# L12: end to end layer

Dina Katabi

6.033 Spring 2007
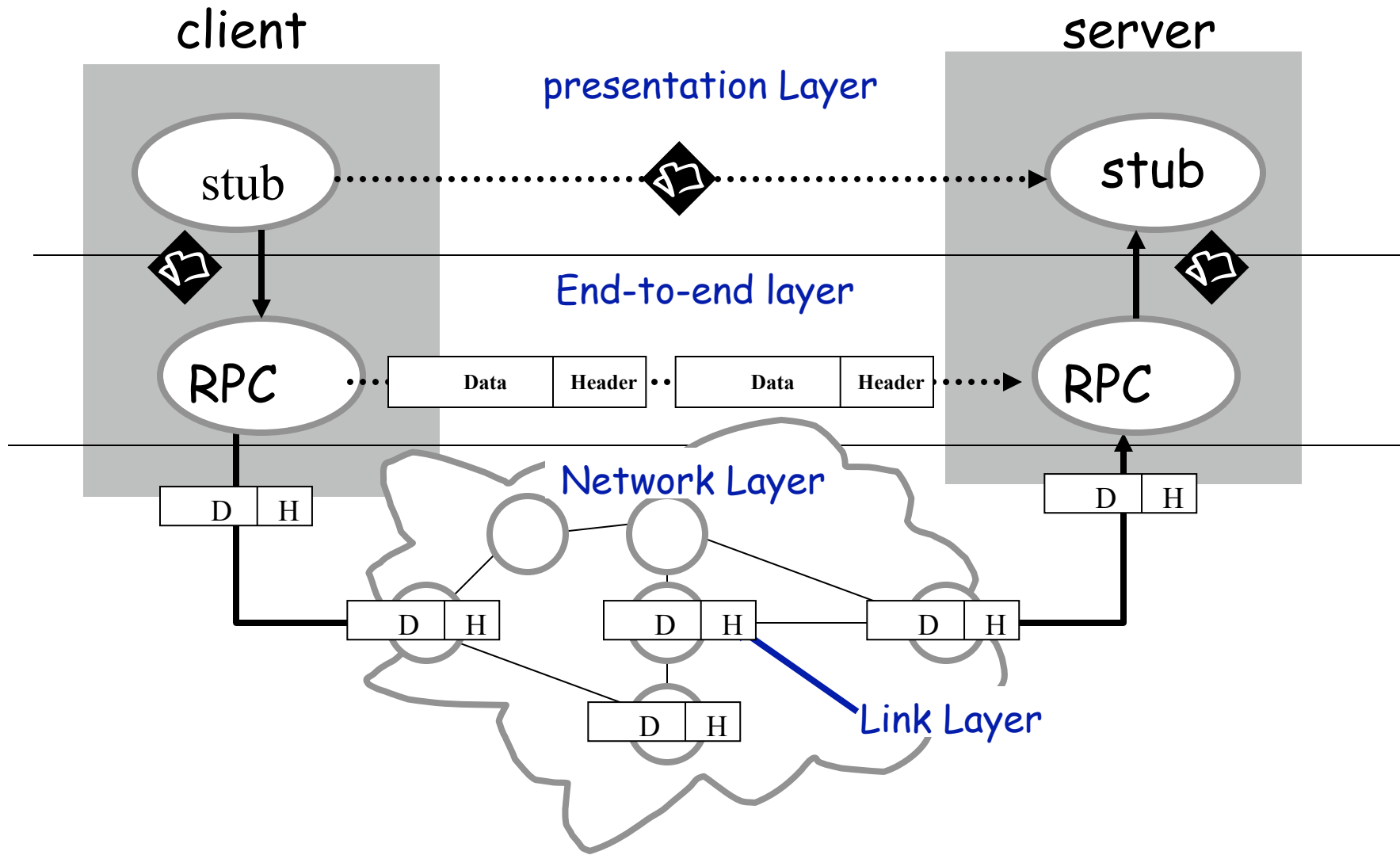
http://web.mit.edu/6.033

Some slides are from lectures by Nick Mckeown, Ion Stoica, Frans Kaashoek, Hari Balakrishnan, Sam Madden, and Robert Morris

