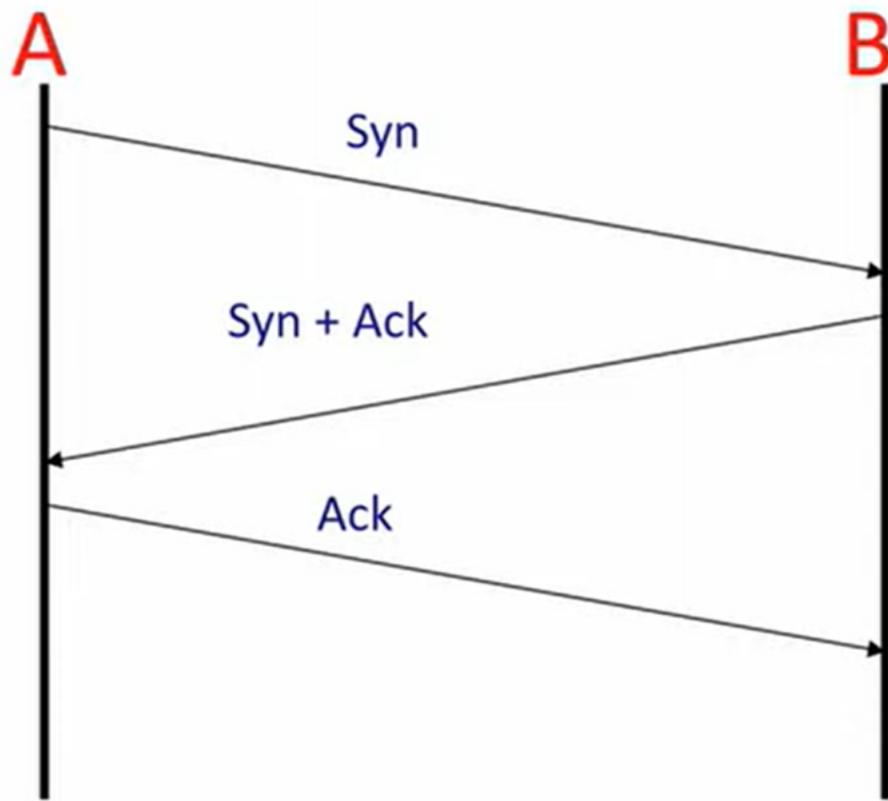


# Transport layer The TCP Service Model

三次握手建立连接

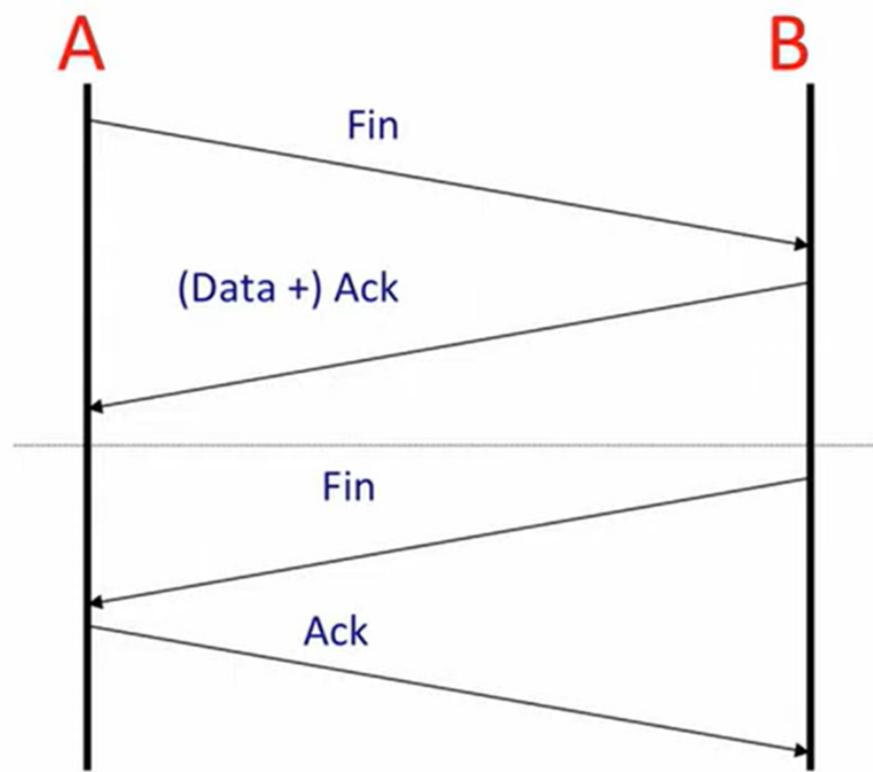
## Connection setup

### 3-way handshake



四次挥手断开连接

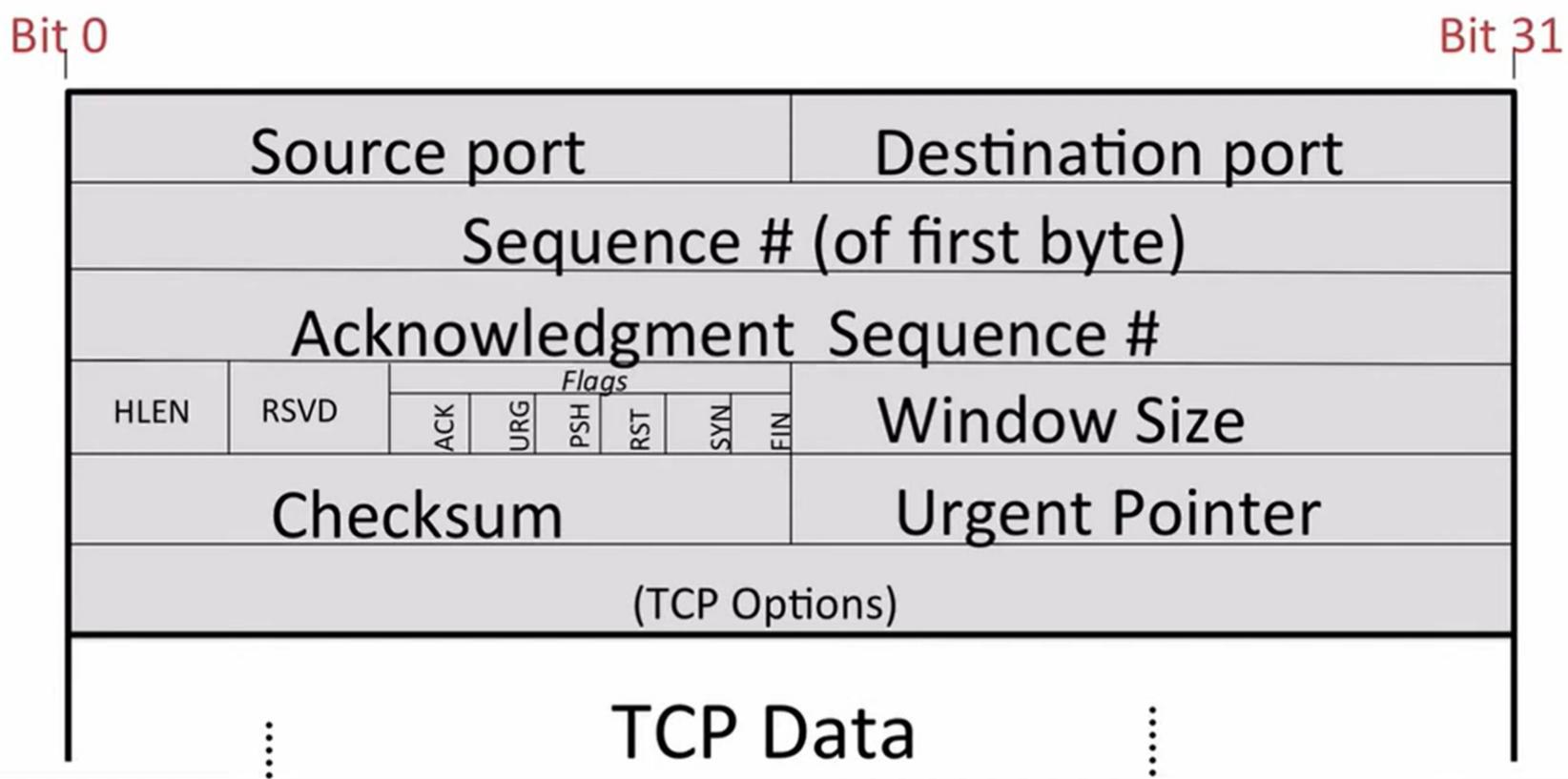
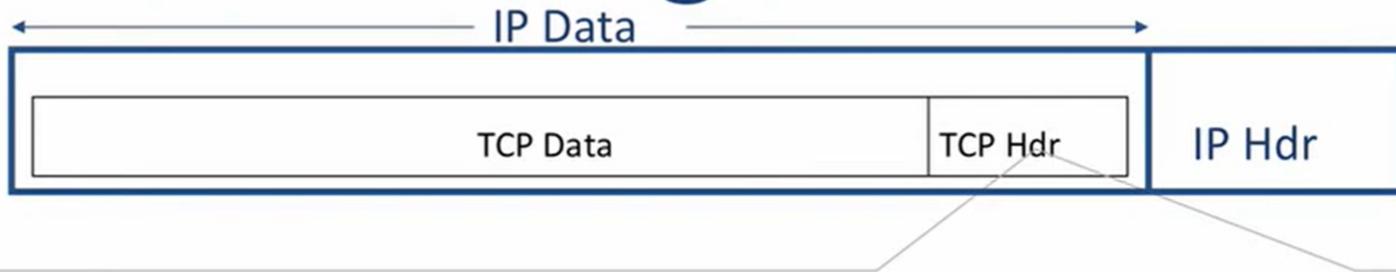
# Connection teardown



