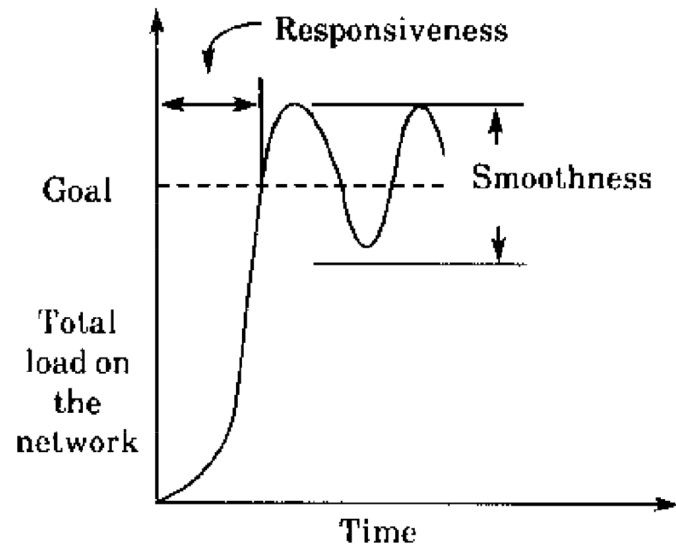
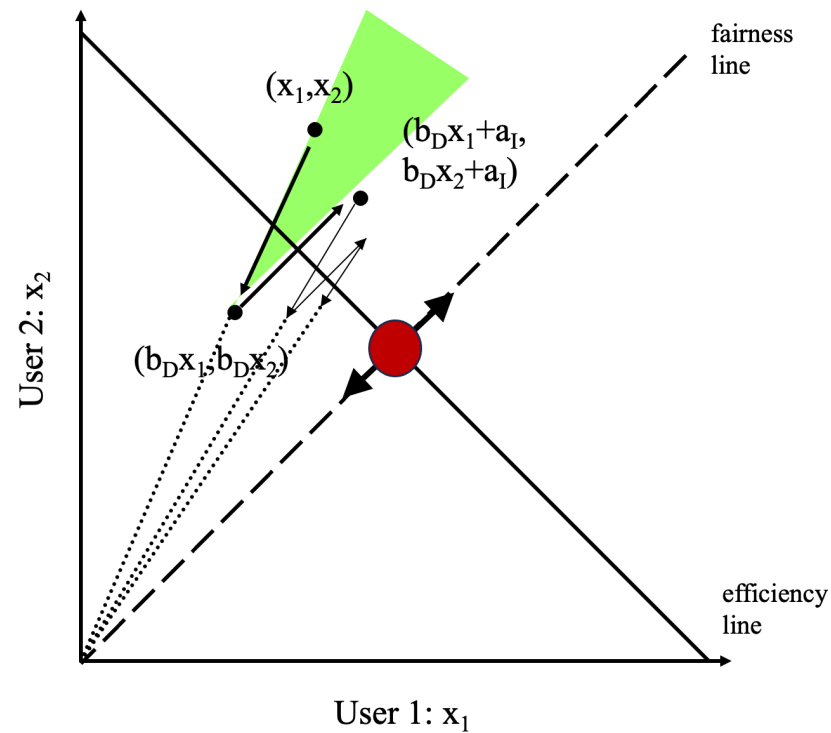
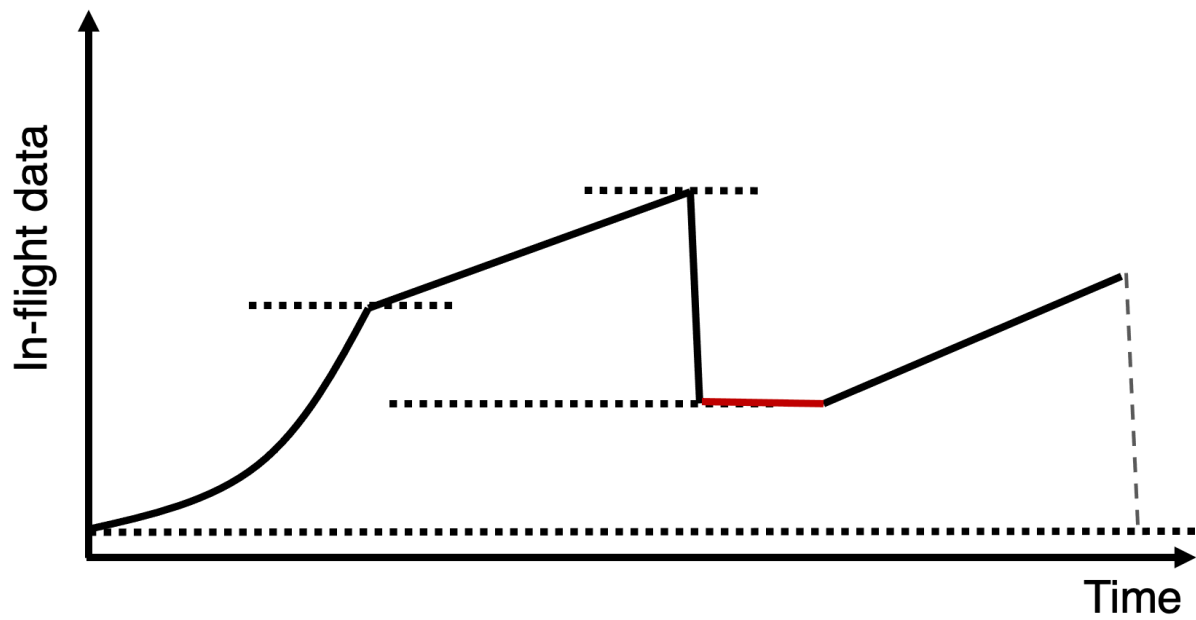
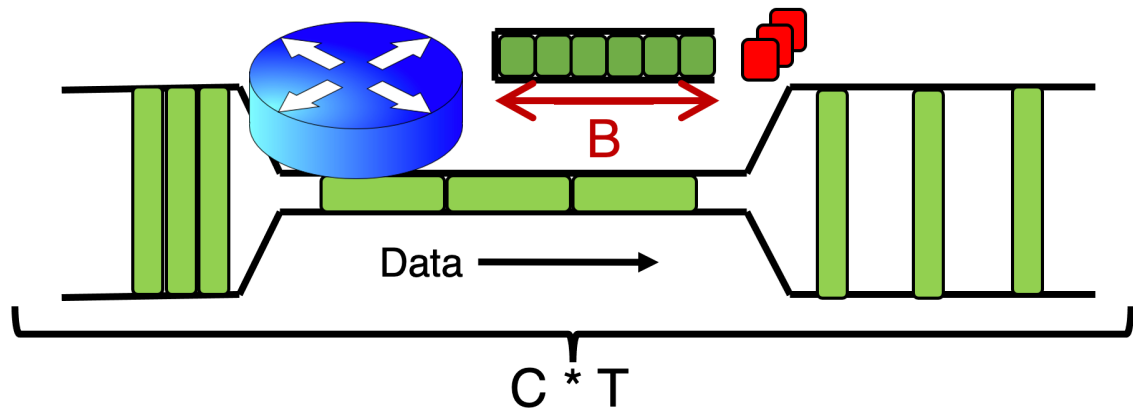


Transport



Simple models are useful

- Chiu and Jain's model isn't indicative of all TCP/AIMD behavior
 - But it's "realistic" enough
 - Stands the test of time: many more sources, much higher bandwidth, ...
- Models should be simple
 - For us to work with
 - For others to understand
 - But they don't have to mimic the "real" thing in every way
- A real, complex model is likely useless, but a realistic, simple model might teach us something

Modeling TCP throughput

Given network characteristics, how quickly can TCP (New Reno) send data?