MIT MASSACHUSETTS INSTITUTE OF TECHNOLOGY

# End-to-end layer

client                  server

presentation Layer

stub              stub

End-to-end layer

RPC     | Data | Header |    | Data | Header |    RPC

| D | H |

Network Layer

| D | H |

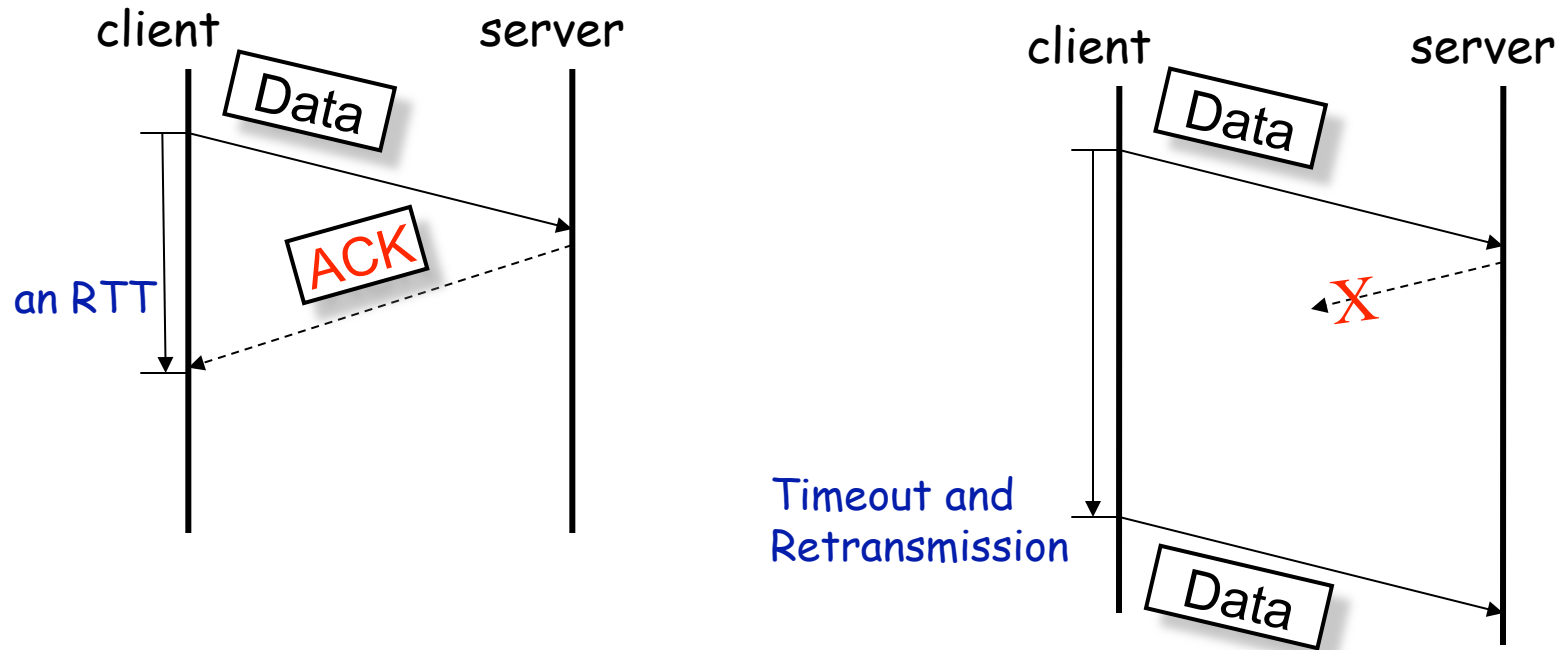| D | H |     | D | H |     | D | H |

| D | H |

Link Layer

# Network layer provides best effort service

- Packets may be:
  - Lossed
  - Delayed (jitter)
  - Duplicated
  - Reordered
  - …

- Problem: Inconvenient service for applications

- Solution: Design protocols for E2E modules
  - Many protocols/modules possible, depending on requirements

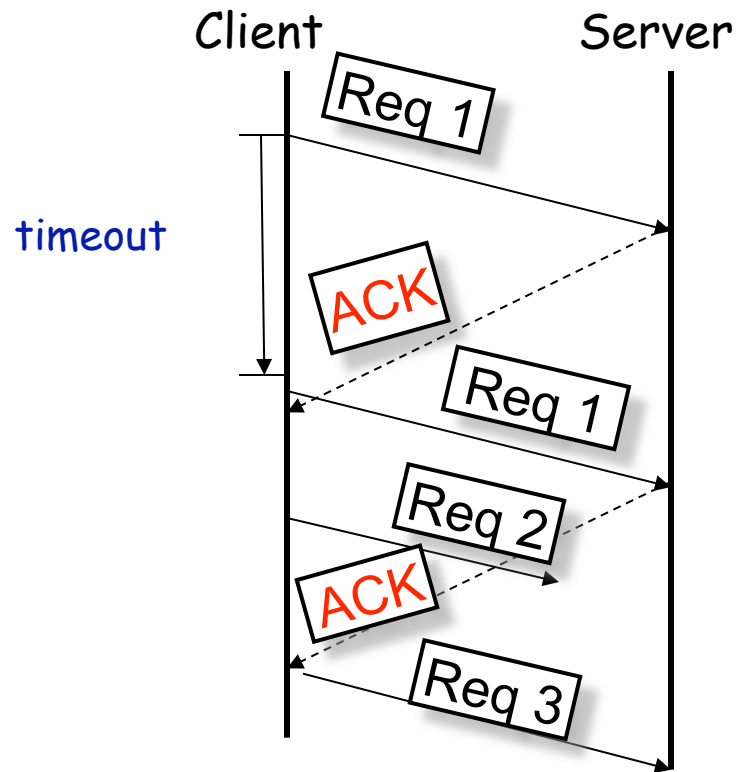# This lecture: some E2E properties

- At most once
- At least once
  - Exactly once?
- Sliding window
- Case study: TCP
- Tomorrow: Network File System (NFS)
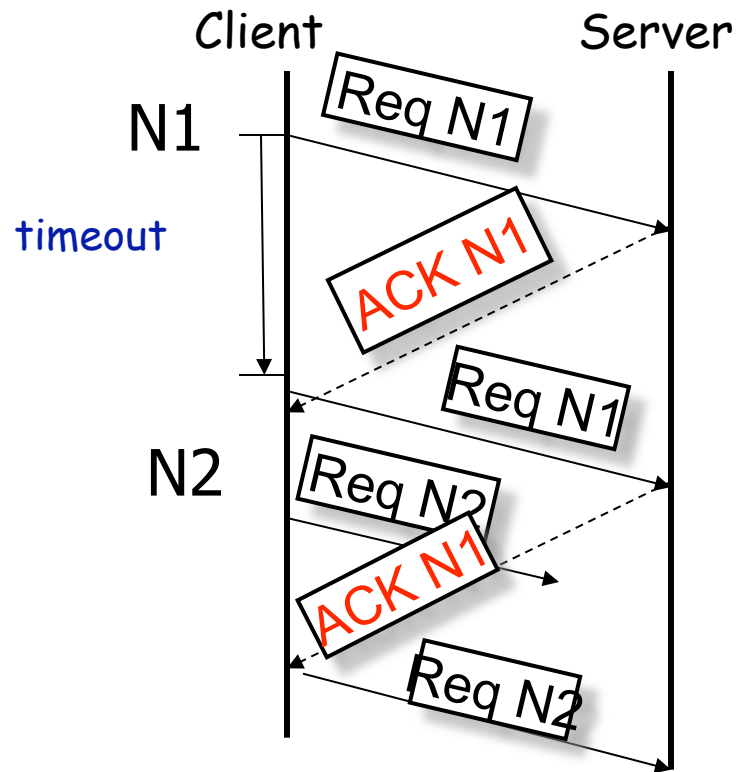
# At Least Once