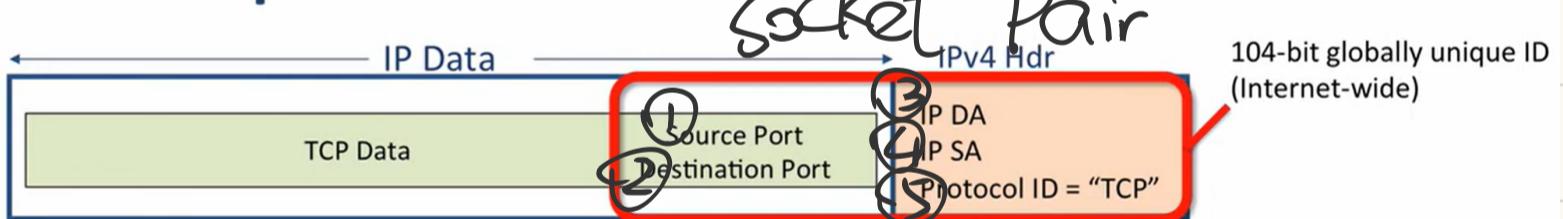
## The TCP Service Model

Property	Behavior
<i>Stream of bytes</i>	Reliable byte delivery service.
<i>Reliable delivery</i>	1. Acknowledgments indicate <u>correct delivery</u> . 2. Checksums detect <u>corrupted data</u> . 3. Sequence numbers detect <u>missing data</u> . 4. Flow-control prevents <u>overrunning receiver</u> .
<i>In-sequence</i>	Data delivered to application in <u>sequence transmitted</u> .
<i>(Congestion Control)</i>	Controls network <u>congestion</u> .

# The TCP Segment Format



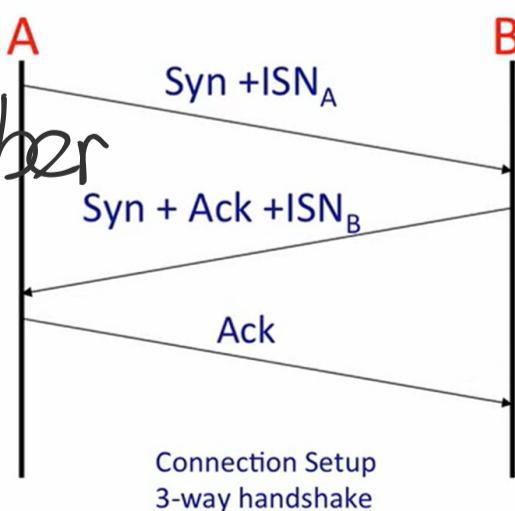
# The Unique ID of a TCP connection



1. Host A increments source port for every new connection

2. TCP picks ISN to avoid overlap with previous connection with same ID.

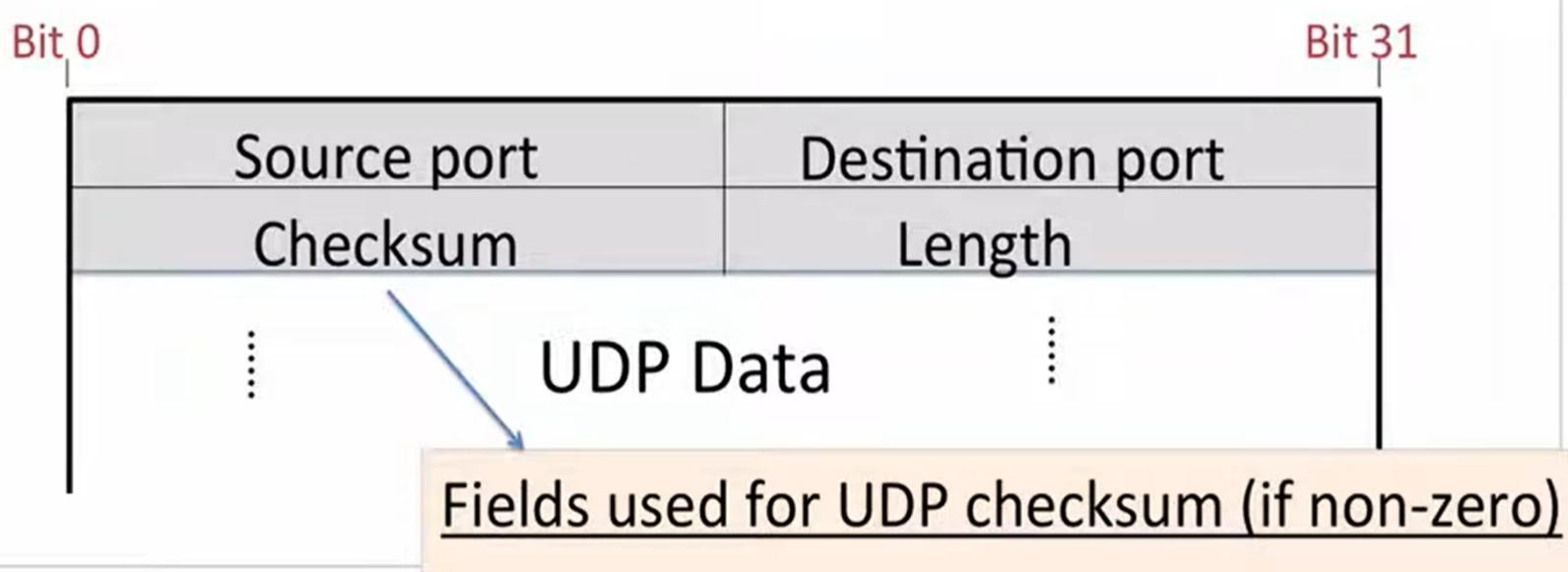
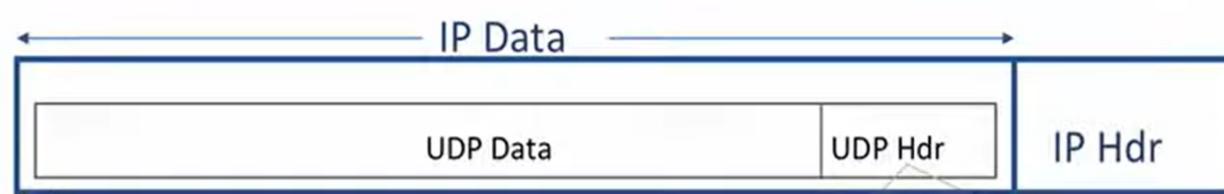
ISN means  
Initial Sequence Number