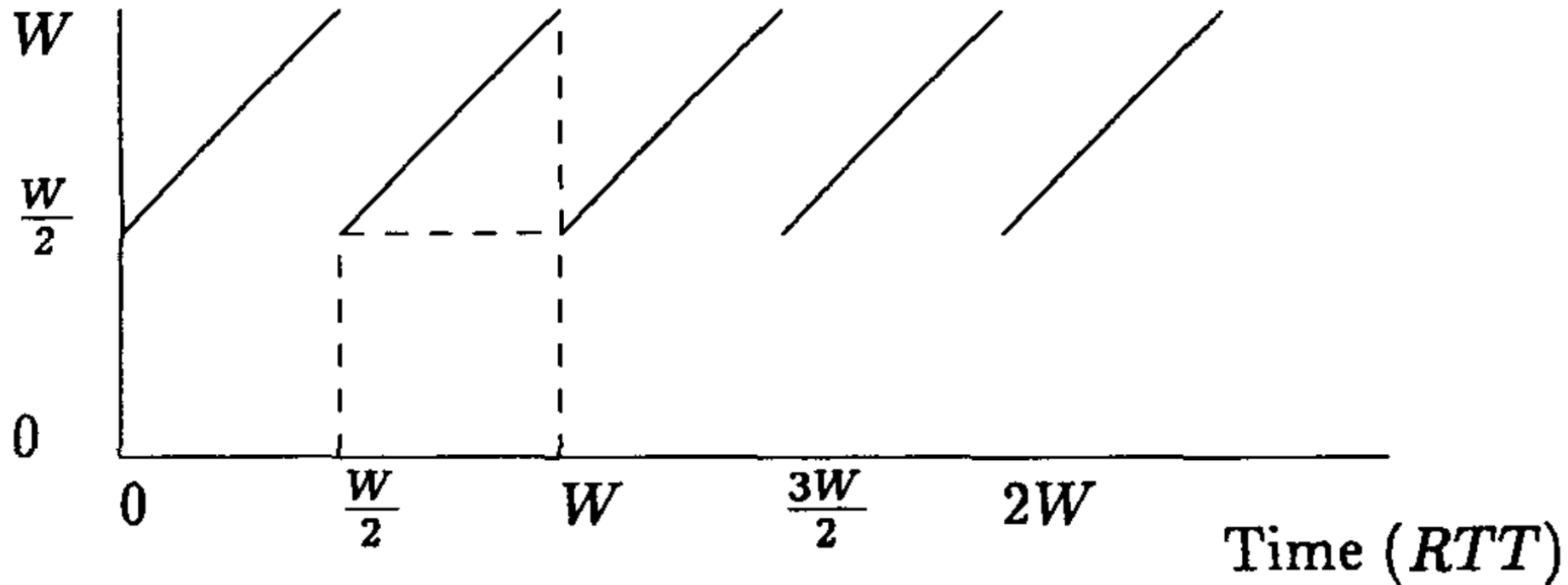
Mathis et al., “Macroscopic behavior of TCP congestion avoidance”

Steady AIMD

$$< \left(\frac{MSS}{RTT} \right) \frac{1}{\sqrt{p}}$$

(Mathis equation)

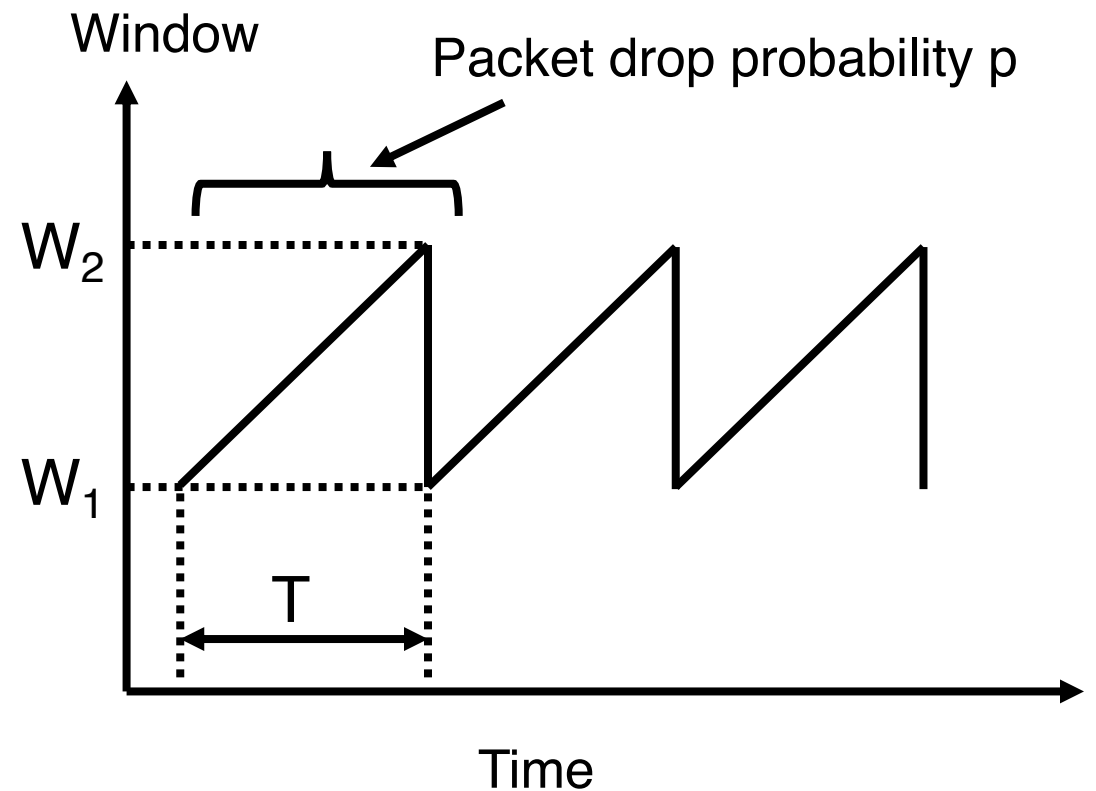
congestion window (packets)



Mathis et al., "Macroscopic behavior of TCP congestion avoidance"

Estimating TCP AIMD throughput

- **Assumptions**
- Single flow, repeating AIMD
- Loss occurs exactly in the last RTT before window reduction
 - Exactly one packet lost
- Assume RTT constant (ignore queueing delay change)
- Relationship between W_1 , W_2 ?
- How many pkts sent over T ?
- Relationship p and # pkts?



