

Data Center Transport

Lecture 10

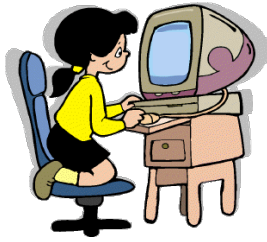
Srinivas Narayana

<http://www.cs.rutgers.edu/~sn624/553-S23>

Parts of this lecture were heavily adapted from slides by Mohammad Alizadeh

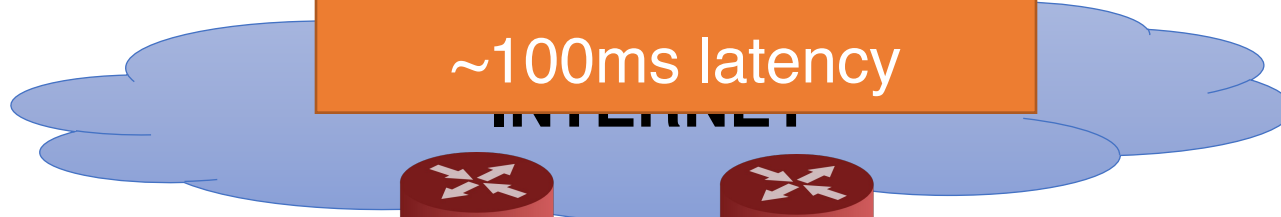
DC Transport Requirements

High throughput, low latency, burst tolerance



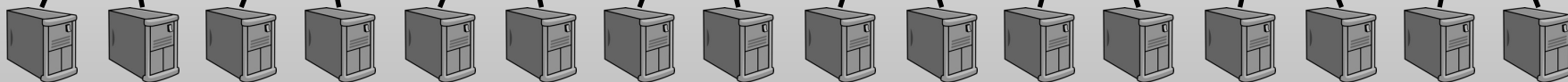
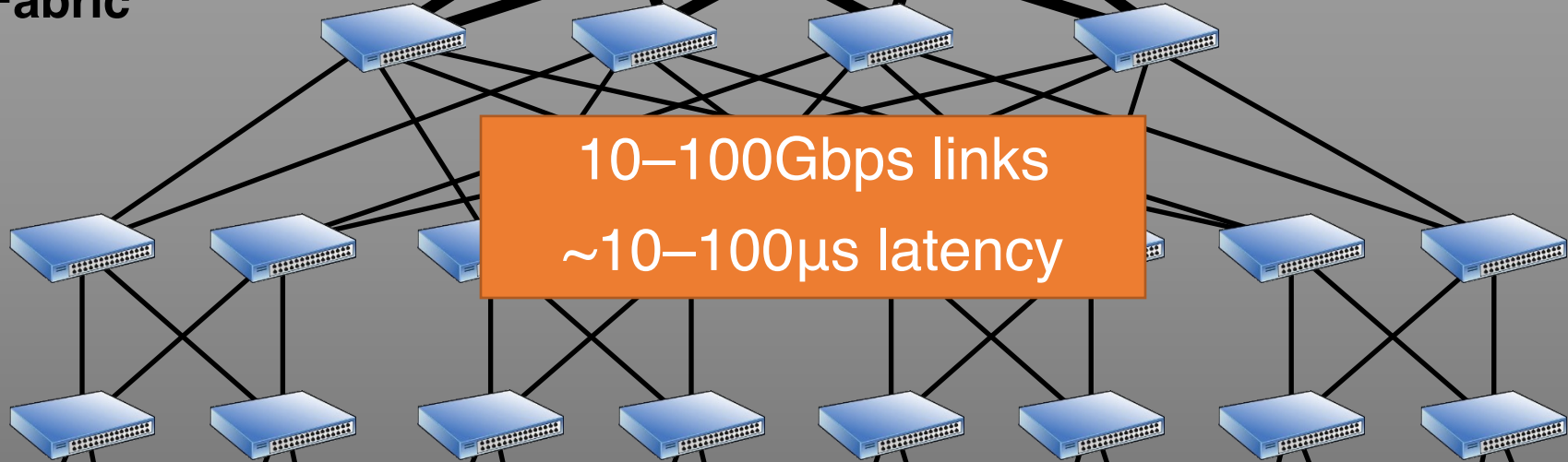
Transport **inside** the DC

100Kbps–100Mbps
links
~100ms latency



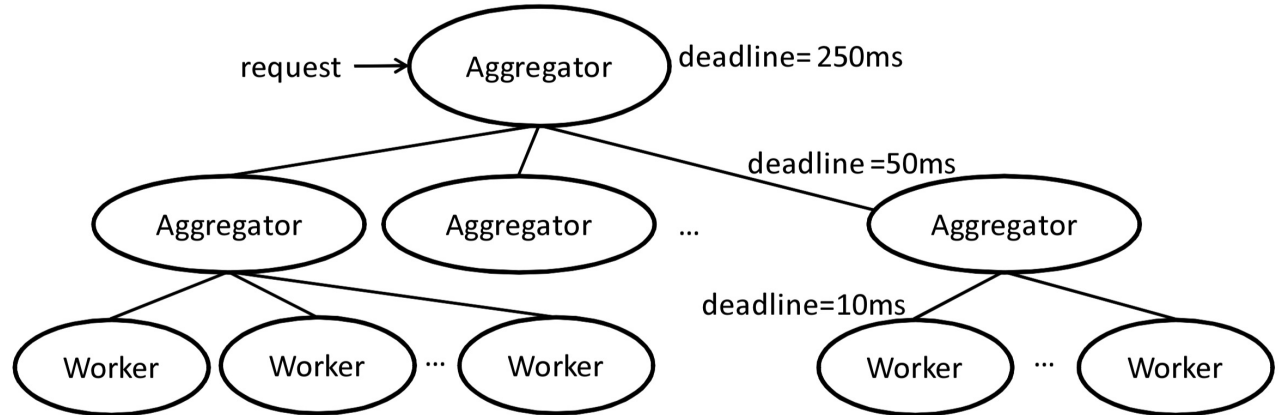
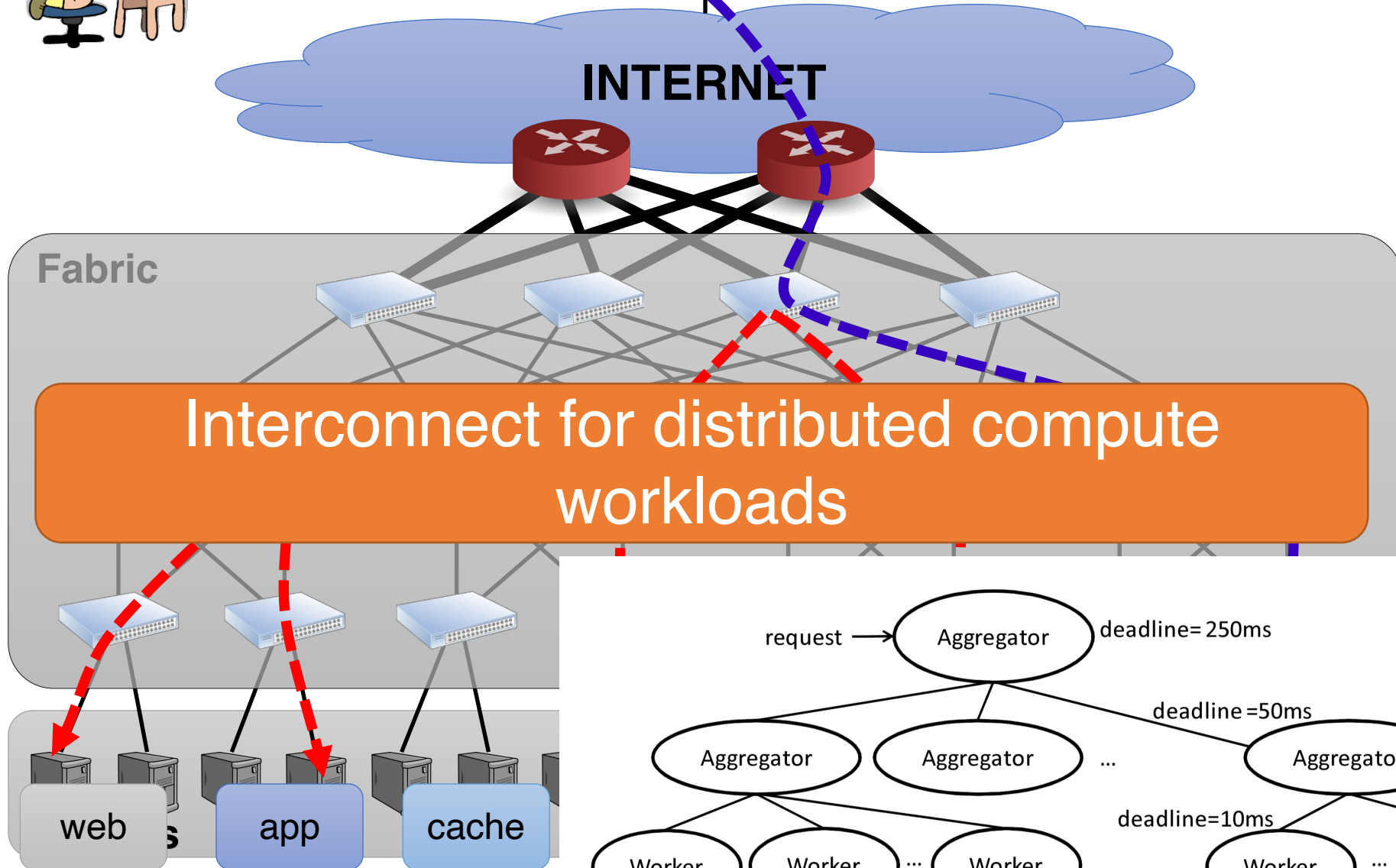
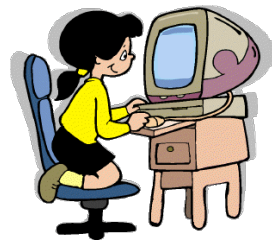
Fabric

