# Lesson 4: End-to-End Protocols

Van-Linh Nguyen

Fall 2024

Online video link:   https://youtu.be/of5CmizjSW8

國立中正大學
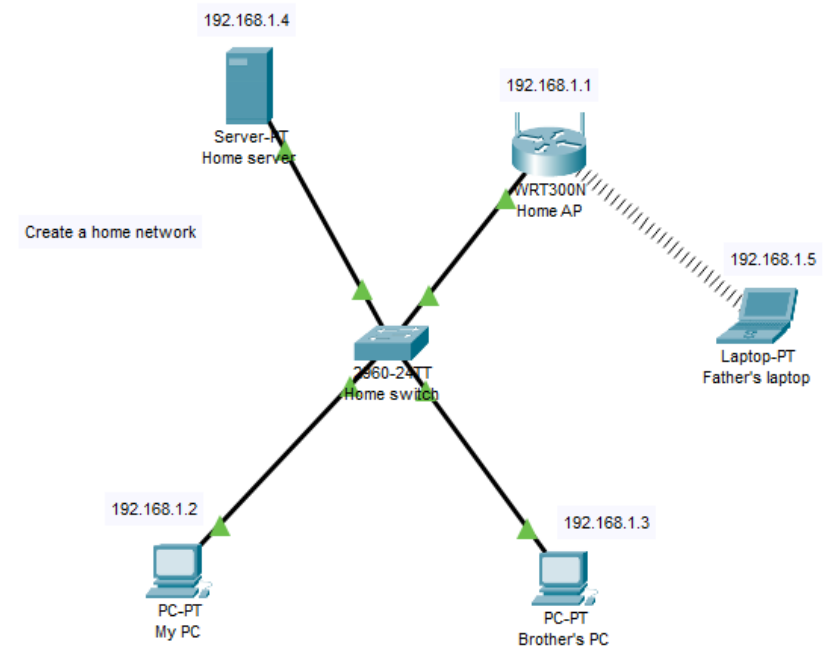Cybersecurity Lab

# Outline

- Lab2: Create a network service/app in Packet Tracer

- Lab3: Create a simple Internet network

- Network protocols
  1. UDP
  2. TCP

# Lab2: Service/application

- This lab is to recall knowledge about application and services

- Based on the first lab, please create a new DNS server with IP 192.168.1.6

- Assign the webserver at 192.168.1.4 with the domain ccu.edu.tw

- Test whether PCs can access Web on the Home server by using the domain name

- Activate FTP services on the Home server and create an account admin/admin → Test FTP connection from laptop
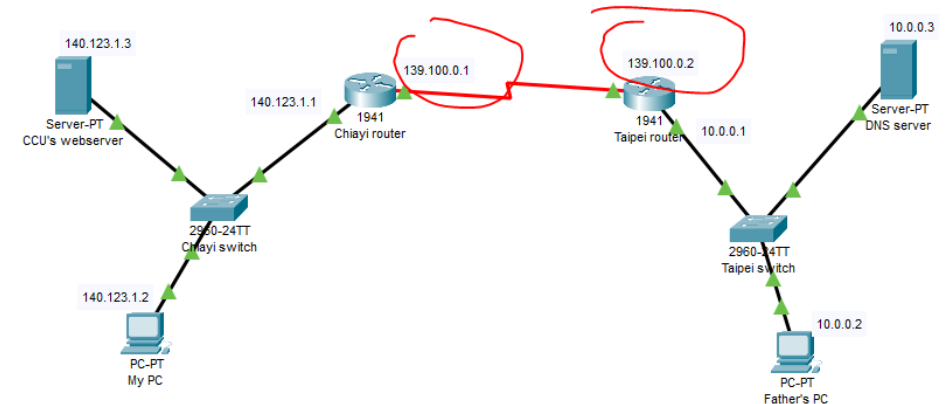


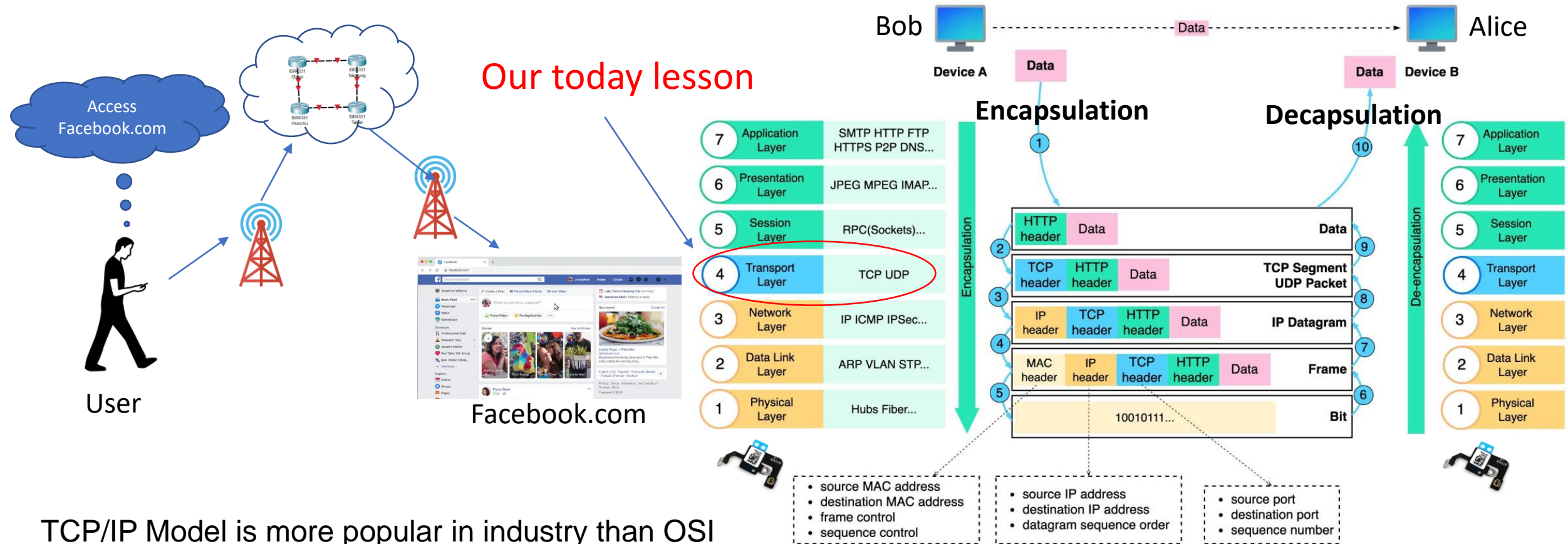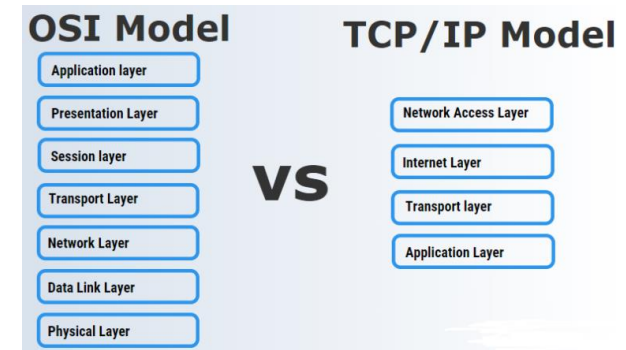國立中正大學
Cybersecurity Lab

# Lab3: Create an Internet network

- Create a simple Internet between Chiayi and Taipei

- The Chiayi network starts with the network IP 140.123.1.0/24

- The Taipei network starts with the network IP 10.0.0.0/16

- Test whether PCs in the Chiayi network can access PCs at the Taipei network

- Create a Webserver at the Chiayi network (CCU's website) and a DNS server at the Taipei network

- Assign the webserver at 140.123.1.3 with the domain ccu.edu.tw

- Create a DNS record for the website at DNS server

- Test whether Father's PC can access CCU's website via domain name



Network IP between Chiayi router and Taipei router must be different
from the connected networks (140.123.1.0, 10.0.0.0)

In the sample lab, we use IP 139.100.0.0 for that network and the IPs for serial port can start with 139.100.0.1 and 139.100.0.2 as
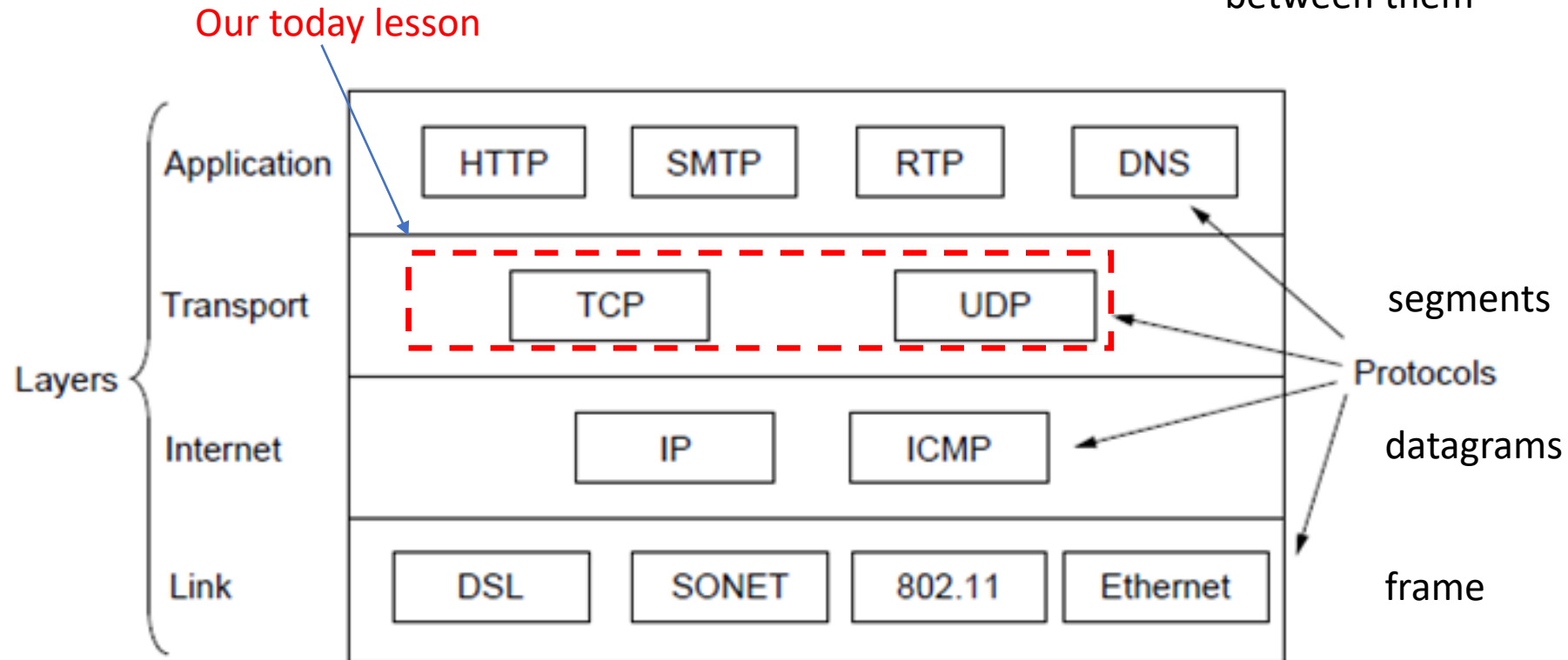Illustrated in the above figure

國立中正大學
Cybersecurity Lab

# Encapsulation/Decapsulation



Our today lesson

Facebook.com

User

Access Facebook.com

Bob — Data — Alice

Encapsulation

Decapsulation

Our today lesson — Transport Layer — TCP UDP

| 7 | Application Layer | SMTP HTTP FTP HTTPS P2P DNS... |
| 6 | Presentation Layer | JPEG MPEG IMAP... |
| 5 | Session Layer | RPC(Sockets)... |
| 4 | Transport Layer | TCP UDP |
| 3 | Network Layer | IP ICMP IPSec... |
| 2 | Data Link Layer | ARP VLAN STP... |
| 1 | Physical Layer | Hubs Fiber... |

- source MAC address
- destination MAC address
- frame control
- sequence control

- source IP address
- destination IP address
- datagram sequence order

- source port
- destination port
- sequence number

TCP/IP Model is more popular in industry than OSI
OSI is for academic research

https://blog.bytebytego.com/

國立中正大學 Cybersecurity Lab

# End-to-end Protocols

The transport layer has responsibility for establishing a temporary communication session between two applications and delivering data between them



https://www.dcs.bbk.ac.uk/~ptw/teaching/IWT/transport-layer/notes.html

# Message encapsulation through layers

Web



Fig. 1: Message transmission through layers.

