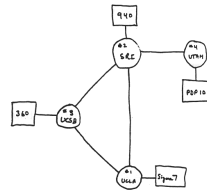


Lecture 3 TCP and Network Architecture



THE ARPA NETWORK
DEC 1969
YALOWS

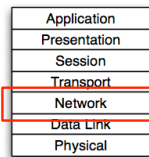
Back to the Network

- The Web uses HTTP, but is built on TCP/IP (like almost all of the Internet)
- Understanding HTTP requires understanding TCP
- No handout on this material yet
- Last networking lecture for awhile

2

Topic 1

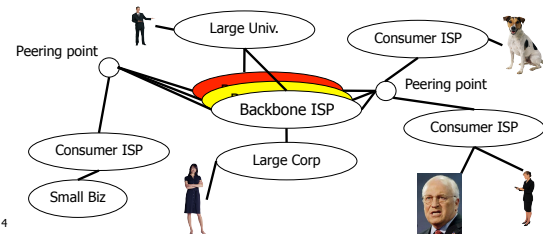
- IP review



3

IP Review

- Best effort packet-switched network
 - Basic unit is *packet*
 - Packets may arrive late, or not at all
 - IP routers form the core of the internet



4

IP Review (2)

- IP addresses are 32 bits
 - E.g., 192.168.1.1
- Forwarding-based networking
 - Each host has a routing table
 - Forward packet to "longest prefix match"

Destination	Gateway
Default	192.168.1.1
192.168.212.111	0:1f:6d:e8:18:0

5

IP Review (3)

- **traceroute** uses debug messages to expose packet's route to destination

```
traceroute to google.com (74.125.95.99), 64 hops max, 52 byte packets
 0 192.168.1.1 (192.168.1.1) 0.037 ms 0.681 ms 1.020 ms (Ann Arbor, MI)
 1 141.212.111.1 (141.212.111.1) 0.566 ms 0.378 ms 0.328 ms
 2 11-cs-bn-arb-2r-bn-arb1-umnet.umich.edu (192.12.80.177) 0.566 ms 0.378 ms 0.328 ms
 3 13-barb-bseb-2r-bn-seb-umnet.umich.edu (192.12.80.11) 0.464 ms 0.495 ms 0.485 ms
 4 13-cs-bn-arb-2r-bn-arb1-umnet.umich.edu (192.12.80.177) 0.566 ms 0.378 ms 0.328 ms
 5 207.72.112.46 (207.72.112.46) 6.552 ms 6.609 ms 6.495 ms (Ann Arbor, MI)
 6 207.72.112.46 (207.72.112.46) 6.552 ms 6.609 ms 6.495 ms (Ann Arbor, MI)
 7 216.239.48.154 (216.239.48.154) 6.785 ms 52.824 ms 6.722 ms (Mountain View, CA)
 8 216.239.48.154 (216.239.48.154) 6.785 ms 52.824 ms 6.722 ms (Mountain View, CA)
 9 209.85.241.29 (209.85.241.29) 17.603 ms (Mountain View, CA)
 10 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 11 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 12 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 13 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 14 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 15 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 16 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 17 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 18 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 19 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 20 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 21 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 22 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 23 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 24 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 25 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 26 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 27 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 28 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 29 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 30 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 31 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 32 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 33 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 34 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 35 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 36 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 37 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 38 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 39 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 40 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 41 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 42 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 43 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 44 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 45 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 46 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 47 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 48 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 49 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 50 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 51 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 52 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 53 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 54 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 55 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 56 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 57 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 58 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 59 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 60 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 61 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 62 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 63 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
 64 209.85.241.29 (209.85.241.29) 18.306 ms 25.116 ms (Mountain View, CA)
```

6

IP Review (4)

- No reason the links have to make sense
 - `traceroute urbanspoon.com`
 - Ann Arbor, MI
 - San Mateo, CA
 - New York, NY
 - <ends>

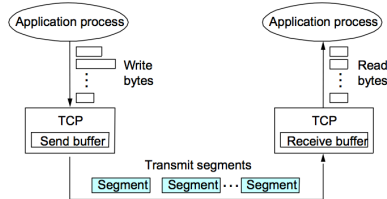
7

TCP

- Transport Control Protocol
 - Reliable, ordered bytestreams
 - Connection-oriented, unlike IP
 - "Virtual circuit" networking built on packet infrastructure
- Processes tied to "host:port" pairs
 - Ports are an OS-level concept
 - Server ports are "well-known" and associated with services
 - Other ports are "ephemeral"