# The UDP Service Model

## User Datagram Protocol

### The UDP Datagram Format



### User Datagram Protocol (UDP)

Property	Behavior
<i>Connectionless Datagram Service</i>	No connection established. Packets may show up in any order.
<i>Self contained datagrams</i>	
<i>Unreliable delivery</i>	<ol style="list-style-type: none"><li>1. No acknowledgments.</li><li>2. No mechanism to detect missing or mis-sequenced datagrams.</li><li>3. No flow control.</li></ol>

# The ICMP Service Model

## Internet Control Message Protocol

### Making the Network Layer Work

三种不同的机制

#### 1. The Internet Protocol (IP)

- The creation of IP datagrams.
- Hop-by-hop delivery from end to end.

#### 2. Routing Tables

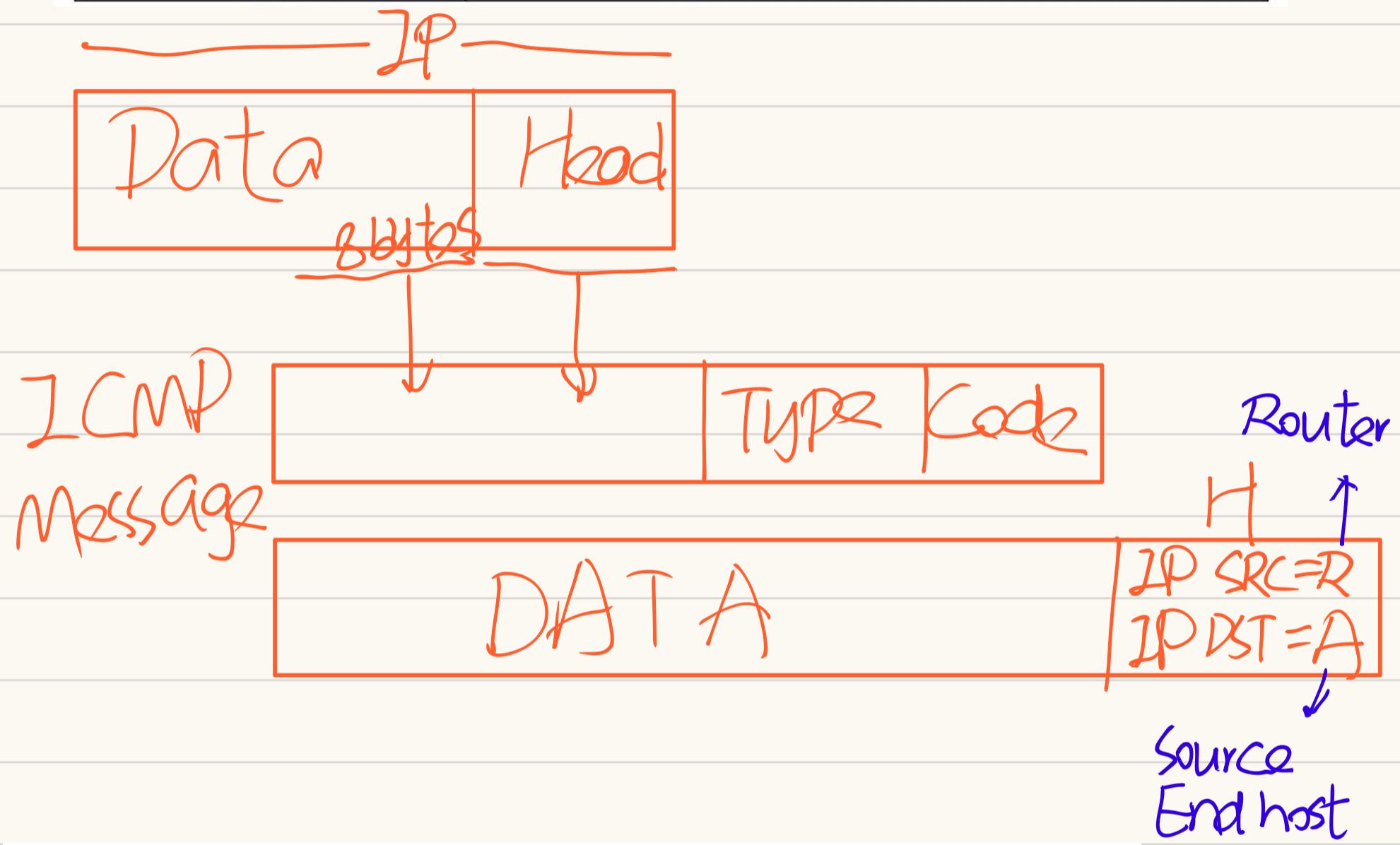
- Algorithms to populate router forwarding tables

#### 3. Internet Control Message Protocol (ICMP)

- Communicates network layer information between end hosts and routers
- Reports error conditions
- Helps us diagnose problems

# The ICMP Service Model

Property	Behavior
<i>Reporting Message</i>	Self-contained message reporting error.
<i>Unreliable</i>	Simple datagram service – no retries.



## (Some) ICMP Message Types

ICMP Type	ICMP Code	Description
0	0	Echo Reply (used by ping)
3	0	Destination Network Unreachable
3	1	Destination Host Unreachable
3	3	Destination Port Unreachable
8	0	Echo Request (used by ping)
11	0	TTL Expired (used by traceroute)

- ① ICMP provides info about the network layer to end hosts and routers
- ② It sits above IP and is therefore strictly a transport layer mechanism.
- ③ The commonly used tools "Ping" and "traceroute" both rely on ICMP.

# The End-to-end Principle

## The End-To-End Principle

The function in question can completely and correctly be implemented only with the knowledge and help of the application standing at the end points of the communication system. Therefore, providing that questioned function as a feature of the communication system itself is not possible. (Sometimes an incomplete version of the function provided by the communication system may be useful as a performance enhancement.) We call this line of reasoning... "the end-to-end argument"

However, this is generally seen as a performance optimization rather than a complete solution

- Saltzer, Reed, and Clark,  
End-to-end Arguments in System Design, 1984

Functions should primarily be implemented at the endpoints of a communication system rather than within the network itself.