10–100Gbps links
~10–100 μ s latency



Servers

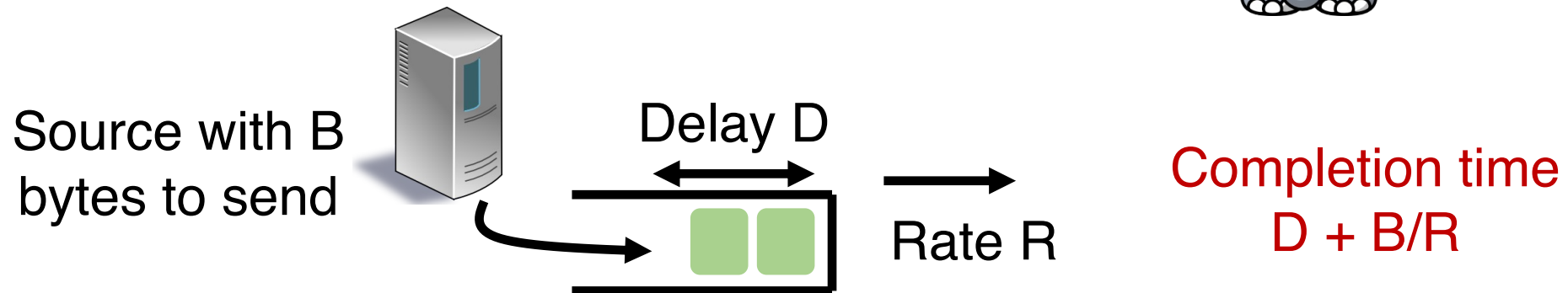
Transport inside the DC



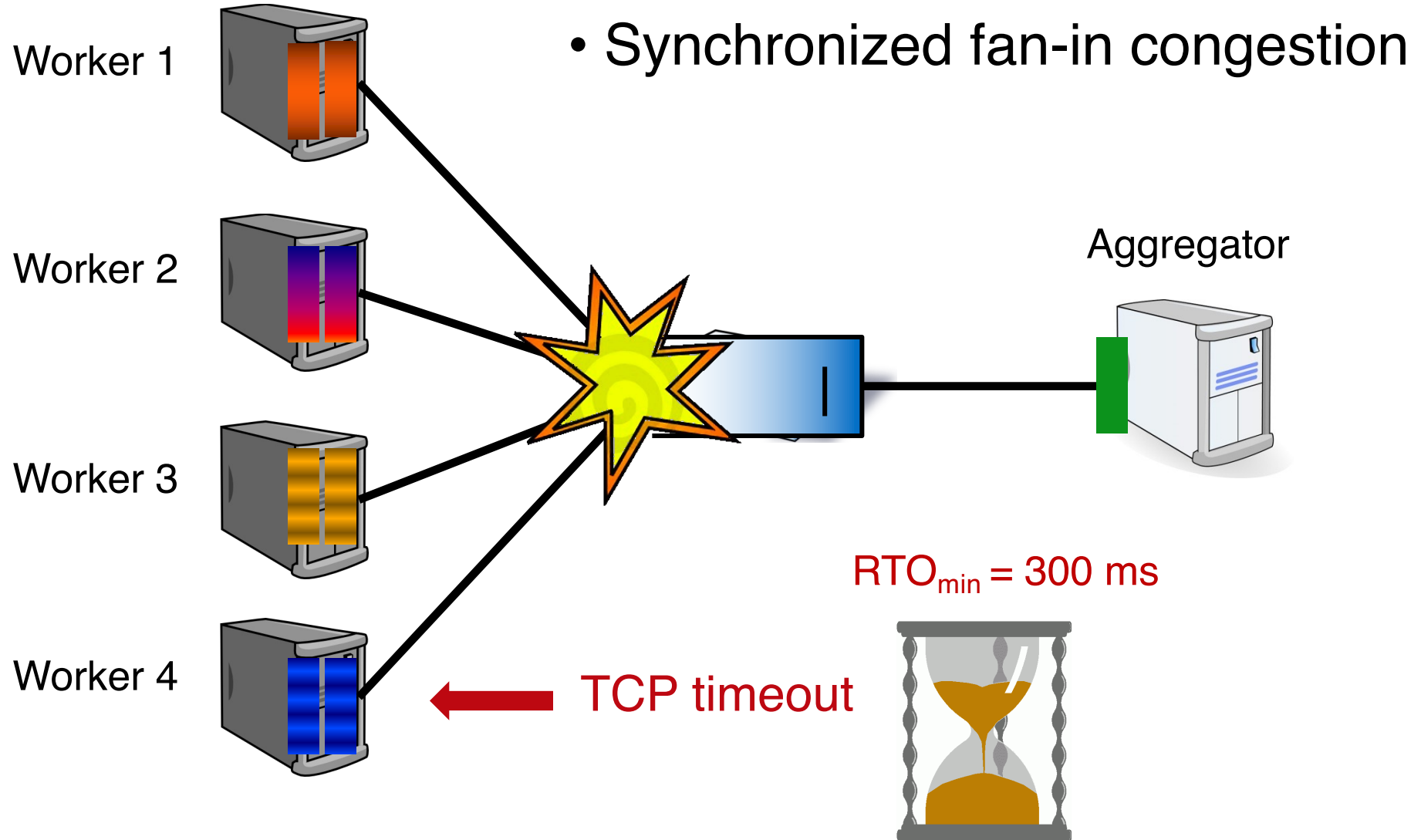
Data center workloads

- Mice and Elephants
- Short messages
(e.g., query, coordination)
- Large flows
(e.g., data update, backup)

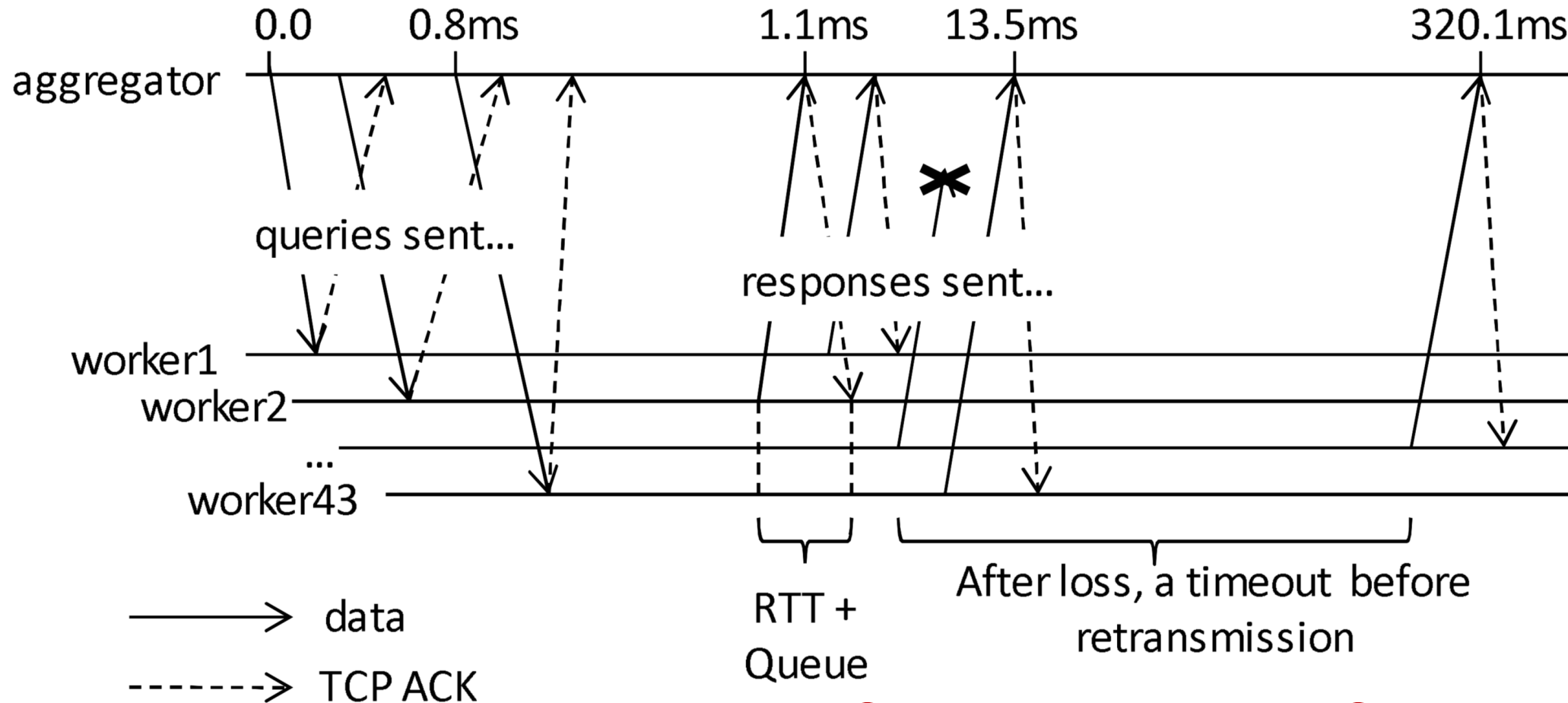
Coexistence creates some performance impairments...