- The application Layer uses the HTTP protocol

- The transport layer uses the TCP (Transmission Control) Protocol

- The network layer uses The IP (Internet) Protocol

- The data link layer of the LANs typically uses the Ethernet protocol

- Protocol Data Units (PDUs)

國立中正大學
Cybersecurity Lab

# End-to-end Protocols

- Below the presentation and session layer is transport layer

- A transport protocol can be expected to provide
  - Guarantees message delivery
  - Delivers messages in the same order they were sent
  - Delivers at most one copy of each message
  - Supports arbitrarily large messages
  - Supports synchronization between the sender and the receiver
  - Allows the receiver to apply flow control to the sender
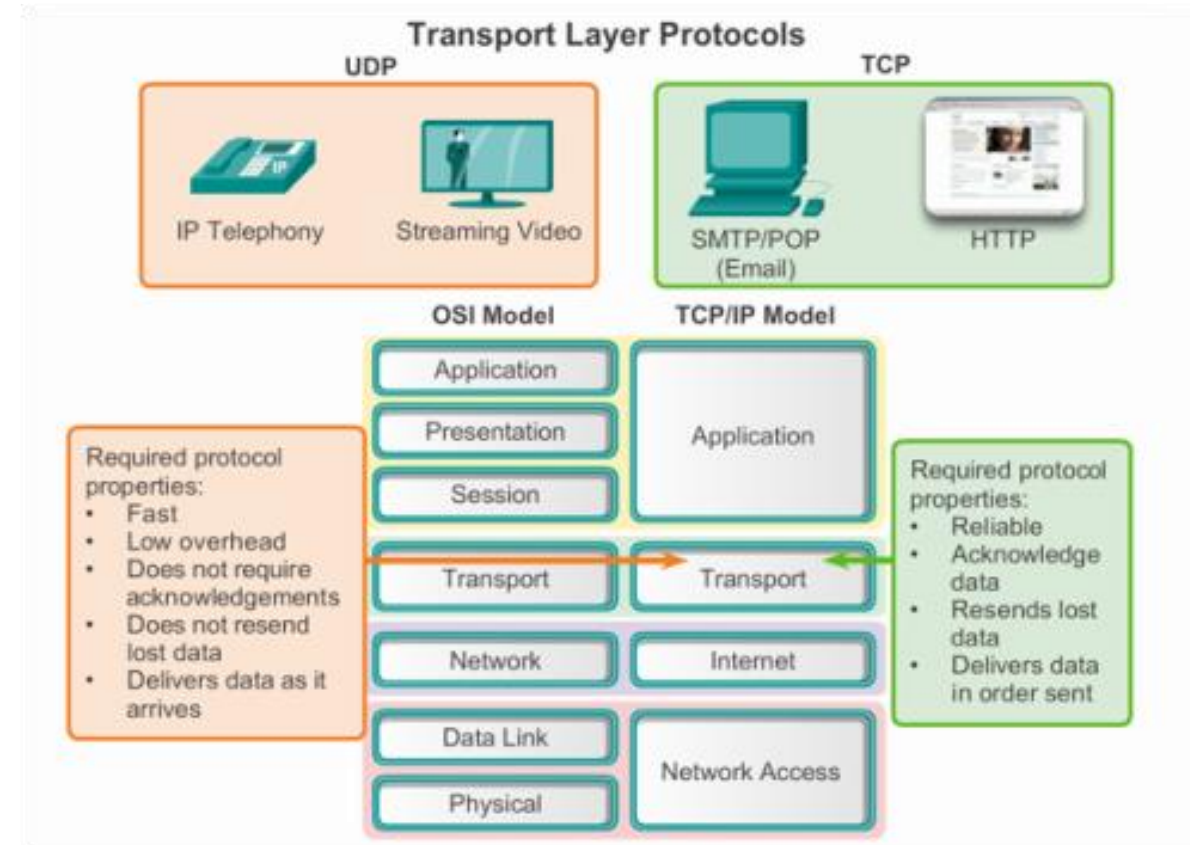  - Supports multiple application processes on each host

# End-to-end Protocols

- Typical limitations of the network on which transport protocol will operate
    - Drop messages
    - Reorder messages
    - Deliver duplicate copies of a given message
    - Limit messages to some finite size
    - Deliver messages after an arbitrarily long delay

# End-to-end Protocols

- Challenge for Transport Protocols
  - Develop algorithms that turn the less-than-desirable properties of the underlying network into the high level of service required by application programs

- Connectionless transport layer
  1. Treat each packet as an individual and delivery to the destination
  2. The receiver doesn't send acknowledgement of the packet

- Connection oriented transport layer
  1. Establish connection between the sender and the receiver before transmitting data
  2. The receiver sends acknowledgement for each received packet

國立中正大學
Cybersecurity Lab

# TCP vs UDP

- Two protocols are used
    1. Connectionless: User Datagram Protocol (UDP)
    2. Connection oriented: Transmission Control Protocol (TCP)
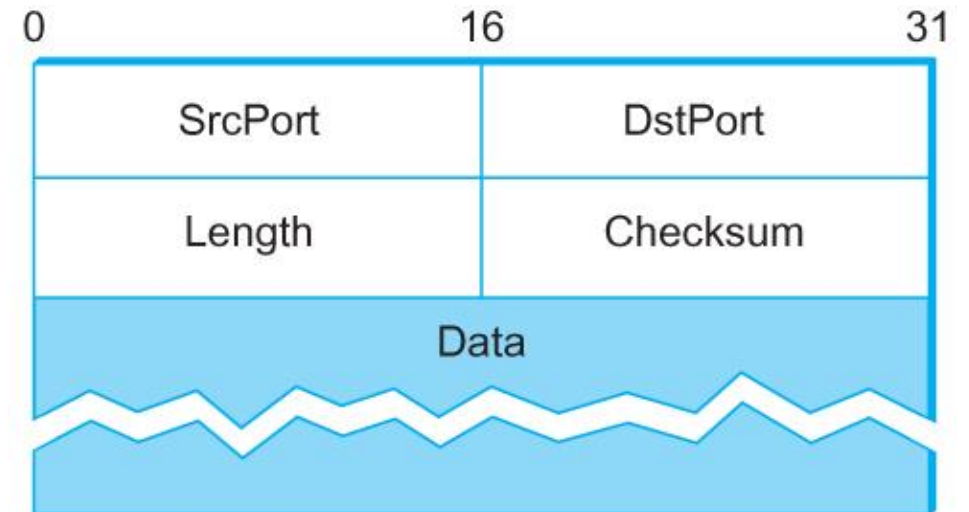
# User Datagram Protocol (UDP)

- UDP is suitable for those applications where <span style="color:red">data loss doesn't affect the perceived quality of the service</span>

- Multiplayer Online Games, VoIP, and Video Live Streaming

- Fast but non-guaranteed transfer. It is also called "best effort" transfer

# UDP Datagram (header)

- Each UDP segment must include UDP header fields identifying the socket connection

- These header fields are the *source port number field* and the *destination port number field*

- Each port number is a 16-bit number: 0 to 65535

- Port numbers below 1024 are called *well-known ports* and are reserved for standard services

| Port number | Application protocol | Description | Transport protocol |
|---|---|---|---|
| 21 | FTP | File transfer | TCP |
| 23 | Telnet | Remote login | TCP |
| 25 | SMTP | E-mail | TCP |
| 53 | DNS | Domain Name System | UDP |
| 79 | Finger | Lookup information about a user | TCP |
| 80 | HTTP | World wide web | TCP |
| 110 | POP-3 | Remote e-mail access | TCP |
| 119 | NNTP | Usenet news | TCP |
| 161 | SNMP | Simple Network Management Protocol | UDP |

Famous ports for some services

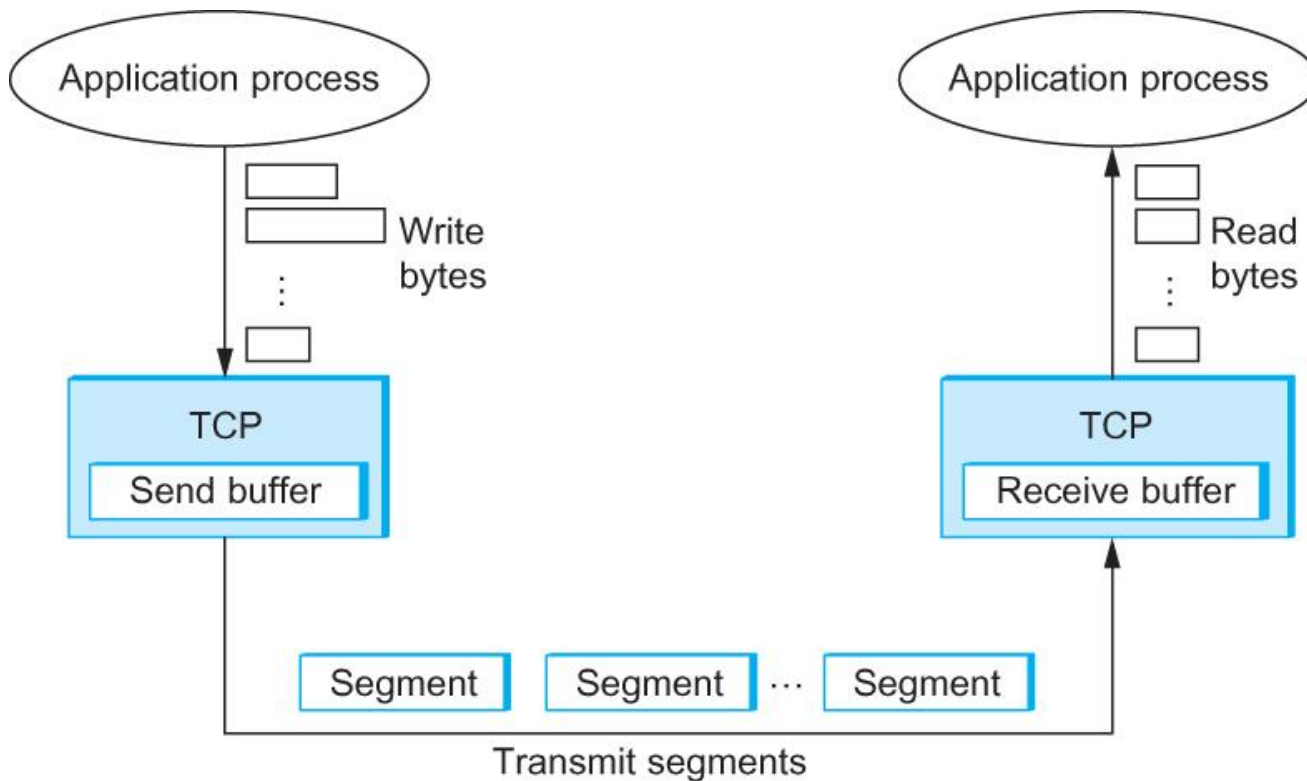| 0 | 16 | 31 |
|---|---|---|
| SrcPort | | DstPort |
| Length | | Checksum |
| Data | | |

Format for UDP header

國立中正大學
Cybersecurity Lab
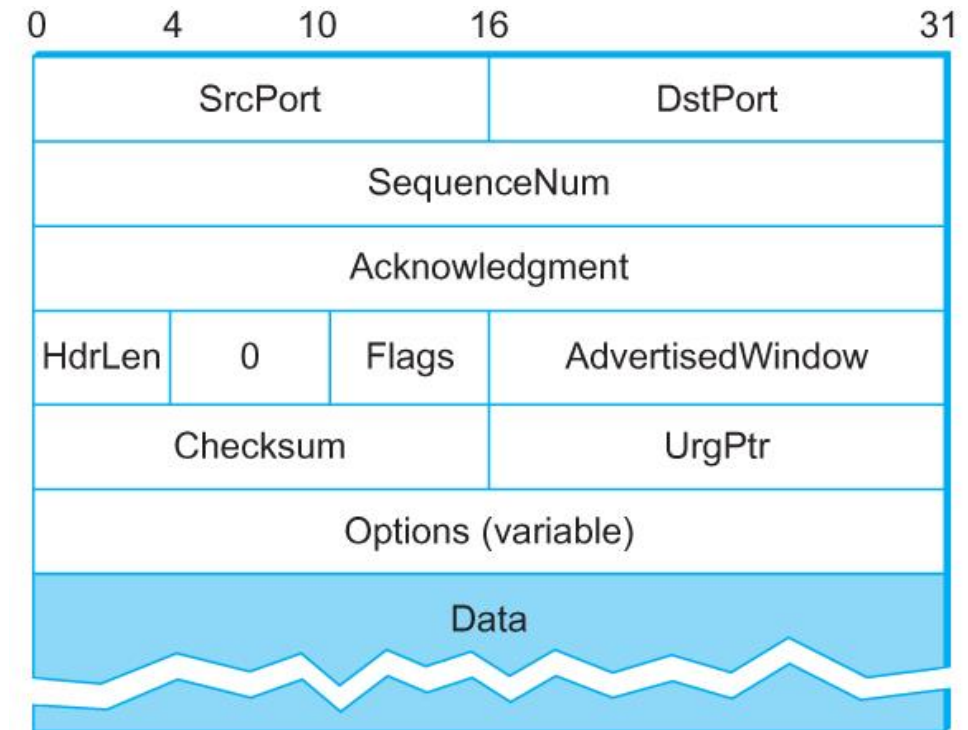
# Transmission Control Protocol (TCP)

- In contrast to UDP, Transmission Control Protocol (TCP) offers the following services
  - Reliable
  - Connection oriented
  - Byte-stream service
- TCP is also known as a three-way handshake protocol

**THREE - WAY HANDSHAKE (TCP)**



(1) SYN

(2) SYN/ACK

(3) ACK ✔

SYN = SYNCHRONIZATION    ACK = ACKNOWLEDGEMENT

# TCP Segment & TCP Header



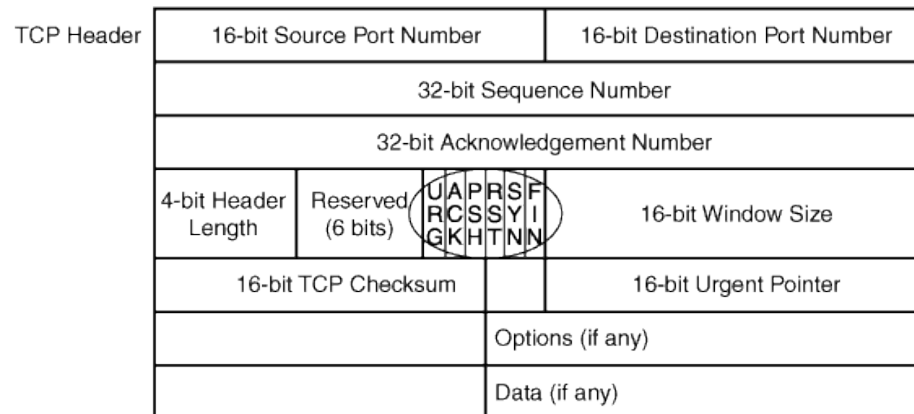How TCP manages a byte stream.



