

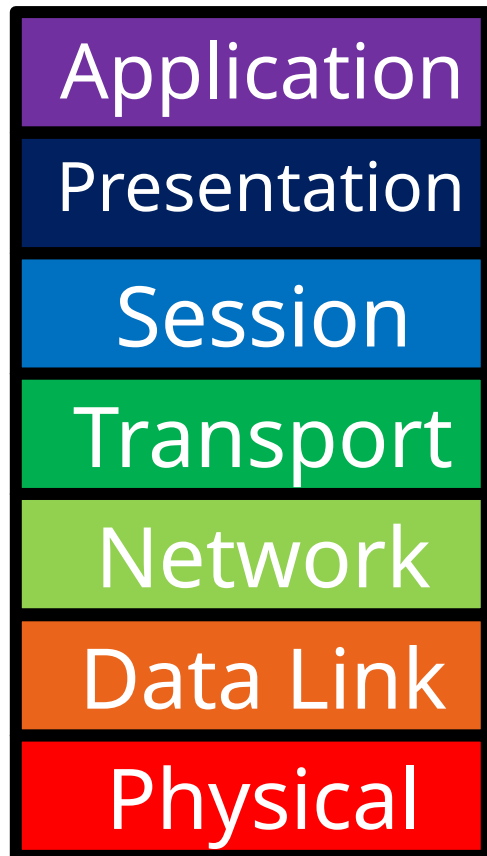
Computer Networks

Lecture 10: Transport layer Part II

Based on slides from D. Choffnes Northeastern U. and P. Gill from StonyBrook University
Revised Autumn 2015 by S. Laki

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

TCP Congestion Control

3

- ❑ **The network is congested if the load in the network is higher than its capacity.**
- ❑ 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

Two Basic Components

4

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)

2. Rate adjustment algorithm

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

Rate Adjustment

5

- 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? !!!!

Implementing Congestion Control

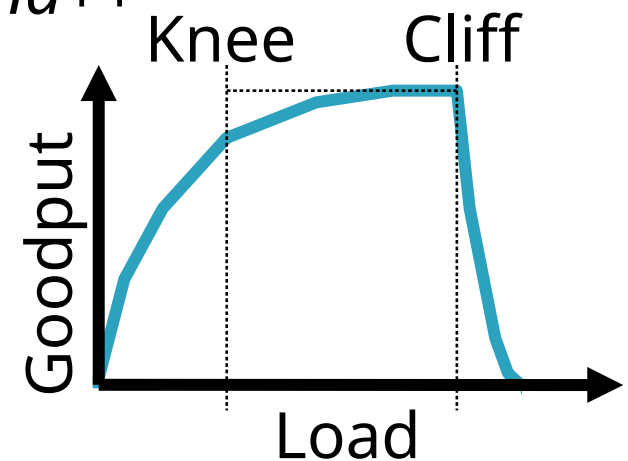
6

- 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

7

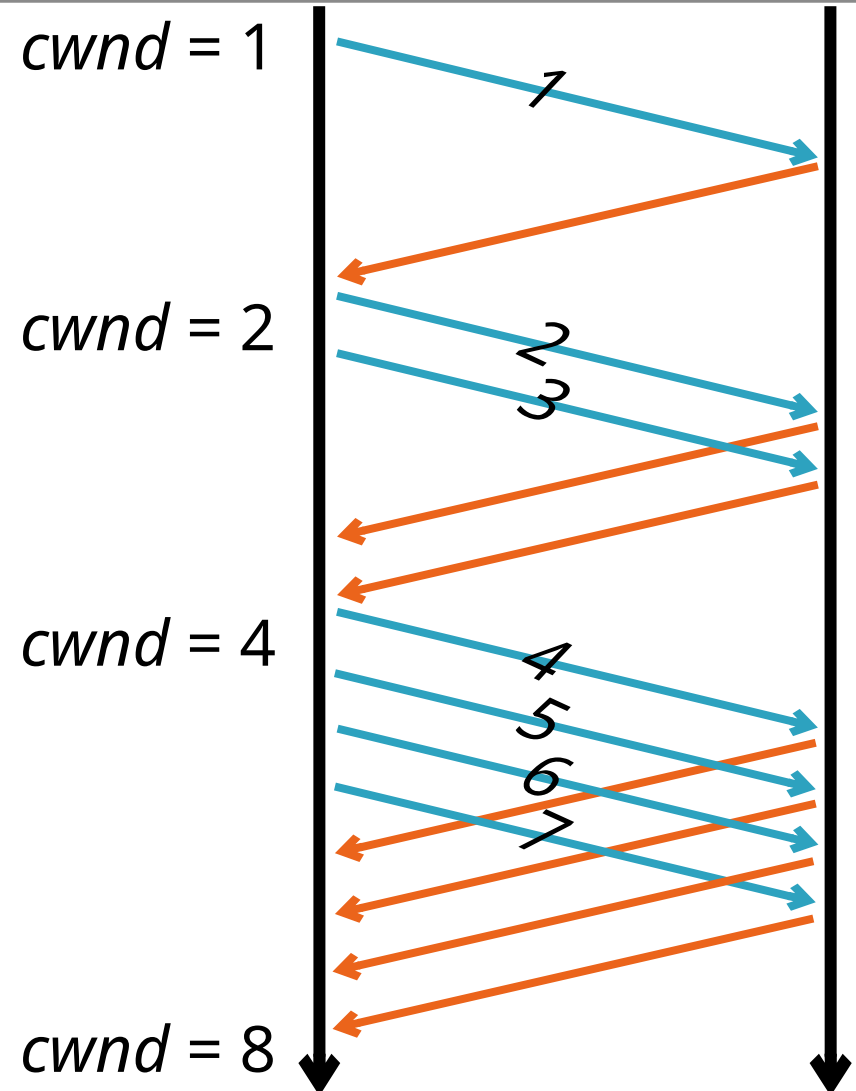
- Goal: reach knee quickly
- Upon starting (or restarting) a connection
 - $cwnd = 1$
 - $ssthresh = adv_wnd$
 - Each time a segment is ACKed, $cwnd++$
- Continues until...
 - $ssthresh$ is reached
 - Or a packet is lost
- Slow Start is not actually slow
 - $cwnd$ increases exponentially



