



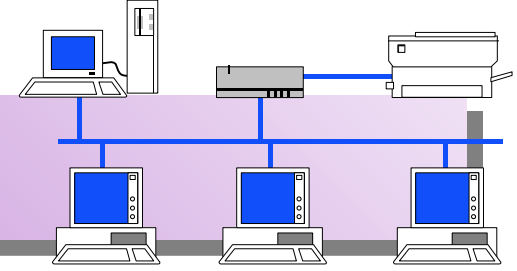
# Unix

## Network Programming (2<sup>nd</sup> Edition)

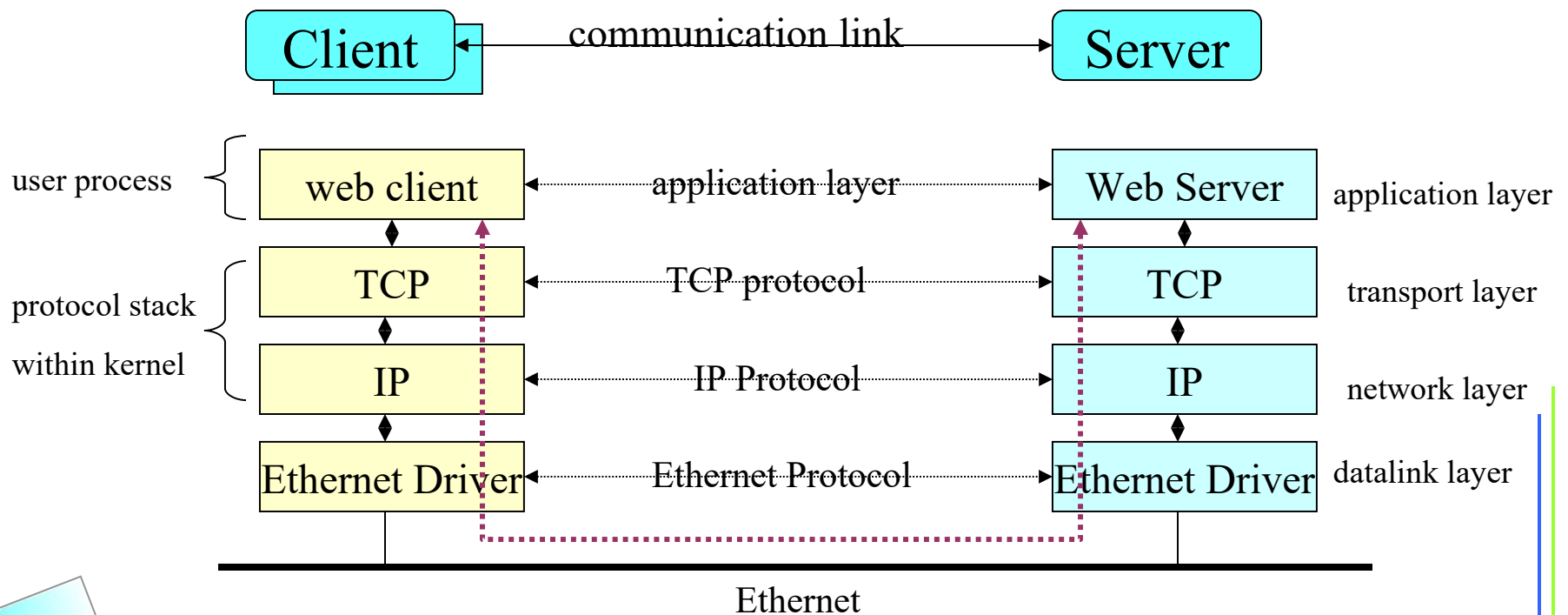
### **Part 1. Introduction and TCP/IP**

#### **Chapter 1. Introduction**

# Introduction



- **Most network application divided into**
  - a **client** and a **server**



# Simple Daytime Client



```

1.  #include "unp.h"
2.
3.  int
4.  main(int argc, char **argv)
5.  {
6.      int      sockfd, n;
7.      char      recvline[MAXLINE + 1];
8.      struct sockaddr_in servaddr;
9.
10.     if (argc != 2)
11.         err_quit("usage: a.out <IPaddress>");
12.
13.     if ( (sockfd = socket(AF_INET, SOCK_STREAM, 0)) < 0)
14.         err_sys("socket error");
15.
16.     bzero(&servaddr, sizeof(servaddr));
17.     servaddr.sin_family = AF_INET;
18.     servaddr.sin_port = htons(13); /* daytime server */
19.     if (inet_pton(AF_INET, argv[1], &servaddr.sin_addr) <= 0)
20.         err_quit("inet_pton error for %s", argv[1]);
21.
22.     if (connect(sockfd, (SA *) &servaddr, sizeof(servaddr)) < 0)
23.         err_sys("connect error");
24.
25.     while ( (n = read(sockfd, recvline, MAXLINE)) > 0) {
26.         recvline[n] = 0; /* null terminate */
27.         if (fputs(recvline, stdout) == EOF)
28.             err_sys("fputs error");
29.     }
30.     if (n < 0)
31.         err_sys("read error");
32.
33.     exit(0);
34. }

```

Create TCP Socket

Specify server's IP  
address and port

Establish connection  
with server

Read and display  
server's reply

# IPv6 Client

```

1.  #include "unp.h"
2.
3.  int
4.  main(int argc, char **argv)
5.  {
6.      int      sockfd, n;
7.      struct sockaddr_in6 servaddr;
8.      char      recvline[MAXLINE + 1];
9.
10.     if (argc != 2)
11.         err_quit("usage: a.out <IPaddress>");
12.
13.     if ( (sockfd = socket(AF_INET6, SOCK_STREAM, 0)) < 0)
14.         err_sys("socket error");
15.
16.     bzero(&servaddr, sizeof(servaddr));
17.     servaddr.sin6_family = AF_INET6;
18.     servaddr.sin6_port = htons(13); /* daytime server */
19.     if (inet_pton(AF_INET6, argv[1], &servaddr.sin6_addr) <= 0)
20.         err_quit("inet_pton error for %s", argv[1]);
21.
22.     if (connect(sockfd, (SA *) &servaddr, sizeof(servaddr)) < 0)
23.         err_sys("connect error");
24.
25.     while ( (n = read(sockfd, recvline, MAXLINE)) > 0) {
26.         recvline[n] = 0; /* null terminate */
27.         if (fputs(recvline, stdout) == EOF)
28.             err_sys("fputs error");
29.     }
30.
31.     if (n < 0)
32.         err_sys("read error");
33.
34.     exit(0);
35. }

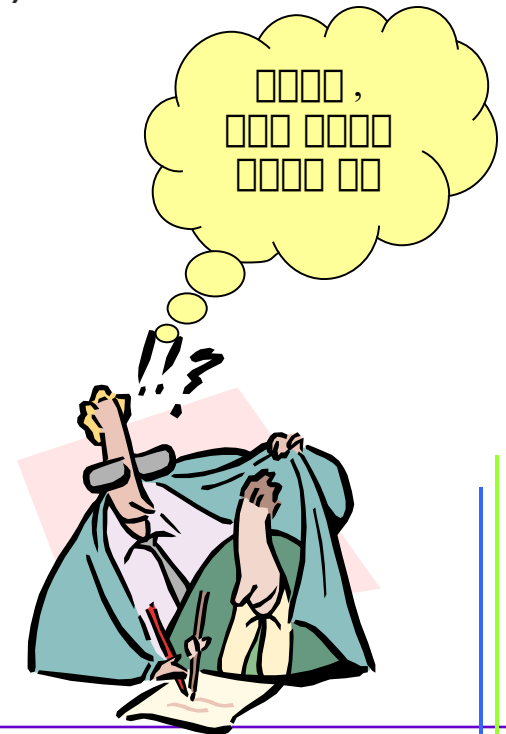
```



# Error Handling (Wrapper Functions)

- **Performs the actual function call, tests the return value, and terminates on an error**
  - `sockfd = Socket(AF_INET, SOCK_STREAM, 0)`

```
int
Socket(int family, int type, int protocol)
{
    int n;
    if (( n=socket(family, type, protocol)) < 0)
        err_sys("Socket error");
    return(n);
}
```



# A Simple Daytime Server

```

1.  #include    "unp.h"
2.  #include    <time.h>
3.  int
4.  main(int argc, char **argv)
5.  {
6.      int                listenfd, connfd;
7.      struct sockaddr_in servaddr;
8.      char                buff[MAXLINE];
9.      time_t             ticks;
10.     listenfd = Socket(AF_INET, SOCK_STREAM, 0);
11.     bzero(&servaddr, sizeof(servaddr));
12.     servaddr.sin_family = AF_INET;
13.     servaddr.sin_addr.s_addr = htonl(INADDR_ANY);
14.     servaddr.sin_port = htons(13); /* daytime
server */
15.     Bind(listenfd, (SA *) &servaddr, sizeof(servaddr));
16.     Listen(listenfd, LISTENQ);
17.     for ( ; ; ) {
18.         connfd = Accept(listenfd, (SA *) NULL,
NULL);
19.         ticks = time(NULL);
20.         sprintf(buff, sizeof(buff), "%.24s\r\n", ctime(&ticks));
21.         Write(connfd, buff, strlen(buff));
22.         Close(connfd);
23.     }
24. }

```

Create Socket

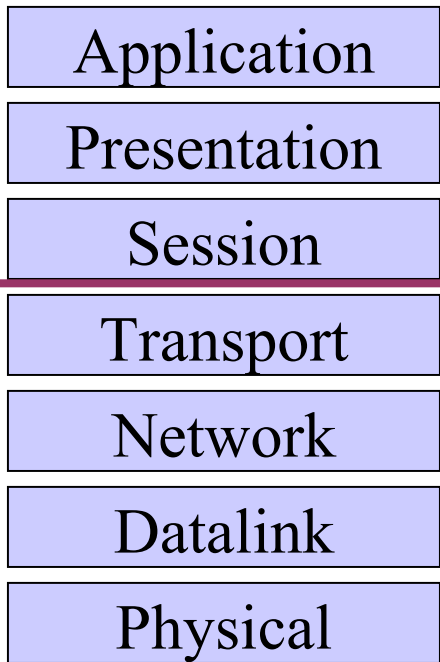
Bind server's well-known port to socket

Convert socket to listening socket

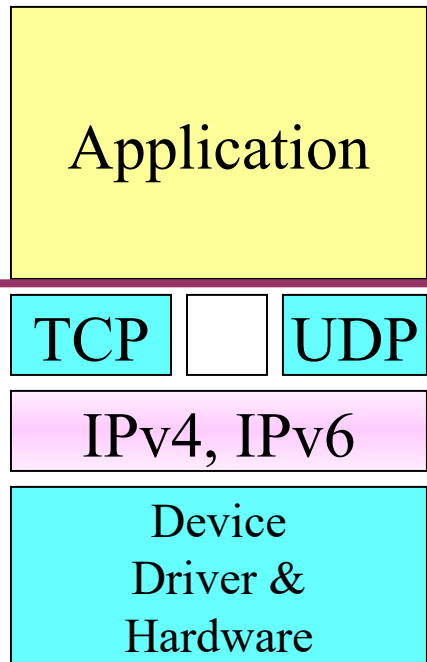
Accept client connection, send reply & close

# OSI Model (Layers in OSI model and Internet protocol suite)

1  
2  
3  
4  
5  
6  
7

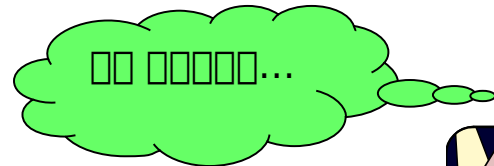
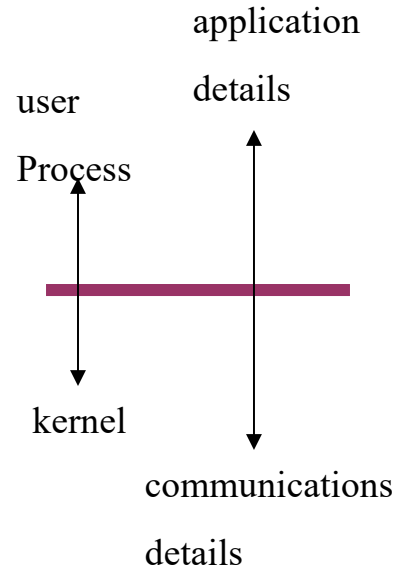


OSI Model



Internet protocol suite

sockets  
XTI

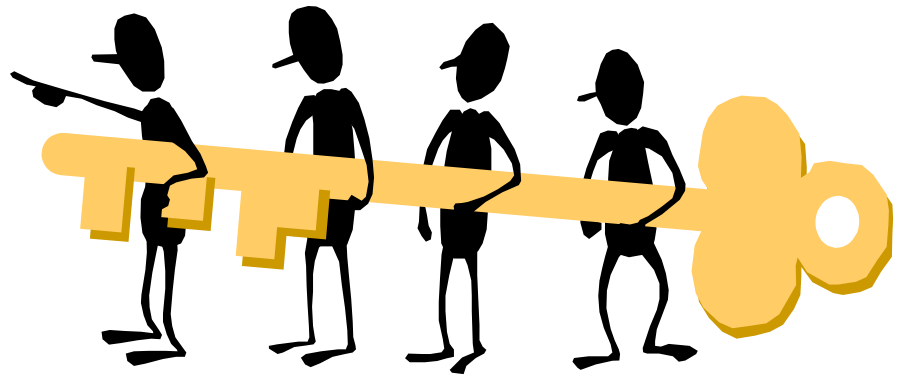


# Unix Standards



- **POSIX**
  - acronym for “Portable Operation System Interface”
  - a family of standards being developed by the IEEE, also adopted as by ISO/IEC
  - <http://www.pasc.org/standing/sd11.html>
- **The Open Group**
  - X/Open Company and Open Software Foundations
  - XPG3(X/Open Portability Guide, Issue 3)
  - <http://www.opengroup.org/public/tech/unix/version2>
- **IETF(Internet Engineering Task Force)**
- **Unix Versions and Portability**
  - focus on Posix.1g





# Unix Network Programming (2<sup>nd</sup> Edition)

## Part 1. Introduction and TCP/IP

### Chapter 2. The Transport Layer: TCP and UDP

# Introduction

- **GOAL**

- to provide enough ***detail to understand*** how to use the protocols from a network programming perspective
- provide ***references to more detailed descriptions*** of the actual design, implementation, and history of the protocols

- **TCP**

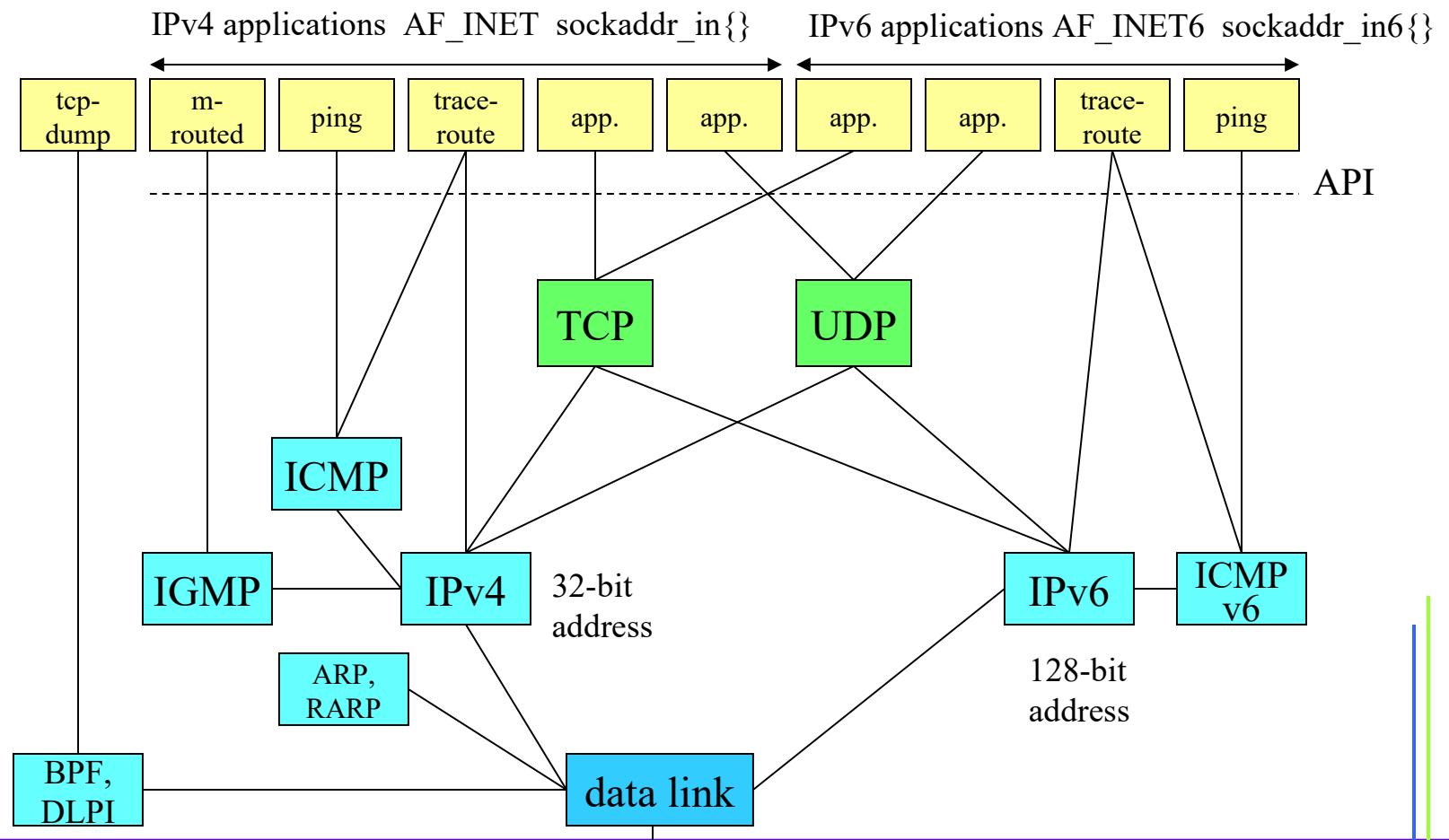
- to write robust clients and server
- sophisticated, byte-stream protocol