TCP Header Format

# TCP Header

- The **SrcPort** and **DstPort**: identify the source and destination ports, respectively.

- The **Acknowledgment**, **SequenceNum**, and **AdvertisedWindow**: involved in TCP's sliding window algorithm.

- Because TCP is a byte-oriented protocol, each byte of data has a sequence number; the **SequenceNum** field contains the <span style="color:red">sequence number</span> for the first byte of data carried in that segment.

- The **Acknowledgment** and **AdvertisedWindow** fields carry information about the flow of data going in the other direction.

| Port number | Application protocol | Description | Transport protocol |
|---|---|---|---|
| 21 | FTP | File transfer | TCP |
| 23 | Telnet | Remote login | TCP |
| 25 | SMTP | E-mail | TCP |
| 53 | DNS | Domain Name System | UDP |
| 79 | Finger | Lookup information about a user | TCP |
| 80 | HTTP | World wide web | TCP |
| 110 | POP-3 | Remote e-mail access | TCP |
| 119 | NNTP | Usenet news | TCP |
| 161 | SNMP | Simple Network Management Protocol | UDP |

國立中正大學
Cybersecurity Lab

# TCP Header

- The 6-bit Flags field is used to relay control information between TCP peers.

- The possible flags include SYN, FIN, RESET, PUSH, URG, and ACK.

- The SYN and FIN flags are used when establishing and terminating a TCP connection, respectively.

- The ACK flag is set any time the Acknowledgment field is valid, implying that the receiver should pay attention to it.

| TCP Header | | |
|---|---|---|
| 16-bit Source Port Number | 16-bit Destination Port Number | |
| 32-bit Sequence Number | | |
| 32-bit Acknowledgement Number | | |
| 4-bit Header Length | Reserved (6 bits) | U A P R S F / R C S S Y I / G K H T N N | 16-bit Window Size |
| 16-bit TCP Checksum | | 16-bit Urgent Pointer |
| Options (if any) | | |
| Data (if any) | | |

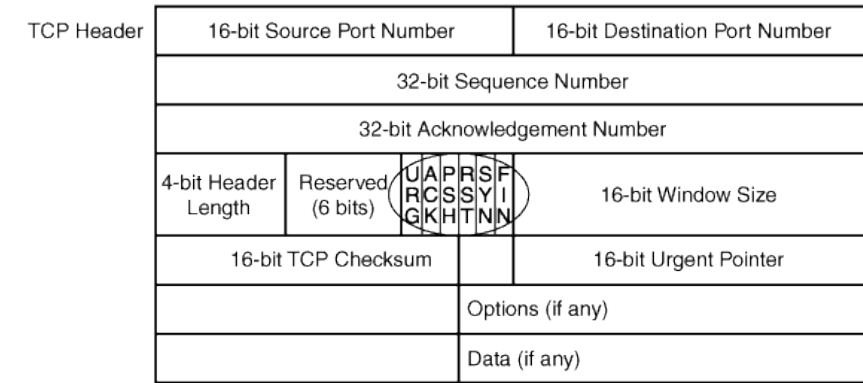The SYN and FIN flags are set.

g041323

國立中正大學
Cybersecurity Lab

# TCP Header

- The **URG** flag signifies that this segment contains urgent data. When this flag is set, the UrgPtr field indicates where the nonurgent data contained in this segment begins.

-

- The urgent data is contained at the front of the segment body, up to and including a value of **UrgPtr** bytes into the segment.

- The **PUSH** flag signifies that the sender invoked the push operation, which indicates to the receiving side of TCP that it should notify the receiving process of this fact.

- Finally, the **RESET** flag signifies that the receiver has become confused

```
▼ Flow 1
    SrcAddr: 192.168.88.102
    DstAddr: 188.188.188.188
    Protocol: TCP (6)
    SrcPort: 27127 (27127)
    DstPort: 23 (23)
  ▼ TCP Flags: 0x12, ACK, SYN
        00.. .... = Reserved: 0x0
        ..0. .... = URG: Not used
        ...1 .... = ACK: Used
        .... 0... = PSH: Not used
        .... .0.. = RST: Not used
        .... ..1. = SYN: Used
        .... ...0 = FIN: Not used
    Source Mac Address: 0c:ad:95:fb:00:00
    Destination Mac Address: 0c:27:76:de:00:00
    InputInt: 1
    Classification Engine ID: PANA-L7 (13)
    Selector ID: 000001
    OutputInt: 0
    Direction: Ingress (0)
    Octets: 84
    Packets: 2
```

國立中正大學
Cybersecurity Lab

# TCP Header

- Finally, the RESET flag signifies that the receiver has become confused, it received a segment it did not expect to receive—and so wants to abort the connection.

- Finally, the Checksum field is used in exactly the same way as for UDP—it is computed over the TCP header, the TCP data, and the pseudo header, which is made up of the source address, destination address, and length fields from the IP header → Validate the packet integrity
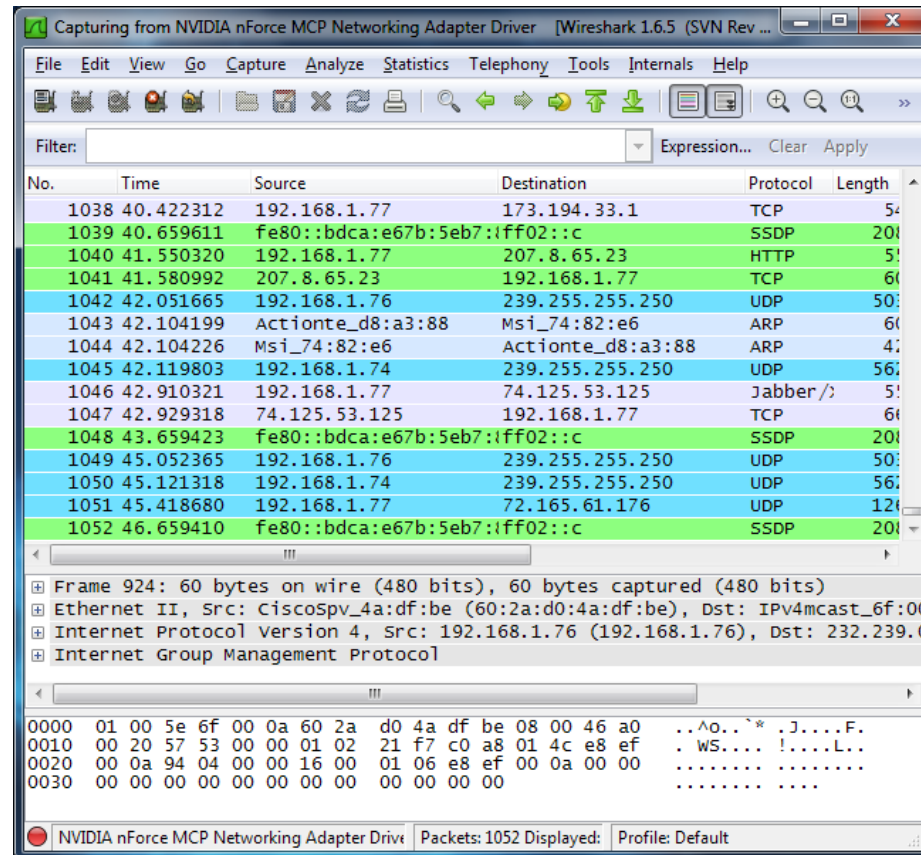


TCP Header

| 16-bit Source Port Number | 16-bit Destination Port Number |
|---|---|
| 32-bit Sequence Number | |
| 32-bit Acknowledgement Number | |
| 4-bit Header Length | Reserved (6 bits) | U R G | A C K | P S H | R S T | S Y N | F I N | 16-bit Window Size |
| 16-bit TCP Checksum | 16-bit Urgent Pointer |
| Options (if any) | |
| Data (if any) | |

The SYN and FIN flags are set.

g041323

國立中正大學
Cybersecurity Lab
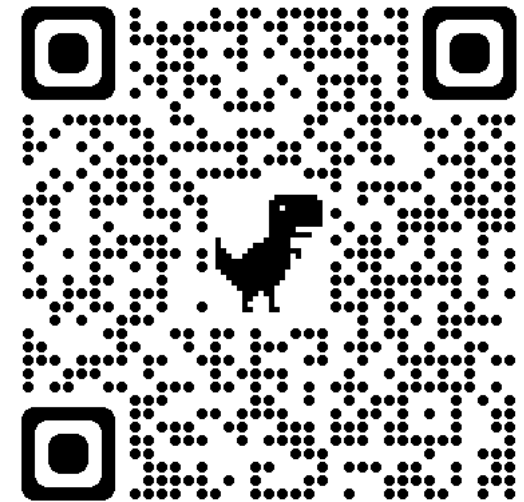
# Packet Analysis in Wireshark

- A free and open-source packet analyzer for network troubleshooting, analysis, software and communications protocol development, and education.



https://www.wireshark.org/



國立中正大學
Cybersecurity Lab

# Wireshark

- Allow us to check TCP/UDP packet in details

- Statistical data
  - ✓ Packet flow (src→ dest)
  - ✓ Endpoints
  - ✓ Packet lengths
  - ✓ I/O graphs



Packet flows

Packet details

# Wireshark

- Allow us to check TCP/UDP packet in details

- Statistic feature
  - ✓ Packet flow (src→ dest)
  - ✓ Endpoints
  - ✓ Packet lengths
  - ✓ I/O graphs

# Wireshark

- Allow us to check TCP/UDP packet in details

- Statistic feature
  - ✓ Packet flow (src→ dest)
  - ✓ Endpoints
  - ✓ Packet lengths
  - ✓ I/O graphs

