or, when bandwidth * delay is large
 - $b * d =$ maximum amount of in-flight data in the network
 - a.k.a. the bandwidth-delay product

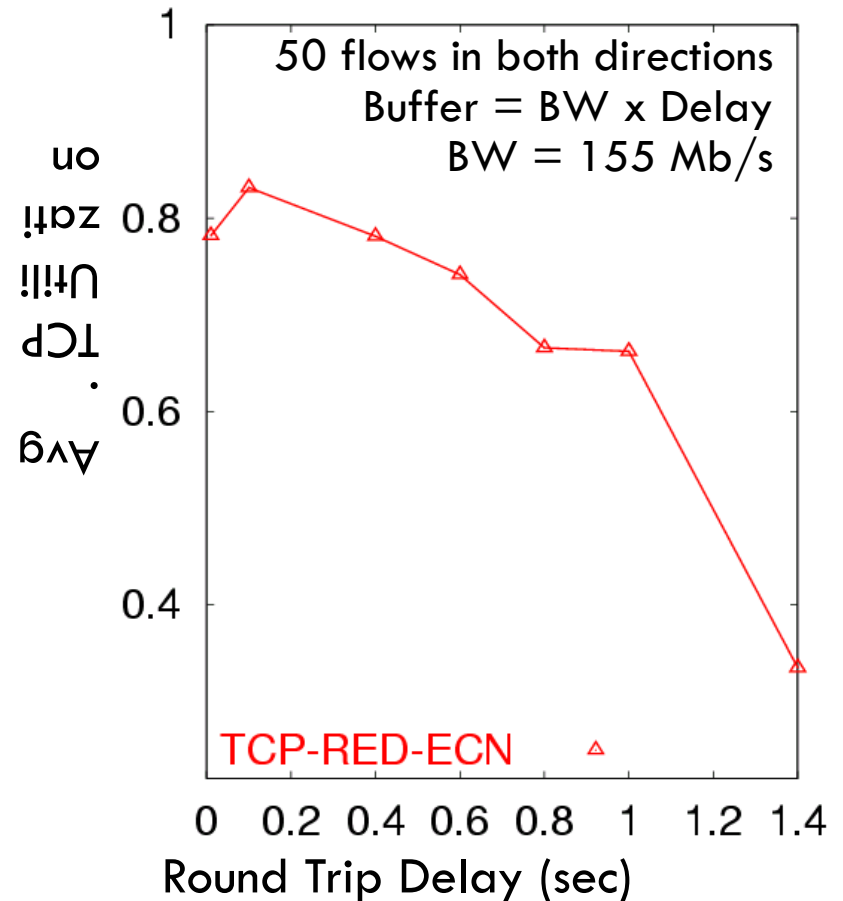
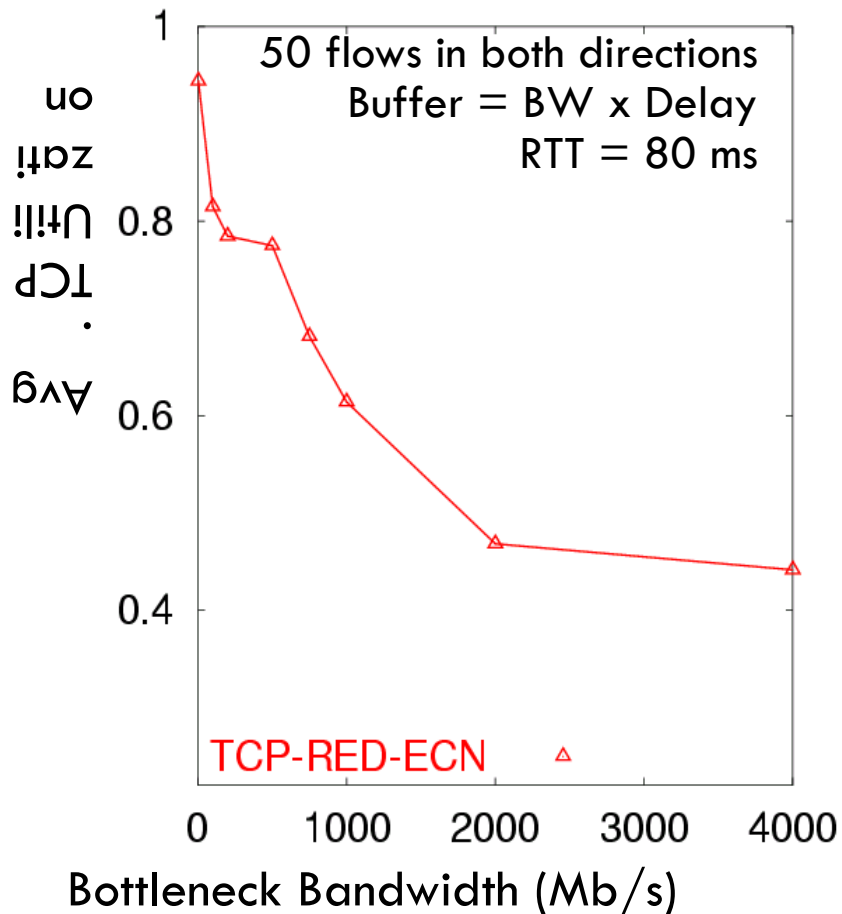
High Bandwidth-Delay Product

34

- Key Problem: TCP performs poorly when
 - The capacity of the network (bandwidth) is large
 - The delay (RTT) of the network is large
 - Or, when bandwidth * delay is large
 - $b * d =$ maximum amount of in-flight data in the network
 - a.k.a. the bandwidth-delay product
- Why does TCP perform poorly?
 - Slow start and additive increase are slow to converge
 - TCP is ACK clocked
 - i.e. TCP can only react as quickly as ACKs are received
 - Large RTT \rightarrow ACKs are delayed \rightarrow TCP is slow to react

Poor Performance of TCP Reno CC

35



Goals

36

- Fast window growth
 - Slow start and additive increase are too slow when bandwidth is large
 - Want to converge more quickly

Goals

36

- Fast window growth
 - ▣ Slow start and additive increase are too slow when bandwidth is large
 - ▣ Want to converge more quickly
- Maintain fairness with other TCP variants
 - ▣ Window growth cannot be too aggressive

Goals

36

- Fast window growth
 - ▣ Slow start and additive increase are too slow when bandwidth is large
 - ▣ Want to converge more quickly
- Maintain fairness with other TCP variants
 - ▣ Window growth cannot be too aggressive
- Improve RTT fairness
 - ▣ TCP Tahoe/Reno flows are not fair when RTTs vary widely

Goals

36

- Fast window growth
 - ▣ Slow start and additive increase are too slow when bandwidth is large
 - ▣ Want to converge more quickly
- Maintain fairness with other TCP variants
 - ▣ Window growth cannot be too aggressive
- Improve RTT fairness
 - ▣ TCP Tahoe/Reno flows are not fair when RTTs vary widely
- Simple implementation

Compound TCP Implementation

37

- Default TCP implementation in Windows
- Key idea: split *cwnd* into two separate windows
 - ▣ Traditional, loss-based window
 - ▣ New, delay-based window

Compound TCP Implementation

37

- Default TCP implementation in Windows
- Key idea: split *cwnd* into two separate windows
 - Traditional, loss-based window
 - New, delay-based window
- $wnd = \min(cwnd + dwnd, adv_wnd)$
 - *cwnd* is controlled by AIMD
 - *dwnd* is the delay window

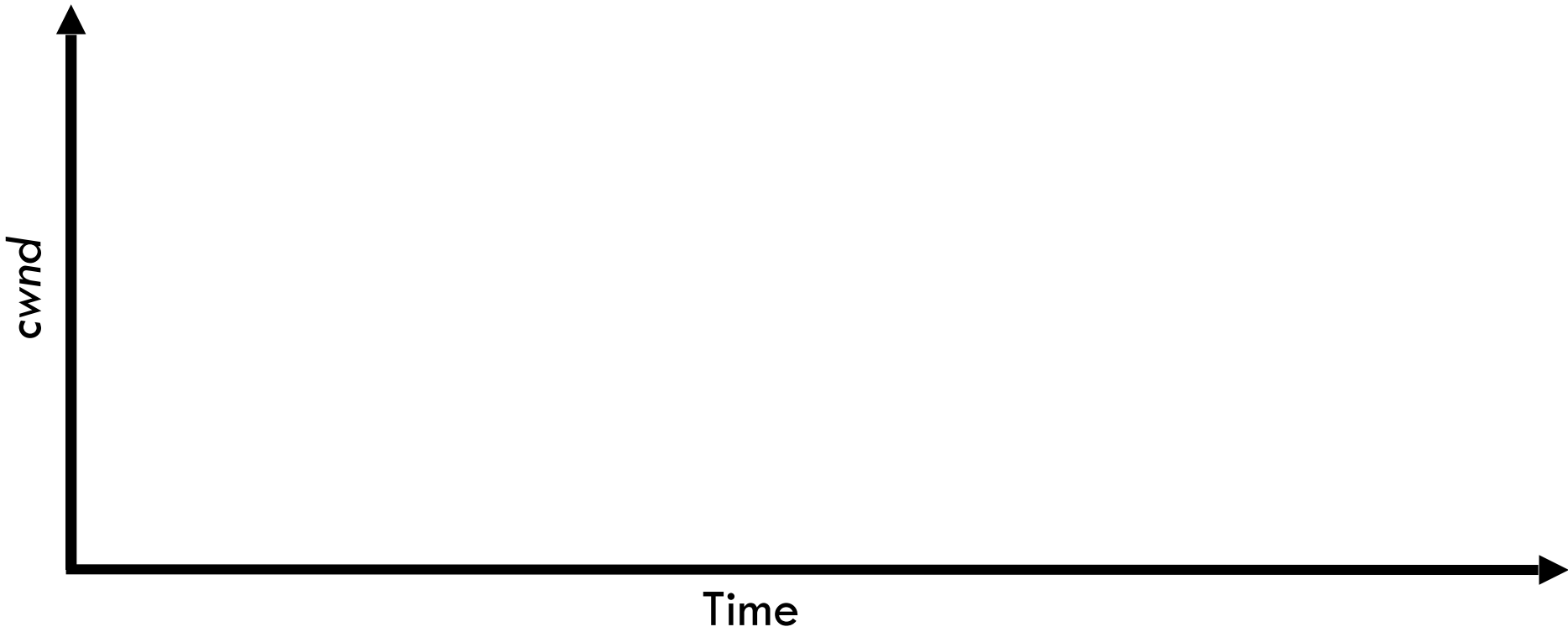
Compound TCP Implementation

37

- Default TCP implementation in Windows
- Key idea: split *cwnd* into two separate windows
 - Traditional, loss-based window
 - New, delay-based window
- $wnd = \min(cwnd + dwnd, adv_wnd)$
 - *cwnd* is controlled by AIMD
 - *dwnd* is the delay window
- Rules for adjusting *dwnd*:
 - If RTT is increasing, decrease *dwnd* ($dwnd \geq 0$)
 - If RTT is decreasing, increase *dwnd*
 - Increase/decrease are proportional to the rate of change