- **UDP**

- simple, unreliable, datagram protocol

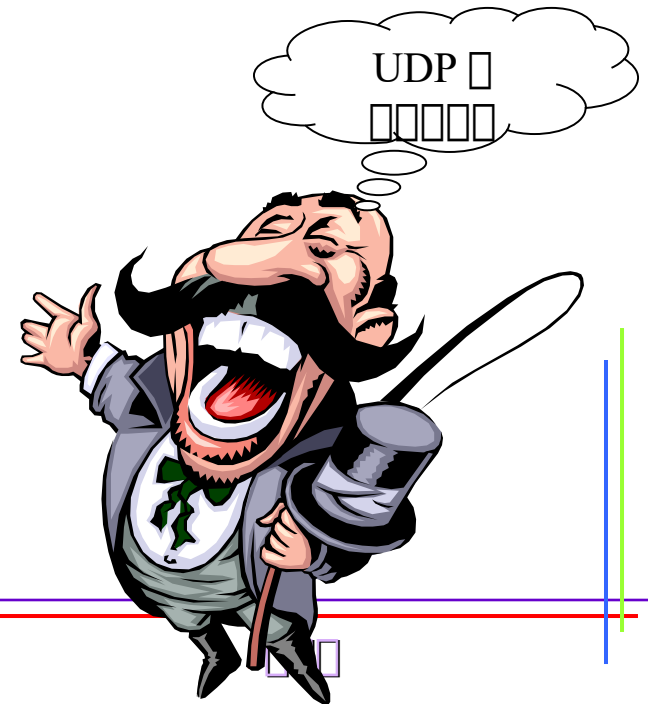


# The Big Picture



# UDP (User Datagram Protocol)

- *simple transport-layer protocol*
- **App writes a datagram to a UDP socket**
- *no guarantee that a UDP datagram ever reaches its final destination*
- *lack of reliability*
- **UDP datagram has a length**
- **a connectionless services**

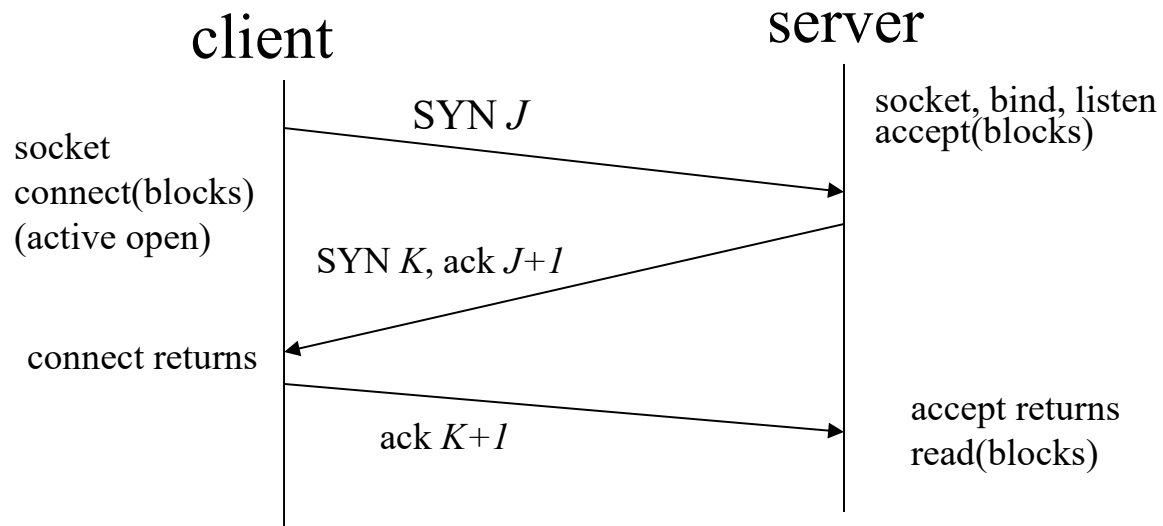


# TCP (Transmission Control Protocol)

- **Provide** *connections* **between clients and servers**
- **Provide** *reliability*
- **also** *sequences* **the data by associating a sequence number with every byte**
- **Provide** *flow control*
- *full-duplex*

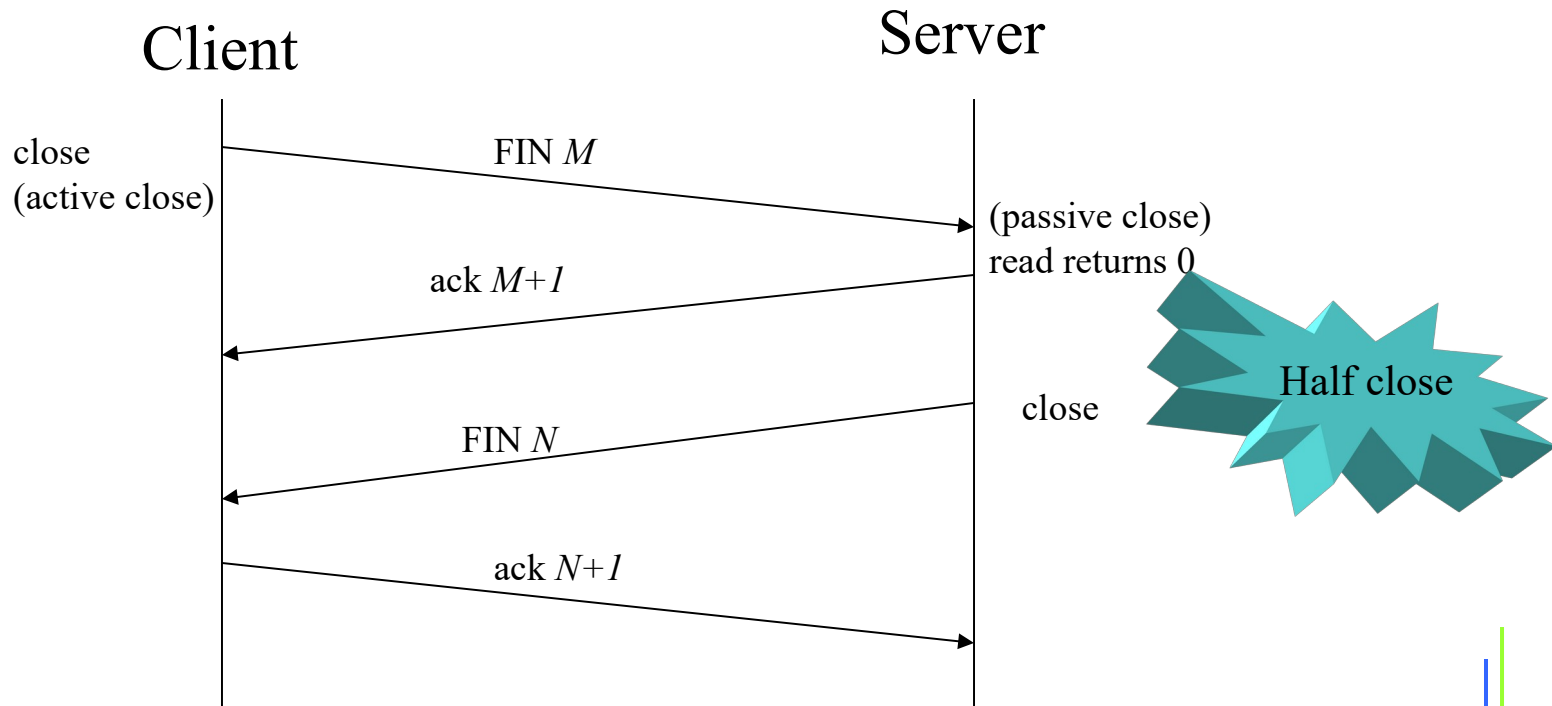


# TCP Connection Establishment



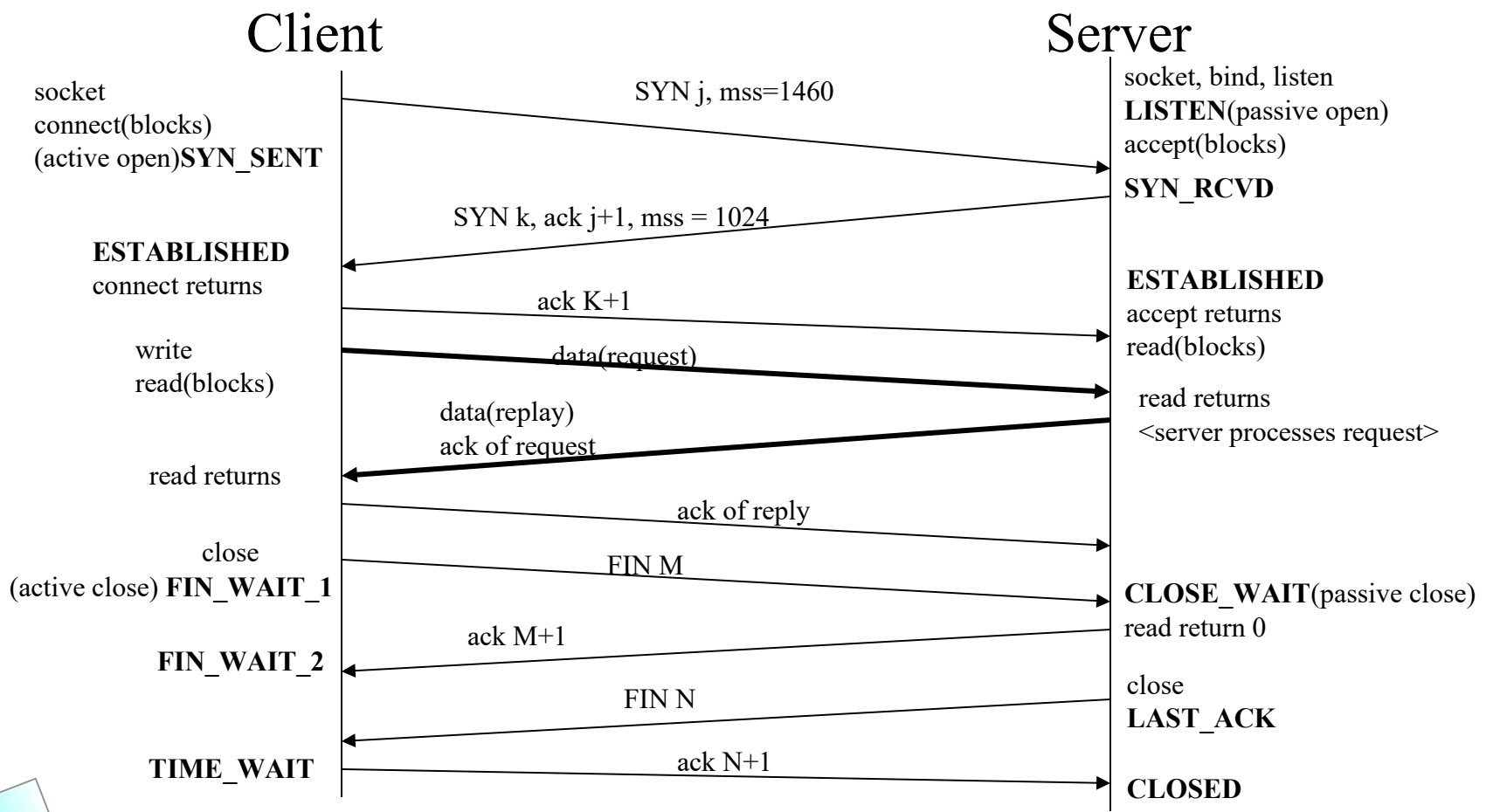
TCP Three-way handshake

# TCP Connection Termination



Packets exchanged when a TCP connection is close

# Watching the Packets

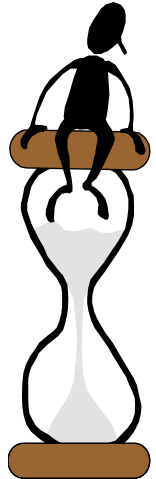




# TIME\_WAIT State & Port Number

- **TIME\_WAIT State**

- MSL(Maximum segment life time)
  - † RFC 1122, about 2 min.
- Two Reasons
  - † to implement TCP's full-duplex connection termination reliably
  - † to allow old duplicate segments to expire in the network



- **Port Number (IANA)**

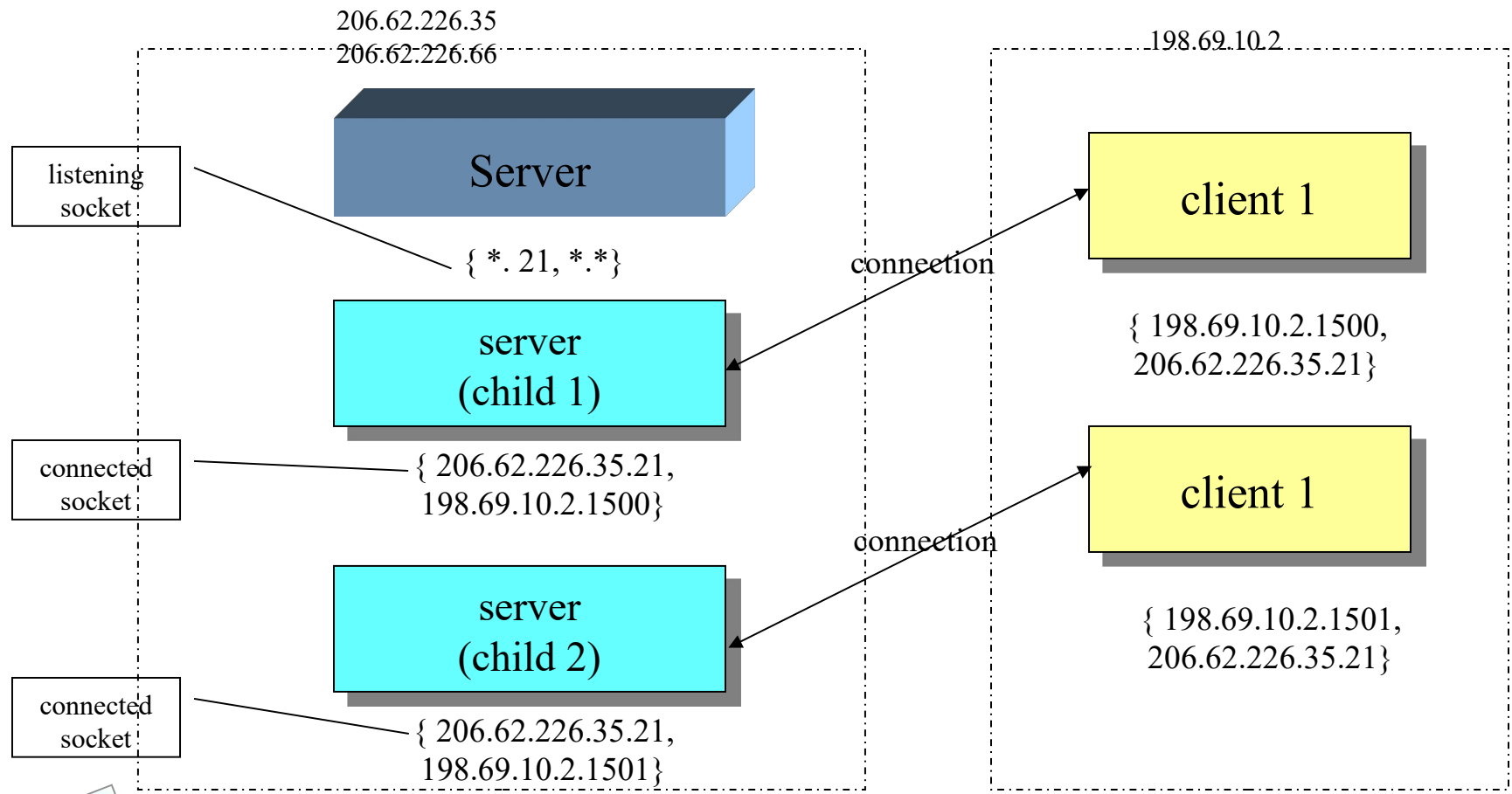
- well-known ports : 0 ~ 1023
- registered port : 1024 ~ 49151
- dynamic or private ports : 49152 ~ 65535

- **Socket Pair**

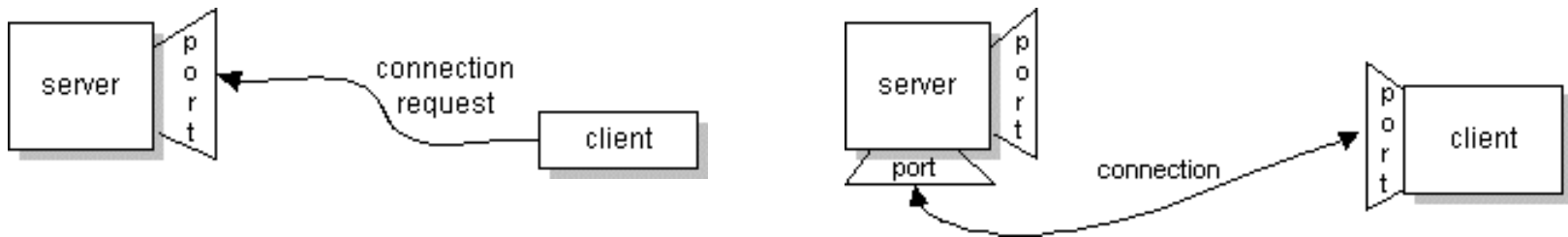
- local IP add., local TCP port, foreign IP add., foreign port