(1) Incast

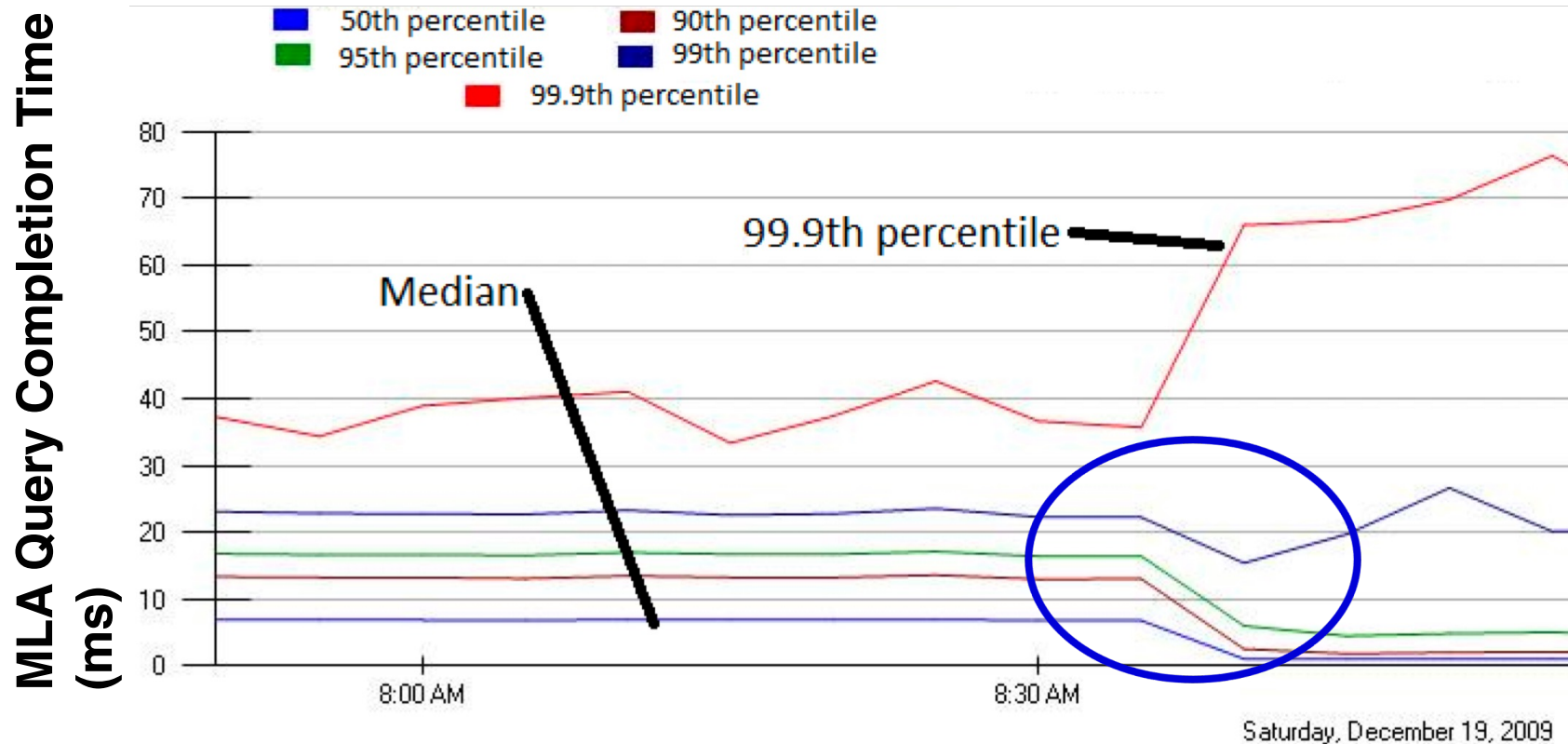


Trace of a real incast event



One option: reduce RTO to mitigate this

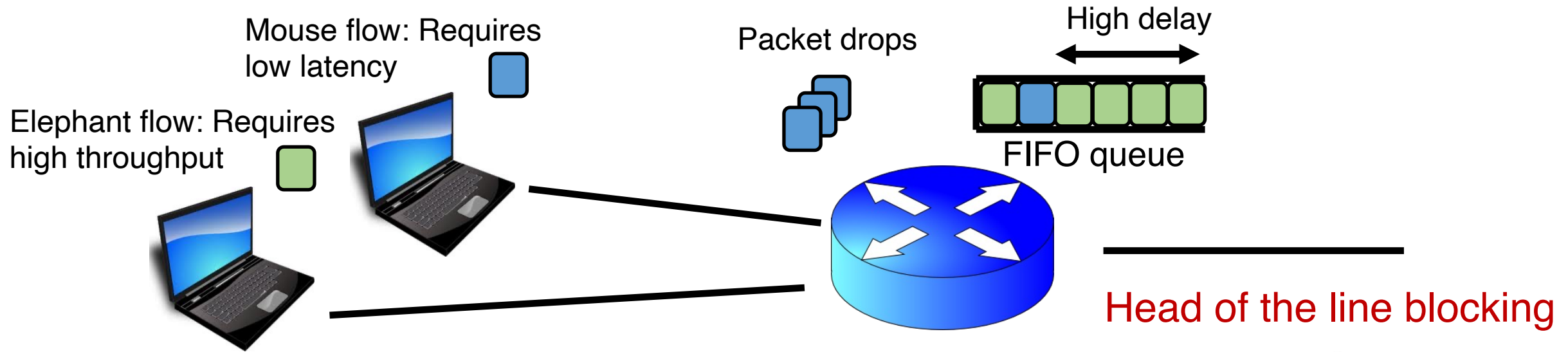
Another option: Jittering to avoid sync





- Jittering switched off around 8:30 am.

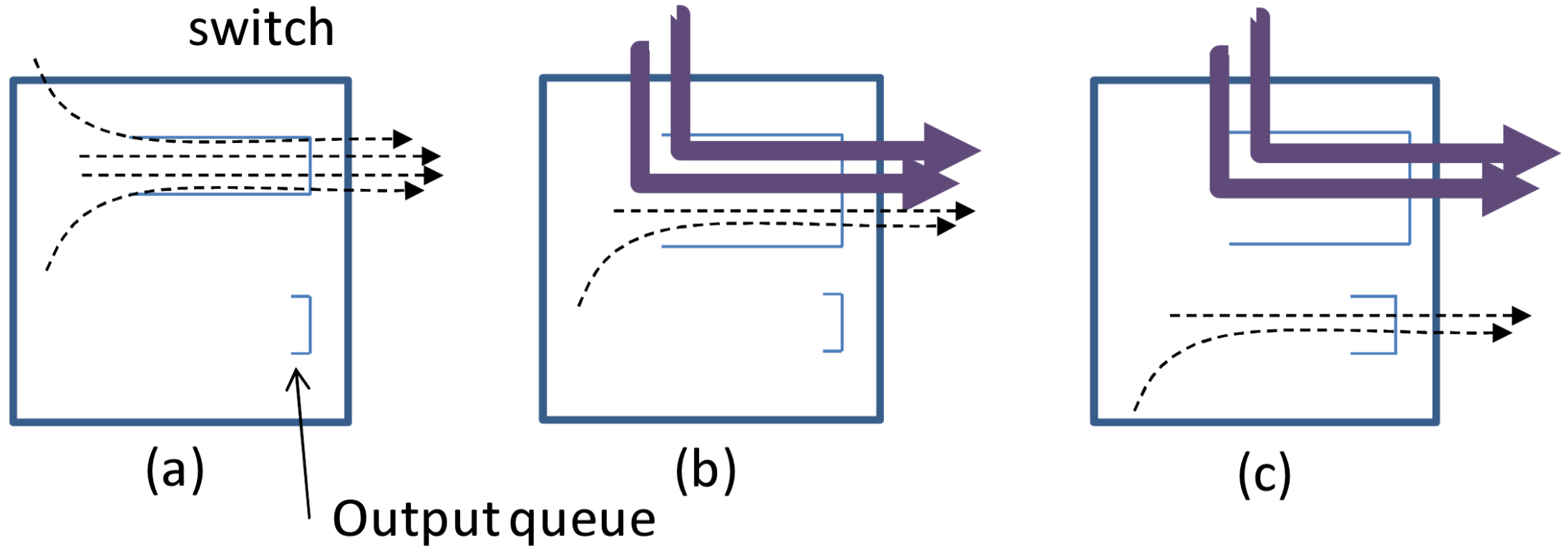
Jittering trades off median for high percentiles

