

Chapter 6

The Transport Layer

陳瑞奇(Rikki)

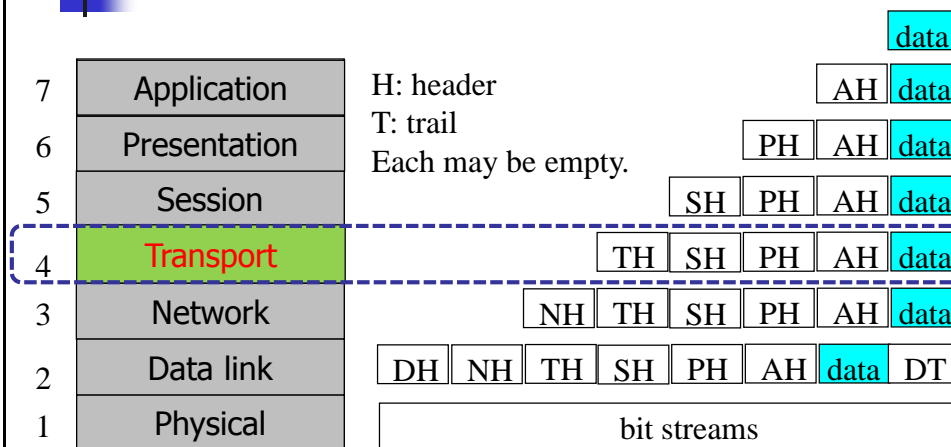
亞洲大學資訊工程學系

Adapted from Computer Networks,
Andrew S. Tanenbaum, Vrije University, Netherlands
& Computer Networking: A Top Down Approach,
Jim Kurose, Keith Ross

Computer Networks, Fifth Edition by Andrew Tanenbaum and David Wetherall, © Pearson Education-Prentice Hall, 2011

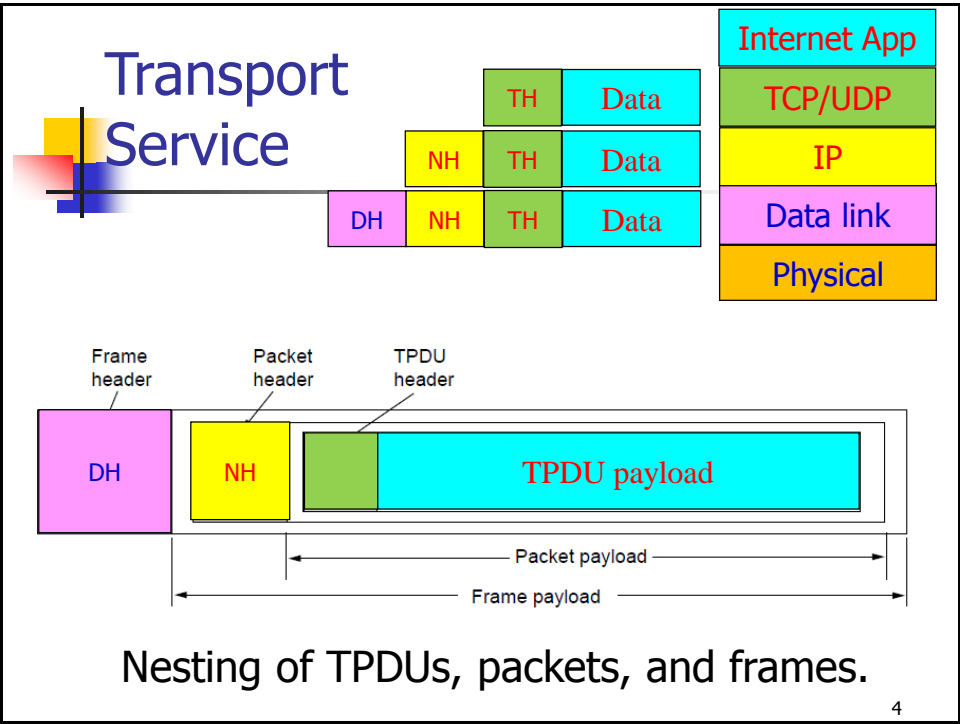
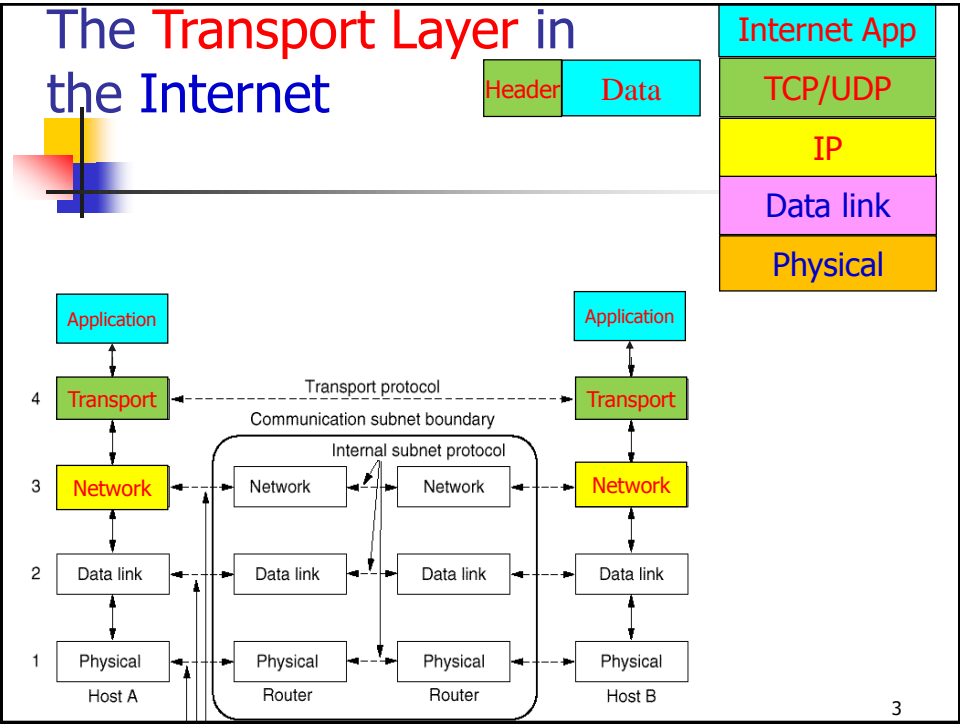
1

The Transport Layer



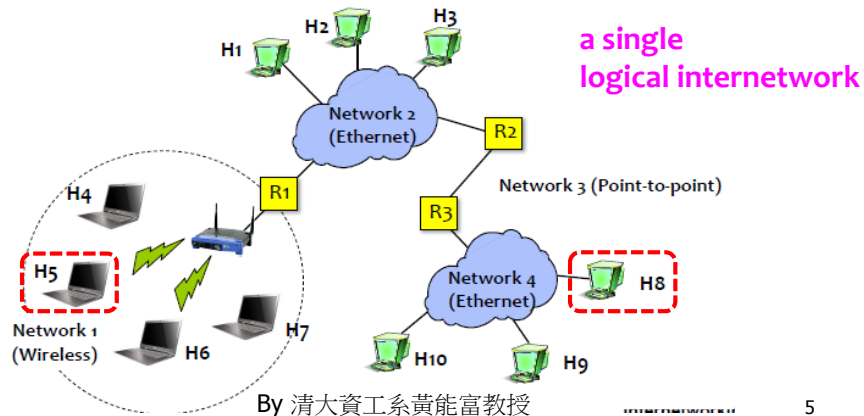
OSI Reference Model

2



Internetworking

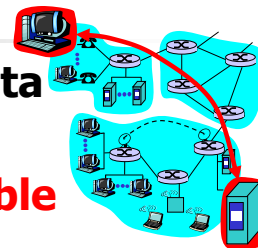
- An arbitrary collection of **networks** **interconnected** to provide some sort of **host-to-host packet delivery service** (**host-to-host transmission**)



5

IP Packet Delivery Model

- **Connectionless** model for data delivery
- **Best-effort** delivery (**unreliable service**)
 - packets are **lost**
 - packets are delivered **out of order**
 - **duplicate copies** of a packet are delivered
 - packets can be **delayed for a long time**



Computer Networking,
Jim Kurose, Keith Ross

6

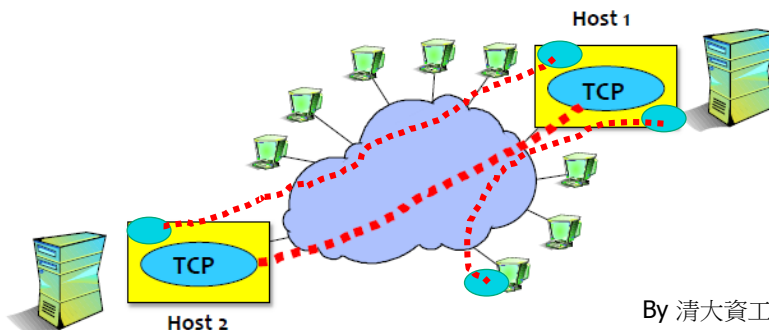
Process-to-process Protocols

- Challenge for Transport Protocols
- Develop algorithms that **turn** the **unreliable service** of the underlying network **into** the **service required by application programs**
 - *Unreliable service* → **Unreliable service (UDP)**
 - *Unreliable service* → **Reliable service (TCP)**

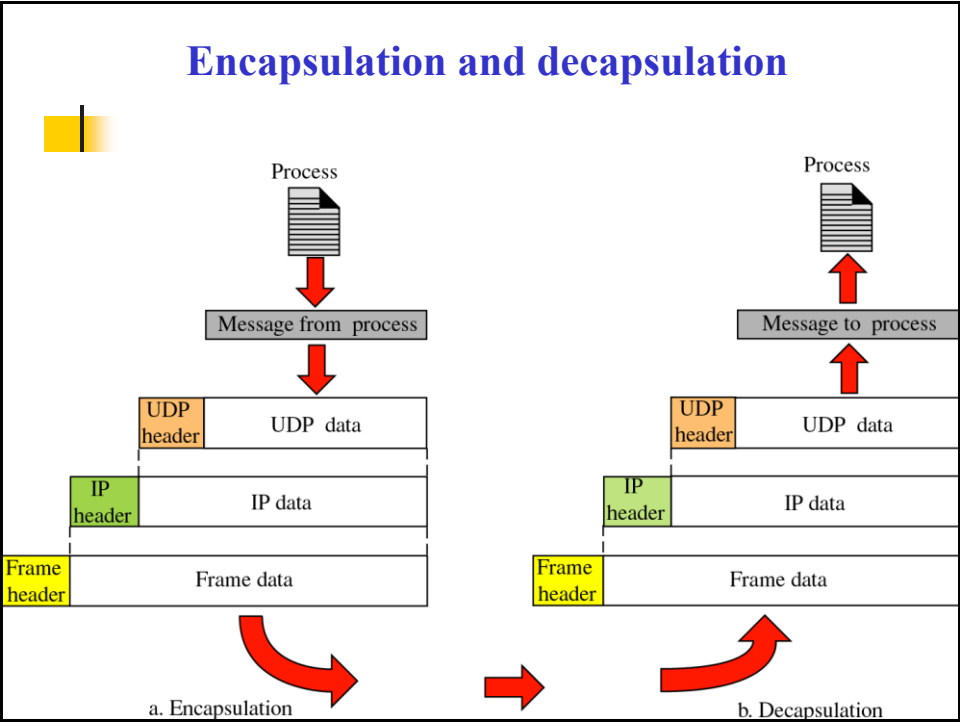
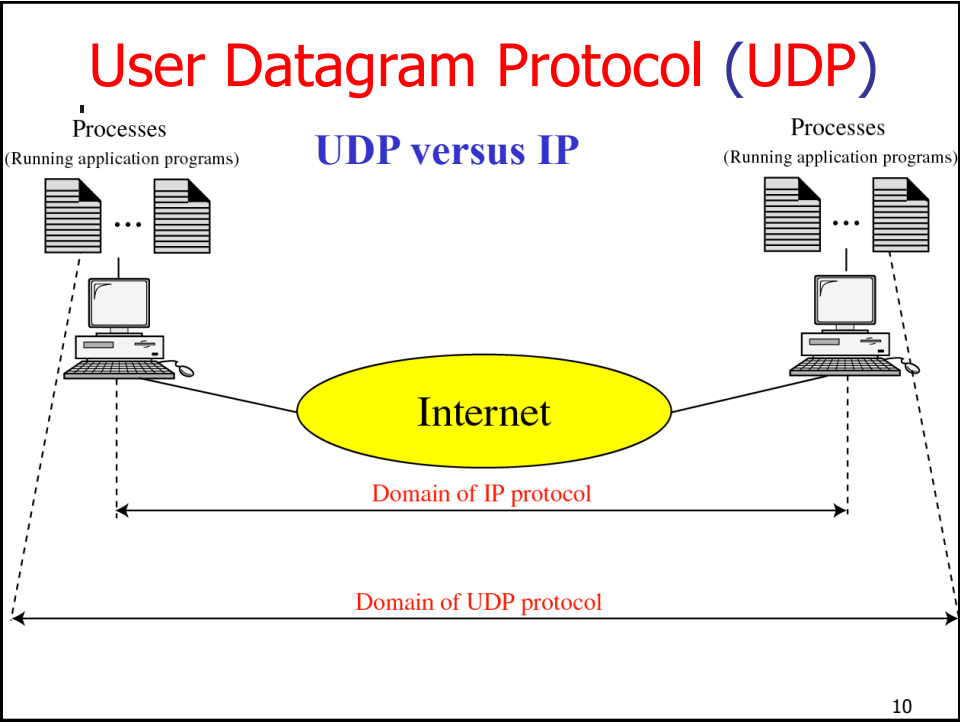
7