Wireshark · Endpoints · Ethernet

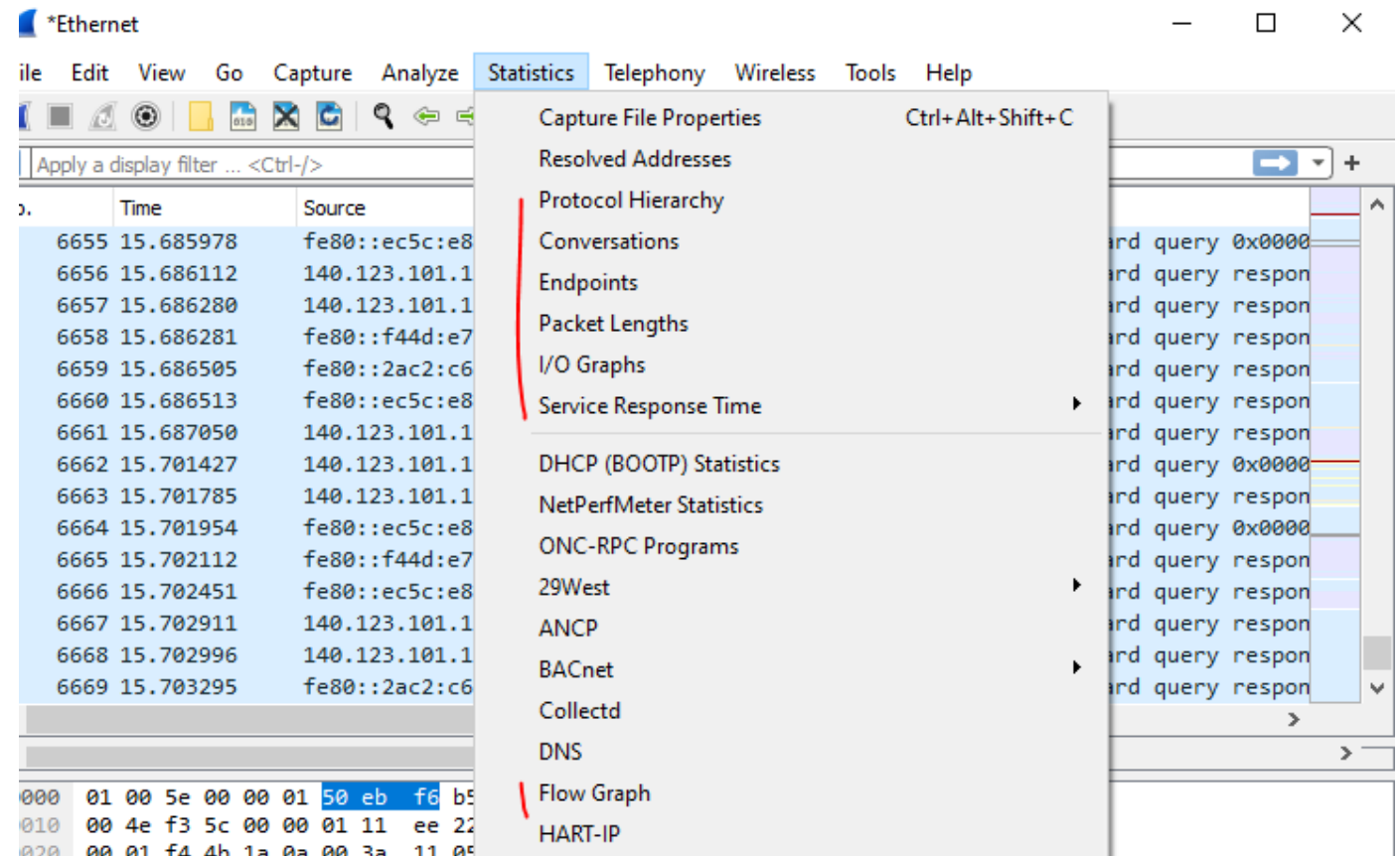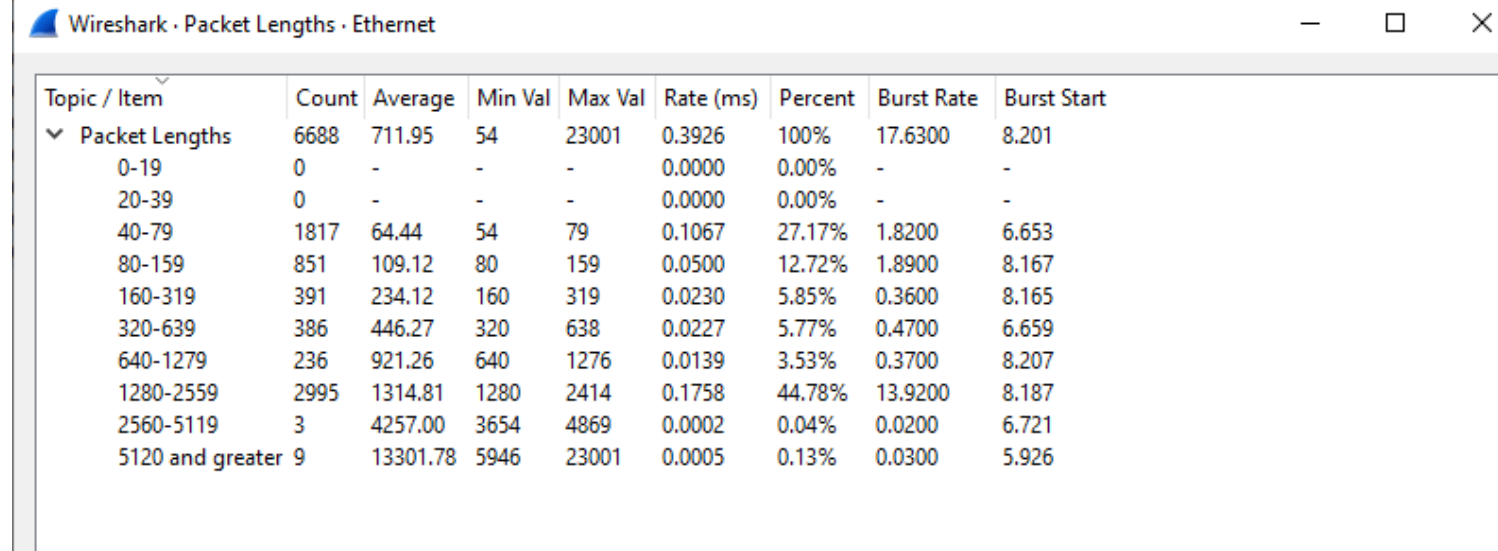| Ethernet · 43 | IPv4 · 92 | IPv6 · 13 | TCP · 167 | UDP · 203 | | |
|---|---|---|---|---|---|---|
| Address | Packets | Bytes | Tx Packets | Tx Bytes | Rx Packets | Rx Bytes |
| 20.78.118.165 | 2 | 114 | 1 | 60 | 1 | 54 |
| 20.89.149.168 | 5 | 328 | 3 | 220 | 2 | 108 |
| 20.108.229.21 | 2 | 115 | 1 | 60 | 1 | 55 |
| 20.110.103.72 | 1 | 60 | 1 | 60 | 0 | 0 |
| 20.126.223.223 | 1 | 60 | 1 | 60 | 0 | 0 |
| 20.189.173.1 | 1 | 60 | 1 | 60 | 0 | 0 |
| 20.189.173.23 | 1 | 60 | 1 | 60 | 0 | 0 |
| 20.190.166.66 | 1 | 60 | 1 | 60 | 0 | 0 |
| 20.210.223.40 | 1 | 60 | 1 | 60 | 0 | 0 |
| 34.36.124.104 | 24 | 8670 | 12 | 6205 | 12 | 2465 |
| 34.64.34.68 | 6 | 402 | 3 | 216 | 3 | 186 |
| 34.200.122.61 | 27 | 13 k | 15 | 10 k | 12 | 3521 |
| 35.168.49.207 | 31 | 16 k | 19 | 8857 | 12 | 7774 |
| 44.218.58.190 | 51 | 20 k | 27 | 14 k | 24 | 5603 |
| 51.104.15.253 | 2 | 120 | 2 | 120 | 0 | 0 |
| 52.72.226.68 | 60 | 20 k | 28 | 11 k | 32 | 9189 |
| 52.86.181.185 | 2 | 114 | 1 | 60 | 1 | 54 |
| 52.92.243.160 | 48 | 24 k | 28 | 22 k | 20 | 2374 |
| 52.111.234.0 | 3 | 260 | 2 | 160 | 1 | 100 |
| 52.113.194.132 | 3 | 180 | 3 | 180 | 0 | 0 |

# Wireshark

- Allow us to check TCP/UDP packet in details

- Statistic feature
  - ✓ Packet flow (src→ dest)
  - ✓ Endpoints
  - ✓ Packet lengths
  - ✓ I/O graphs



Wireshark · Packet Lengths · Ethernet

| Topic / Item | Count | Average | Min Val | Max Val | Rate (ms) | Percent | Burst Rate | Burst Start |
|---|---|---|---|---|---|---|---|---|
| ⌄ Packet Lengths | 6688 | 711.95 | 54 | 23001 | 0.3926 | 100% | 17.6300 | 8.201 |
| 0-19 | 0 | - | - | - | 0.0000 | 0.00% | - | - |
| 20-39 | 0 | - | - | - | 0.0000 | 0.00% | - | - |
| 40-79 | 1817 | 64.44 | 54 | 79 | 0.1067 | 27.17% | 1.8200 | 6.653 |
| 80-159 | 851 | 109.12 | 80 | 159 | 0.0500 | 12.72% | 1.8900 | 8.167 |
| 160-319 | 391 | 234.12 | 160 | 319 | 0.0230 | 5.85% | 0.3600 | 8.165 |
| 320-639 | 386 | 446.27 | 320 | 638 | 0.0227 | 5.77% | 0.4700 | 6.659 |
| 640-1279 | 236 | 921.26 | 640 | 1276 | 0.0139 | 3.53% | 0.3700 | 8.207 |
| 1280-2559 | 2995 | 1314.81 | 1280 | 2414 | 0.1758 | 44.78% | 13.9200 | 8.187 |
| 2560-5119 | 3 | 4257.00 | 3654 | 4869 | 0.0002 | 0.04% | 0.0200 | 6.721 |
| 5120 and greater | 9 | 13301.78 | 5946 | 23001 | 0.0005 | 0.13% | 0.0300 | 5.926 |

# Wireshark

- Allow us to check TCP/UDP packet in details

- Statistic data
  - ✓ Packet flow (src→ dest)
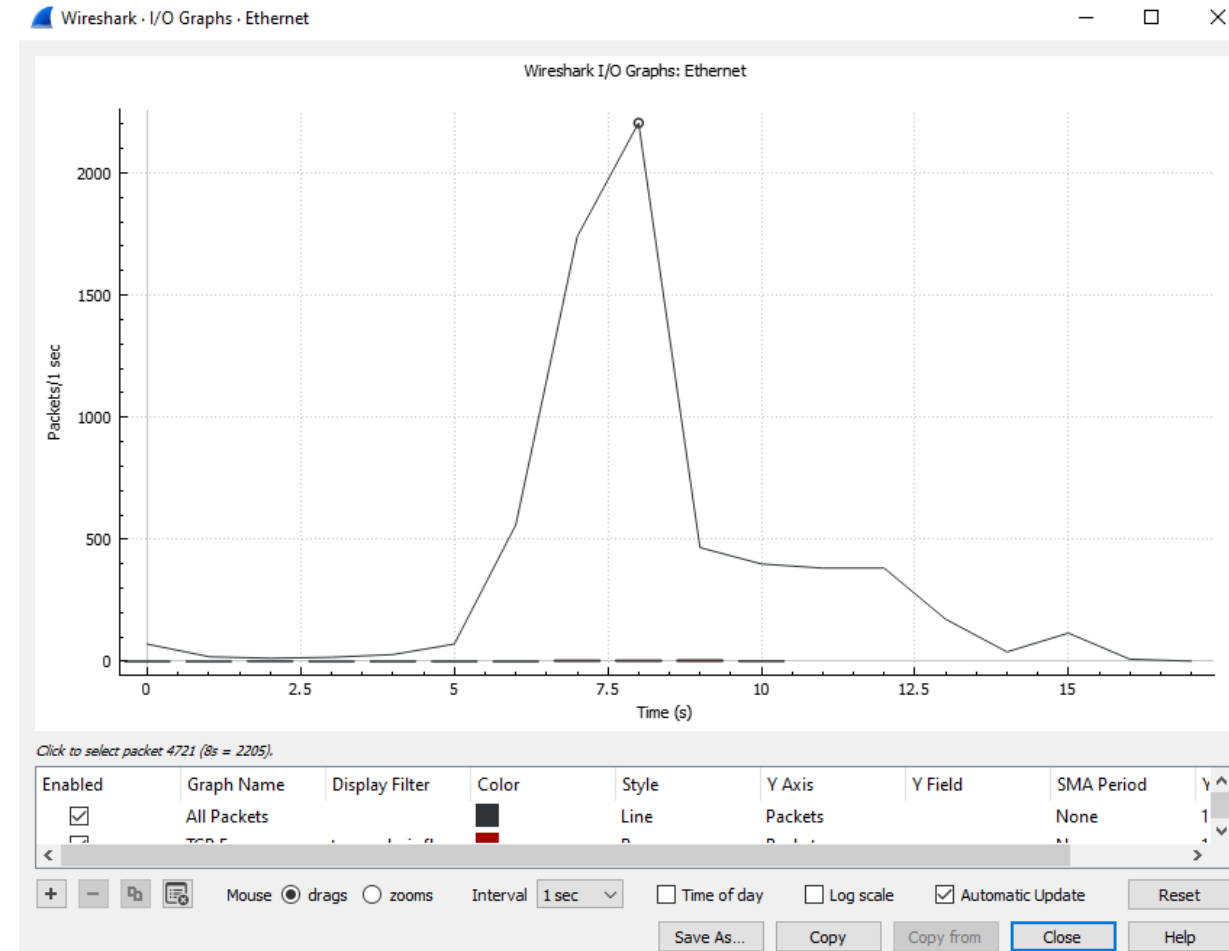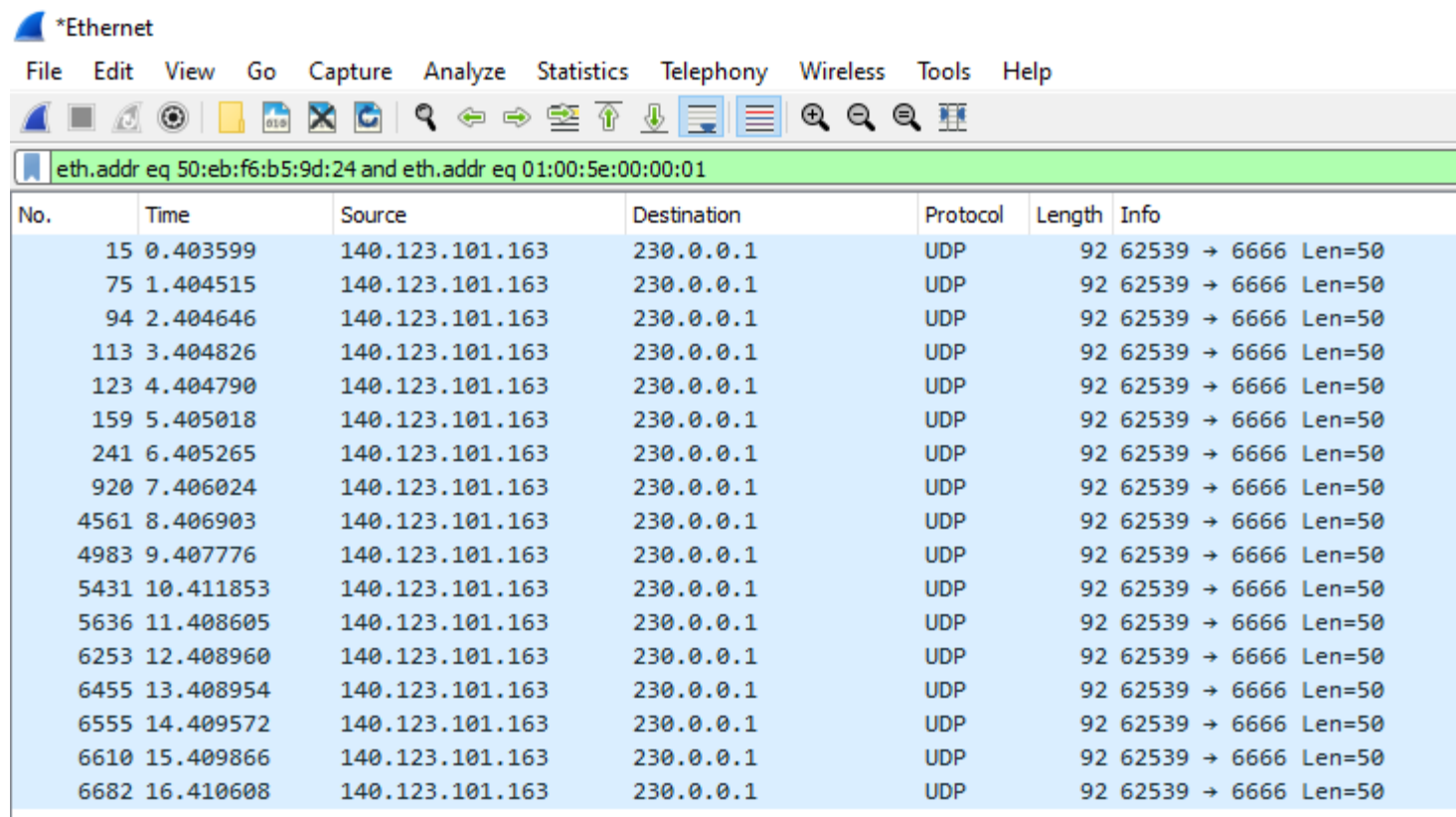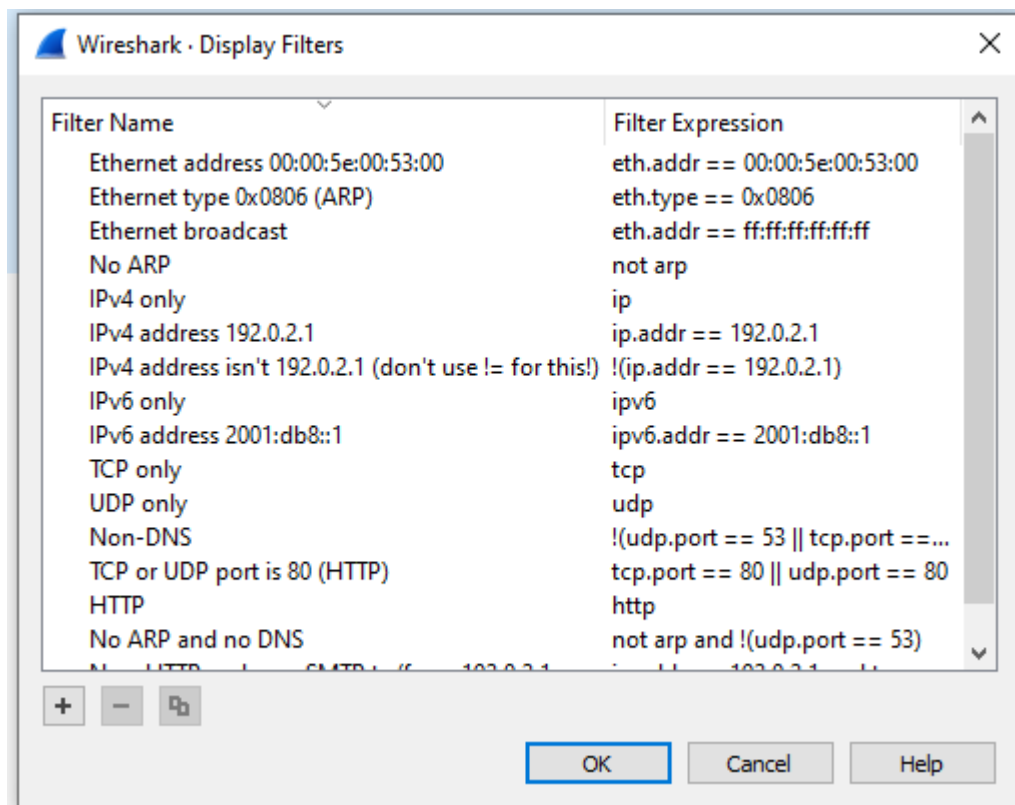  - ✓ Endpoints
  - ✓ Packet lengths
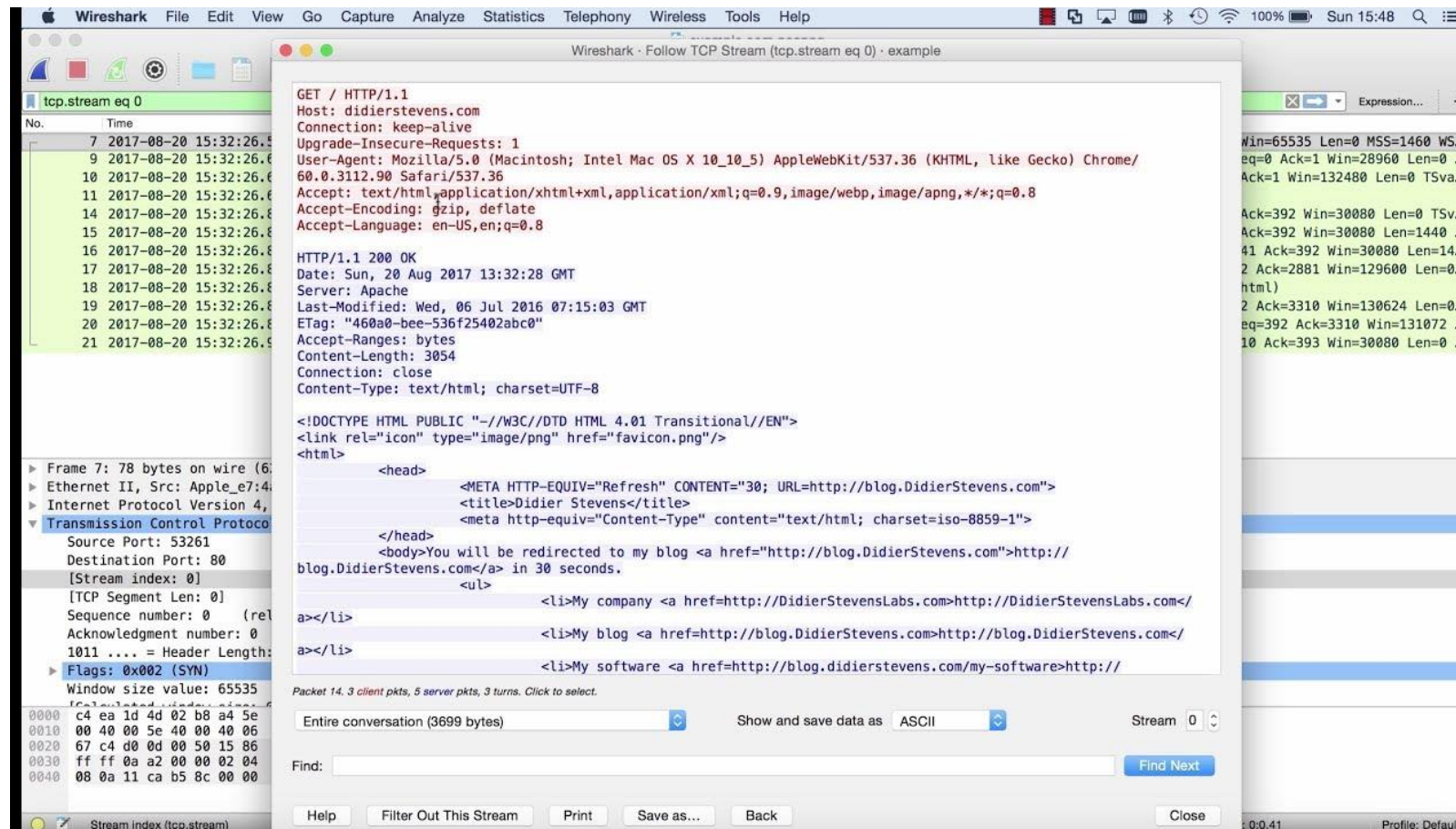  - ✓ I/O graphs