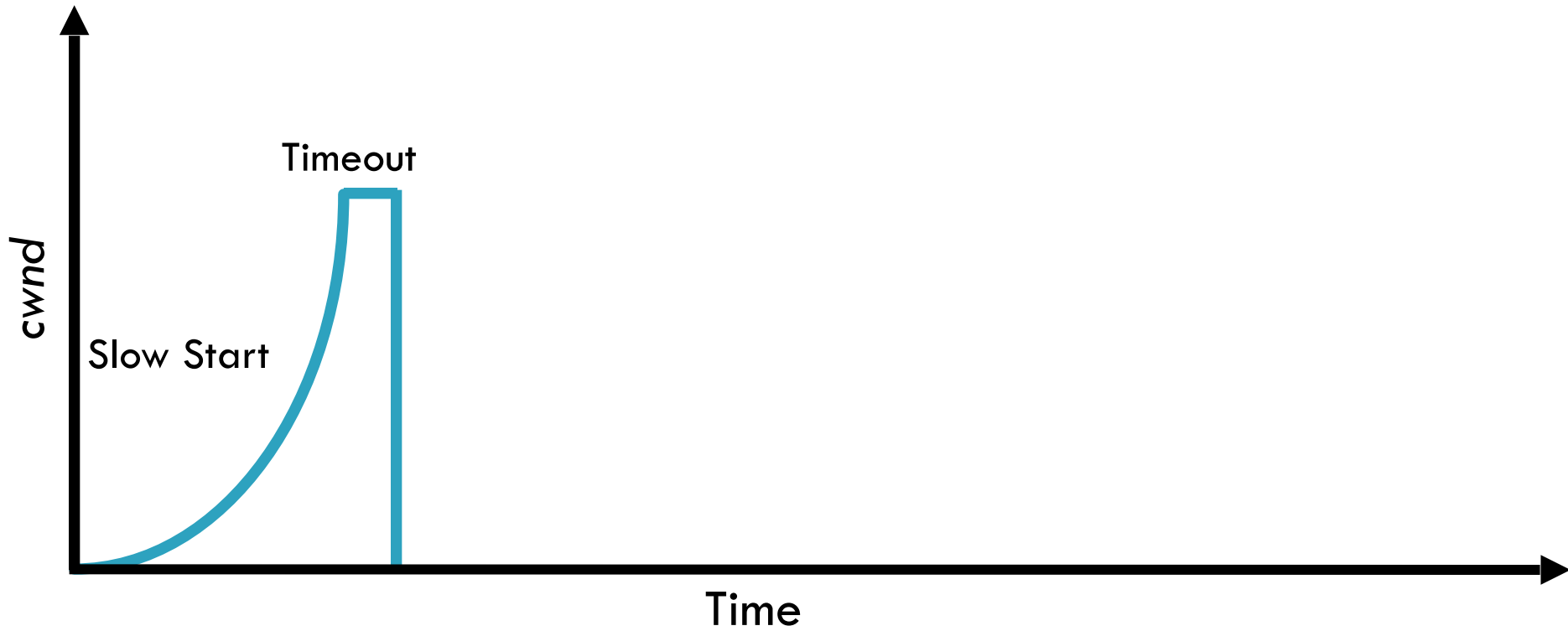
Compound TCP Example

38



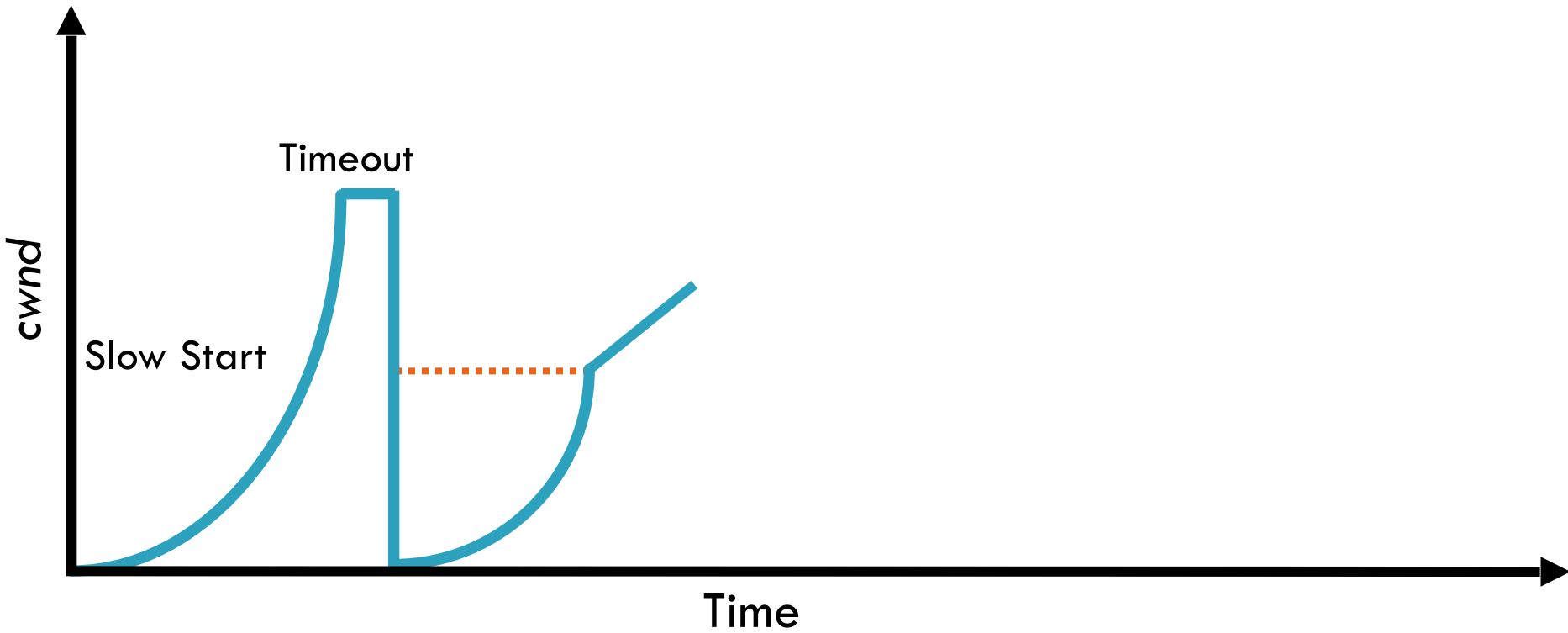
Compound TCP Example

38



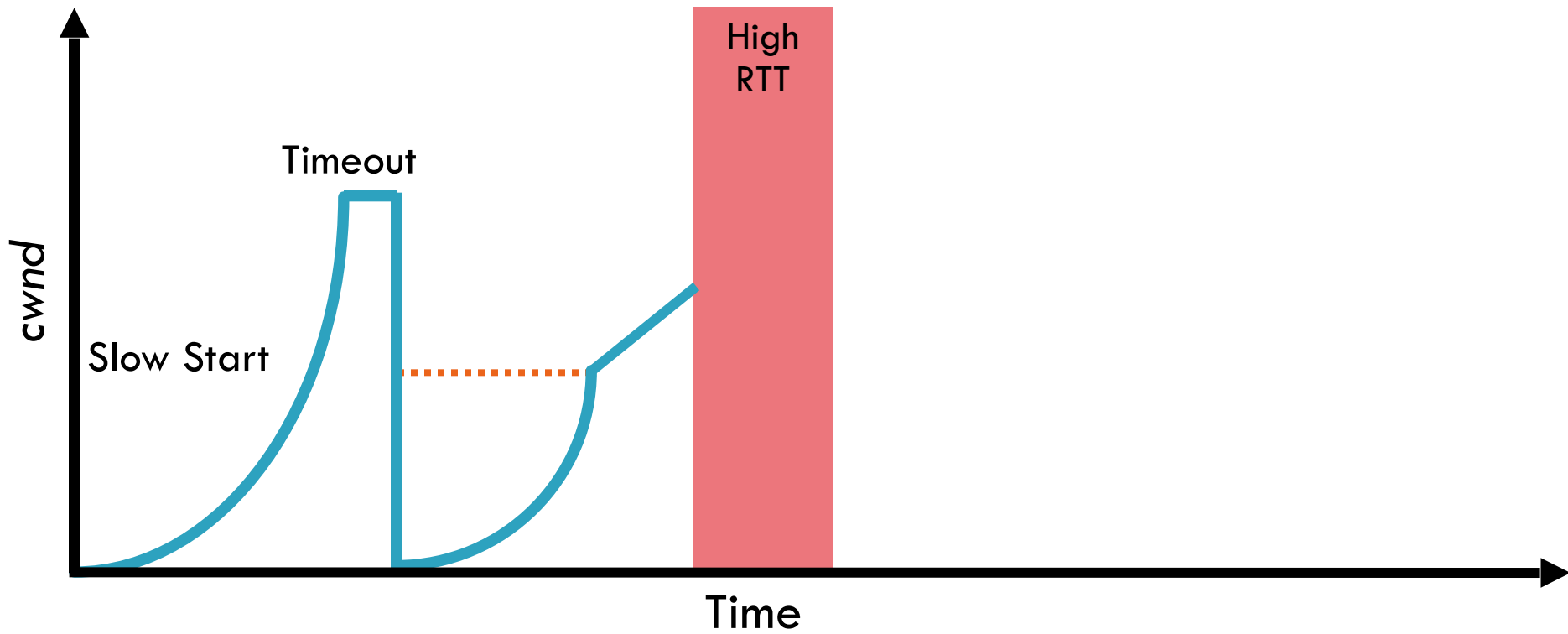
Compound TCP Example

38



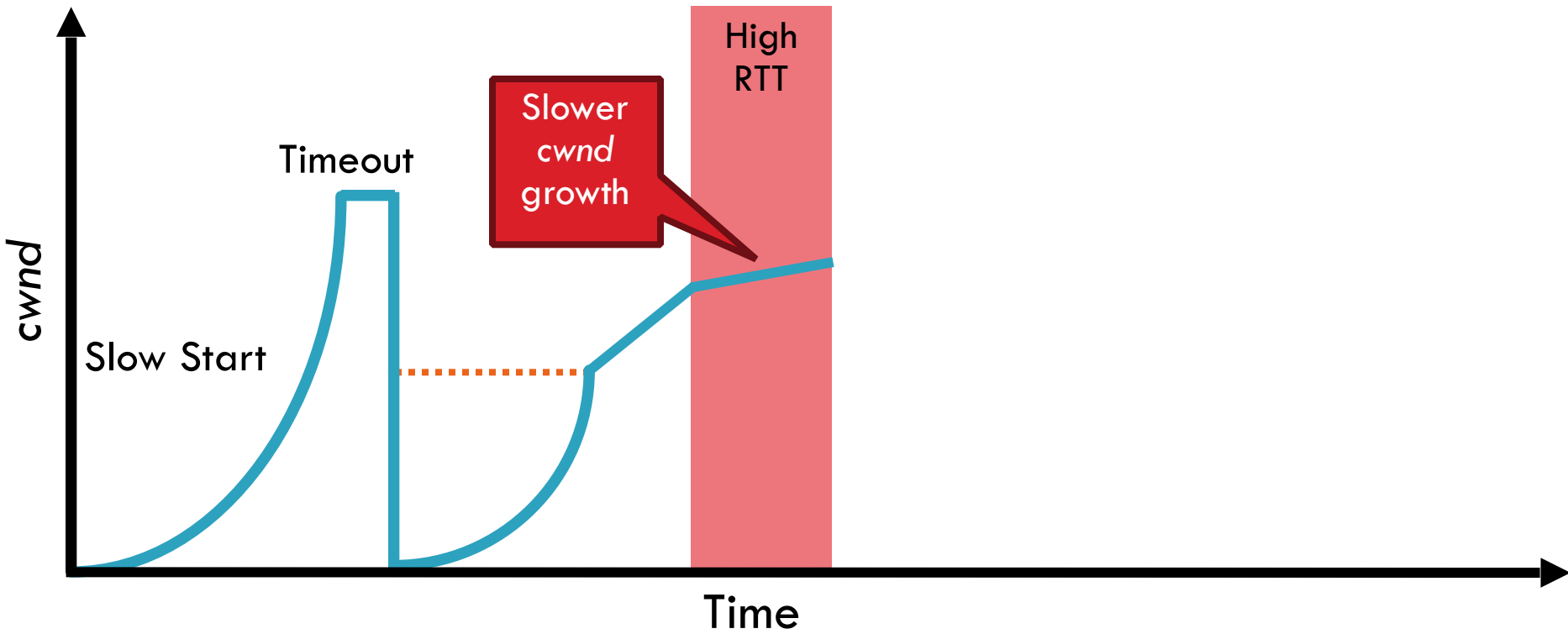
Compound TCP Example

38



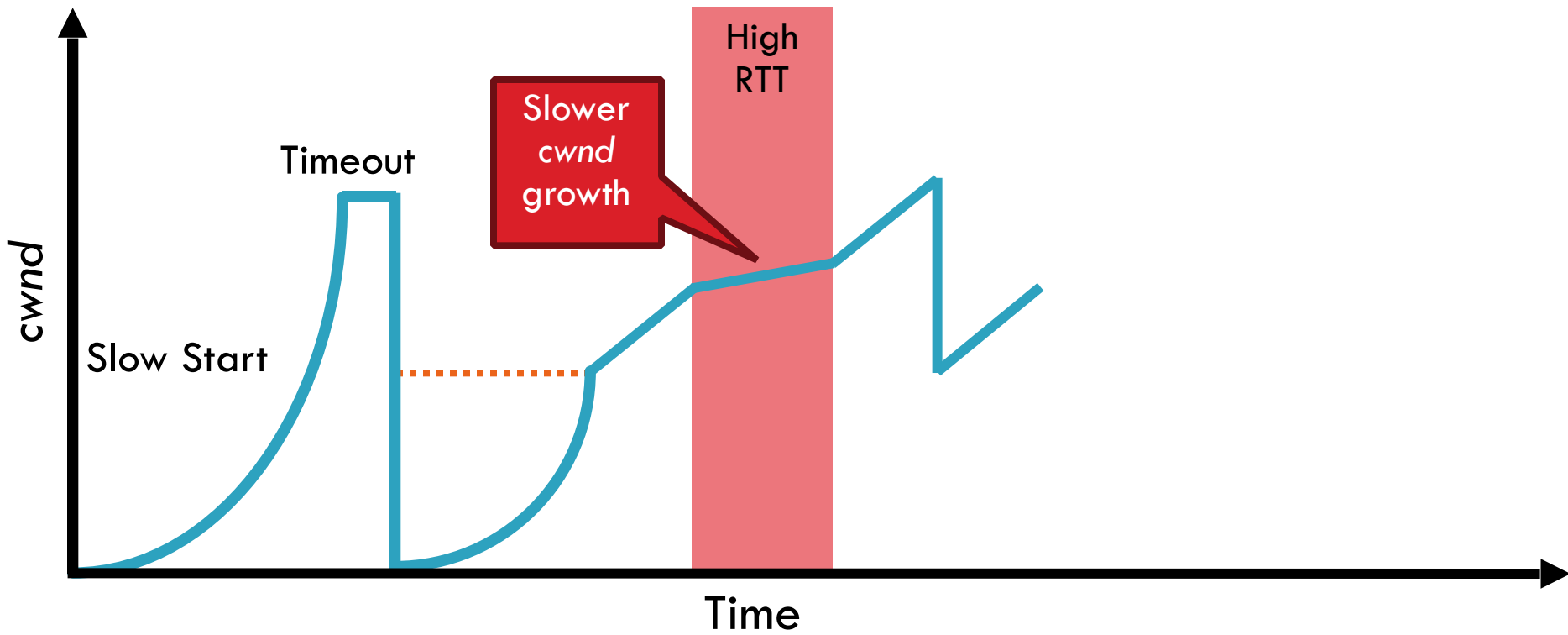
Compound TCP Example

38



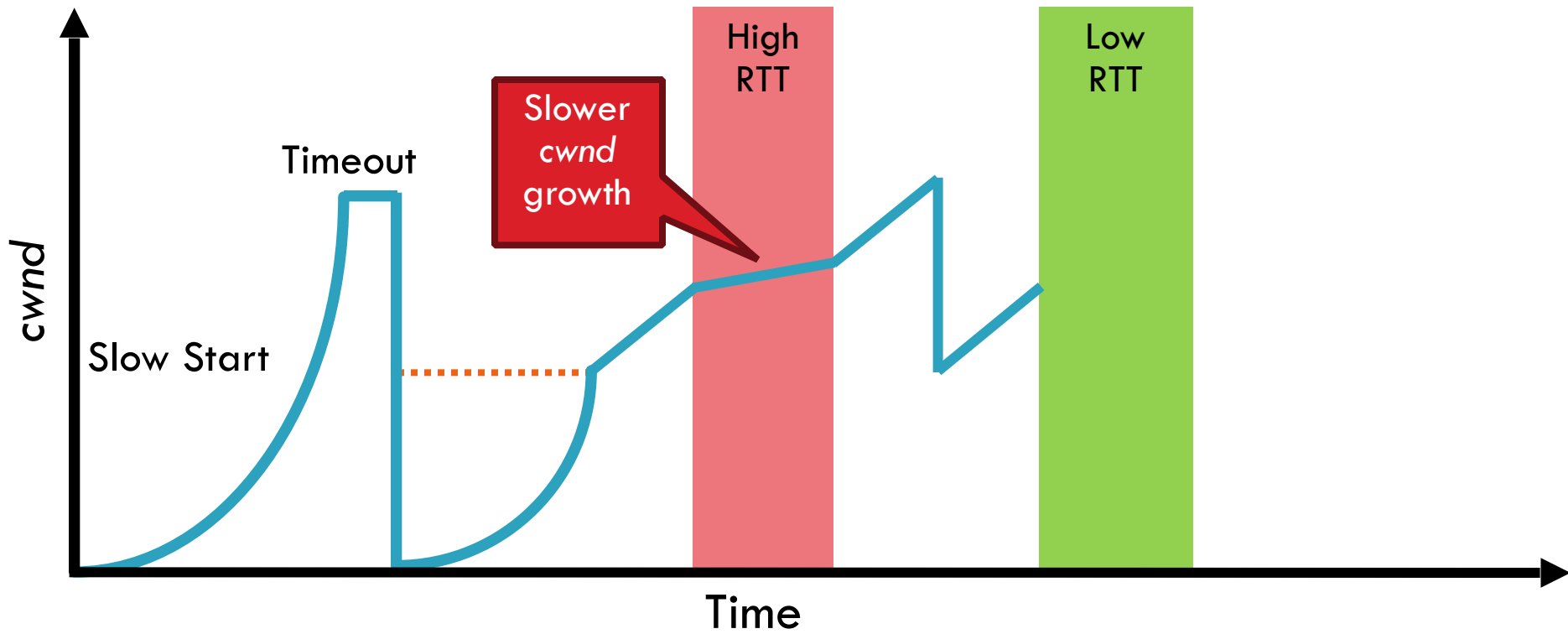
Compound TCP Example

38



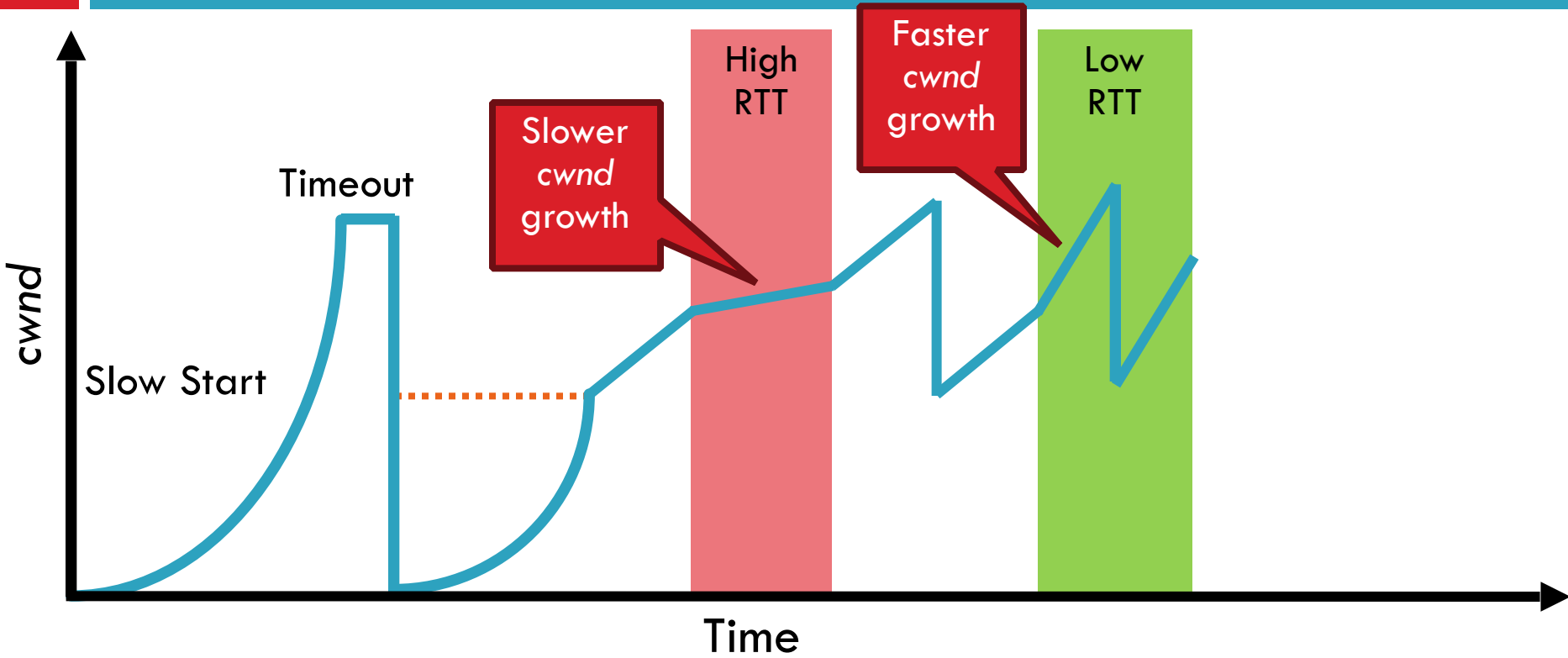
Compound TCP Example

38



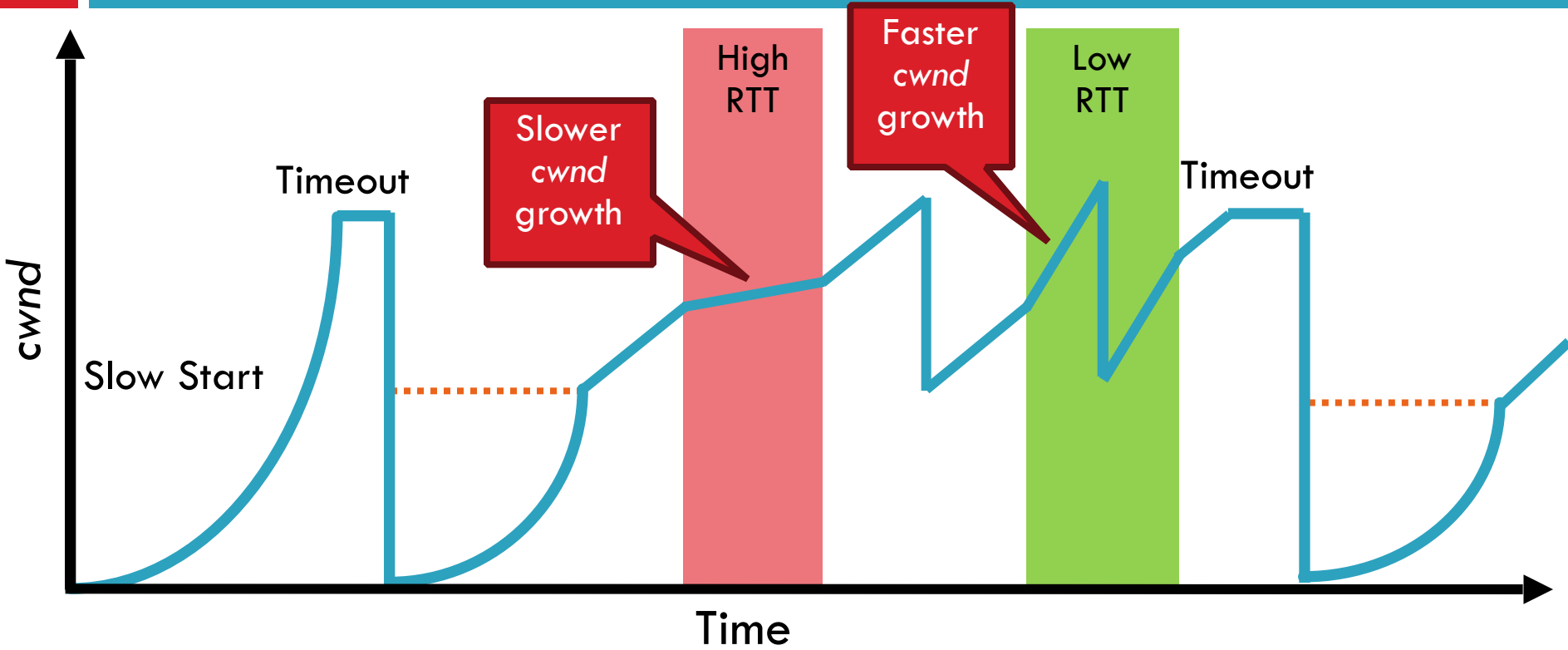
Compound TCP Example

38



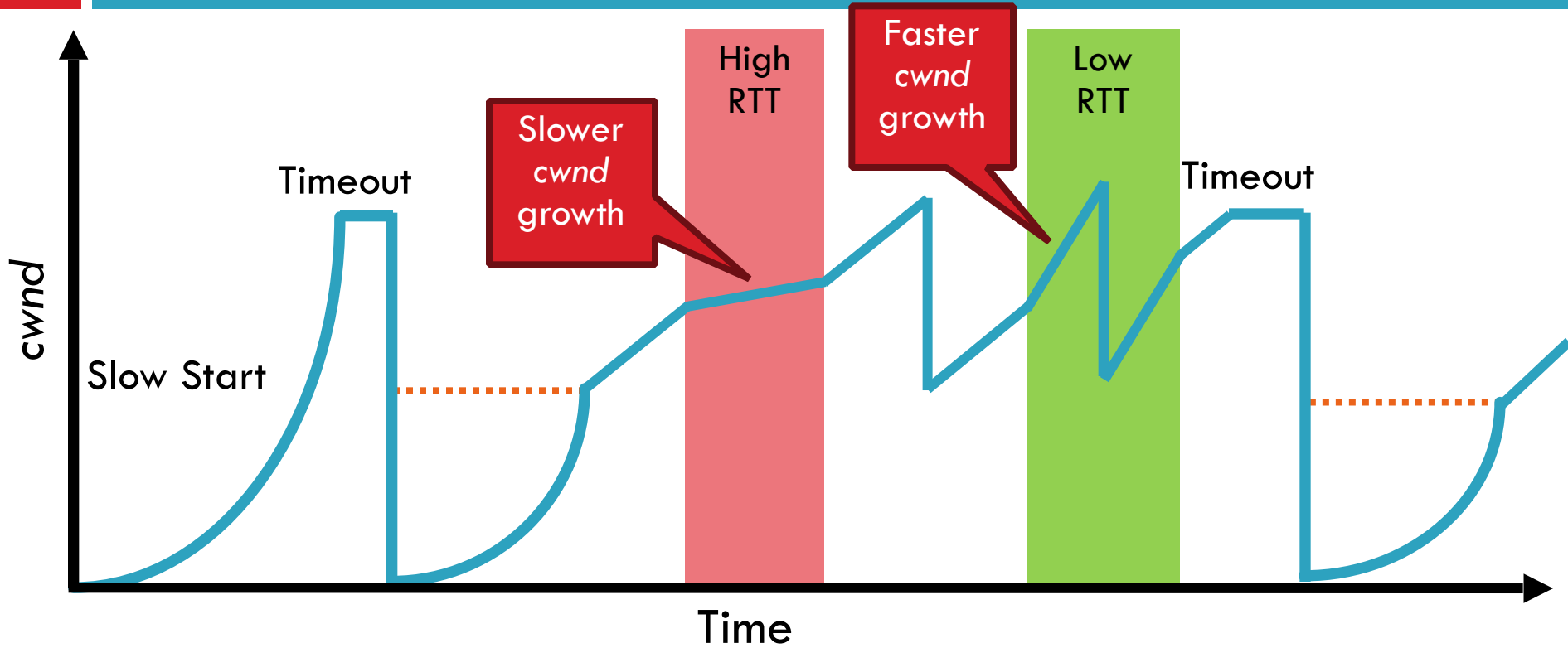
Compound TCP Example

38



Compound TCP Example

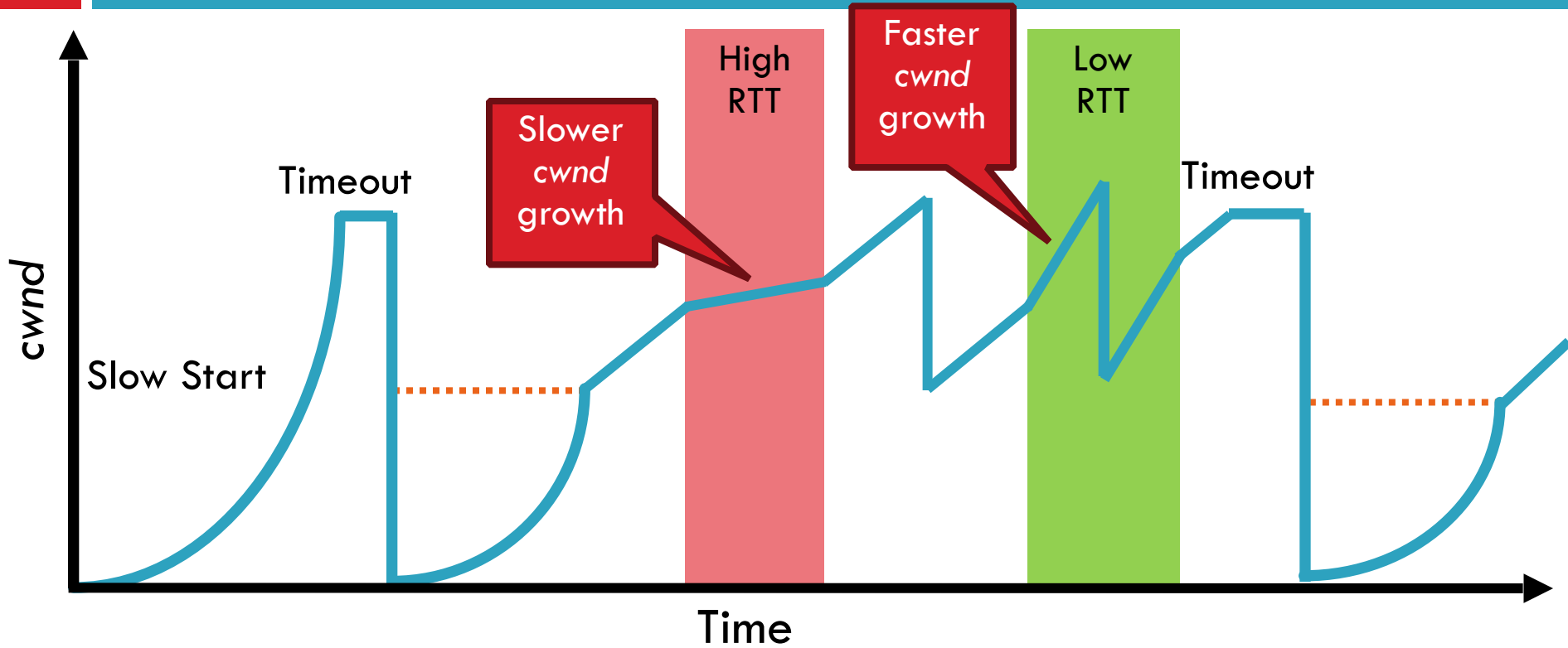
38



- Aggressiveness corresponds to changes in RTT

Compound TCP Example

