## eecs 489 COMPUTER NETWORKS

Lecture 29:
TCP Connection Establishment

## Internet Protocol Stack

**application:** supporting network applications
• HTTP, SMTP, FTP, etc.

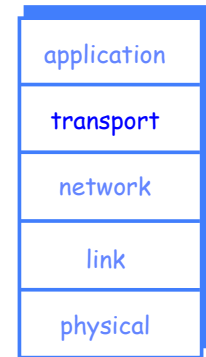**transport:** endhost-endhost data transfer
• TCP, UDP

**network:** routing of datagrams from source to destination
• IP, routing protocols

**link:** data transfer between neighboring network elements
• Ethernet, WiFi

**physical:** bits "on the wire"

| application |
| --- |
| transport |
| network |
| link |
| physical |

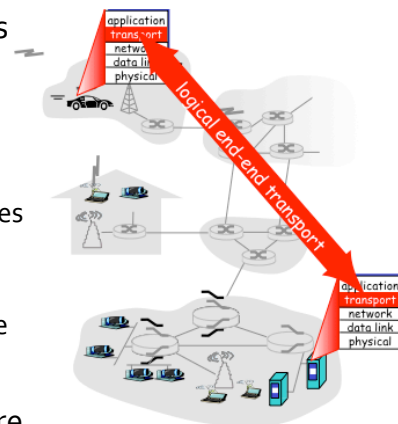## Transport Protocols

Provide logical communication between application processes running on different hosts

Run on end hosts

• sender: breaks each application message into segments, and passes them onto the network layer

• receiver: reassembles segments into messages, passes them to the application layer

Multiple transport protocols are available to applications

## Internet Transport Protocols
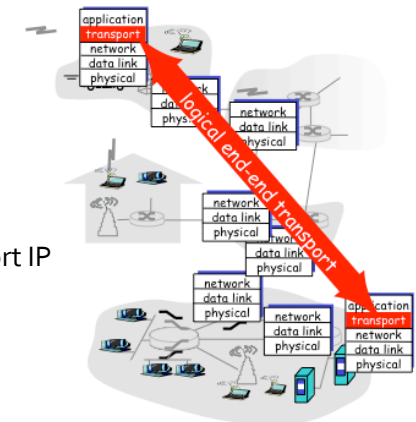
Reliable, in-order delivery (TCP)
• connection setup
• flow control
• congestion control

Unreliable, unordered delivery (UDP)
• no-frills extension of best-effort IP

Services not available
• delay guarantees
• bandwidth guarantees

# Why Would Anyone Use UDP?

Lightweight communication between processes:

- faster than TCP, no connection establishment/tear-down stages (1 vs. 2.5 rtts)
  - simply send messages to and receive them from a socket

- no connection state at server and client
  - avoid overhead and delays of ordered, reliable delivery
  - no allocation of buffers, parameters, sequence #s, etc.
  - making it easier to handle many active clients at once

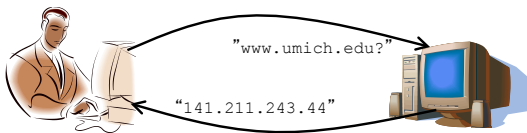# Why Would Anyone Use UDP?

Lightweight communication between processes:

- small segment header
  - UDP header is only eight-bytes long

- fine control over what data is sent and when
  - as soon as an application process writes into the socket
  - UDP packages the data and sends the packet

- no congestion control
  - UDP can blast away as fast as the network can handle

- broadcast & multicast can only use UDP (Why?)

# Popular Applications That Use UDP

Simple query protocols like DNS
- overhead of connection establishment is overkill
- easier to have the application retransmit if necessary



"www.umich.edu?"

"141.211.243.44"

Multimedia streaming
- loss tolerant, rate sensitive
  - by the time a lost packet is retransmitted, it's too late
  - congestion control introduces too much jitter
  - e.g., calls, video streaming, gaming

Reliable UDP: add reliability at application layer
- application-specific error recovery!