Slow Start Example

8

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



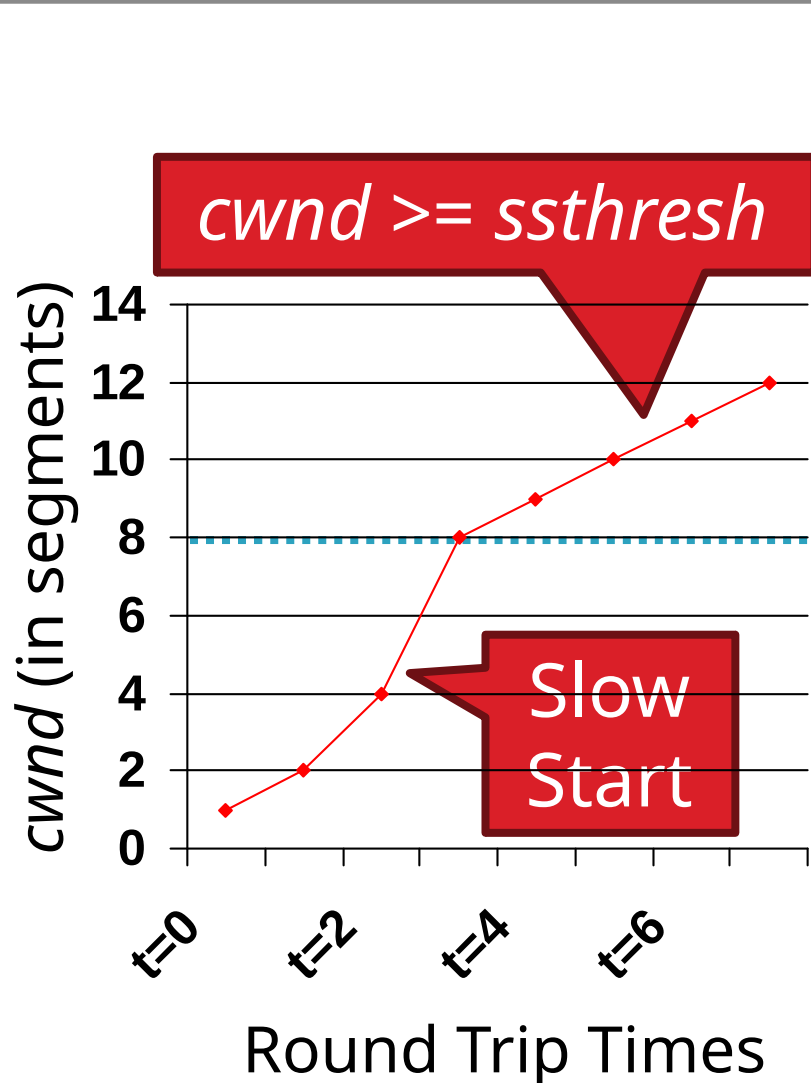
Congestion Avoidance

9

- Additive Increase Multiplicative Decrease (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