# Concurrent Server



# □□□ (from Java Tutorial)



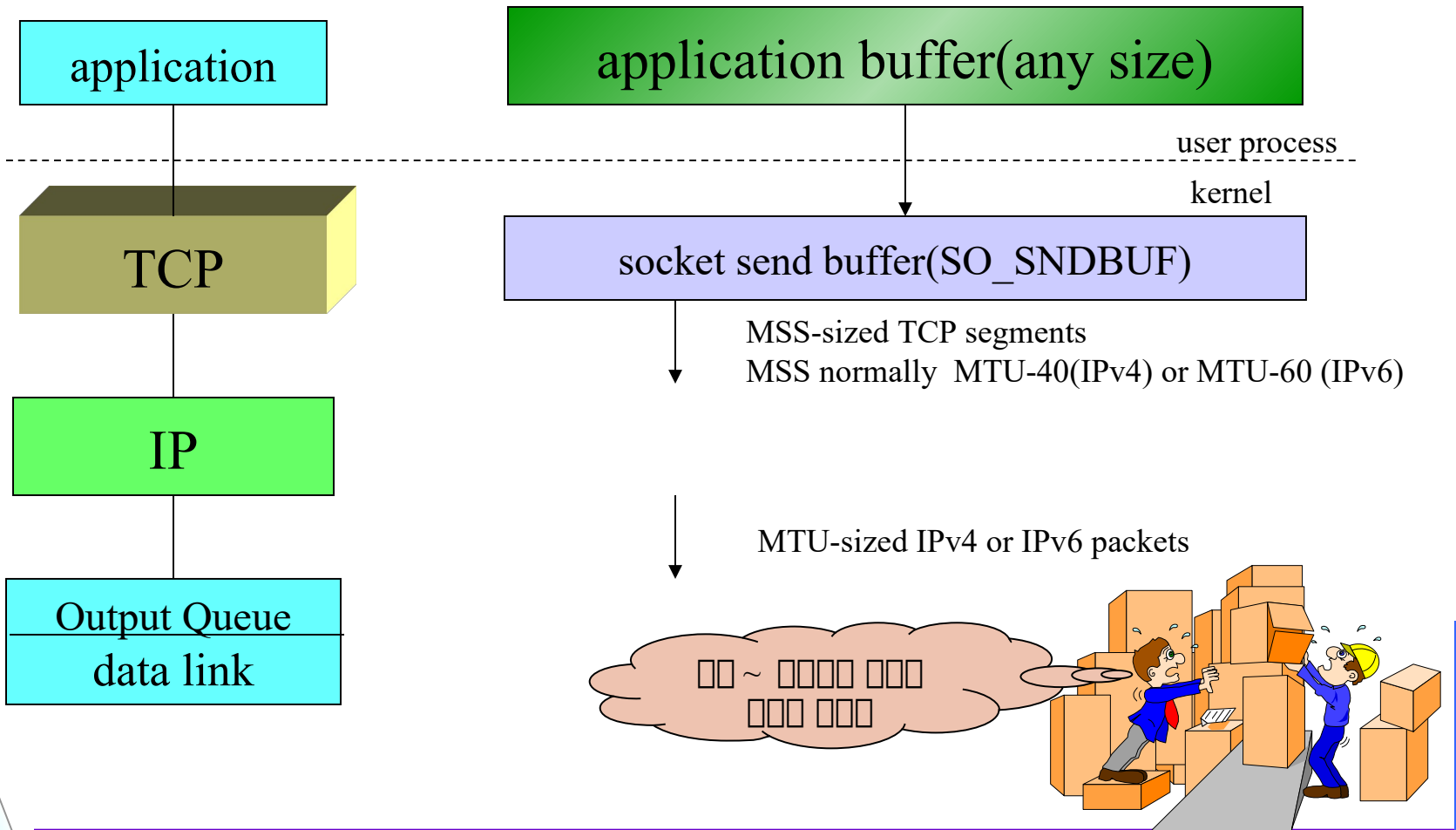
**Normally, a server runs on a specific computer and has a socket that is bound to a specific port number. The server just waits, listening to the socket for a client to make a connection request. On the client-side: The client knows the hostname of the machine on which the server is running and the port number to which the server is connected. To make a connection request, the client tries to rendezvous with the server on the server's machine and port.**

**If everything goes well, the server accepts the connection. Upon acceptance, *the server gets a new socket bound to a different port. It needs a new socket (and consequently a different port number) so that it can continue to listen to the original socket for connection requests while tending to the needs of the connected client.***

**On the client side, if the connection is accepted, a socket is successfully created and the client can use the socket to communicate with the server. Note that the socket on the client side is not bound to the port number used to rendezvous with the server. Rather, the client is assigned a port number local to the machine on which the client is running. The client and server can now communicate by writing to or reading from their sockets.**

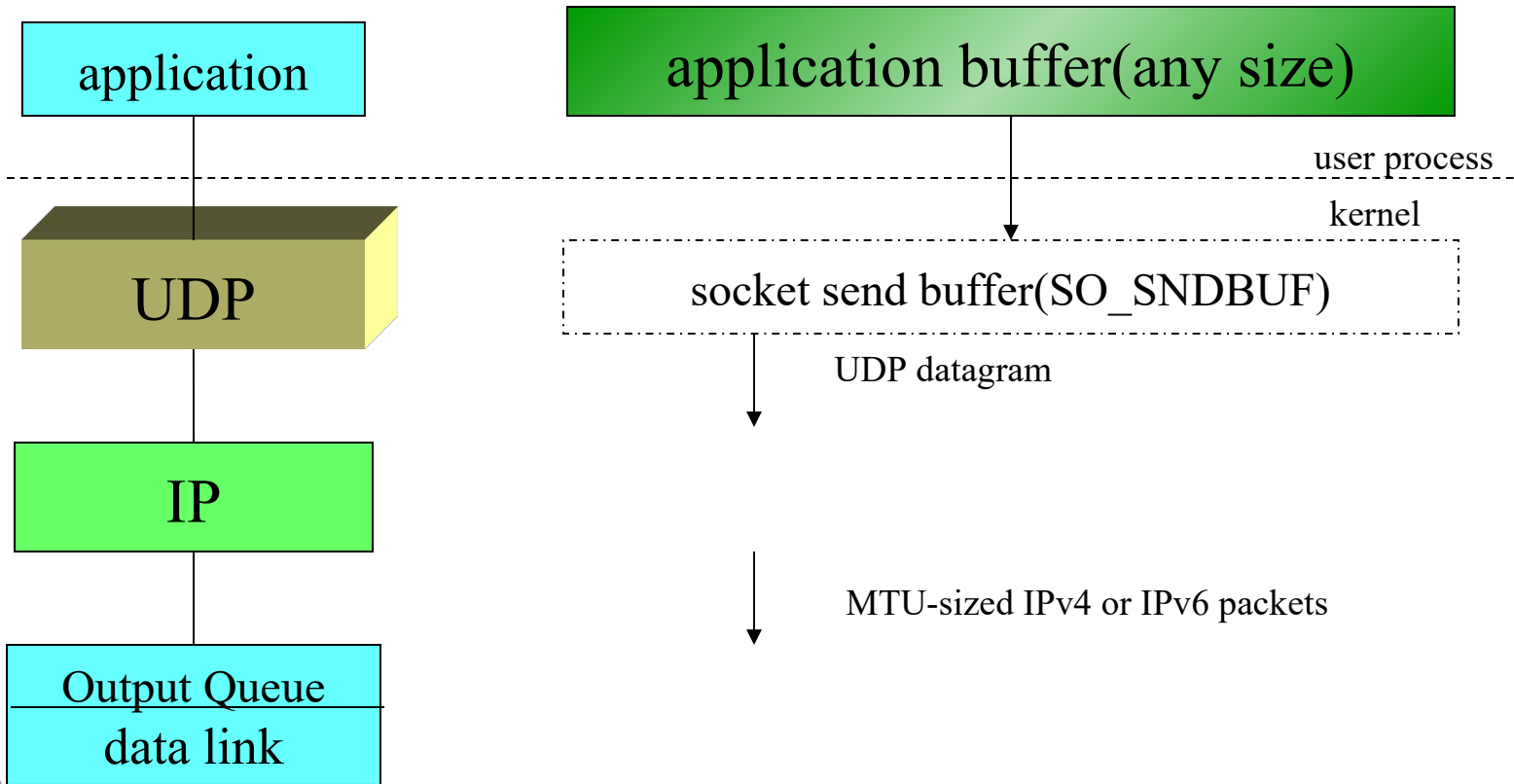
# TCP Output

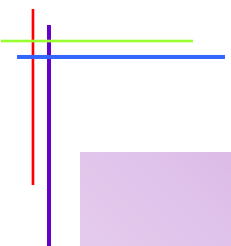
(Steps and buffers involved when application writes to a TCP socket)



# UDP Output

(Steps and buffers involved when application writes to a UDP socket)





# UNIX Network Programming (chapter 3 ~ 5)

# Agenda

- ◆ **Socket Address Structures**
- ◆ **Socket Functions**
- ◆ **TCP Socket**
- ◆ **POSIX □□□□**
- ◆ **Server Source Example**
- ◆ **Client Source Example**

# Socket Address Structures

## Basic

```
struct sockaddr {
    uint8_t    sa_len;        /* for variable socket address structure */
    sa_family_t sa_family;    /* address family */
    char       sa_data[14];   /* protocol-specific address */
};
```

## IPv4

```
struct in_addr {
    in_addr_t    s_addr;     /* 32 bit IPv4 */
};

struct sockaddr_in {
    uint8_t      sin_len;     /* value = 16 */
    sa_family_t  sin_family;
    in_port_t    sin_port;    /* 16 bit */
    struct in_addr sin_addr;   /* 32 bit IPv4 */
    char         sin_zero[8]; /* unused */
};
```

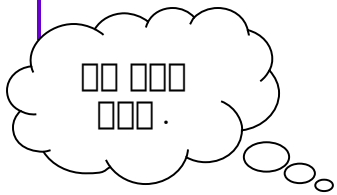
## IPv6

```
struct in6_addr {
    uint8_t      s6_addr[16]; /* 128 bit IPv6 */
};

struct sockaddr_in6 {
    uint8_t      sin6_len;    /* 16 */
    sa_family_t  sin6_family;
    in_port_t    sin6_port;
    uint32_t     sin6_flowinfo; /* priority & flow */
    struct in6_addr sin6_addr; /* 128 bit IPv6 */
};
```



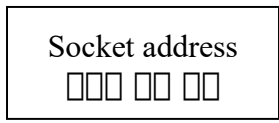
# Socket Functions (1)



```

int socket( int family, int type, int protocol );
int bind( int sockfd, struct sockaddr *addrptr, int addrlen );
int connect( int sockfd, struct sockaddr *addrptr, int addrlen );
int listen( int sockfd, int backlog );
int accept( int sockfd, struct sockaddr *addrptr, int *addrlen );
int close( int sockfd );

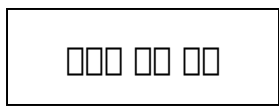
```



```

int getsocketname( int sockfd, struct sockaddr *loacladdr, socklen_t *addrlen );
int getpeername ( int sockfd, struct sockaddr *peeraddr, socklen_t *addrlen );

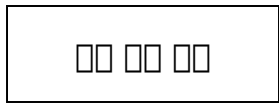
```



```

void bzero( void #dest, size_t nbytes );
void bcopy( const void *src, void *dest, size_t nbytes );
int bcmp( const void *ptr1, const void *ptr2, size_t nbytes );

```



```

int inet_aton( const char *strptr, struct in_addr *addrptr );
char *inet_ntoa( struct in_addr inaddr );
int inet_pton( int family, const char *strptr, void *addrptr );
const char *inet_ntop( int family, const void *addrptr, char *strptr, size_t len );

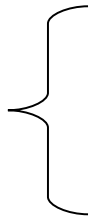
```

IPv4

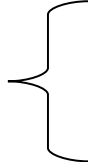
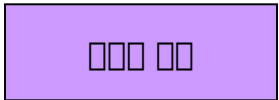
IPv4 and IPv6



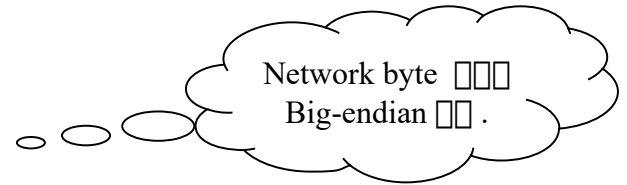
# Socket Functions (2)



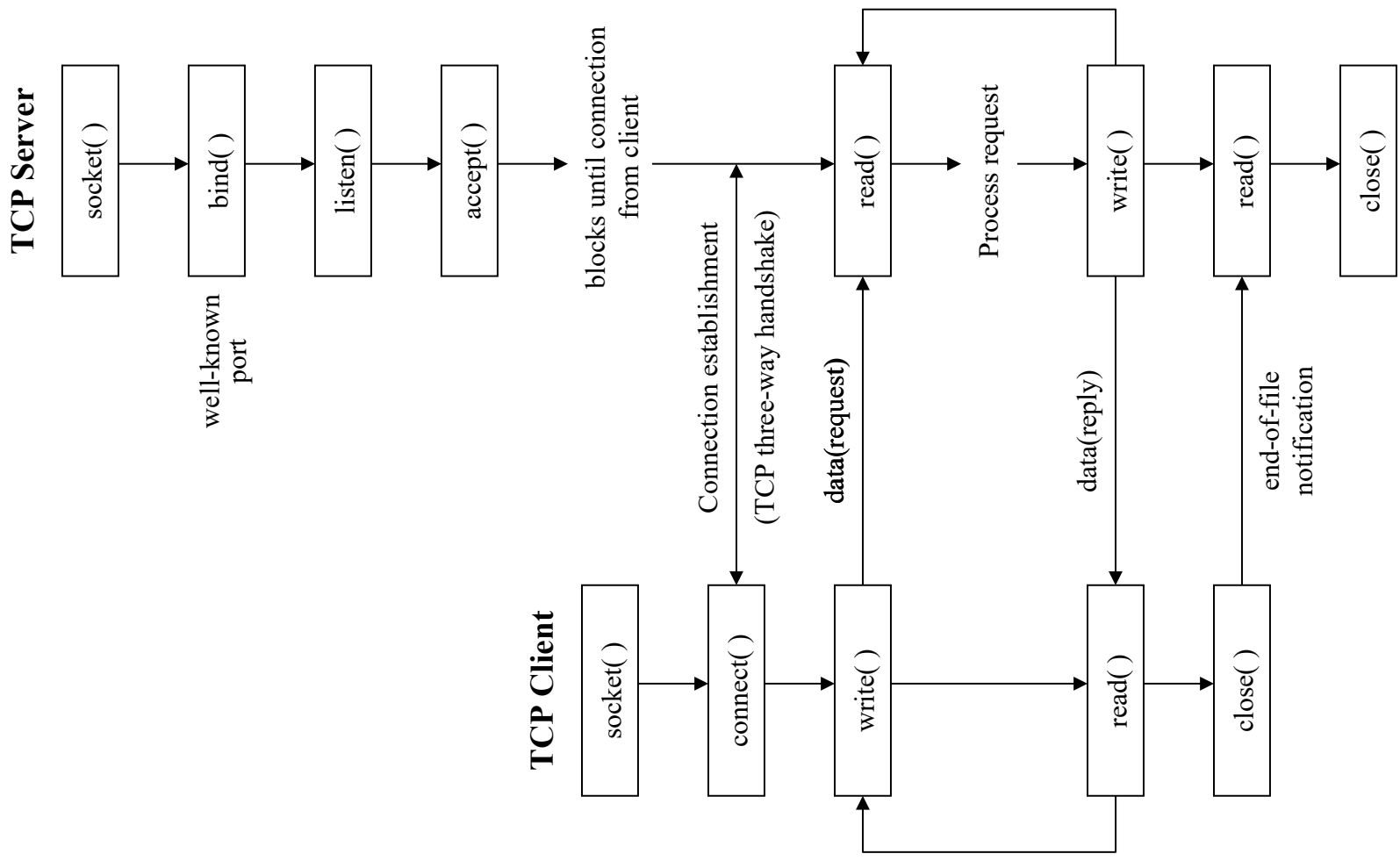
```
uint32_t htonl( uint32_t hostlong );
uint32_t ntohl( uint32_t netlong );
uint16_t htons( uint16_t hostshort );
uint16_t ntohs( uint16_t hostshort );
```



```
ssize_t readn( int filedes, void *buff, size_t nbytes );
ssize_t writen( int filedes, const void *buff, size_t nbytes );
ssize_t readline( int filedes, void buff, size_t maxlen );
```



# TCP Socket



# socket()

**int socket( int family, int type, int protocol );**

family	Description
AF_INET	IPv4 protocols
AF_INET6	IPv6 protocols
AF_LOCAL	Unix domain protocols
AF_ROUTE	Routing sockets
AF_KEY	Key sockets

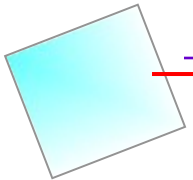
< Address family >

type	Description
SOCK_STREAM	stream socket
SOCK_DGRAM	datagram socket
SOCK_RAW	raw socket

< Socket Type >

	AF_INET	AF_INET6	AF_LOCAL	AF_ROUTE	AF_KEY
SOCK_STREAM	TCP	TCP	Yes		
SOCK_DGRAM	UDP	UDP	Yes		
SOCK_RAW	IPv4	IPv6		Yes	Yes

- ☞ protocol argument ☐ raw socket ☐ 0 ☐ ☐ .
- ☞ key socket : kernel ☐ ☐ ☐ key table ☐ ☐ interface ☐ ☐ .
- ☞ raw socket : TCP ☐ UDP ☐ header ☐ ☐ IP header ☐ ☐ ☐ ☐ .



# connect(), bind()

**int connect( int sockfd, const struct sockaddr \*servaddr, socklen\_t addrlen );**

- ☞ client → remote socket → server → → → → → → → →
- ☞ three-way handshake → →

**int bind( int sockfd, const struct sockaddr \*myaddr, socklen\_t addrlen );**

- ☞ → → → socket → myaddr → → → → → → → →
- ☞ → → server → → → → → service port → → → → → → → →
- ☞ client → → bind() → → → → → → → → kernel → → → → → port → → → → → → → →  
→ → → → binding → → → → .