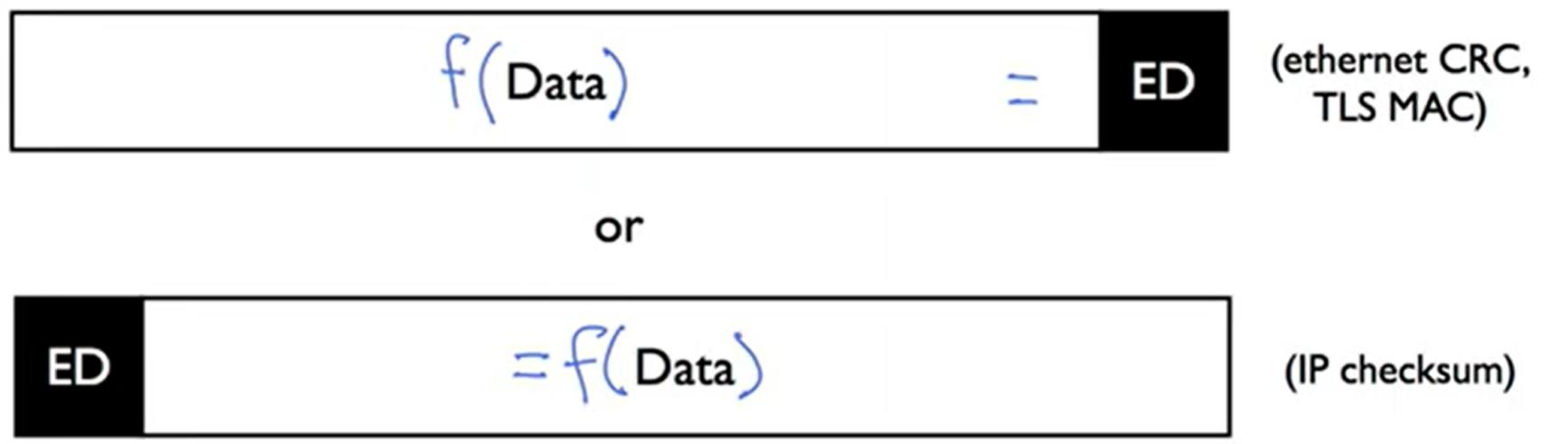
## “Strong” End to End

The network's job is to transmit datagrams as efficiently and flexibly as possible. Everything else should be done at the fringes...

Should not be implemented – [RFC 1958]  
in the middle

(Strictly)

# Error Detection : 3 schemes checksum, CRC, MAC



- Checksum adds up values in packet (IP,TCP)
  - ▶ Very fast, cheap to compute even in software
  - ▶ Not very robust
- Cyclic redundancy code computes remainder of a polynomial (Ethernet)
  - ▶ More expensive than checksum (easy today, easy in hardware)
  - ▶ Protects against any 2 bit error, any burst  $\leq c$  bits long, any odd number of errors
- Message authentication code: cryptographic transformation of data (TLS)
  - ▶ Robust to malicious modifications, but not errors
  - ▶ If strong, any 2 messages have a  $2^{-c}$  chance of having the same code

# IP Checksum

- IP, UDP, and TCP use one's complement checksum algorithm:
  - ▶ Set checksum field to 0, sum all 16-bit words in packet
  - ▶ Add any carry bits back in:  $0x8000 + 0x8000 = 0x0001$
  - ▶ Flip bits (0xc379 becomes 0x3c86), unless 0xffff, then checksum is 0xffff
  - ▶ To check: sum whole packet including checksum, should be 0xffff
- Benefits: fast, easy to compute and check
  - ▶ Motivated by earliest software implementations
- Drawbacks: poor error detection
  - ▶ Only guarantees detecting a single bit error
  - ▶ Can detect other errors, but actual guarantees are both weak and complex

# Cyclic Redundancy Check (CRC)

- Cyclic Redundancy Check (CRC): distill  $n$  bits of data into  $c$  bits,  $c \ll n$ 
  - ▶ Can't detect all errors:  $2^{-c}$  chance another packet's CRC matches
- CRC designed to detect certain forms of errors: stronger than checksum
  - ▶ Any message with an odd number of bit errors
  - ▶ Any message with 2 bits in error
  - ▶ Any message with a single burst of errors  $\leq c$  bits long
- Link layers typically use CRCs
  - ▶ Fast to compute in hardware (details in a moment)
  - ▶ Can be computed incrementally
  - ▶ Good error detection for physical layer burst errors

# Diversion: CRC Mathematical Basis

- Cyclic Redundancy Check (CRC): distill  $n$  bits of data into  $c$  bits,  $c \ll n$
- Uses polynomial long division
  - ▶ Consider the message  $M$  a polynomial with coefficients 0 or 1 (pad with  $c$  zeroes)
    - E.g.,  $M = 10011101 = x^7 + x^4 + x^3 + x^2 + 1$
  - ▶ Use a generator polynomial  $G$  of degree  $c$  also with coefficients 0 or 1
    - Pad first term (always 1) for frustrating historical reasons
    - E.g.  $G = 1011 = x^4 + x^3 + x + 1$
    - USB (CRC-16) =  $0x8005 = x^{16} + x^{15} + x^2 + 1$
  - ▶ Divide  $M$  by  $G$ , the remainder is the CRC: pick  $G$  carefully!
- Append CRC to message  $M$ :  $M' = M + CRC$ 
  - ▶ Long division of  $M'$  with  $G$  has a remainder of 0

## MAC

- Message Authentication Code (MAC)
  - ▶ Not to be confused with Media Access Control (MAC)!
- Uses cryptography to generate  $c = MAC(M, s)$ ,  $|c| \ll |M|$ 
  - ▶ Using  $M$  and secret  $s$ , can verify  $c = MAC(M, s)$
  - ▶ If you don't have  $s$ , very very hard to generate  $c$
  - ▶ Very very hard to generate an  $M$  whose MAC is  $c$
  - ▶  $M + c$  means the other person probably has the secret (or they're replayed!)
- Cryptographically strong MAC means flipping one bit of  $M$  causes every bit in the new  $c$  to be randomly 1 or 0 (no information)
  - ▶ Not as good for error detection as a CRC!
  - ▶ But protects against adversaries