10



$cwnd = 1$

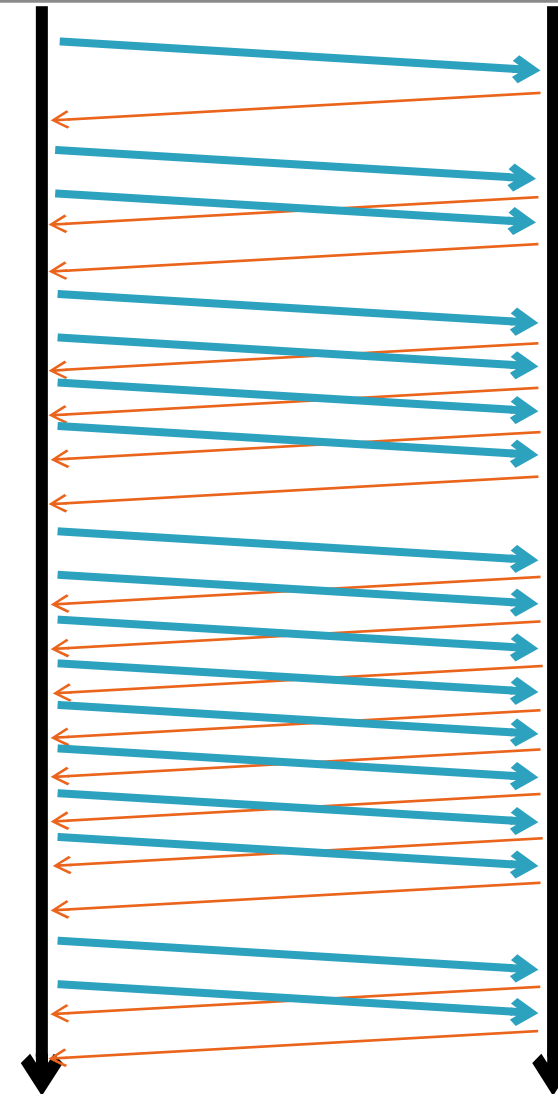
$cwnd = 2$

$cwnd = 4$

$ssthresh = 8$

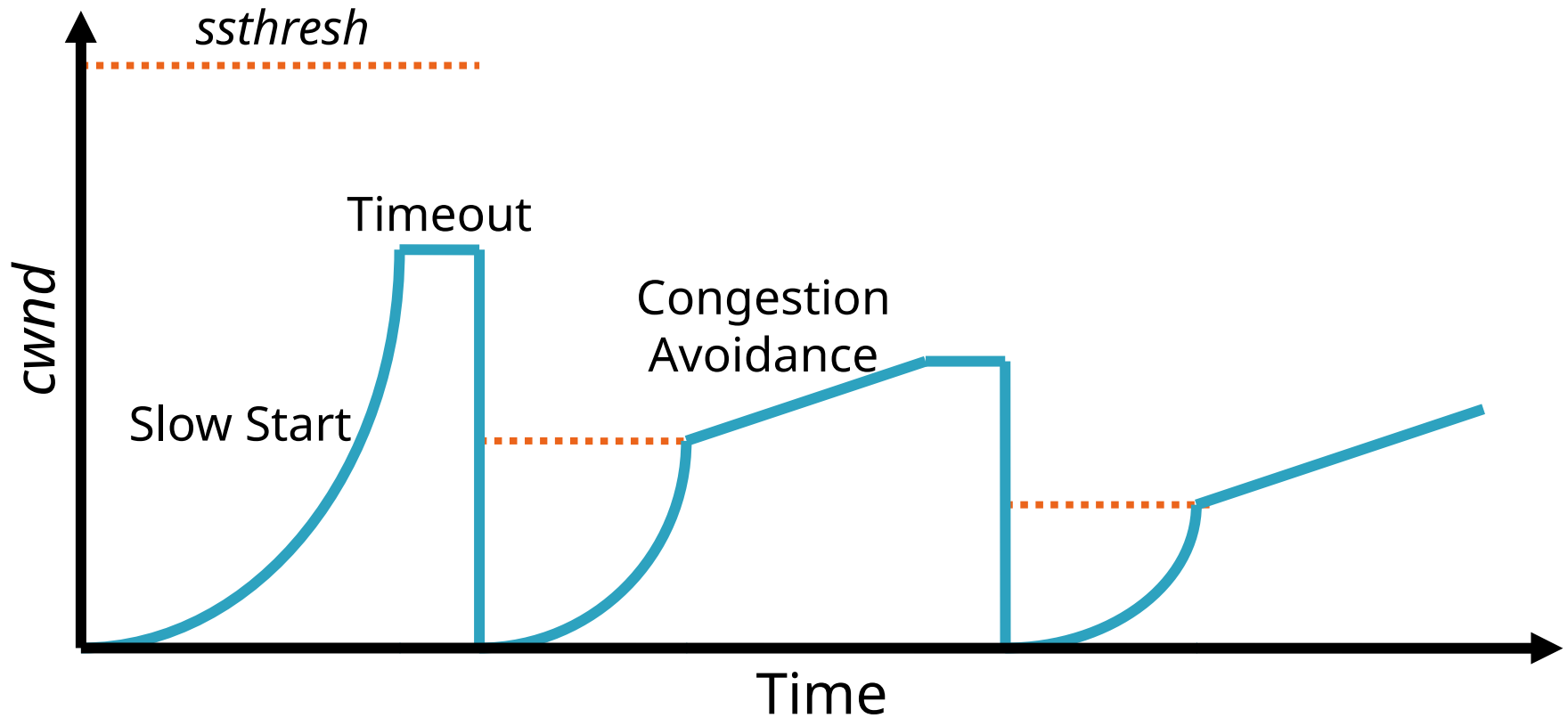
$cwnd = 8$

$cwnd = 9$



The Big Picture – TCP Tahoe (the original TCP)