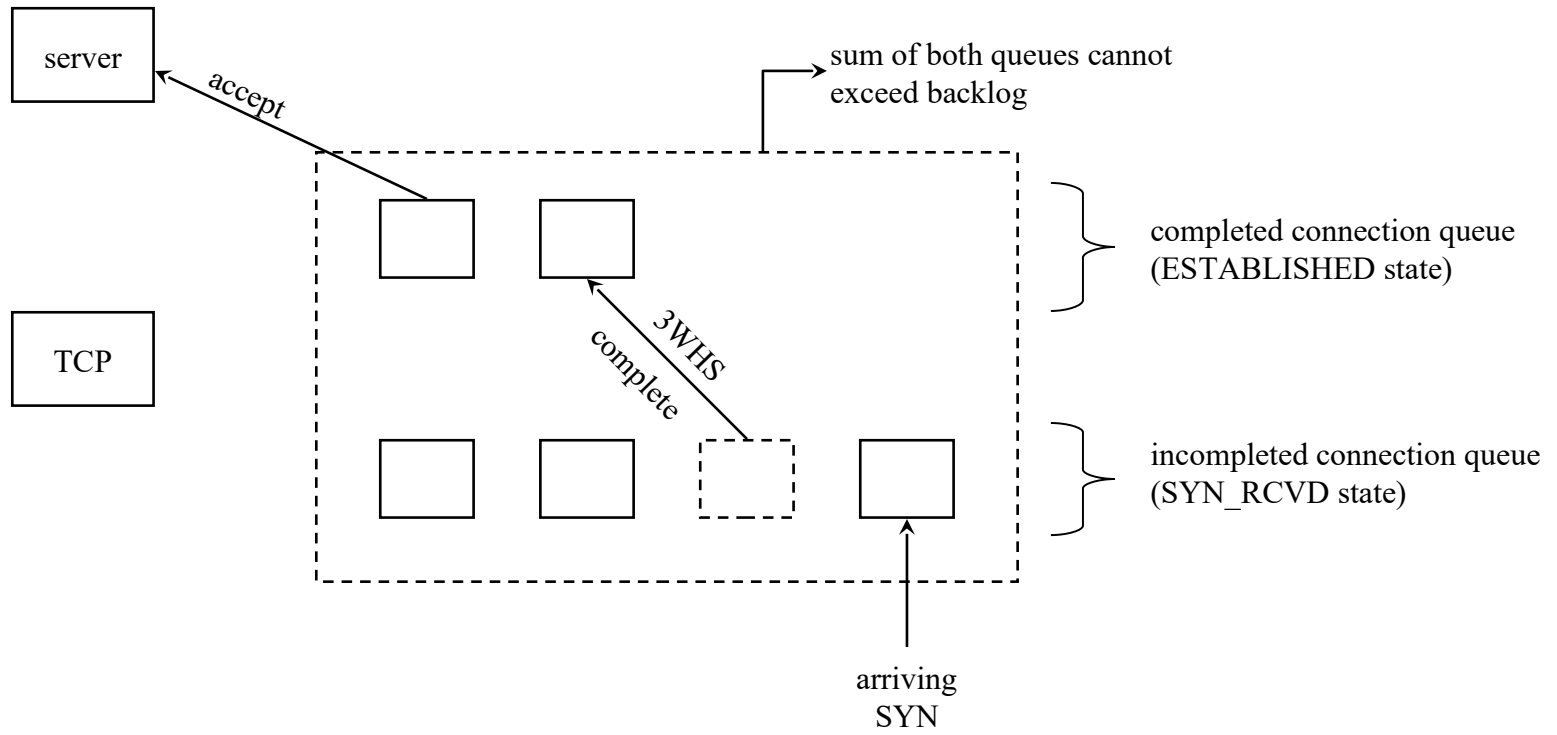
Process specifies		Result
IP address	port	
wildcard	0	kernel chooses IP address and port
wildcard	nonzero	kernel chooses IP address, process specifies port
local IP address	0	process specifies IP address, kernel chooses port
local IP address	nonzero	process specifies IP address and port

# listen()

```
int listen( int sockfd, int backlog );
```

☞ server □ □□□ □□

☞ server program □ □□ □□□ server □ □□□□ socket □ □□□□ □□□ listen() system call □ □□□□ kernel □ □□□ .



# accept(), close()

```
int accept( int sockfd, struct sockaddr *cliaddr, socklen_t *addrlen );
```

☞ incoming connection request □ □

☞ cliaddr □ accept □ client □ □□□ □□□□ .

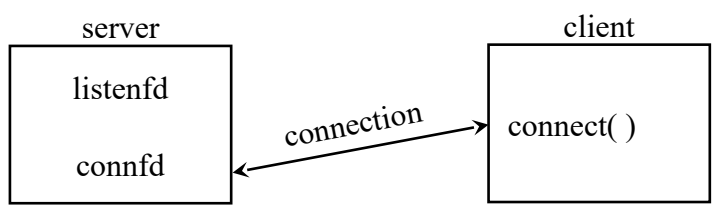
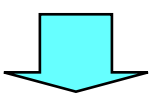
```
int close( int sockfd );
```

☞ socket □ □□□□ □□ □□□□ □□□□ .

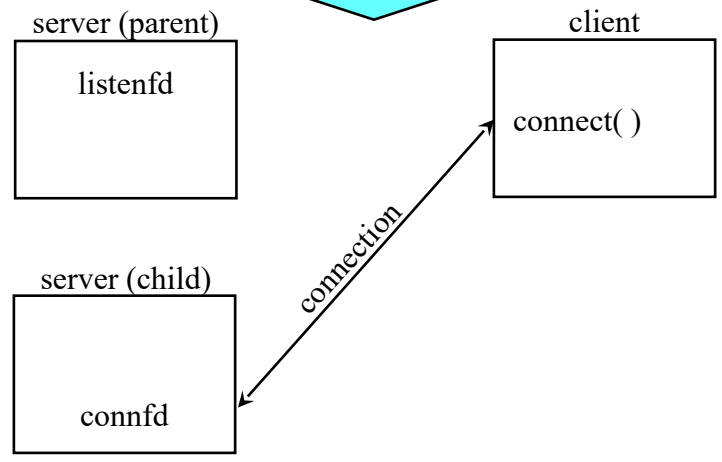
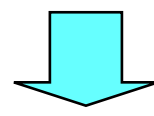
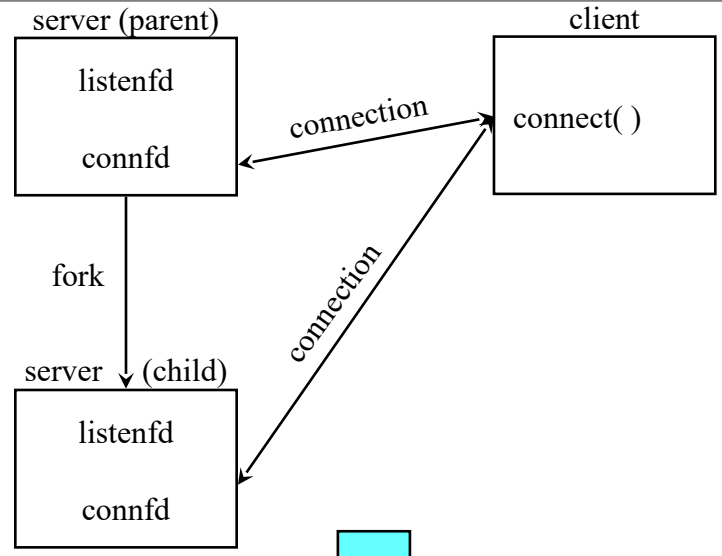
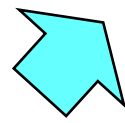
# fork ()

```
pid_t fork( void );
```

process 0 000 0000 0000 00 00 00 000000  
00 .



\* listenfd : listening socket  
connfd : connected socket





# POSIX

☞ `sigaction` disposition `act` , `oact` action `act` .

☞ `sigaction` disposition `act` `sigaction` `act` .

1) `SIGKILL` `SIGSTOP` `act` `act` .

2) `SIG_IGN` `act` `act` . `SIGKILL` `SIGSTOP` `act` .

3) `SIG_DFL` `act` default disposition `act` . `act` `act` process `act` , `act` process `act` memory `act` working directory `act` .

```
typedef void Sigfunc( int );
Sigfunc *signal( int signo, Sigfunc *func );
{
    struct sigaction act, oact;

    act.sa_handler = func;
    sigemptyset( &act.sa_mask );
    act.sa_flags = 0;

    if ( sigaction( signo, &act, &oact ) < 0 )
        return ( SIG_ERR );
    return ( oact.sa_handler );
}
signo = signal ( SIG_IGN , SIG_DFL
```

# SIGCHLD signal

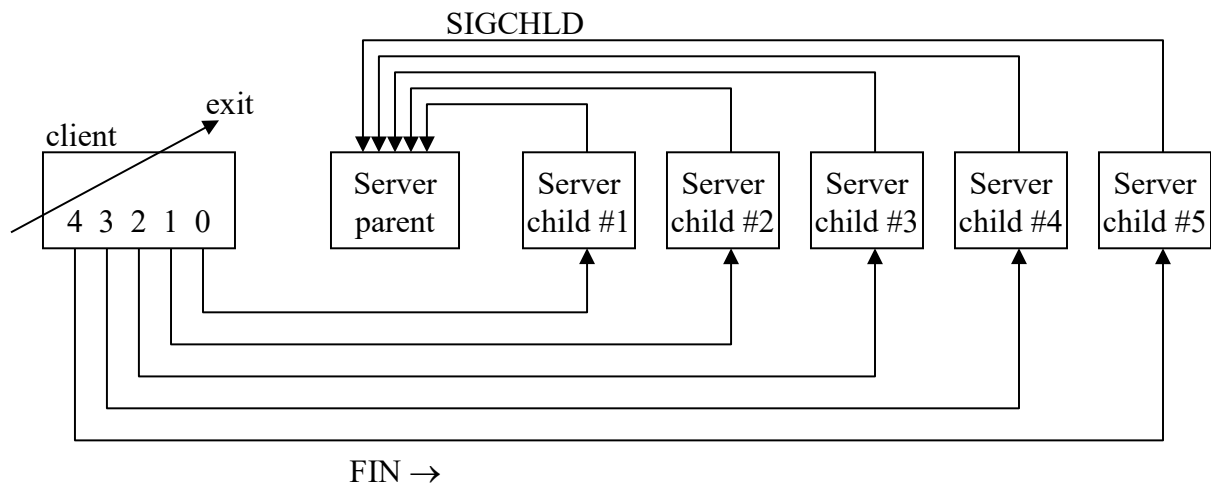
Process is terminated by kernel process parent is notified.

zombie child is notified wait server is notified waitpid blocking.

```
pid_t wait( int *statloc );
```

```
pid_t waitpid( pid_t pid, int *statloc, int options );
```

option "WNOHANG" kernel child process blocking, child process waitpid blocking.



< wait >

PID	TT	STAT	TIME	COMMAND
21282	p1	S	0:00.09	./tcpserv03
21284	p1	Z	0:00.00	(tcpserv03)
21285	p1	Z	0:00.00	(tcpserv03)
21286	p1	Z	0:00.00	(tcpserv03)
21287	p1	Z	0:00.00	(tcpserv03)



# SIGPIPE signal

- ☞ sever □ □ □ □ □ □ server TCP □ client □ □ □ data □ □ □ process □ □ □ □ □ □ RST □ □ □ □ .
- ☞ process □ RST □ □ □ socket □ data □ □ □ write □ EPIPE □ □ □ □ □ .
- ☞ □ □ □ □ □ □ □ □ □ □ “ server terminated prematurely” □ □ □ error message □ □ □ □ □ □ ..

# Server source example (1)

```

int main( int argc, char **argv )
{
    int listenfd, connfd;
    pid_t childpid;
    socklen_t clien;
    struct sockaddr_in cliaddr, servaddr;
    void sig_chld( int );

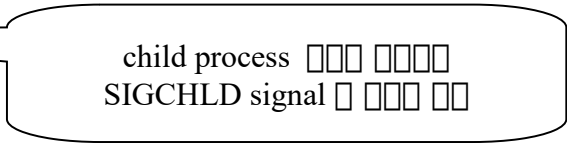
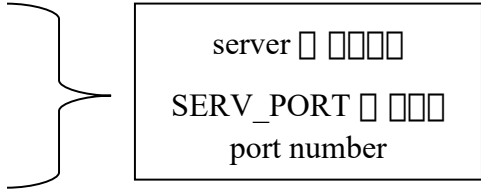
    listenfd = Socket( AF_INET, SOCK_STREAM, 0 );

    bzero( &servaddr, sizeof(servaddr) );
    servaddr.sin_family = AF_INET;
    servaddr.sin_addr.s_addr = htonl( INADDR_ANY );
    servaddr.sin_port = htons( SERV_PORT );

    Bind( listenfd, (SA *)&servaddr, sizeof(servaddr) );
    Listen( listenfd, LISTENQ );
    Signal( SIGCHLD, sig_chld );

    for ( ; ; ) {
        clien = sizeof( cliaddr );

```



# Server source example (2)

```

if ( (connfd = accept( listenfd, (SA *)&cliaddr, &clilen ) ) < 0 ) {
    if ( errno == EINTR )
        continue;
    else
        err_sys( "accept error" );
}
if ( (childpid = Fork() ) == 0 ) {
    Close( listenfd );          /* child closes listening socket */
    str_echo( connfd );        /* do it all */
    exit(0);
}
Close( connfd );              /* parent closes connected socket */
}
}

void sig_chld( int signo )
{
    pid_t pid;
    int stat;

    while ( ( pid = waitpid( -1, &stat, WNOHANG ) ) > 0 )
        printf( "child %d terminated\n", pid );
    return;
}

```

slow system call( `accept` )  
`EINTR`( `int` `errnum` )  
`int` `errcode`

# Client source example

```

int main( int argc, char **argv )
{
    int i, sockfd[5];
    struct sockaddr_in seraddr;

    if ( argc != 2 )
        err_quit( "usage: tcpcli <IPaddress>" );

    for ( i=0; i<5; i++ ) {
        sockfd[i] = Socket( AF_INET, SOCK_STREAM, 0 );

        bzero( &servaddr, sizeof(servaddr) );
        servaddr.sin_family = AF_INET;
        servaddr.sin_port = htons( SERV_PORT );
        Inet_pton( AF_INET, argv[1], &servaddr.sin_addr );

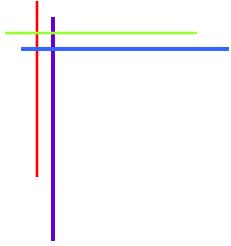
        Connect( sockfd[i], (SA *)&servaddr, sizeof( servaddr ) );
    }

    str_cli( stdin, sockfd[0] );          /* do it all */

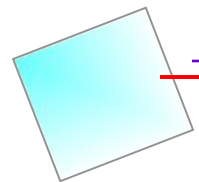
    exit(0);
}

```

□□□ server □□□  
□□□□ .



# Chapter 6 *I/O Multiplexing: The select and poll Functions*



# Introduction

- **I/O Multiplexing** □ **select, poll, pselect** □□□ □□ □□□
- **I/O Multiplexing** □ □□□□ □
  - client □ □□ □□ descriptor( □□ □□ , interactive input □ network socket) □ □□□□ □□
  - client □ □□□ □□ □□ socket □ □□□□ □□ □□
  - TCP □□□□ listening socket □ connected socket □ □□□□ □□
  - □□□□ TCP □ UDP □ □□□□ □□□□□ □□
  - □□□□ □□□□ □□ □□ □□□□□ □□□□□ □□ □□ ( □ : inetd)



# I/O Models

- **blocking I/O**
- **nonblocking I/O**
- **I/O multiplexing (select and poll)**
- **signal-driven I/O (SIGIO)**
- **asynchronous I/O (Posix.1 aio\_ functions)**

✱ □ I/O □□□ □□ □□□ **Text** □ **Fig 6.1~6.6** □ □□

# select Function

- **Allows the process to instruct the kernel to wait for any one of multiple events to occur and to wake up the process only when one or more of these events occurs or when a specified amount of time has passed.**

```
#include <sys/select.h>
#include <sys/time.h>

int select(int maxfdp1, fd_set *readset, fd_set *writerset, fd_set *exceptset,
           const struct timeval *timeout);
```

Returns: positive count of ready descriptors, 0 on timeout, -1 on error

1. Wait forever: timeout = NULL
2. Wait up to a fixed amount time: timeout specified in the timeval structure
3. Don't wait at all: timeout = 0

# select Function (cont'd)

- **fd\_set**
  - void FD\_ZERO(fd\_set \*fdset);
  - void FD\_SET(int fd, fd\_set \*fdset);
  - void FD\_CLR(int fd, fd\_set \*fdset);
  - int FD\_ISSET(int fd, fd\_set \*fdset);
- **Under What Conditions Is a Descriptor Ready?**
  - Ready for reading if...
    - † receive buffer  $\geq$  low-water mark for receive buf.  
(`read` return  $> 0$ )
    - † read-half `read` return  $= 0$
    - † listening socket, `accept` return  $\neq 0$
    - † pending error (`read` return  $= -1$ )

## select Function (cont'd)

- Ready for writing if...
  - † socket send buffer  $\geq$  low-water mark for send buf. (2048) bytes, UDP (write  $> 0$ )
  - † write-half (write  $\neq$  SIGPIPE)
  - † pending error (write  $= -1$ )
- Exception condition if...
  - † Out-of-band data (Chap. 21 ..)