(2) Head of line blocking: Queue buildup



- Resource contention occurs in the **core** of the network 
- Congestion control will react, but may be too little & too late:
 - Congestion control can't prevent packet drops “now”
 - Congestion control won't prevent high-sending-rate flows from inflicting large delays or recurring drops 

(2) HOL Blocking on a switch port



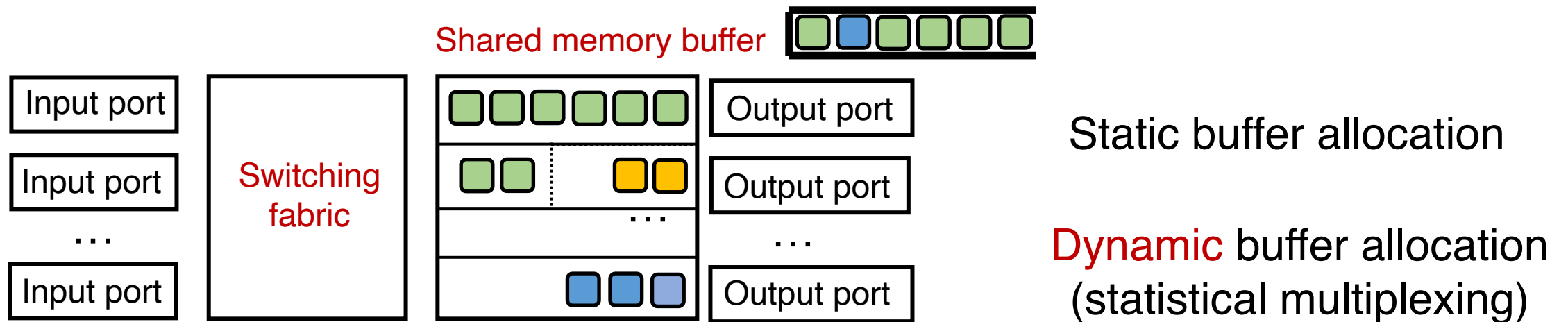
Key:  = large flow  = small flow

Queue buildup on a port increases delays for all flows using that port
(head of the line blocking)

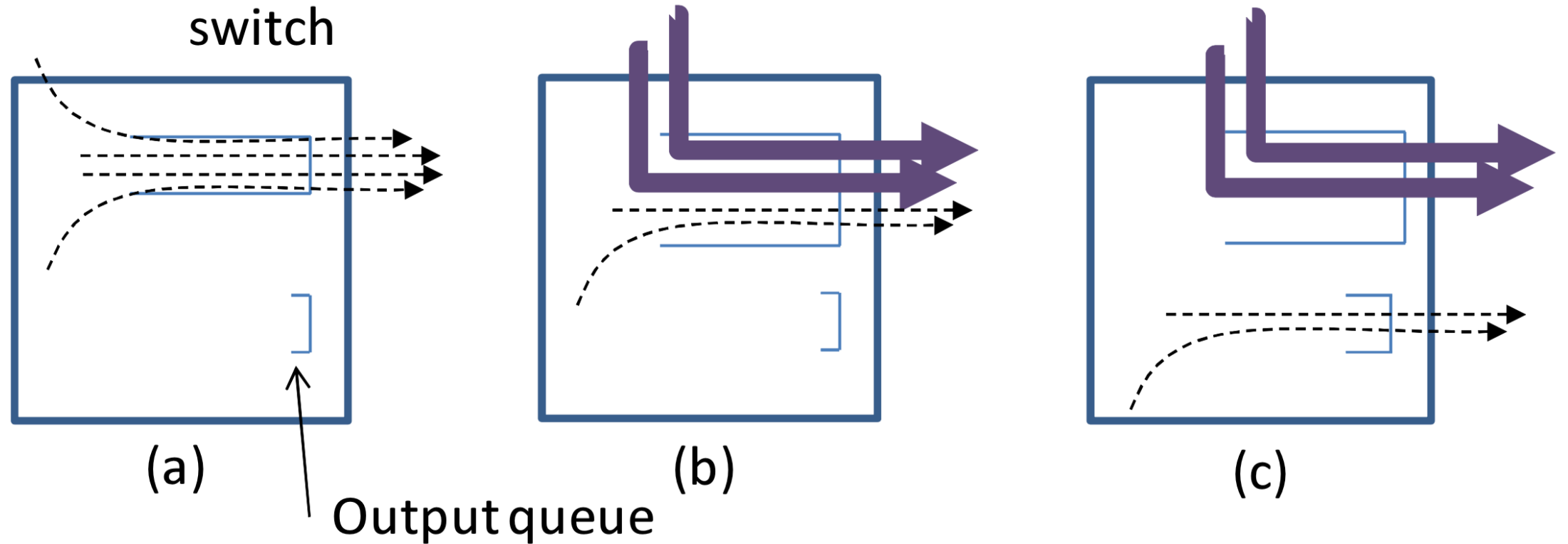
Reducing RTO doesn't help latency when there are no drops

(3) Shared memory buffering

- Where should the packets not currently serviced wait?
- Input-queued vs. **output-queued** (preferable design)
- **Buffer management**: how to put packets into the buffer
- **Scheduling**: how to schedule packets leaving the buffer



(3) Buffer Pressure



Key:  = large flow  = small flow

Shared memory buffering means that queues building up on a **different** port can also impact flows.

Need to keep queues small. Use delay-based CC?

- Keep just a few packets in queues by observing **delays**

$$\text{queue_use} = \text{cwnd} - \text{BWE} \times \text{RTT}_{\text{noLoad}} = \text{cwnd} \times (1 - \text{RTT}_{\text{noLoad}} / \text{RTT}_{\text{actual}})$$

- Adjust window such that only a few packets are in queue

$$\alpha \leq \text{queue_use} \leq \beta$$

- **RTT estimates need to be very accurate and precise**
 - Difficult in low-RTT data centers.
 - Challenges: Software queueing & scheduling delays. Timer tick res

Data Center TCP (DCTCP)

Design of the congestion control algorithm

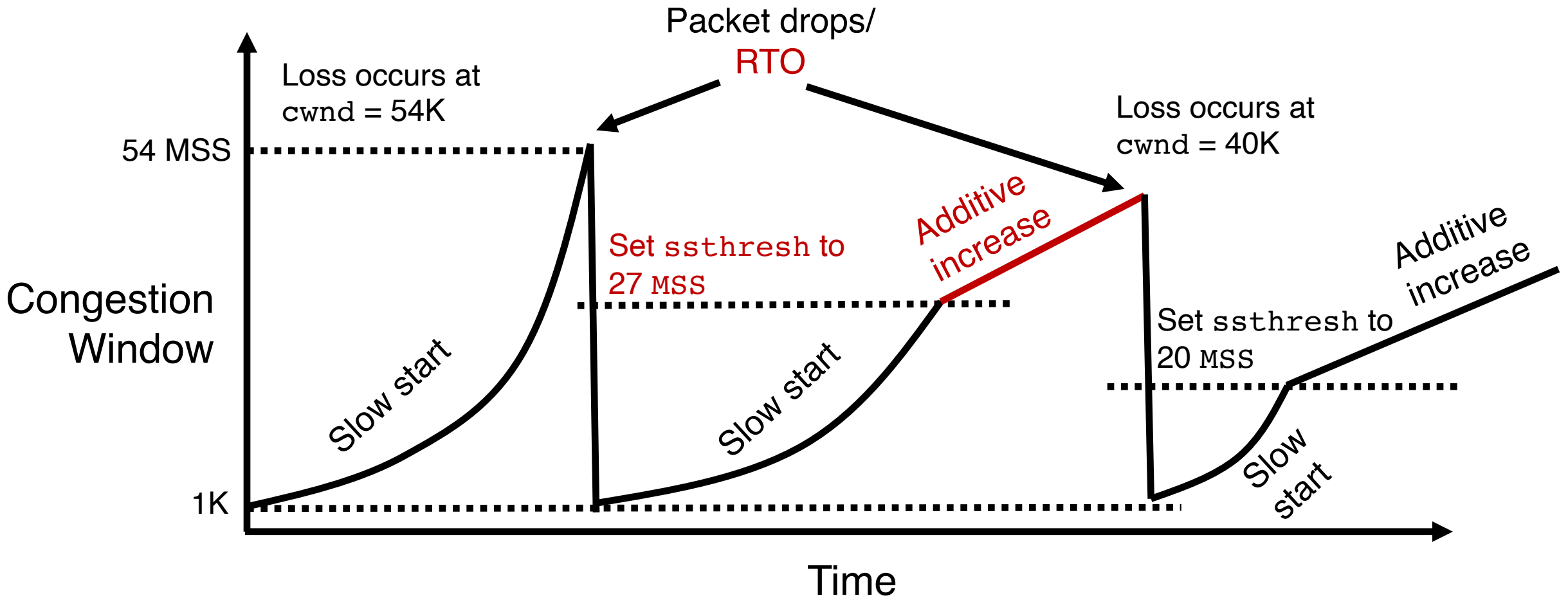
Review: TCP congestion control

- Keep some **in-flight** (un-ACK'ed) packets: **congestion window**
- Adjust window based on several algorithms: TCP New Reno:
 - Startup: slow start
 - Steady state: AIMD
 - Loss: fast retransmission, fast recovery
- Main question for this lecture:
 - **(How) should this design change for data centers?**