# TCP: Transmission Control Protocol

Connection oriented
- explicit set-up and tear-down of TCP session

Stream-of-bytes service
- sends and receives a stream of bytes, not discrete messages

Flow control
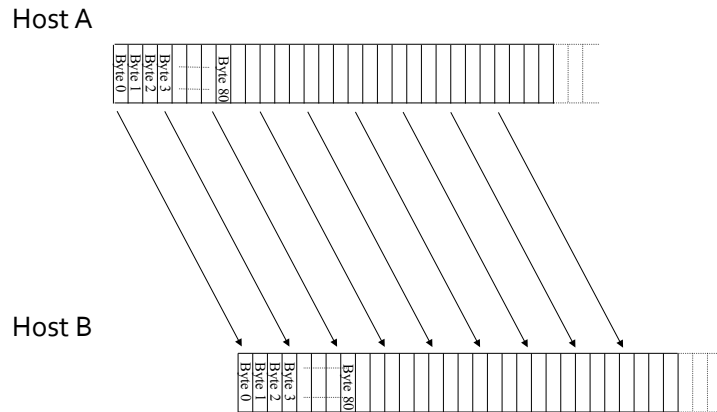- prevent overflow of the receiver's buffer space

Reliable delivery
- retransmission of lost packets

Congestion control
- adapt to network congestion for the greater good

Does not provide timing guarantee nor minimum bandwidth guarantee

# TCP "Stream of Bytes" Service

Host A

Host B

# Emulated Using TCP "Segments"

Host A

TCP Data

Segment sent when:
1. segment full (Max Segment Size),
2. not full, but times out waiting for more data from app, or
3. "pushed" by application

TCP Data

Host B

# TCP: Transmission Control Protocol

Provides reliability on datagram network
What does reliable delivery entail?

- 
- 
- 
- 

Link layer already provides reliable delivery
Why do we need provide it again at the transport layer?

- 
- 

# Need for E2E Reliability

pkt 3 corrupted
or dropped

# E2E Reliability

Lack of reliable delivery due to:
- re-routed packets
- bit error
- dropped/lost packets (due to congestion)
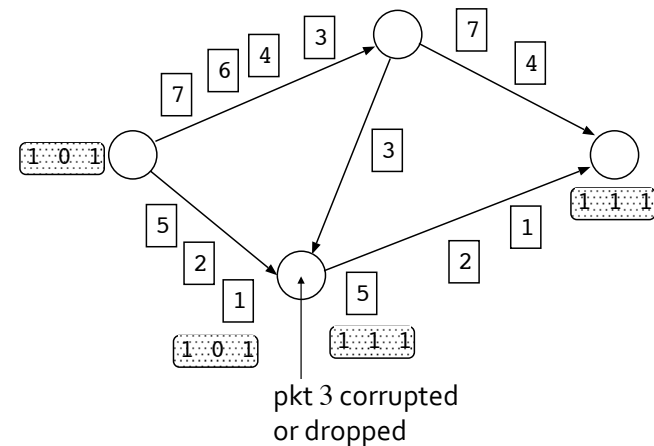- system reboots

What are some of the tools available to us to achieve reliability at the transport layer, given unreliable network layer?
- 
- 
- 
- 

# Sequence Number

What are the uses of sequence number in providing reliability?
- 
- 
- 

# TCP's Cumulative ACK

ACKs the last byte received in-order

Tells sender the next-expected seq#

If bytes $0$ to $n$ have been received, ACK says $n+1$

subsequent out-of-order packets generate the same cumulative ACK:

```
                    Sender          Destination
                      S                  D
    8 bytes: seq# 0
   10 bytes: seq# 8
   10 bytes: seq# 18                    ack 8
   10 bytes: seq# 28        x           ack 18

                                        ack 18
```

Advantage: lost ACK can be "covered" by later ACKs

Disadvantage: size of gap between two packets not known to sender

# TCP Cumulative ACK



Host A

ISN (initial sequence number)

Sequence number = 1st byte

TCP Data | TCP HDR

ACK sequence number = next expected byte

TCP Data | TCP HDR

Host B

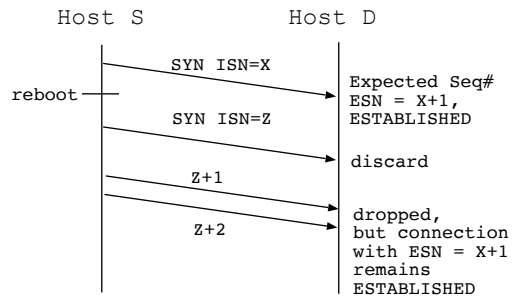# Connection Establishment

TCP SYNchronization packet to establish a connection carries the Initial Sequence Number (ISN)