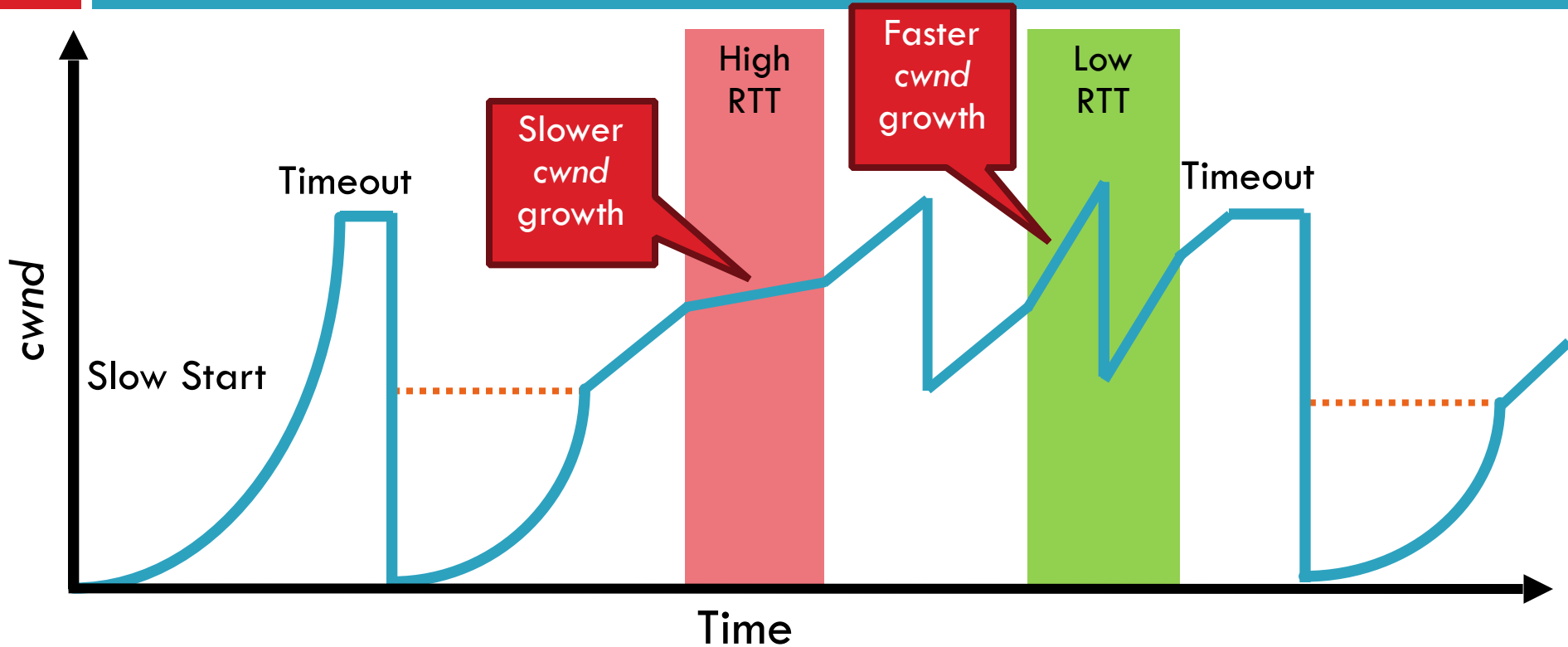
38



- Aggressiveness corresponds to changes in RTT
- Advantages: fast ramp up, more fair to flows with different RTTs

Compound TCP Example

38



- Aggressiveness corresponds to changes in RTT
- Advantages: fast ramp up, more fair to flows with different RTTs
- Disadvantage: must estimate RTT, which is very challenging

TCP CUBIC Implementation

39

- Default TCP implementation in Linux
- Replace AIMD with cubic function

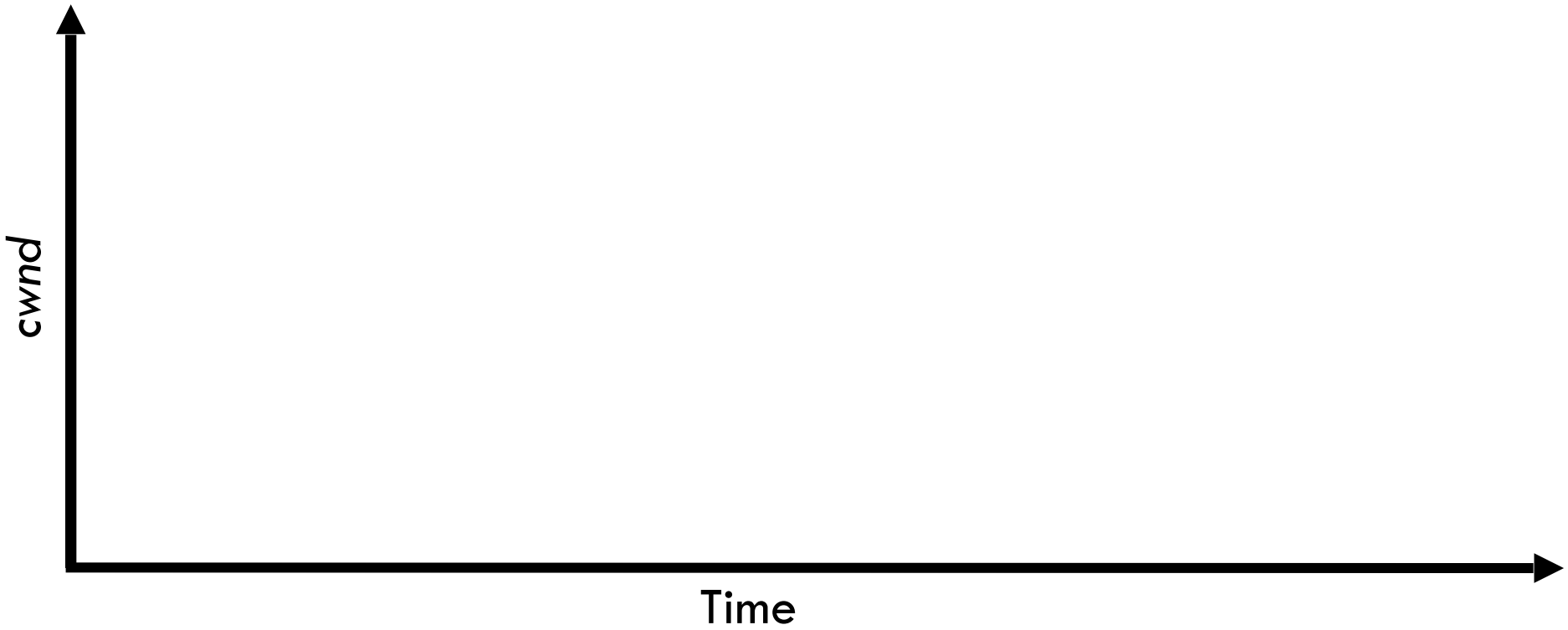
$$W_{cubic} = C(T - K)^3 + W_{max} \quad (1)$$

C is a scaling constant, and $K = \sqrt[3]{\frac{W_{max}\beta}{C}}$

- $\beta \rightarrow$ a constant fraction for multiplicative increase
- $T \rightarrow$ time since last packet drop
- $W_{max} \rightarrow$ cwnd when last packet dropped