I/O graphs

# Packet flow filter in Wireshark

• Find all packets of a flow

**Wireshark · Display Filters**

| Filter Name | Filter Expression |
|---|---|
| Ethernet address 00:00:5e:00:53:00 | eth.addr == 00:00:5e:00:53:00 |
| Ethernet type 0x0806 (ARP) | eth.type == 0x0806 |
| Ethernet broadcast | eth.addr == ff:ff:ff:ff:ff:ff |
| No ARP | not arp |
| IPv4 only | ip |
| IPv4 address 192.0.2.1 | ip.addr == 192.0.2.1 |
| IPv4 address isn't 192.0.2.1 (don't use != for this!) | !(ip.addr == 192.0.2.1) |
| IPv6 only | ipv6 |
| IPv6 address 2001:db8::1 | ipv6.addr == 2001:db8::1 |
| TCP only | tcp |
| UDP only | udp |
| Non-DNS | !(udp.port == 53 \|\| tcp.port ==... |
| TCP or UDP port is 80 (HTTP) | tcp.port == 80 \|\| udp.port == 80 |
| HTTP | http |
| No ARP and no DNS | not arp and !(udp.port == 53) |

[ + ] [ − ] [ ⧉ ]

[ OK ] [ Cancel ] [ Help ]

**\*Ethernet**

File   Edit   View   Go   Capture   Analyze   Statistics   Telephony   Wireless   Tools   Help

`eth.addr eq 50:eb:f6:b5:9d:24 and eth.addr eq 01:00:5e:00:00:01`

| No. | Time | Source | Destination | Protocol | Length | Info |
|---|---|---|---|---|---|---|
| 15 | 0.403599 | 140.123.101.163 | 230.0.0.1 | UDP | 92 | 62539 → 6666 Len=50 |
| 75 | 1.404515 | 140.123.101.163 | 230.0.0.1 | UDP | 92 | 62539 → 6666 Len=50 |
| 94 | 2.404646 | 140.123.101.163 | 230.0.0.1 | UDP | 92 | 62539 → 6666 Len=50 |
| 113 | 3.404826 | 140.123.101.163 | 230.0.0.1 | UDP | 92 | 62539 → 6666 Len=50 |
| 123 | 4.404790 | 140.123.101.163 | 230.0.0.1 | UDP | 92 | 62539 → 6666 Len=50 |
| 159 | 5.405018 | 140.123.101.163 | 230.0.0.1 | UDP | 92 | 62539 → 6666 Len=50 |
| 241 | 6.405265 | 140.123.101.163 | 230.0.0.1 | UDP | 92 | 62539 → 6666 Len=50 |
| 920 | 7.406024 | 140.123.101.163 | 230.0.0.1 | UDP | 92 | 62539 → 6666 Len=50 |
| 4561 | 8.406903 | 140.123.101.163 | 230.0.0.1 | UDP | 92 | 62539 → 6666 Len=50 |
| 4983 | 9.407776 | 140.123.101.163 | 230.0.0.1 | UDP | 92 | 62539 → 6666 Len=50 |
| 5431 | 10.411853 | 140.123.101.163 | 230.0.0.1 | UDP | 92 | 62539 → 6666 Len=50 |
| 5636 | 11.408605 | 140.123.101.163 | 230.0.0.1 | UDP | 92 | 62539 → 6666 Len=50 |
| 6253 | 12.408960 | 140.123.101.163 | 230.0.0.1 | UDP | 92 | 62539 → 6666 Len=50 |
| 6455 | 13.408954 | 140.123.101.163 | 230.0.0.1 | UDP | 92 | 62539 → 6666 Len=50 |
| 6555 | 14.409572 | 140.123.101.163 | 230.0.0.1 | UDP | 92 | 62539 → 6666 Len=50 |
| 6610 | 15.409866 | 140.123.101.163 | 230.0.0.1 | UDP | 92 | 62539 → 6666 Len=50 |
| 6682 | 16.410608 | 140.123.101.163 | 230.0.0.1 | UDP | 92 | 62539 → 6666 Len=50 |

國立中正大學
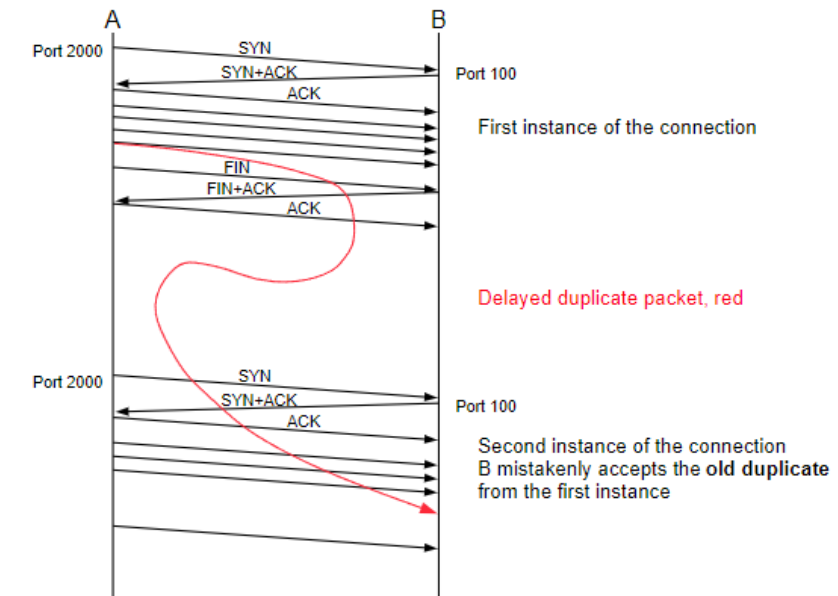Cybersecurity Lab

# TCP stream follow in Wireshark

- Read all packets of http requests

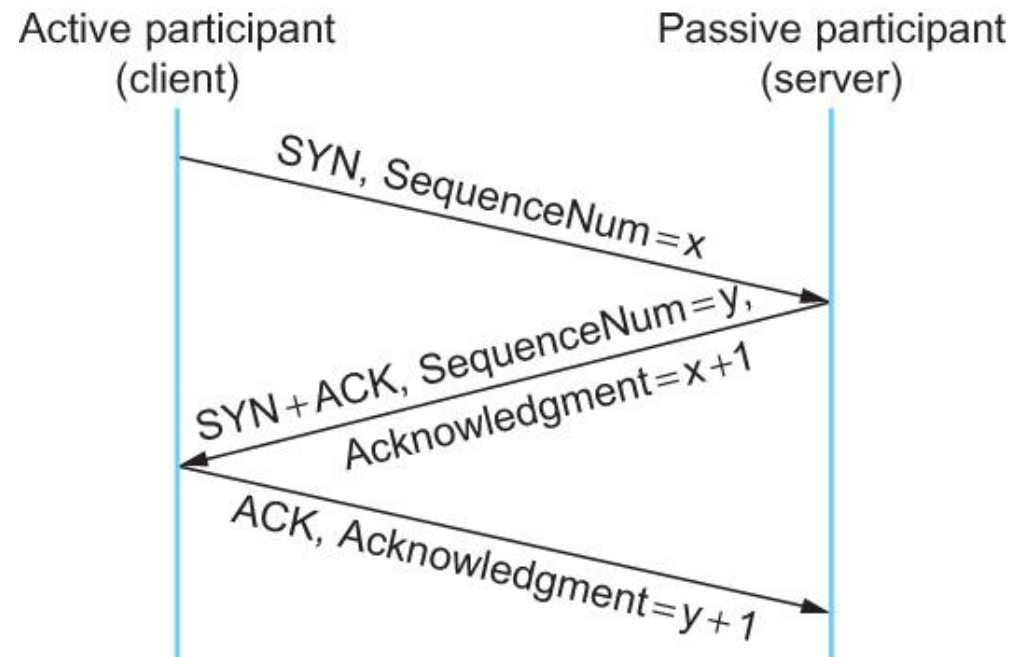# Network issues

# TCP Issues

- At the heart of TCP is the sliding window algorithm

- As TCP runs over the Internet rather than <span style="color:red">a point-to-point</span> link, the following issues need to be addressed by the sliding window algorithm