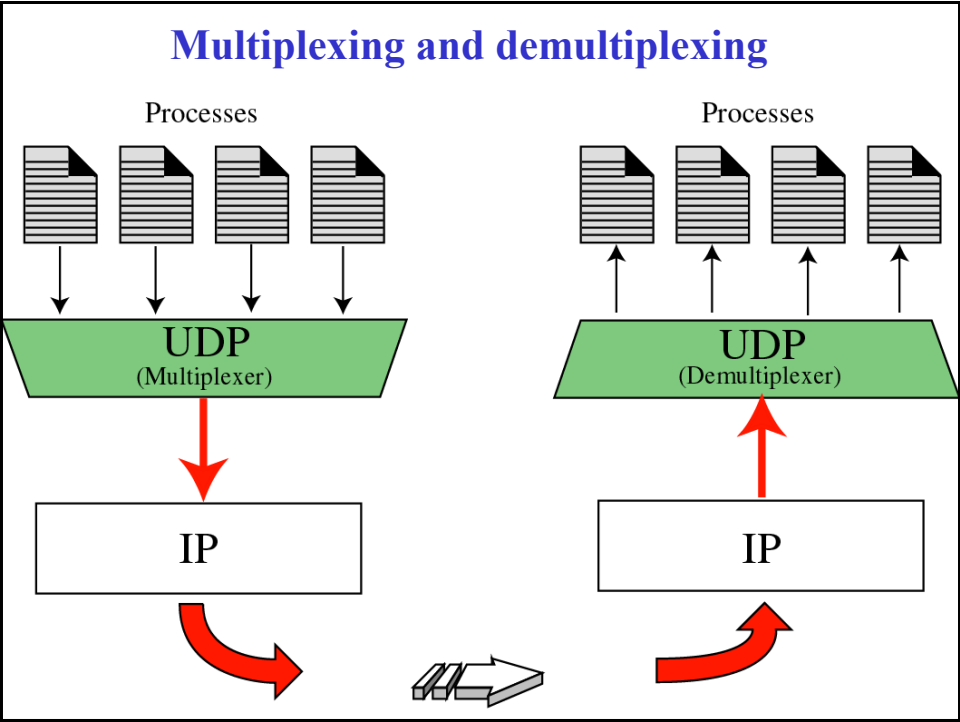
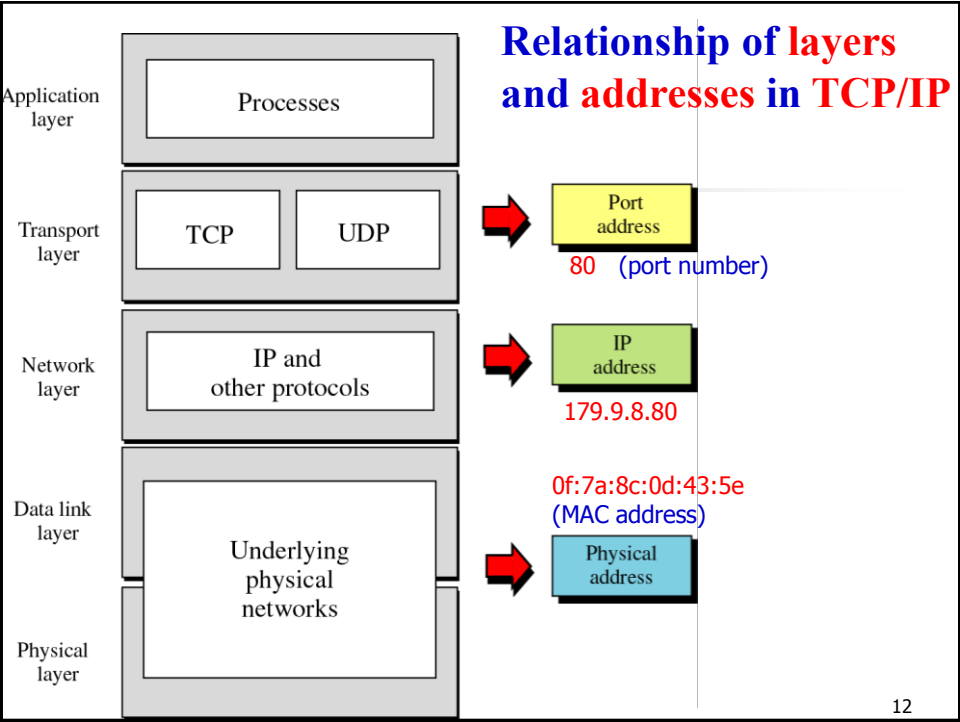
Description of Layers

- Transport Layer (提供不同主機 processes 之間的資料傳送)
 - Implements a process-to-process channel
 - Unit of data exchanges in this layer is called a message **segment**
 - **UDP (User Datagram Protocol)** – Unreliable service
 - **TCP (Transmission Control Protocol)** – Reliable service



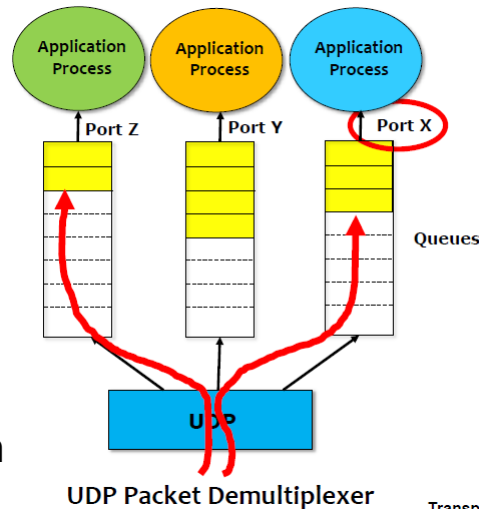
By 清大資工系黃能富教授





User Datagram Protocol (UDP)

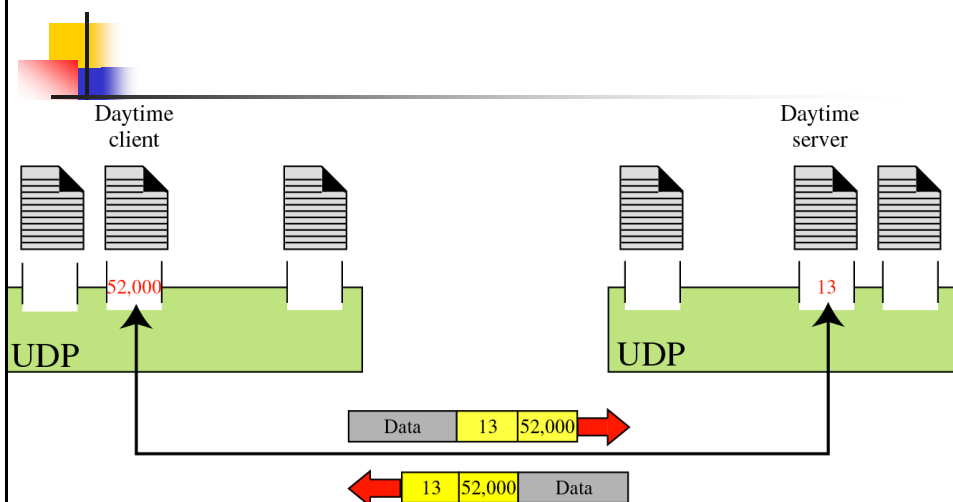
- Simple Demultiplexer
- Extends host-to-host delivery service of the underlying network into a **process-to-process** communication service
- Adds a level of **demultiplexing** which allows multiple application processes on each host to share the network



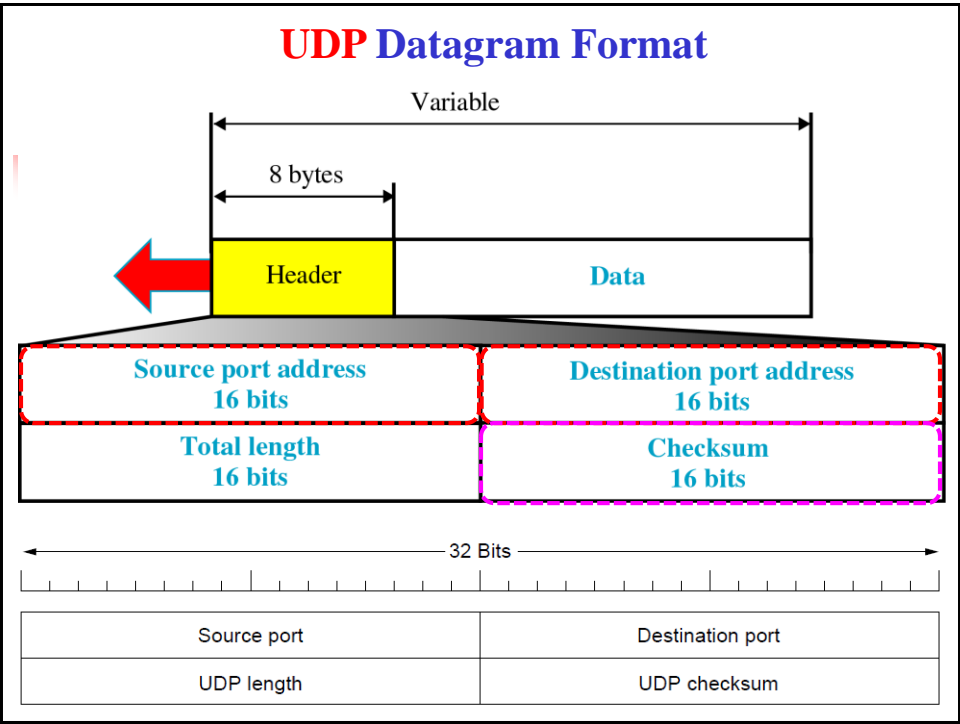
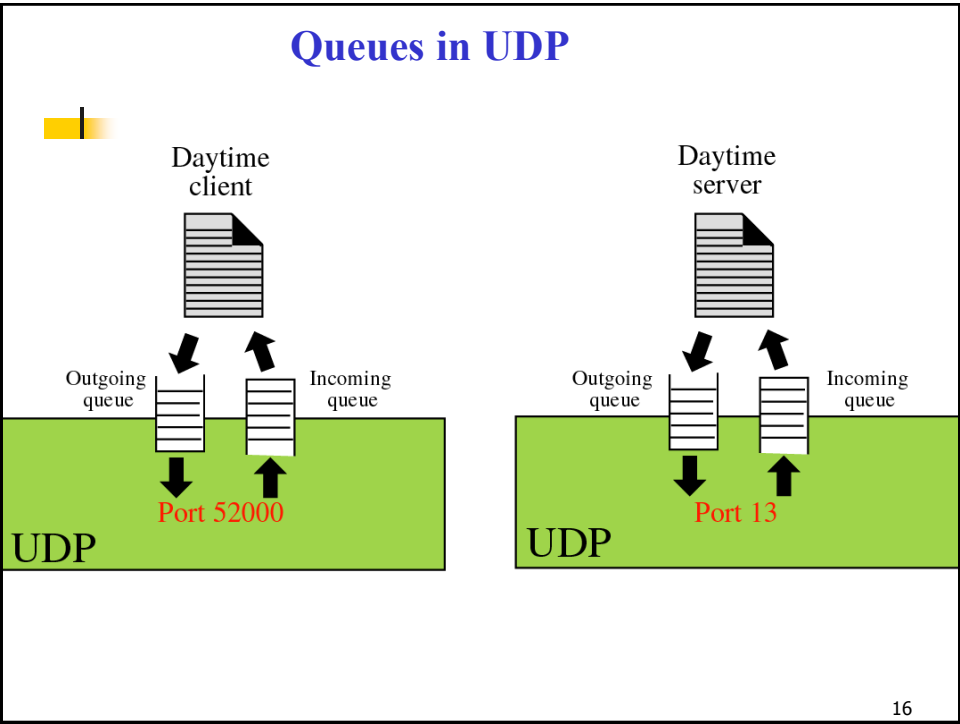
By 清大資工系黃能富教授

14

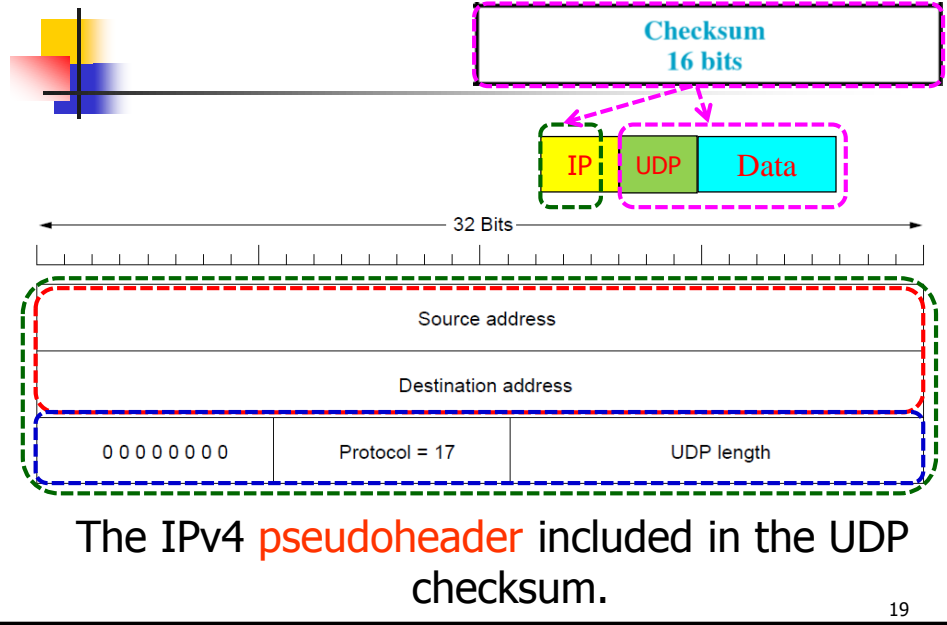
UDP Port numbers



15

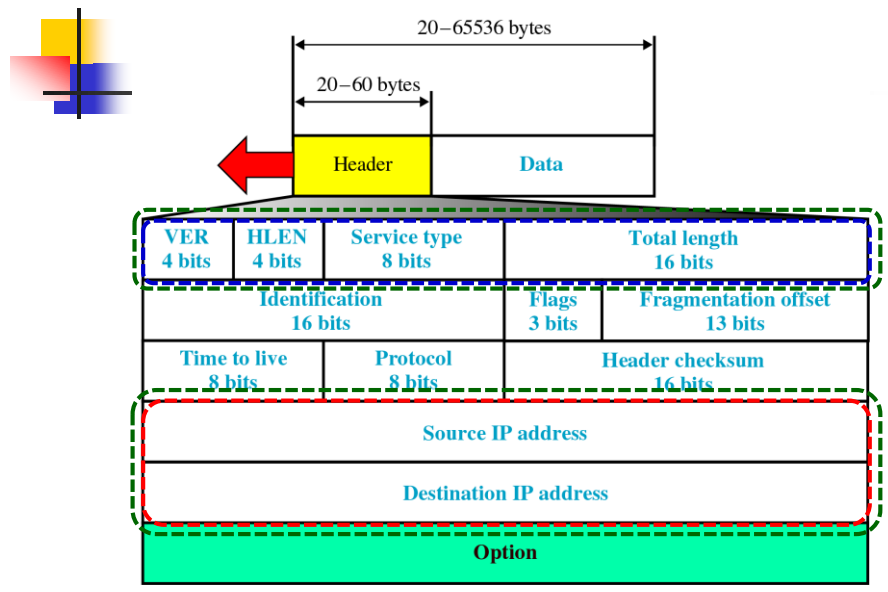


Introduction to UDP (2)