11



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

The Evolution of TCP

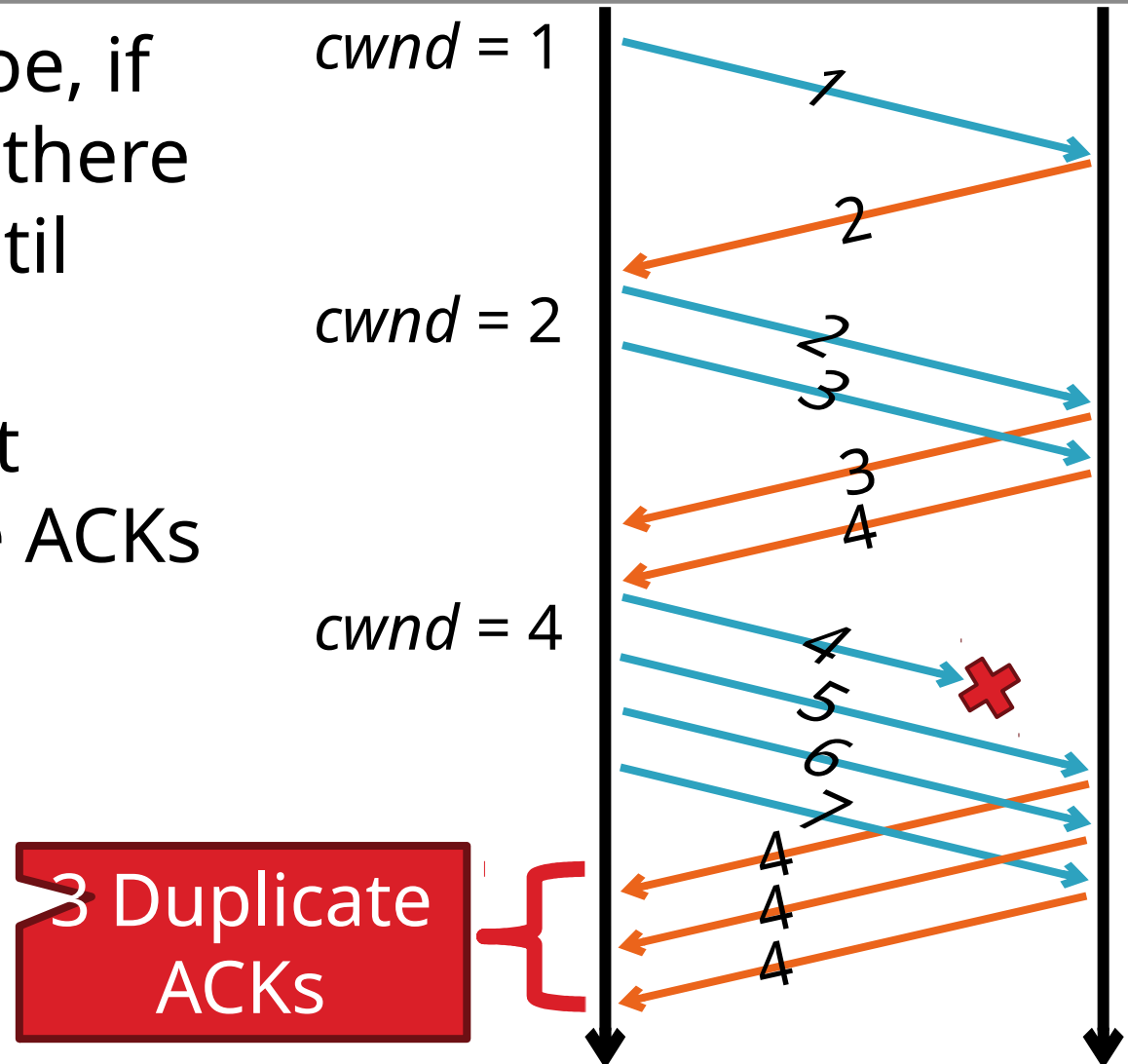
13

- 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
 - 3 duplicate ACKs? -> retransmit (don't wait for RTO)
 - Fast recovery
 - On loss: set $cwnd = cwnd/2$ ($ssthresh = \text{new } cwnd \text{ value}$)

TCP Reno: Fast Retransmit

14

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

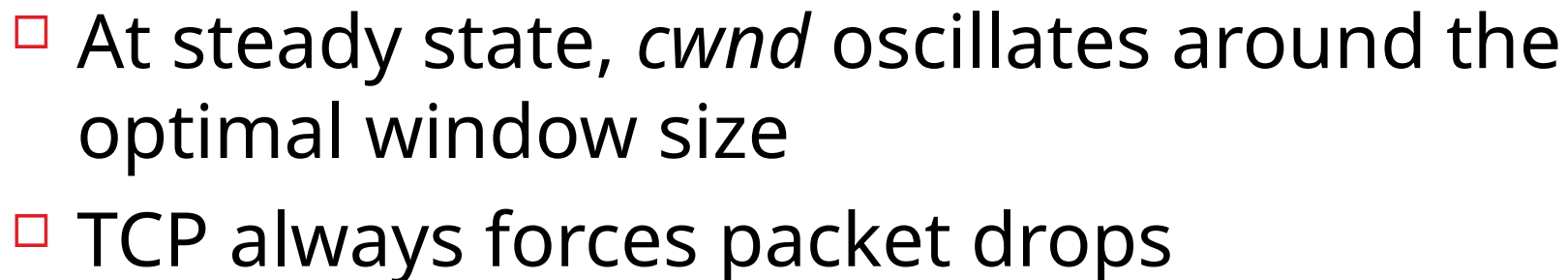


TCP Reno: Fast Recovery

15

- After a fast-retransmit set *cwnd* to $cwnd/2$
 - Also reset *ssthresh* to the new halved *cwnd* value
 - 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

16



Many TCP Variants...

17

- Tahoe: the original
 - Slow start with AIMD
 - Dynamic RTO based on RTT estimate
- Reno:
 - fast retransmit (3 dupACKs)
 - fast recovery ($\text{cwnd} = \text{cwnd}/2$ on loss)
- 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