First try:

```
        Host S              Host D

              SYN ISN=X
reboot ───────────────────▶  Expected Seq#
                             ESN = X+1,
              SYN ISN=Z      ESTABLISHED
       ───────────────────▶
                             discard
                Z+1
       ───────────────────▶
                             dropped,
                Z+2          but connection
       ───────────────────▶  with ESN = X+1
                             remains
                             ESTABLISHED
```
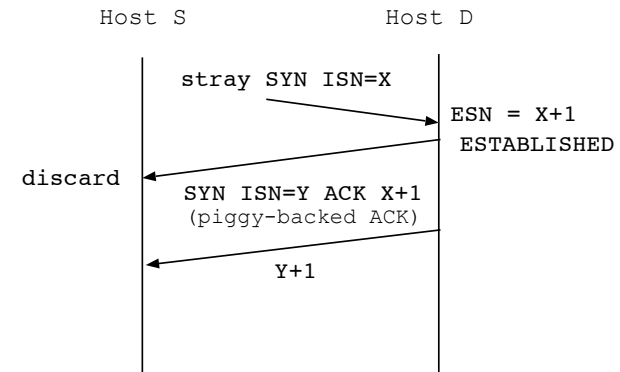
Lesson: connection request must be ACKed

# Connection Establishment

Second try:

```
        Host S              Host D

              stray SYN ISN=X
       ───────────────────▶  ESN = X+1
                             ESTABLISHED

discard ◀──── SYN ISN=Y ACK X+1
              (piggy-backed ACK)

              Y+1
       ◀───────────────────
```

Lesson: connection ACK must be ACKed or rejected

# TCP Connection Establishment

Three-way handshake:

```
     Host S                  Host D

           SYN ISN=X
     ───────────────────▶  ESN = X+1
                             ESTABLISHED
     SYN ISN=Y
     ACK X+1
     ◀───────────────────
ESN = Y+1
ESTABLISHED
           ACK Y+1
     ───────────────────▶

             X+1
     ───────────────────▶
```

SYN uses a seq#

# Three-Way Handshake

How three-way handshake solves the original problems:

```
   Host S      Host D          Host S        Host D

        SYN ISN=X                  stray SYN ISN=X
reboot ─────────────▶  SYN ISN=Y  ────────────────▶ ESN = X+1
SYN ISN=Z             ACK X+1                        ESTABLISHED
                                  discard ◀─ SYN ISN=Y ACK X+1
RST Y ◀──╳──▶ RST Z
                                              RST Y
```

What if the SYN packet is lost?
• since there's no good way to gauge RTT
• some TCP sender times out after 3-6 seconds

# TCP Segment



## IP packet
• no bigger than Maximum Transmission Unit (MTU)
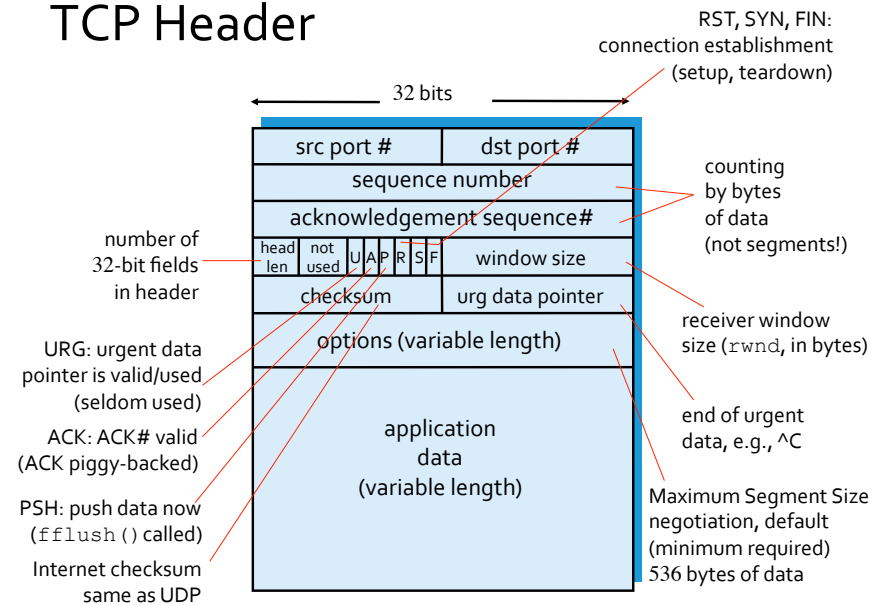• e.g., up to 1500 bytes on an Ethernet