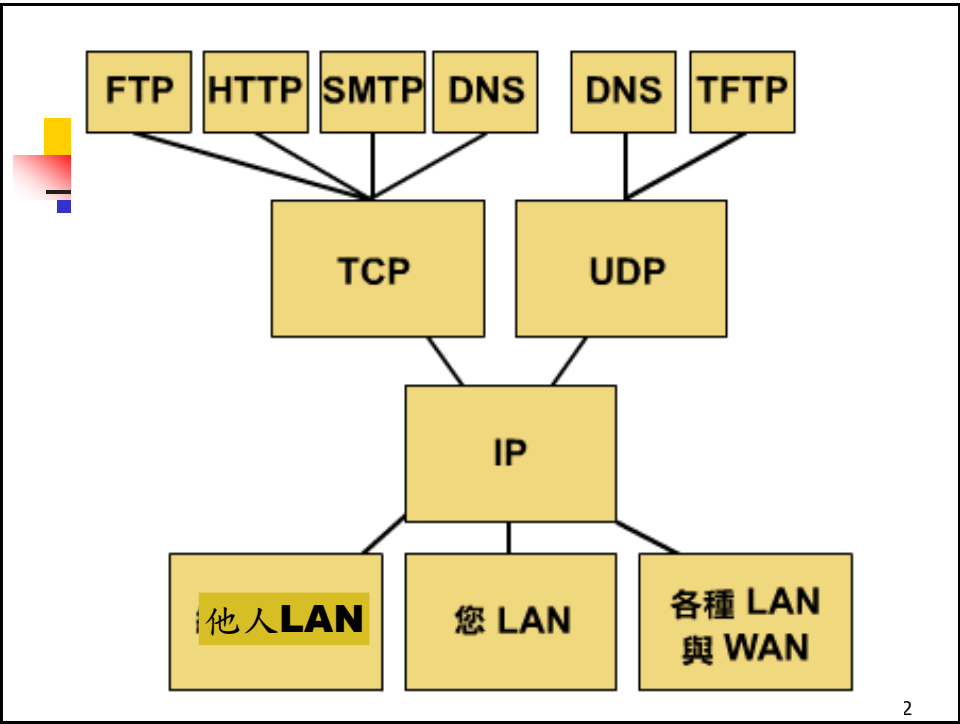
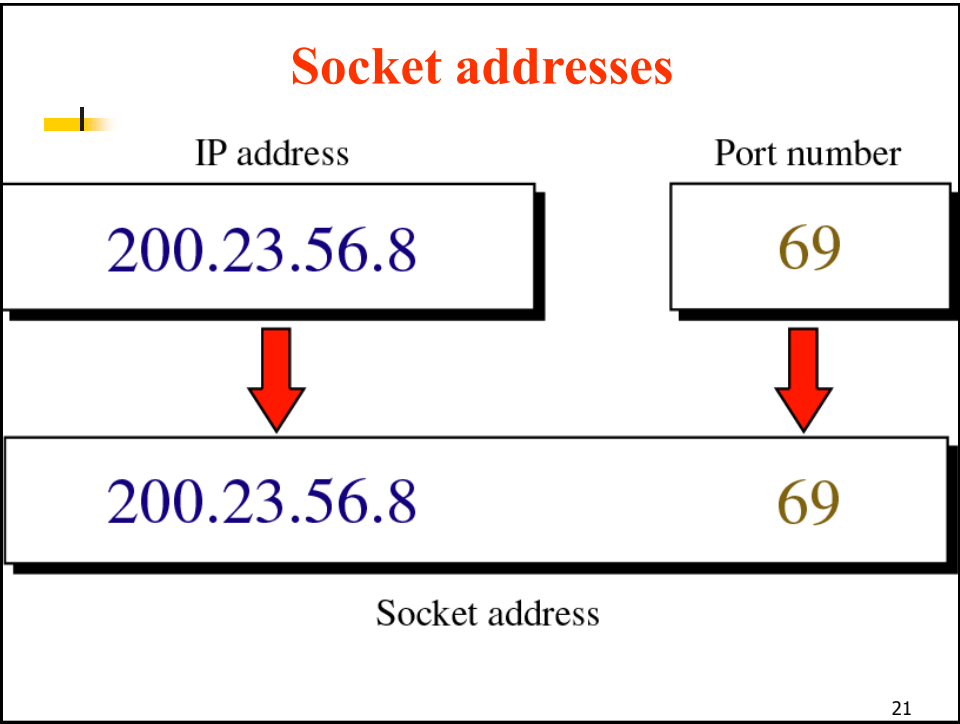


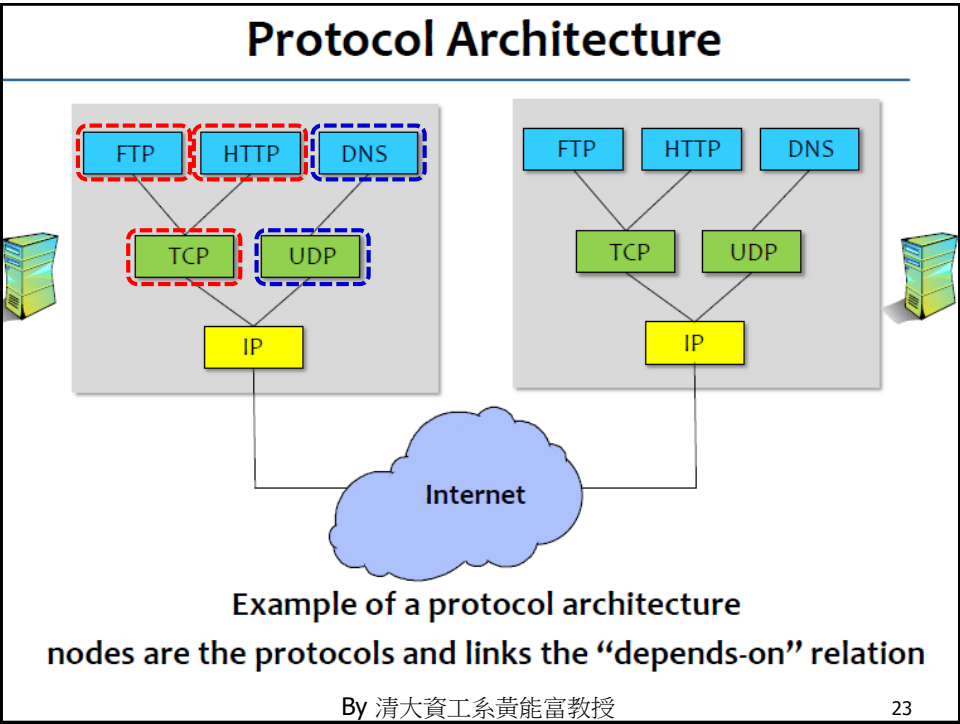
19

IP Datagram

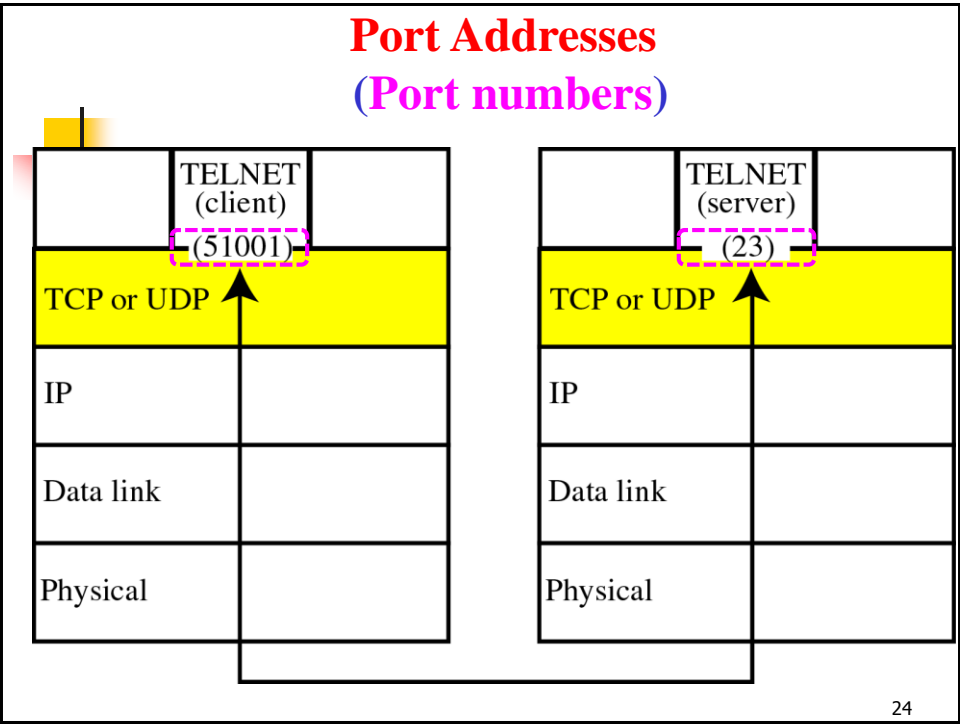


20

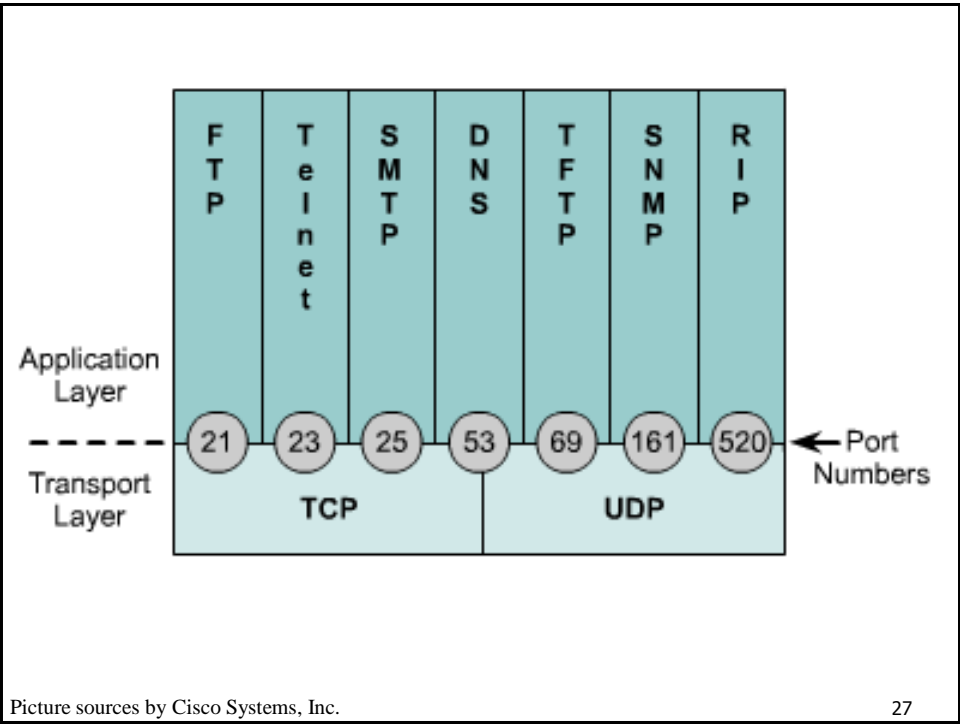
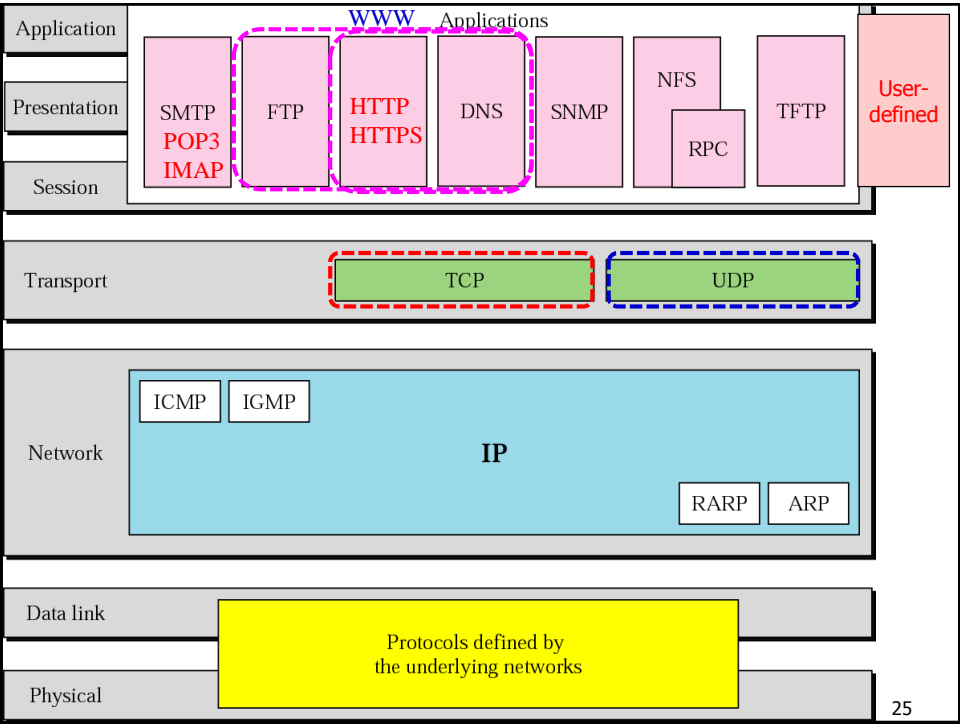




23



24



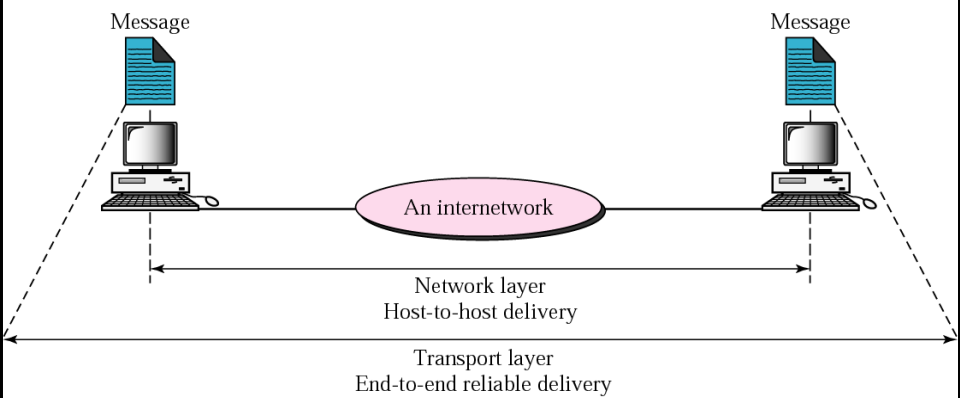
The TCP Service Model (1)

Port	Protocol	Use
20, 21	FTP	File transfer
22	SSH	Remote login, replacement for Telnet
25	SMTP	Email
80	HTTP	World Wide Web
110	POP-3	Remote email access
143	IMAP	Remote email access
443	HTTPS	Secure Web (HTTP over SSL/TLS)
543	RTSP	Media player control
631	IPP	Printer sharing