Implications

- Throughput has a $1/\sqrt{p}$ dependence on packet loss rate
- Getting full bottleneck throughput requires loss rate $1/(\text{BDP})^2$
- RTT unfairness
 - Flows with a smaller RTT get better throughput (ramp up faster)
- Engineering implications:
 - Split TCP (CDNs, data center frontends, ...)
 - Special considerations for long-distance connections

Widely Deployed TCPs

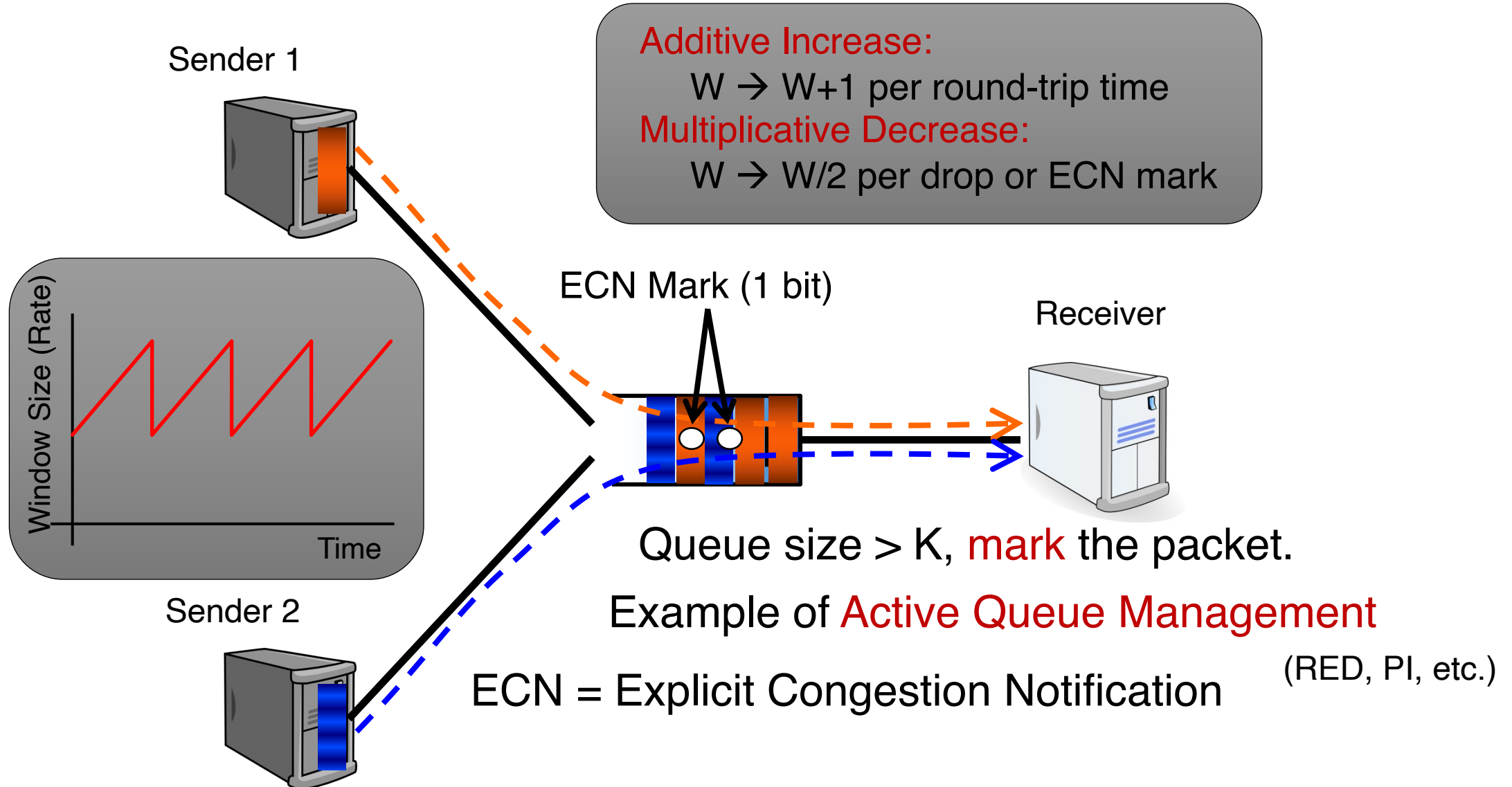
Data Center TCP

Alizadeh et al.

Context

- Regular TCP: window evolution with all signals measured end to end
- Data centers: hardware under single administrative control
 - Could network switches do better?
- What if switches provided better feedback than loss?