## TCP packet
• IP packet with a TCP header and data inside
• TCP header is typically 20 bytes long

## TCP segment
• no more than Maximum Segment Size (MSS) bytes
• e.g., up to 1460 consecutive bytes from the stream
• (note: PA3's MSS includes IP header)

[Rexford]

# TCP Header

RST, SYN, FIN:
connection establishment
(setup, teardown)



32 bits

src port #    dst port #
sequence number
acknowledgement sequence#
head len | not used | U A P R S F | window size
checksum | urg data pointer
options (variable length)
application data (variable length)

number of 32-bit fields in header

URG: urgent data pointer is valid/used (seldom used)

ACK: ACK# valid (ACK piggy-backed)

PSH: push data now (`fflush()` called)

Internet checksum same as UDP

counting by bytes of data (not segments!)

receiver window size (`rwnd`, in bytes)

end of urgent data, e.g., ^C

Maximum Segment Size negotiation, default (minimum required) 536 bytes of data

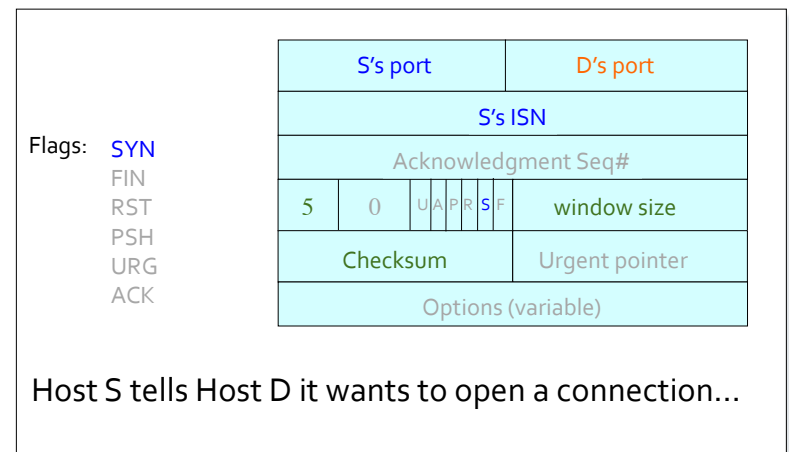# TCP Header Fields

**Sequence number:**
• sequence numbers count bytes sent
• seq# of a packet is the seq# of the first byte it carries

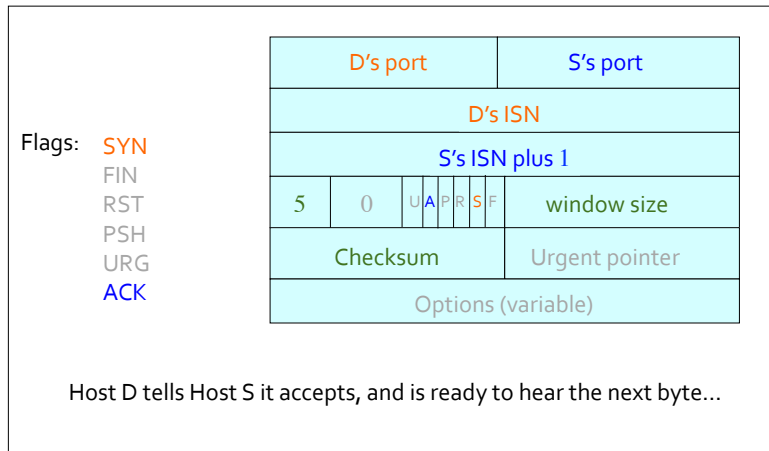**Acknowledgement sequence number:**

allows for piggy-backed ACK
• TCP traffic is often bidirectional
  • both data and ACKs travel in both directions
• ACK packets have high overhead
  • 40 bytes for the TCP/IP headers, carrying no data
• piggy-backing allows a host D to send its ACK to host S along with its data for host S
• delayed ACK: TCP allows the receiver to delay sending of ACKs to increase chances of piggy-backing and to reduce number of ACKs (since they're cumulative)