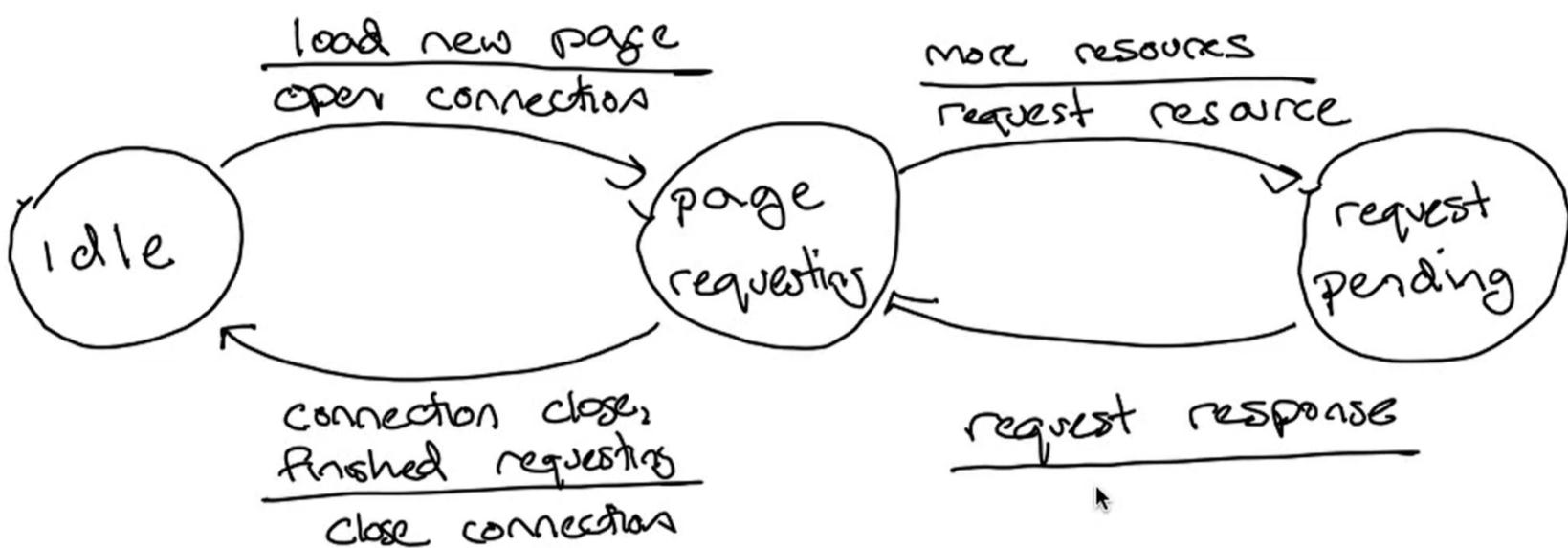
For each error detection algorithm, mark which errors it can guarantee catching, if any. The CRCs use a good generator polynomial and the MAC algorithm is cryptographically strong.

Algorithm	Single bit error	Run of 2 bit errors	Run of 9 bit errors	Two bit errors 100 bits apart
8-bit checksum	X	Z	Z	Z
16-bit checksum	X	Z	Z	Z
8-bit CRC	Y	X	Z	X
16-bit CRC	Y	Z	X	Z
32-bit MAC	Z	Z	Z	Z