8

TCP Delivery



9

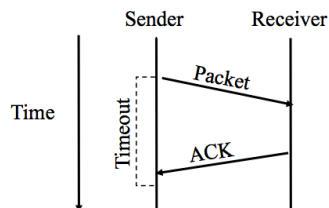
Reliable connections

- Basic principle of reliable TCP connection is retransmission
 - Each packet has 32-bit Sequence Number
 - Every SeqNo is ACKed by receiver
 - When timeout expires, sender emits again
 - Simple!
- OK, maybe not quite that easy

10

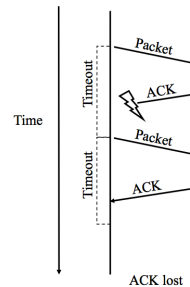
Stop and Wait

- Send a packet, wait for ACK
- Receiver ACKs everything

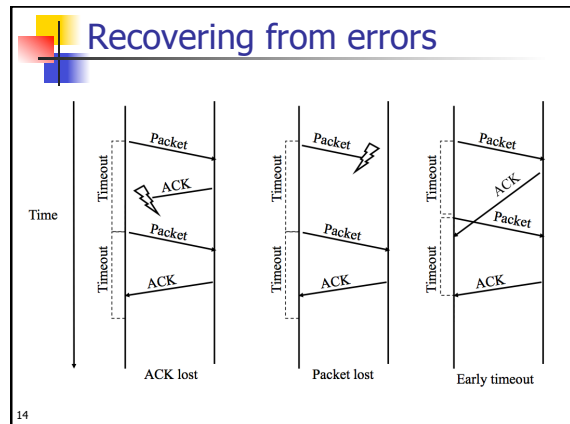
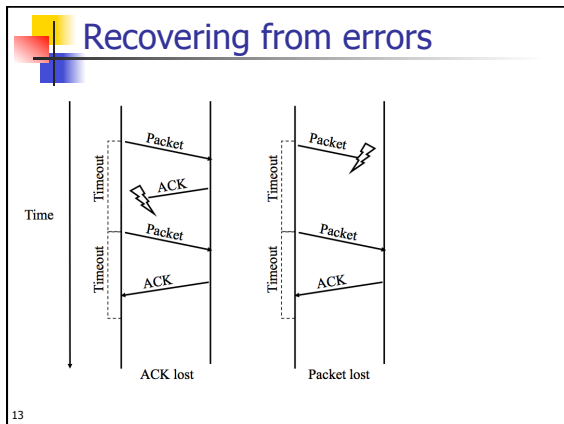


11

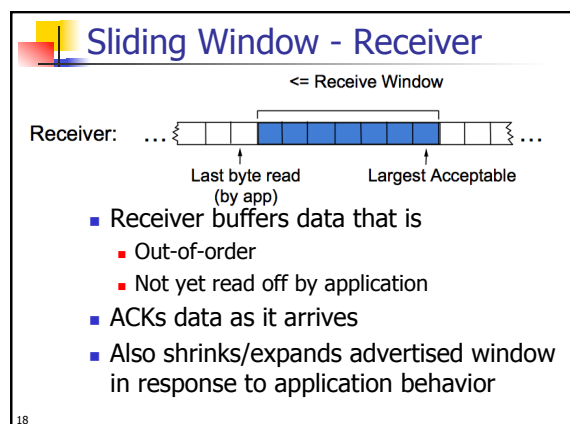
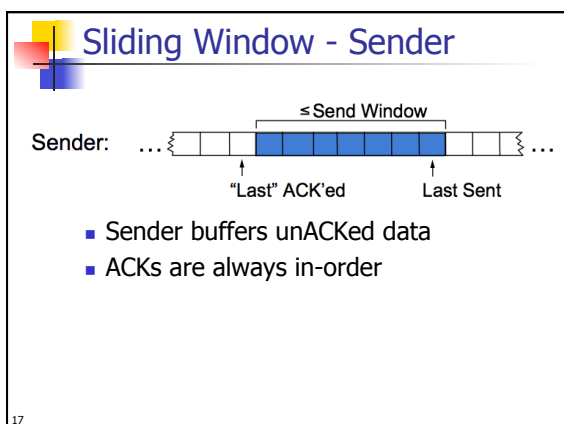
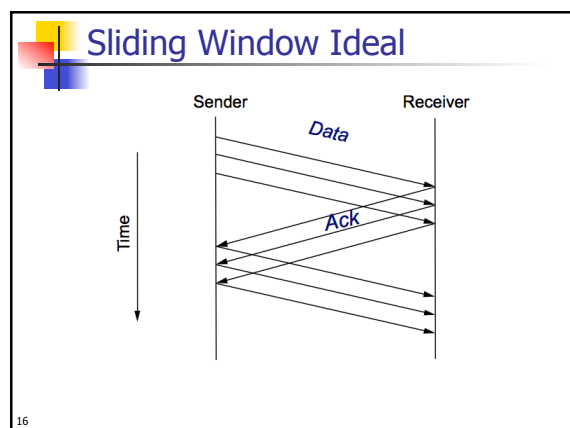
Recovering from errors