18

- 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

High Bandwidth-Delay Product

19

- 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

Goals

20

- 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

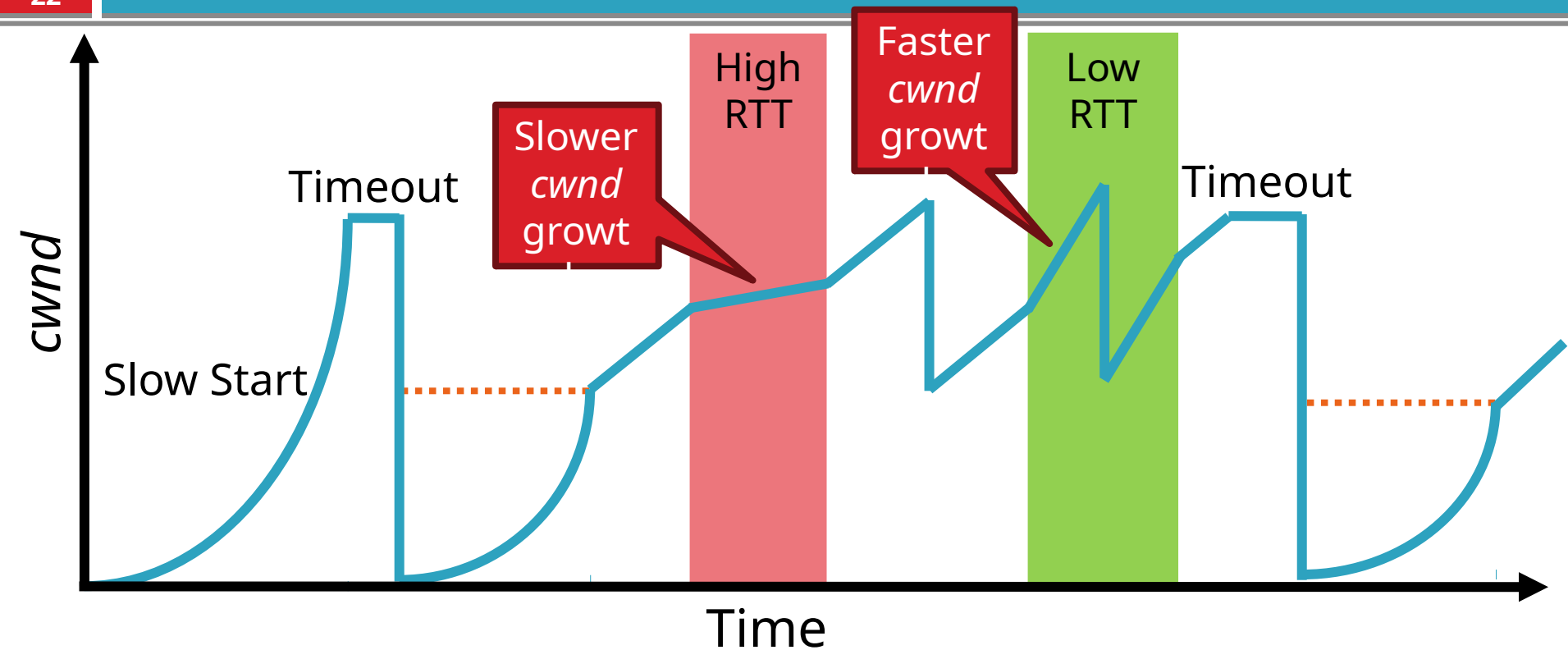
Compound TCP Implementation

21

- 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

22



- 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

23

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

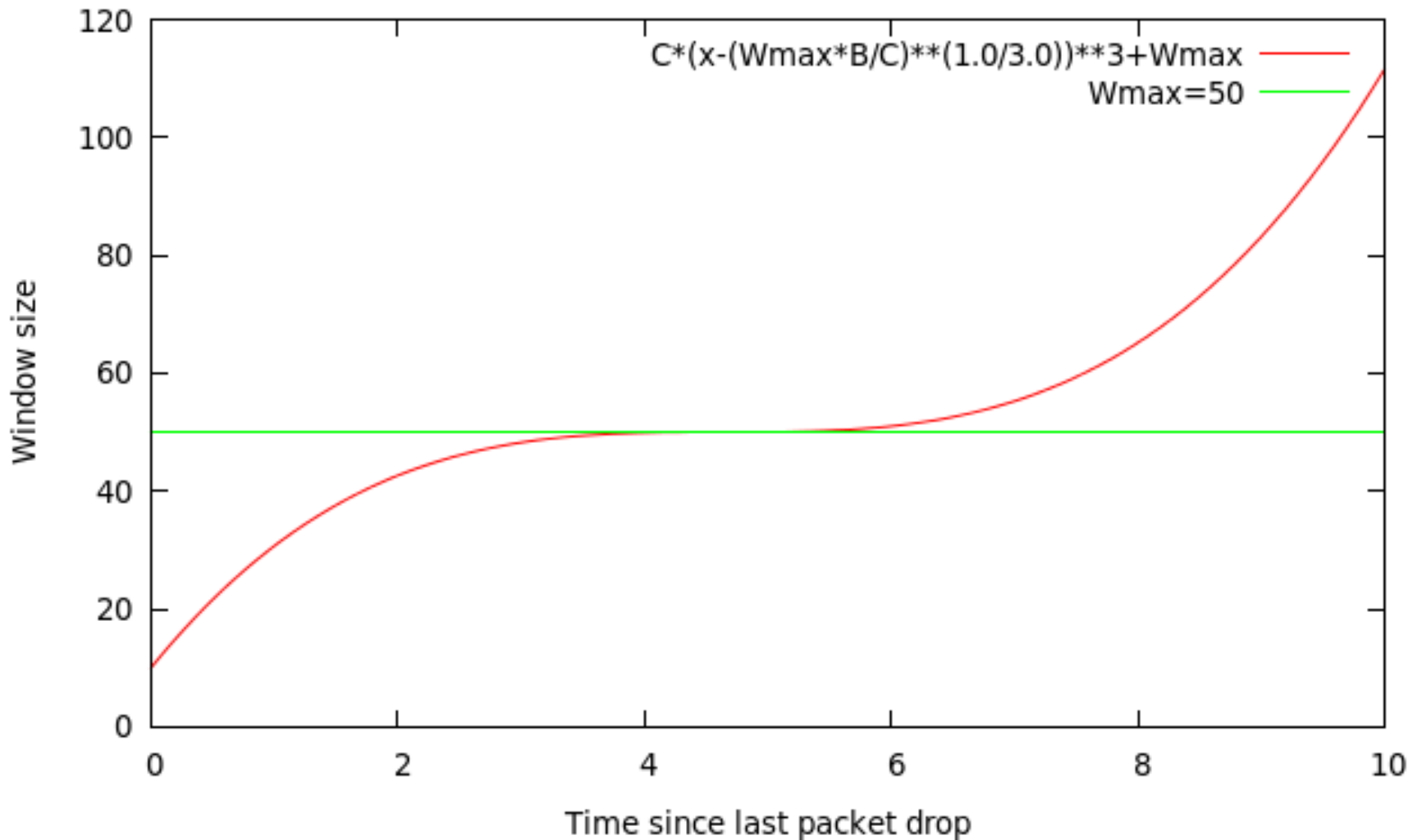
$$W_{cubic} = C(T - K)^3 + W_{max} \quad (1)$$