Behavior of Additive Increase

Say MSS = 1 KByte

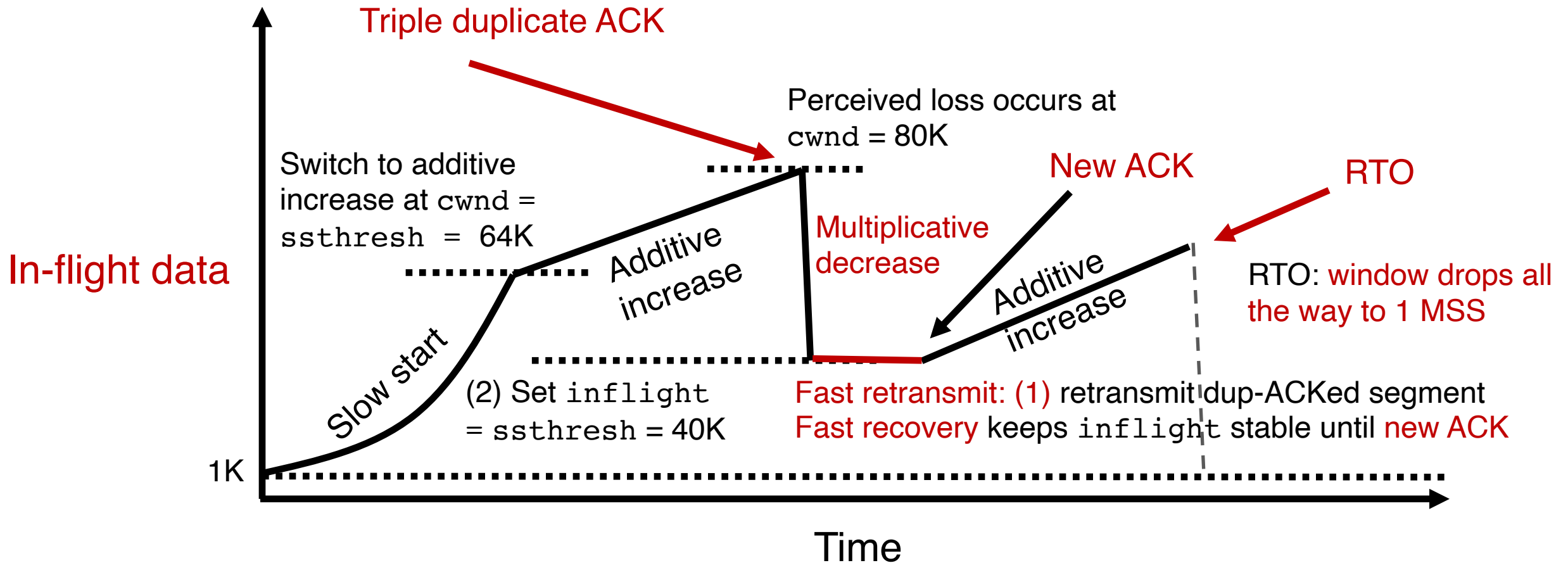
Default ssthresh = 64KB = 64 MSS



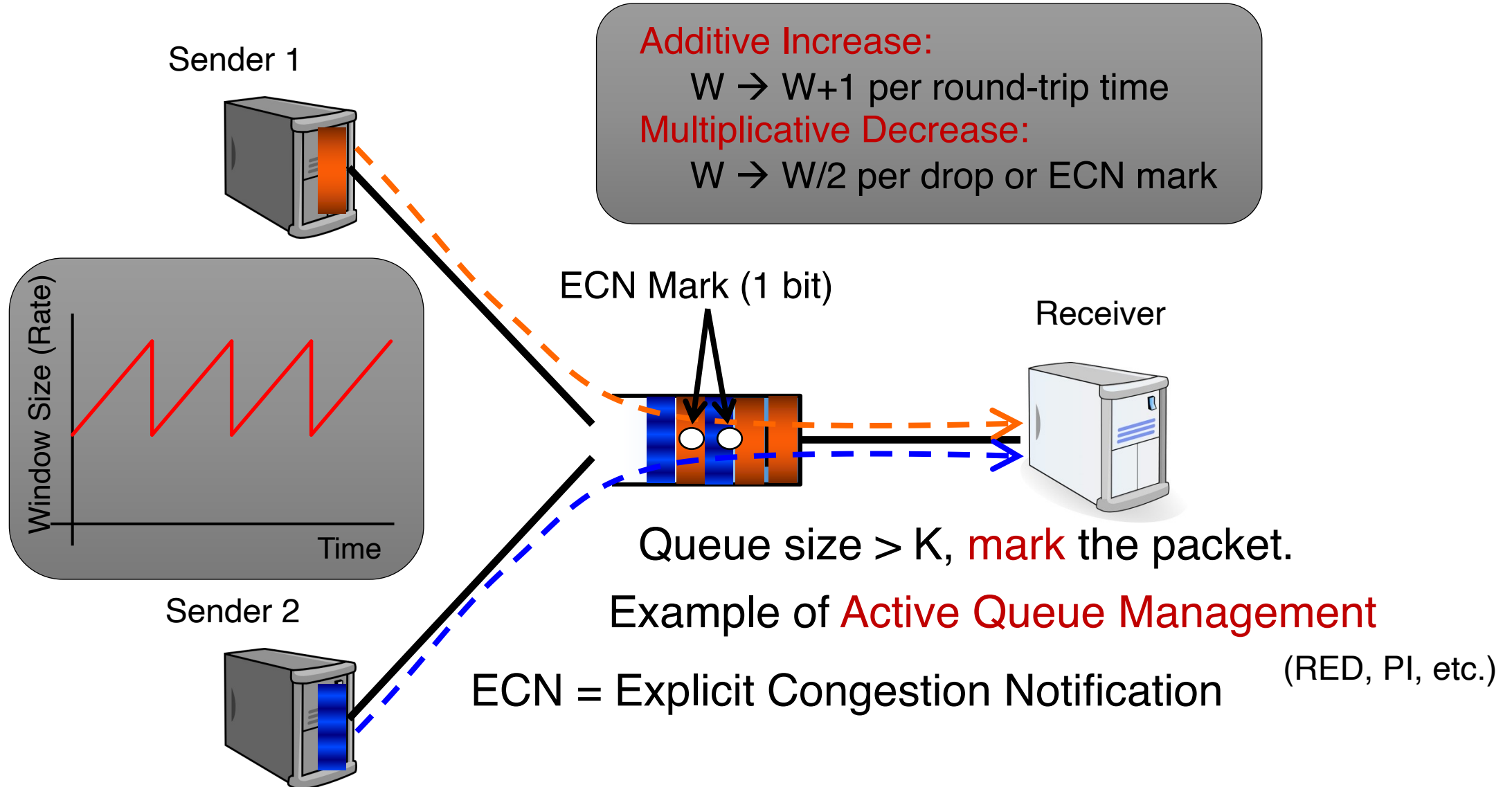
Additive Increase/Multiplicative Decrease