Explicit Congestion Notification



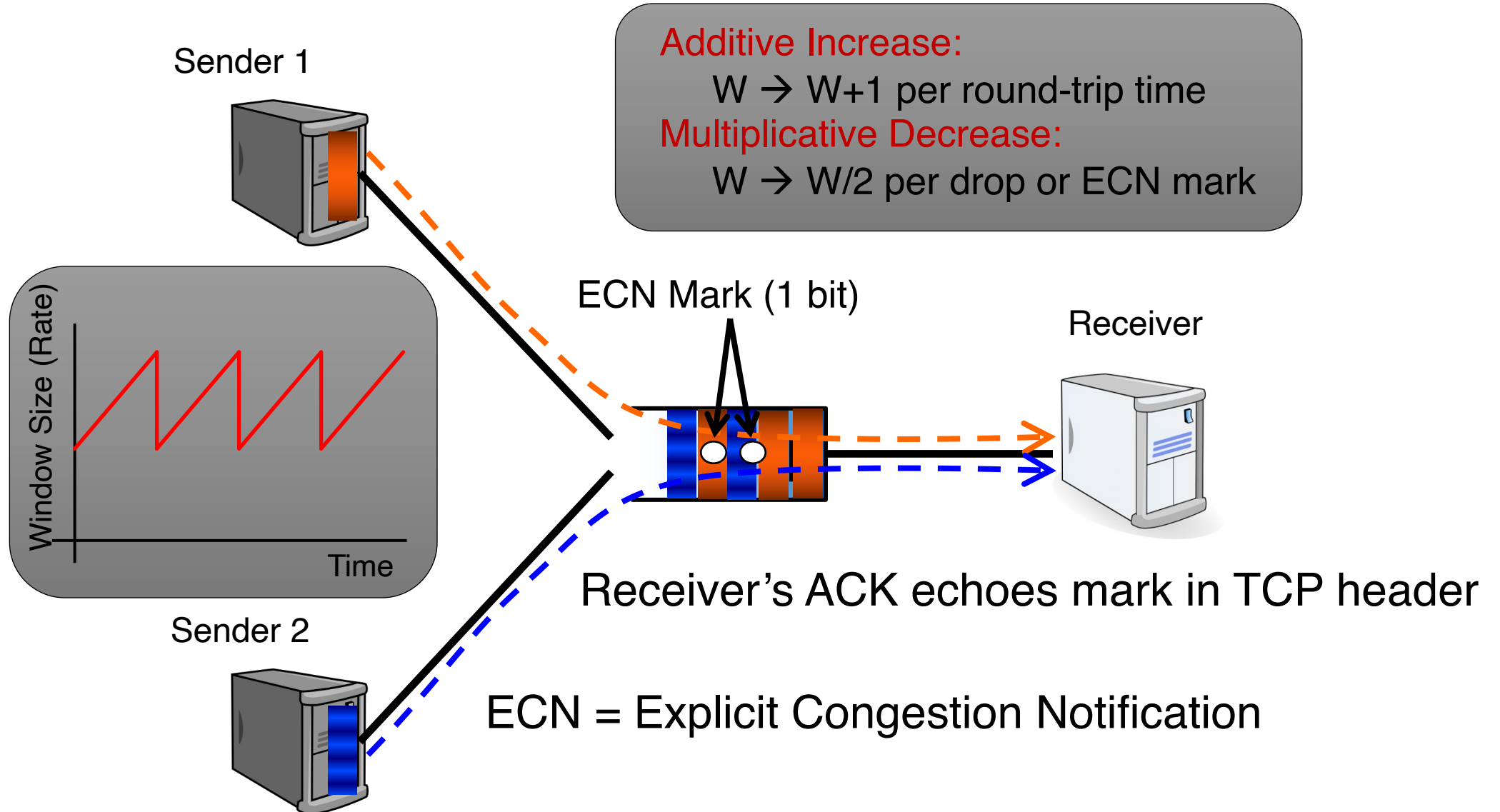
ECN set on the IP header by routers

- 00 – Not ECN-Capable Transport, Not-ECT
- 01 – ECN Capable Transport(1), ECT(1)
- 10 – ECN Capable Transport(0), ECT(0)
- 11 – Congestion Experienced, CE. → **Dropped** by router if TCP sender is not ECN enabled

IPv4 header format

Offsets	Octet	0				1				2				3																			
Octet	Bit	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31
0	0	Version				IHL				DSCP				ECN				Total Length															
4	32	Identification								Flags				Fragment Offset																			
8	64	Time To Live				Protocol				Header Checksum																							
12	96	Source IP Address																															
16	128	Destination IP Address																															
20	160	Options (if IHL > 5)																															
⋮	⋮																																
56	448																																

Explicit Congestion Notification



ECN on the TCP header

TCP segment header

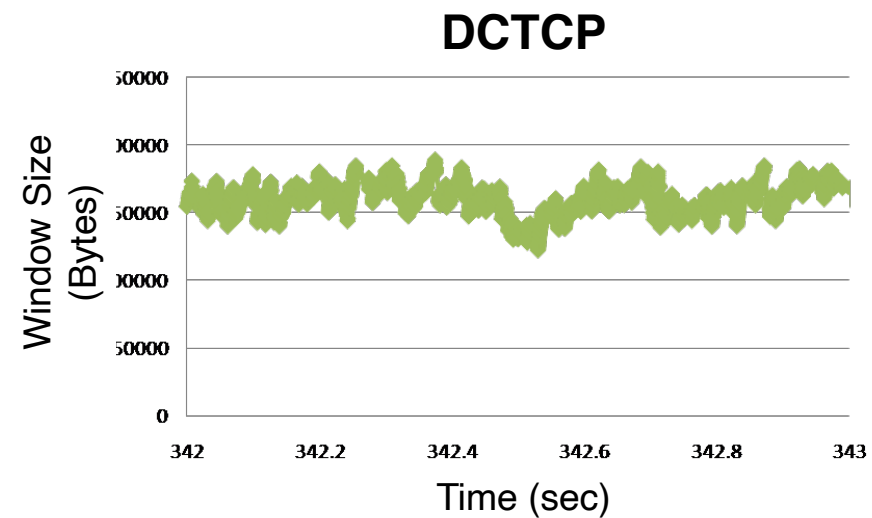
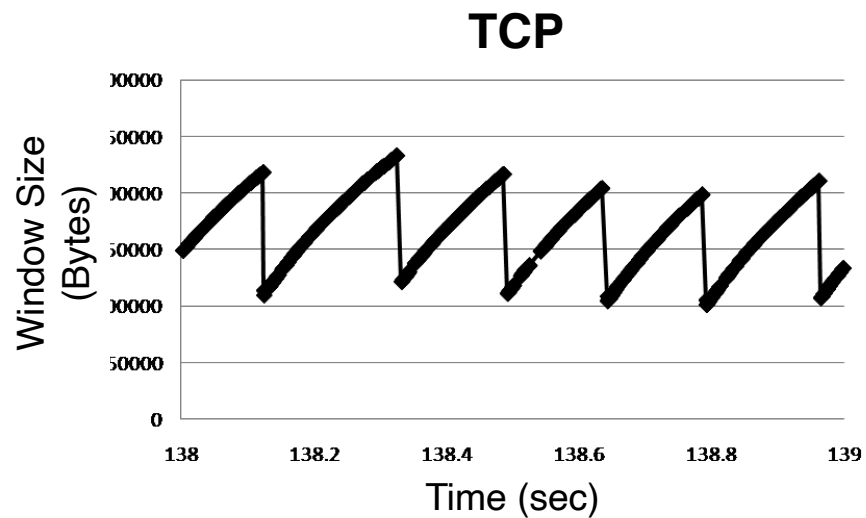
Offsets	Octet	0								1								2								3							
Octet	Bit	7	6	5	4	3	2	1	0	7	6	5	4	3	2	1	0	7	6	5	4	3	2	1	0	7	6	5	4	3	2	1	0
0	0	Source port																Destination port															
4	32	Sequence number																															
8	64	Acknowledgment number (if ACK set)																															
12	96	Data offset	Reserved 0000	C W R	E C E	U R G	A R K	P S H	R S T	S S T	F I N	Window Size																					
16	128	Checksum																Urgent pointer (if URG set)															
20	160	Options (if <i>data offset</i> > 5. Padded at the end with "0" bits if necessary.)																															
:	:																																
56	448																																

DCTCP: Main idea

- Extract multi-bit feedback from single-bit stream of ECN marks
 - Reduce window size based on **fraction** of marked packets

DCTCP: Main idea

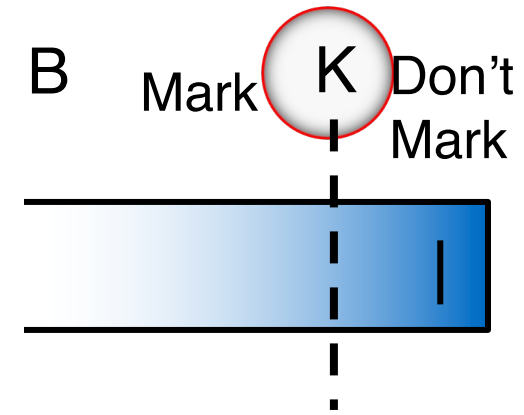
ECN Marks	TCP	DCTCP
1 0 1 1 1 1 0 1 1 1	Cut window by 50%	Cut window by 40%
0 0 0 0 0 0 0 0 0 1	Cut window by 50%	Cut window by 5%



DCTCP algorithm

Router side:

- Mark packets when Queue Length $> K$.