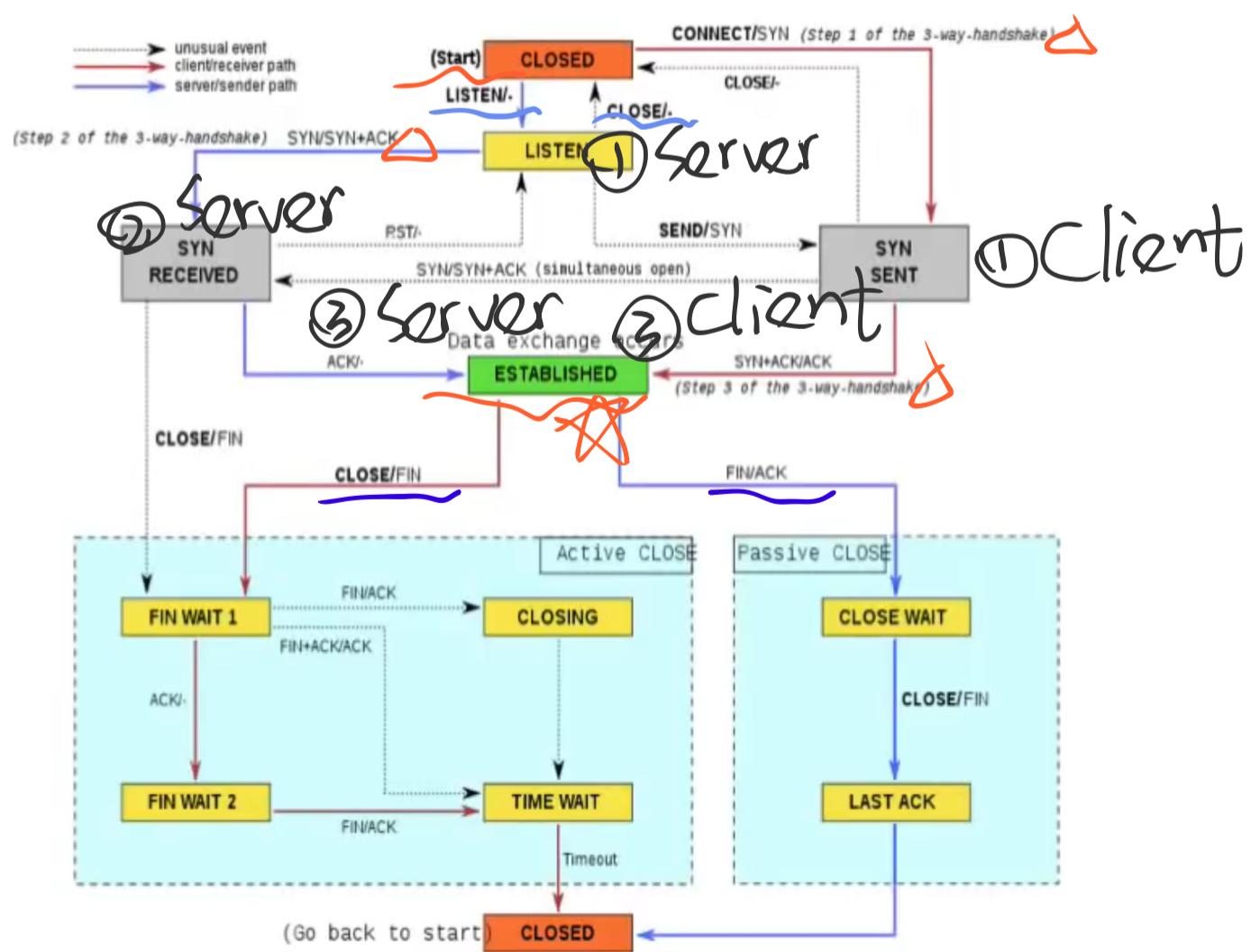
# Finite State Machines

## Protocol Specification

# FSM Example: HTTP Request



# FSM Example: TCP Connection

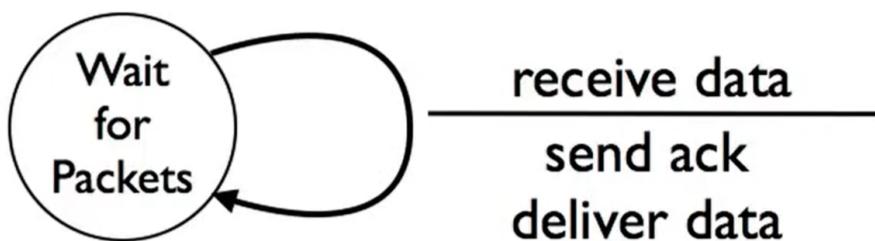


# Flow Control

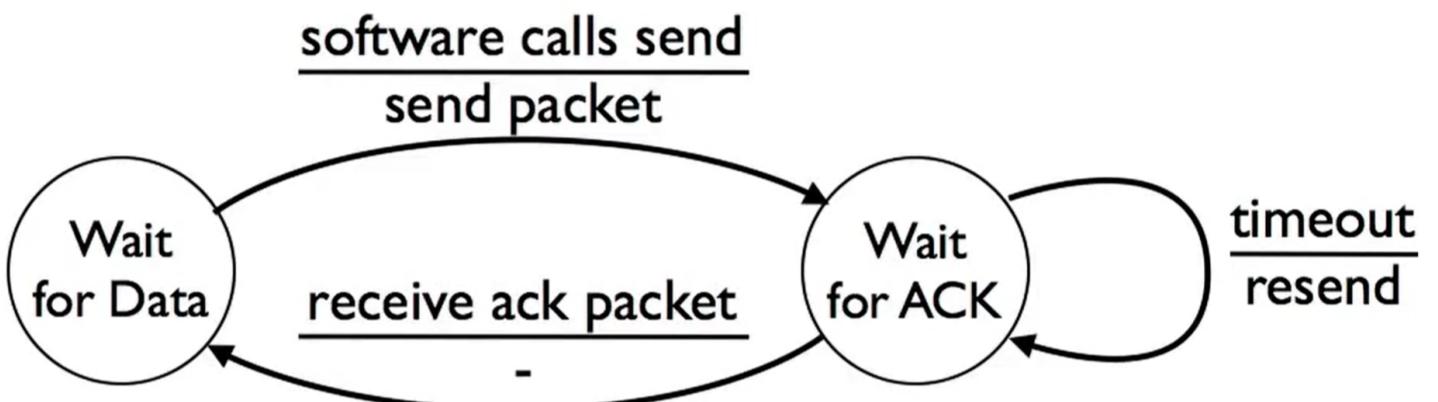
## Stop-and-Wait

### Stop and Wait FSM

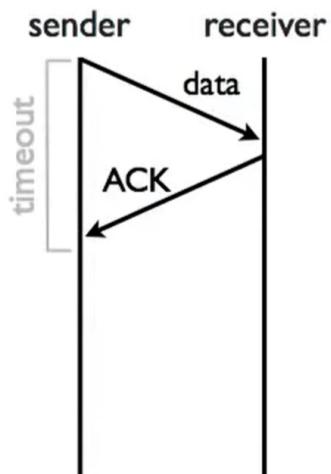
Receiver FSM



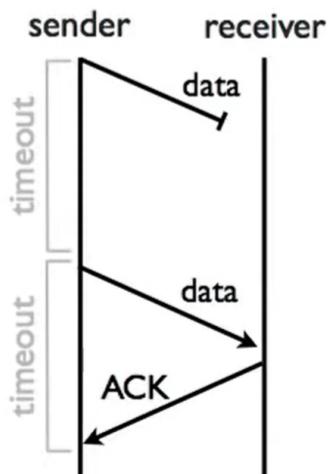
Sender FSM



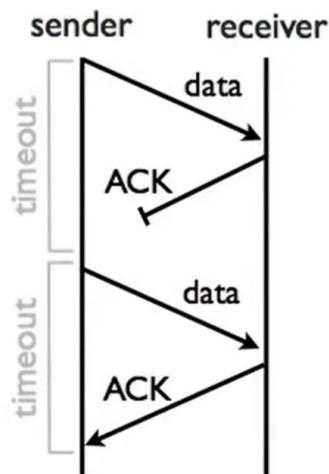
### Example Executions



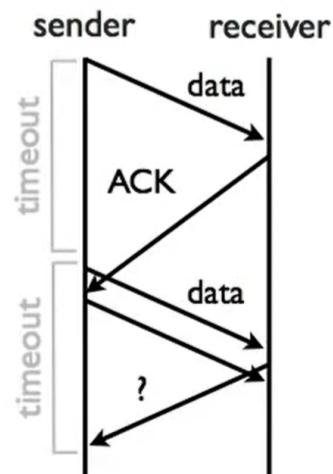
No Loss



Data Loss



ACK Loss

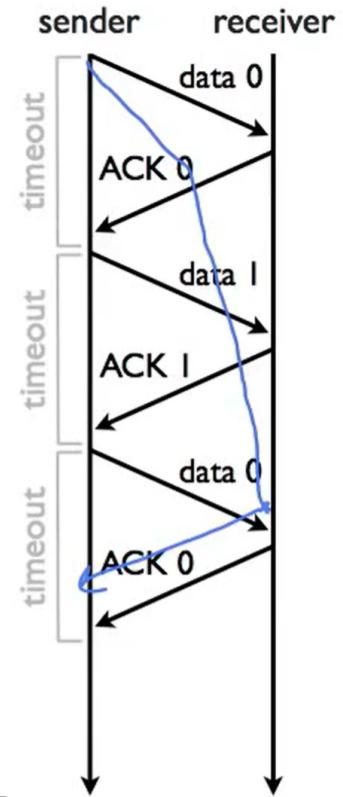


ACK Delay



# To Solve the ACK Delay Duplicates

- Use 1-bit counter in data and acknowledgments
  - Receiver can tell if new data or duplicate
- Some simplifying assumptions
  - Network does not duplicate packets
  - Packets not delayed multiple timeouts



## Conclusion:

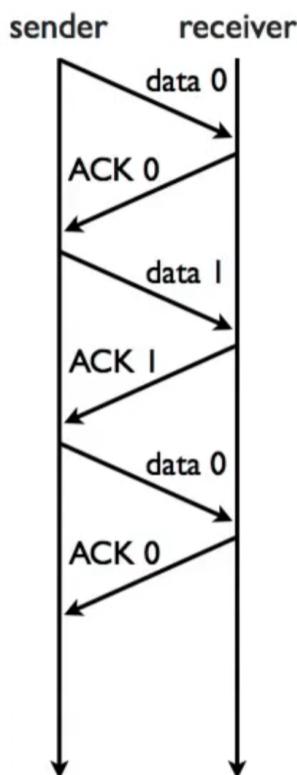
- ① At most one packet in flight at any time.
- ② Sender sends one packet
- ③ Receiver sends ACK packet when it receives data.
- ④ On receiving ACK, sender sends new data
- ⑤ On timeout, sender sends current data
- ⑥ Use 1-bit counter to detect duplicates.

# Stop and Wait & Sliding Window

↓ Generate

## Example Execution

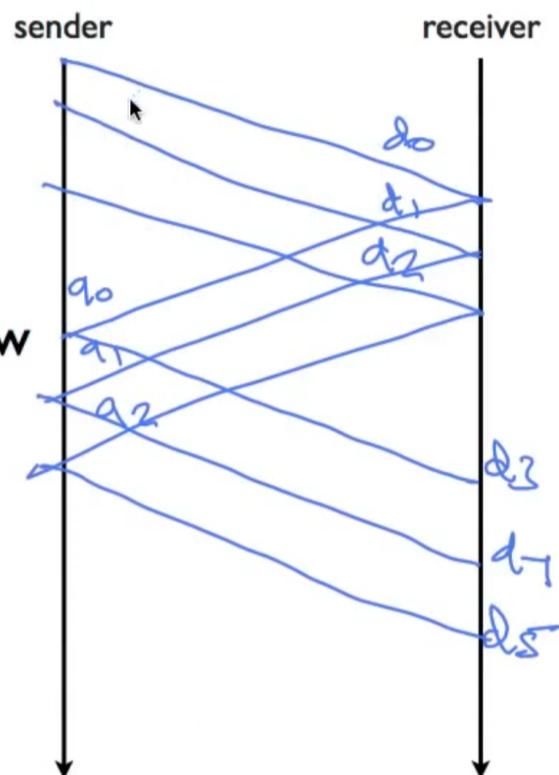
Stop and Wait



Sliding Window

$w = 3$

5



## Sliding Window Sender

- Every segment has a sequence number (SeqNo)
- Maintain 3 variables
  - ▶ Send window size (SWS)
  - ▶ Last acknowledgment received (LAR)
  - ▶ Last segment sent (LSS)
- Maintain invariant:  $(LSS - LAR) \leq SWS$
- Advance LAR on new acknowledgment
- Buffer up to SWS segments