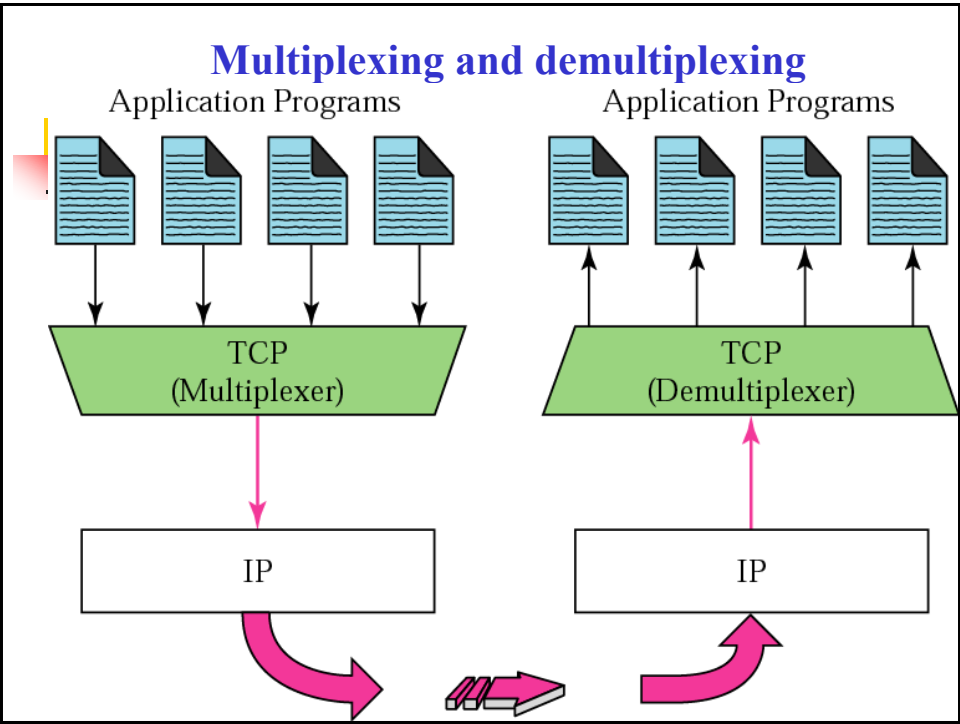
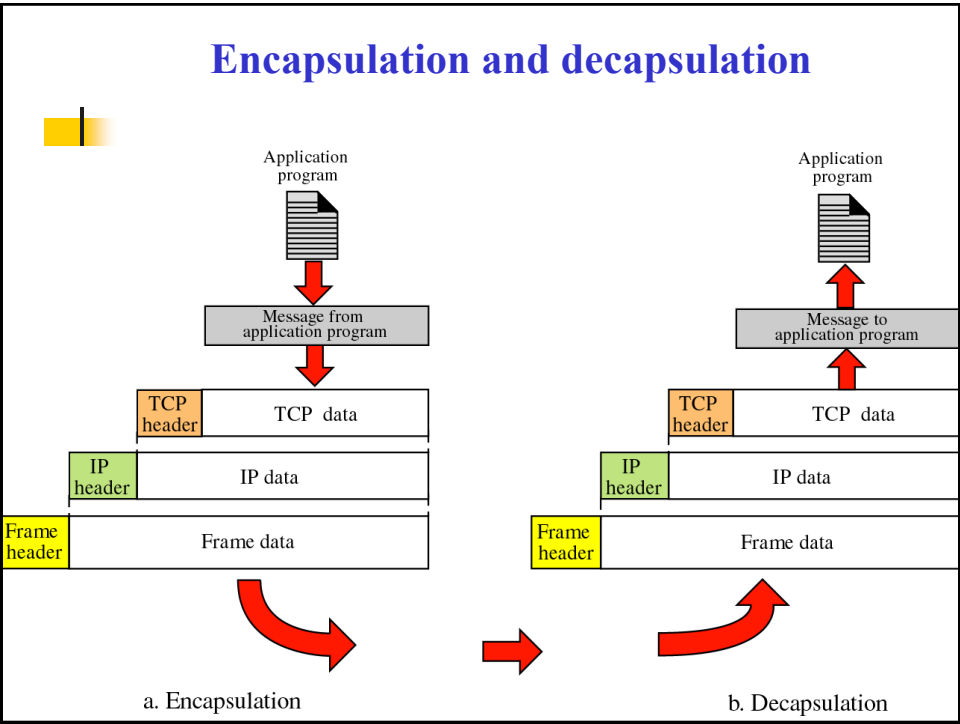
Some assigned ports

28

Reliable end-to-end delivery of a message TCP (Transmission Control Protocol)



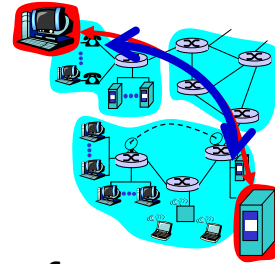
29



Reliable Byte Stream protocol (TCP)

- TCP offers the following services

- Reliable
- Connection oriented
- Byte-stream service



- Flow control: preventing senders from **overrunning** the capacity of **the receivers**
- Congestion control: preventing too much data from being injected into **the network**, thereby causing the Internet to become **overloaded**

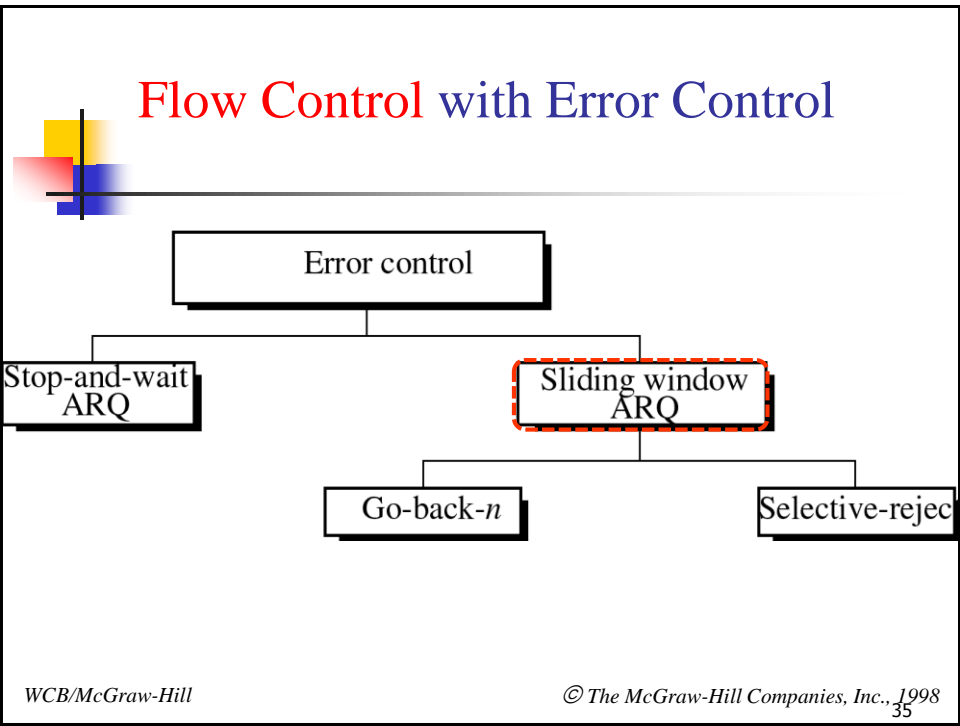
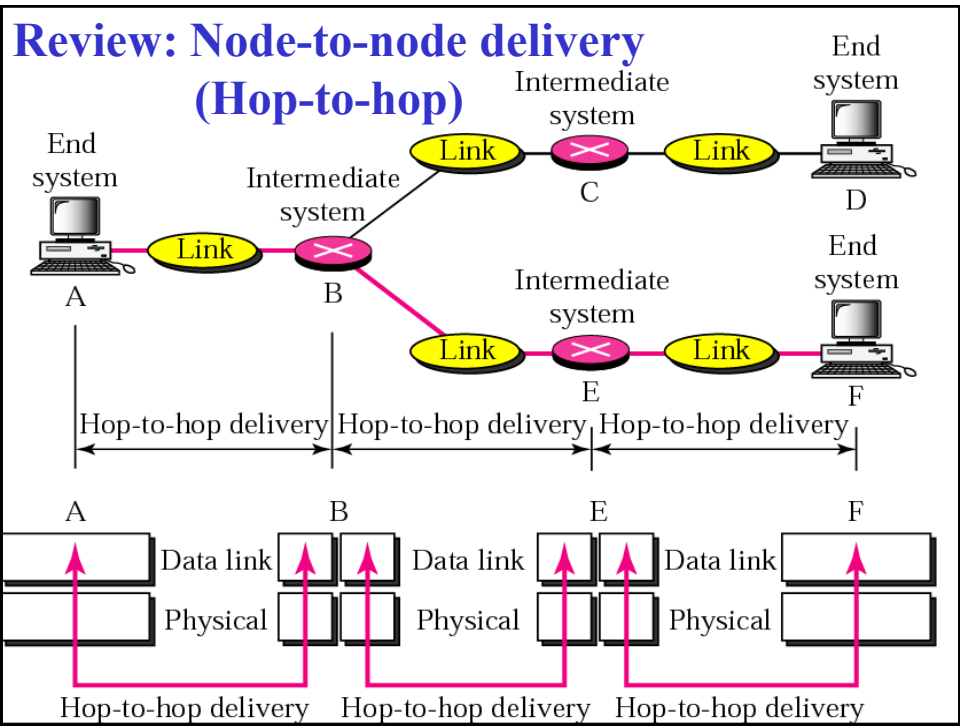
32

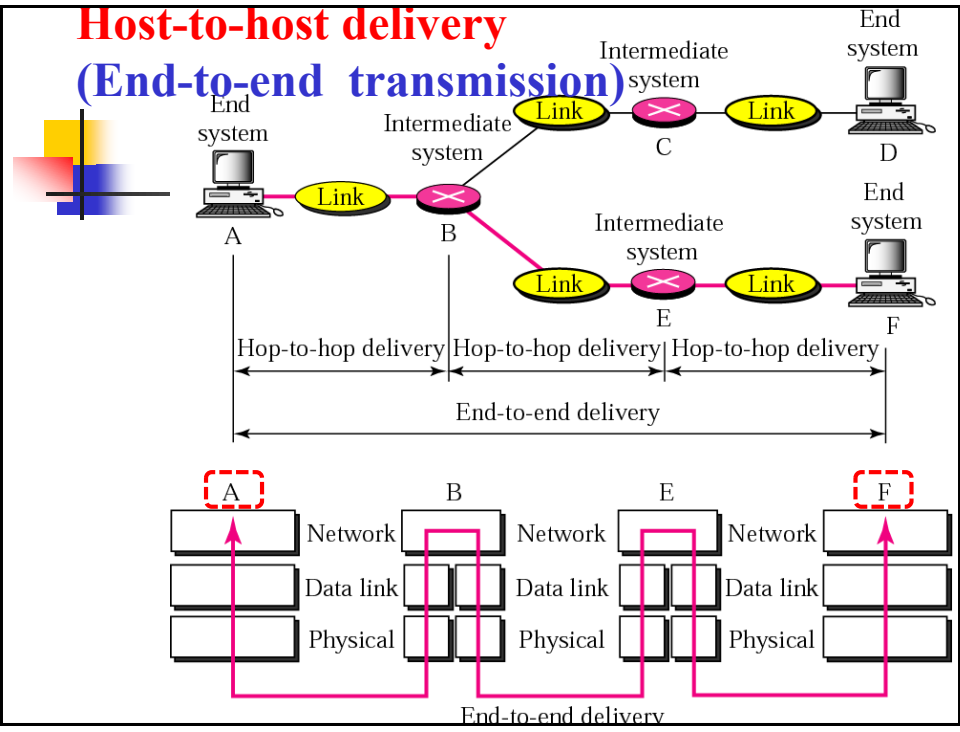
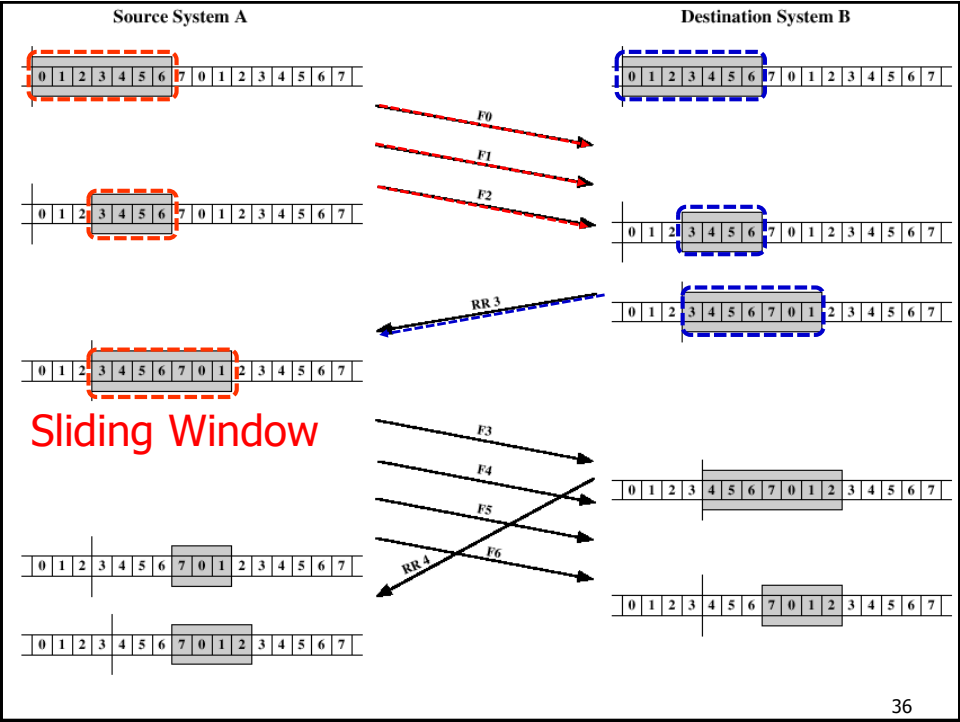
End-to-end Issues

- TCP runs **over the Internet** rather than a **point-to-point link**
- The **TCP sliding window algorithm** need to consider:
 - TCP supports **logical connections** between processes that are running on two different computers in **the Internet**
 - TCP connections are likely to have widely **different RTT times**
 - Packets may get **reordered** in the Internet

By 清大資工系黃能富教授

33





TCP Segment:

- TCP is a **byte-oriented** protocol
- The **sender writes bytes** into a TCP connection and the **receiver reads bytes** out of the TCP connection.
- The source TCP **buffers enough bytes from the sending process** to fill a **reasonably sized packet** and then sends this packet to its peer on the destination host.
- The destination TCP then puts the contents of the packet into a **receive buffer**, and the receiving process reads from this buffer.

By 清大資工系黃能富教授

38

TCP Segment

Write bytes

Read bytes

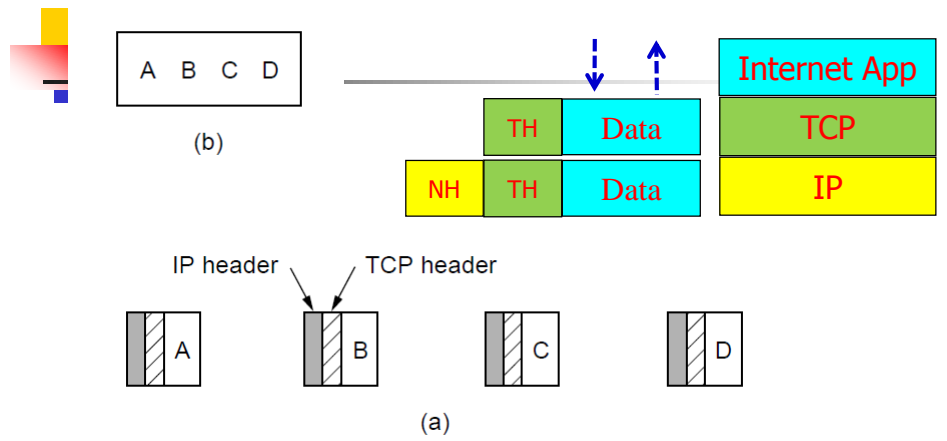
Transmit segments

How TCP manages a byte stream.