- Sender persistently sends until it receives an ack
- Challenges:
  - Duplicate ACKs
  - What value for timer

# Duplicate ACK problem



- Problem: Request 2 is not delivered
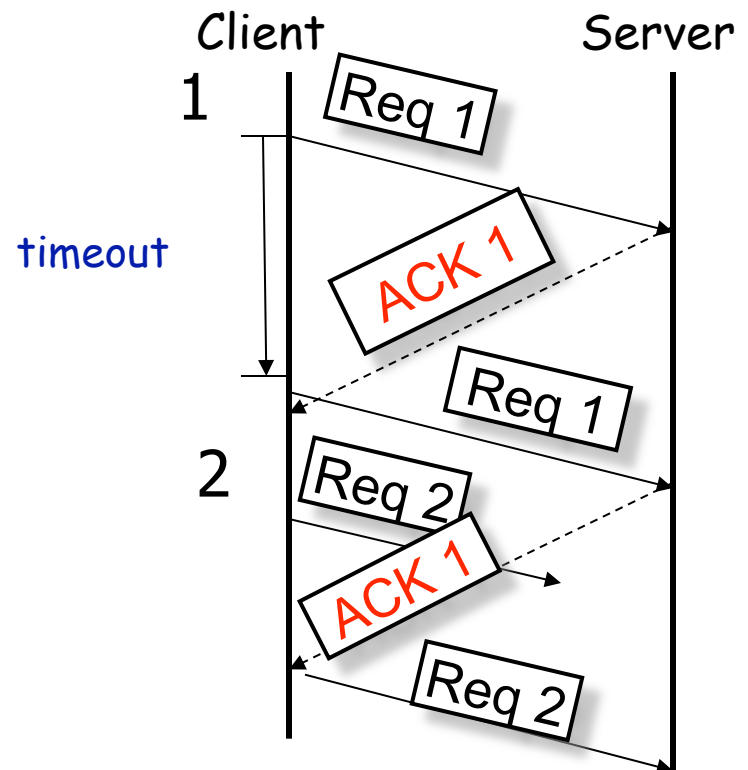  - violates at-least once delivery

# Solution: nonce



- Label request and ack with unique identifier that is never re-used

# Engineering a nonce

- Use sequence numbers
- Challenges:
  - Wrap around?
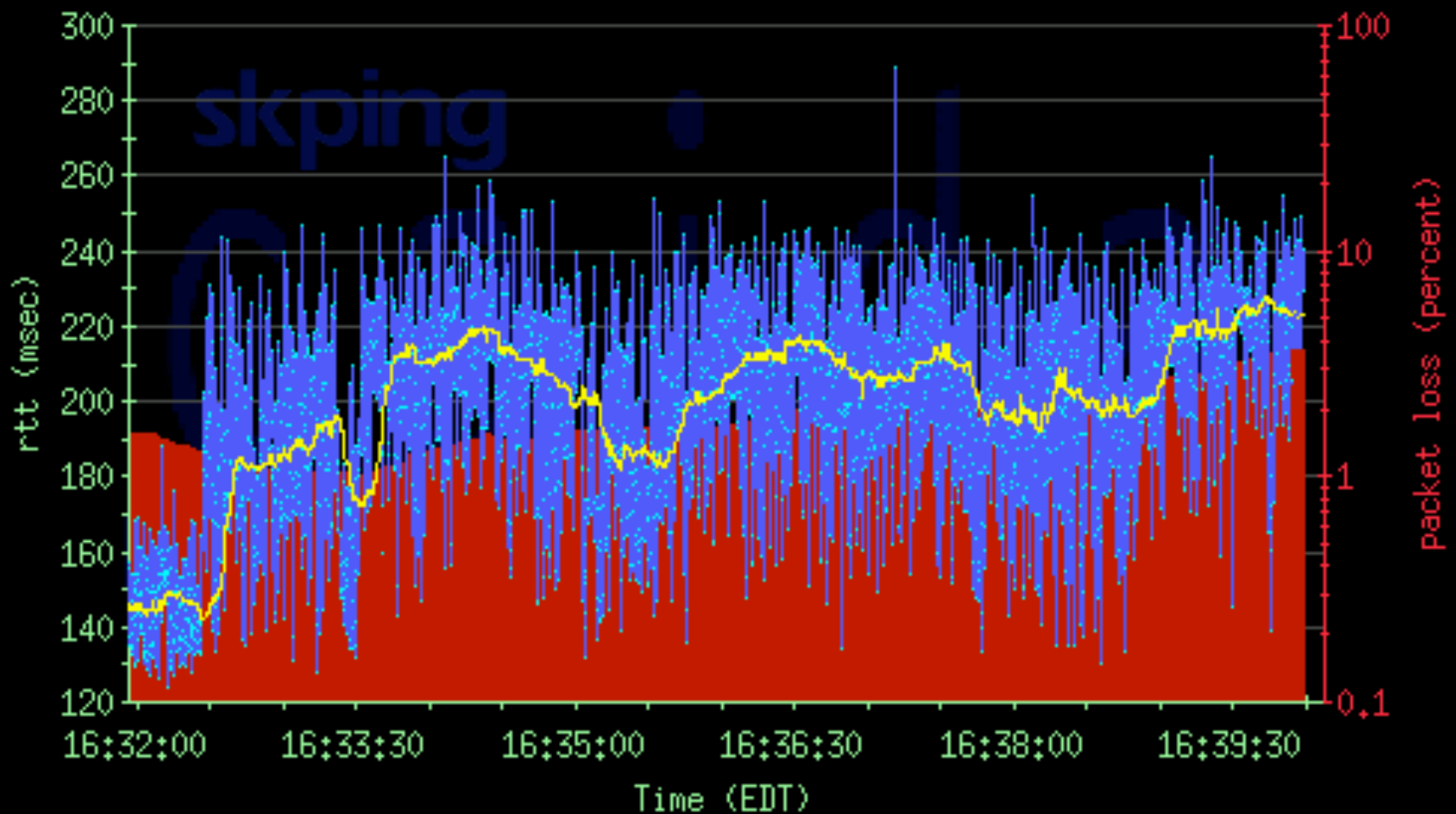  - Failures?