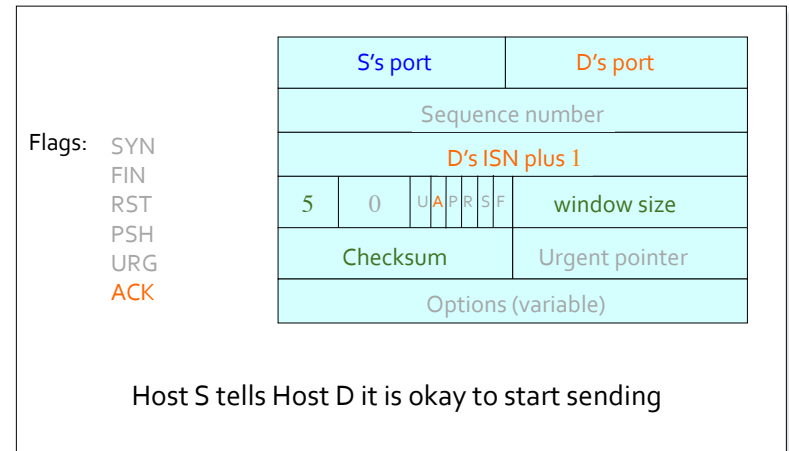
# Step 1: S's Initial SYN Packet



Flags: **SYN**
FIN
RST
PSH
URG
ACK

| S's port | D's port |
| S's ISN | |
| Acknowledgment Seq# | |
| 5 | 0 | U A P R S F | window size |
| Checksum | Urgent pointer |
| Options (variable) | |

Host S tells Host D it wants to open a connection...

[Rexford]

# Step 2: D's SYNACK Packet

Flags:
SYN
FIN
RST
PSH
URG
ACK

| D's port | S's port |
|---|---|
| D's ISN | |
| S's ISN plus 1 | |

| 5 | 0 | U A P R S F | window size |
|---|---|---|---|

| Checksum | Urgent pointer |
|---|---|
| Options (variable) | |

Host D tells Host S it accepts, and is ready to hear the next byte...

Upon receiving Host D's SYNACK packet, Host S can start sending data

[Rexford]

# Step 3: S's ACK of the SYNACK

Flags:
SYN
FIN
RST
PSH
URG
ACK

| S's port | D's port |
|---|---|
| Sequence number | |
| D's ISN plus 1 | |

| 5 | 0 | U A P R S F | window size |
|---|---|---|---|

| Checksum | Urgent pointer |
|---|---|
| Options (variable) | |

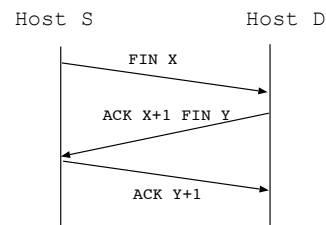Host S tells Host D it is okay to start sending

Upon receiving Host S's packet, Host D can start sending data

[Rexford]

# Connection Tear-down

When to release a connection?
How do you know the other side is done sending and all sent packets have arrived?
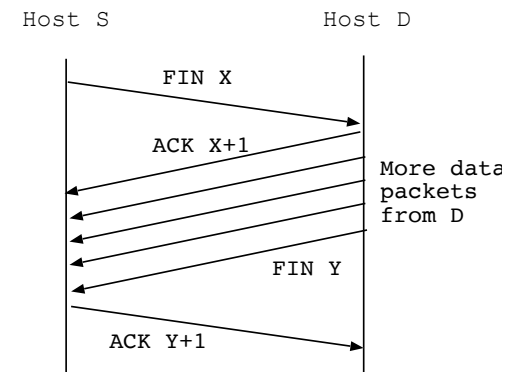
Use 3-way handshake to tear-down connection:

```
Host S              Host D

          FIN X
        --------->

      ACK X+1 FIN Y
        <---------

          ACK Y+1
        --------->
```