By 清大資工系黃能富教授

39

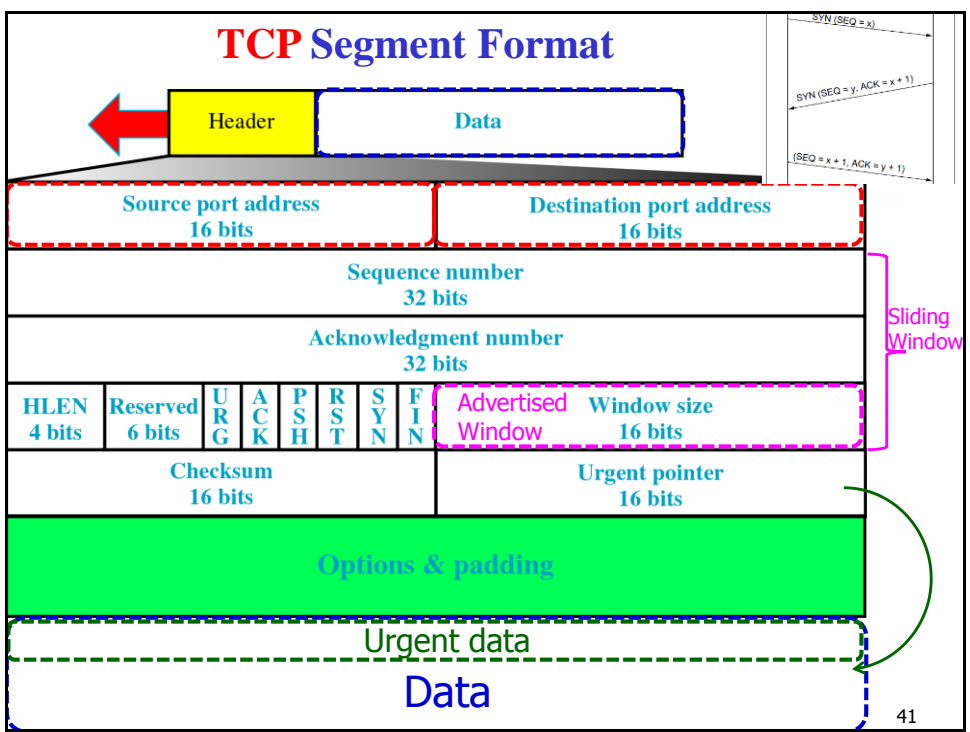
The TCP Service Model (2)



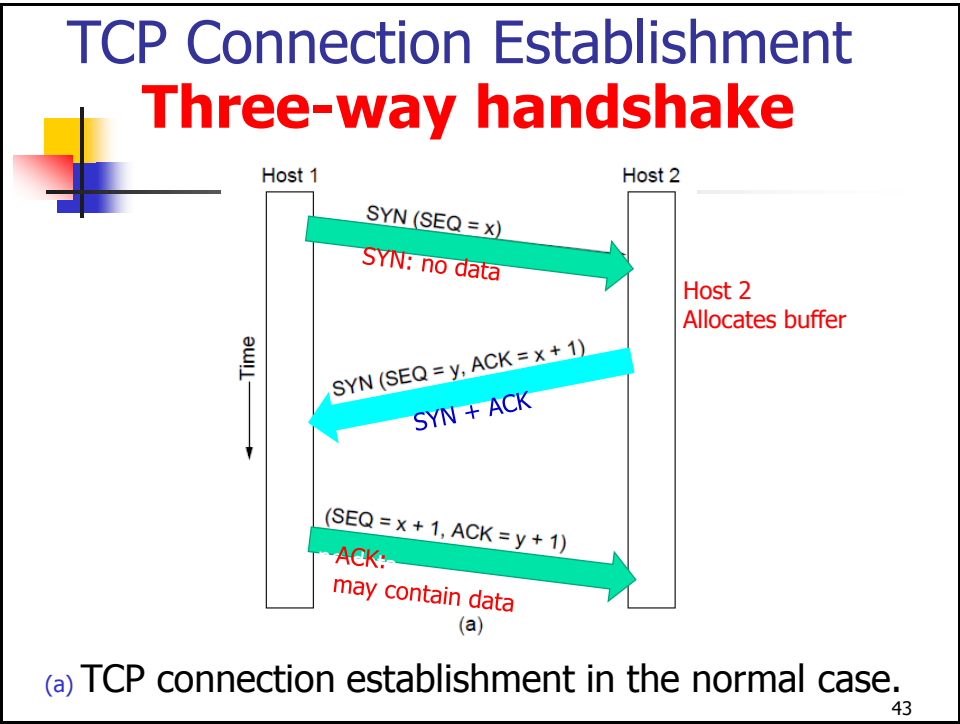
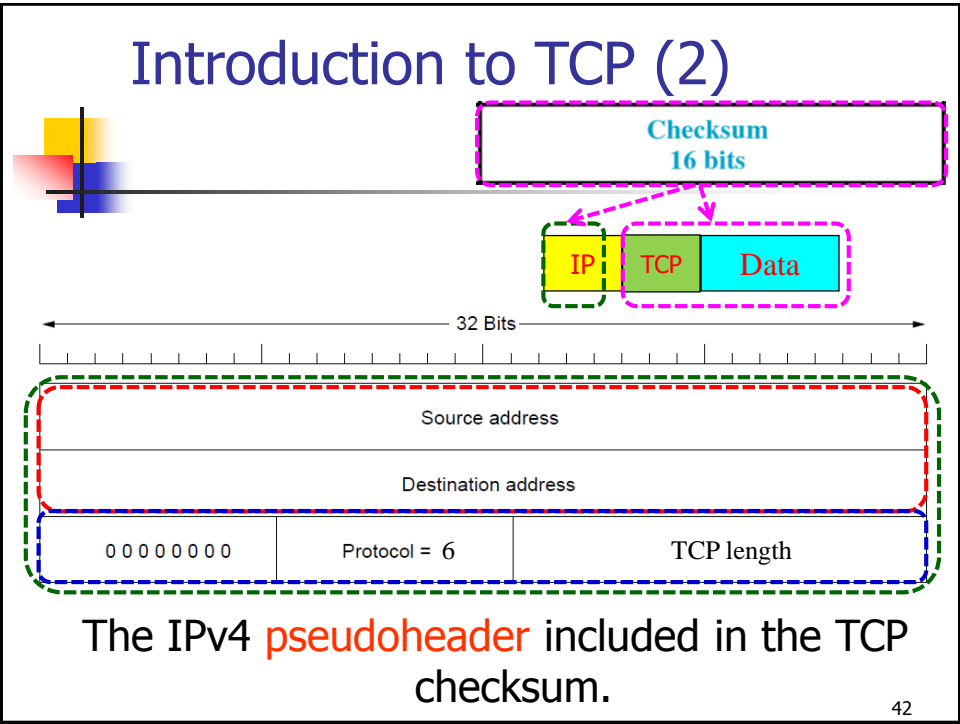
- (a) Four 512-byte segments sent as separate IP diagrams
- (b) The 2048 bytes of data delivered to the application in a single READ call

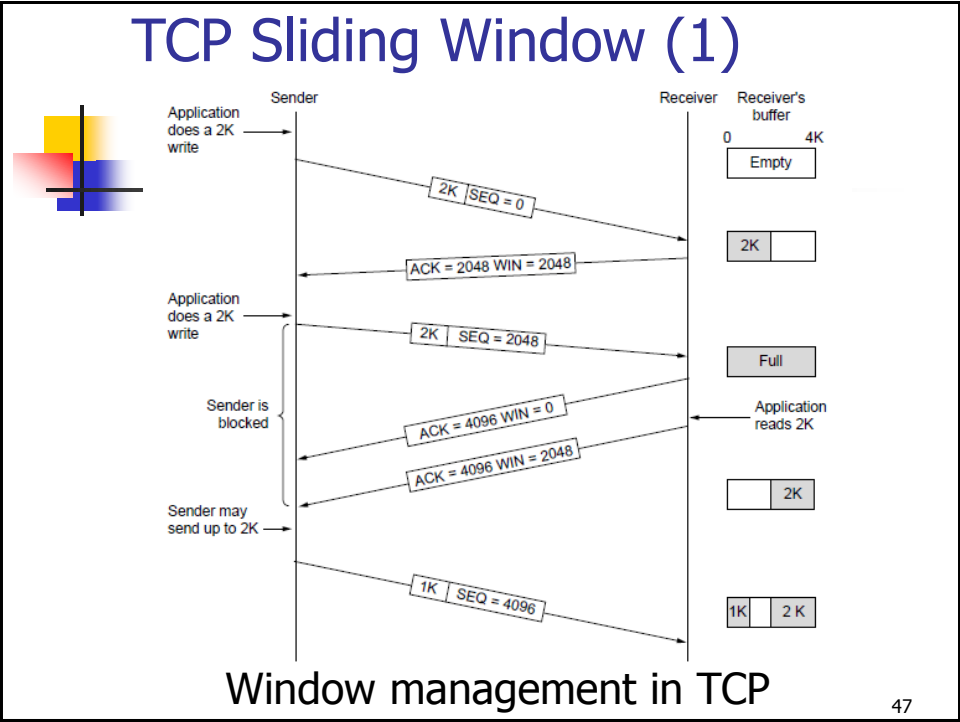
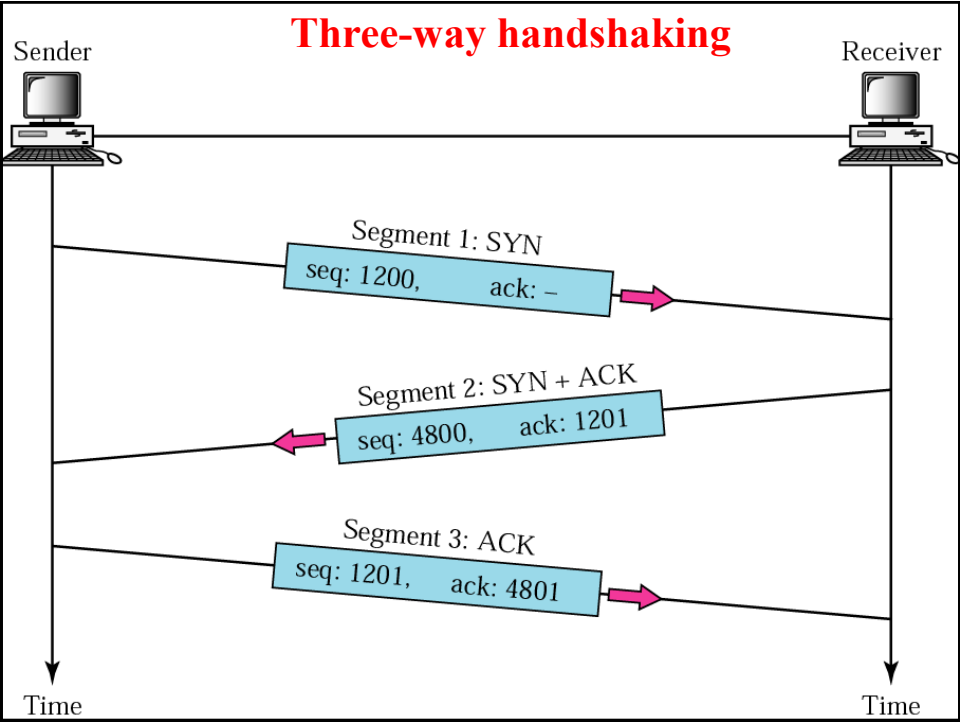
40

TCP Segment Format



41





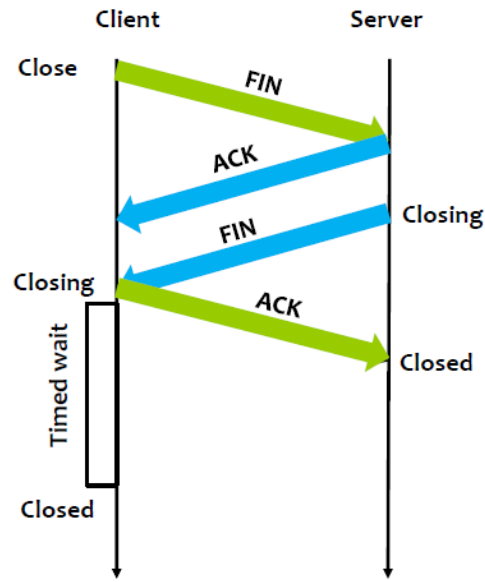
TCP Connection Management (cont.)

Closing a connection:

client closes socket:
`clientSocket.close()`
`;`

Step 1: Client sends TCP FIN control segment to server

Step 2: Server receives FIN, replies with ACK. Closes connection, sends FIN.



By 清大資工系黃能富教授

48

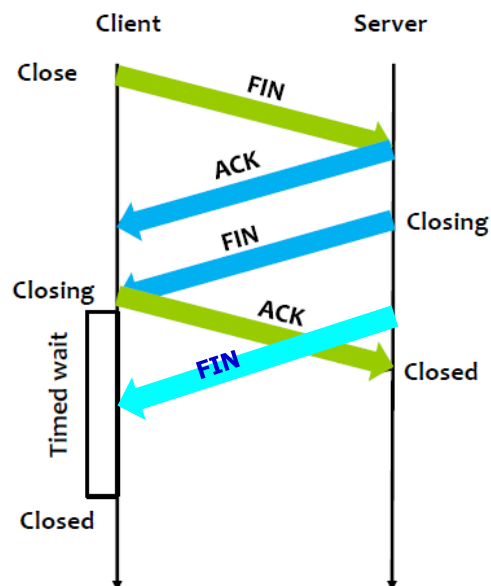
TCP Connection Management (cont.)

Step 3: Client receives FIN, replies with ACK.

- Enters “timed wait” - will respond with ACK to received FINs

Step 4: Server receives ACK. Connection closed.

Note: with small modification, can handle simultaneous FINs.



By 清大資工系黃能富教授

72



TCP Connection Management Modeling

(1)



The states used in the TCP connection management finite state machine.