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

- B □ a constant fraction for multiplicative increase
- T □ time since last packet drop
- W_max □ cwnd when last packet dropped

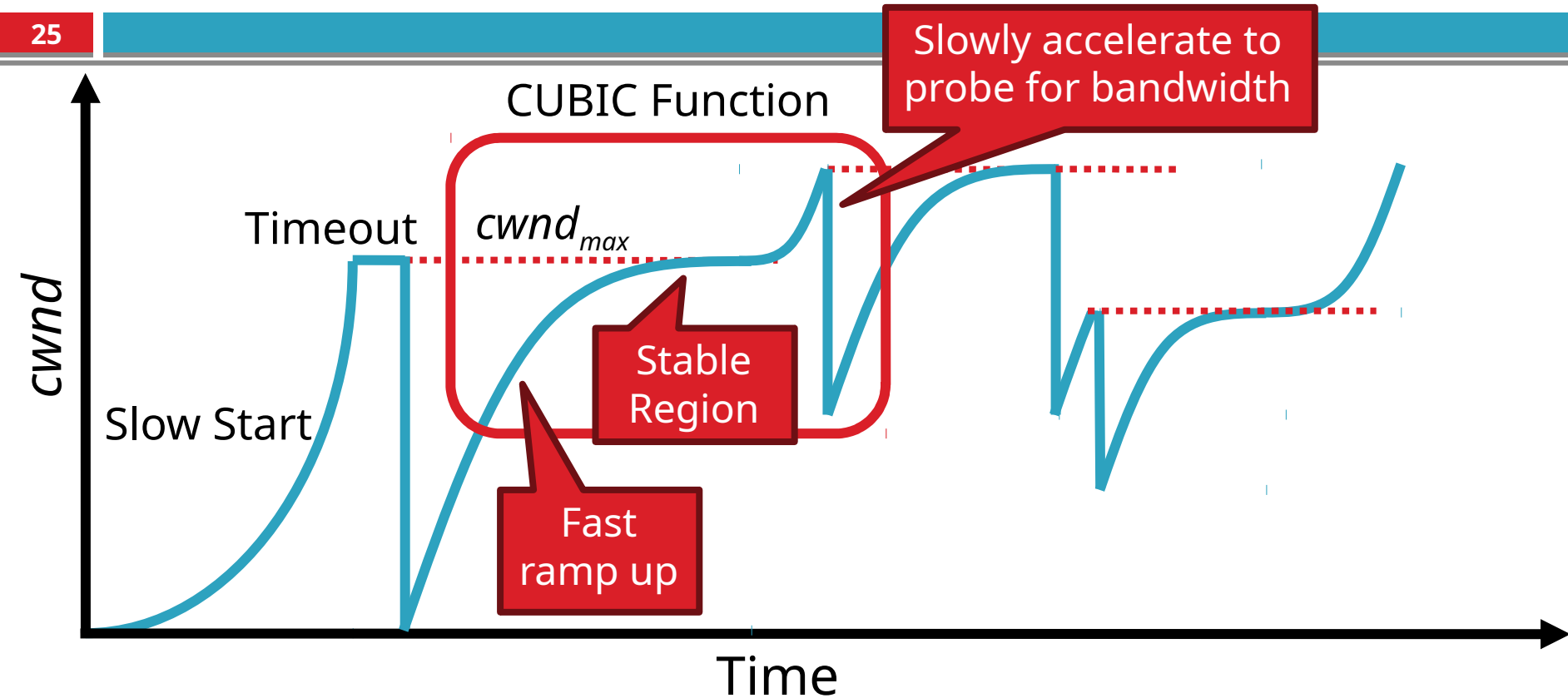
TCP CUBIC Implementation

24



TCP CUBIC Example

25



- 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

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

Issues with TCP

27

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

Small Flows

28

- 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

29

- 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

30

- ❑ 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 (sequence number field)
 - Client will reflect the state back to the server