  - TCP supports logical connections between processes that are running on two different computers in the Internet
  - TCP connections are likely to have widely different Round Trip (RTT) times
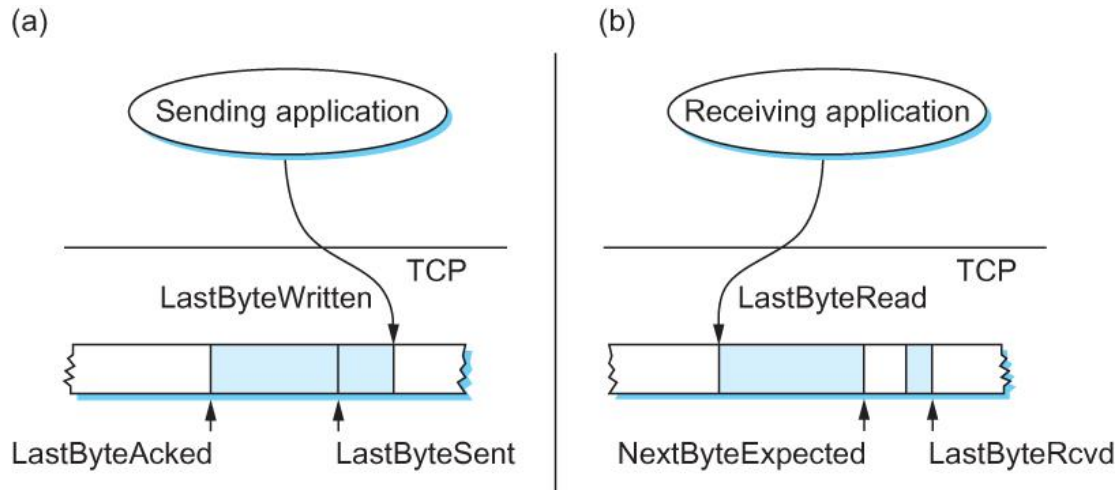  - Packets may get reordered in the Internet

# Connection Establishment/Termination in TCP



Timeline for three-way handshake algorithm

# Sliding Window Revisited

- TCP's variant of the sliding window algorithm, which serves several purposes:
  - (1) it guarantees the reliable delivery of data,
  - (2) it ensures that data is delivered in order, and
  - (3) it enforces flow control between the sender and the receiver.



- Sending Side
  LastByteAcked ≤ LastByteSent
  LastByteSent ≤ LastByteWritten
- Receiving Side
  LastByteRead < NextByteExpected
  NextByteExpected ≤ LastByteRcvd + 1

Relationship between TCP send buffer (a) and receive buffer (b).

# End-to-end Issues

- TCP is a byte-oriented protocol, which means that the sender writes bytes into a TCP connection and the receiver reads bytes out of the TCP connection.

- Although "byte stream" describes the service TCP offers to application processes, TCP does not, itself, transmit individual bytes over the Internet.

- TCP needs a mechanism using which each side of a connection will learn what resources the other side is able to apply to the connection → TCP flow/congestion control

- TCP needs a mechanism using which the sending side will learn the capacity of the network

# Flow control VS Congestion control

- **Flow control** involves preventing senders from overrunning the capacity of the receivers

- **Congestion control** involves preventing too much data from being injected into the network, thereby causing switches or links to become overloaded *bandwidth limited*

- **Flow Control:** Algorithms to prevent that the sender overruns the receiver with information

- **Error Control:** Algorithms to recover or conceal the effects from packet losses

- **Congestion Control:** Algorithms to prevent that the sender overloads the network

→ The goal of each of the control mechanisms are different.

→ In TCP, the implementation of these algorithms is combined

Opposite objectives
- End-system
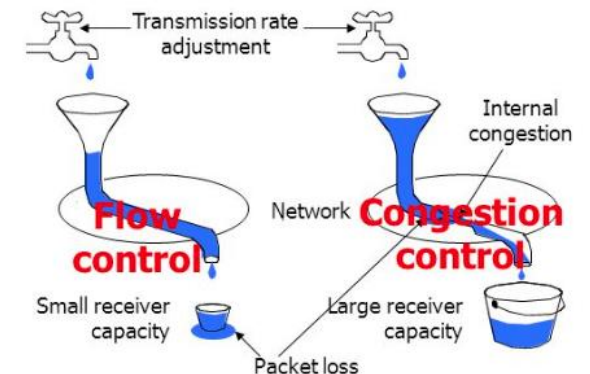  - Optimize its own throughput
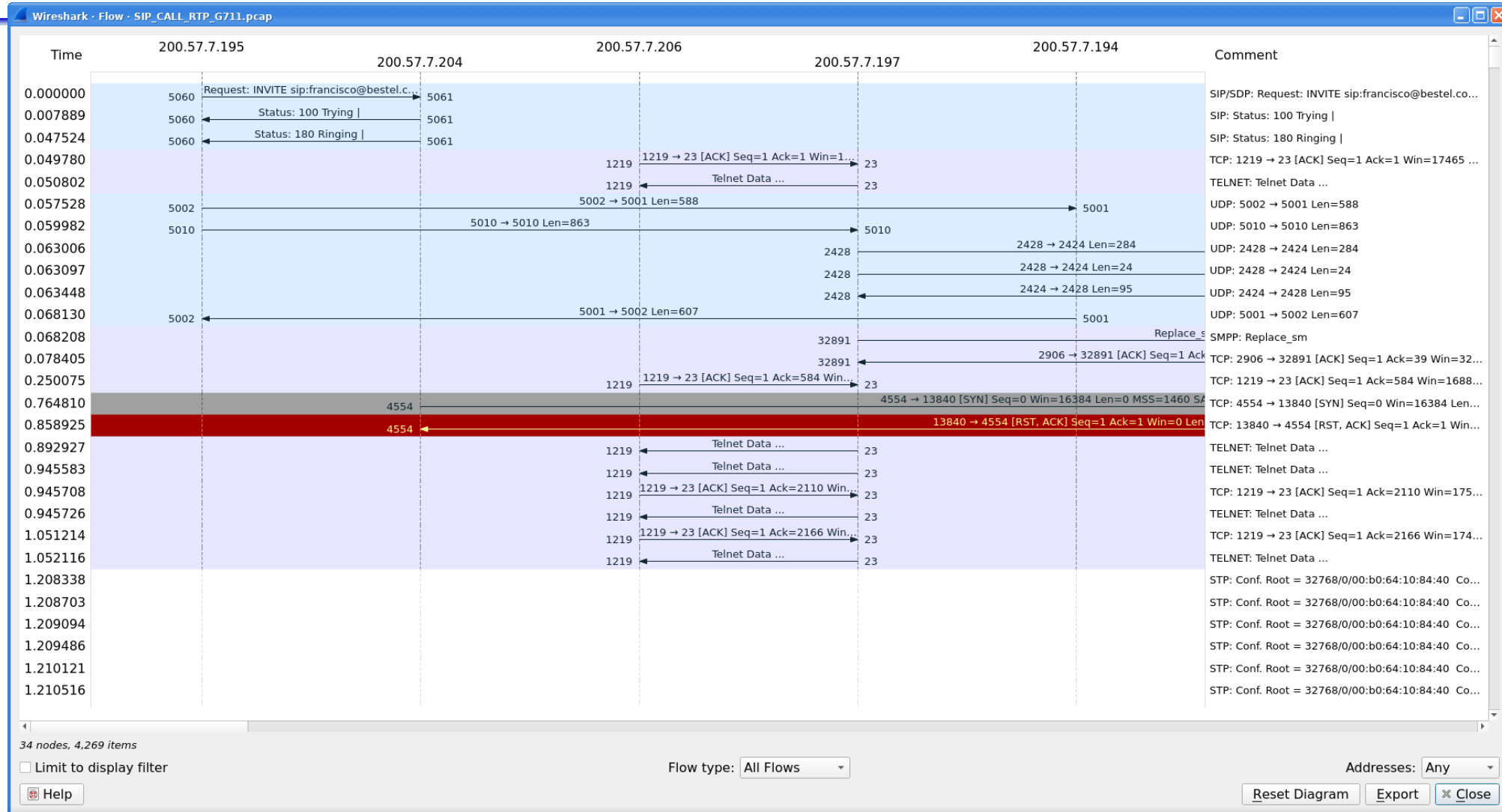  - Possibly at the expense of other end-systems

- Two different problems
  - Receiver capacity
  - Network capacity
  - Cannot be distinguished easily at all places
  - Should be differentiated

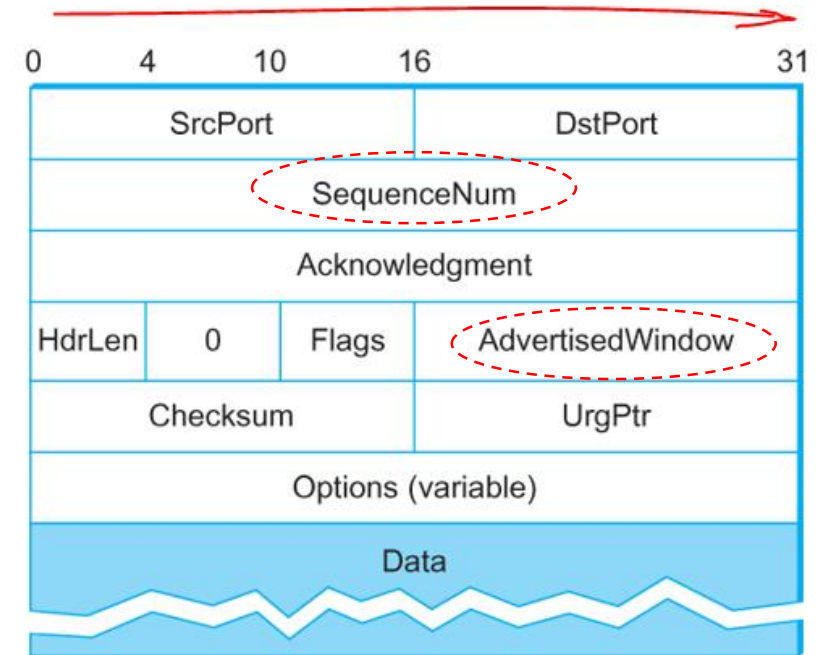Opposite objectives
- Network
  - Optimize overall throughput

Transmission rate adjustment

**Flow control**    Network  **Congestion control**

Internal congestion

Small receiver capacity    Large receiver capacity

Packet loss

# TCP Flow Control

- LastByteRcvd − LastByteRead ≤ MaxRcvBuffer

- AdvertisedWindow = MaxRcvBuffer − ((NextByteExpected − 1) − LastByteRead)

- LastByteSent − LastByteAcked ≤ AdvertisedWindow

- EffectiveWindow = AdvertisedWindow − (LastByteSent − LastByteAcked)

- LastByteWritten − LastByteAcked ≤ MaxSendBuffer

- If the sending process tries to write y bytes to TCP, but

  (LastByteWritten − LastByteAcked) + y > MaxSendBuffer

  then TCP blocks the sending process and does not allow it to generate more data.

# TCP flow control in Wireshark

# Consequence number's limit: Wraparound

- **SequenceNum**: 32 bits longs     $2^{32} -1 = 4GB$

- Every packet's sequence number must be unique → the Sequence Number will be exhausted ( > 4GB).

- When accessible, the sequence numbers that were previously utilized can be reused as needed

- Reusing of sequence numbers is known as the Wraparound

- AdvertisedWindow: 16 bits long

  - TCP has satisfied the requirement of the sliding

  - window algorithm that is the sequence number

  - space be twice as big as the window size

  - $2^{32} >> 2 \times 2^{16}$



TCP Header Format

# Protecting against Wraparound

- Relevance of the 32-bit sequence number space

  - The sequence number used on a given connection might wraparound

  - A byte with sequence number $x$ could be sent at one time, and then at a later time a second byte with the same sequence number $x$ could be sent

  - Packets cannot survive in the Internet for longer than the **MSS**

  - **MSS** is set to 120 sec

  - We need to make sure that the sequence number does not wrap around within a 120-second period of time

  - Depends on how fast data can be transmitted over the Internet

| Bandwidth | Time until Wraparound |
|---|---|
| T1 (1.5 Mbps) | 6.4 hours |
| Ethernet (10 Mbps) | 57 minutes |
| T3 (45 Mbps) | 13 minutes |
| Fast Ethernet (100 Mbps) | 6 minutes |
| OC-3 (155 Mbps) | 4 minutes |
| OC-12 (622 Mbps) | 55 seconds |
| OC-48 (2.5 Gbps) | 14 seconds |

Time until 32-bit sequence number space wraps around.

# Keeping the Pipe Full

- 16-bit AdvertisedWindow field must be big enough to allow the sender to keep the pipe full

- Clearly the receiver is <span style="color:red">free not to open the window as large as</span> the AdvertisedWindow field allows

- If the receiver has enough buffer space
  - The window needs to be opened far enough to allow a full
  - Delay × bandwidth product's worth of data
  - Assuming an RTT of 100 ms

| Bandwidth | Delay × Bandwidth Product |
|---|---|
| T1 (1.5 Mbps) | 18 KB |
| Ethernet (10 Mbps) | 122 KB |
| T3 (45 Mbps) | 549 KB |
| Fast Ethernet (100 Mbps) | 1.2 MB |
| OC-3 (155 Mbps) | 1.8 MB |
| OC-12 (622 Mbps) | 7.4 MB |
| OC-48 (2.5 Gbps) | 29.6 MB |

Required window size for 100-ms RTT.