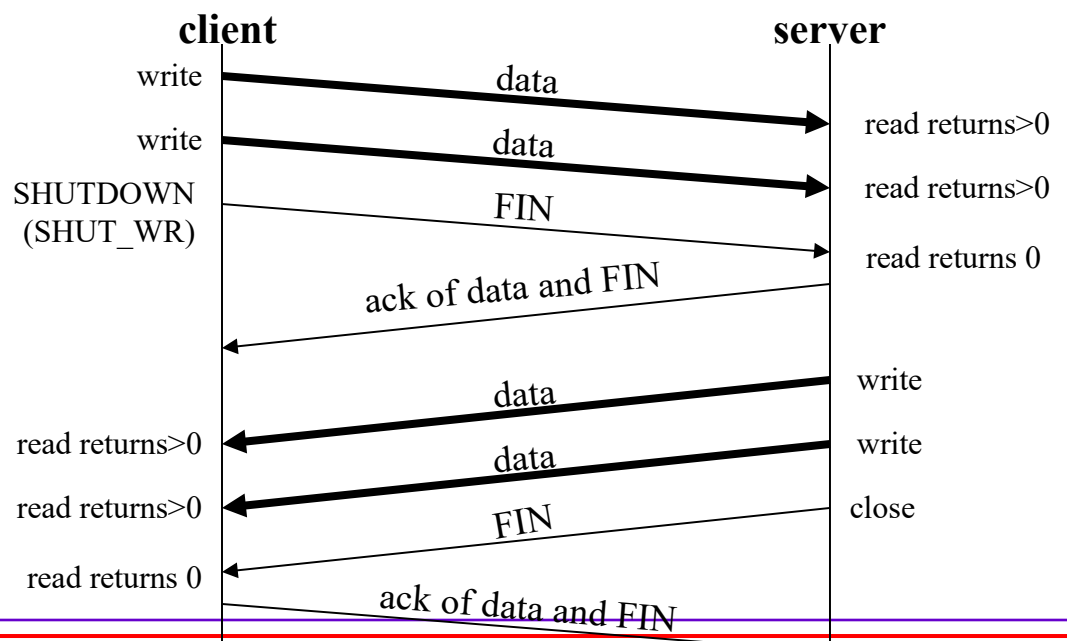
# shutdown Function

```
#include <sys/socket.h>
```

```
int shutdown (int sockfd, int howto);
```

Returns: 0 if OK, -1 on error

✳ **howto** : SHUT\_RD / SHUT\_WR / SHUT\_RDWR



## Example (str\_cli Function–Revisited)

```
1 #include "unp.h"
2
3 void
4 str_cli(FILE *fp, int sockfd)
5 {
6     int     maxfdp1, stdineof;
7     fd_set  rset;
8     char    sendline[MAXLINE], recvline[MAXLINE];
9
10    stdineof = 0;
11    FD_ZERO(&rset);
12    for ( ; ; ) {
13        if (stdineof == 0)
14            FD_SET(fileno(fp), &rset);
15        FD_SET(sockfd, &rset);
16        maxfdp1 = max(fileno(fp), sockfd) + 1;
17        Select(maxfdp1, &rset, NULL, NULL, NULL);
18
19        if (FD_ISSET(sockfd, &rset)) { /* socket is readable */
20            if (Readline(sockfd, recvline, MAXLINE) == 0) {
21                if (stdineof == 1)
22                    return; /* normal termination */
23                else
24                    err_quit("str_cli: server terminated prematurely");
25            }
26        }
27    }
28 }
```

## Example (str\_cli Function-Revisited)

```
26
27     Fputs(recvline, stdout);
28 }
29
30 if (FD_ISSET(fileno(fp), &rset)) { /* input is readable */
31     if (Fgets(sendline, MAXLINE, fp) == NULL) {
32         stdineof = 1;
33         Shutdown(sockfd, SHUT_WR); /* send FIN */
34         FD_CLR(fileno(fp), &rset);
35         continue;
36     }
37
38     Writen(sockfd, sendline, strlen(sendline));
39 }
40 }
41 }
```

# pselect and poll Functions

```
#include <sys/select.h>
#include <signal.h>
#include <time.h>
```

```
int pselect (int maxfdp1, fd_set *readset, fd_set *writerset, fd_set *exceptset, const struct timespec
             *timeout, const sigset_t *sigmask);
```

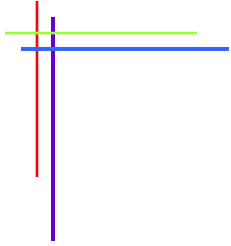
Returns: count of ready descriptors, 0 on timeout, -1 on error

```
#include <poll.h>
```

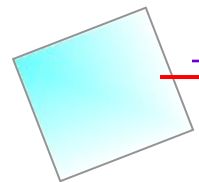
```
int poll (struct pollfd *fdarray, unsigned long nfds, int timeout);
```

Returns: count of ready descriptors, 0 on timeout, -1 on error





# Chapter 7 *Socket Options*



# Socket Option □ get/set □□ □□

- **getsockopt and setsockopt functions**

```
#include <sys/socket.h>
```

```
int getsockopt(int sockfd, int level, int optname, void *optval, socklen_t *optlen);
```

```
int setsockopt(int sockfd, int level, int optname, const void *optval, socklen_t *optlen);
```

Both return: 0 if OK, -1 on error

- **fnctl function**
- **ioctl function**

✱ **getsockopt, setsockopt** □□□ □□ □□□ □□ □□ □□  
**(Fig. 7.1 of text)**

# Generic Socket Options

- **SO\_BROADCAST**

- enable/disable the ability of the process to send broadcast messages
- supported for only datagram sockets
- broadcast message `sockaddr_in` `sockaddr_in` (Ethernet, Token ring `0`)  
`0` `0`
- `0` `0` `0` `0` `0`, destination address `0` broadcast address `0` `EADDRRESOLVE` `0`

- **SO\_DEBUG**

- TCP `0` `0`
- `0` TCP `0` `0` / `0` `0` `0` `0` `0` `0` circular buffer `0` `0`
- trpt program `0` `0` `0` `0`...

# Generic Socket Options (cont'd)

- **SO\_DONTROUTE**

- 0: 0x00000000, 1: 0x00000001
- destination is point-to-point link, shared network returns ENETUNREACH
- datagram MSG\_DONTROUTE flag send/sendto/sendmsg

- **SO\_ERROR**

- pending error (Exxx) returned
- pending error 2
  - † select call block -1 return
  - † signal-driven I/O SIGIO
- getsockopt socket pending error, 1 so\_error 0



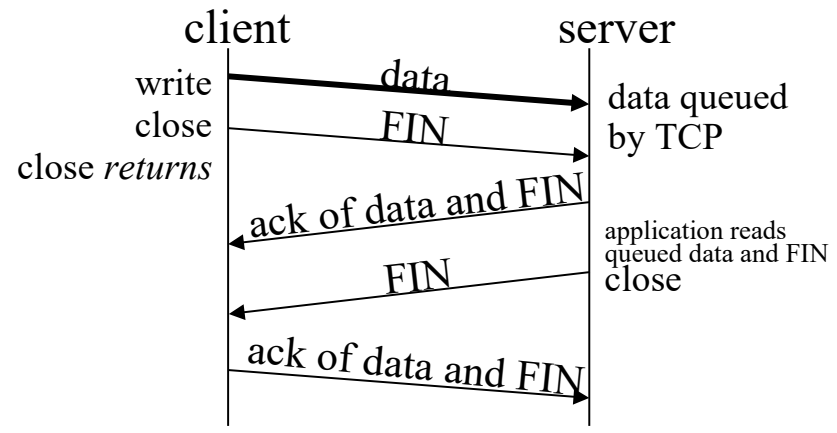
# Generic Socket Options (cont'd)

## • SO\_LINGER

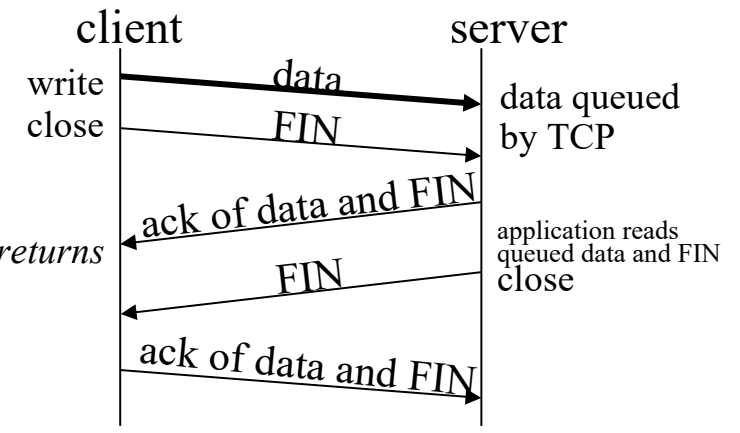
- Connection oriented protocol `close` function `int close(int sockfd)`

```

struct linger {
    int l_onoff; /* 0=off, nonzero=on */
    int l_linger; /* linger time, units as seconds */
}
  
```



Default operation of close



close with SO\_LINGER option



# Generic Socket Options (cont'd)

- **SO\_OOINLINE**

- out-of-band data `int opt`: `int` `int` `int` OOB data `int` normal input queue `int`

- **SO\_RCVBUF / SO\_SNDBUF**

- `int` `int` send buffer `int` receive buffer `int` `int` `int`
- **receive buffer** `int` `int` **TCP** `int` **flow control** `int` `int`

- **SO\_RCVLOWAT / SO\_SNDLOWAT**

- select function `int` `int` receive/send low-water mark `int` `int`
- receive low-water mark `int` default `int` 1
- send low-water mark `int` default `int` 2048

- **SO\_RCVTIMEO / SO\_SNDTIMEO**

- send/receive `int` `int` `int` `int`
- receive timeout affects: read, readv, recv, recvfrom, recvmsg
- send timeout affects: write, writev, send, sendto, sendmsg

# Generic Socket Options (cont'd)

- **SO\_REUSEADDR / SO\_REUSEPORT**

- SO\_REUSEADDR 0000 4 0000 00 0000 0000

- 1. 0 0000 00 0000 0000 00000 listening server 0 0 0000 binding 0 0 00 0

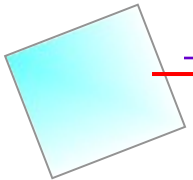
0000

- (a) listening server 0 00
- (b) 00 0000 00000, 0 client 0 0000 00 000000 0000
- (c) listening server 0 00000, 00 000000 00 0000 00000 0000
- (d) listening server 0 0000

00 00 00 SO\_REUSEADDR 0000 00000 00 0000, (d) 00000 bind 0 0000

- 2. 0000 instance 0 00 00 **IP** 0000 bind 00 00, 00 0000 00 00 instance 0 0000 0 0000 00

- ※ TCP 00 completely duplicate binding 0 00000 00.





# Generic Socket Options (cont'd)

3. `bind` `IP` `binding`
4. UDP `completely duplicate binding`

## - `SO_REUSEPORT` (4.4BSD)

- † multicasting `SO_REUSEPORT` (4.4BSD)
- † `IP` `binding` `SO_REUSEPORT` `completely duplicate binding`
- † `IP` `multicast address` `SO_REUSEADDR` `SO_REUSEPORT`

## • **`SO_TYPE`**

- `SOCK_STREAM`, `SOCK_DGRAM`, etc)

## • **`SO_USELOOPBACK` (`AF_ROUTE`)**

- `SO_USELOOPBACK`

# fnctl / ioctl Functions

- **fnctl, ioctl, routing socket operation** □ □ □

□ □ Operation	fnctl	ioctl	Routing socket	Posix.1g
set socket for nonblocking I/O	F_SETFL, O_NONBLOCK	FIONBIO		fnctl
set socket for signal-driven I/O	F_SETFL, O_ASYNC	FIOASYNC		fnctl
set socket owner	F_SETOWN	SIOCSPGRP or FIOSETOWN		fnctl
get socket owner	F_GETOWN	SIOCGPGRP or FIOGETOWN		fnctl
get #bytes in socket receive buffer		FIONREAD		
test for socket at out-of-bank mark		SIOCATMARK		socketmark
obtain interface list		SIOCGIFCONF	sysctl	
interface operations		SIOC[GS]Ifxxx		
ARP cache operations		SIOCxARP	RTM_xxx	
routing table operations		SIOCxxxRT	RTM_xxx	

# fnctl □□ □□ □

- **Nonblocking I/O** □ □□□□ □

```
int flags;
```

```
        /* Set socket nonblocking */
if ( (flags = fnctl (fd, F_GETFL, 0)) < 0 )
    err_sys("F_GETFL error");
flags |= O_NONBLOCK;
if ( (fnctl(fd, F_SETFL, flags) < 0 )
    err_sys("F_SETFL error");
```

- **Nonblocking I/O** □ □□□□ □

```
.....
flags &= ~O_NONBLOCK;
if ( (fnctl(fd, F_SETFL, flags) < 0 )
    err_sys("F_SETFL error");
```



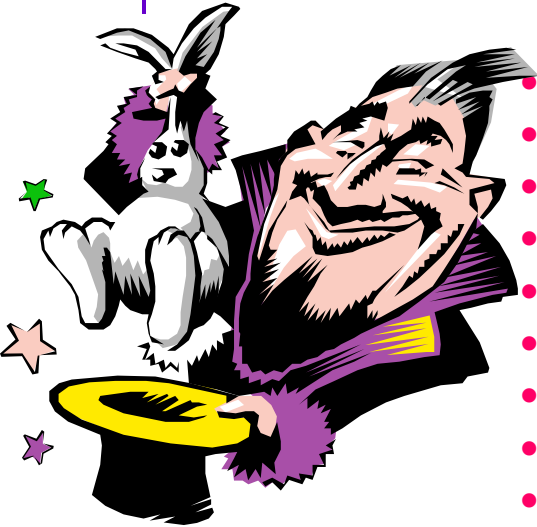
# Unix Network Programming

## Chapter 8 Elementary UDP Socket



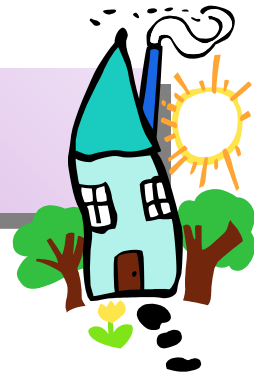
UDP □ □

# Contents

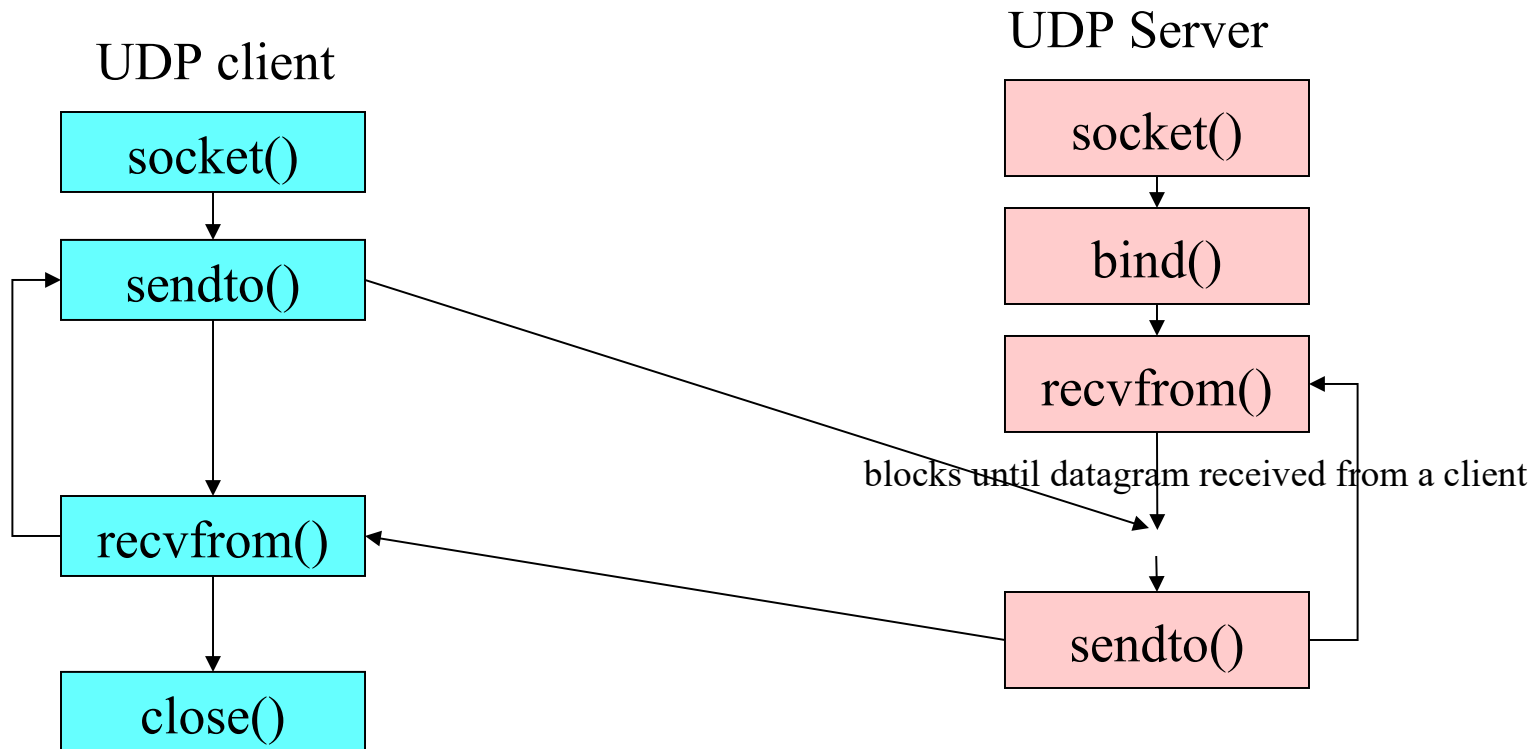


- **Introduction**
- *recvfrom* and *sendto* functions
- **UDP Echo Server : main Function**
- **UDP Echo Server : dg\_echo Function**
- **UDP Echo Client : main Function**
- **UDP Echo Client : dg\_cli Function**
- **Lost Diagrams**
- **Verifying Received Response**
- **Server Not Running**
- **Summary of UDP example**
- **connect Function with UDP**
- **dg\_cli Function(Revisited)**
- **Lack of Flow Control with UDP**
- **Determining Outgoing Interface with UDP**
- **TCP and UDP Echo Server Using select**
- **Summary**

# Introduction



- **connectionless, unreliable, datagram protocol**



# recvfrom and sendto Functions

- **#include <sys/socket.h>**

13 □□□ □□ . recv, send,  
recvmsg & sendmsg  
□□ 0 □□□ □□

- **ssize\_t recvfrom(int sockfd, void \*buff, size\_T nbytes, int flags, struct sockaddr \*from, socklen\_t \*addrlen);**

accept

- **ssize\_t sendto(int sockfd, const void \*buff, size\_t nbytes, int flags,**

connect

**const struct sockaddr \*to, socklen\_t addrlen);**

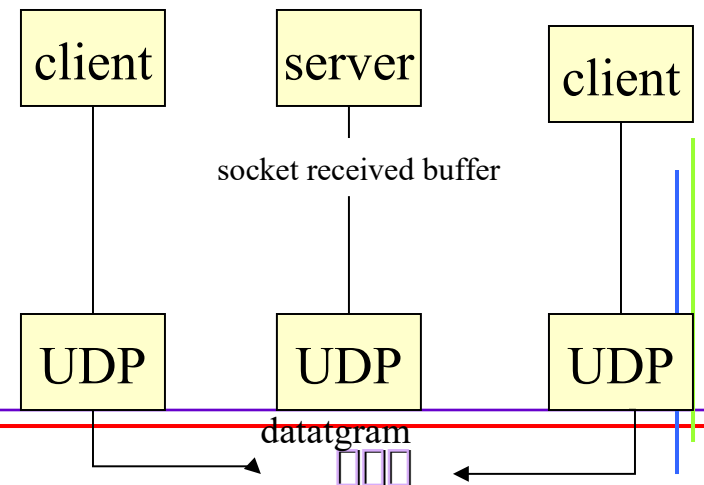
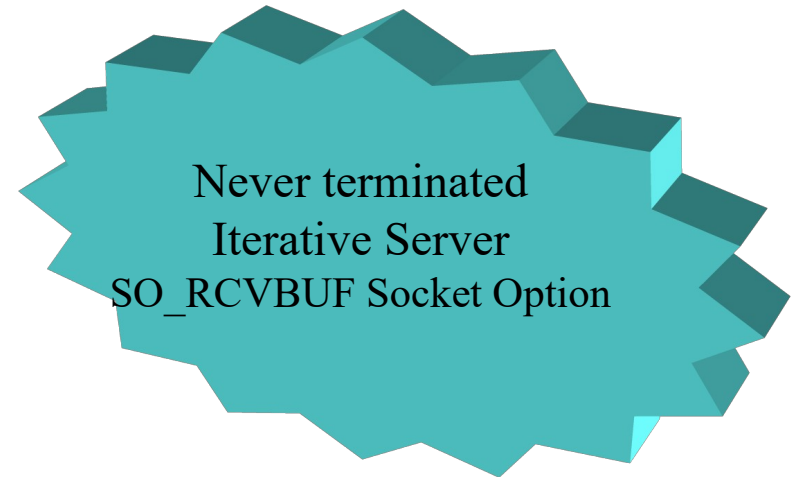
**both return : # of byte read or written if OK, -1 on error**