Sender side:

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

$$\text{each RTT: } F = \frac{\# \text{ of marked ACKs}}{\text{Total \# of ACKs}} \Rightarrow \alpha \leftarrow (1 - g)\alpha + gF$$

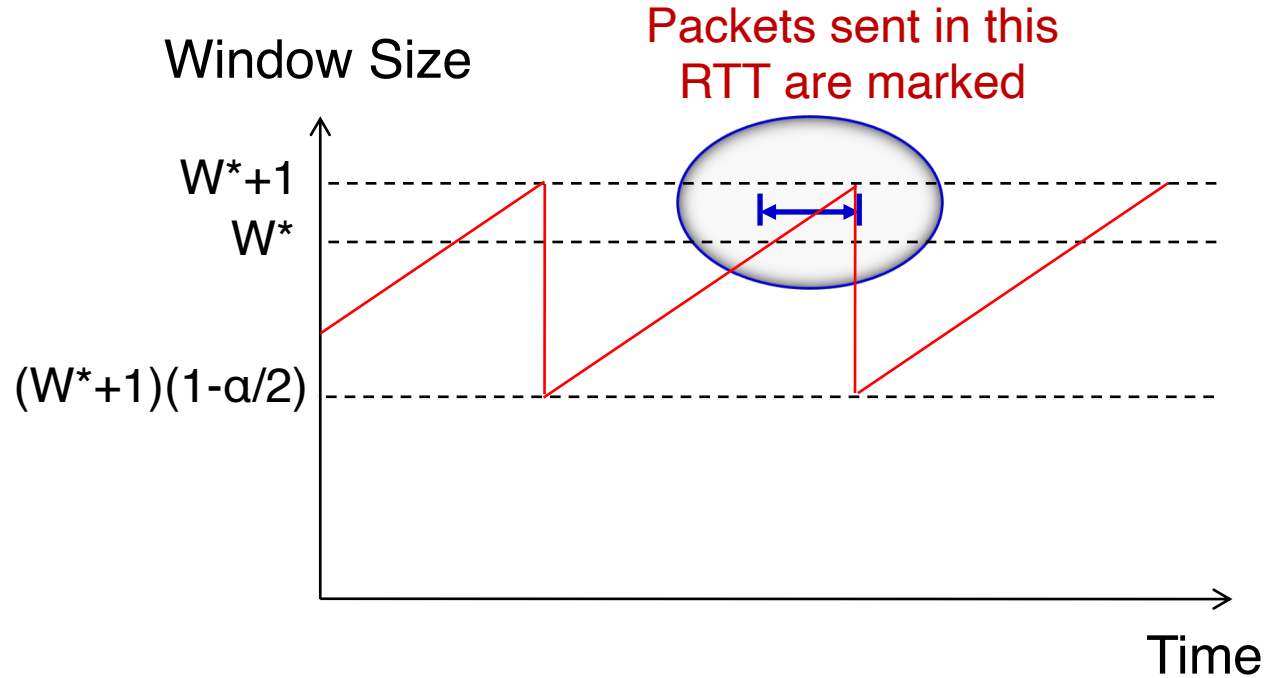
- **Adaptive window decreases:** $W \leftarrow (1 - \frac{\alpha}{2})W$
 - Note: decrease factor between 1 and 2.

Reacting to and controlling queue size **distribution**, specifically, the region above K .

Setting protocol parameters

- When should router start marking? **K**:
 - Mark too late: higher queueing delay (and maybe loss)
 - Mark too early: queues too small, **lose throughput**
 - **Want min queue size > 0 even when TCP windows drop**
- What is the ideal buffer size for DCTCP?
 - Regular TCP: Bandwidth-delay product
 - **Want buffer $>$ max queue size**

Use model to set parameters



$$W^* = (C \times RTT + K) / N$$

$$\alpha = \frac{\text{\# of pkts in last RTT of Period}}{\text{\# of pkts in Period}}$$