Also means:

$$LSS \leq LART + SWS$$

For example

$$SWS = 5 \quad LAR = 11$$

Then only 12, 13, 14, 15, 16 allow to be sent



# Sliding Window Receiver

Also means

$$LAS \leq RWS + LSR$$

For exampk: if  $RWS = 5$

$$LSR = 3$$

Then it can

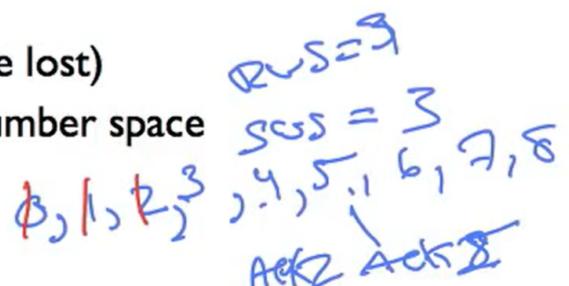
only accept 4, 5, 6, 7, 8

means that I've received 3 and everything before it, not just 3

## RWS, SWS, and Sequence Space

ensure no buffer area

- $RWS \geq 1, SWS \geq 1, RWS \leq SWS$  is wasted
- If  $RWS = 1$ , "go back N" protocol, need  $SWS+1$  sequence numbers
- If  $RWS = SWS$ , need  $2SWS$  sequence numbers
- Generally need  $RWS+SWS$  sequence numbers
  - ▶ RWS packets in unknown state (ACK may/may not be lost)
  - ▶ SWS packets in flight must not overflow sequence number space



## Sliding Window Flow Control

1. Allow a "window" of unacknowledged packets in flight
2. When acknowledgement arrives, advance window
3. Required sequence number space size depends on window sizes.

# Retransmission Strategies

Go Back N

Selective Repeat

Pessimistic

Optimistic

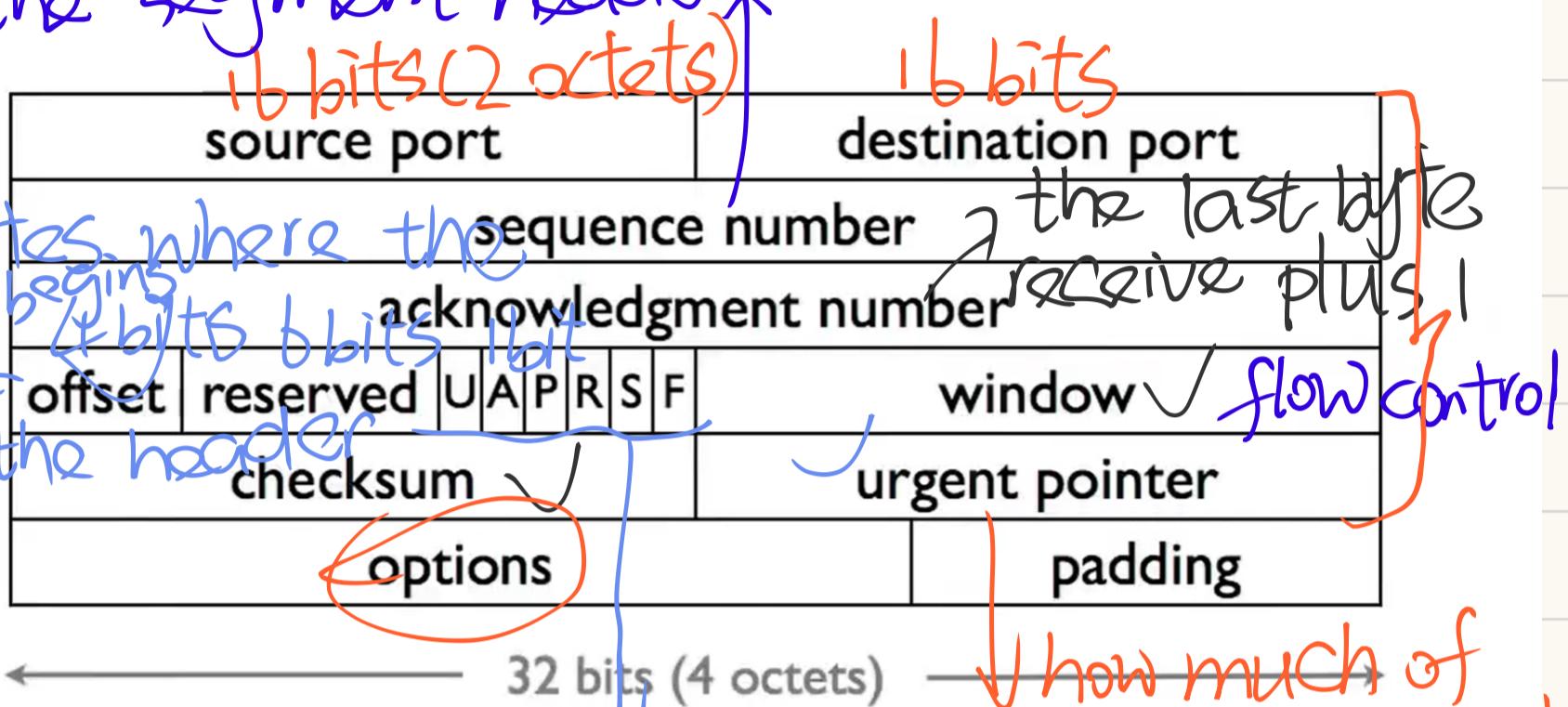
Sender has a window of  $N$   
Receiver has a window of 1

Sender and Receiver  
both have a  
window of  $N$

# TCP Header

## TCP Header

is of the first byte of the data region which follows the segment header.