# Timer value

- Fixed is bad. RTT changes depending on congestion
  - Pick a value that's too big, wait too long to retransmit a packet
  - Pick a value too small, generates a duplicate (retransmitted packet).
- Adapt the estimate of RTT → adaptive timeout

# RTT Measurements
**(collected by Caida)**

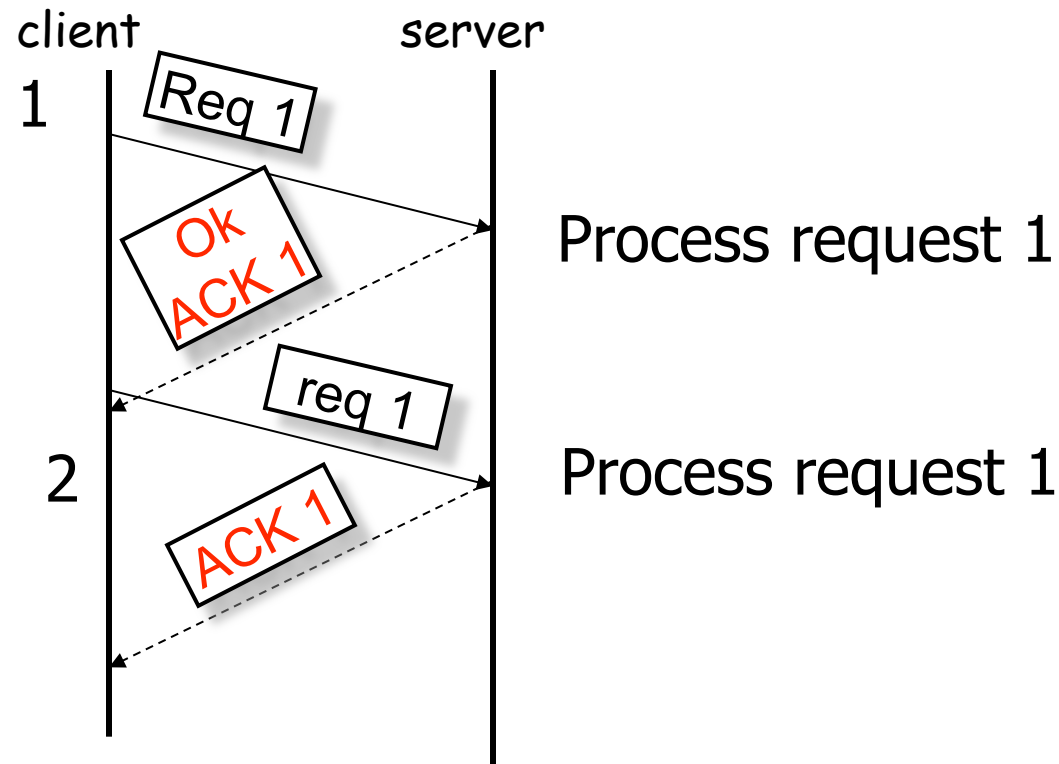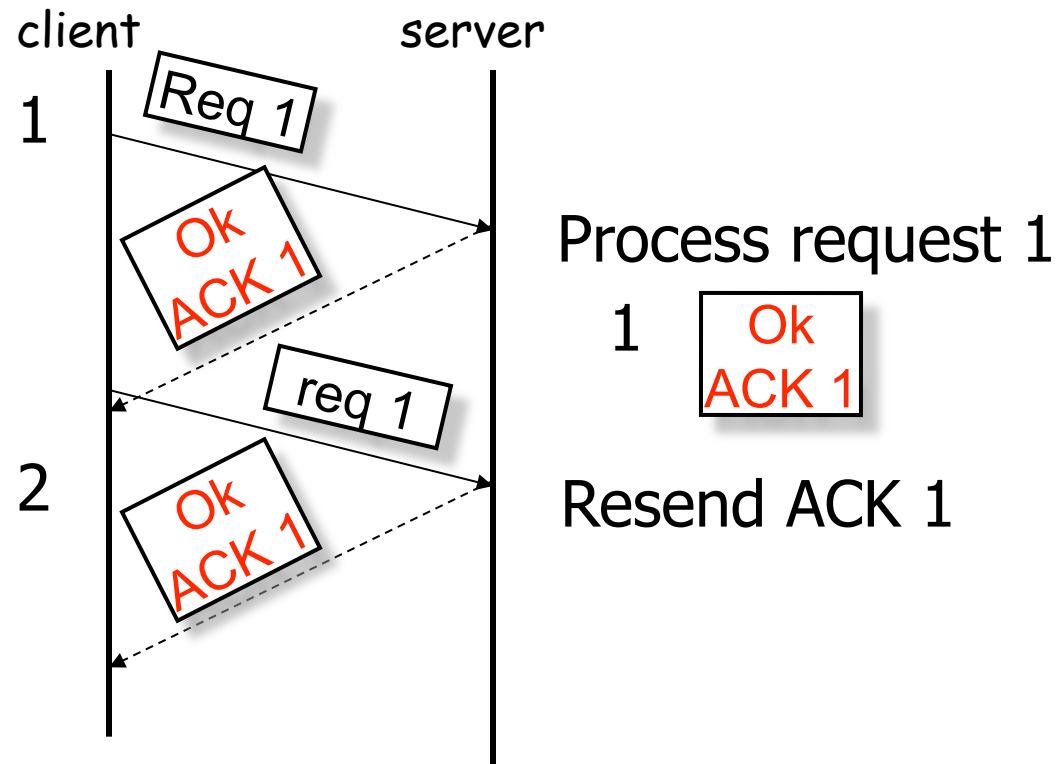# Adaptive Timeout: Exponential weighted moving averages