State	Description
CLOSED	No connection is active or pending
LISTEN	The server is waiting for an incoming call
SYN RCVD	A connection request has arrived; wait for ACK
SYN SENT	The application has started to open a connection
ESTABLISHED	The normal data transfer state
FIN WAIT 1	The application has said it is finished
FIN WAIT 2	The other side has agreed to release
TIME WAIT	Wait for all packets to die off
CLOSING	Both sides have tried to close simultaneously
CLOSE WAIT	The other side has initiated a release
LAST ACK	Wait for all packets to die off

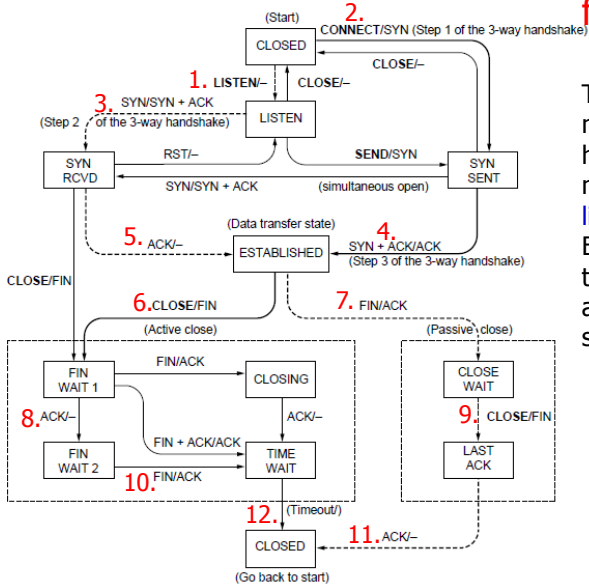
53

TCP Connection Management Modeling

(2)



TCP connection management
finite state machine.

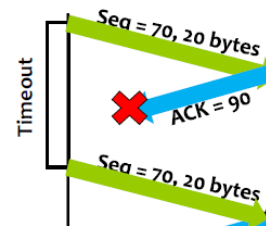


The heavy solid line is the normal path for a client. The heavy dashed line is the normal path for a server. The light lines are unusual events. Each transition is labeled by the event causing it and the action resulting from it, separated by a slash.

54

Timeout value for Retransmission

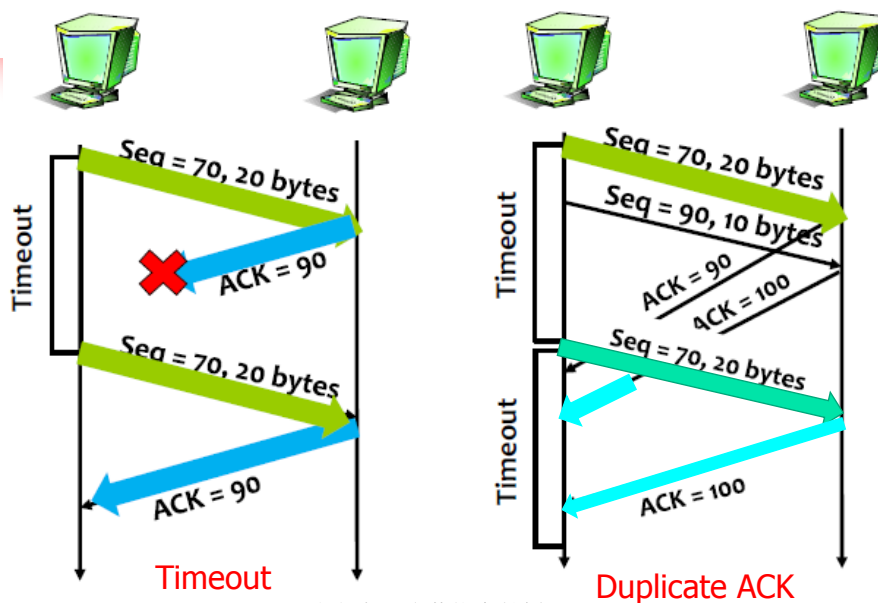
- Measure **SampleRTT** for each **segment/ACK** pair
- Compute **weighted average of RTT**
 - $\text{EstRTT} = a \times \text{EstRTT} + (1 - a) \times \text{SampleRTT}$
 - a between 0.8 and 0.9
- Set timeout based on **EstRTT**
 - $\text{TimeOut} = 2 \times \text{EstRTT}$



By 清大資工系黃能富教授

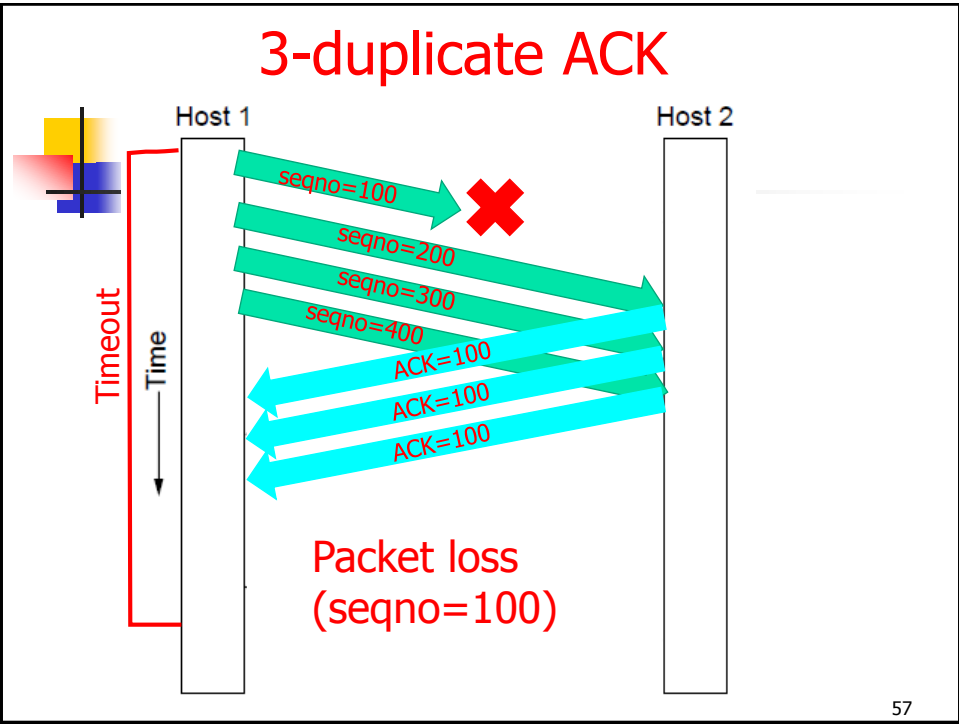
55

TCP retransmission scenarios





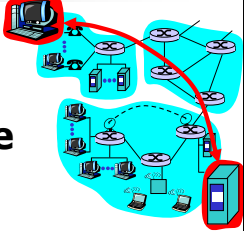
By 清大資工系黃能富教授

56



TCP Congestion Control

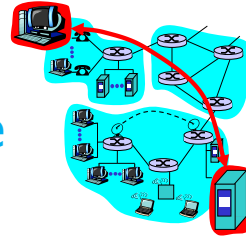
- Flow control
 - Advertised Window
- Congestion control
 - **Congestion window: self-clocking** by using **ACKs** to pace the transmission of packets



Source: <https://cc0.wfublog.com/> or <http://photopin.com/> or Google

8

TCP Congestion Control



Additive Increase Multiplicative Decrease (AIMD)

- TCP's **effective window** is revised as follows:
 - $\text{MaxWindow} = \text{MIN}(\text{CongestionWindow}, \text{AdvertisedWindow})$
 - **EffectiveWindow** = $\text{MaxWindow} - (\text{LastByteSent} - \text{LastByteAcked})$.
- A TCP source is allowed to send **no faster than the slowest component** can accommodate
 - the **network** or the **destination host**

By 清大資工系黃能富教授

59

Additive Increase Multiplicative Decrease (AIMD)



A packet =
maximum segment size
(MSS)
bytes



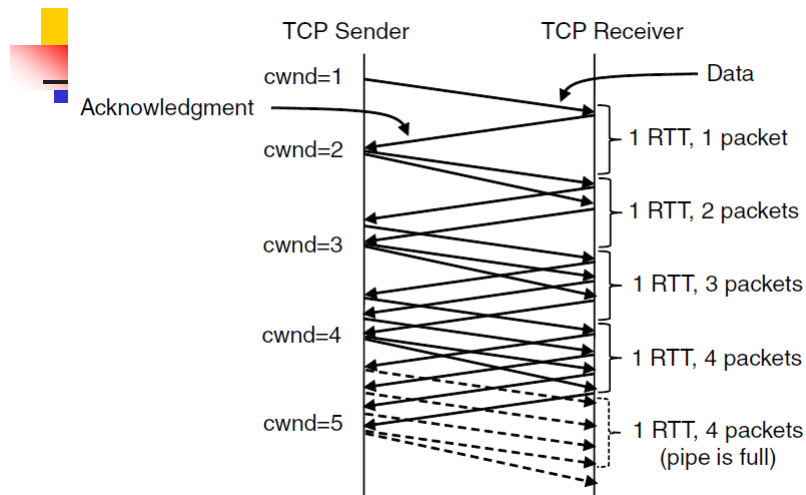
Additive increase

One more packet is added for each RTT

By 清大資工系黃能富教授

60

TCP Congestion Control (1)

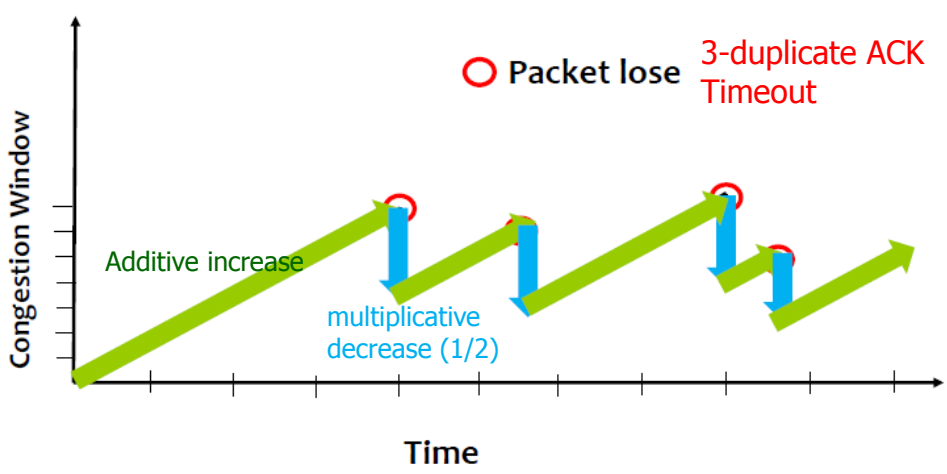


Additive increase from an initial congestion window of 1 segment.

61

Additive Increase Multiplicative Decrease (AIMD)

Trace: Sawtooth behavior



By 清大資工系黃能富教授

62

TCP Congestion Control

- **Slow start**: to increase the congestion window **rapidly from a cold start**.
- Slow start effectively increases the congestion window **exponentially, rather than linearly**.



http://auto.ferrari.com/zh_CN/wp-content/uploads/sites/11/2015/02/news-ferrari-580x362.jpg
http://news.xinhuanet.com/auto/2009-04/15/xinsrc_36204061509371092322728.jpg

63

Slow Start

- Initially, the CongestionWindow= 1 packet.
- Example: MSS = 500 bytes, RTT = 200 msec, initial rate = 20 kbps
- 每收到一個ACK就加1 packet (MSS)
- Upon **receiving** the corresponding **two ACKs**, TCP increments CongestionWindow by 2—one for each ACK—and next **sends four packets**.
- TCP effectively **doubles the number of packets** it has in transit **every RTT**.

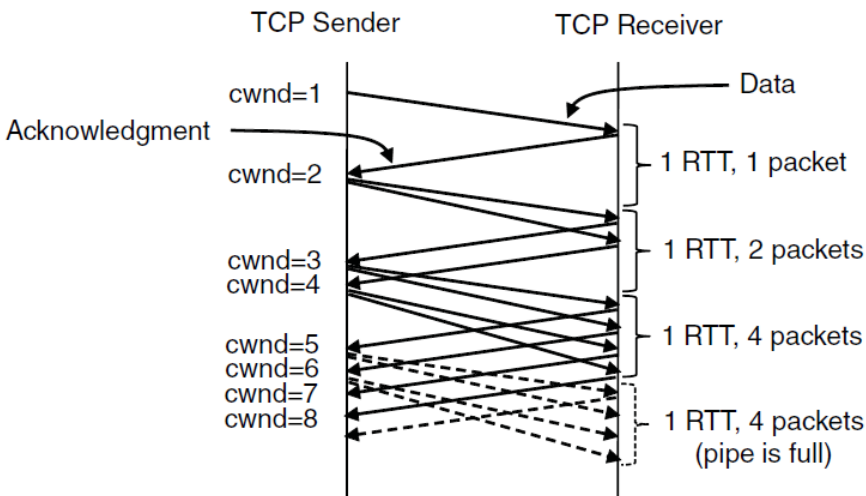


By 清大資工系黃能富教授

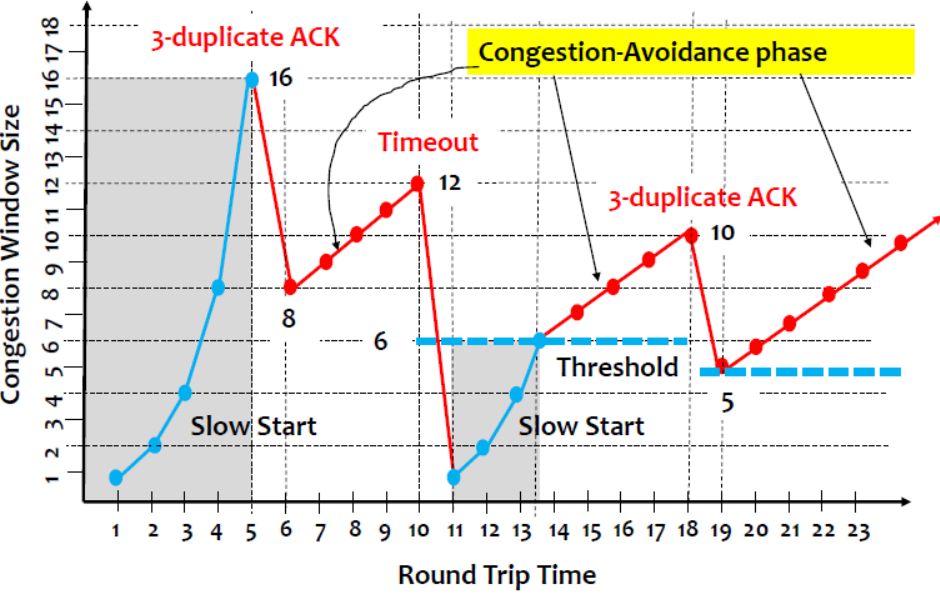
64

TCP Congestion Control (2)