Further topics

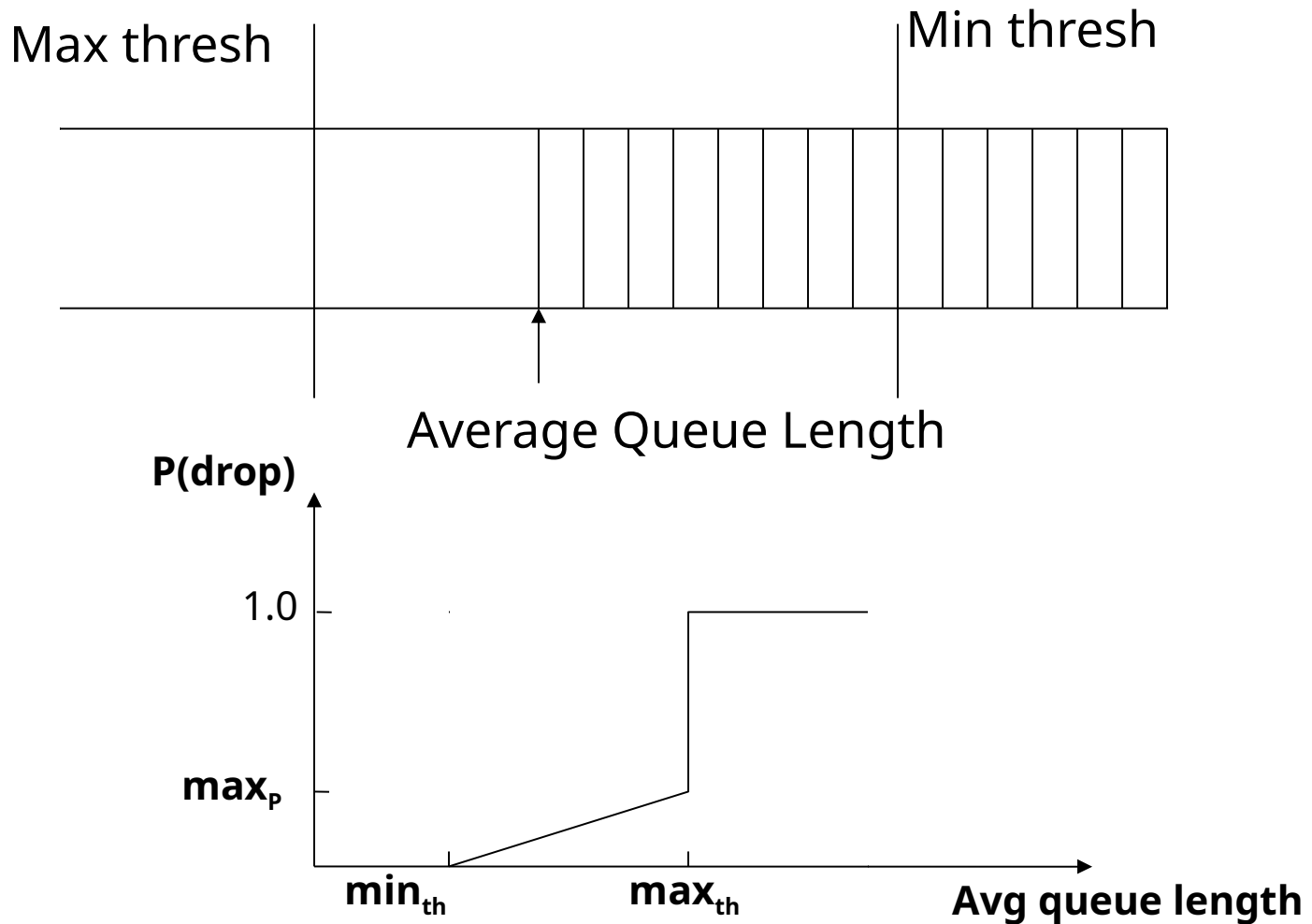
Typical Internet Queuing

- ❑ 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

RED Algorithm

- Maintain running average of queue length
- If $\text{avgq} < \text{min}_{\text{th}}$ do nothing
 - Low queuing, send packets through
- If $\text{avgq} > \text{max}_{\text{th}}$, drop packet
 - Protection from misbehaving sources
- Else mark packet in a manner proportional to queue length
 - Notify sources of incipient congestion
 - E.g. by ECN IP field or dropping packets with a given probability

RED Operation



RED Algorithm

- Maintain running average of queue length
- For each packet arrival
 - Calculate average queue size (avg)
 - If $\min_{th} \leq avg < \max_{th}$
 - Calculate probability P_a
 - With probability P_a
 - Mark the arriving packet: drop or set-up ECN
 - Else if $\max_{th} \leq avg$
 - Mark the arriving packet: drop, ECN

Data Center TCP: DCTCP

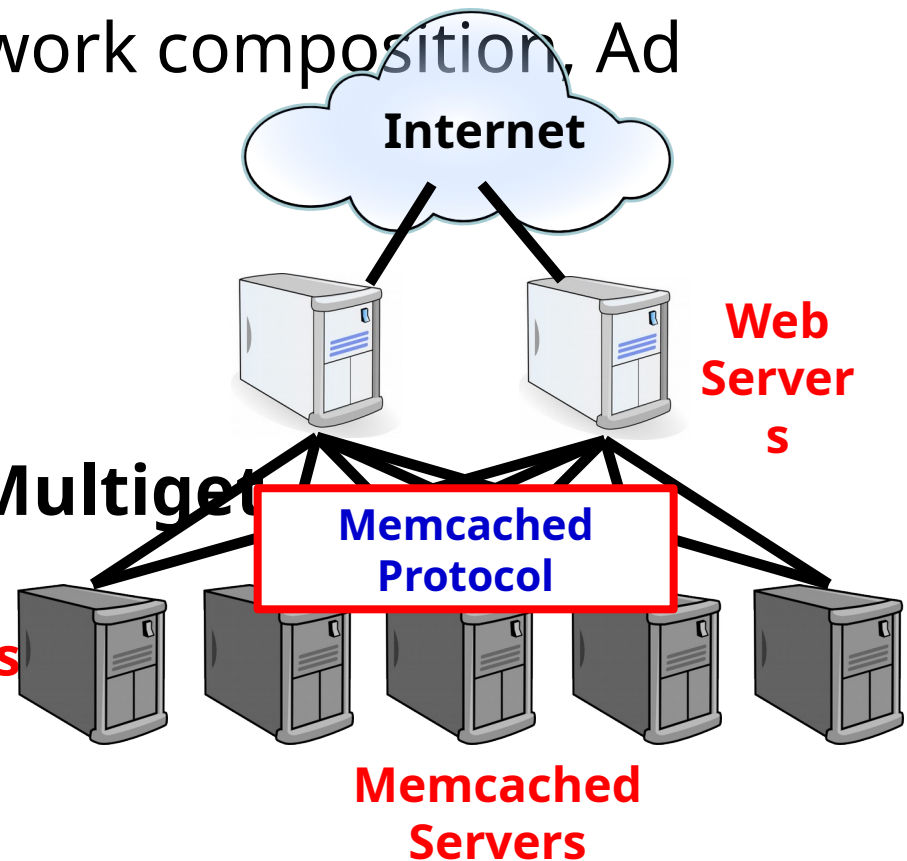
Generality of Partition/Aggregate

- The foundation for many large-scale web applications.
 - Web search, Social network composition, Ad selection, etc.