# UDP Echo Server

```

1.  #include "unp.h"
2.  int main(int argc, char **argv)
3.  { int    sockfd;
4.    struct sockaddr_in servaddr, cliaddr;
5.    sockfd = Socket(AF_INET, SOCK_DGRAM, 0);
6.    bzero(&servaddr, sizeof(servaddr));
7.    servaddr.sin_family    = AF_INET;
8.    servaddr.sin_addr.s_addr = htonl(INADDR_ANY);
9.    servaddr.sin_port      = htons(SERV_PORT);
10.   Bind(sockfd, (SA *) &servaddr, sizeof(servaddr));
11.   dg_echo(sockfd, (SA *) &cliaddr, sizeof(cliaddr)); }
12. void dg_echo(int sockfd, SA *pcliaddr, socklen_t clien){
13.     int  n;
14.     socklen_t len;
15.     char  mesg[MAXLINE];
16.     for ( ; ; ) {
17.         len = clien;
18.         n = Recvfrom(sockfd, mesg, MAXLINE, 0, pcliaddr, &len);
19.         Sendto(sockfd, mesg, n, 0, pcliaddr, len);
20.     }

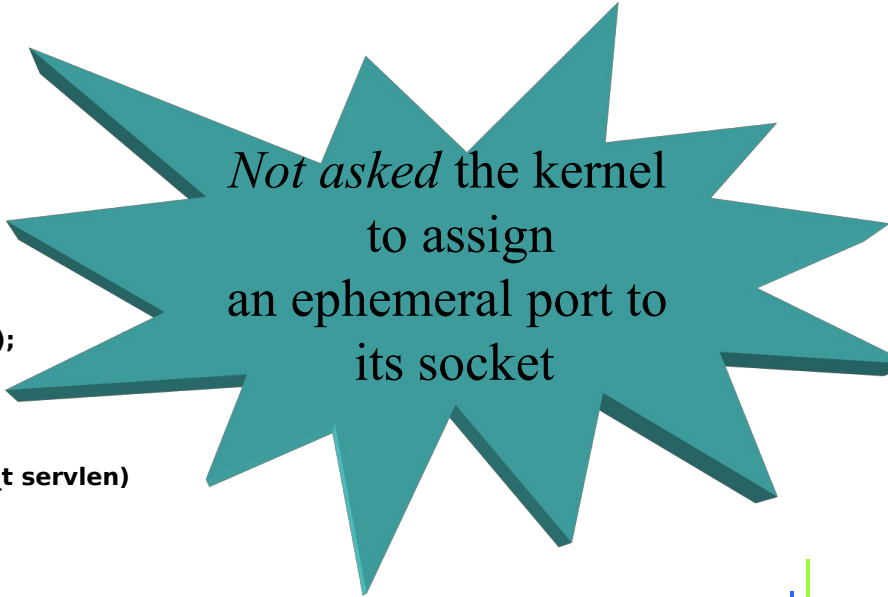
```





# UDP Echo Client

```
1. #include "unp.h"
2. int main(int argc, char **argv)
3. {
4.     int sockfd;
5.     struct sockaddr_in servaddr;
6.     if (argc != 2)
7.         err_quit("usage: udpcli <IPaddress>");
8.     bzero(&servaddr, sizeof(servaddr));
9.     servaddr.sin_family = AF_INET;
10.    servaddr.sin_port = htons(SERV_PORT);
11.    Inet_pton(AF_INET, argv[1], &servaddr.sin_addr);
12.    sockfd = Socket(AF_INET, SOCK_DGRAM, 0);
13.    dg_cli(stdin, sockfd, (SA *) &servaddr, sizeof(servaddr));
14.    exit(0);
15. }
16. void
17. dg_cli(FILE *fp, int sockfd, const SA *pservaddr, socklen_t servlen)
18. {
19.     int n;
20.     char sendline[MAXLINE], recvline[MAXLINE + 1];
21.     while (Fgets(sendline, MAXLINE, fp) != NULL) {
22.         Sendto(sockfd, sendline, strlen(sendline), 0, pservaddr, servlen);
23.         n = Recvfrom(sockfd, recvline, MAXLINE, 0, NULL, NULL);
24.         recvline[n] = 0; /* null terminate */
25.         Fputs(recvline, stdout);
26.     } }
```



*Not asked* the kernel  
to assign  
an ephemeral port to  
its socket

# Lost Datagrams

- **Our UDP client-server example is *not reliable*.**
- **If a client is lost, the client will *block forever* in its call to *recvfrom* in the function `dg_cli`, waiting for a server reply.**
- **Only way to prevent this is to *place a timeout* on the client's call to *recvfrom*.(13.2)**
  - connect with a Timeout Using SIGALRM
  - recvfrom with a Timeout Using SIGALRM
- **Reliability to a UDP client-server in 20.5 p542**

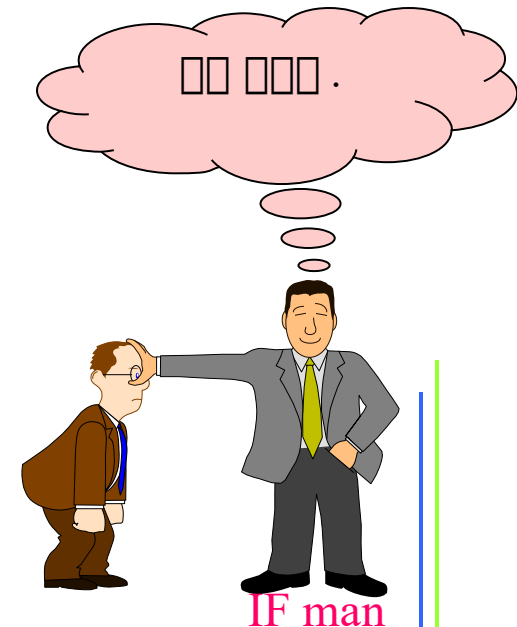


# Verifying Received Response

```

1.  #include "unp.h"
2.  void dg_cli(FILE *fp, int sockfd, const SA *pservaddr, socklen_t servlen)
3.  {
4.      int    n;
5.      char   sendline[MAXLINE], recvline[MAXLINE + 1];
6.      socklen_t len;
7.      struct sockaddr *preply_addr;
8.      preply_addr = Malloc(servlen);
9.      while (Fgets(sendline, MAXLINE, fp) != NULL) {
10.         Sendto(sockfd, sendline, strlen(sendline), 0, pservaddr, servlen);
11.         len = servlen;
12.         n = Recvfrom(sockfd, recvline, MAXLINE, 0, preply_addr, &len);
13.         if (len != servlen || memcmp(pservaddr, preply_addr, len) != 0) {
14.             printf("reply from %s (ignored)\n",
15.                 Sock_ntop(preply_addr, len));
16.             continue;
17.         }
18.         recvline[n] = 0; /* null terminate */
19.         Fputs(recvline, stdout);
20.     }
21. }

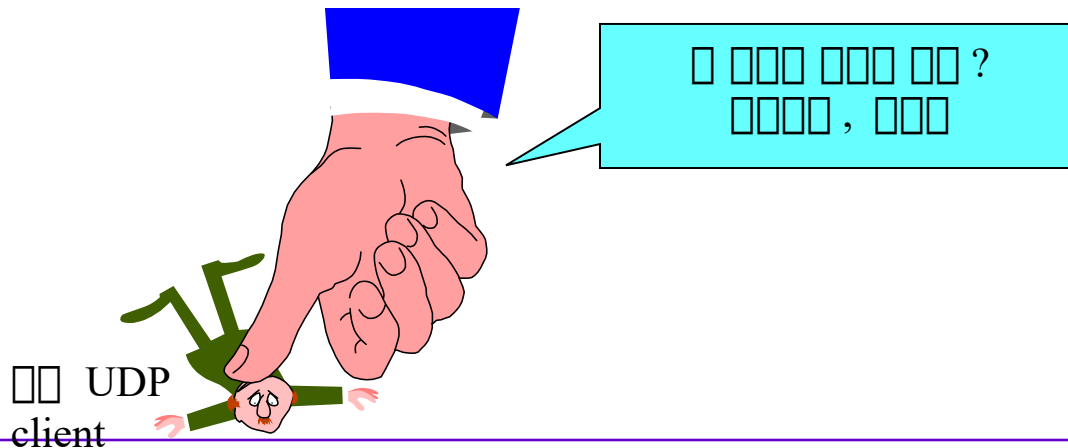
```



# Server Not Running

- **Asynchronous error**

- errors that are reported some time **after** the packet was sent.
- ICMP port unreachable
- TCP : always report these errors to the application.



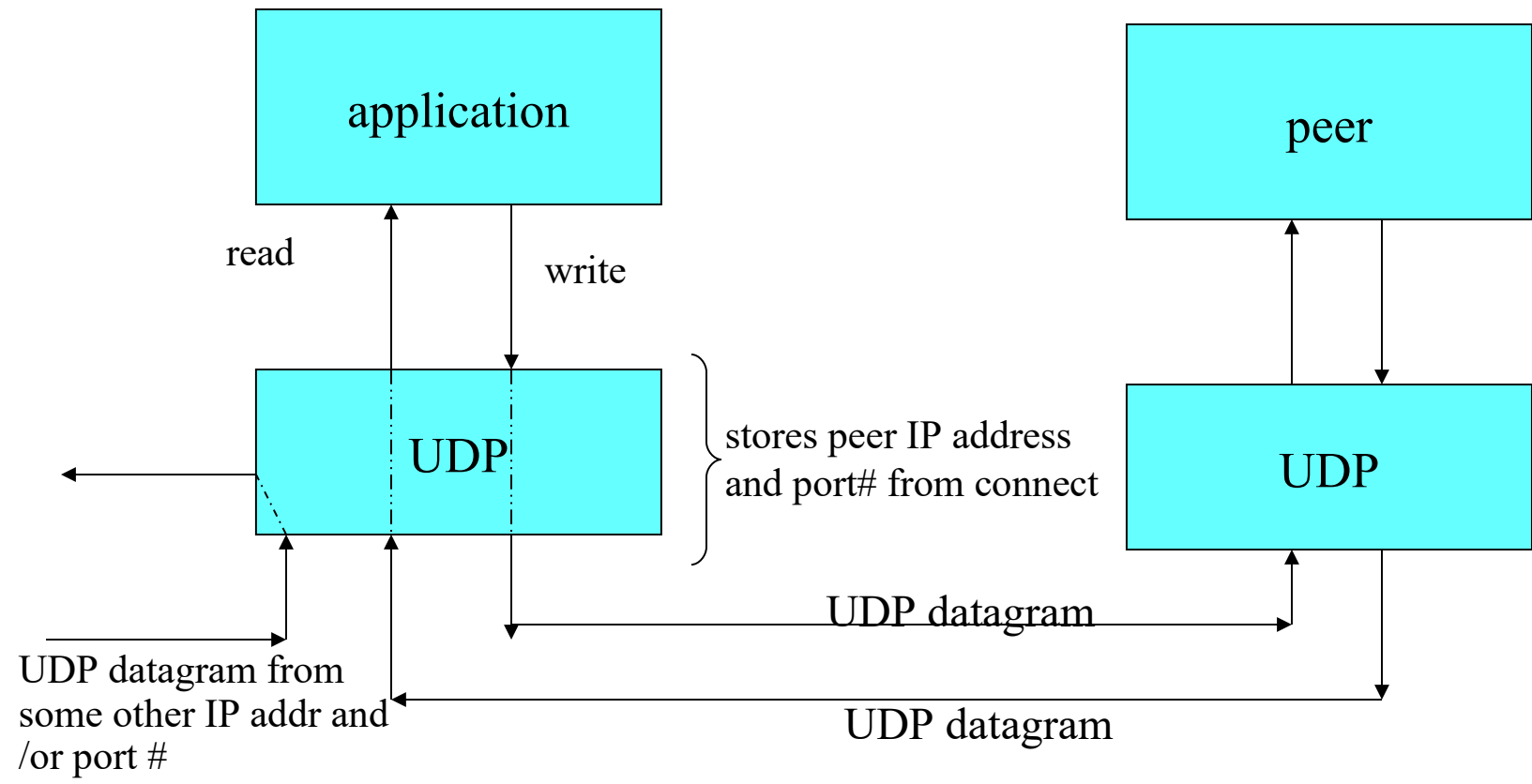
# Connect Function with UDP(1)

- **unconnected UDP Socket**
  - the default when we create a UDP socket
- **connected UDP Socket**
  - the result of calling connect on a UDP socket



1. Does *not* use **sendto** but use **write** or **send** instead.
2. Does *not* use **recvfrom** but use **read** or **recv** instead.
3. **Asynchronous errors** are returned to the process for a connected UDP socket. The corollary, as we previously described, is that an unconnected UDP socket does not receive any asynchronous errors

# connect Function with UDP(2)

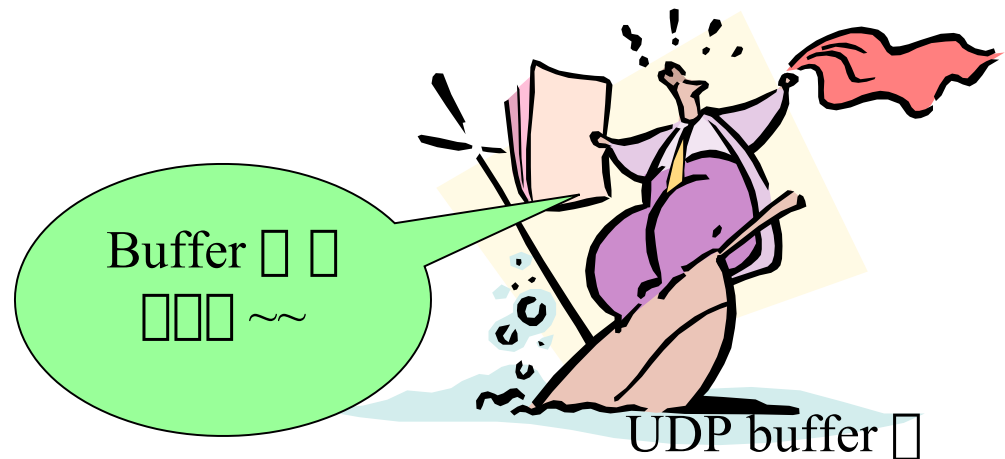


# dg\_client Function

```
1.  #include "unp.h"
2.  void
3.  dg_cli(FILE *fp, int sockfd, const SA *pservaddr, socklen_t servlen)
4.  {
5.      int n;
6.      char sendline[MAXLINE], recvline[MAXLINE + 1];
7.      Connect(sockfd, (SA *) pservaddr, servlen);
8.      while (Fgets(sendline, MAXLINE, fp) != NULL) {
9.          Write(sockfd, sendline, strlen(sendline));
10.         n = Read(sockfd, recvline, MAXLINE);
11.         recvline[n] = 0; /* null terminate */
12.         Fputs(recvline, stdout);
13.     }
14. }
```

# Lack of Flow Control with UDP

- **UDP has no flow control and unreliable**
  - ex) UDP sender to overrun the receiver.
- **UDP Socket Receive Buffer**
  - can change receive Buffer : `SO_RCVBUF`
  - default : 41,600 bytes(BSD/OS)
    - † actual limit : 246,723 bytes





# TCP and UDP Echo Server Using select

```

1. int main(int argc, char **argv) {
2.     int     listenfd, connfd, udpfd, nready, maxfdp1;
3.     char     mesg[MAXLINE];
4.     pid_t    childpid;
5.     fd_set   rset;
6.     ssize_t  n;
7.     socklen_t len;
8.     const int on = 1;
9.     struct sockaddr_in cliaddr, servaddr;
10.    void     sig_chld(int);
11.    /* create listening TCP socket */
12.    listenfd = Socket(AF_INET, SOCK_STREAM, 0);
13.    bzero(&servaddr, sizeof(servaddr));
14.    servaddr.sin_family = AF_INET;
15.    servaddr.sin_addr.s_addr = htonl(INADDR_ANY);
16.    servaddr.sin_port = htons(SERV_PORT);
17.    Setsockopt(listenfd, SOL_SOCKET, SO_REUSEADDR,
18.    &on, sizeof(on));
19.    Bind(listenfd, (SA *) &servaddr, sizeof(servaddr));
20.    Listen(listenfd, LISTENQ);
21.    /* create UDP socket */
22.    udpfd = Socket(AF_INET, SOCK_DGRAM, 0);
23.    bzero(&servaddr, sizeof(servaddr));
24.    servaddr.sin_family = AF_INET;
25.    servaddr.sin_addr.s_addr = htonl(INADDR_ANY);
26.    servaddr.sin_port = htons(SERV_PORT);
27.    Bind(udpfd, (SA *) &servaddr, sizeof(servaddr));
28.    /* include udpserverselect02 */
29.    Signal(SIGCHLD, sig_chld); /* must call waitpid() */
30.    FD_ZERO(&rset);
31.    maxfdp1 = max(listenfd, udpfd) + 1;
32.    for ( ; ; ) {
33.        FD_SET(listenfd, &rset);
34.        FD_SET(udpfd, &rset);
35.        if ( (nready = select(maxfdp1, &rset, NULL, NULL, NULL)) < 0)
36.        {
37.            if (errno == EINTR) continue; /* back to for() */
38.            else err_sys("select error"); }
39.        if (FD_ISSET(listenfd, &rset)) {
40.            len = sizeof(cliaddr);
41.            connfd = Accept(listenfd, (SA *) &cliaddr, &len);
42.            if ( (childpid = Fork()) == 0) { /* child process */
43.                Close(listenfd); /* close listening socket */
44.                str_echo(connfd); /* process the request */
45.                exit(0); }
46.            Close(connfd); /* parent closes connected socket */
47.        }
48.        if (FD_ISSET(udpfd, &rset)) {
49.            len = sizeof(cliaddr);
50.            n = Recvfrom(udpfd, mesg, MAXLINE, 0, (SA *) &cliaddr,
51.            &len);
52.            Sendto(udpfd, mesg, n, 0, (SA *) &cliaddr, len);
53.        }
54.    }