FIN also uses a seq#

# Connection Tear-down

If the other side still has data to send:

```
Host S              Host D

          FIN X
        --------->

        ACK X+1
        <---------
                    More data
                    packets
                    from D
          FIN Y
        <---------

          ACK Y+1
        --------->
```
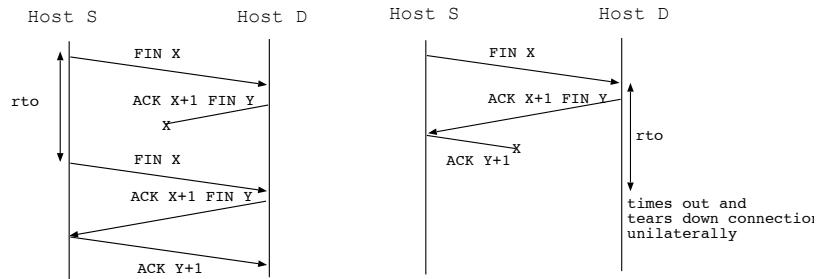
Why not delay `ACK X+1` until `FIN Y`?

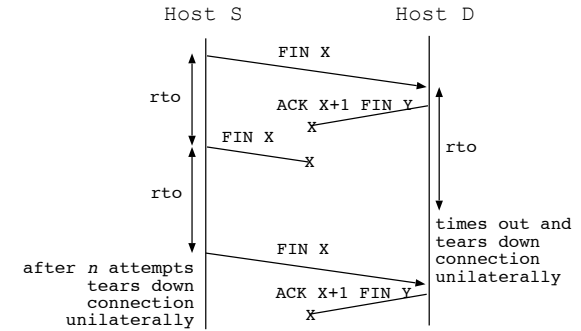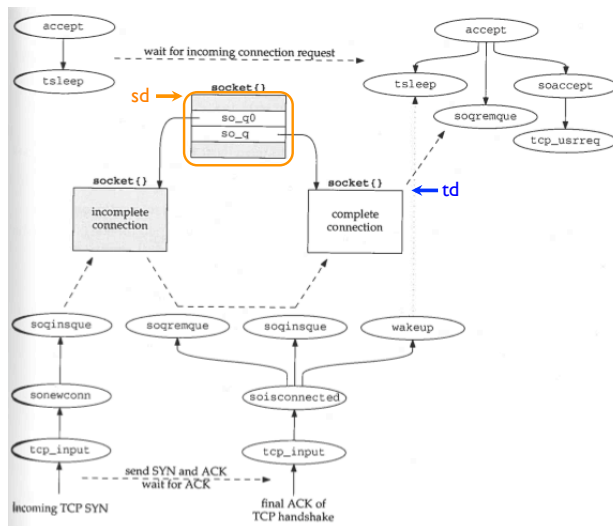# Connection Tear-down

Still depends on timeout for correctness:

TCP connection tear-down depends on timers for correctness, but uses 3-way handshake for performance improvement
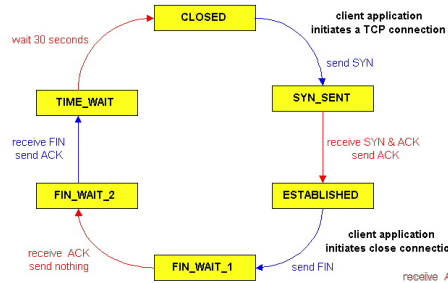
# Socket Connection Queues



Stevens TCP/IP Illustrated v. 2 pp. 441, 461
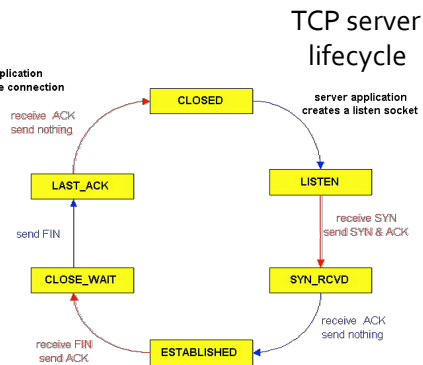
# TCP Connection Establishment Demo

```
% ifconfig –a
% sudo tcpdump –i en0 –S host web.eecs.umich.edu
```