- Example: **Facebook**

Partition/Aggregate ~ Multiget

- Aggregators: **Web Servers**
- Workers: **Memcached Servers**

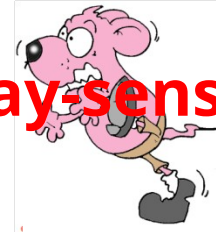


Workloads

38

- Partition/Aggregate
(Query)

→ **Delay-sensitive**



- Short messages [50KB-1MB]
(Coordination, Control state)

→ **Delay-sensitive**



- Large flows [1MB-50MB]
(Data update)

→ **Throughput-sensitive**



Impairments

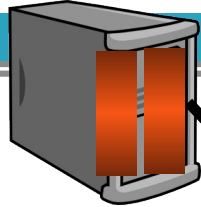
39

- Incast
- Queue Buildup
- Buffer Pressure

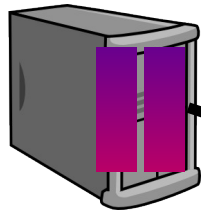
Incast

40

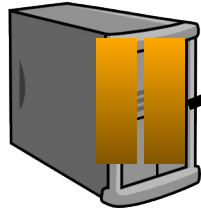
Worker
1



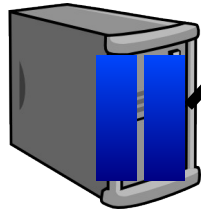
Worker
2



Worker 3



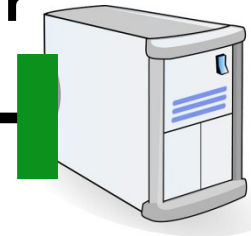
Worker
4



- Synchronized mice collide.

➤ **Caused by
Partition/Aggregate.**

Aggregato
r



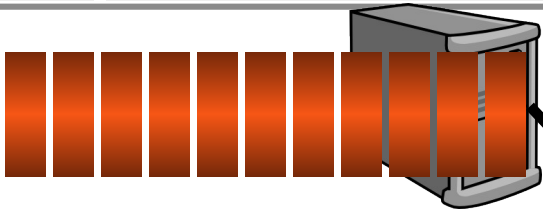
$RTO_{min} = 300 \text{ ms}$



← **TCP
timeout**

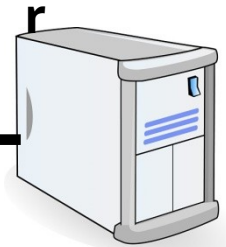
Queue Buildup

Sender 1

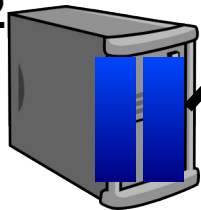


- Big flows buildup queues.
 - **Increased latency for short flows.**

Receiver



Sender 2



- Measurements in Bing cluster
 - **For 90% packets: $RTT < 1ms$**
 - **For 10% packets: $1ms < RTT < 15ms$**

Data Center Transport Requirements

42

1. High Burst Tolerance

- Incast due to Partition/Aggregate is common.

2. Low Latency

- Short flows, queries

3. High Throughput

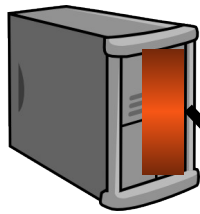
- Continuous data updates, large file transfers

The challenge is to achieve these three together.

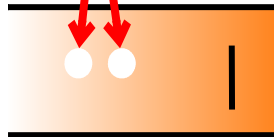
DCTCP: The TCP/ECN Control Loop

Sender 1

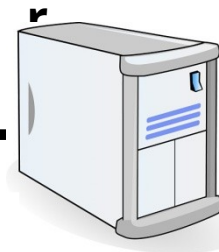
ECN = Explicit Congestion Notification



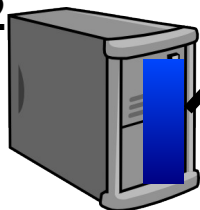
ECN Mark (1 bit)



Receive



**Sender
2**



DCTCP: Two Key Ideas

18

1. React in proportion to the **extent** of congestion, not its **presence**.

✓ Reduces **variance** in sending rates, lowering queuing

ECN Marks	TCP	DCTCP
1 0 1 1 1 1 0 1 1 1	Cut window by 50%	Cut window by 40%
0 0 0 0 0 0 0 0 0 1	Cut window by 50%	Cut window by 5%

2. Mark based on **instantaneous** queue length.

✓ Fastfeedback to better deal with bursts.

Data Center TCP Algorithm

19

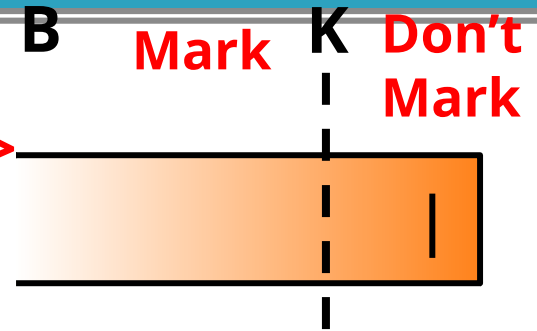
Switch side:

- Mark packets when **Queue Length** >

K.

Sender side:

- Maintain running average of ***fraction*** of packets marked (**α**).



In each RTT:

$$F = \frac{\# \text{ of marked ACKs}}{\text{Total \# of ACKs}}$$

$$\alpha \leftarrow (1 - g)\alpha + gF$$

$$Cwnd \leftarrow \left(1 - \frac{\alpha}{2}\right)Cwnd$$

➤ Adaptive window decreases:

- Note: decrease factor between 1 and 2.