```

# /etc/services □□□

•	<b>tcpmux</b>	<b>1/tcp</b>	
•	<i>echo</i>	<i>7/tcp</i>	
•	<i>echo</i>	<i>7/udp</i>	
•	<b>discard</b>	<b>9/tcp</b>	<b>sink null</b>
•	<b>discard</b>	<b>9/udp</b>	<b>sink null</b>
•	<b>sysstat</b>	<b>11/tcp</b>	<b>users</b>
•	<i>daytime</i>	<i>13/tcp</i>	
•	<i>daytime</i>	<i>13/udp</i>	
•	<b>netstat</b>	<b>15/tcp</b>	
•	<b>chargen</b>	<b>19/tcp</b>	<b>ttytst source</b>
•	<b>chargen</b>	<b>19/udp</b>	<b>ttytst source</b>
•	<b>ftp-data</b>	<b>20/tcp</b>	
•	<b>ftp</b>	<b>21/tcp</b>	
•	<b>telnet</b>	<b>23/tcp</b>	
•	<i>smtp</i>	<i>25/tcp</i>	<b>mail</b>
•	<i>time</i>	<i>37/tcp</i>	<i>timserver</i>
•	<i>time</i>	<i>37/udp</i>	<b>timserver</b>
•	<b>name</b>	<b>42/udp</b>	<b>nameserver</b>
•	<b>whois</b>	<b>43/tcp</b>	<b>nicname</b>
•	<i>domain</i>	<i>53/udp</i>	
•	<i>domain</i>	<i>53/tcp</i>	
•	<b>bootps</b>	<b>67/udp</b>	<b># BOOTP/DHCP server</b>
•	<b>bootpc</b>	<b>68/udp</b>	<b># BOOTP/DHCP client</b>

1 port □  
 UDP/TCP □□□  
 □□□□ □□ □□□□  
 □□



# usually to sri-nic

# BOOTP/DHCP server  
 # BOOTP/DHCP client



# Unix Network Programming

## Chapter 9.

### Elementary Name and Address Conversions



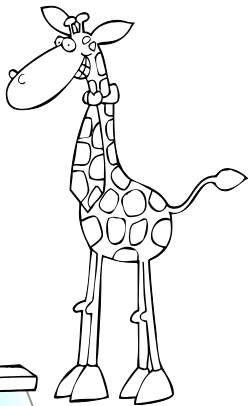
# Contents



- Introduction
- Domain Name System
- *gethostbyname* Functions
- RES\_USE\_INET6 Resolver Option
- *gethostbyname2* Function and IPV6 Support
- *gethostbyaddr* Function
- *uname* Function
- *gethostname* Function
- *getservbyname* and *getservbyport* Functions
- Other Networking Information
- Summary

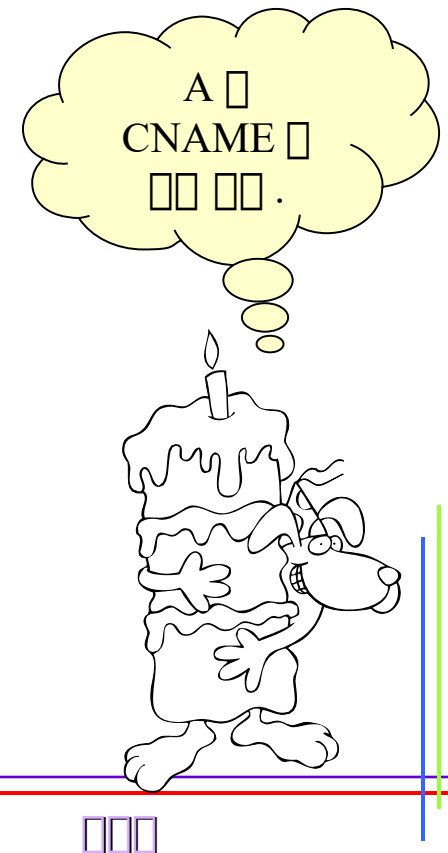
# Introduction

- **So far, we used numeric addresses for the host**
- **User names instead of numbers for numerous reasons**
  - names are easier to remember
    - † gethostbyname, gethostbyaddr (hostnames -> IP addr)
    - † getservbyname, getservbyport(service names -> port )

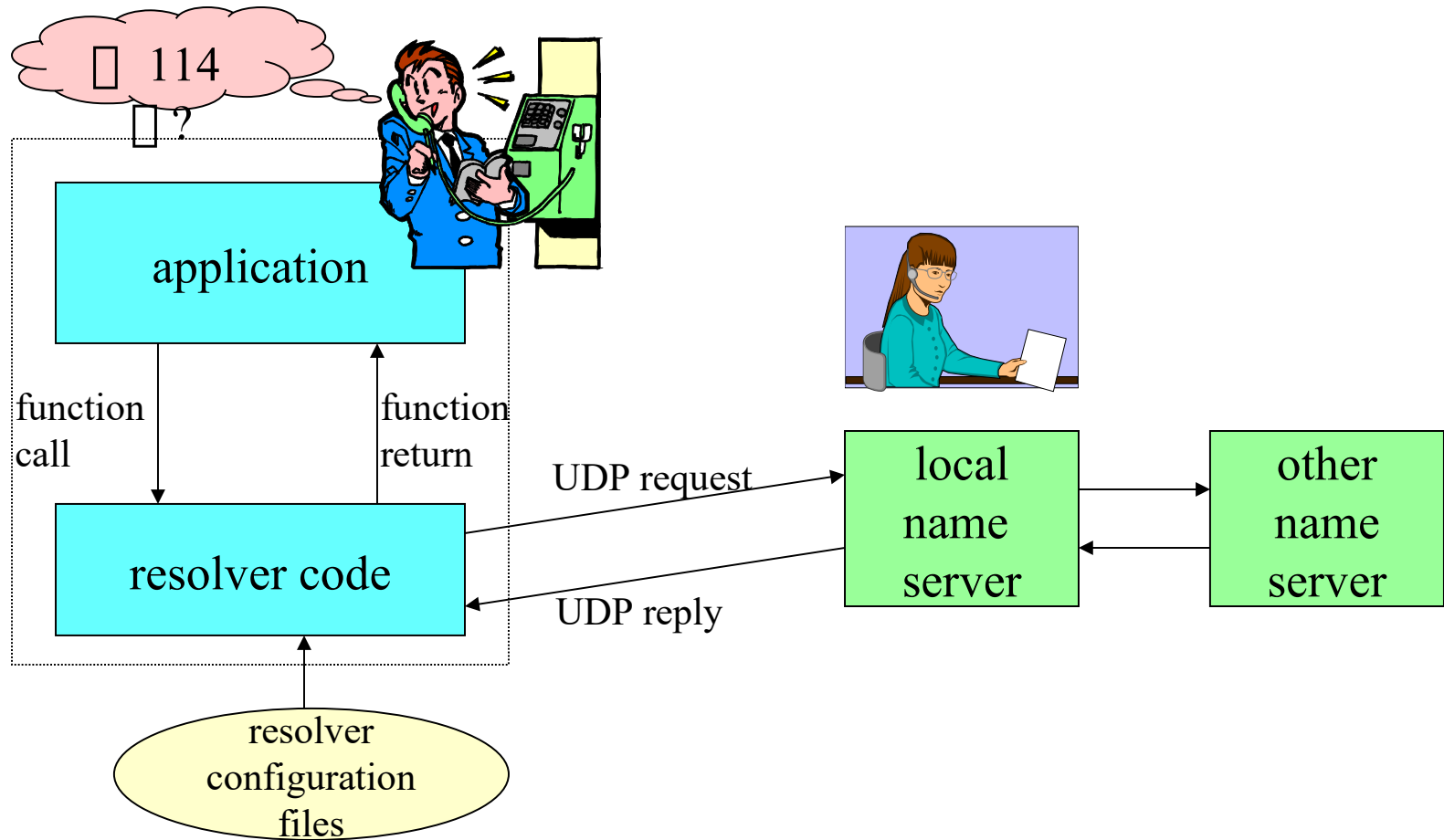


# Domain Name System

- **DNS is used primarily to map between hostnames and IP addresses.**
- **Resource Records**
  - A : 32bit IPv4 address
  - AAAA : 128bit IPv6 address
  - PTR : pointer records
  - MX : mail exchanger
  - CNAME : canonical name (ftp, www,...)



# Typical arrangement of clients, resolvers, and name servers



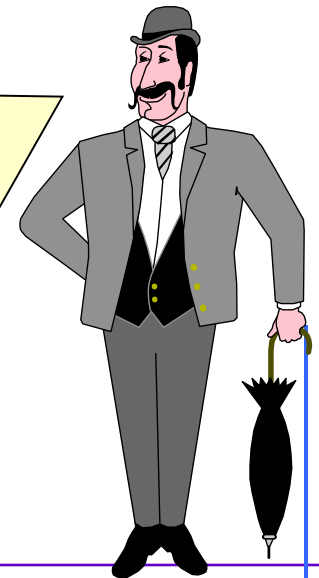
# gethostbyname function

- **#include <netdb.h>**
- **struct hostent \*gethostbyname(const char \*hostname);**

**Returns : nonull pointer if OK, NULL on error with h\_errno set**

```
struct hostent {
    char    *h_name;
    char    **h_aliases;
    int     h_addrtype;
    int     h_length;
    char    **h_addr_list;
};
#define h_addr h_addr_list[0];
```

page 242 □ □□□ ,  
 □□□□ □□□  
 □□□□ □ □ □□□ .  
 IPv6 □ 9.4, 9.5 □□  
 □□□□ □□  
 □□□□□□  
 □□□□□□ .





# gethostbyaddr Function / uname Function

AF\_INET  
or  
AF\_INET6

- **#include <netdb.h>**
- **struct hostent \*gethostbyaddr(const char \*addr, size\_t len, int family)**

**Returns : nonull pointer if OK, NULL on error with h\_errno set**

- **#include <sys/utsname.h>**
- **int uname(struct utsname \*name)**

**Returns : nonnegative value is OK, -1 on error**

```
#define _UTS_NAMESIZE 16
#define _UTS_NODESIZE 256
struct utsname{
    char sysname[_UTS_NAMESIZE]; /* os name */
    char nodename[_UTS_NODESIZE]; /* node name */
    char release[_UTS_NAMESIZE]; /* O.S. release level */
    char version[_UTS_NAMESIZE]; /* O.S. version level */
    char machine[_UTS_NAMESIZE]; /*hardware type */
}
```

uname() □□□  
local IP □  
□□□□



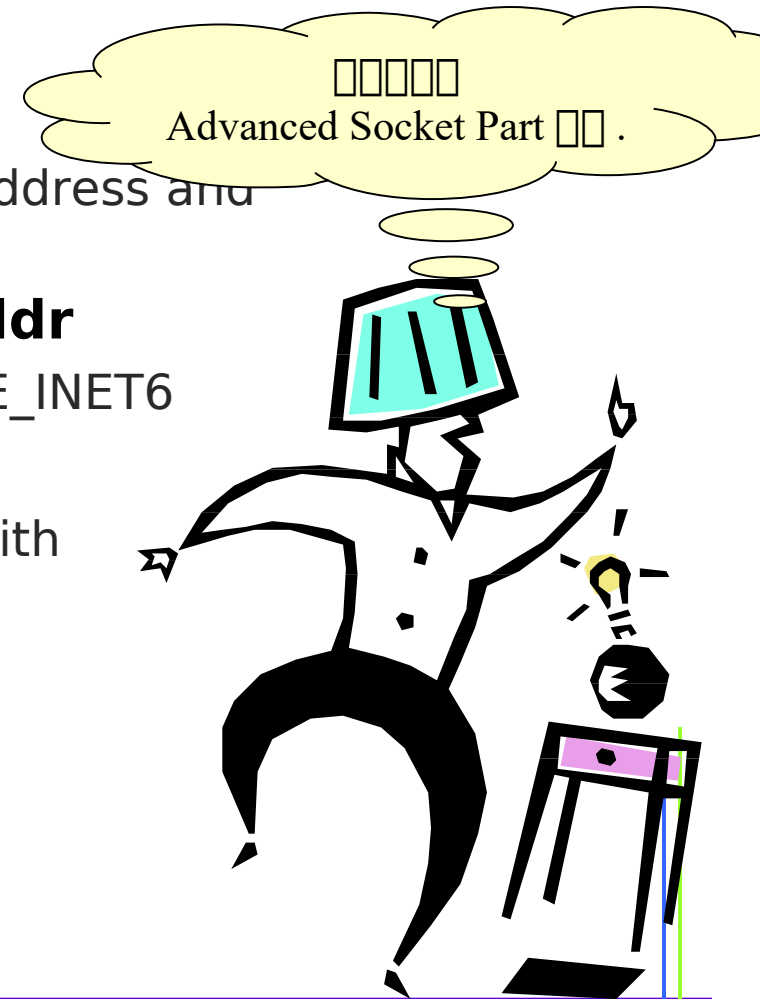
# getservbyname / getservbyport

- **/etc/services** □ □□□□ , **local machine** □ □□□□ **server** □ □□□□ , □□□□ □□□□□ **port** □ □□□□□□ .

Information	Data file	Structure	keyed lookup function
<b>hosts</b>	<b>/etc/hosts</b>	<b>hostent</b>	<b>gethostbyaddr, gethostbyname</b>
<b>networks</b>	<b>/etc/networks</b>	<b>netent</b>	<b>getnetbyaddr, getnetbyname</b>
<b>protocols</b>	<b>/etc/protocols</b>	<b>protoent</b>	<b>getprotobyname, getprotobynumber</b>
<b>services</b>	<b>/etc/services</b>	<b>servent</b>	<b>getservbyname, getservbyport</b>

# Summary

- **resolver**
  - to convert a hostname into an IP address and vice versa
- **gethostbyname / getthostbyaddr**
  - IPV6 □ : gethostname2(), RES\_USE\_INET6
- **getservbyname**
  - commonly used function dealing with services names and port





# UNIX Network Programming

( chap 10 - chap 11 )

# Agenda

## ■ IPv4 and IPv6 Interoperability ( chap.10 )

- ☞ **IPv6 Server on dual-stack host**
- ☞ **Processing of sever on dual-stack host**
- ☞ **Dual stack host**
- ☞ **Processing of client requests**
- ☞ **IPv4-mapped IPv6**

# Agenda

## ■ Advanced Name and Address Conversions ( chap.11 )

☞ Why use *getaddrinfo()* & *getnameinfo()* ?