# TCP Connection Management FSM

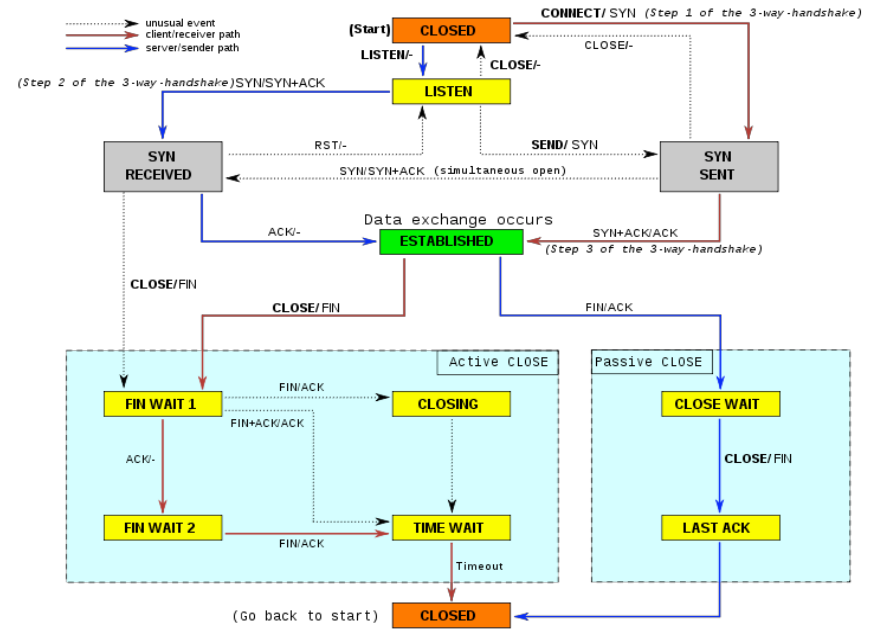Put together . . .

**TCP client lifecycle**

CLOSED
wait 30 seconds
client application initiates a TCP connection
send SYN
TIME_WAIT
SYN_SENT
receive FIN send ACK
receive SYN & ACK send ACK
FIN_WAIT_2
ESTABLISHED
receive ACK send nothing
client application initiates close connection
FIN_WAIT_1
send FIN

**TCP server lifecycle**

CLOSED
receive ACK send nothing
server application creates a listen socket
LAST_ACK
LISTEN
send FIN
receive SYN send SYN & ACK
CLOSE_WAIT
SYN_RCVD
receive FIN send ACK
receive ACK send nothing
ESTABLISHED

---

(Start) CLOSED — CONNECT/SYN (Step 1 of the 3-way-handshake)
CLOSE/-
LISTEN/-
CLOSE/-
(Step 2 of the 3-way-handshake) SYN/SYN+ACK — LISTEN
SYN RECEIVED
RST/-
SEND/SYN
SYN SENT
SYN/SYN+ACK (simultaneous open)
Data exchange occurs
ACK/-
ESTABLISHED
SYN+ACK/ACK (Step 3 of the 3-way-handshake)
CLOSE/FIN
CLOSE/FIN
FIN/ACK

Active CLOSE
FIN WAIT 1 — FIN/ACK — CLOSING
FIN+ACK/ACK
ACK/-
FIN WAIT 2
FIN/ACK — TIME WAIT
Timeout

Passive CLOSE
CLOSE WAIT
CLOSE/FIN
LAST ACK

(Go back to start) CLOSED

- - -> unusual event
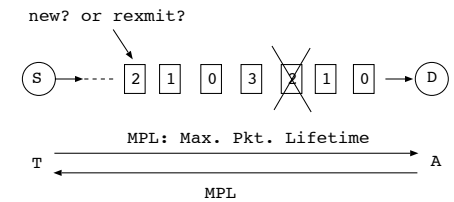——> client/receiver path
——> server/sender path

---

# Finite Sequence Number Space

Issues arising from having finite sequence number space:
1. choice of sequence space size
2. sequence number wrap around
3. initial sequence number (ISN) choice

---

# Sequence Number Space Size

If we had only 2 bits to keep track of sequence numbers:

new? or rexmit?