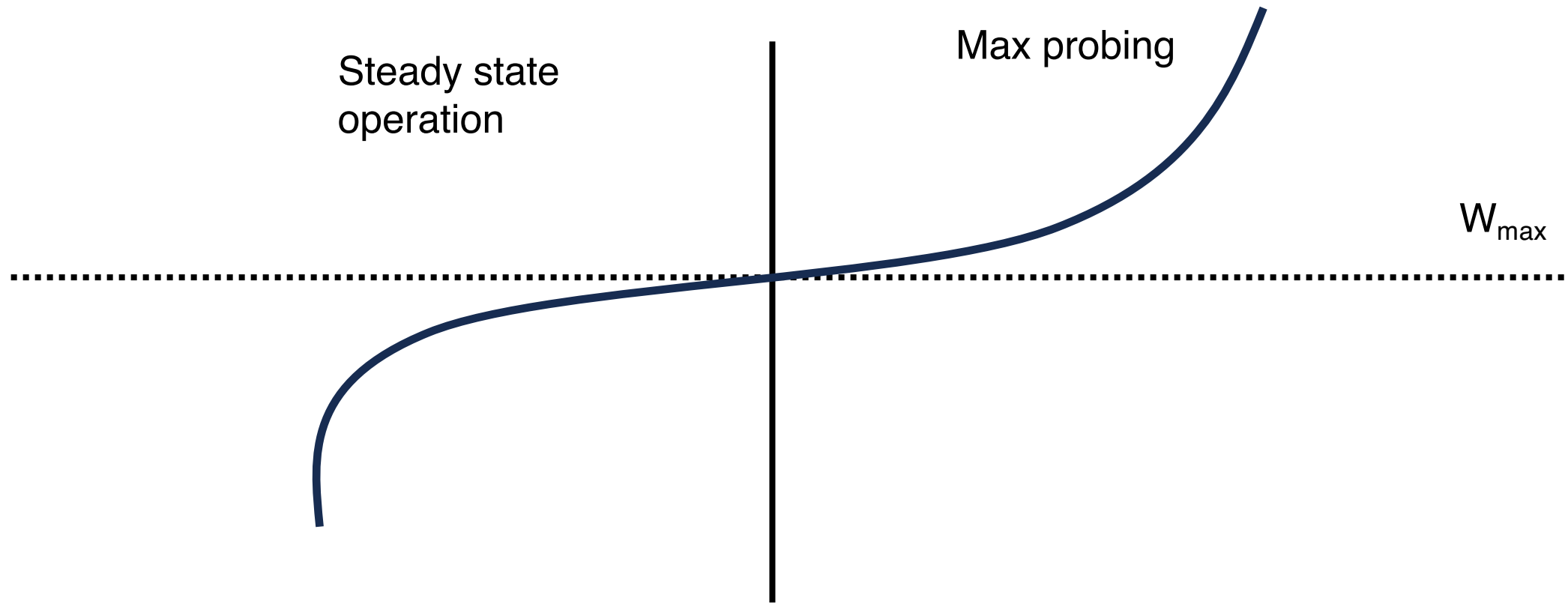
$$\alpha \approx \sqrt{2/W^*}$$

$$Q_{max} = K + N$$

$$Q_{min} = K + N - \frac{1}{2} \sqrt{2N(C \times RTT + K)} \quad K > (C \times RTT) / 7$$

TCP Cubic

Support faster window growth



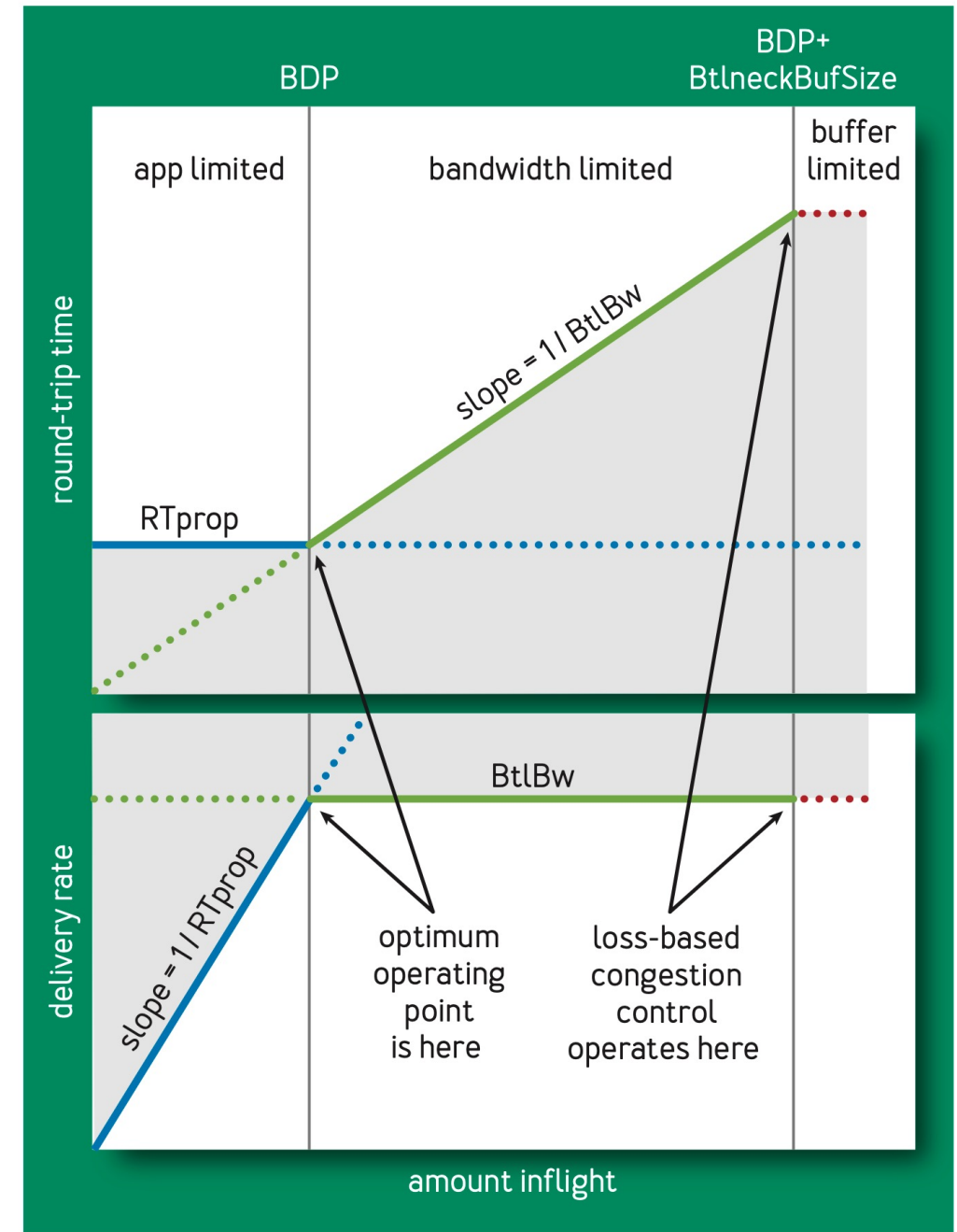
$$W(t) = C(t - K)^3 + W_{max}$$

TCP BBR

Cardwell et al., Google

Principle of Operation

- Find optimal operating point in terms of bottleneck bandwidth and delays
- Cannot simultaneously probe bottleneck bandwidth and propagation delay
- → Occasionally drop window to estimate propagation RTT



Endpoint algorithms alone are
insufficient

The approach that the Internet takes to allocate resources in the network core is to use a **distributed algorithm (congestion control)** running at endpoints.

This allows the Internet to **scale** to a large # of endpoints.

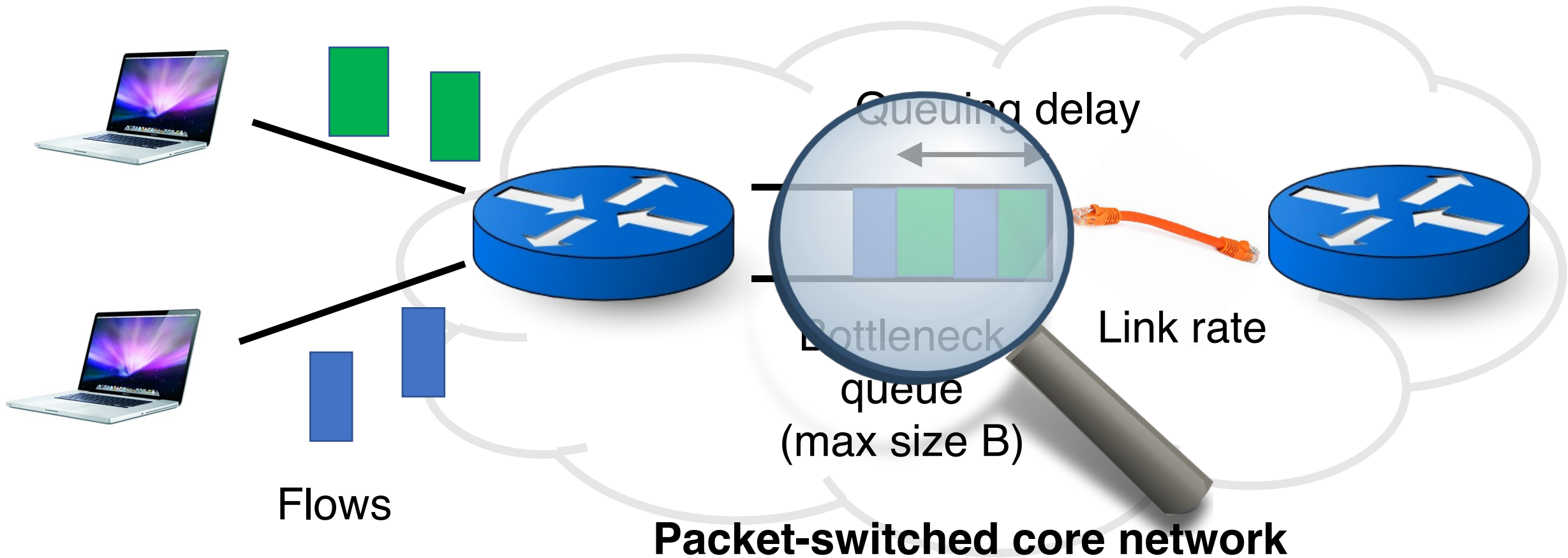
However, it also places **trust in endpoints**.