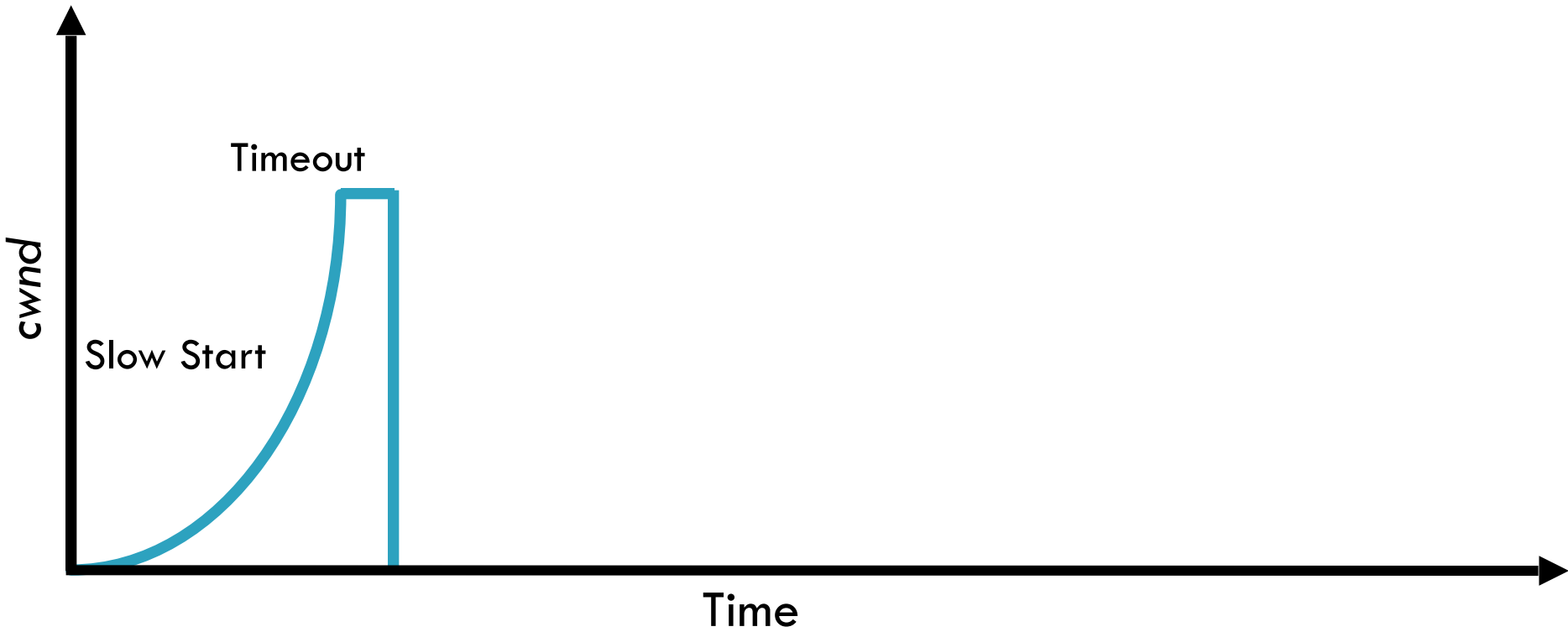
TCP CUBIC Example

40



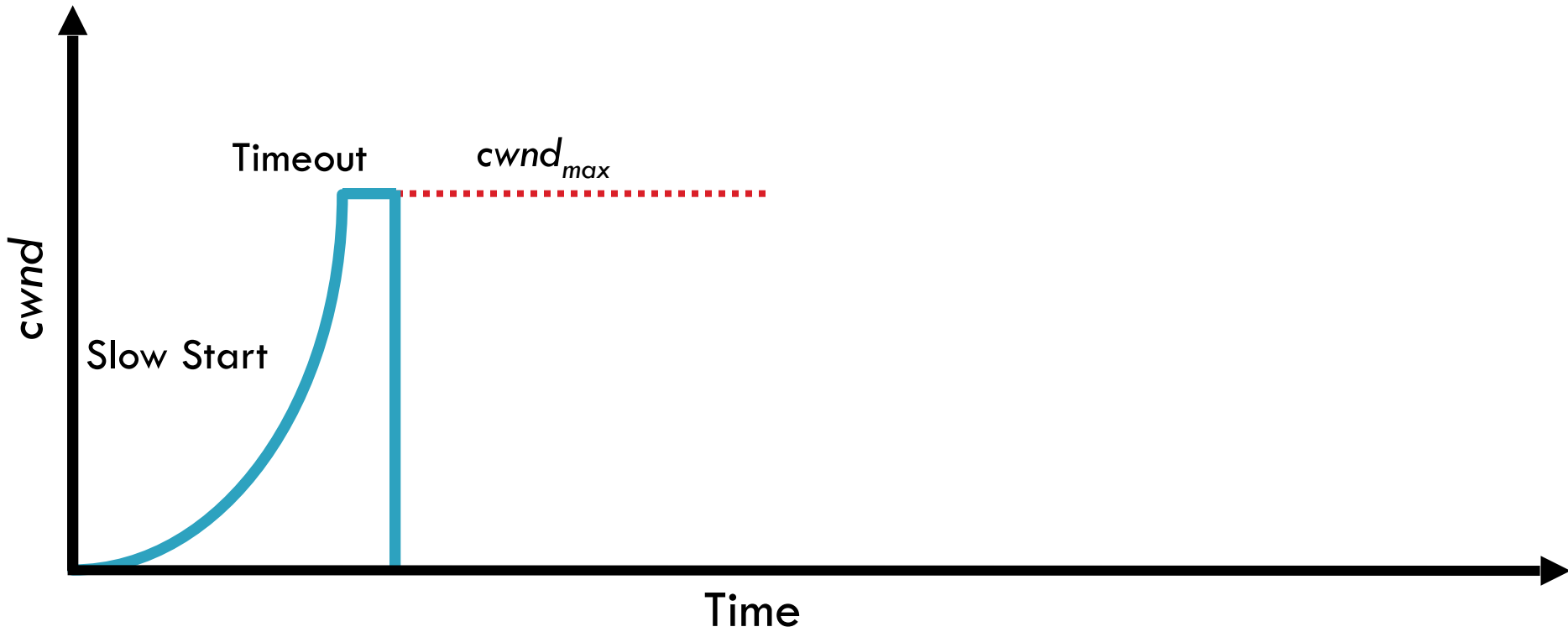
TCP CUBIC Example

40



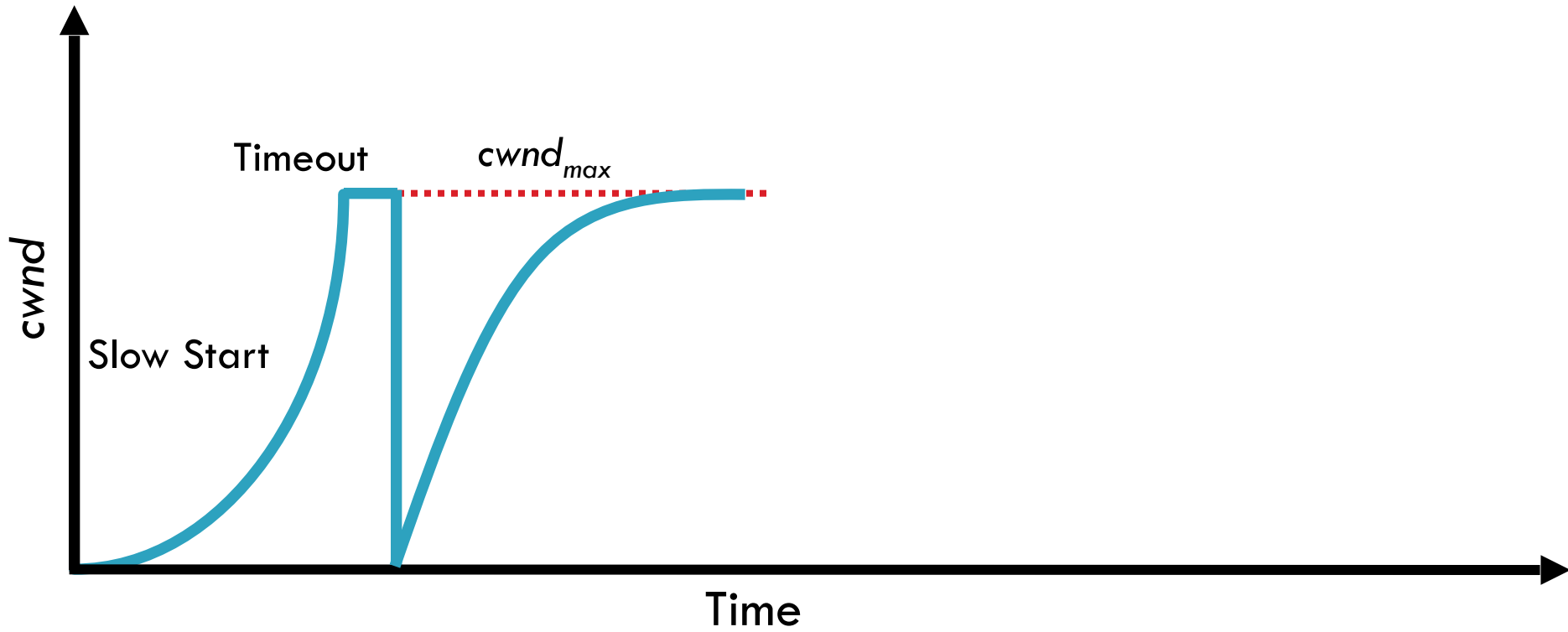
TCP CUBIC Example

40



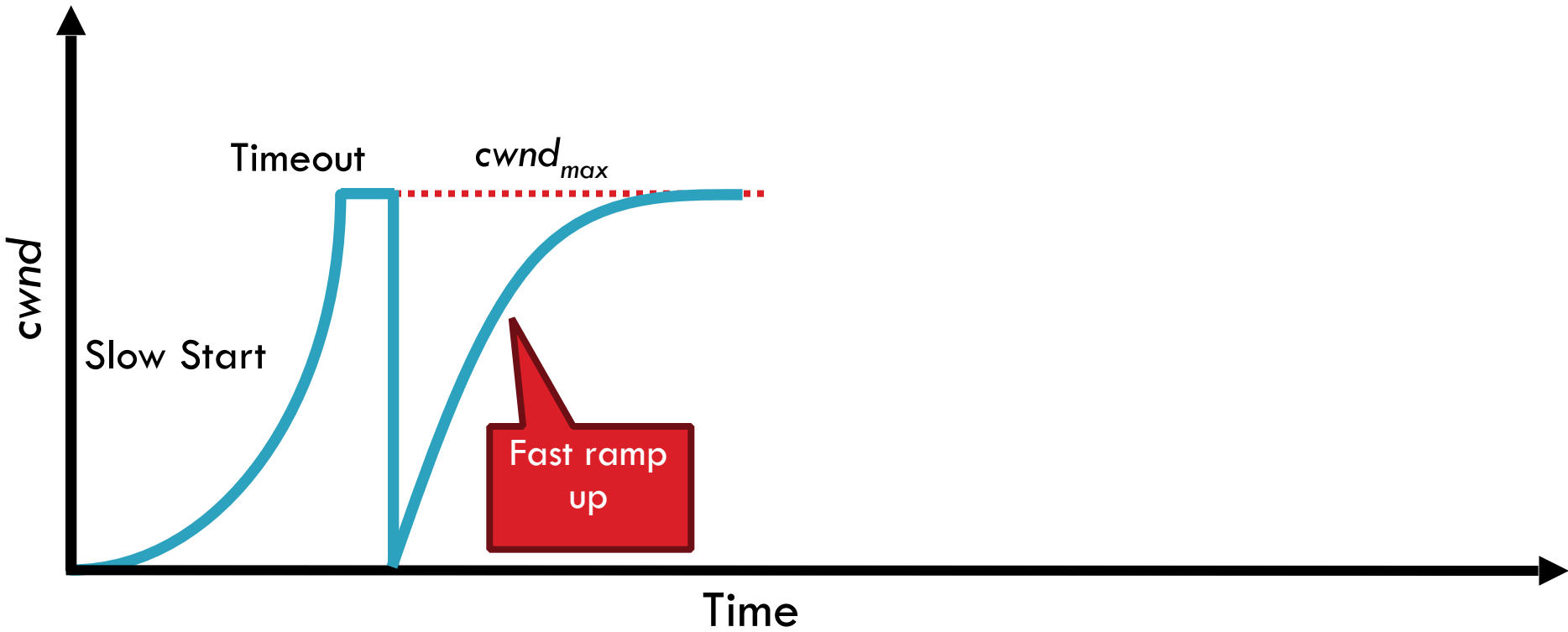
TCP CUBIC Example

40



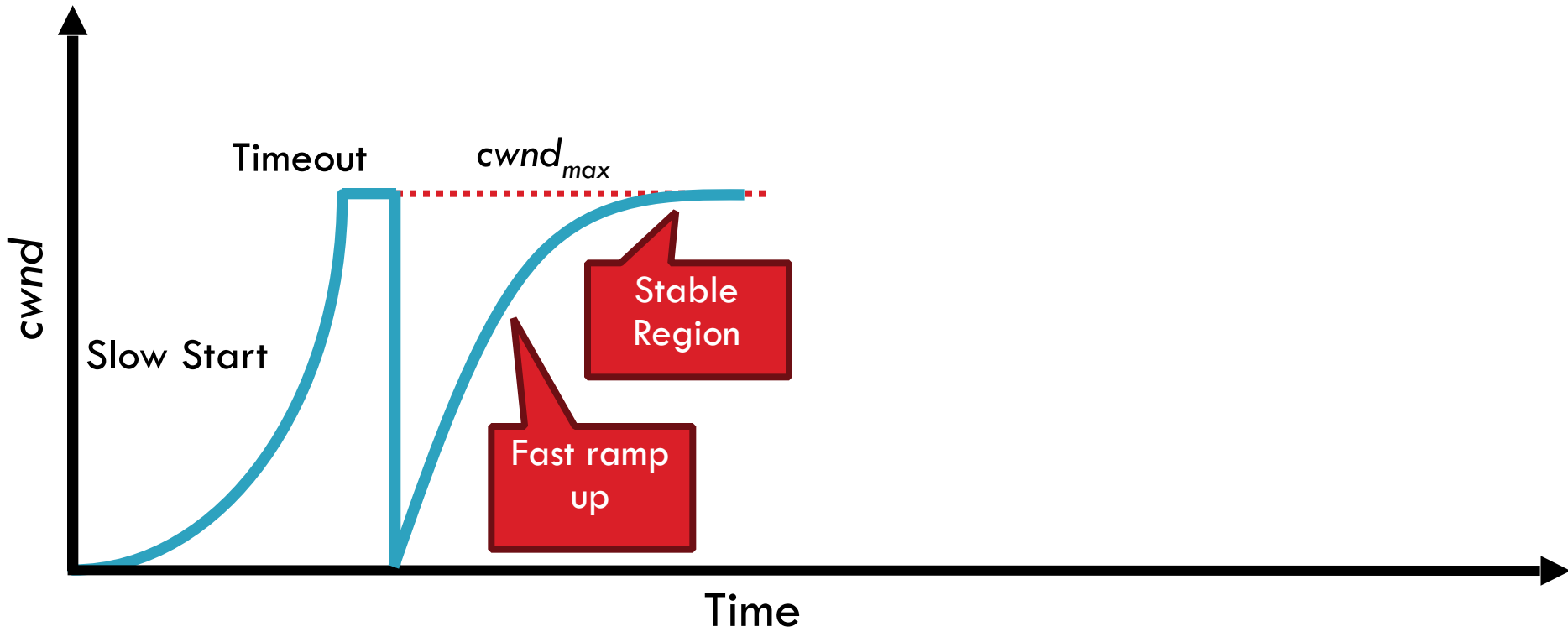
TCP CUBIC Example

40



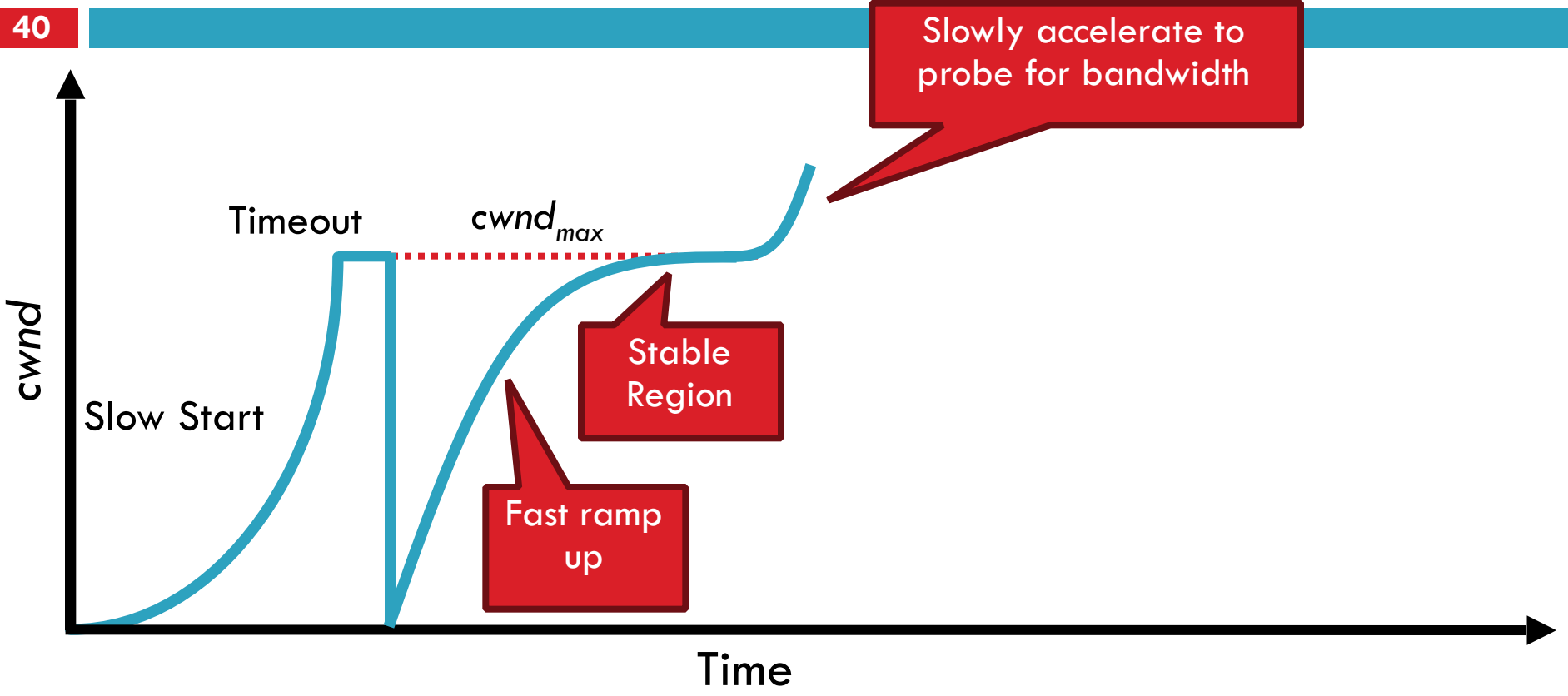
TCP CUBIC Example

40



TCP CUBIC Example

40



Slowly accelerate to probe for bandwidth

Timeout

$cwnd_{max}$

Slow Start

Stable Region

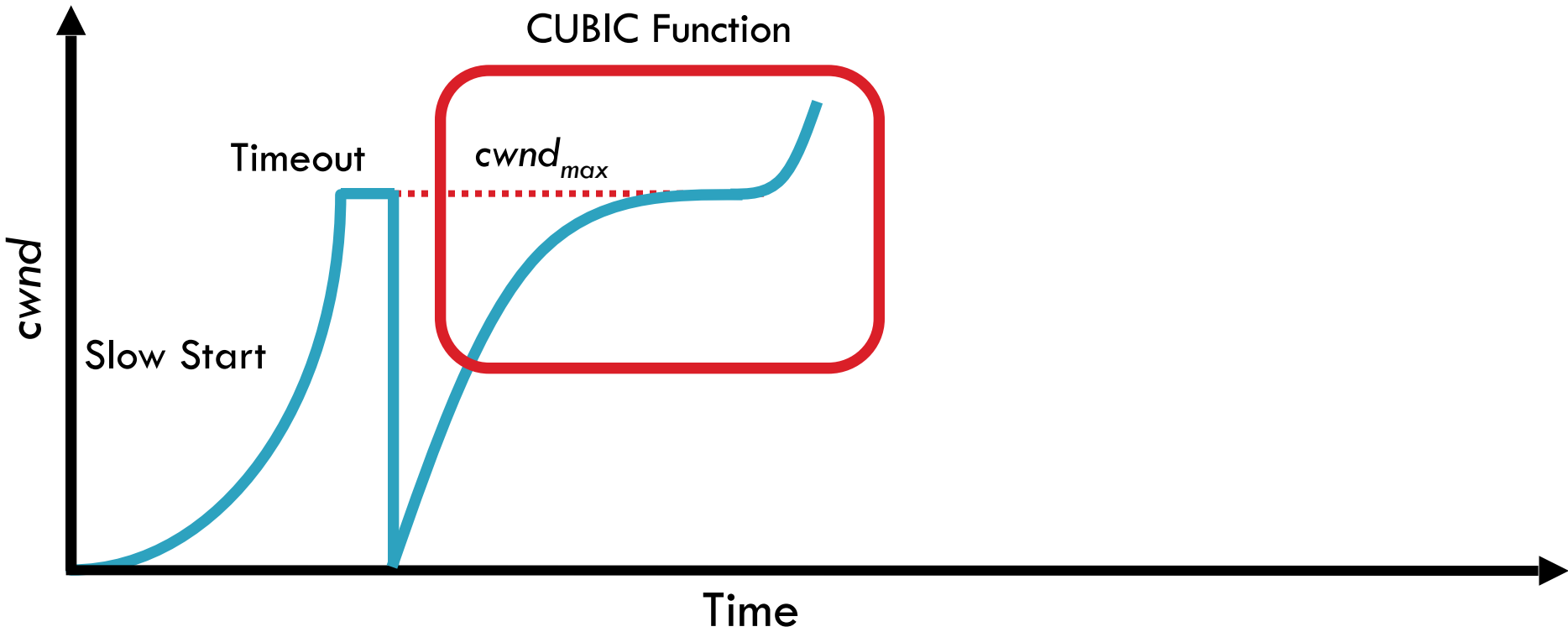
Fast ramp up

Time

$cwnd$

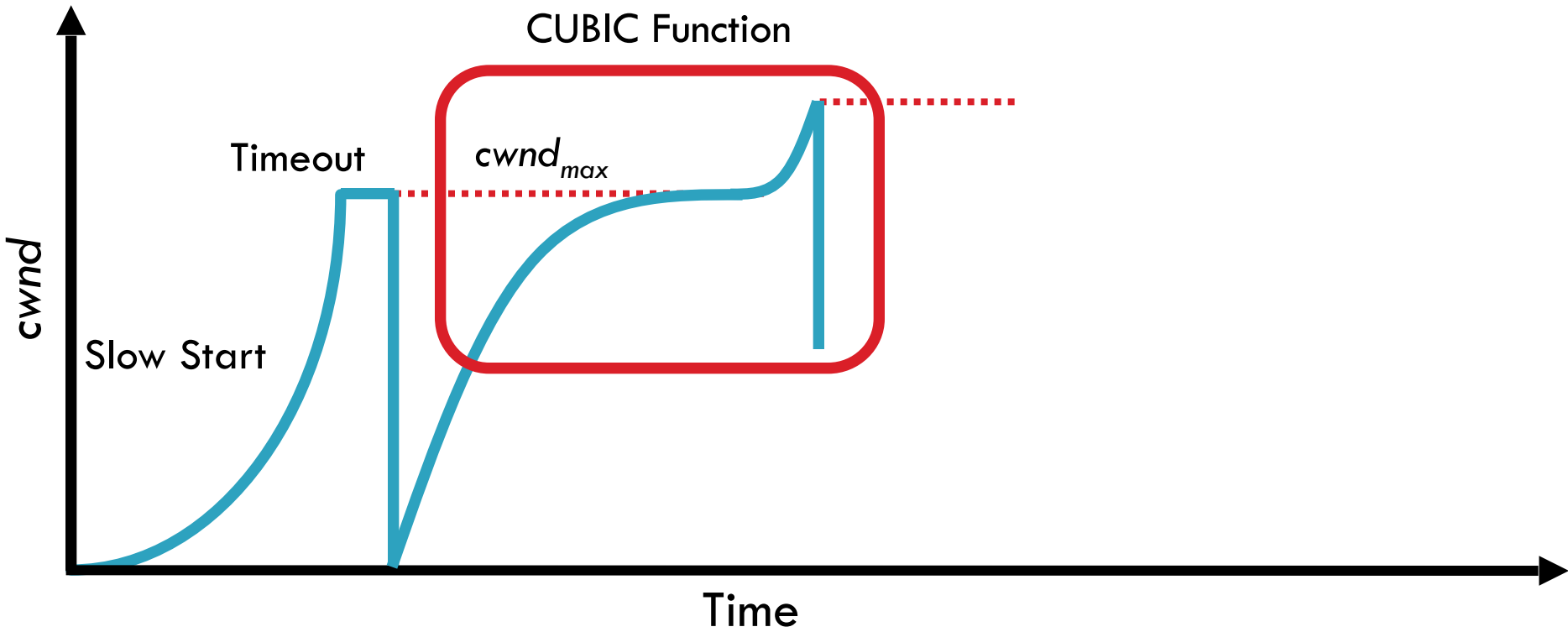
TCP CUBIC Example

40



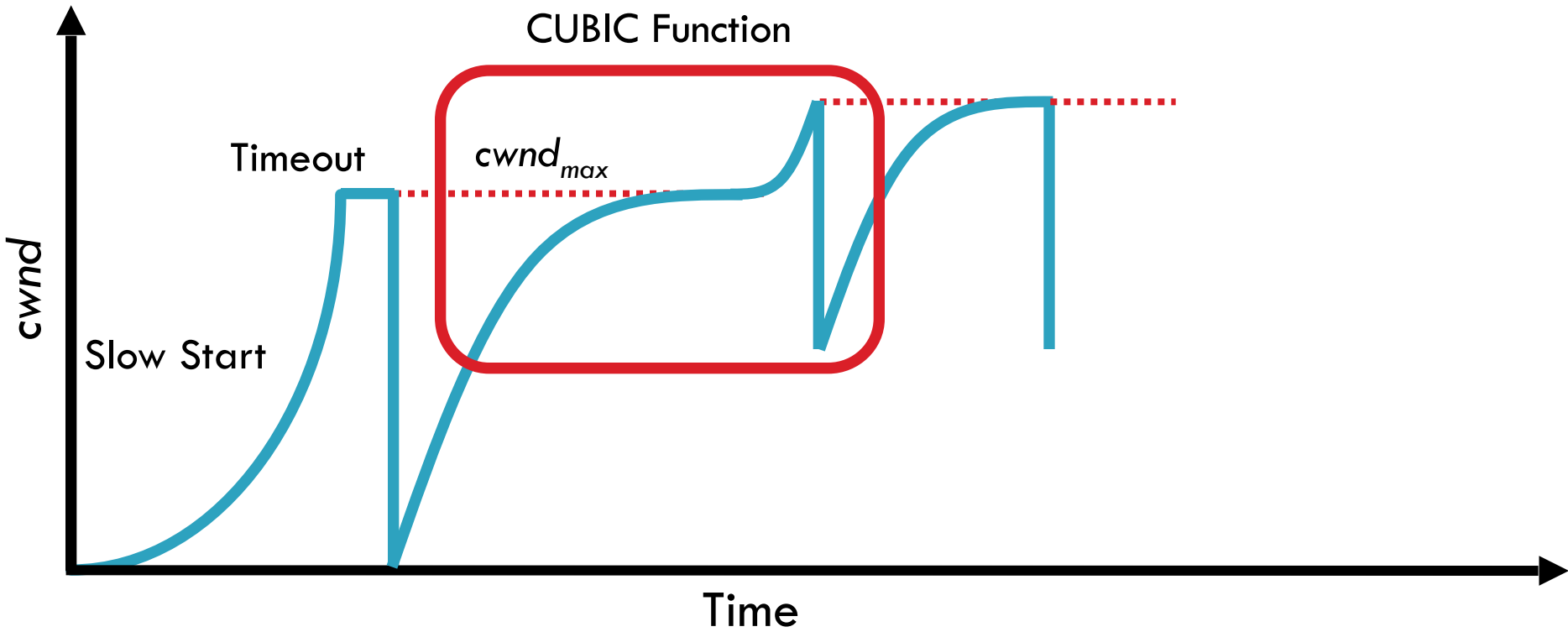
TCP CUBIC Example

40



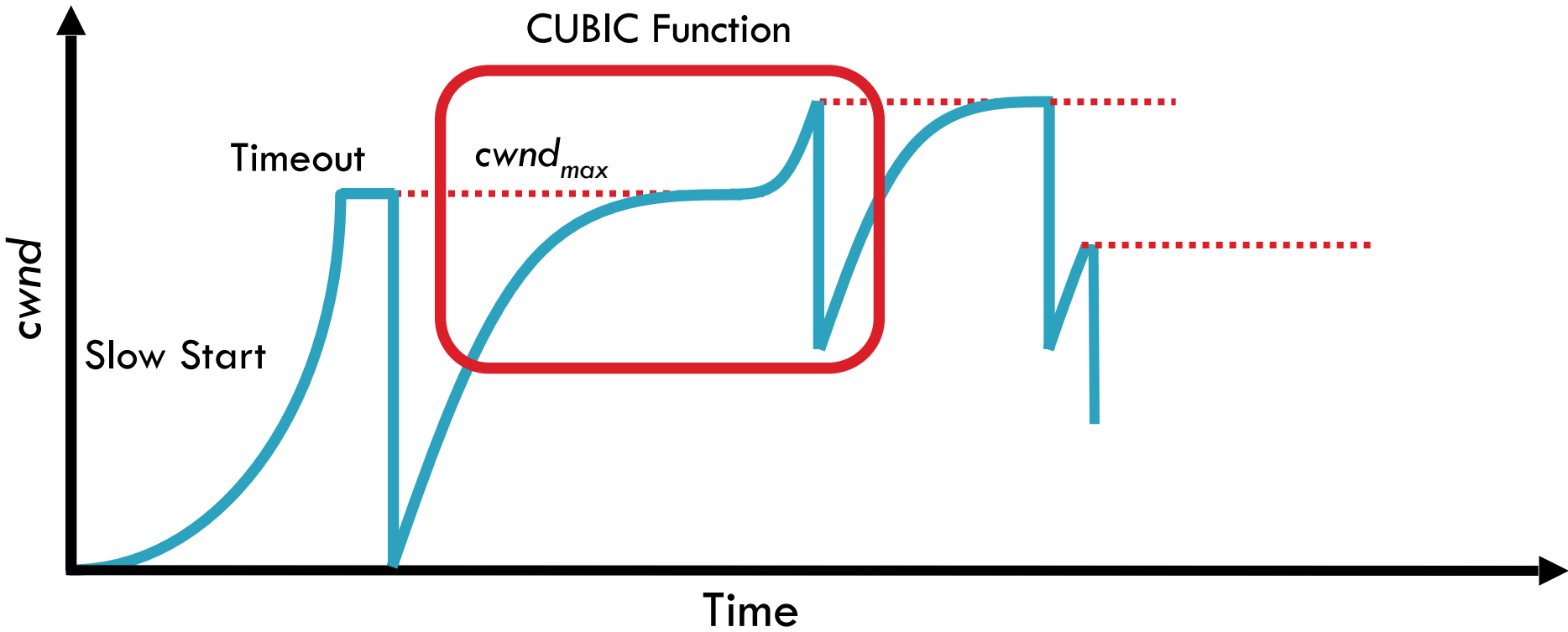
TCP CUBIC Example

40