- Samples $S_1$, $S_2$, $S_3$, ..
- Algorithm
  - EstimatedRTT = $T_0$
  - EstimatedRTT = $\alpha$ S + (1- $\alpha$) EstimatedRTT
  - where $0 \leq \alpha \leq 1$
- What values should one pick for $\alpha$ and $T_0$?
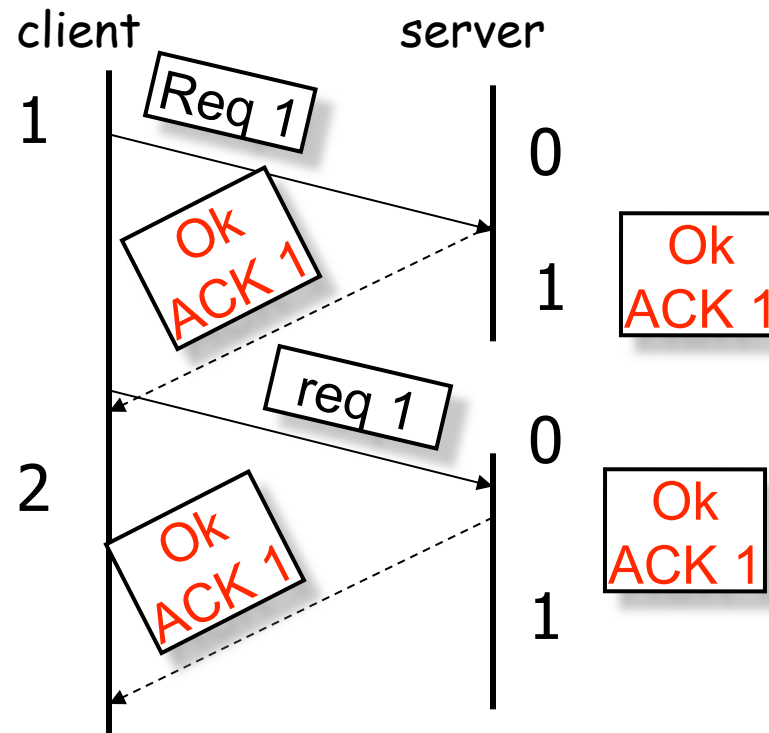  - Adaptive timeout is also hard

# At Most Once Challenges



- Server shouldn't process req 1
- Server should send result preferably

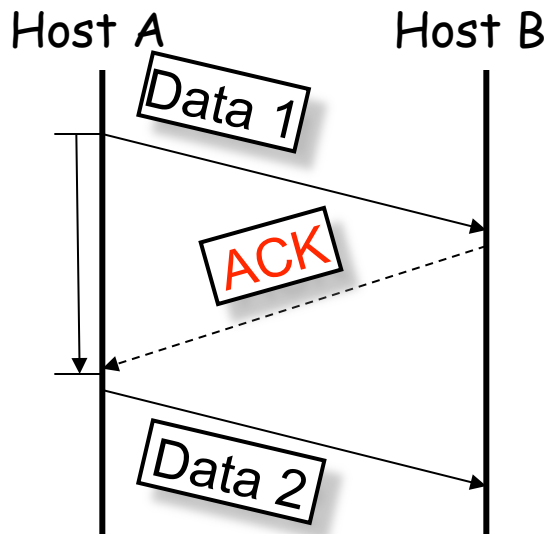# Idea: remember sequence number



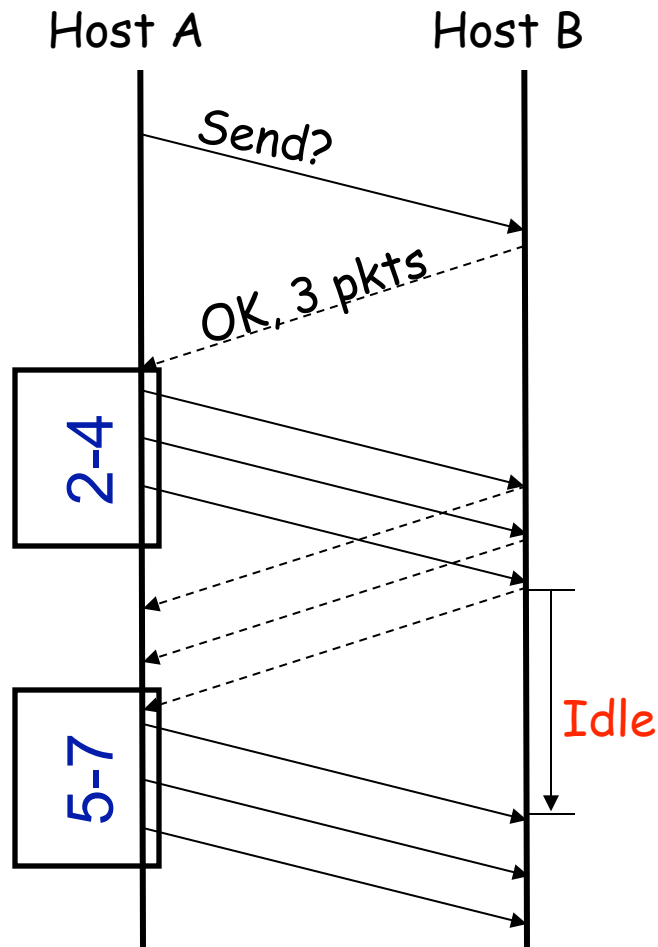- Server remembers also last few responses

# Problem: failures



- Performed request 1 twice!
- How to maintain the last nonce per sender (tombstone)?
  - Write to non-volatile storage?
  - Move the problem?  (e.g., different port number)
  - Make probability of mistake small?
- How about exactly once?  (Need transactions)