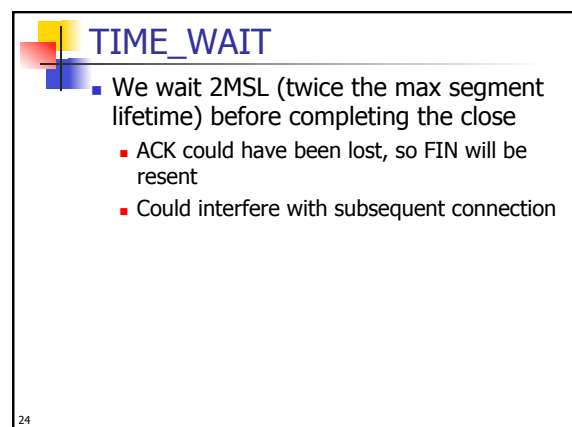
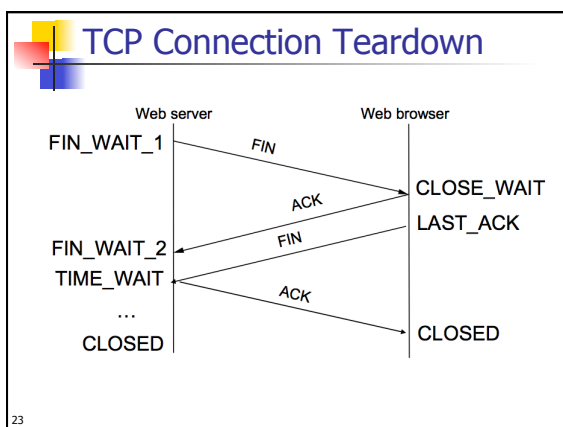
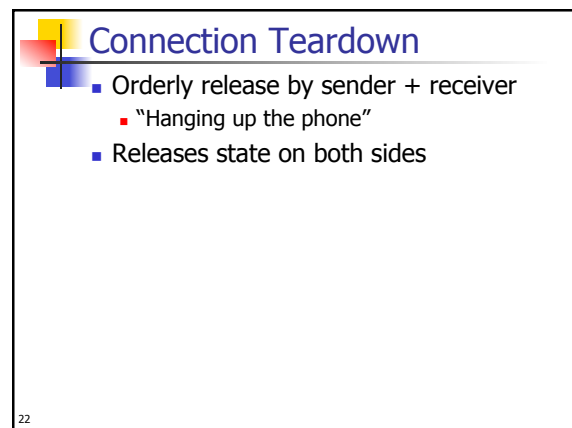
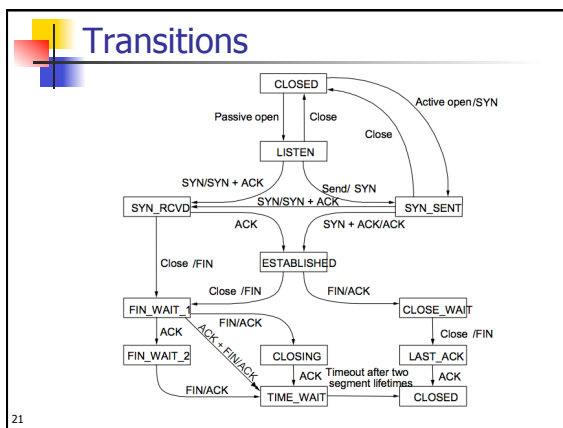
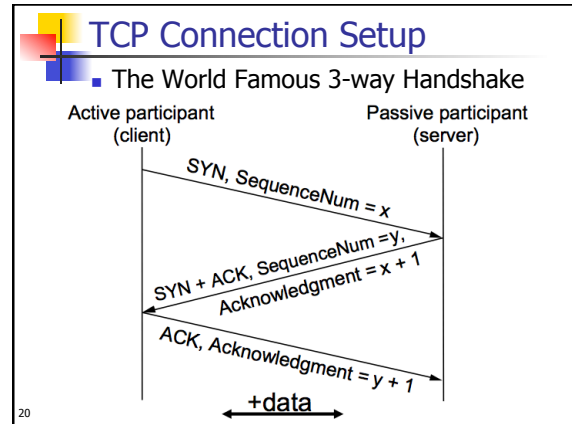
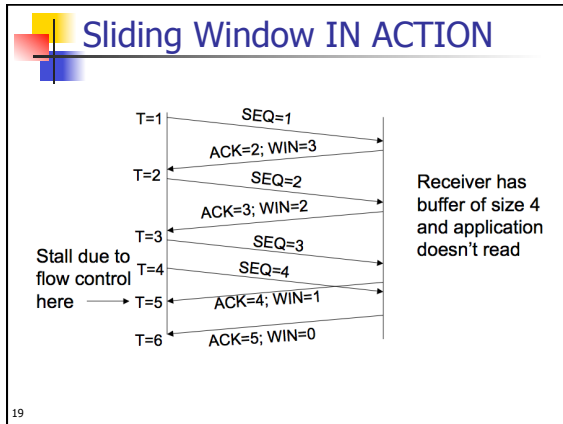