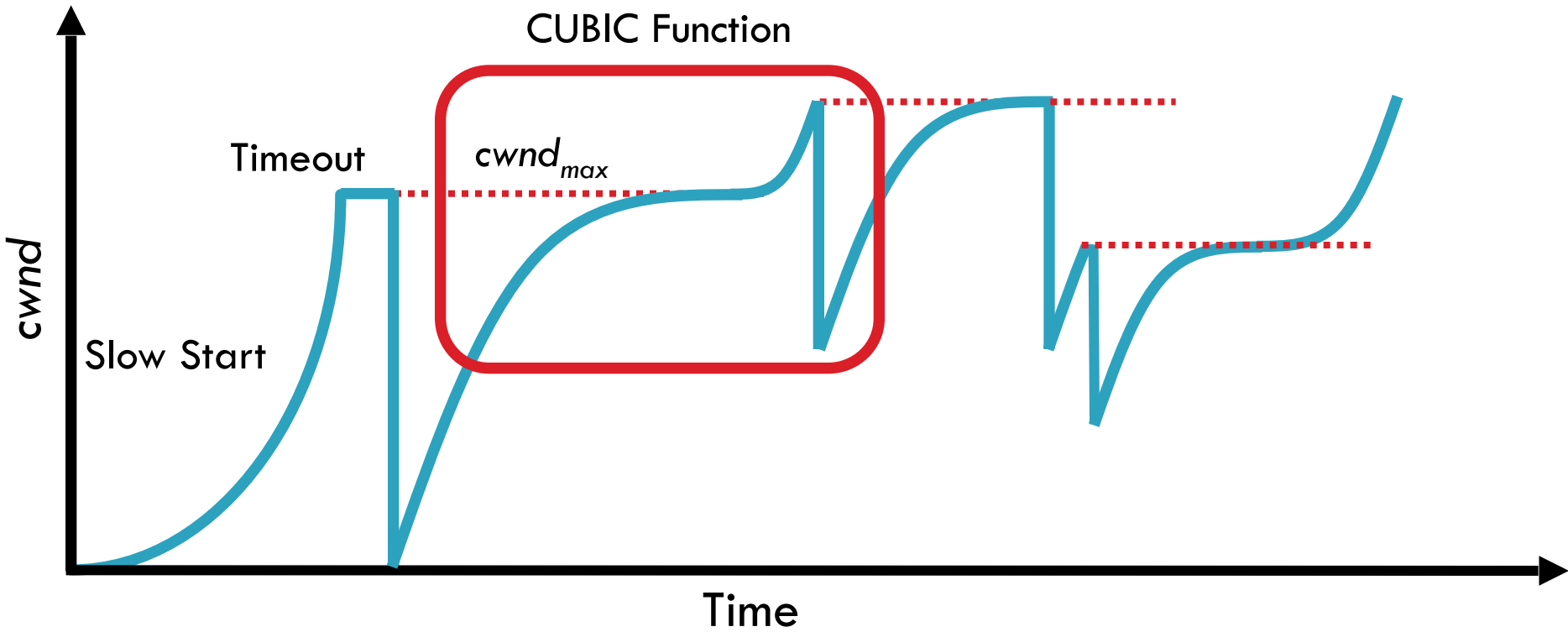
TCP CUBIC Example

40



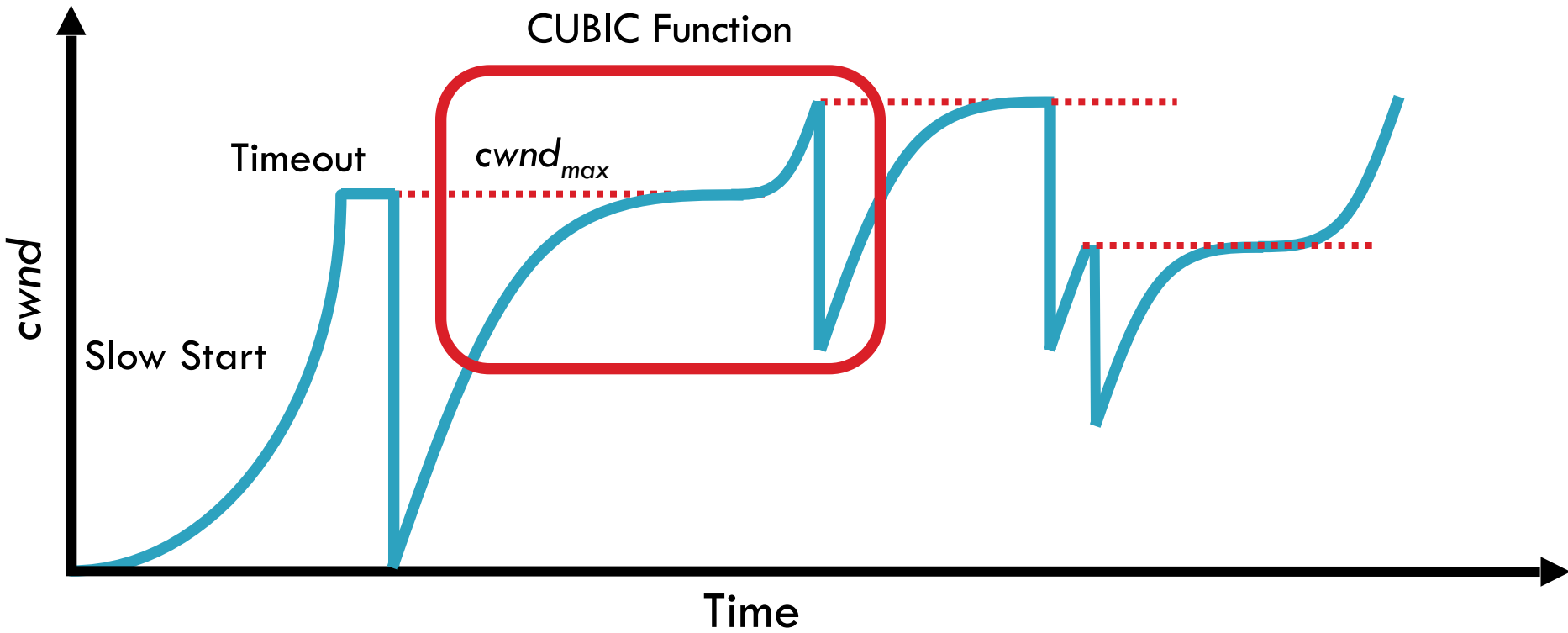
TCP CUBIC Example

40



TCP CUBIC Example

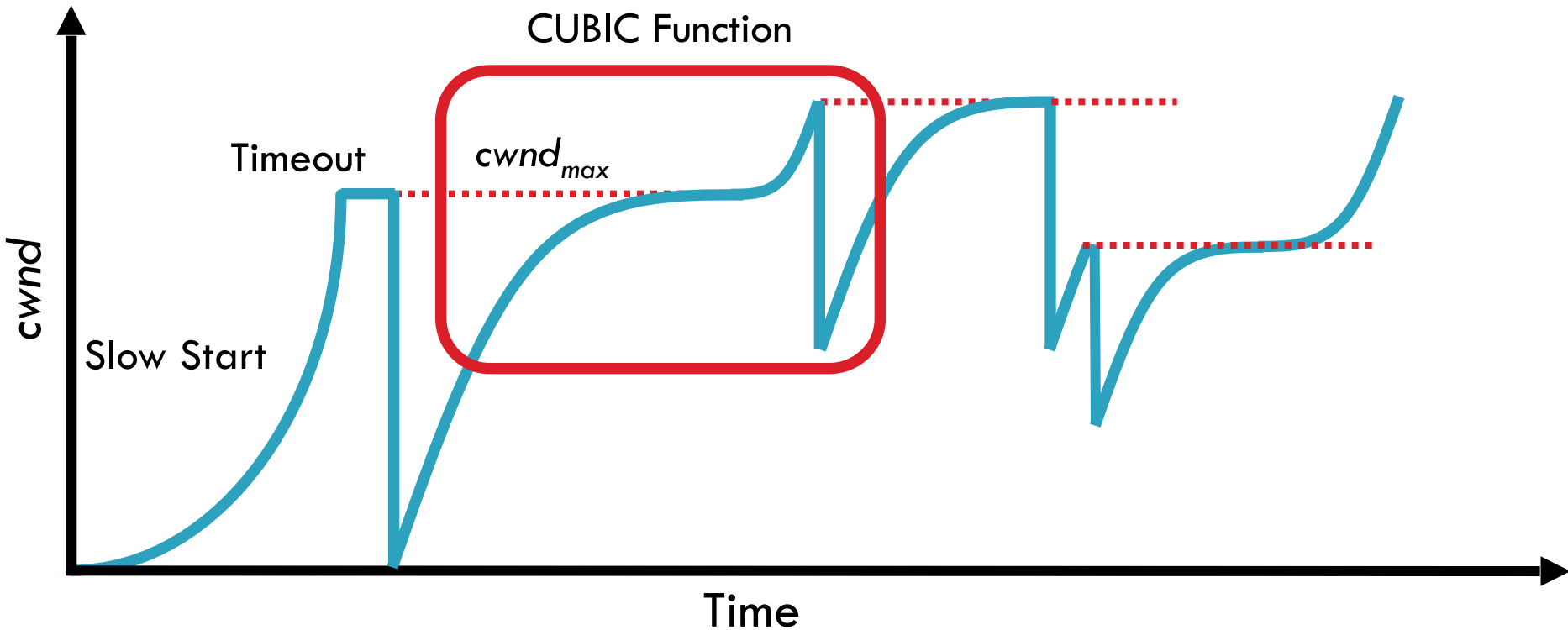
40



- Less wasted bandwidth due to fast ramp up

TCP CUBIC Example

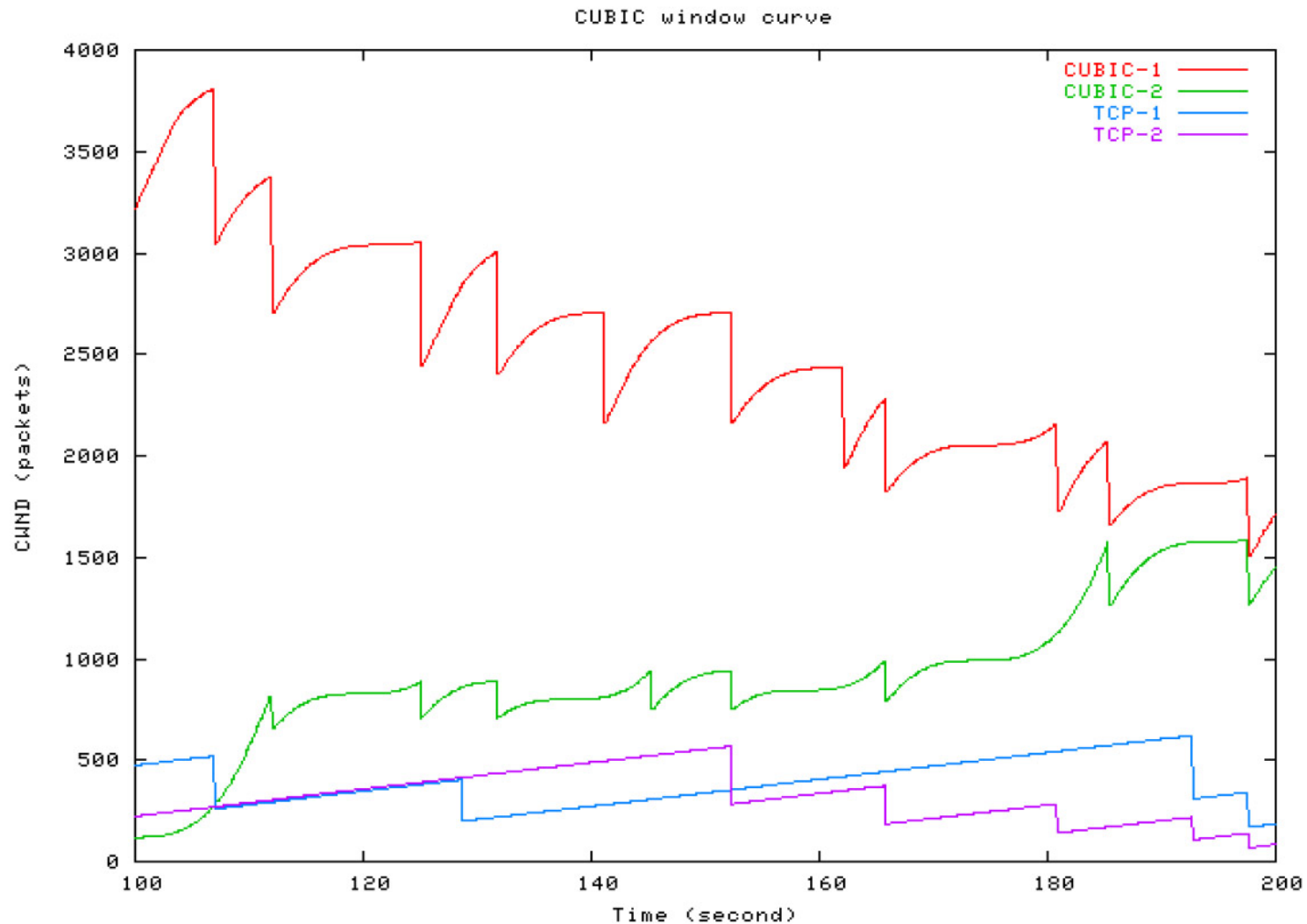
40



- Less wasted bandwidth due to fast ramp up
- Stable region and slow acceleration help maintain fairness
 - ▣ Fast ramp up is more aggressive than additive increase
 - ▣ To be fair to Tahoe/Reno, CUBIC needs to be less aggressive

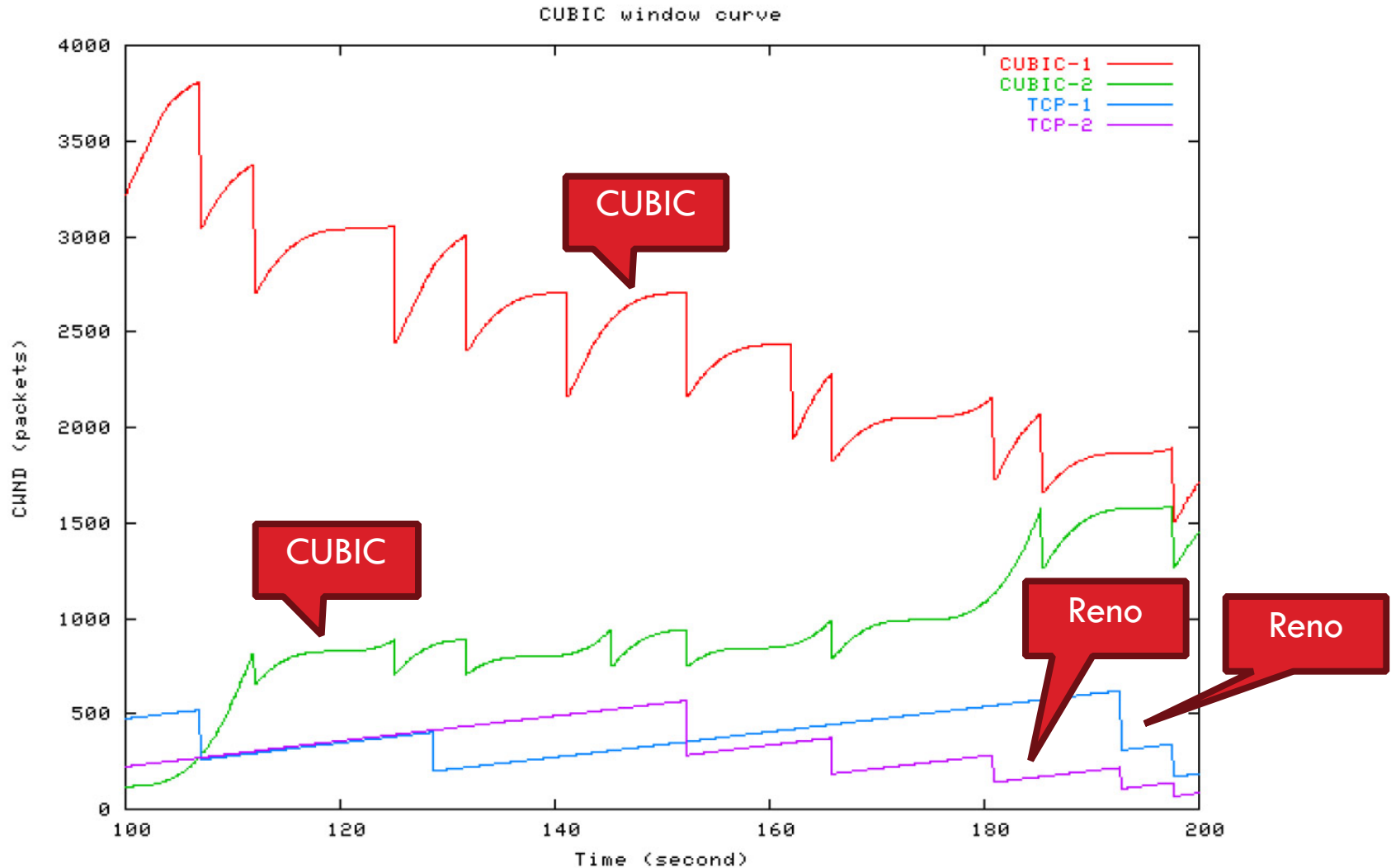
Simulations of CUBIC Flows

41



Simulations of CUBIC Flows

41



Deploying TCP Variants



- TCP assumes all flows employ TCP-like congestion control
 - ▣ TCP-friendly or TCP-compatible
 - ▣ Violated by UDP :(