# How fast should the sender sends?
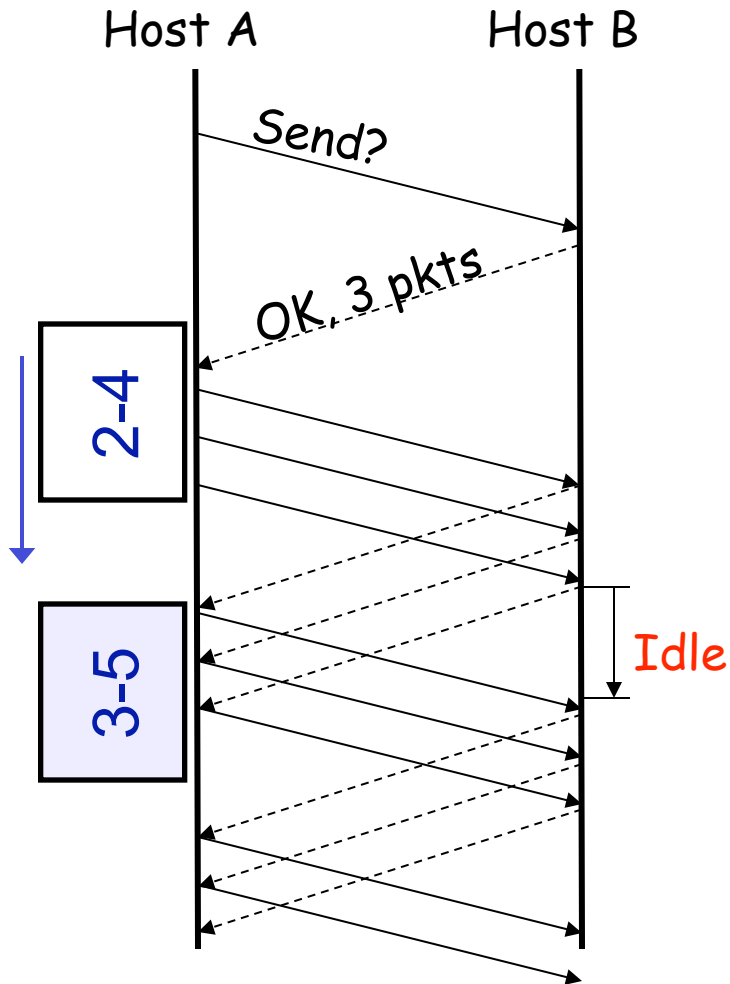


Host A          Host B

Data 1

ACK

Data 2

- Waiting for acks is too slow
- Throughput is one packet/ RTT
  - Say packet is 500 bytes
  - RTT 100ms
  - → Throughput = 40Kb/s, Awful!
- Overlap pkt transmission

# Send a window of packets



- Assume the receiver is the bottleneck
  - Maybe because the receiver is a slow machine

- Receiver needs to tell the sender when and how much it can send

- The window advances once all previous packets are acked → too slow

# Sliding Window

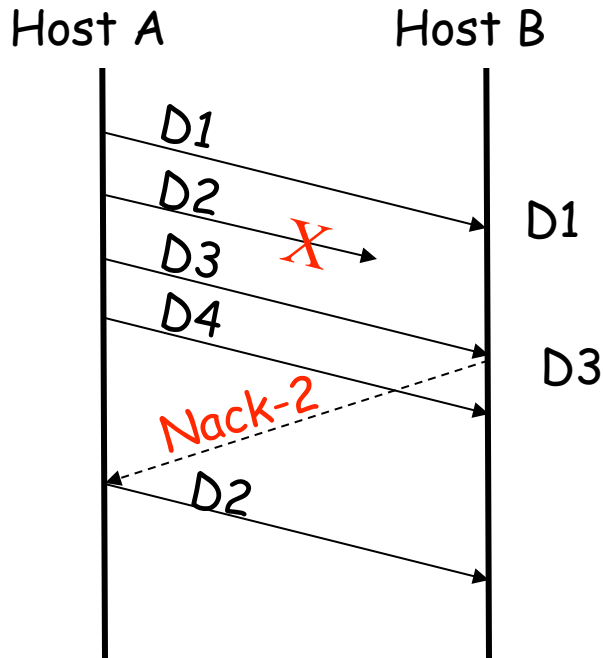Host A          Host B

Send?

OK, 3 pkts

2-4

3-5

Idle

- Senders advances the window whenever it receives an ack → sliding window

- But what is the right value for the window?