12



- ### Flow Control
- Any downside to Stop-and-Wait?
 - Sliding window** technique for managing send/receive capacity
 - Receiver indicates "receive window" it is willing to buffer
 - Sender cannot have more packets in transit (unACKed) than what can fit in the window
 - Ideally, window is $\text{bandwidth} * \text{RT delay}$
 - Keeps as much data in transit as possible
- 15





Congestion Control

- Buffers in middle of network can become overloaded
- Not part of window negotiation

Packets queued here

25

Congestion Control

- When has packet been lost?
 - Too long a timer, you're waiting pointlessly
 - Too short, you're adding needless load
- Many retransmits can induce *congestion collapse* (it happened in late 1980s!)
 - Retransmits just add to congestion
 - Capacity of network falls dramatically
- **Congestion control** is the art of dynamically sizing the sender's window to avoid slamming in-network routers
 - Ideally based on round-trip time

26

Probing

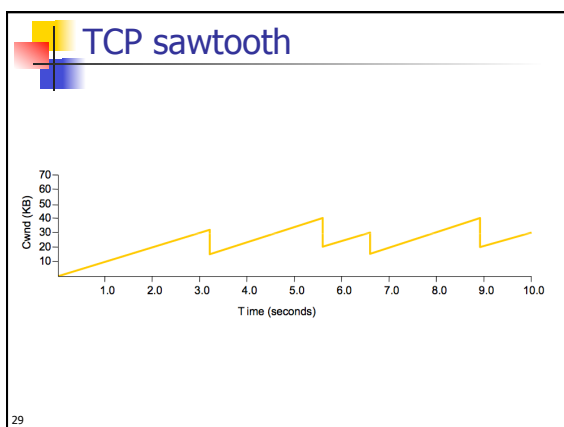
- Hosts use *packet-probing* to figure out bandwidth to remote side
- A lost packet is assumed to be dropped due to router buffer congestion
 - Good assumption on internet
 - What about on wireless?

27

One approach: AIMD

- Additive Increase, Multiplicative Decrease
- Increase slowly if there is bandwidth
- Back down fast in event of loss

28



Another: Slow Start

- AIMD can take ages to find best level
- Instead, SS *doubles* every RTT
- Use Slow Start at first; then AIMD at SS threshold

30

HTTP on TCP

- TCP Summary
 - A roundtrip for each window of bytes
 - Takes some time to find best rate
 - OS buffer overhead for each connection
- HTTP Summary
 - Connections short-lived
 - Small amounts of data per connection
 - Many connections per HTML page
- Naïve HTTP on TCP yields *awful* performance

31

HTTP 1.0 on TCP

- Every HTML-embedded item requires a GET

32

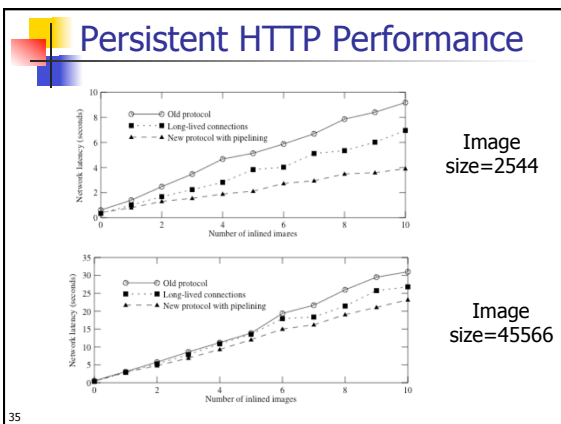
HTTP on TCP

33

HTTP/1.1

- Persistent Connections
 - One TCP connection, many HTTP methods
 - Good for small objects, not large
 - Makes parallel downloads difficult

34



HTTP/1.1

- Persistent Connections
 - One TCP connection, many HTTP methods
 - Good for small objects, not large
 - Makes parallel downloads difficult
- Caching
 - GET + IF-MODIFIED-SINCE <timestamp>
 - When combined with browser cache, eliminates a lot of unneeded data transfer
 - Same number of roundtrips
 - Lots of different caches possible
 - Browser, department proxy, Akamai

36