Say MSS = 1 KByte


Default ssthresh = 64KB = 64 MSS



Explicit Congestion Notification



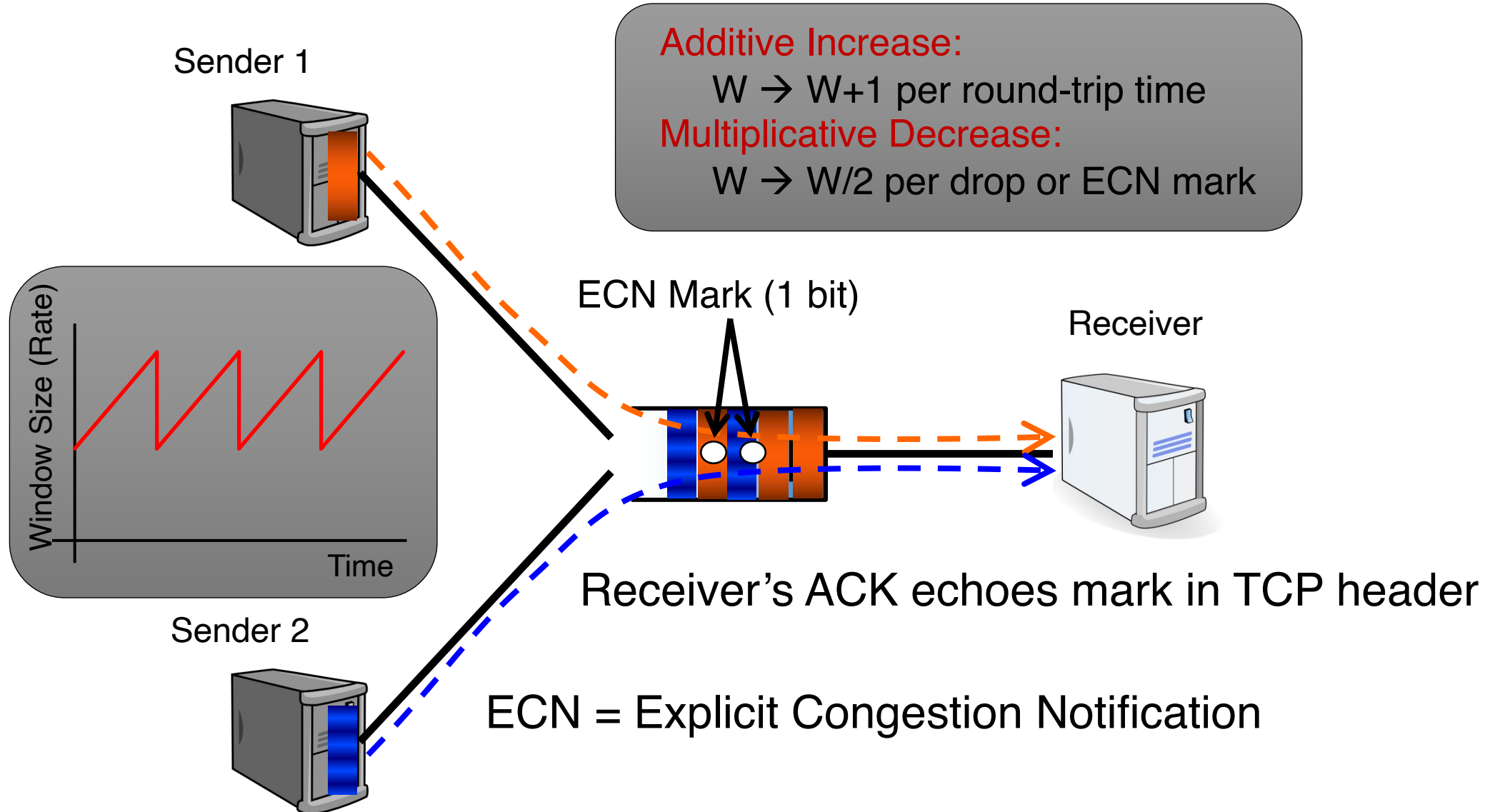
ECN set on the IP header by routers

- 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** 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

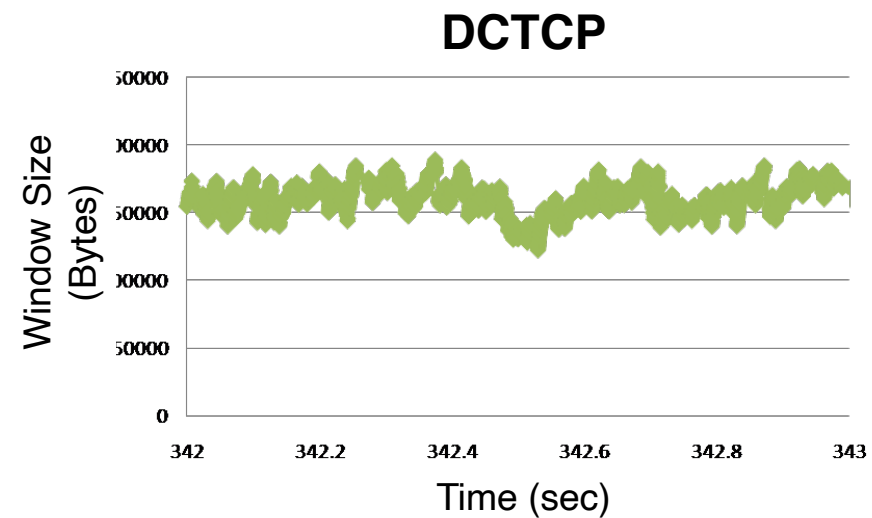
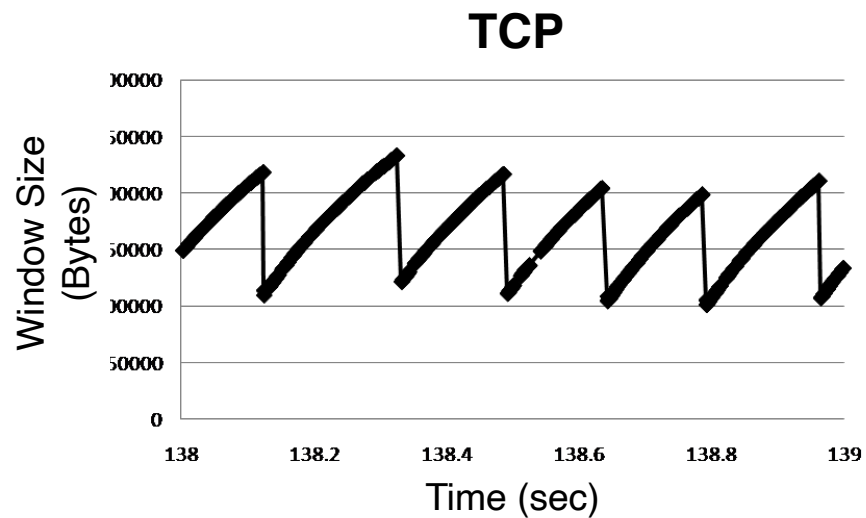
Offsets	Octet	0								1								2								3							
		7	6	5	4	3	2	1	0	7	6	5	4	3	2	1	0	7	6	5	4	3	2	1	0	7	6	5	4	3	2	1	0
0	0	Source port																Destination port															
4	32	Sequence number																															
8	64	Acknowledgment number (if ACK set)																															
12	96	Data offset	Reserved 0000	C W R	E C R E	U R G	A C 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

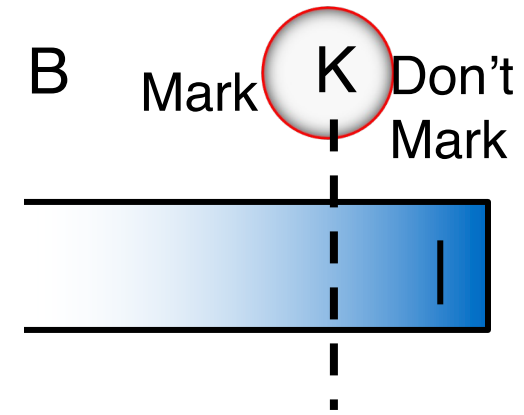
ECN Marks	TCP	DCTCP
1 0 1 1 1 1 0 1 1 1	Cut window by 50%	Cut window by 40%
0 0 0 0 0 0 0 0 0 1	Cut window by 50%	Cut window by 5%



DCTCP algorithm

Switch 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:** $W \leftarrow (1 - \frac{\alpha}{2})W$

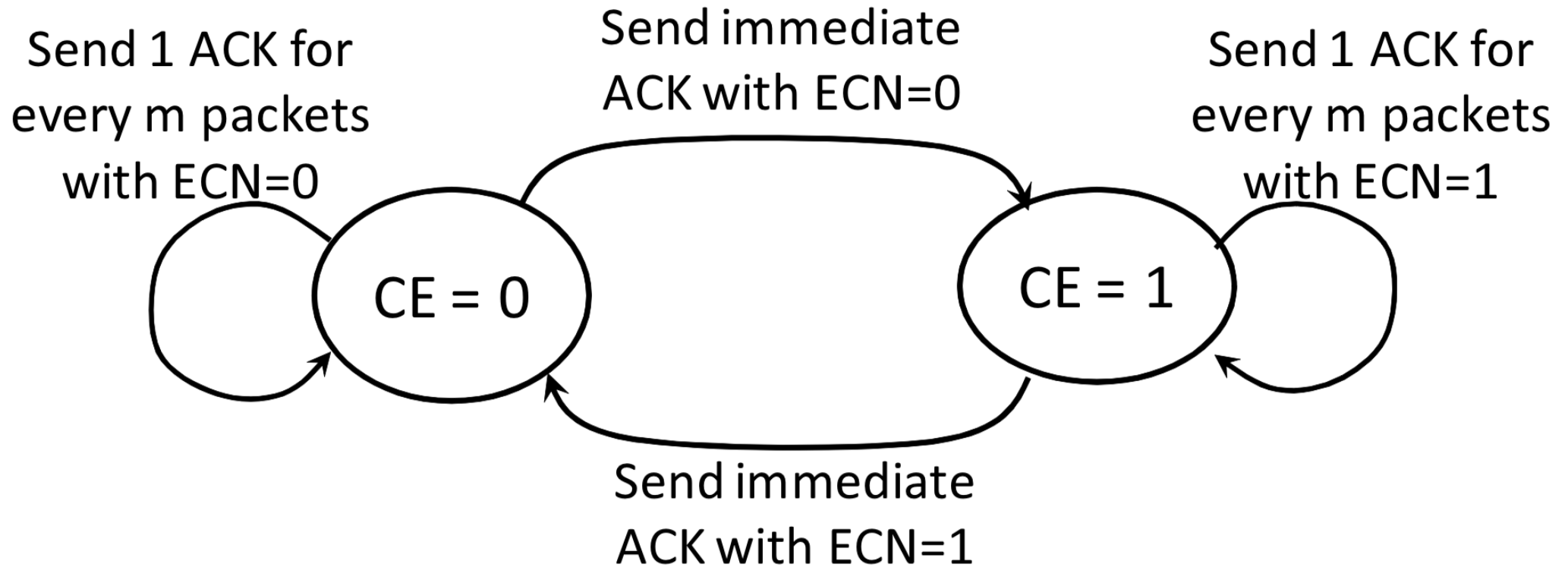
- Note: decrease factor between 1 and 2.