Uncooperative, buggy, malicious endpoints

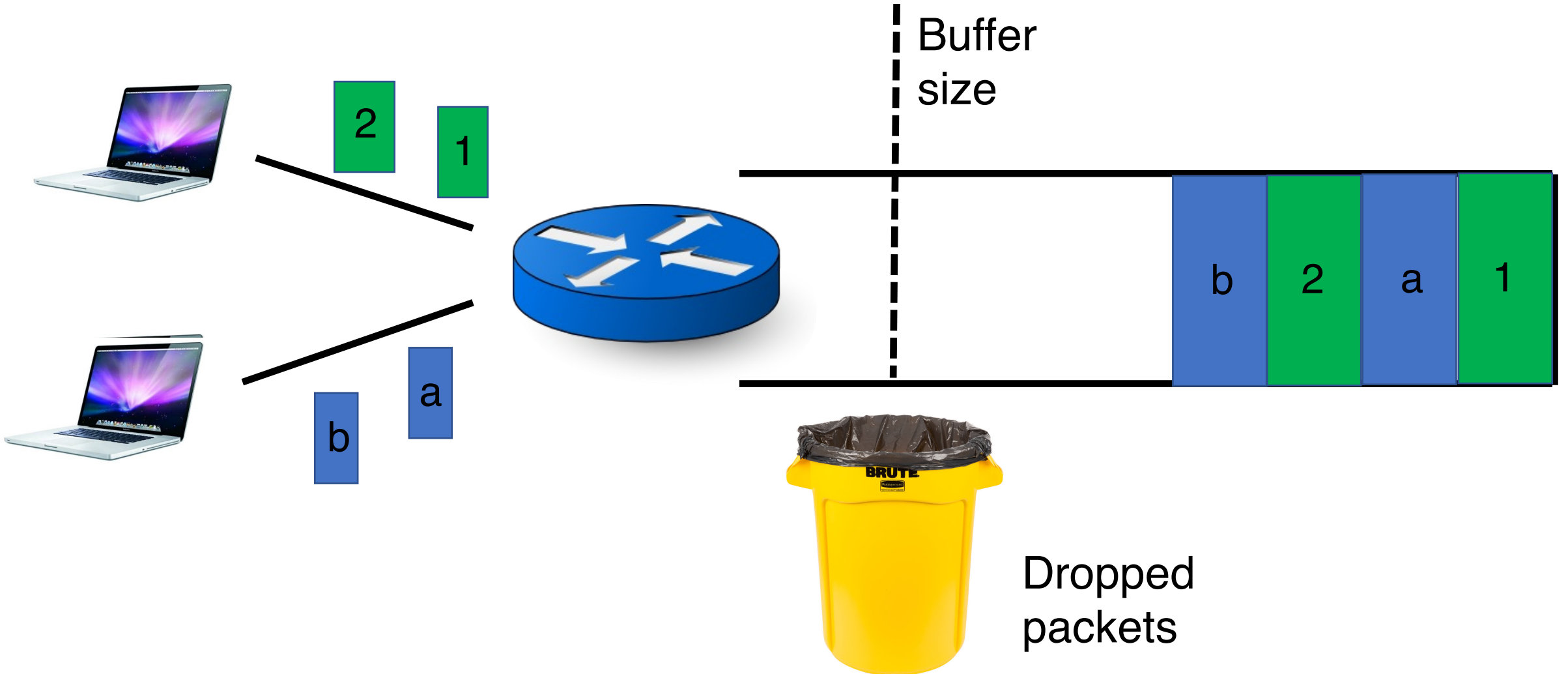


- What if an endpoint is buggy, or malicious?
- We'd like the network core to do something better than best-effort