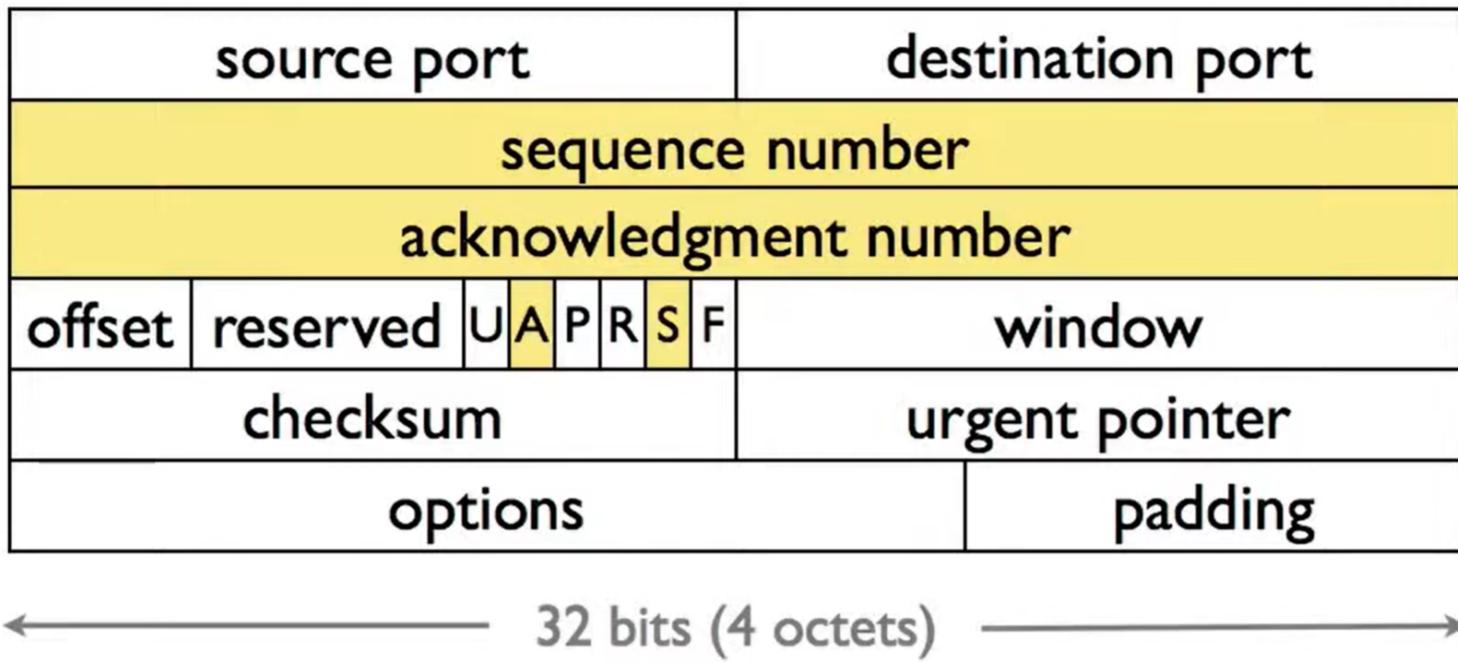
(whether urgent) (whether buffer)  
U → urgent, P → push  
A → ACK, R → RST, S → SYN, F → FIN

The TCP header can have a variable length due to options field, that is the reason why "offset" is needed.

The smallest TCP header (without any options) is 20 bytes (5 words of 4 bytes each)

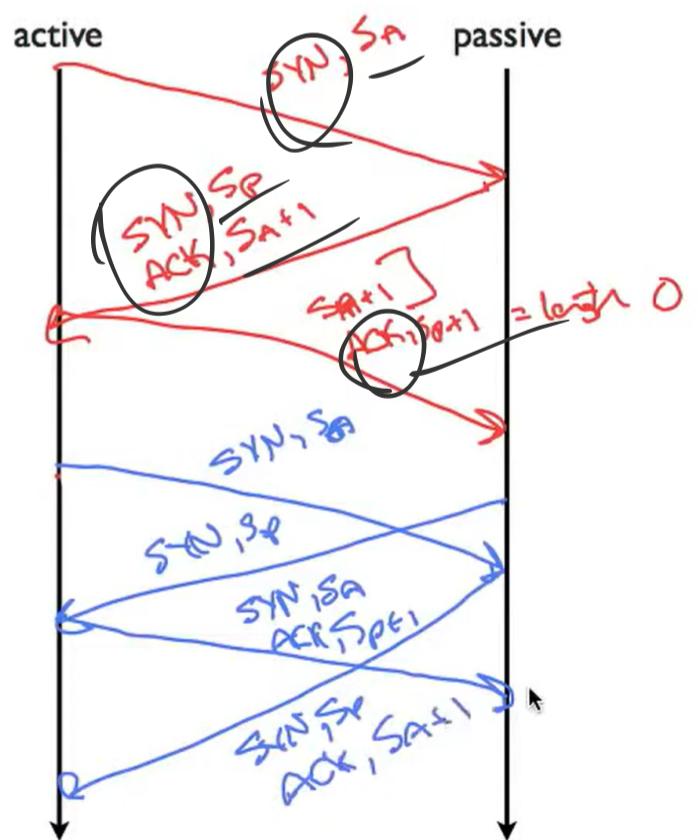
# TCP Set Up and Teardown

## Connection Setup

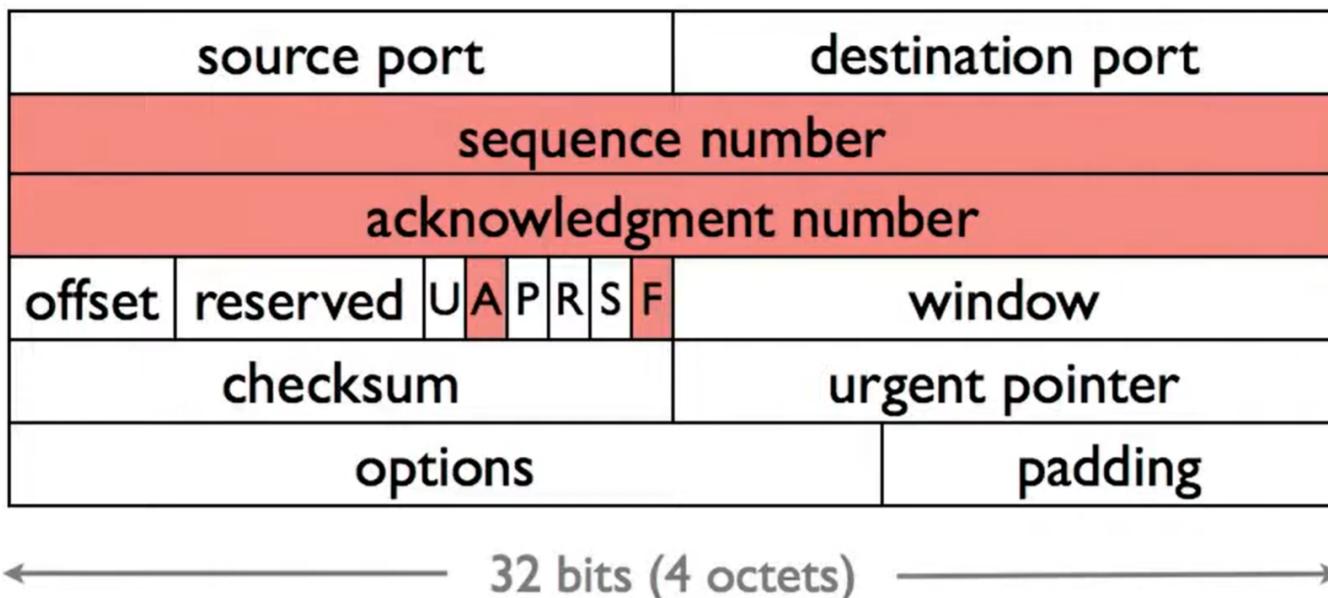


## 3-way Handshake

- Active opener sends first packet
  - ▶ SYN with sequence number
- Passive opener responds
  - ▶ SYN with sequence number
  - ▶ ACKs active opener's SYN packet
- Active opener responds
  - ▶ ACKs passive opener's SYN packet
- Also support “simultaneous open”
  - ▶ Two SYNs pass each other
  - ▶ Each side ACKs the other



# Connection Teardown



# Connection Teardown

- FIN bit says no more data to send
    - ▶ Caused by `close()` or `shutdown()` on other end
  - Both sides must send FIN to terminate a connection
  - Typical teardown exchange:
    - ▶ A → B: **FIN**, seq  $S_A$ , ack  $S_B$
    - ▶ B → A: ack  $S_A + 1$
    - ▶ B → A: **FIN**, seq  $S_B$ , ack  $S_A + 1$
    - ▶ A → B: ack  $S_B + 1$
  - Can also have simultaneous close
  - Can A and B forget about closed socket after final message?
- bidirectional connection*