Deploying TCP Variants



- TCP assumes all flows employ TCP-like congestion control
 - ▣ TCP-friendly or TCP-compatible
 - ▣ Violated by UDP :(
- If new congestion control algorithms are developed, they must be TCP-friendly

Deploying TCP Variants

- TCP assumes all flows employ TCP-like congestion control
 - ▣ TCP-friendly or TCP-compatible
 - ▣ Violated by UDP :(
- If new congestion control algorithms are developed, they must be TCP-friendly
- Be wary of unforeseen interactions
 - ▣ Variants work well with others like themselves
 - ▣ Different variants competing for resources may trigger unfair, pathological behavior

TCP Perspectives



- Cerf/Kahn
 - ▣ Provide flow control
 - ▣ Congestion handled by retransmission

TCP Perspectives



- Cerf/Kahn
 - ▣ Provide flow control
 - ▣ Congestion handled by retransmission
- Jacobson / Karels
 - ▣ Need to avoid congestion
 - ▣ RTT estimates critical
 - ▣ Queuing theory can help

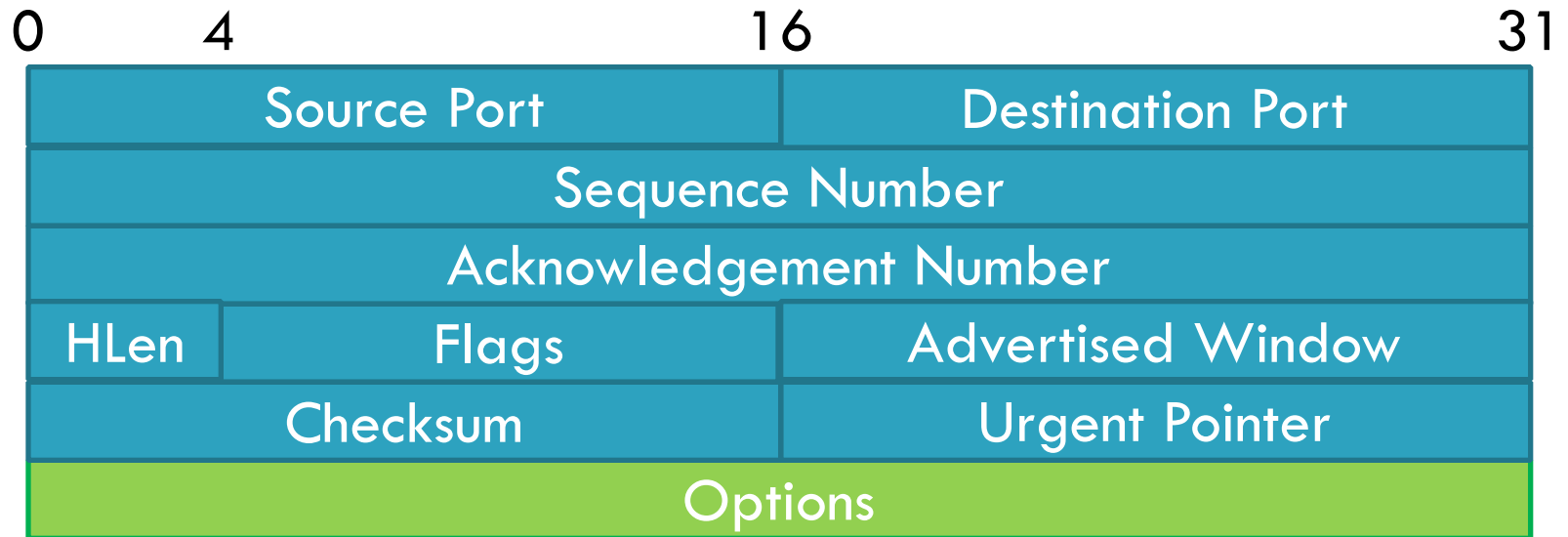
TCP Perspectives

- Cerf/Kahn
 - ▣ Provide flow control
 - ▣ Congestion handled by retransmission
- Jacobson / Karels
 - ▣ Need to avoid congestion
 - ▣ RTT estimates critical
 - ▣ Queuing theory can help
- Winstein/Balakrishnan
 - ▣ TCP is maximizing an objective function
 - Fairness/efficiency
 - Throughput/delay
 - ▣ Let a machine pick the best fit for your environment

- ❑ Congestion Control
- ❑ Evolution of TCP
- ❑ Problems with TCP

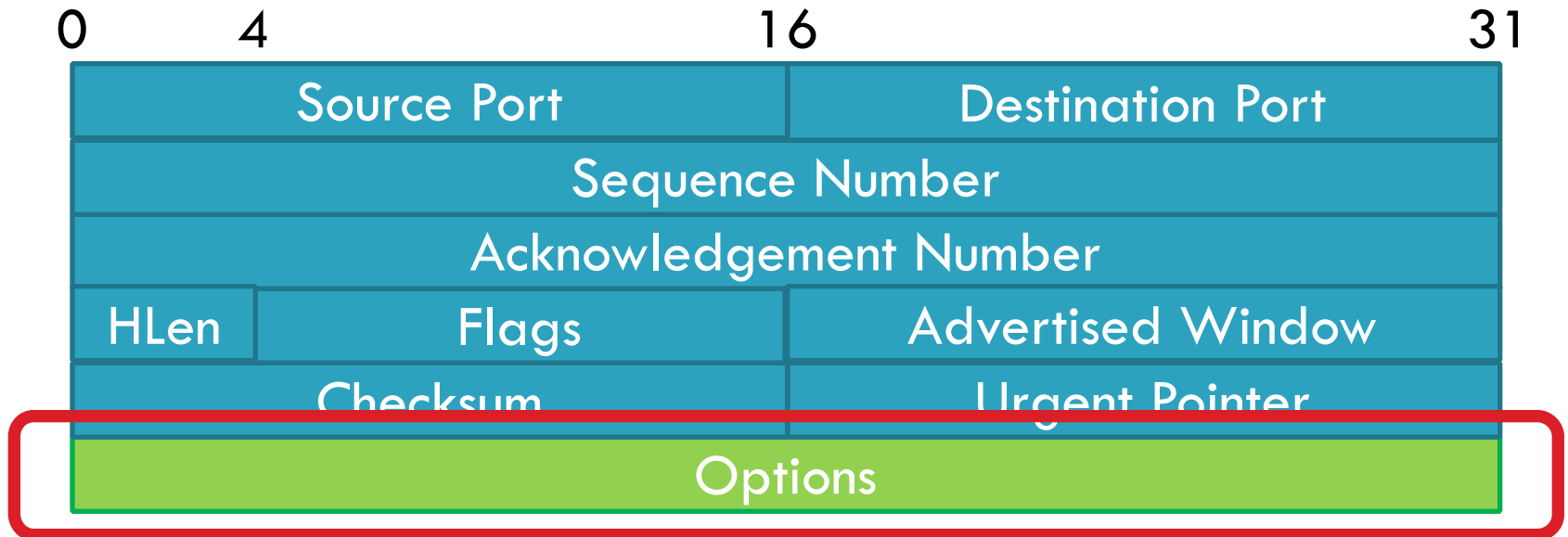
Common TCP Options

45



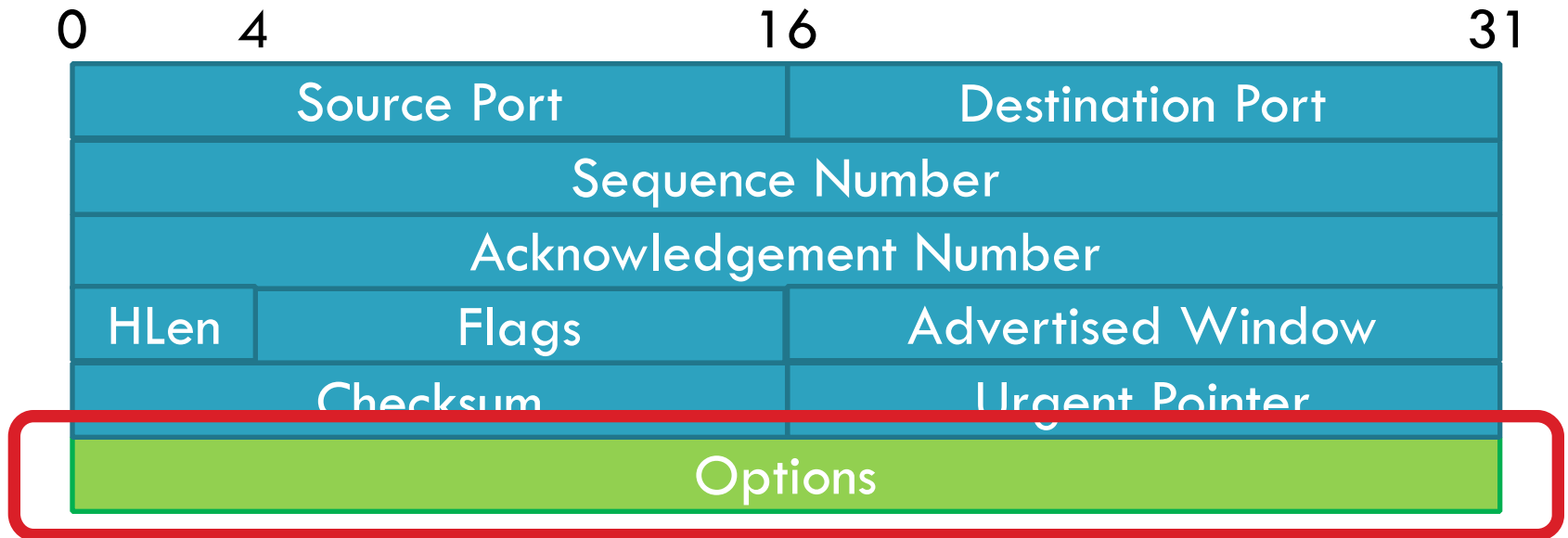
Common TCP Options

45



Common TCP Options

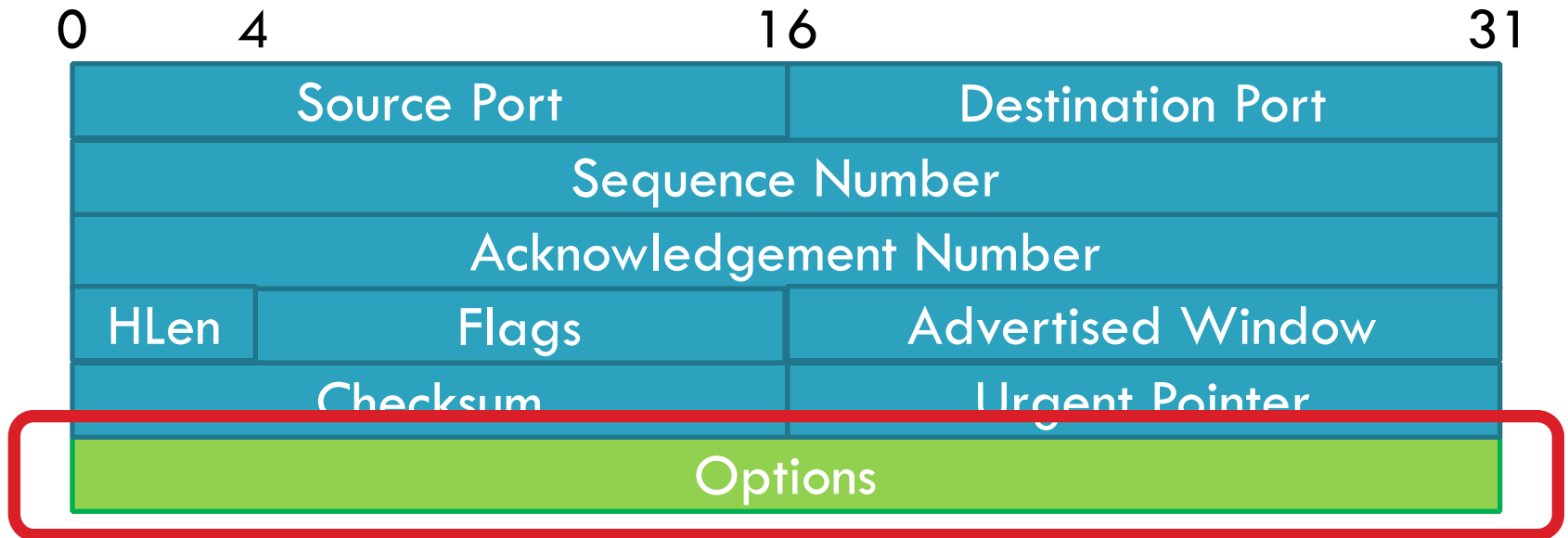
45



- Window scaling

Common TCP Options

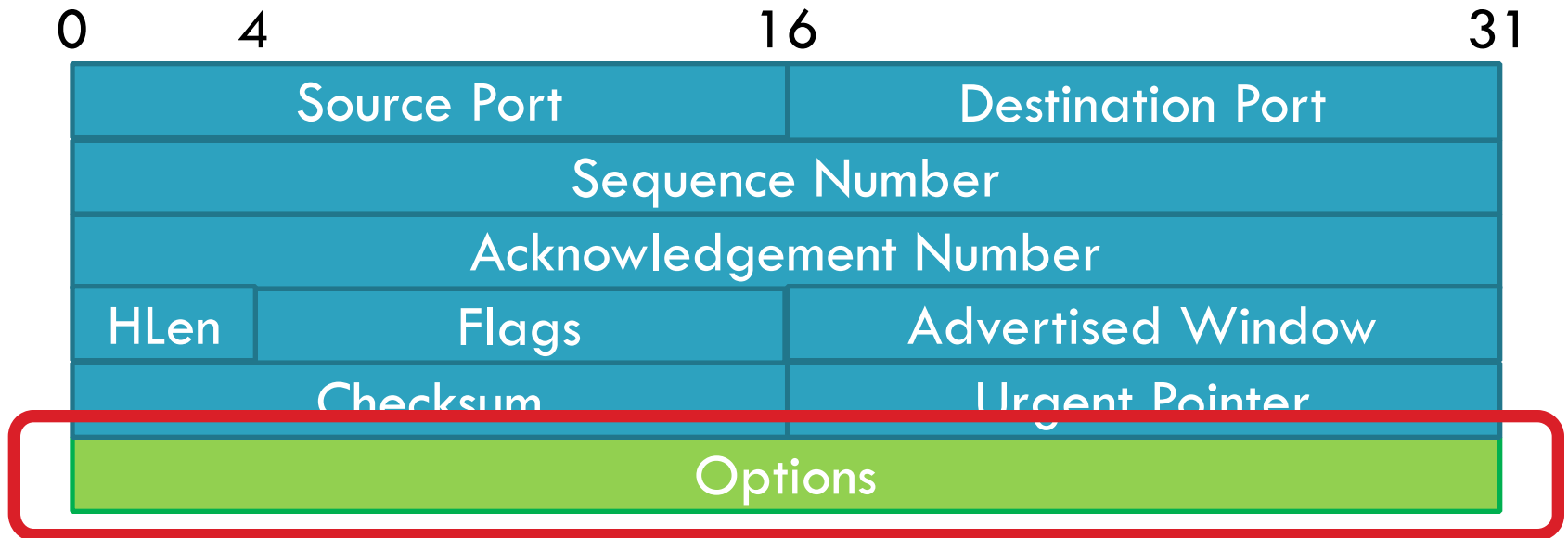
45



- Window scaling
- SACK: selective acknowledgement

Common TCP Options

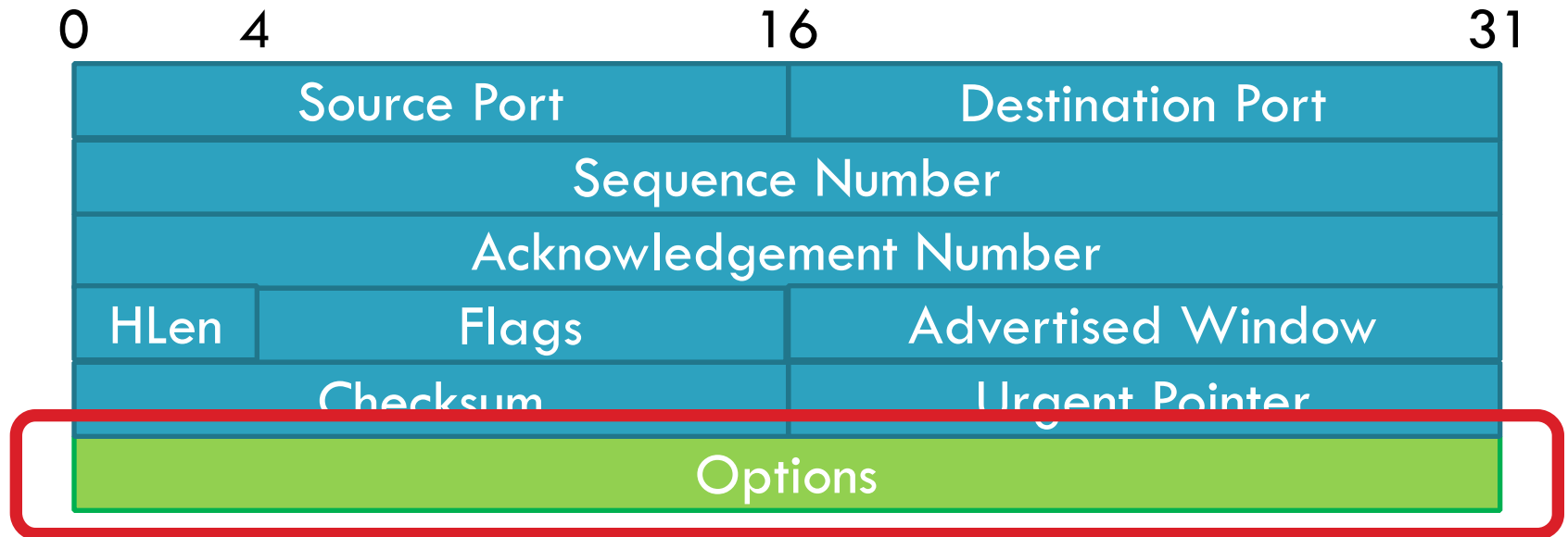
45



- Window scaling
- SACK: selective acknowledgement
- Maximum segment size (MSS)

Common TCP Options

45



- Window scaling
- SACK: selective acknowledgement
- Maximum segment size (MSS)
- Timestamp

Window Scaling

46

- Problem: the advertised window is only 16-bits
 - Effectively caps the window at 65536B, 64KB
 - Example: 1.5Mbps link, 513ms RTT

Window Scaling

46

- Problem: the advertised window is only 16-bits
 - Effectively caps the window at 65536B, 64KB
 - Example: 1.5Mbps link, 513ms RTT

$$(1.5\text{Mbps} * 0.513\text{s}) = 94\text{KB}$$

$$64\text{KB} / 94\text{KB} = 68\% \text{ of maximum possible speed}$$

Window Scaling

46

- Problem: the advertised window is only 16-bits
 - Effectively caps the window at 65536B, 64KB
 - Example: 1.5Mbps link, 513ms RTT