S - - - - 2  1  0  3  2  1  0 → D

MPL: Max. Pkt. Lifetime
T ←——————————————→ A
MPL

Let:
$A$: time taken by receiver to ACK packet
$T$: time sender continues retransmitting if an ACK is not received

Maximum Segment Lifetime (MSL): $2\text{MPL} + T + A$

## Duplicated Sequence Number

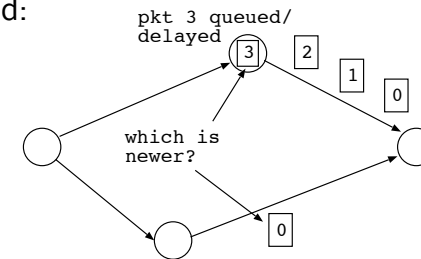Maximum Segment Lifetime (MSL): $2\mathrm{MPL} + T + A$

Want: no sequence number may be duplicated within an MSL

What would cause a sequence number to be duplicated within an MSL?
- sequence number wraps around (within a single connection)
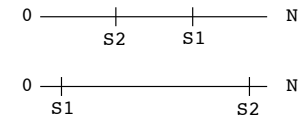- sequence number re-used (across connections)

## Sequence Number Wrap Around

Sequence space ($N$) is finite and sequence number can wrap around:



Assuming $s_1$ and $s_2$ are not more than $N/2$ apart, $s_1 > s_2$ if either:
1. $|s_1 - s_2| < N/2$ and $s_1 > s_2$, or
2. $|s_1 - s_2| > N/2$ and $s_1 < s_2$



## Required Sequence Number Size

To prevent duplicate sequence number: for $s_1 > s_2$, $s_1$ and $s_2$ cannot be more than $N/2$ apart ($|s_1 - s_2| < N/2$) within an MSL

For TCP, $N = 2^{32}-1$, $N/2$ maximum separation requirement means that only $n = 31$ bits are usable

Let $\mu$ be the transmission bandwidth, worst case, assuming sender can "fill the pipe", want:
$\mu < (N/2)/\mathrm{MSL}$ or, $2^n > \mu*\mathrm{MSL}$

- example: SF-NY MPL is 25 msec.
  let $\mathrm{MSL} = 2$ min, for $n = 31$ bits,
  $\mu$ must be $< 17.8$ MB/s (143 Mbps)

## Required Sequence Number Size

But TCP transmission is constrained by the receiver's advertised window (`rwnd`), which is of 16-bit size

Only 64 KB can be outstanding at any one time, which takes less than ⅓ of a second to clear the network even at "slow" T1 speed (1.5 Mbps)

So we don't have to worry about sequence number re-use due to wrap around

# Initial Sequence Number (ISN)

Sequence number for the very first byte

Why not always start with an ISN of 0?

IP addresses and port #s uniquely identify a
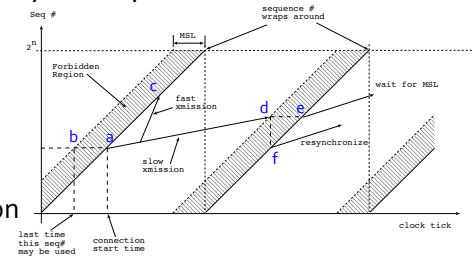connection, but port numbers get reused when:
•
•
and connections get "reincarnated"

We want ISNs that will not let packets from an old
connection that are still in flight to be mistaken for
packets from the new connection

# ISN from System Clock

Assume clock keeps ticking even when machine is down
Want: no seq. number may be duplicated within an MSL

Forbidden region: using
sequence number currently
in forbidden region could
cause duplicate ISN chosen
by "reincarnated" connection



What to do on hitting forbidden region (d)?
d) nothing: seqno duplicated by "reincarnated" connection
e) wait for MSL before resuming transmission
f) resynchronize sequence number
either (e) or (f), connection stalled

# TCP's Handling of ISN

TCP requires changing the ISN over time
• set from a 32-bit clock that ticks every $4\ \mu$seconds
• which only wraps around once every 4.55 hours
  • unlikely for reincarnated connection to share seqno

Connection cannot be reused for MSL time
• on connection tear-down, wait for 2MSL
  (TIME-WAIT state, 30 secs)
  `bind: Address already in use`
• on reboot, do not create connection for MSL
  (2 minutes boot time, so not a problem)