國立中正大學
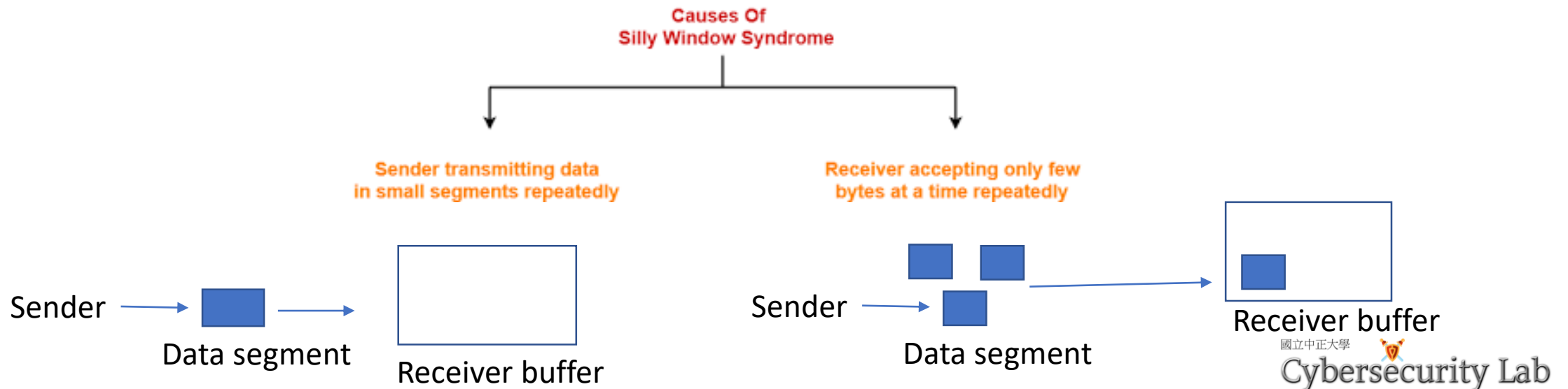Cybersecurity Lab
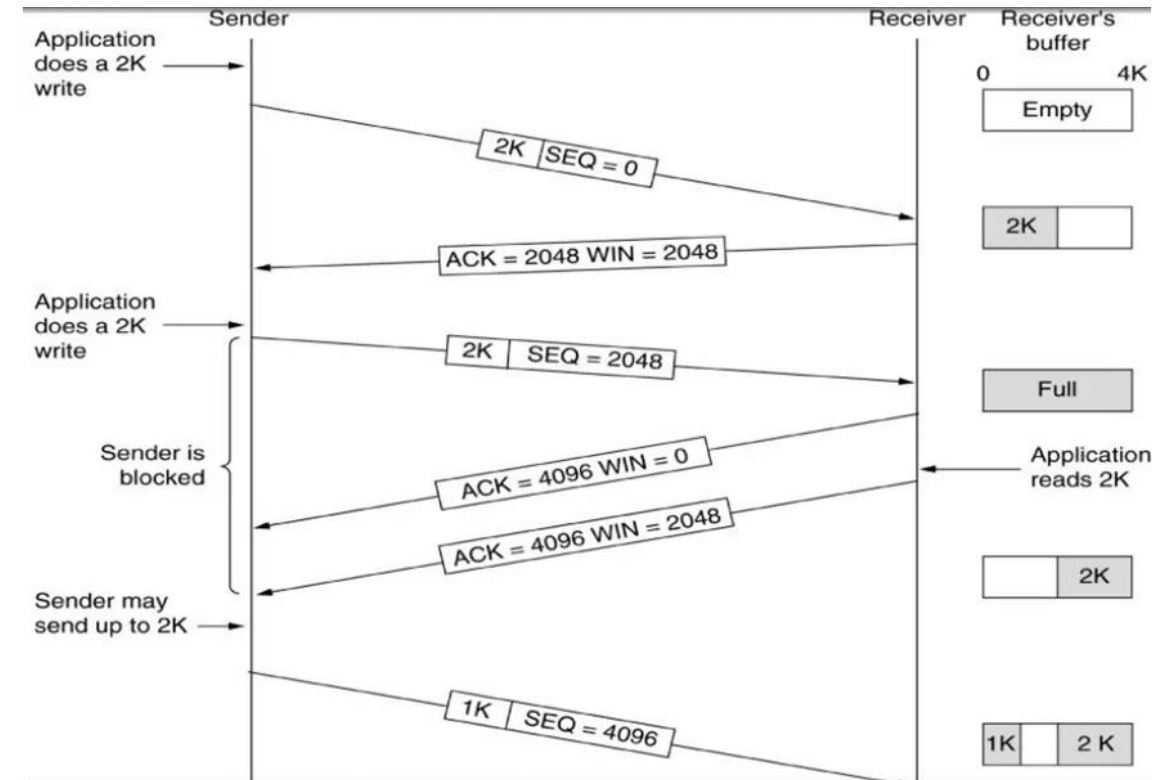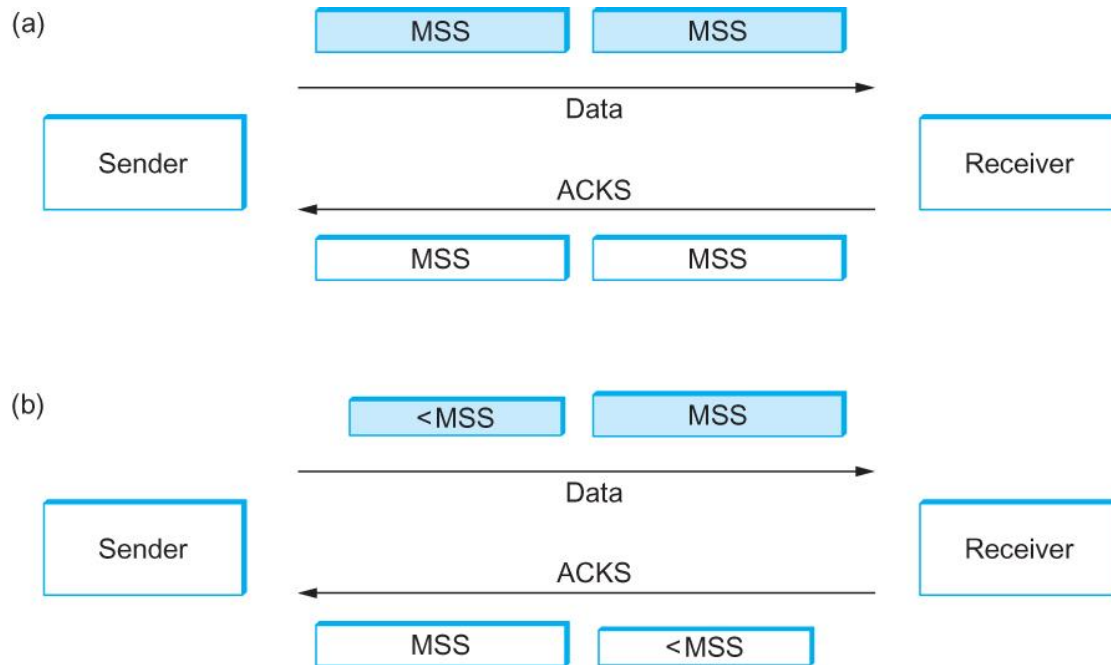
# Triggering Transmission

- How does TCP decide to transmit a segment?
  - TCP supports a byte stream abstraction
  - Application programs write bytes into streams
  - It is up to TCP to decide that it has enough bytes to send a segment

- What factors governs this decision
  - Ignore flow control: window is wide open, as would be the case when the connection starts
  - TCP has three mechanism to trigger the transmission of a segment
    - 1) TCP maintains a variable MSS and sends a segment as soon as it has collected MSS bytes from the sending process
      - MSS is usually set to the size of the largest segment TCP can send without causing local IP to fragment.
      - MSS: MTU of directly connected network – (TCP header + and IP header)
    - 2) Sending process has explicitly asked TCP to send it
      - TCP supports push operation
    - 3) When a timer fires
      - Resulting segment contains as many bytes as are currently buffered for transmission

# Silly Window Syndrome

- A problem in computer networking caused by poorly implemented TCP flow control

- This problem occurs when the sending application program creates data slowly, the receiving application program consumes data slowly, or both

**Causes Of Silly Window Syndrome**

**Sender transmitting data in small segments repeatedly**

**Receiver accepting only few bytes at a time repeatedly**

Sender → Data segment → Receiver buffer

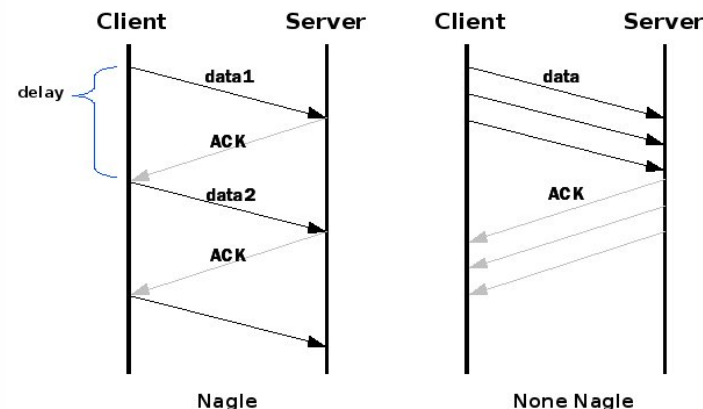Sender → Data segment → Receiver buffer

# Silly Window Syndrome



Silly Window Syndrome

Source: Gateoverflow

# Nagle's Algorithm

- If there is data to send but the window is open less than MSS, then we may want to wait some amount of time before sending the available data

- But how long?

- If we wait too long, then we hurt interactive applications like Telnet

- If we don't wait long enough, then we risk sending a bunch of tiny packets and falling into the *silly window* syndrome
  - The solution is to introduce a timer and to transmit when the timer expires

# Nagle's Algorithm

- We could use a clock-based timer, for example one that fires every 100 ms

- Nagle introduced an elegant self-clocking solution

- Key Idea
  - As long as TCP has any data in flight, the sender will eventually receive an ACK
  - This ACK can be treated like a timer firing, triggering the transmission of more data

# Nagle's Algorithm

When the application produces data to send

    if both the available data and the window ≥ MSS

        send a full segment

    else

        if there is unACKed data in flight

            buffer the new data until an ACK arrives

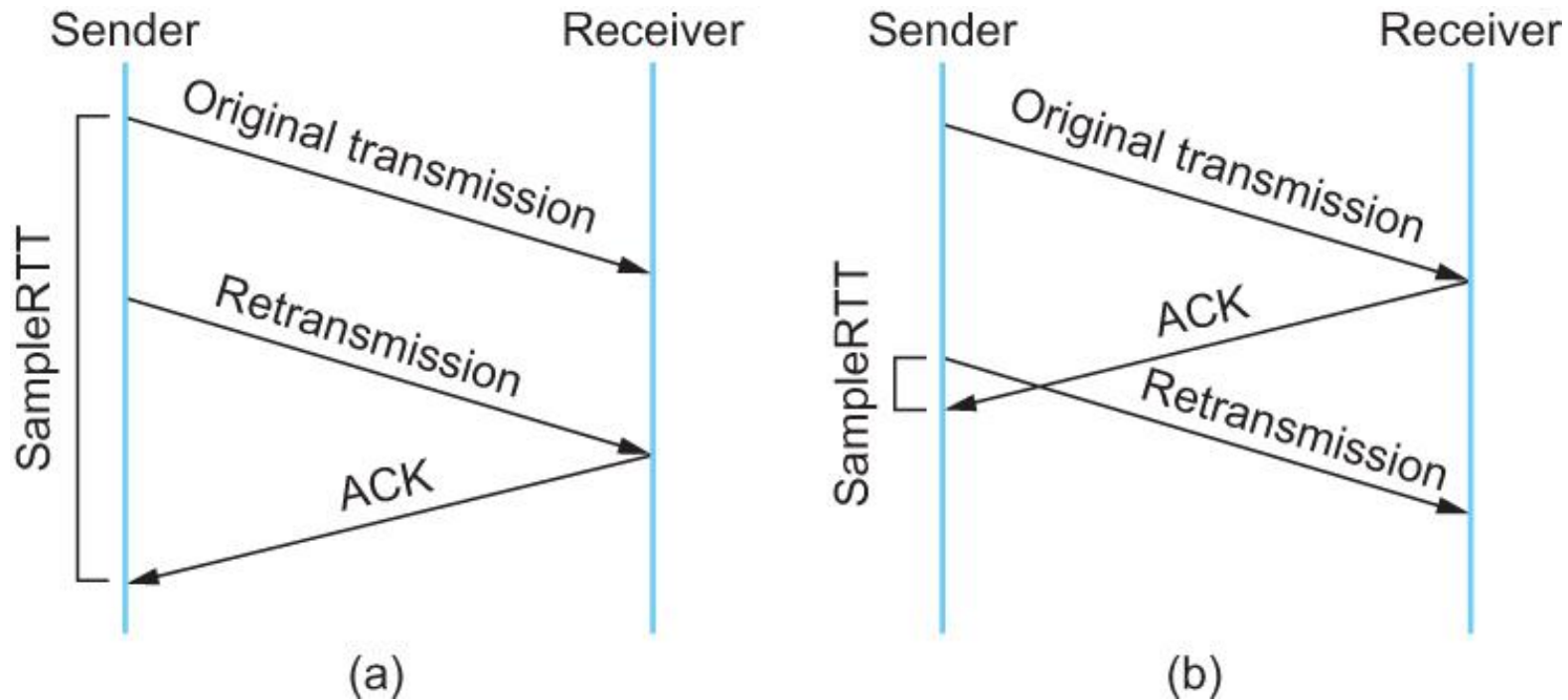        else

            send all the new data now

# Adaptive Retransmission

- Original Algorithm
  - Measure **SampleRTT** for each segment/ ACK pair
  - Compute weighted average of RTT
    - **EstRTT** = $\alpha$ x **EstRTT** + (1 - $\alpha$)x **SampleRTT**
    - $\alpha$ between 0.8 and 0.9
  - Set timeout based on **EstRTT**
    - **TimeOut** = 2 x **EstRTT**

# Original Algorithm

- Problem
  - ACK does not really acknowledge a transmission
    - It actually acknowledges the receipt of data
  - When a segment is retransmitted and then an ACK arrives at the sender
    - It is impossible to decide if this ACK should be associated with the first or the second transmission for calculating RTTs

# Karn/Partridge Algorithm



Associating the ACK with (a) original transmission versus (b) retransmission

# Karn/Partridge Algorithm

- Do not sample RTT when retransmitting

- Double timeout after each retransmission

- Karn-Partridge algorithm was an improvement over the original approach, but it does not eliminate congestion

- We need to understand how timeout is related to congestion
  - If you timeout too soon, you may unnecessarily retransmit a segment which adds load to the network

# Karn/Partridge Algorithm

- Main problem with the original computation is that it does not take variance of Sample RTTs into consideration.