Slow start from an initial congestion window of 1 segment

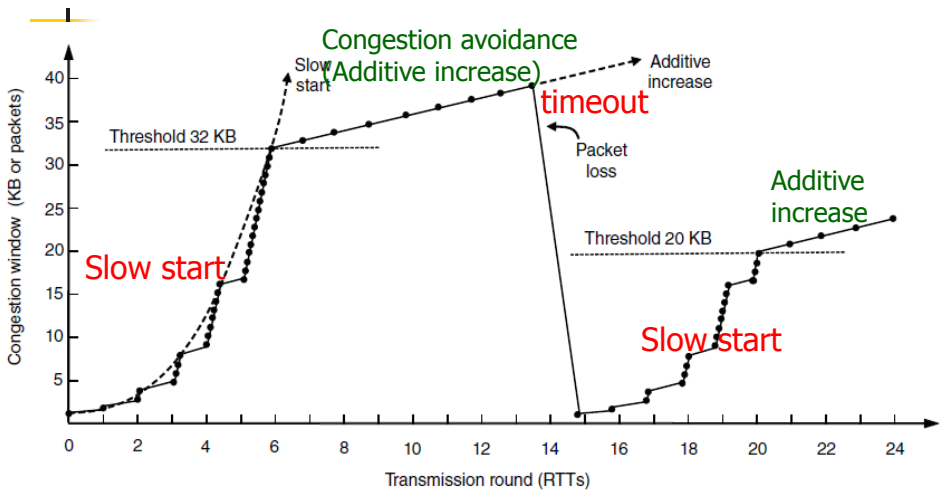


TCP Congestion Control



By 清大資工系黃能富教授

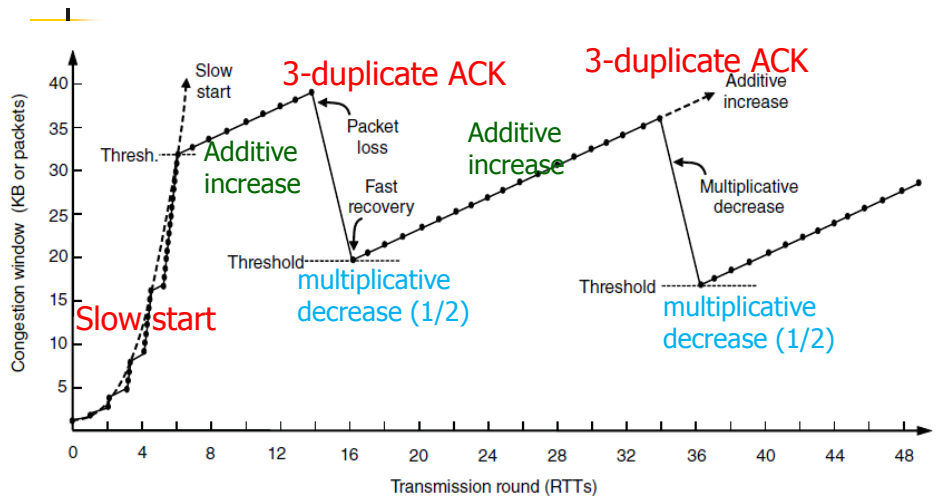
TCP Congestion Control (3)



Slow start followed by additive increase in TCP Tahoe.

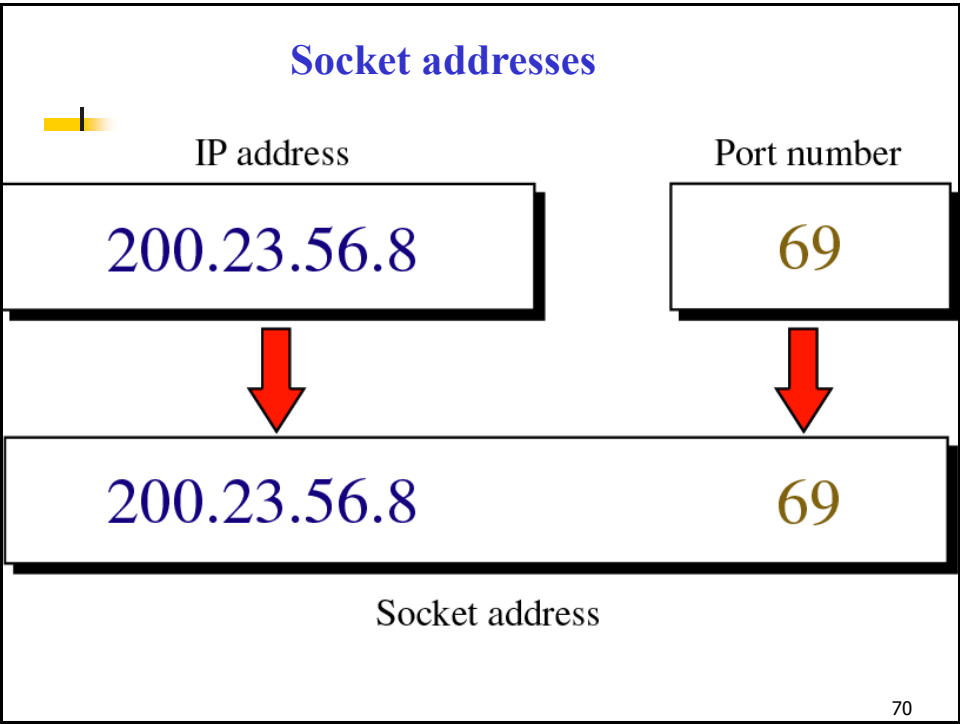
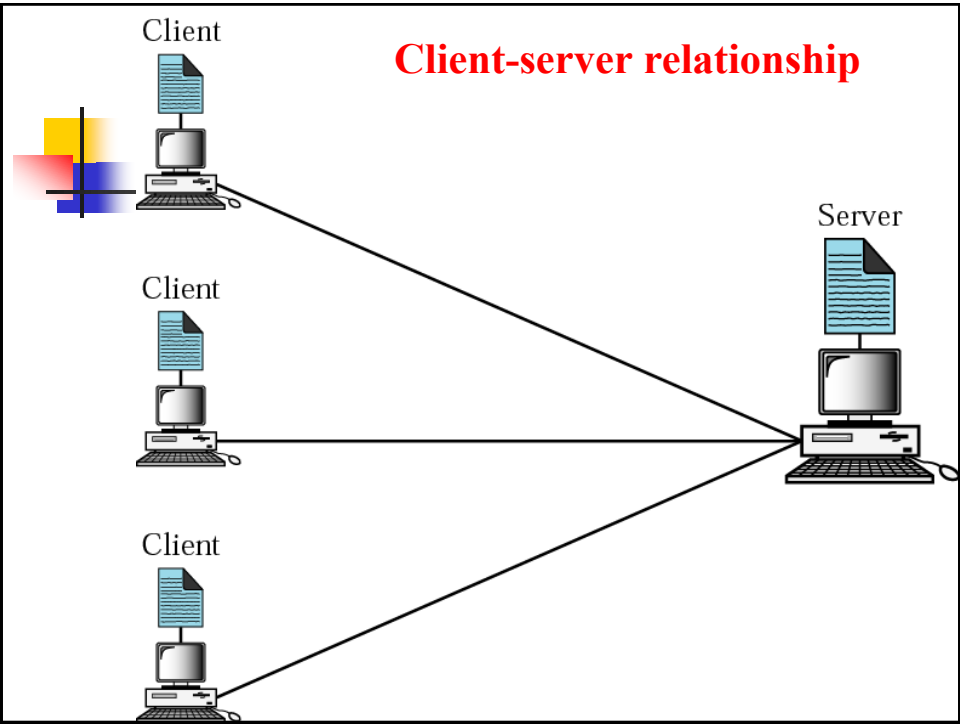
67

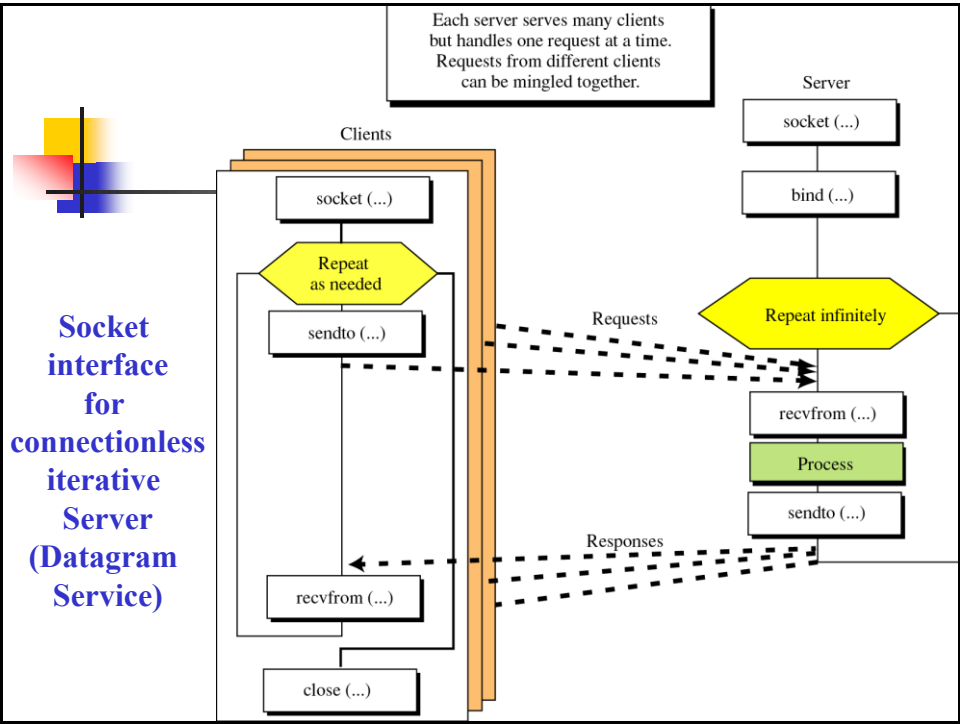
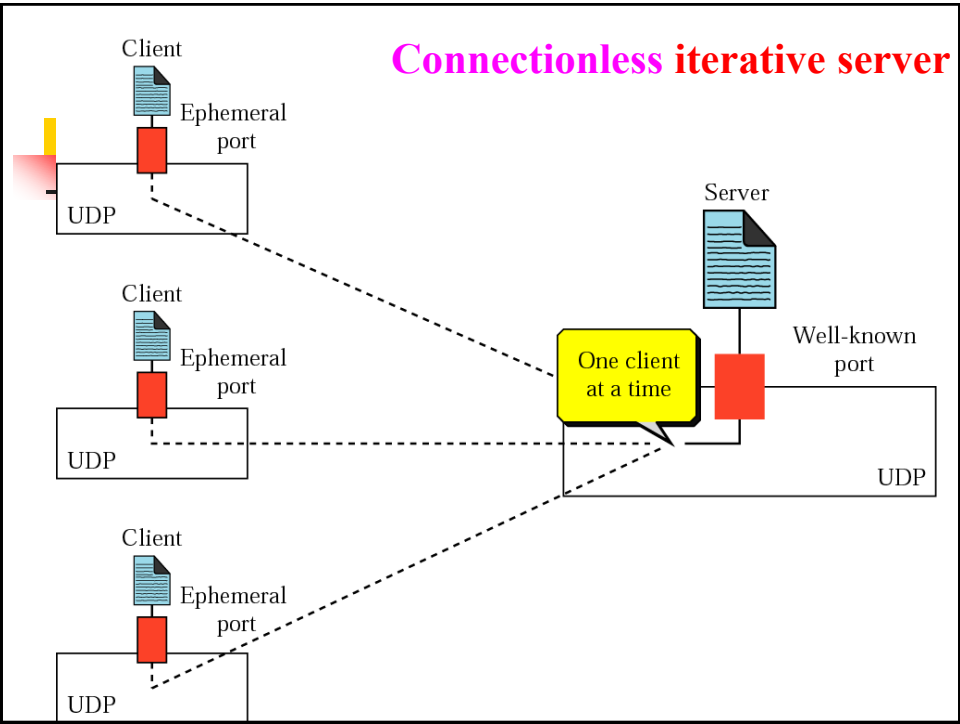
TCP Congestion Control (4)

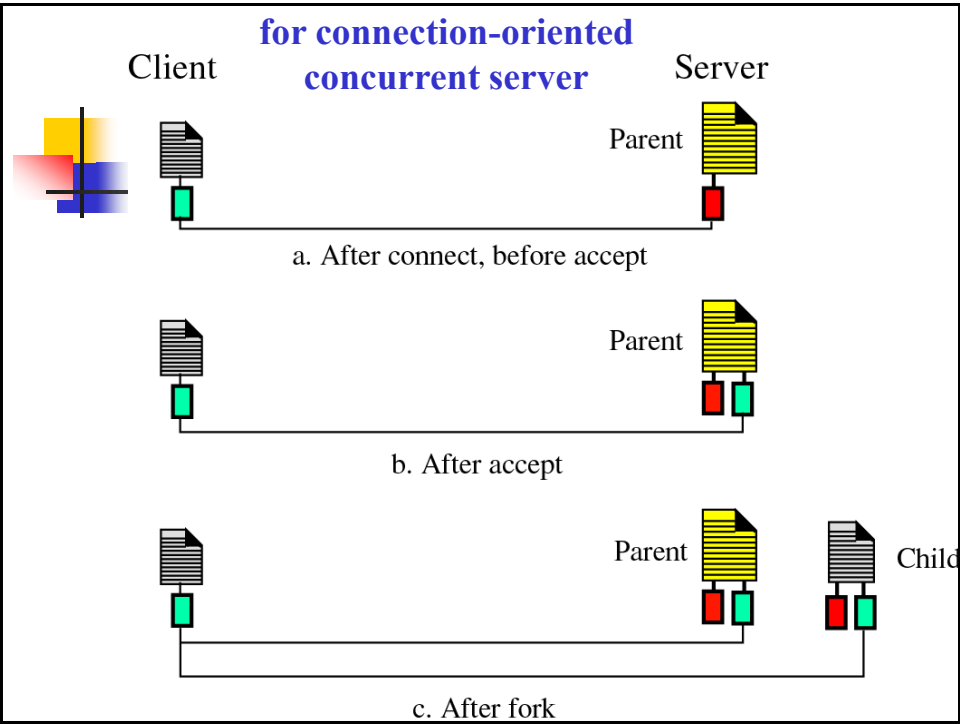
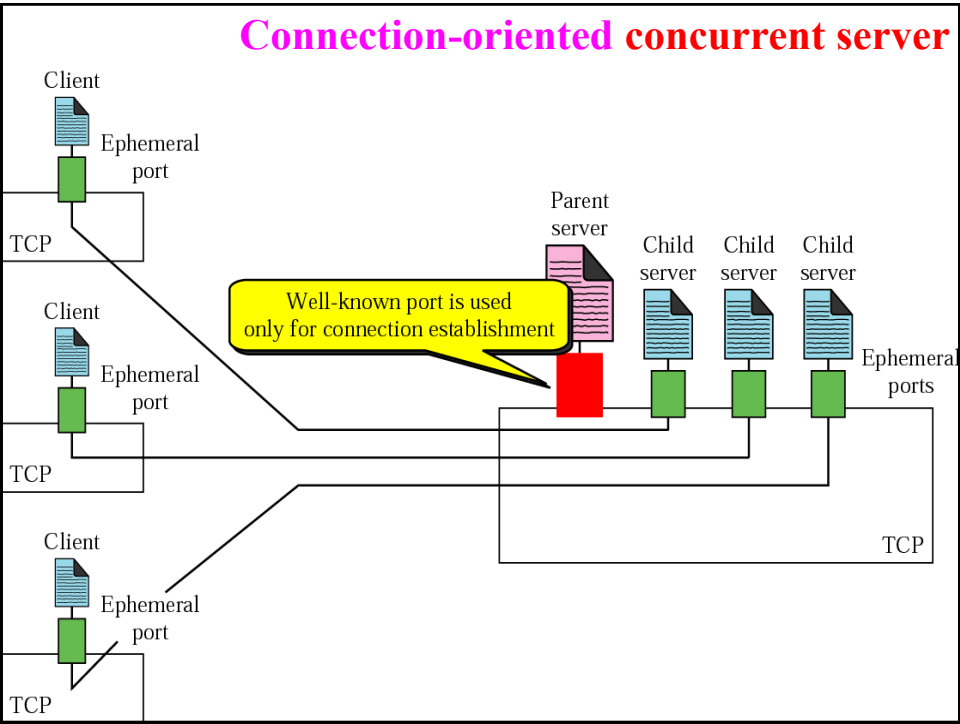