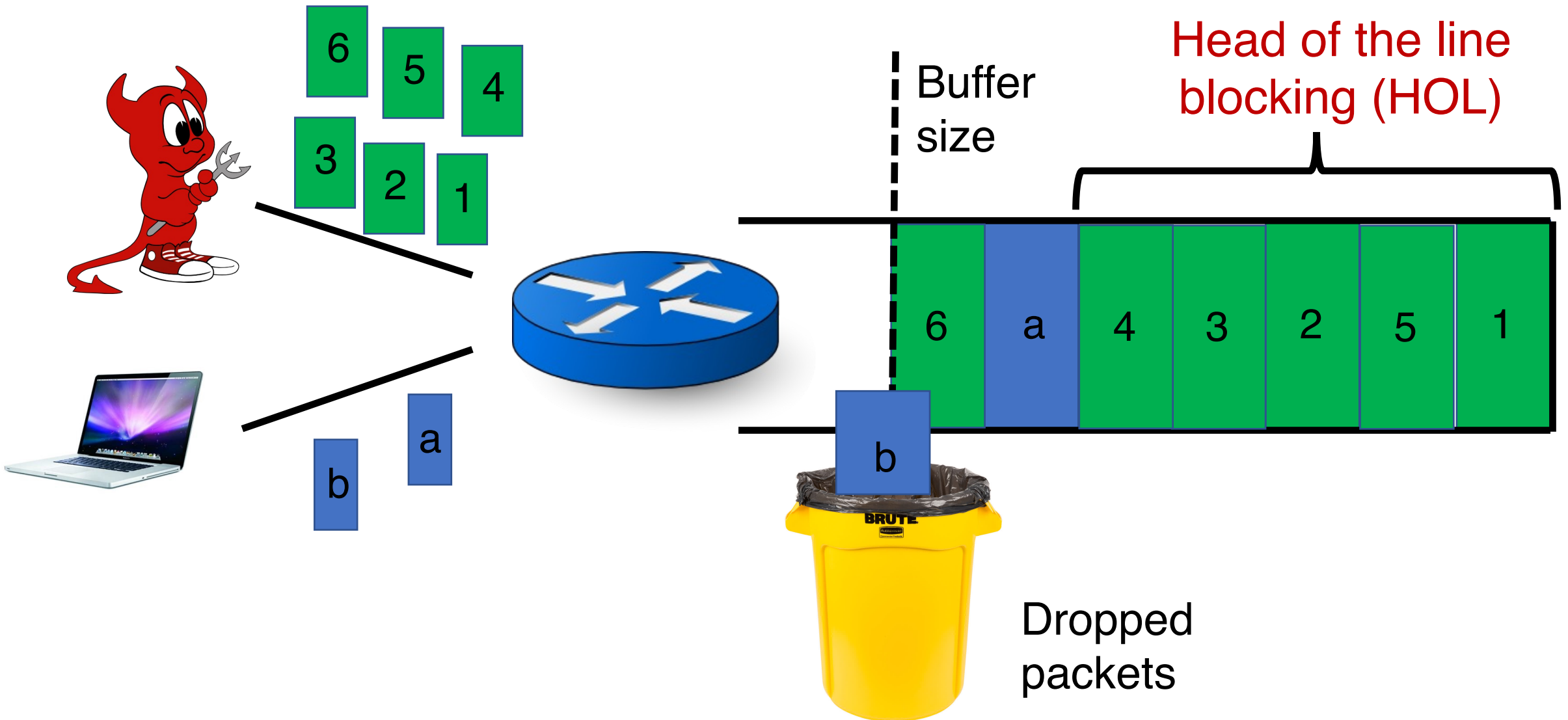
Simplified model of bottleneck link



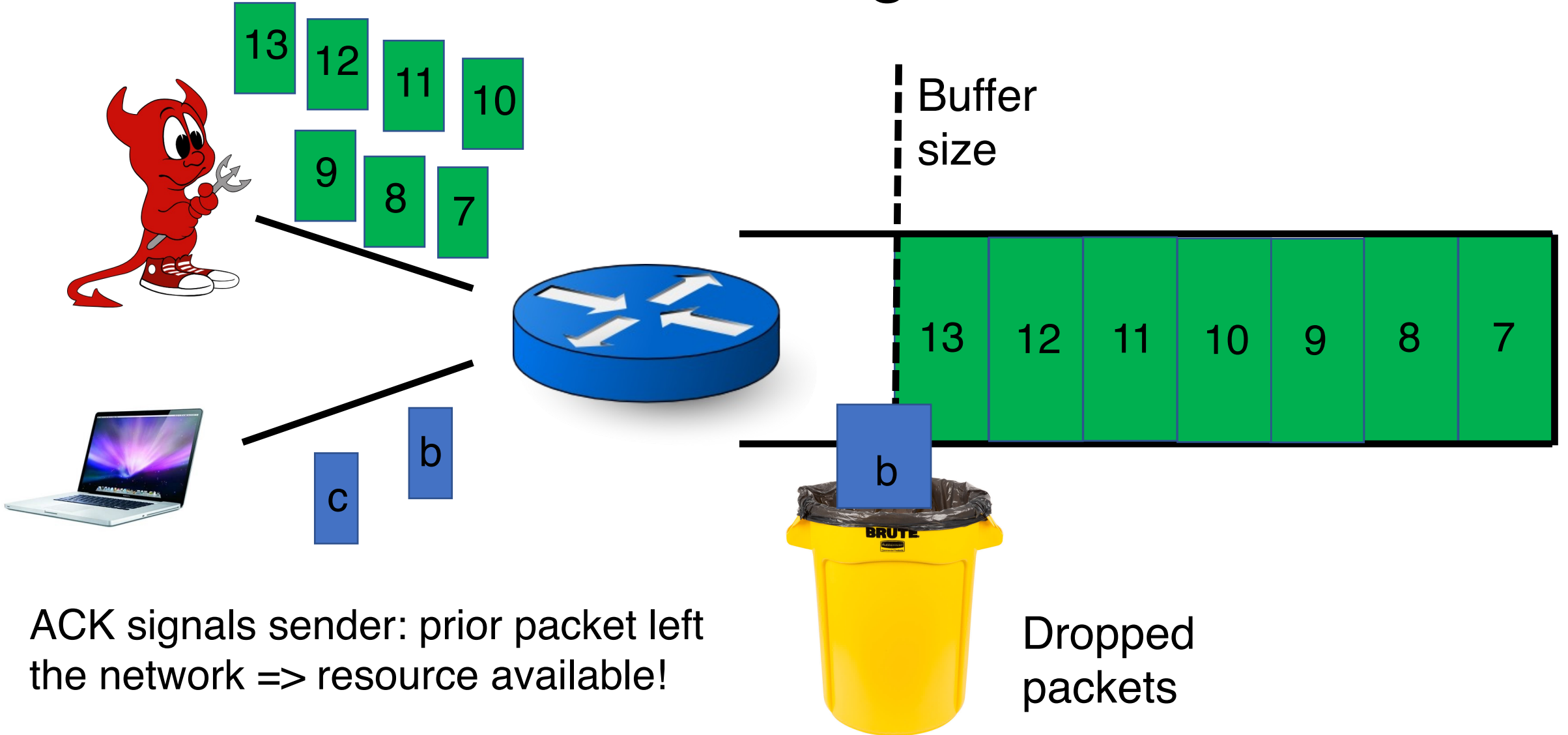
FIFO scheduling + Tail-drop buffer mgmt



FIFO scheduling + Tail-drop buffer mgmt



Next RTT: ACK-clocking makes it worse



Network can be monopolized by bad endpoints

- ACK clocking **synchronizes** senders to when resource is available
 - Conversely, packet losses desynchronize the sender
- Contending packet arrivals may not be random enough
 - e.g., Blue flow can't capture buffer space for *a few* round-trips
- Can observe this effect when many TCP flows compete
 - Some TCP flows can never get off the ground
- A FIFO tail-drop queue incentivizes sources to misbehave

Packet scheduling disciplines @ routers

- Significantly influences how packets are treated regardless of the endpoint's transmissions
 - Implementations of **Quality of Service (QoS)** within large networks
 - Implications for **net neutrality** debates
- Intellectually interesting and influential question
 - Important connections to job scheduling in systems
- Just like in life, how you schedule work is highly impactful

Relationship: scheduling & transport

- Packet scheduling is dealing with things transport has already put into the network
- Transport requires a few round-trip times to react; scheduling does something “immediately”
- If you could schedule transmission out of the endpoint, you could get them to zoom through the network without waiting anywhere in queues
- In modern clusters, goal of transport is to often act as if a centralized scheduler directly chose the packet to transmit from the endpoint (leaving other packets waiting at the endpoint)