Reacting to and controlling queue size **distribution**

Delayed ACKs

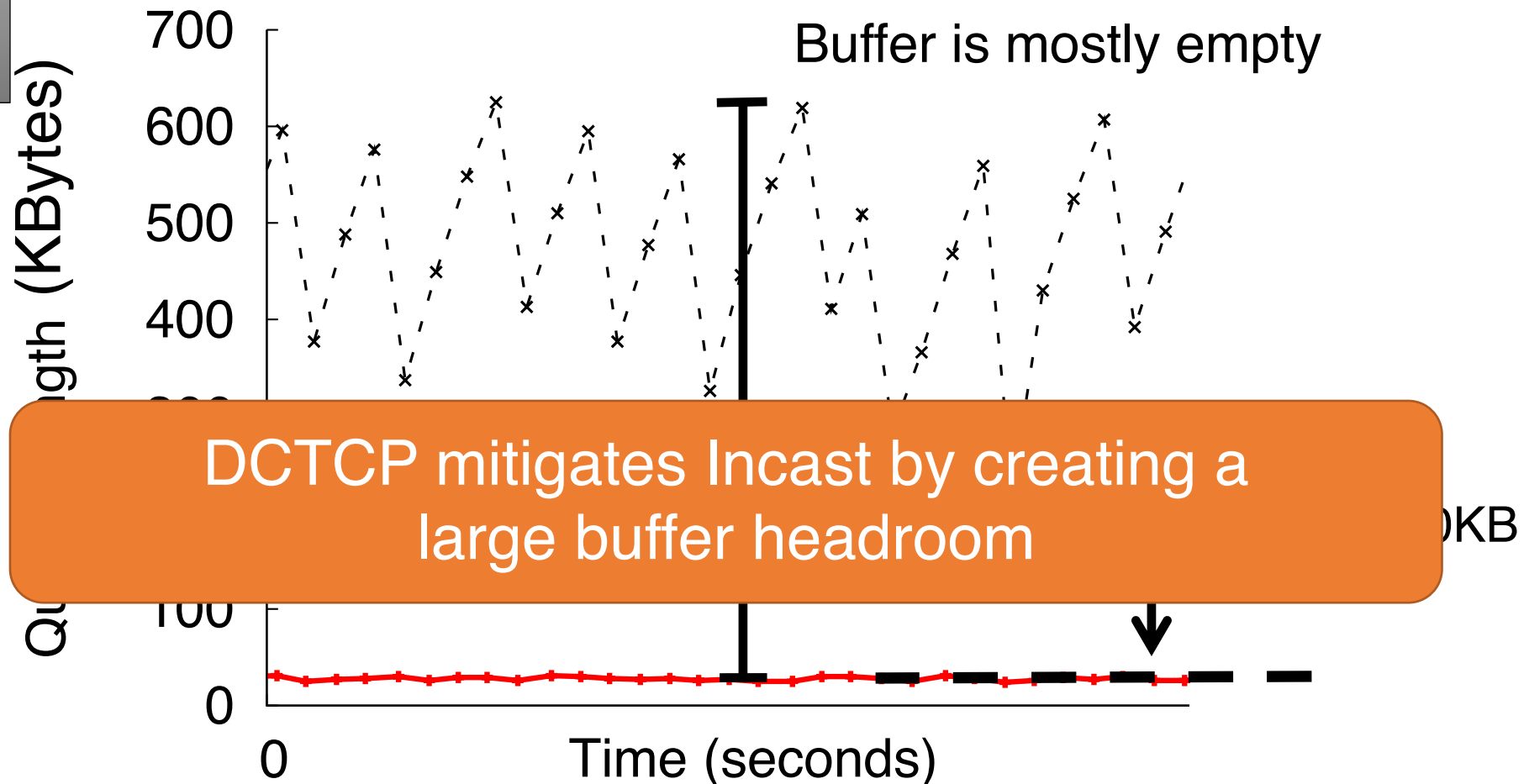
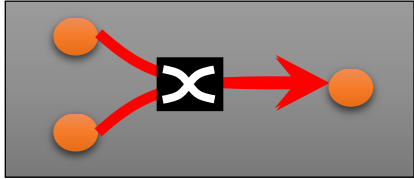
- Not every packet is ACKnowledged by receiver
- Too many ACKs: increase packet processing load
 - Typical policy: ACK every m packets, or after sender has paused transmitting for a delayed ACK timeout
- How to allow the sender to see the full stream of ECN marks?

Efficient and “lossless” ACK generation



DCTCP vs TCP

Experiment: 2 flows (Win 7 stack), Broadcom 1Gbps Switch



Why it works

1. Low Latency

- ✓ Small buffer occupancies → low queuing delay

2. High Throughput

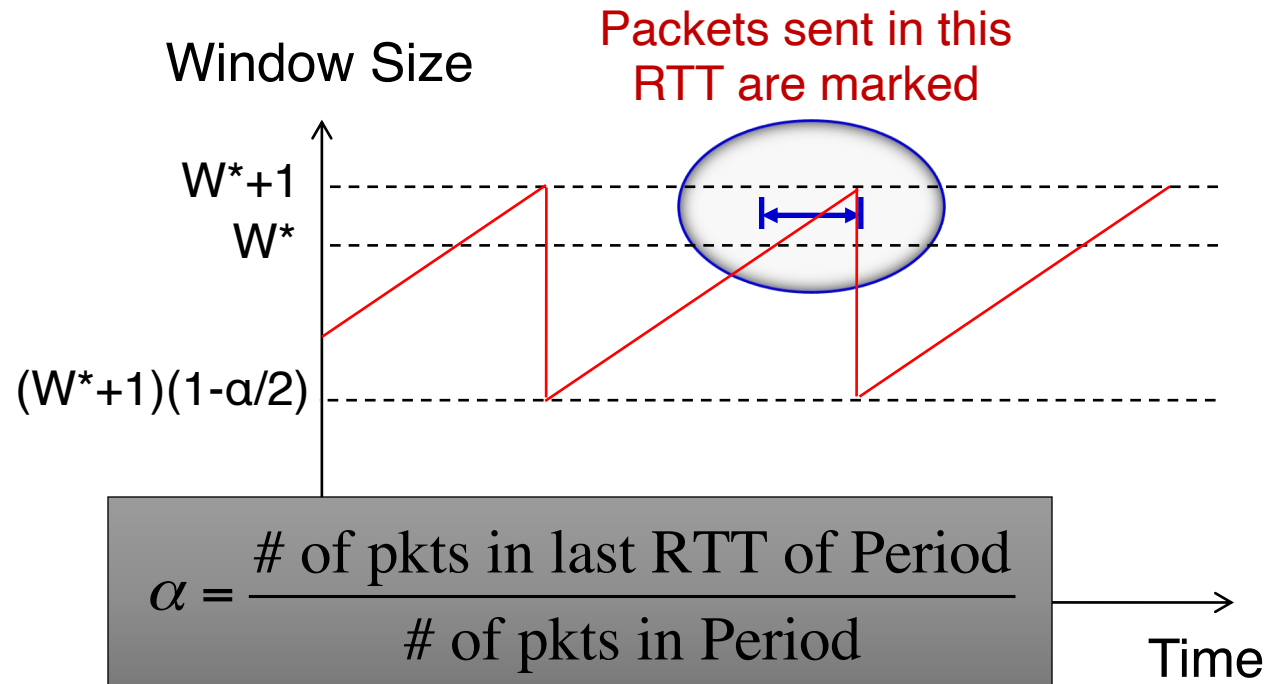
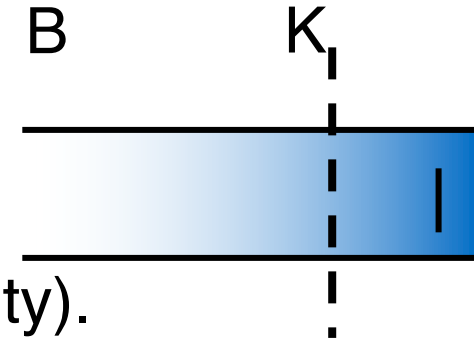
- ✓ ECN averaging → smooth rate adjustments, low variance

3. High Burst Tolerance

- ✓ Large buffer headroom → bursts fit
- ✓ Aggressive marking → sources react before packets are dropped

Setting parameters: A bit of analysis

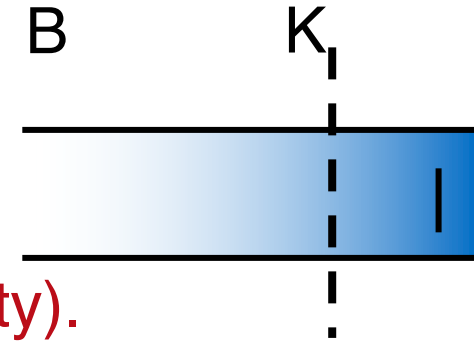
- How much buffering does DCTCP need for 100% throughput?
- Need to quantify queue size **oscillations** (stability).