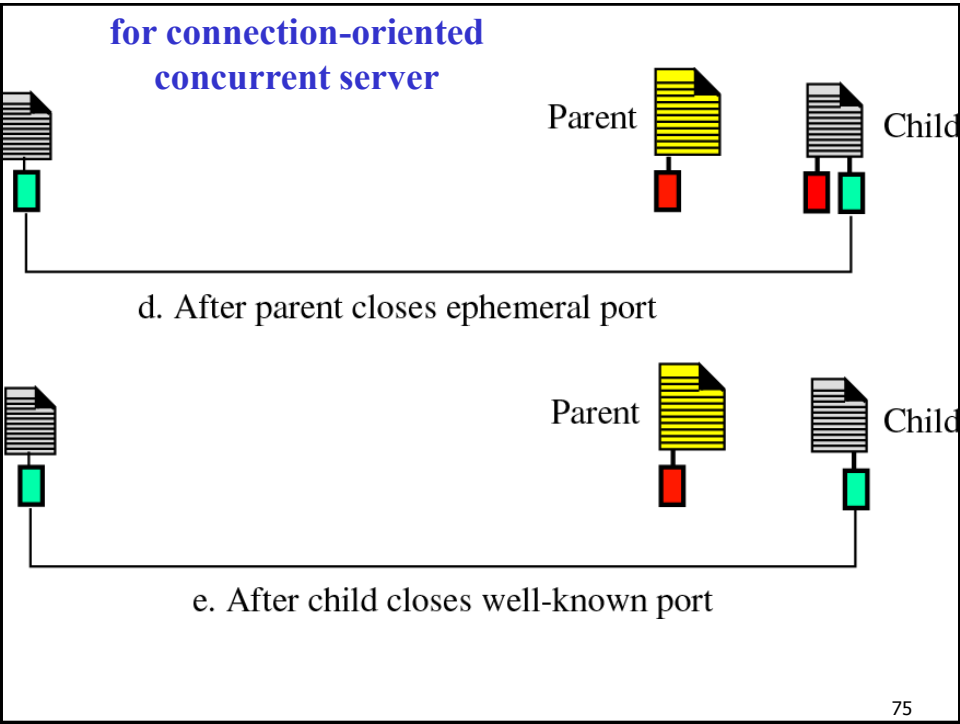
Fast recovery and the sawtooth pattern of TCP Reno.

68



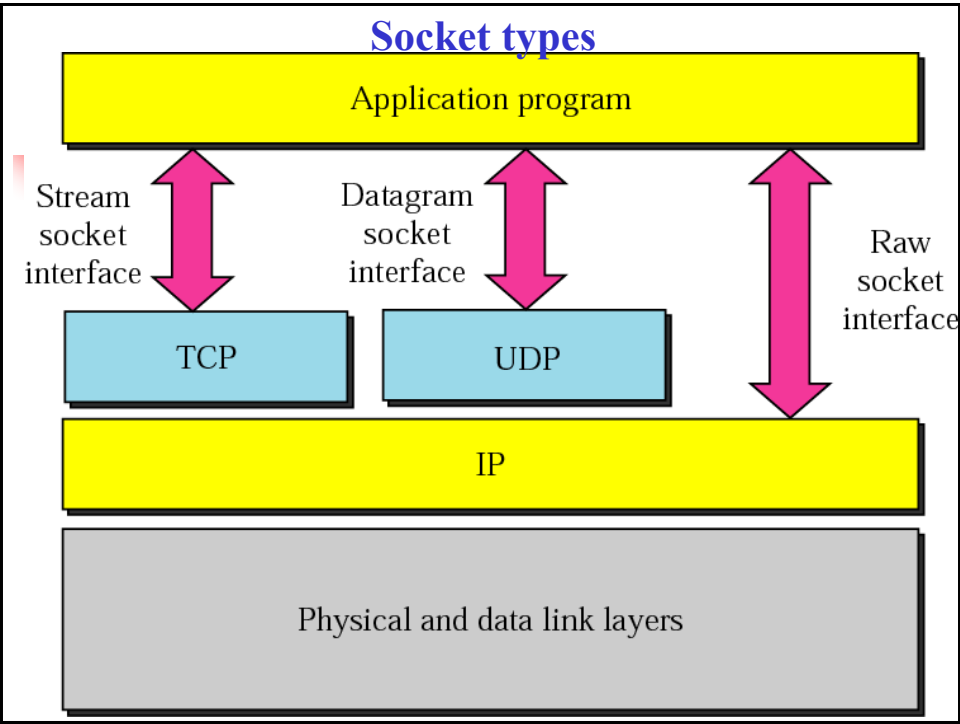
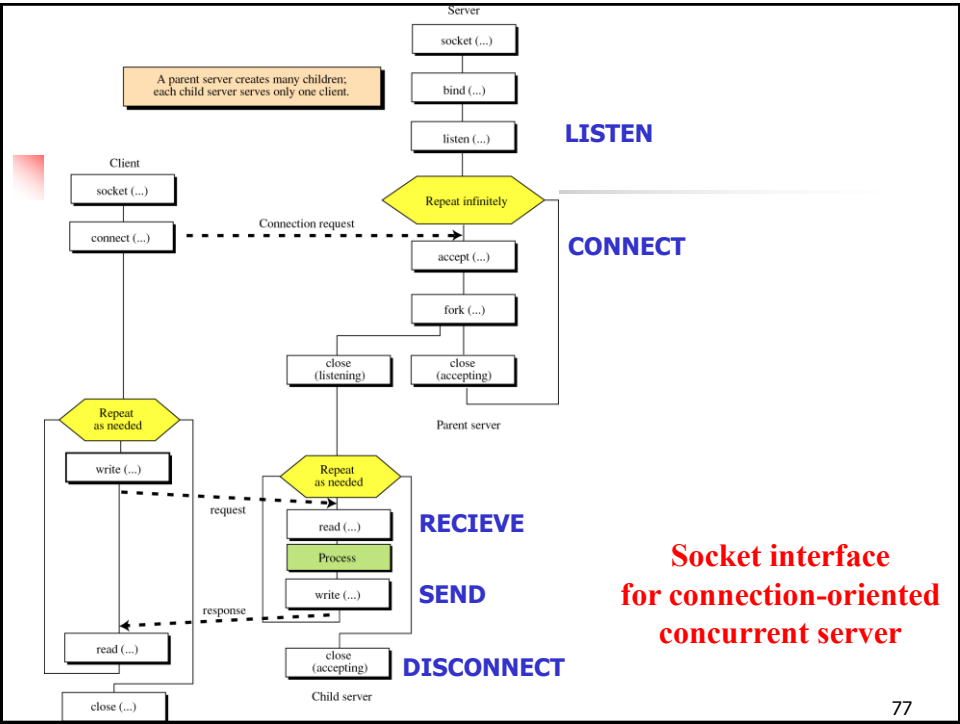






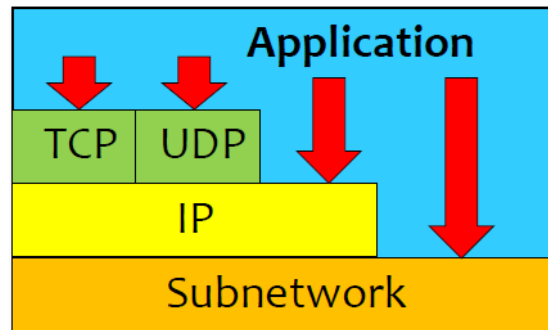
**Example
of a
server
program
with
parent
and
child
processes**

```
.....  
void main (void)  
{  
    ..... ;  
    ..... ;  
    pid_t pid ;  
    for ( ; ; )  
    {  
        Connection from client  
        pid = fork () ;  
        if (pid != 0 )  
        {  
            Code for parent  
        }  
        else  
        {  
            Code for child  
        }  
    }  
}
```



Internet Architecture

- Defined by IETF
- Three main features
 - **Does not imply strict layering.** The application is **free to bypass** the defined transport layers and to directly use IP or other underlying networks



By 清大資工系黃能富教授

79

Example of Socket Programming: An Internet File Server (1)

/* This page contains a client program that can request a file from the server program
* on the next page. The server responds by sending the whole file.
*/

```
#include <sys/types.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <netdb.h>
```

```
#define SERVER_PORT 12345
#define BUF_SIZE 4096
```

```
/* arbitrary, but client & server must agree */
/* block transfer size */
```

```
int main(int argc, char **argv)
{
```

```
    int c, s, bytes;
```

```
    char buf[BUF_SIZE];
```

```
    struct hostent *h;
```

```
    struct sockaddr_in channel;
```

```
/* buffer for incoming file */
```

```
/* info about server */
```

```
/* holds IP address */
```

```
    . . .
```

Client code using sockets

80

Example of Socket Programming: An Internet File Server (2)

```

if (argc != 3) fatal("Usage: client server-name file-name");
h = gethostbyname(argv[1]);          /* look up host's IP address */
if (!h) fatal("gethostbyname failed");

s = socket(PF_INET, SOCK_STREAM, IPPROTO_TCP);
if (s < 0) fatal("socket");
memset(&channel, 0, sizeof(channel));
channel.sin_family = AF_INET;
memcpy(&channel.sin_addr.s_addr, h->h_addr, h->h_length);
channel.sin_port = htons(SERVER_PORT);

c = connect(s, (struct sockaddr *) &channel, sizeof(channel));
if (c < 0) fatal("connect failed");

```

...

Client code using sockets

81

Example of Socket Programming: An Internet File Server (3)

```

c = connect(s, (struct sockaddr *) &channel, sizeof(channel));
if (c < 0) fatal("connect failed");

/* Connection is now established. Send file name including 0 byte at end. */
write(s, argv[2], strlen(argv[2])+1);

/* Go get the file and write it to standard output. */
while (1) {
    bytes = read(s, buf, BUF_SIZE);    /* read from socket */
    if (bytes <= 0) exit(0);           /* check for end of file */
    write(1, buf, bytes);              /* write to standard output */
}

fatal(char *string)
{
    printf("%s\n", string);
    exit(1);
}

```

Client code using sockets

82

Example of Socket Programming: An Internet File Server (4)

```
#include <sys/types.h>                /* This is the server code */
#include <sys/fcntl.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <netdb.h>

#define SERVER_PORT 12345              /* arbitrary, but client & server must agree */
#define BUF_SIZE 4096                 /* block transfer size */
#define QUEUE_SIZE 10

int main(int argc, char *argv[])
{
    int s, b, l, fd, sa, bytes, on = 1;
    char buff[BUF_SIZE];               /* buffer for outgoing file */
    struct sockaddr_in channel;        /* holds IP address */

    . . .
```

Server code

83

Example of Socket Programming: An Internet File Server (5)

```
/* Build address structure to bind to socket. */
memset(&channel, 0, sizeof(channel)); /* zero channel */
channel.sin_family = AF_INET;
channel.sin_addr.s_addr = htonl(INADDR_ANY);
channel.sin_port = htons(SERVER_PORT);

/* Passive open. Wait for connection. */
s = socket(AF_INET, SOCK_STREAM, IPPROTO_TCP); /* create socket */
if (s < 0) fatal("socket failed");
setsockopt(s, SOL_SOCKET, SO_REUSEADDR, (char *) &on, sizeof(on));

b = bind(s, (struct sockaddr *) &channel, sizeof(channel));
if (b < 0) fatal("bind failed");

l = listen(s, QUEUE_SIZE);             /* specify queue size */
if (l < 0) fatal("listen failed");

. . .
```

Server code

84

Example of Socket Programming: An Internet File Server (6)



```
/* Socket is now set up and bound. Wait for connection and process it. */
while (1) {
    sa = accept(s, 0, 0);          /* block for connection request */
    if (sa < 0) fatal("accept failed");

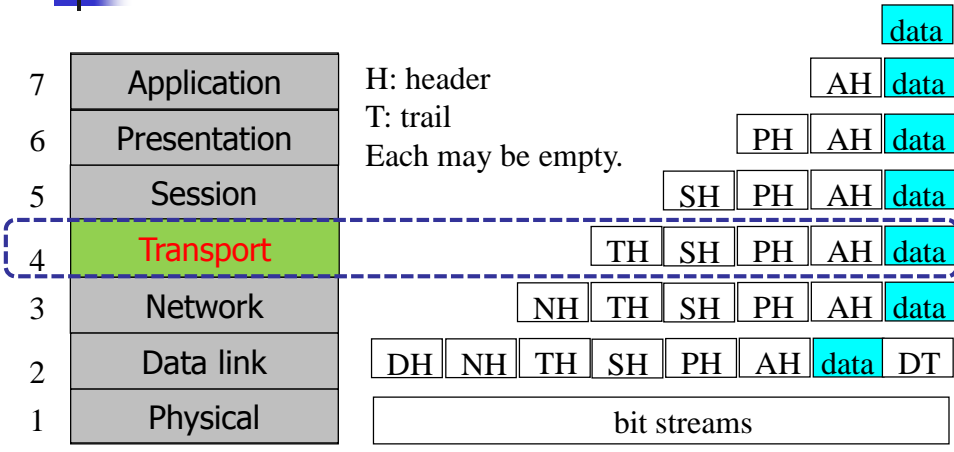
    read(sa, buf, BUF_SIZE);      /* read file name from socket */

    /* Get and return the file. */
    fd = open(buf, O_RDONLY);      /* open the file to be sent back */
    if (fd < 0) fatal("open failed");


    while (1) {
        bytes = read(fd, buf, BUF_SIZE); /* read from file */
        if (bytes <= 0) break;          /* check for end of file */
        write(sa, buf, bytes);          /* write bytes to socket */
    }
    close(fd);                       /* close file */
    close(sa);                       /* close connection */
}
```

Server code

The Transport Layer



OSI Reference Model



End of Chapter 6

Questions?
Thank you!

Computer Networks, Fifth Edition by Andrew Tanenbaum and David Wetherall, © Pearson Education-Prentice Hall, 2011

87