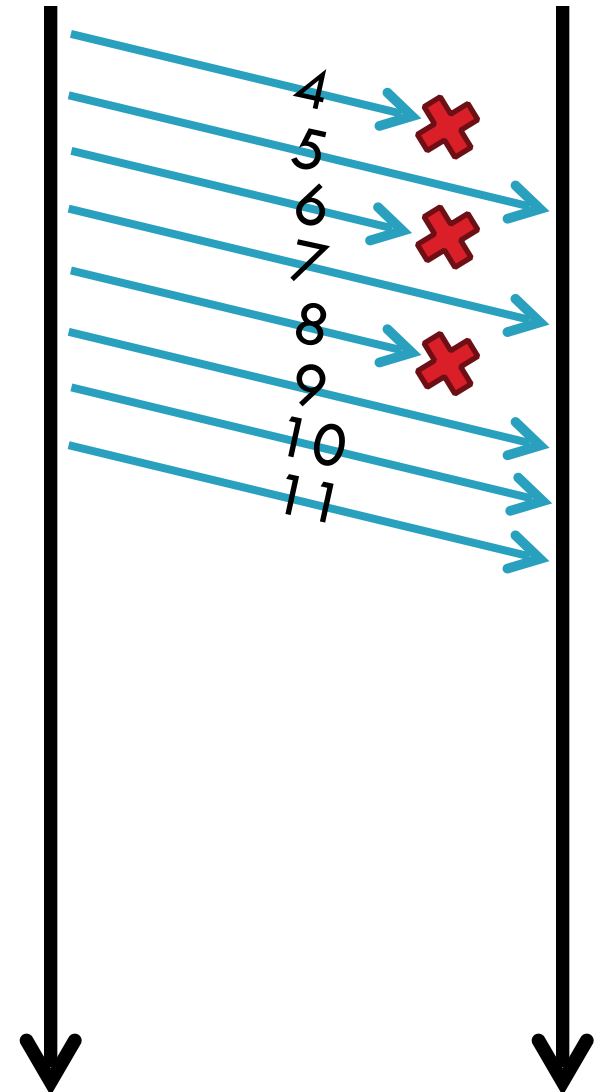
$$(1.5\text{Mbps} * 0.513\text{s}) = 94\text{KB}$$

$$64\text{KB} / 94\text{KB} = 68\% \text{ of maximum possible speed}$$

- Solution: introduce a window scaling value
 - $wnd = adv_wnd \ll wnd_scale;$
 - Maximum shift is 14 bits, 1GB maximum window

SACK: Selective Acknowledgment

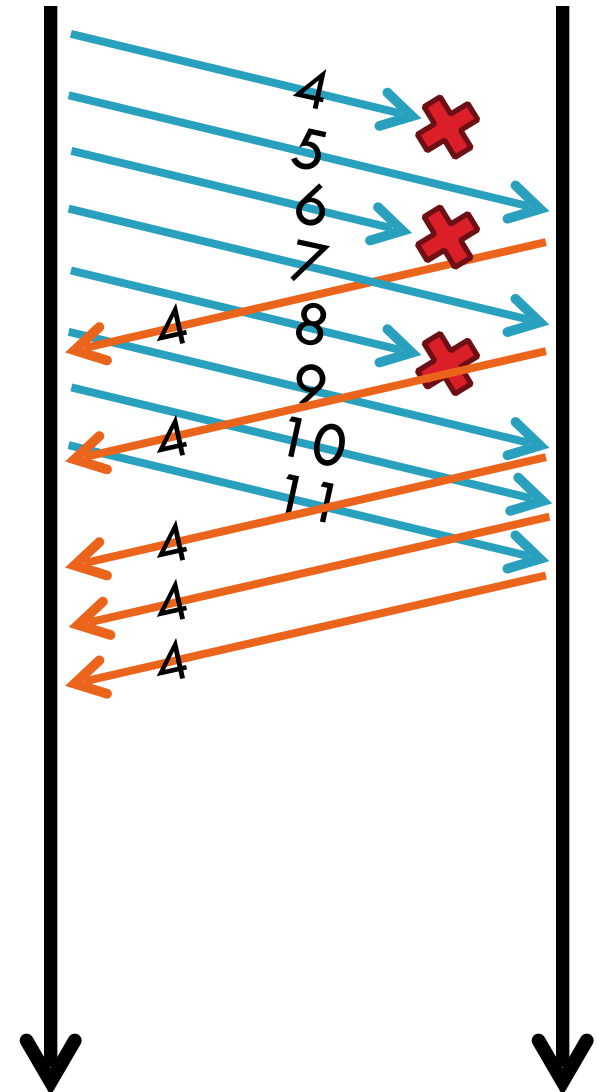
47



SACK: Selective Acknowledgment

47

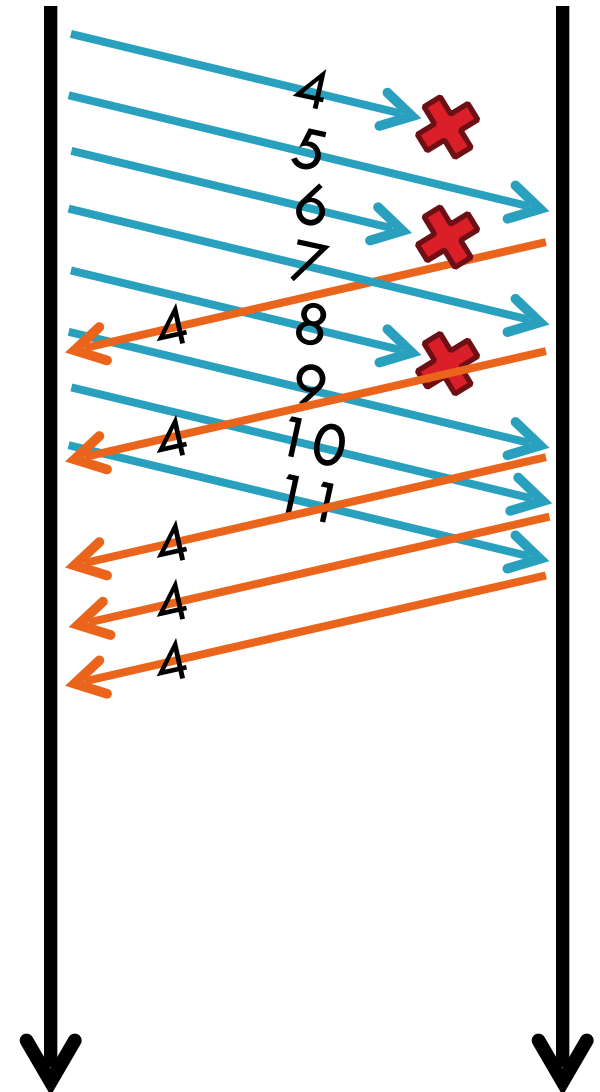
- Problem: duplicate ACKs only tell us about 1 missing packet
 - ▣ Multiple rounds of dup ACKs needed to fill all holes



SACK: Selective Acknowledgment

47

- Problem: duplicate ACKs only tell us about 1 missing packet
 - ▣ Multiple rounds of dup ACKs needed to fill all holes
- Solution: selective ACK
 - ▣ Include received, out-of-order sequence numbers in TCP header
 - ▣ Explicitly tells the sender about holes in the sequence



Other Common Options

48

- Maximum segment size (MSS)
 - Essentially, what is the hosts MTU
 - Saves on path discovery overhead

Other Common Options

48

- Maximum segment size (MSS)
 - Essentially, what is the hosts MTU
 - Saves on path discovery overhead
- Timestamp
 - When was the packet sent (approximately)?
 - Used to prevent sequence number wraparound
 - PAWS algorithm

Issues with TCP

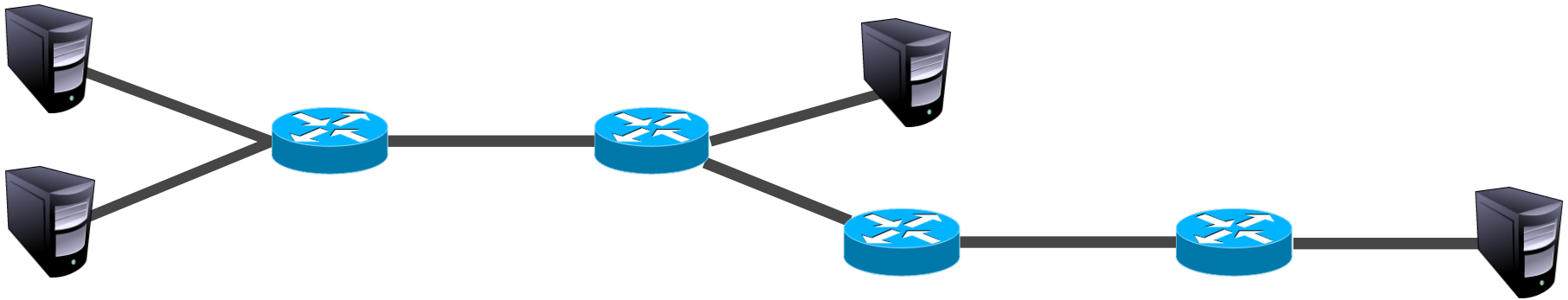
49

- The vast majority of Internet traffic is TCP
- However, many issues with the protocol
 - Lack of fairness
 - Synchronization of flows
 - Poor performance with small flows
 - Really poor performance on wireless networks
 - Susceptibility to denial of service

Fairness

50

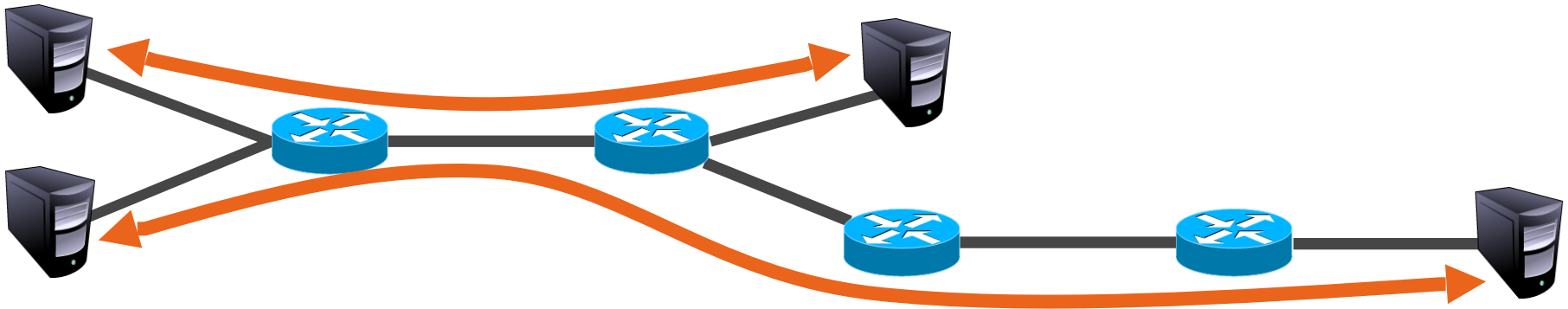
- Problem: TCP throughput depends on RTT



Fairness

50

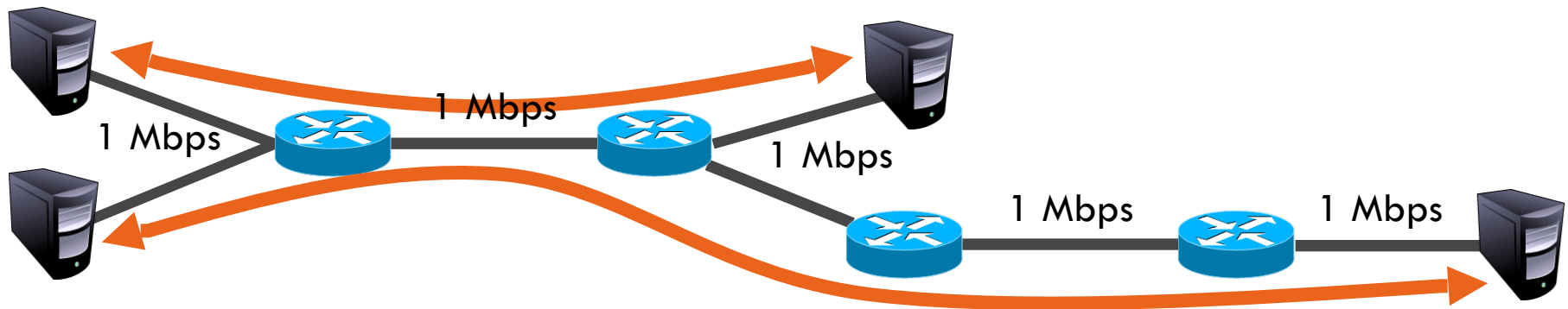
- Problem: TCP throughput depends on RTT



Fairness

50

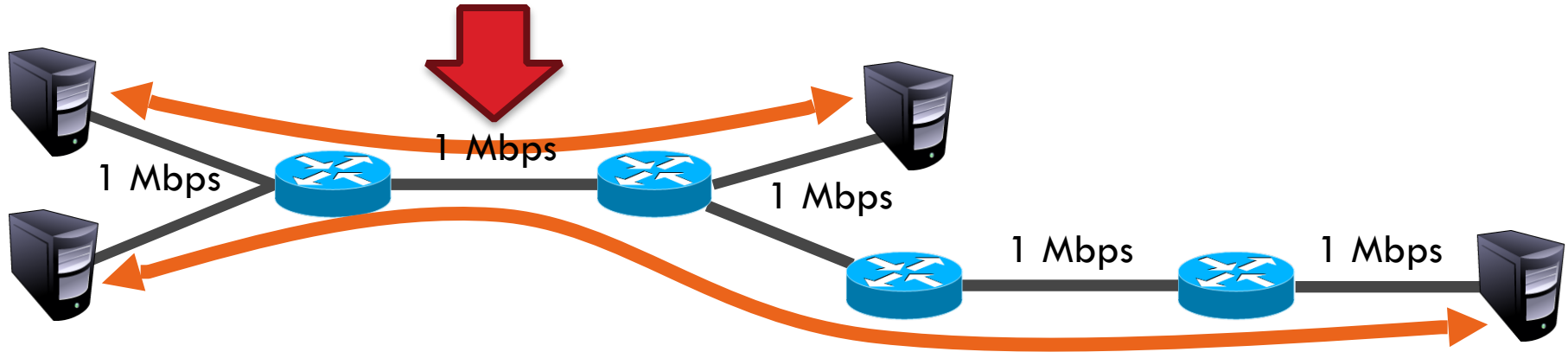
- Problem: TCP throughput depends on RTT



Fairness

50

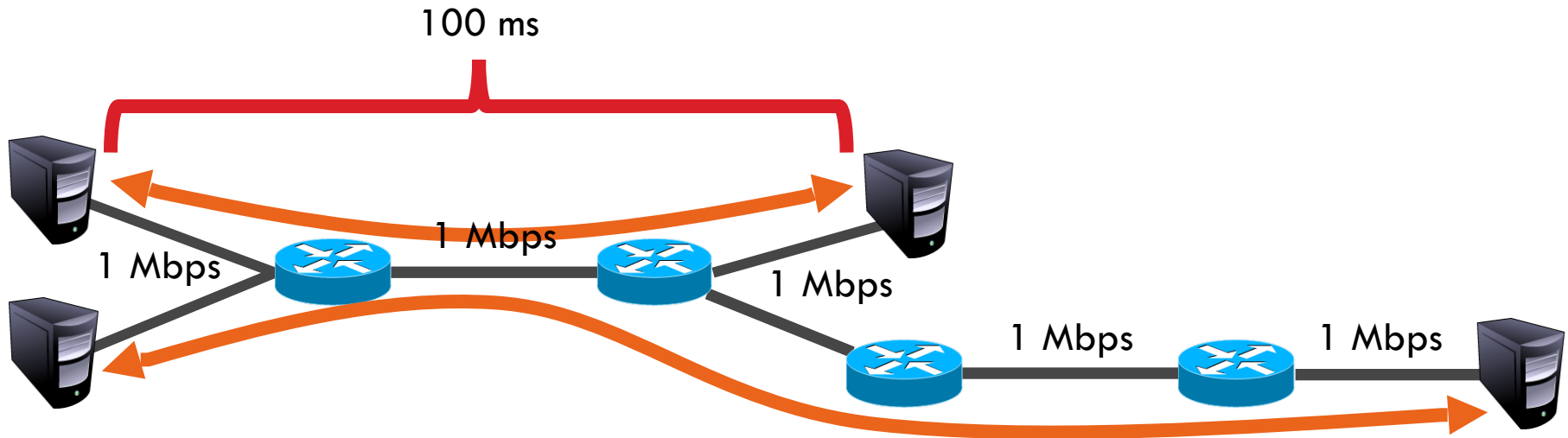
- Problem: TCP throughput depends on RTT



Fairness

50

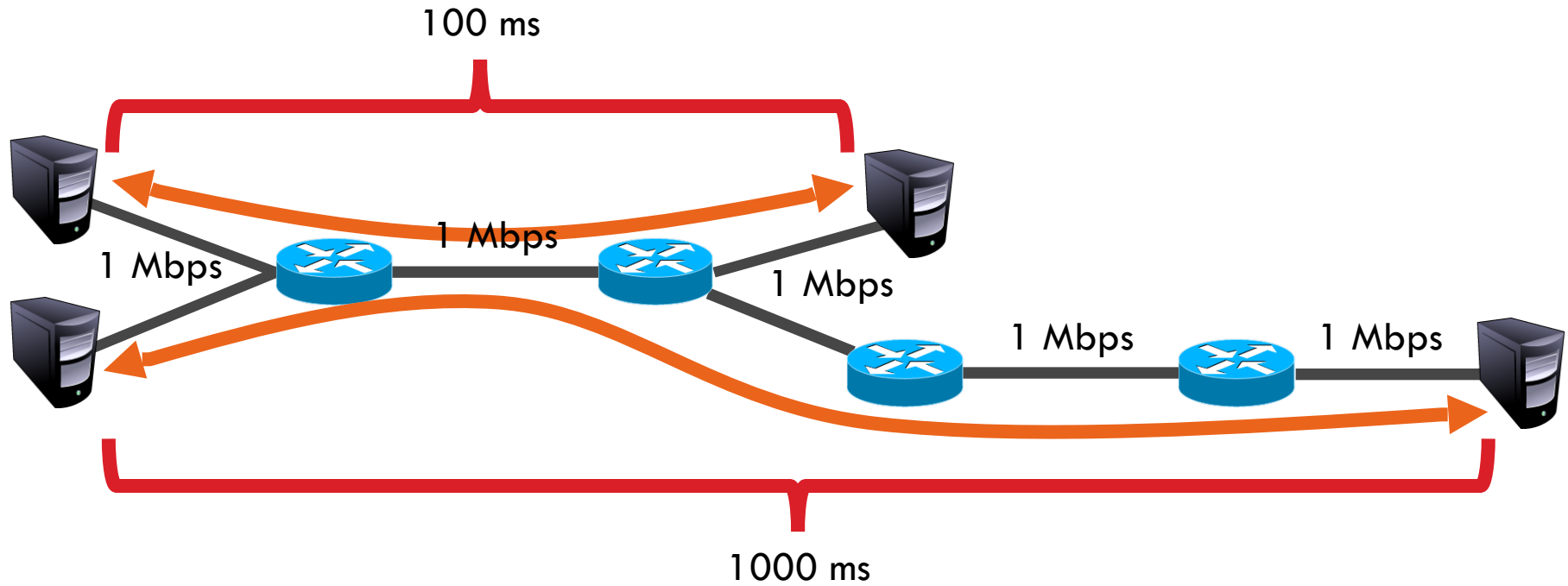
- Problem: TCP throughput depends on RTT