# The Right Window Size

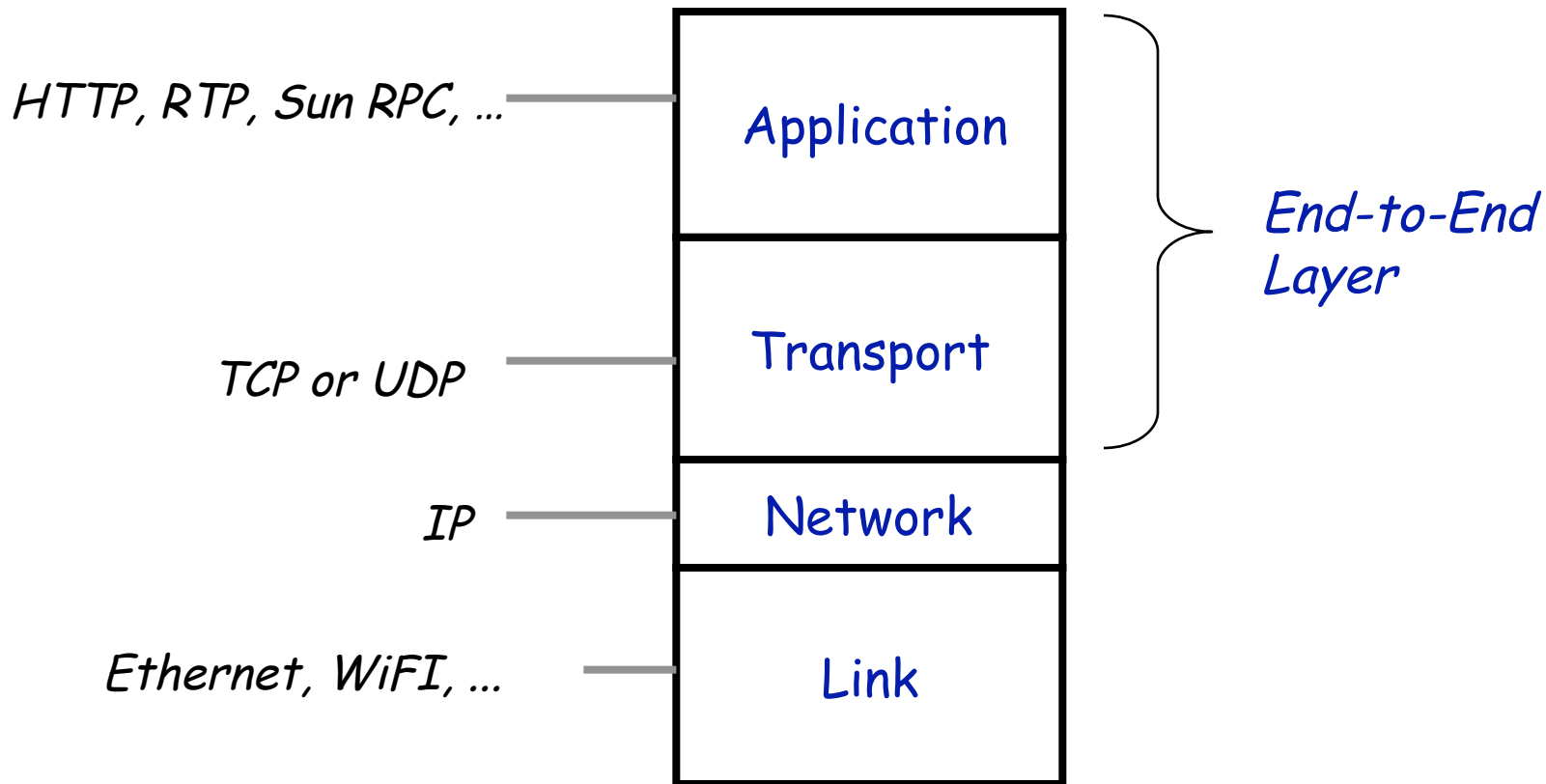- Assume server is bottleneck
  - Goal: make idle time on server zero
  - Assume: server rate is B bytes/s
  - Window size = B x RTT
  - Danger: sequence number wrap around

- What if network is bottleneck?
  - Many senders?
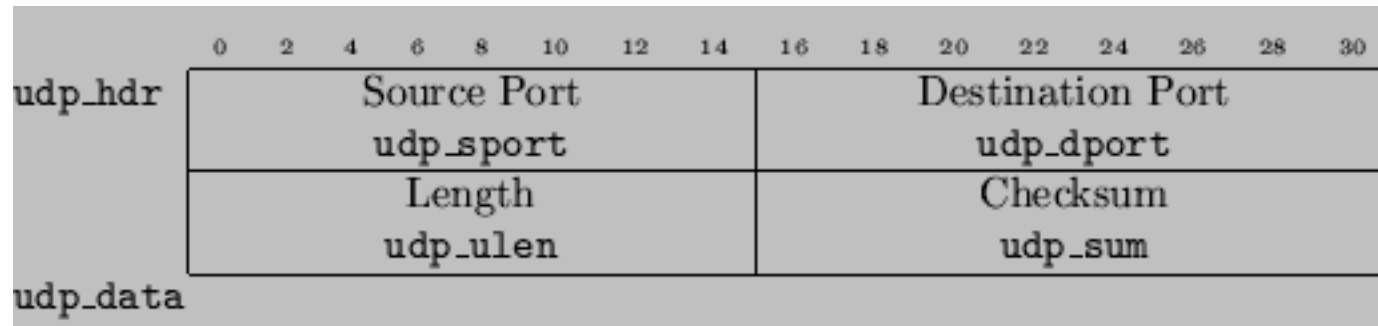  - Sharing?
  - Next lecture

# "Negative" ACK



- Minimize reliance on timer
- Add sequence numbers to packets
- Send a Nack when the receiver finds a hole in the sequence numbers
- Difficulties
  - Reordering
  - Cannot eliminate acks, because we need to ack the last packet