- If the variance among Sample RTTs is small
  - Then the Estimated RTT can be better trusted
  - There is no need to multiply this by 2 to compute the timeout

- On the other hand, a large variance in the samples suggest that timeout value should not be tightly coupled to the Estimated RTT

- Jacobson/Karels proposed a new scheme for TCP retransmission

# Jacobson/Karels Algorithm

- Difference = SampleRTT − EstimatedRTT

- EstimatedRTT = EstimatedRTT + ( × Difference)

- Deviation = Deviation + (|Difference| − Deviation)

- TimeOut = μ × EstimatedRTT +  × Deviation
    - where based on experience, μ is typically set to 1 and  is set to 4. Thus, when the variance is small, TimeOut is close to EstimatedRTT; a large variance causes the deviation term to dominate the calculation.

# TCP congestion control

- The mechanism that <span style="color:red">prevents congestion</span> from happening or removes it after congestion takes place

- Congestion window state of TCP that <span style="color:red">limits the amount of data to be sent</span> by the sender into the network <span style="color:red">even before receiving the acknowledgment</span>
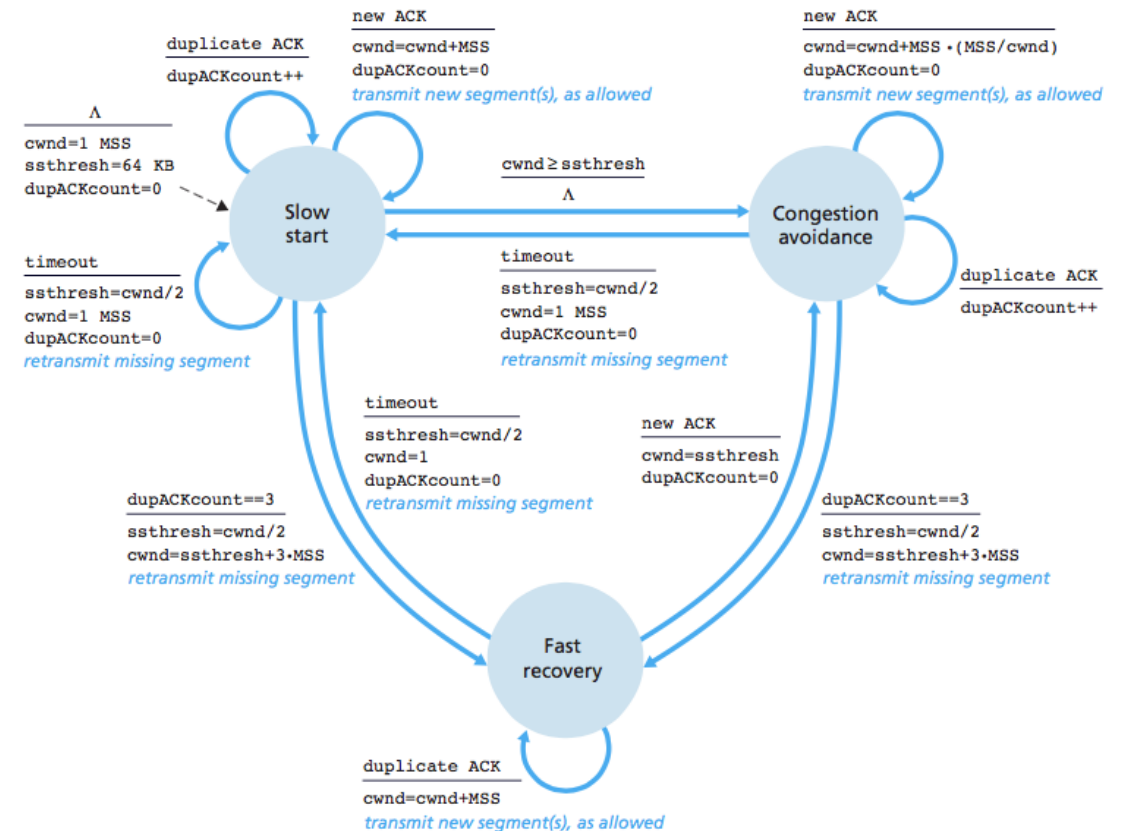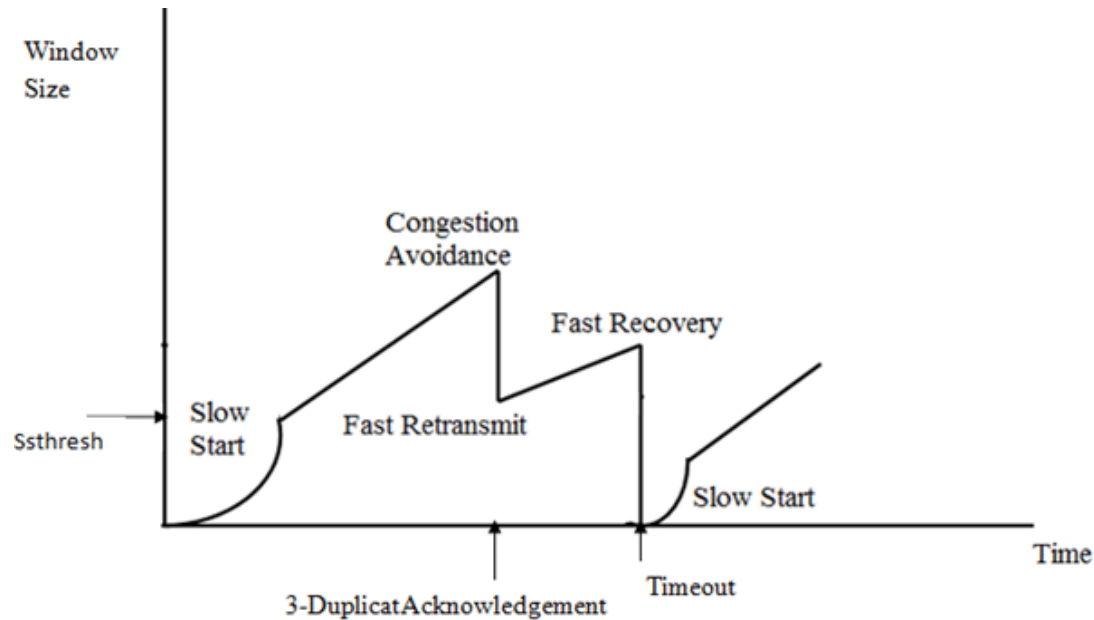
# TCP congestion control list

- The naming convention for congestion control algorithms (CCAs) may have originated in a 1996 paper by Kevin Fall and Sally Floyd

- There are many variants of congestion algorithms: TCP Reno/TCP Vegas….

- Efficient TCP congestion control is still a open problem for research

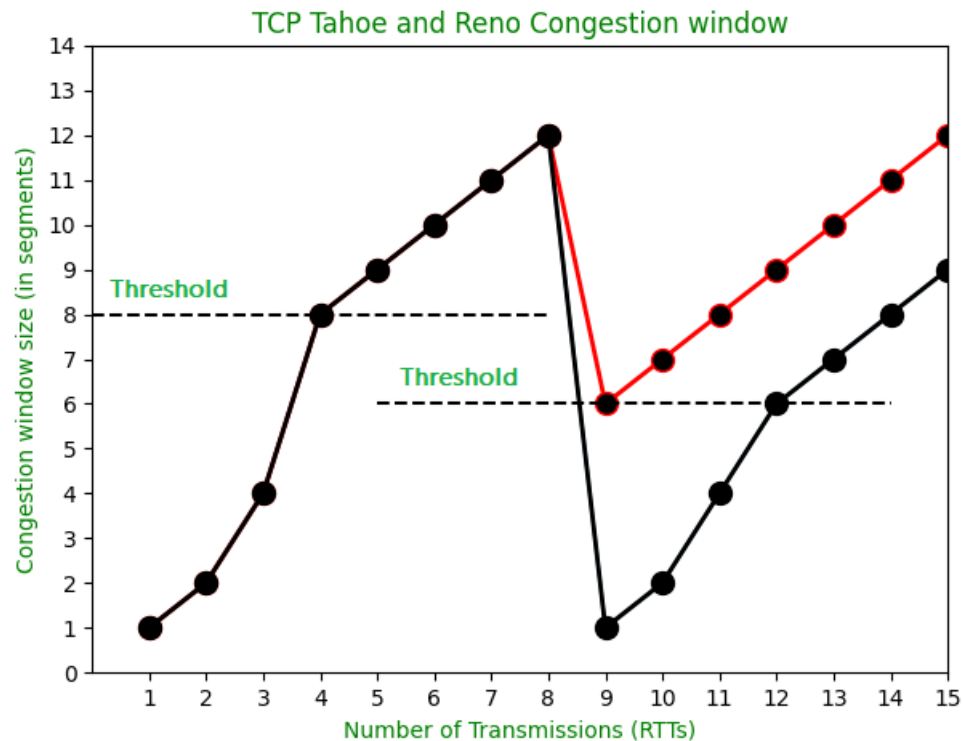| Variant | Feedback | Required changes | Benefits | Fairness |
|---|---|---|---|---|
| (New) Reno | Loss | — | — | Delay |
| Vegas | Delay | Sender | Less loss | Proportional |
| High Speed | Loss | Sender | High bandwidth | |
| BIC | Loss | Sender | High bandwidth | |
| CUBIC | Loss | Sender | High bandwidth | |
| C2TCP[11][12] | Loss/Delay | Sender | Ultra-low latency and high bandwidth | |
| NATCP[13] | Multi-bit signal | Sender | Near Optimal Performance | |
| Elastic-TCP | Loss/Delay | Sender | High bandwidth/short & long-distance | |
| Agile-TCP | Loss | Sender | High bandwidth/short-distance | |
| H-TCP | Loss | Sender | High bandwidth | |
| FAST | Delay | Sender | High bandwidth | Proportional |
| Compound TCP | Loss/Delay | Sender | High bandwidth | Proportional |
| Westwood | Loss/Delay | Sender | Lossy links | |
| Jersey | Loss/Delay | Sender | Lossy links | |
| BBR[14] | Delay | Sender | BLVC, Bufferbloat | |
| CLAMP | Multi-bit signal | Receiver, Router | Variable-rate links | Max-min |
| TFRC | Loss | Sender, Receiver | No Retransmission | Minimum delay |
| XCP | Multi-bit signal | Sender, Receiver, Router | BLFC | Max-min |
| VCP | 2-bit signal | Sender, Receiver, Router | BLF | Proportional |
| MaxNet | Multi-bit signal | Sender, Receiver, Router | BLFSC | Max-min |
| JetMax | Multi-bit signal | Sender, Receiver, Router | High bandwidth | Max-min |
| RED | Loss | Router | Reduced delay | |
| ECN | Single-bit signal | Sender, Receiver, Router | Reduced loss | |

# TCP congestion control

- Some famous congestion control mechanisms

# TCP congestion control

TCP Reno may be the best choice due to its ability to adapt quickly and avoid congestion

**TCP Tahoe and Reno Congestion window**

(Graph: Congestion window size (in segments) vs Number of Transmissions (RTTs), showing Threshold lines at 8 and 6)

## TCP Congestion Control

| Algoritms | condition | Design | action |
|---|---|---|---|
| Slow Start | cwnd <= ssthres; | cwnd doubles per RTT | cwnd+=1MSS per ACK |
| Congestion Avoidance | cwnd > ssthres | cwnd++ per RTT (additive increase) | cwnd+=1/cwnd * MSS per ACK |
| fast retx | 3 duplicate ACK | reduce the cwnd by half (multicative decreasing) | ssthres = max(cwnd/2,2) cwnd = ssthres + 3 MSS; retx the lost packet |
| fast recovery | receiving a new ACK after fast retx | finish the 1/2 reduction of cwnd in fast retx/fast recovery | cwnd = ssthres; tx if allowed by cwnd |
| | upon a dup ACK after fast retx before fast recovery | ("transition phrase") | cwnd +=1MSS; Note: it is different from slow start. |
| RTO timeout | time out | reset | ssthres = max(cwnd/2,2) cwnd = 1; retx the lost packet |

國立中正大學
Cybersecurity Lab

# Summary

- We have discussed how to transmit data via transport layer.

- We have discussed UDP

- We have discussed TCP/TCP flow control/congestion control

- We have discussed to use Wireshark to analyze packets

- We have discussed several algorithms to optimize TCP transmission