Fairness

50

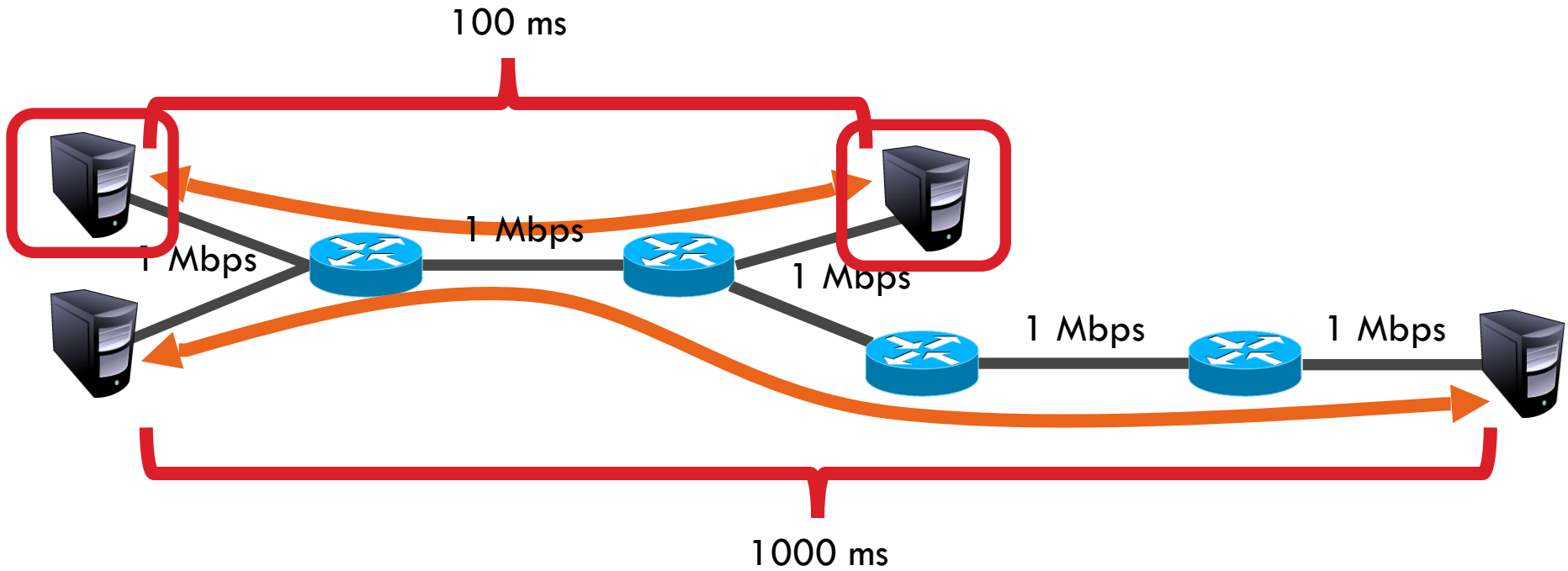
- Problem: TCP throughput depends on RTT



Fairness

50

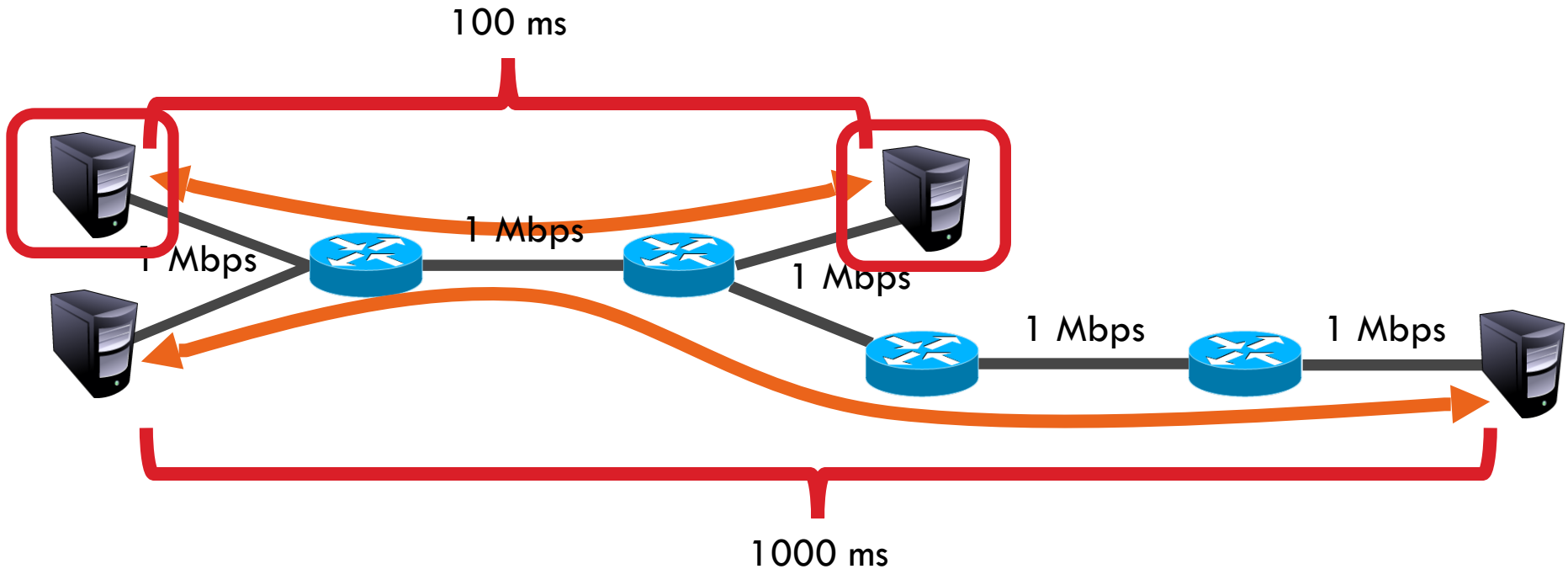
- Problem: TCP throughput depends on RTT



Fairness

50

- Problem: TCP throughput depends on RTT

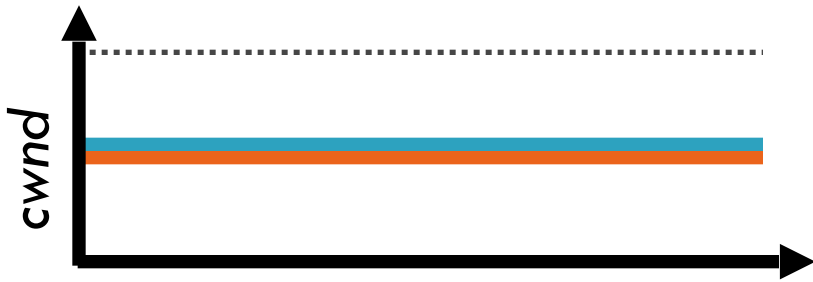


- ACK clocking makes TCP inherently unfair
- Possible solution: maintain a separate delay window
 - ▣ Implemented by Microsoft's Compound TCP

Synchronization of Flows

51

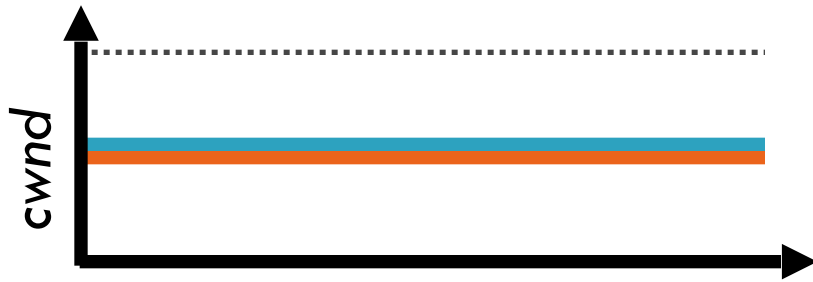
- Ideal bandwidth sharing



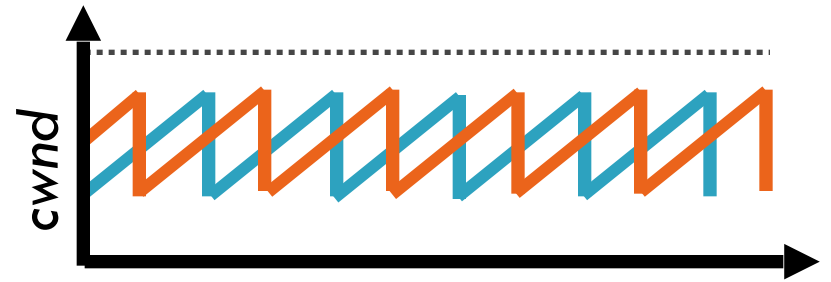
Synchronization of Flows

51

- Ideal bandwidth sharing



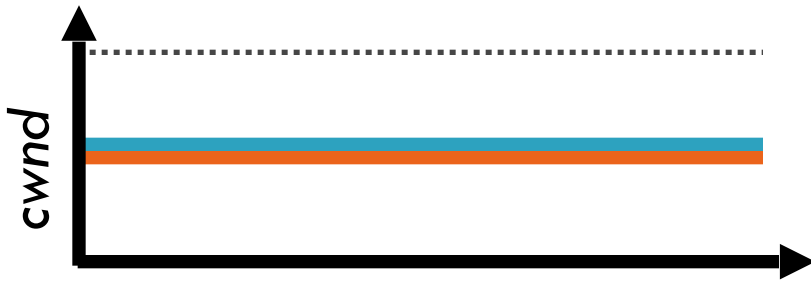
- Oscillating, but high overall utilization



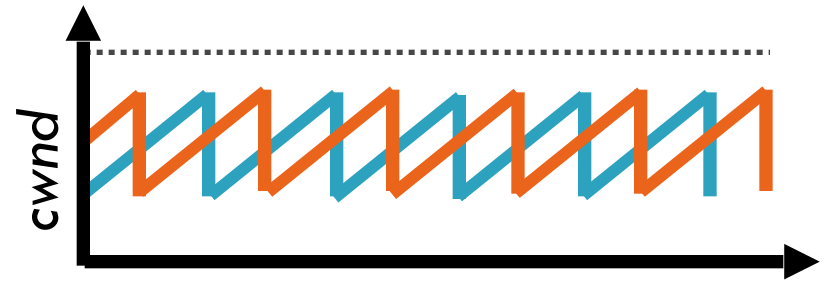
Synchronization of Flows

51

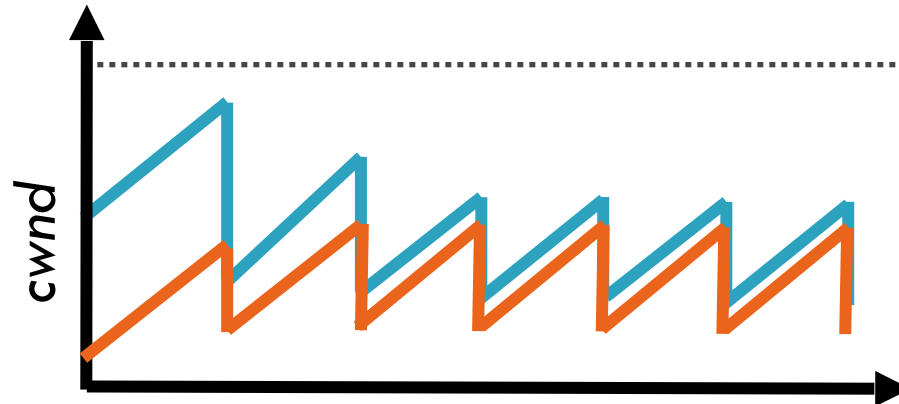
- Ideal bandwidth sharing



- Oscillating, but high overall utilization



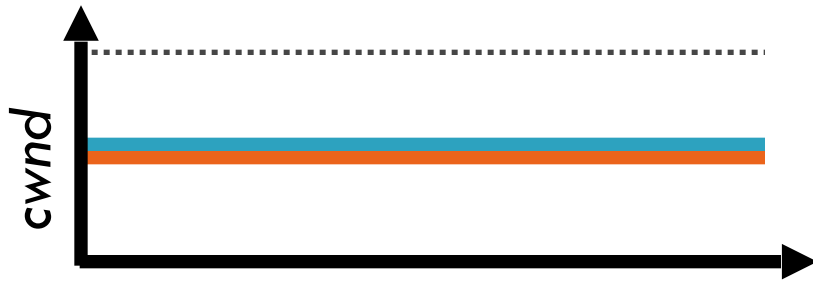
- In reality, flows synchronize



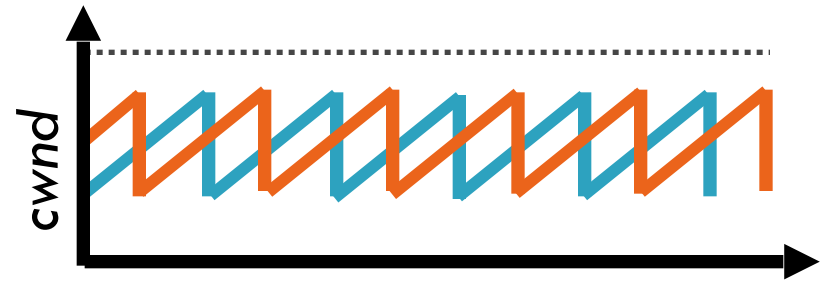
Synchronization of Flows

51

- Ideal bandwidth sharing

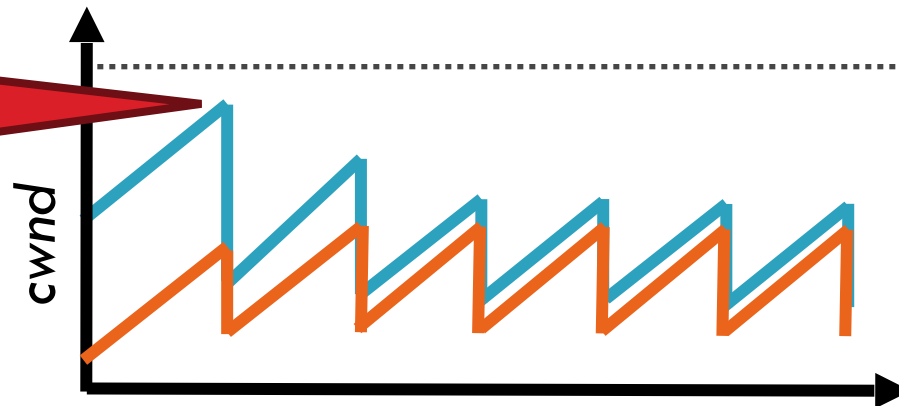


- Oscillating, but high overall utilization



- In reality, flows synchronize

One flow causes all flows to drop packets



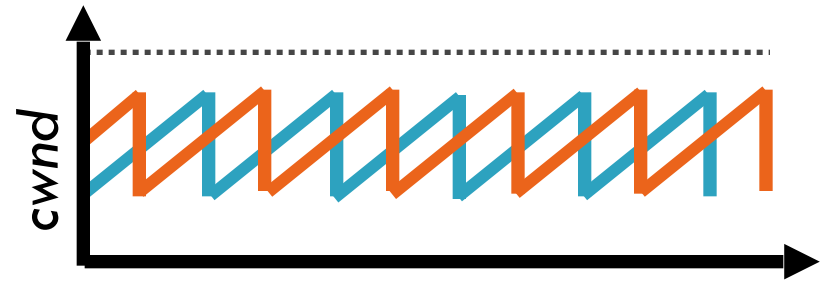
Synchronization of Flows

51

- Ideal bandwidth sharing

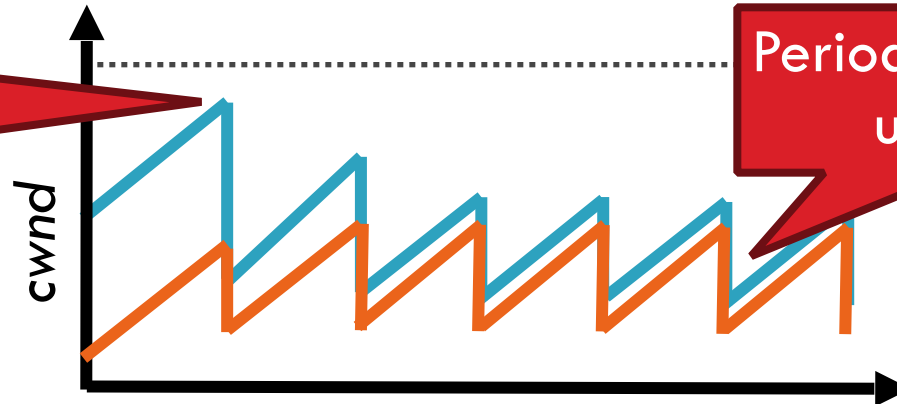


- Oscillating, but high overall utilization



- In reality, flows synchronize

One flow causes all flows to drop packets



Periodic lulls of low utilization

Small Flows

52

- Problem: TCP is biased against short flows
 - 1 RTT wasted for connection setup (SYN, SYN/ACK)
 - *cwnd* always starts at 1

Small Flows

52

- Problem: TCP is biased against short flows
 - 1 RTT wasted for connection setup (SYN, SYN/ACK)
 - *cwnd* always starts at 1
- Vast majority of Internet traffic is short flows
 - Mostly HTTP transfers, <100KB
 - Most TCP flows never leave slow start!

Small Flows

52

- Problem: TCP is biased against short flows
 - ▣ 1 RTT wasted for connection setup (SYN, SYN/ACK)
 - ▣ *cwnd* always starts at 1
- Vast majority of Internet traffic is short flows
 - ▣ Mostly HTTP transfers, <100KB
 - ▣ Most TCP flows never leave slow start!
- Proposed solutions (driven by Google):
 - ▣ Increase initial *cwnd* to 10
 - ▣ TCP Fast Open: use cryptographic hashes to identify receivers, eliminate the need for three-way handshake

Wireless Networks

53

- Problem: Tahoe and Reno assume loss = congestion
 - ▣ True on the WAN, bit errors are very rare
 - ▣ False on wireless, interference is very common

Wireless Networks

53

- Problem: Tahoe and Reno assume loss = congestion
 - ▣ True on the WAN, bit errors are very rare
 - ▣ False on wireless, interference is very common
- TCP throughput $\sim 1/\sqrt{\text{drop rate}}$
 - ▣ Even a few interference drops can kill performance

Wireless Networks

53

- Problem: Tahoe and Reno assume loss = congestion
 - ▣ True on the WAN, bit errors are very rare
 - ▣ False on wireless, interference is very common
- TCP throughput $\sim 1/\sqrt{\text{drop rate}}$
 - ▣ Even a few interference drops can kill performance
- Possible solutions:
 - ▣ Break layering, push data link info up to TCP
 - ▣ Use delay-based congestion detection (TCP Vegas)
 - ▣ Explicit congestion notification (ECN)

Denial of Service

54

- Problem: TCP connections require state
 - Initial SYN allocates resources on the server
 - State must persist for several minutes (RTO)

Denial of Service

54

- Problem: TCP connections require state
 - ▣ Initial SYN allocates resources on the server
 - ▣ State must persist for several minutes (RTO)
- SYN flood: send enough SYNs to a server to allocate all memory/meltdown the kernel

Denial of Service

54

- Problem: TCP connections require state
 - ▣ Initial SYN allocates resources on the server
 - ▣ State must persist for several minutes (RTO)
- SYN flood: send enough SYNs to a server to allocate all memory/meltdown the kernel
- Solution: SYN cookies
 - ▣ Idea: don't store initial state on the server
 - ▣ Securely insert state into the SYN/ACK packet
 - ▣ Client will reflect the state back to the server

SYN Cookies

55

0

Sequence Number

SYN Cookies

55



SYN Cookies

55



- Did the client really send me a SYN recently?
 - ▣ Timestamp: freshness check
 - ▣ Cryptographic hash: prevents spoofed packets

SYN Cookies

55



- Did the client really send me a SYN recently?
 - ▣ Timestamp: freshness check
 - ▣ Cryptographic hash: prevents spoofed packets
- Maximum segment size (MSS)
 - ▣ Usually stated by the client during initial SYN
 - ▣ Server should store this value...
 - ▣ Reflect the clients value back through them

SYN Cookies in Practice

56

- Advantages
 - Effective at mitigating SYN floods
 - Compatible with all TCP versions
 - Only need to modify the server
 - No need for client support

SYN Cookies in Practice

56

□ Advantages

- Effective at mitigating SYN floods
- Compatible with all TCP versions
- Only need to modify the server
- No need for client support

□ Disadvantages

- MSS limited to 3 bits, may be smaller than clients actual MSS
- Server forgets all other TCP options included with the client's SYN
 - SACK support, window scaling, etc.