Setting parameters: A bit of analysis

- How small can queues be without loss of throughput?

➤ Need to quantify queue size oscillations (Stability).



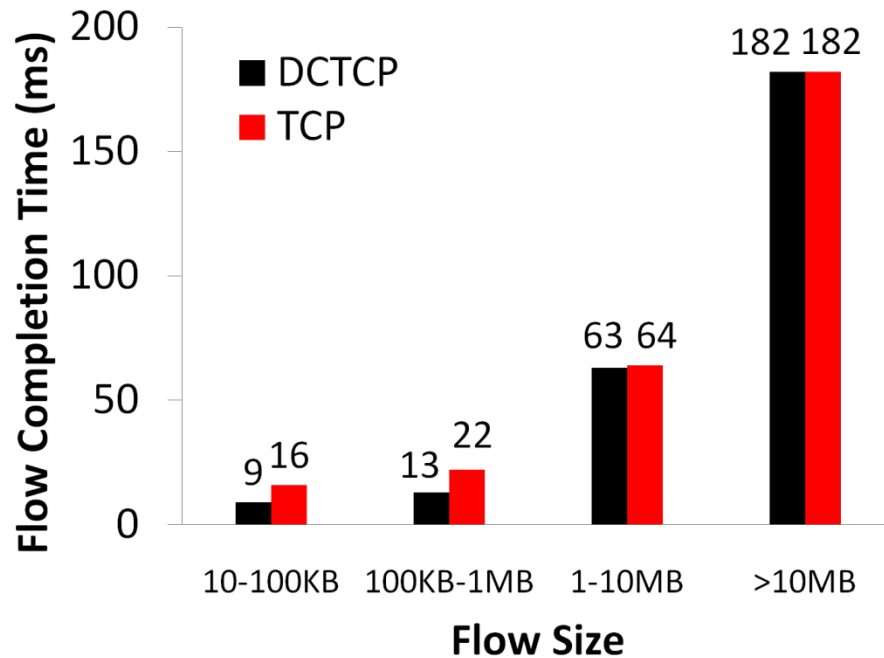
$$K > (1/7) C \times RTT$$



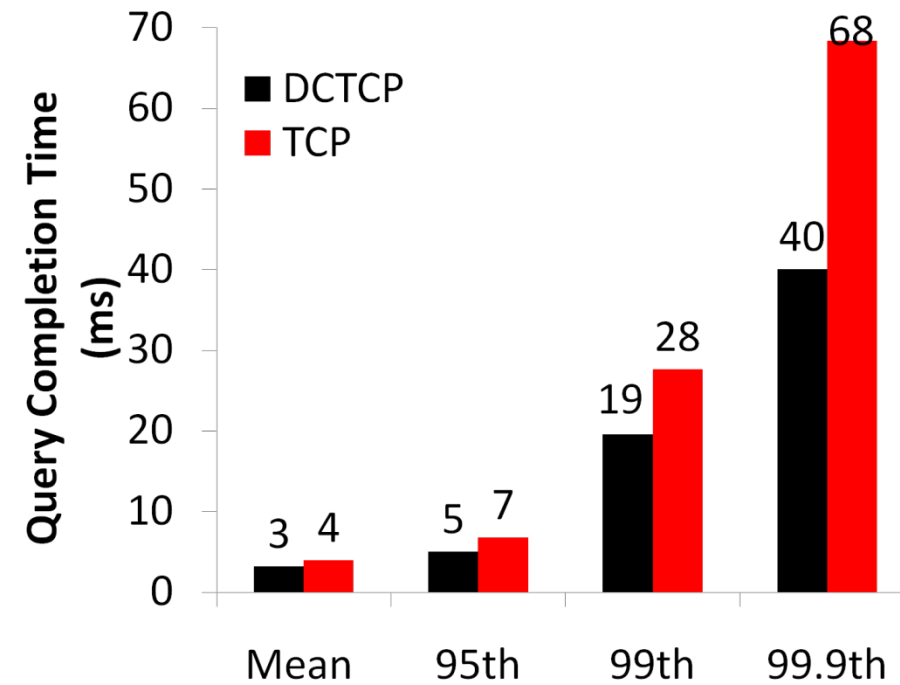
$$\text{for TCP:} \\ K > C \times RTT$$

Bing benchmark (baseline)

Background Flows

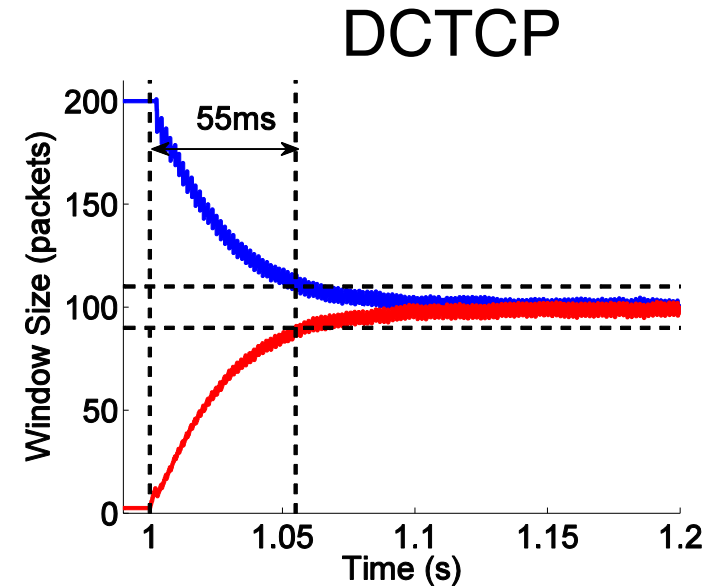
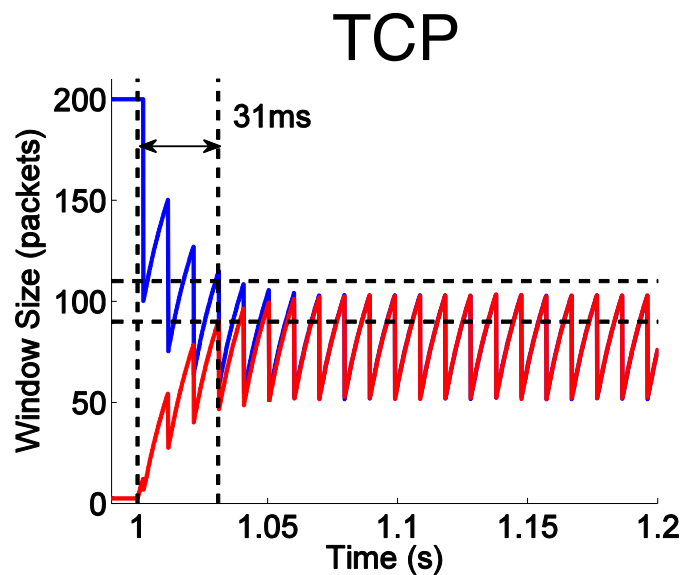


Query Flows



Convergence time

- DCTCP takes at most ~40% more **RTTs** than TCP
 - “Analysis of DCTCP”, SIGMETRICS 2011
- **Intuition:** DCTCP makes smaller adjustments than TCP, but makes them much more frequently



CC evaluation: many aspects to consider

- Throughput, delays, flow completion times
- Fairness, convergence times
- Specific impairments:
 - incast (many to one, all to all)
 - Queue buildup
 - Buffer pressure
 - Collateral damage from incast
- Multi-hop versus single-hop bottlenecks
- Comparison against existing TCPs and AQMs
- How deployable: app awareness, hardware compatibility, ...

CC Deployment Concerns

Life isn't easy in the fast lane

Practical deployment concerns in DCs

- **Coexistence with legacy protocols like TCP Cubic**
 - Application code can't be upgraded in one shot
- **Minimum window size matters during heavy incast events**
 - e.g., 2 packets versus 1 packet: no reactive scheme can work if buffers are so small that drop in one RTT
- **Enabling appropriate options at senders, receivers, and routers**
 - Non “ECN-capable” flagged packets will be dropped when $Q > K$
 - ... including the SYN packets of any connection
- **Receive-side buffer tuning**
 - Reacting to increasing buffer demands at the endpoints takes time
 - Static: Usually, receive buffer must be at least BDP; also influenced significantly by queueing