# E2E layer in Internet



HTTP, RTP, Sun RPC, ... — Application

TCP or UDP — Transport

IP — Network

Ethernet, WiFI, ... — Link

End-to-End Layer

The 4-layer Internet model

# UDP

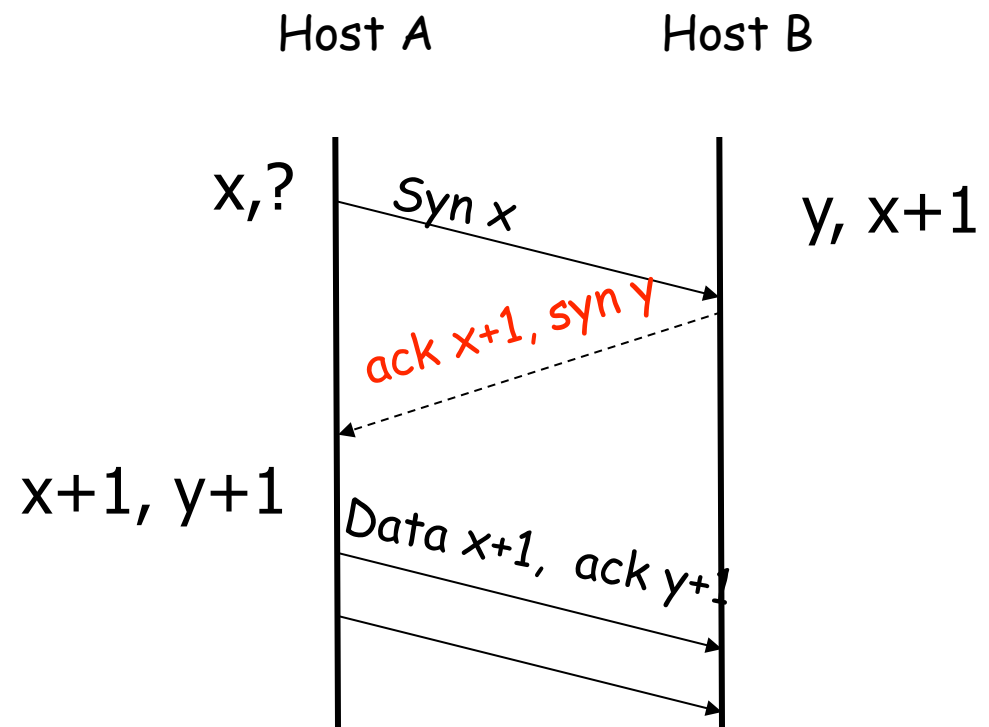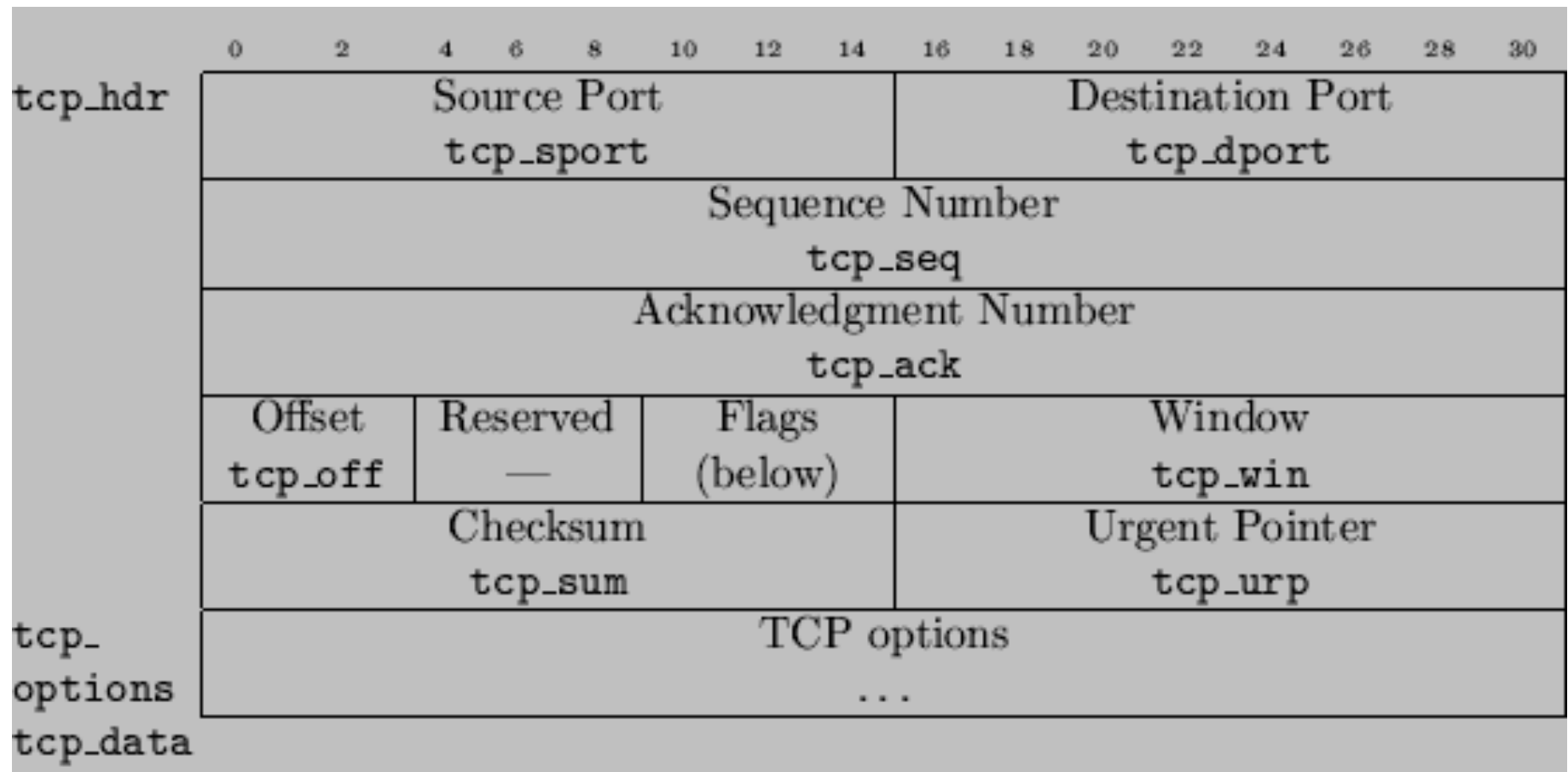| | 0 2 4 6 8 10 12 14 | 16 18 20 22 24 26 28 30 |
|---|---|---|
| udp_hdr | Source Port<br>udp_sport | Destination Port<br>udp_dport |
| | Length<br>udp_ulen | Checksum<br>udp_sum |
| udp_data | | |

# Transmission Control Protocol (TCP)

- Connection-oriented
- Delivers bytes at-most-once
- Bidirectional
  - ACKs are piggybacked

Host A          Host B

x,?     Syn x                    y, x+1

        ack x+1, syn y

x+1, y+1    Data x+1, ack y+1

# TCP header

# Closing a TCP connection



Host A          Host B

x,y      fin x

         ack x+1

         fin y

timed wait    ack y+1

timeout

closed              closed

y, x