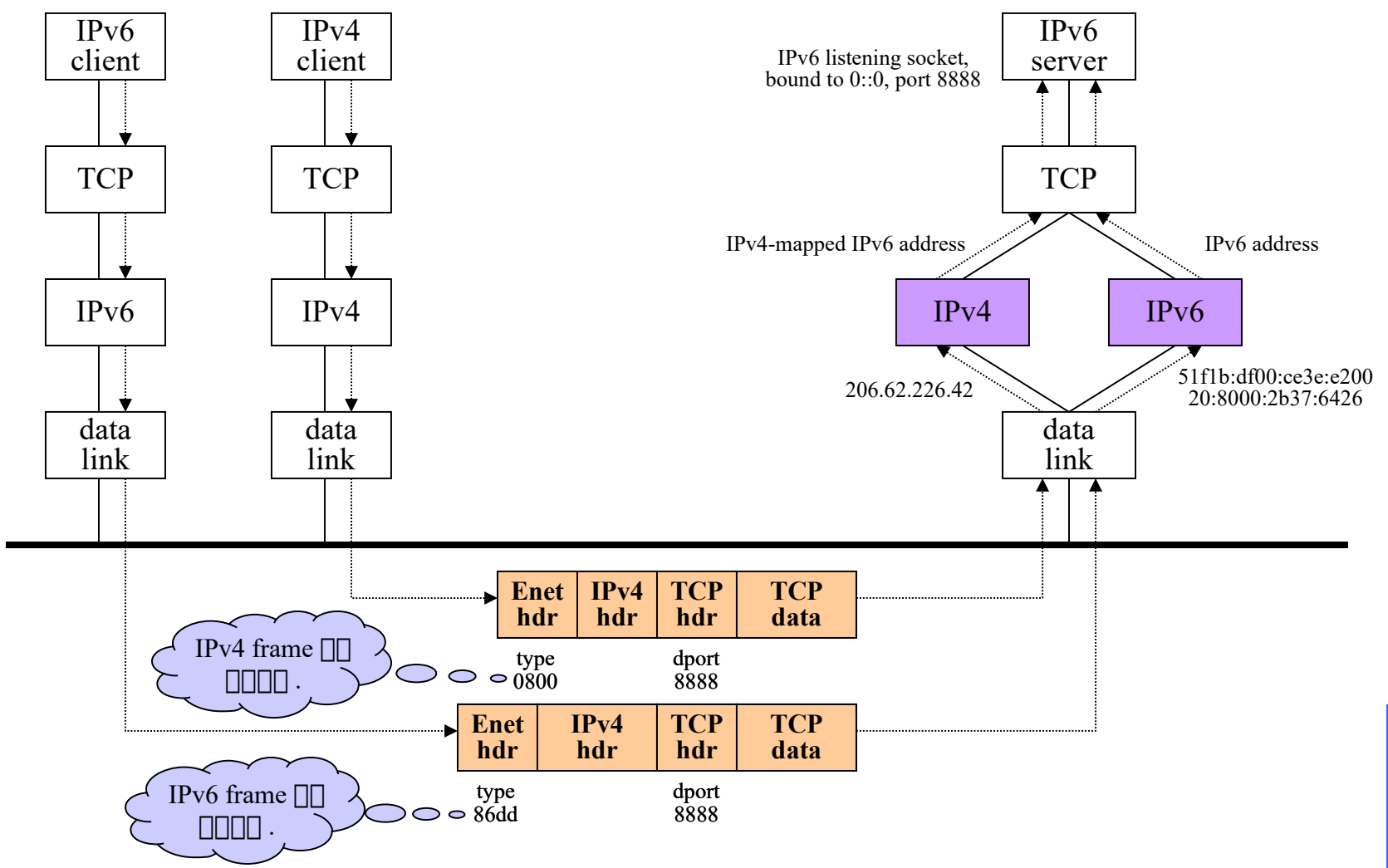
☞ *getaddrinfo()*

☞ Action and Result of *getaddrinfo()*

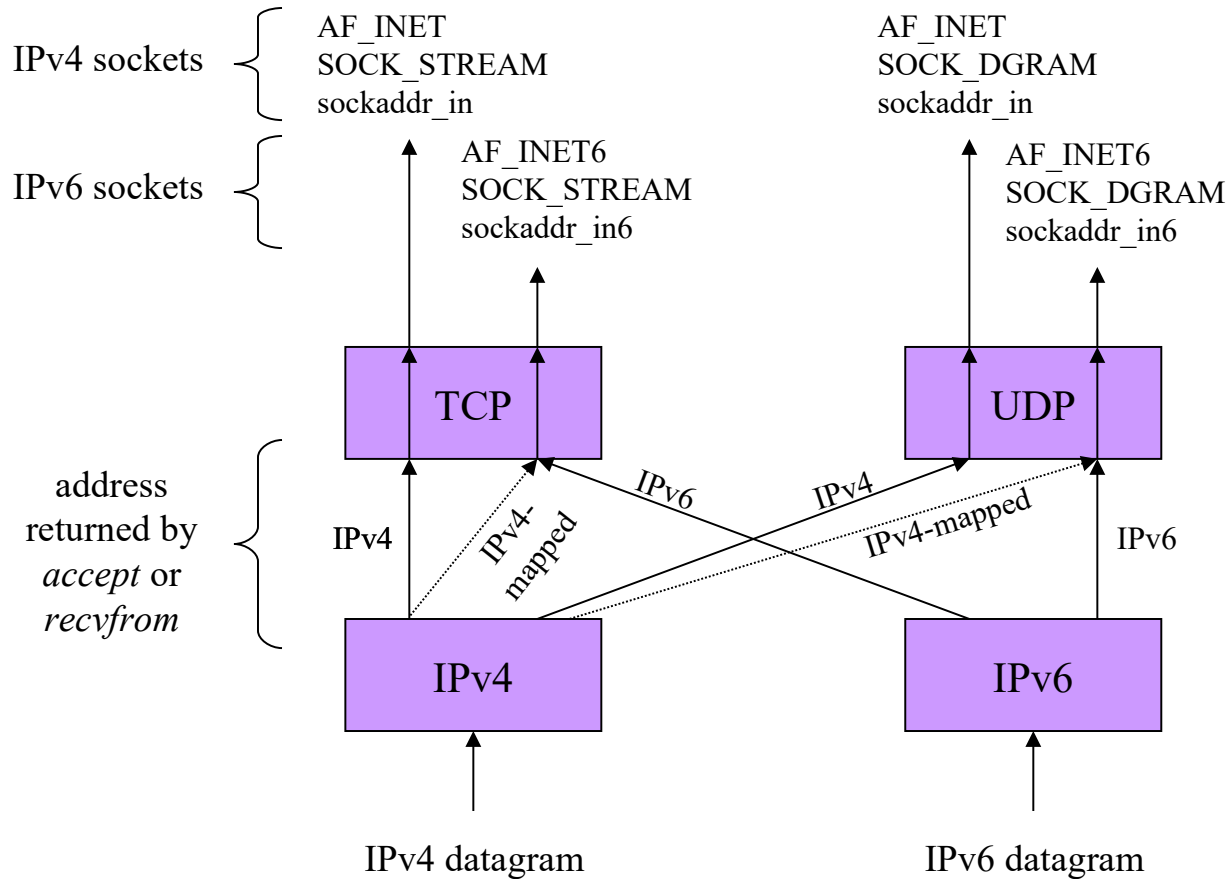
☞ *getnameinfo()*

☞ Reentrant functions

# IPv6 Server on dual-stack host



# Processing of server on dual-stack host



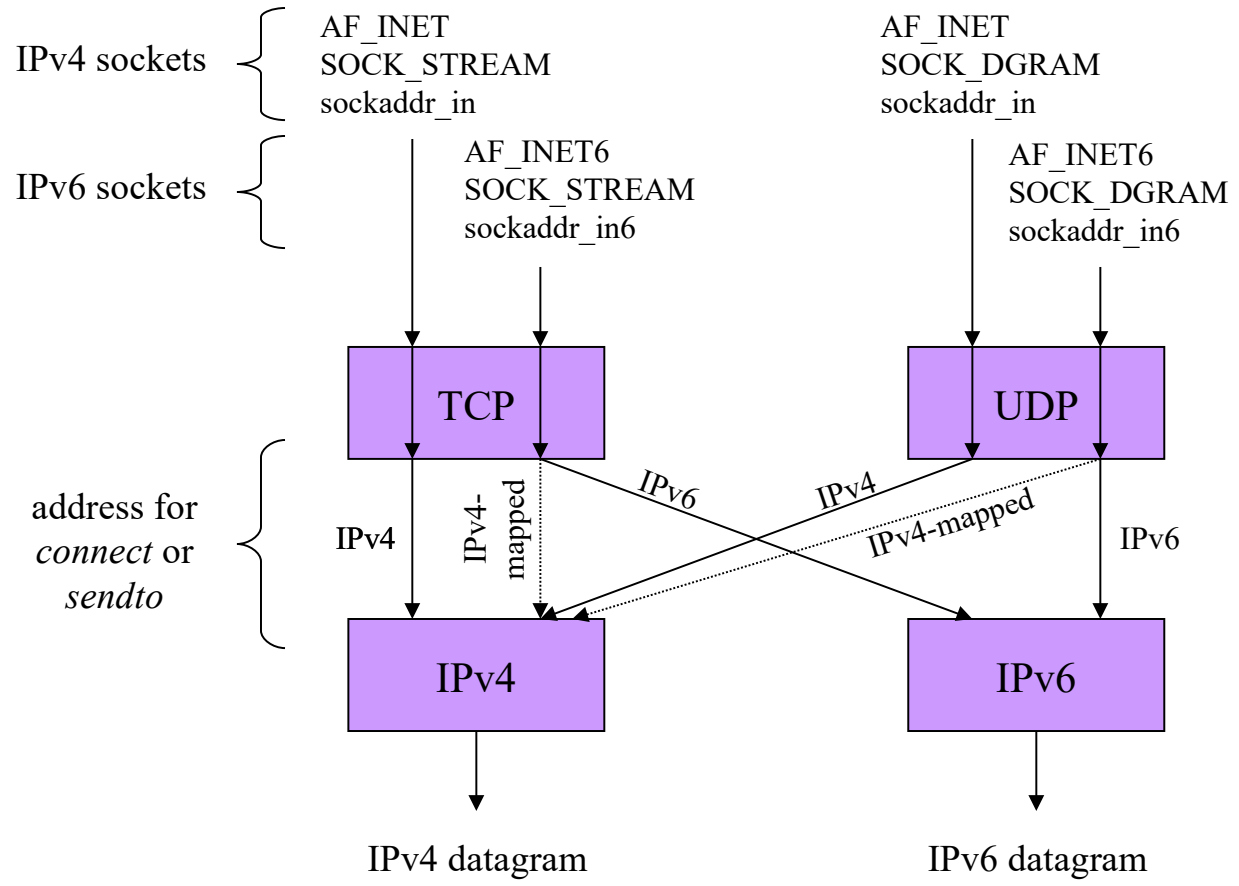


# Dual-stack host

## Listening socket rules

1. A listening IPv4 socket can accept incoming connections from only IPv4 clients.
2. If a server has a listening IPv6 socket that has bound the wildcard address, that socket can accept incoming connections from either IPv4 clients or IPv6 clients. For a connection from an IPv4 client the server's local address for the connections will be the corresponding IPv4-mapped IPv6 address.
3. If a server has a listening IPv6 socket that has bound an IPv6 address other than an IPv4-mapped IPv6 address, that socket can accept incoming connections from IPv6 clients only

# Processing of client requests



# IPv4-mapped IPv6

## IPv6 server

IPv4 → IPv4-mapped IPv6 : kernel `accept` `recvfrom` application `sendto` .

## IPv6 client

IPv6 → IPv4-mapped IPv6 : resolver `gethostbyname` `gethostbyaddr` etc )  
 application `connect` `sendto` .

### DNS records example

solaris	IN A	206.62.226.33
	IN AAAA	5f1b:df00:c33e:e200:0020:0800:2078:e3e3
aix	IN A	206.62.226.43
	IN A	5f1b:df00:c33e:e200:0020:0800:5afc:2b36
bsdi2	IN A	206.62.226.34

### IPv4-mapped IPv6 address example

return address - `0::ffff:206.62.226.34`

# Summary of interoperability

	IPv4 sever IPv4-only host (A only)	IPv6 sever IPv6-only host (AAAA only)	IPv4 sever dual-stack host (A and AAAA)	IPv6 sever dual-stack host (A and AAAA)
IPv4 client, IPv4-only host	IPv4	X	IPv4	IPv4
IPv6 client, IPv6-only host	X	IPv6	X	IPv6
IPv4 client, dual-stack host	IPv4	X	IPv4	IPv4
IPv6 client, dual-stack host	IPv4	IPv6	X*	IPv6

## Why use `getaddrinfo()` & `getnameinfo()` ?

`gethostbyname()` □ `gethostbyaddr()` □ protocol □ □□□□□ .  
 ▶ address family □ □□□□ □□□ .

Because of !!  
 protocol independence  
 of our application

```
int getaddrinfo( const char *hostname, const char *service,
                const struct addrinfo *hints, struct addrinfo **result );

int getnameinfo( const struct sockaddr *sockaddr, socklen_t addrlen,
                 char *host, size_t hostlen,
                 char *serv, size_t servlen, int flags );
```

# getaddrinfo()

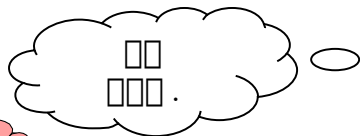
hostname : host name or address string

service : service name or decimal port number string

```

struct addrinfo {                                // define at <netdb.h>
    int     ai_flags;                            /* AI_PASSIVE, AI_CANONNAME */
    int     ai_family;                          /* AF_XXX */
    int     ai_socktype;                        /* SOCK_XXX */
    int     ai_protocol;                        /* 0 or IPPROTO_XXX, TCP or UDP */
    size_t  ai_addrlen;                         /* length of ai_addr, IPv4:16, IPv6:24 */
    char    *ai_canonname;                      /* ptr to canonical name for host */
    struct  sockaddr *ai_addr;                  /* ptr to socket address structure */
    struct  addrinfo *ai_next;                 /* ptr to next structure in linked list */
};

```



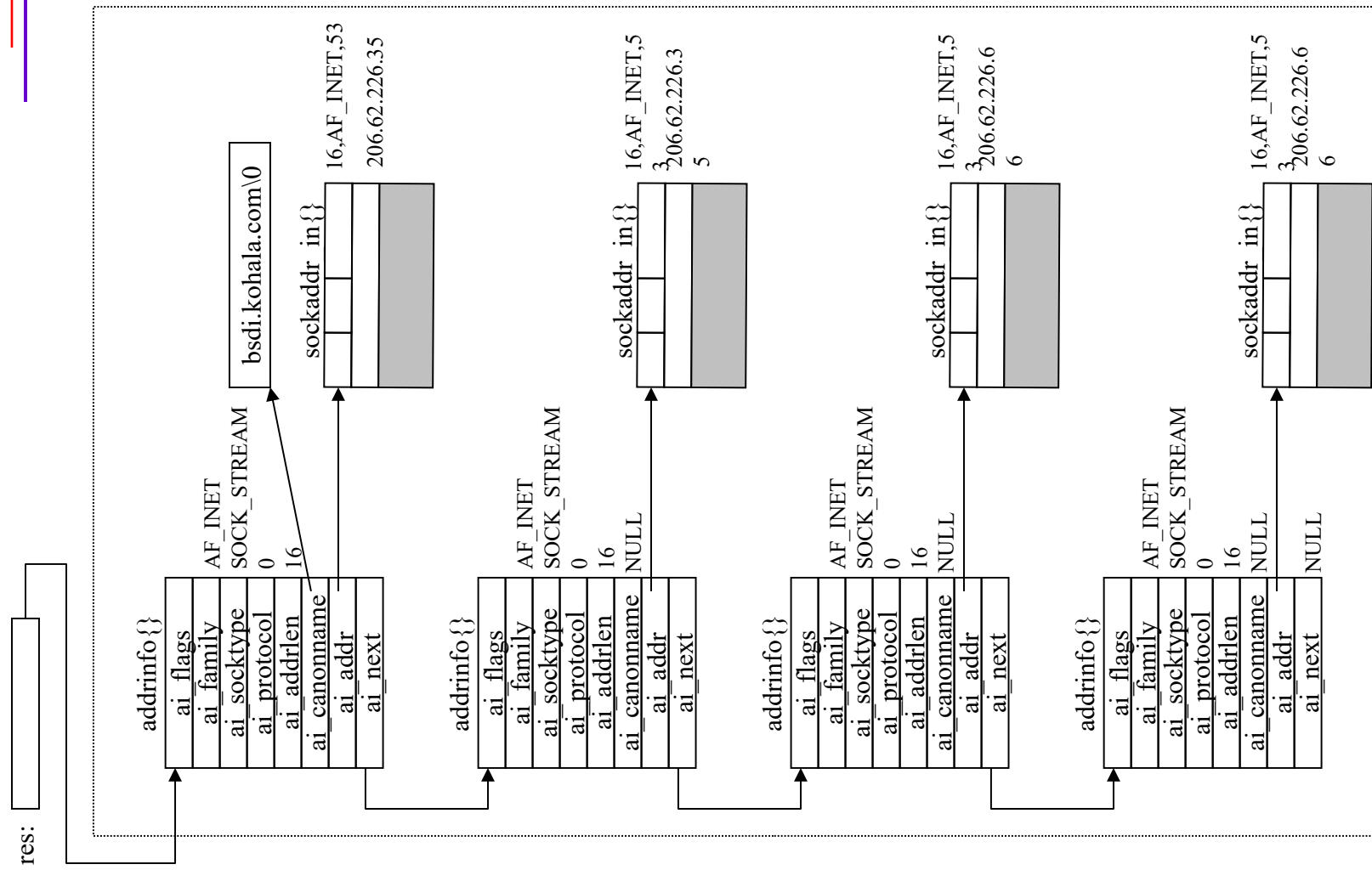
**Ex)**

```

struct addrinfo  hints, *res;
bzero( &hints, sizeof(hints) );
hints.ai_flags = AI_CANONNAME;
hints.ai_family = AF_INET;
getaddrinfo( "bdsi", "domain", &hints, &res );

```

# IP address



dynamically allocated by *getaddrinfo*



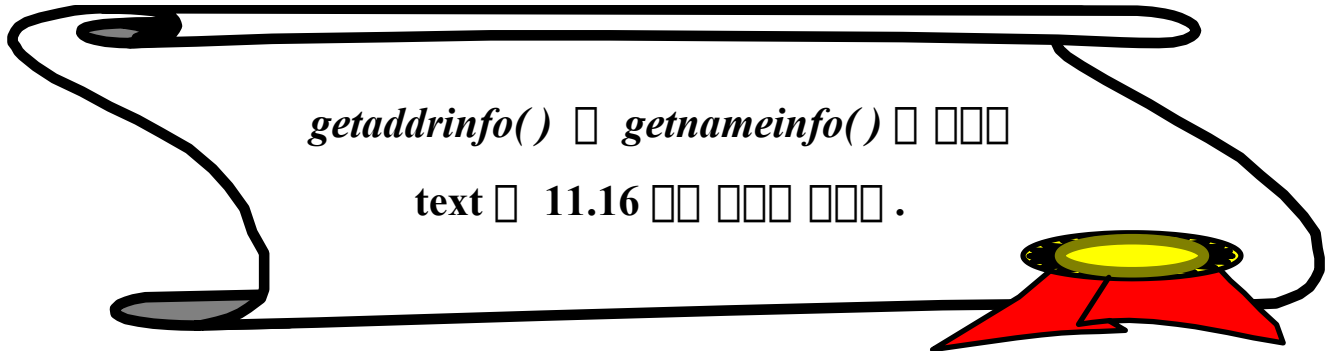
# Action and Result of *getaddrinfo()*

Hostname specified by caller	Address family specified by caller	Hostname string contains	Result	Action
nonnull hostname string; active or passive	AF_UNSPEC	hostname	all AAAA records returned as <code>sockaddr_in6{}s</code> and all A records returned as <code>sockaddr_in{}s</code>	two DNS searches : gethostbyname2(AF_INET6) with RES_USE_INET6 off gethostbyname2(AF_INET) with RES_USE_INET6 off
		hex string	one <code>sockaddr_in6{}s</code>	<code>inet_pton(AF_INET6)</code>
		dotted decimal	one <code>sockaddr_in{}s</code>	<code>inet_pton(AF_INET)</code>
	AF_INET6	hostname	all AAAA records returned as <code>sockaddr_in6{}s</code> <b>else</b> all A records returned as IPv4-mapped IPv6 as <code>sockaddr_in{}s</code>	gethostbyname() with RES_USE_INET6 <b>on</b>
		hex string	one <code>sockaddr_in6{}s</code>	<code>inet_pton(AF_INET6)</code>
		dotted decimal	error : EAI_ADDRFAMILY	
	AF_INET	hostname	all A records returned as <code>sockaddr_in{}s</code>	gethostbyname() with RES_USE_INET6 off
		hex string	error : EAI_ADDRFAMILY	
		dotted decimal	one <code>sockaddr_in{}s</code>	<code>inet_pton(AF_INET)</code>

Continue ...



Hostname specified by caller	Address family specified by caller	Hostname string contains	Result	Action
null hostname string; passive	AF_UNSPEC	implied 0::0 implied 0.0.0.0	one sockaddr_in6{}s and one sockaddr_in{}s	inet_pton(AF_INET6) inet_pton(AF_INET)
	AF_INET6	implied 0::0	one sockaddr_in6{}	inet_pton(AF_INET6)
	AF_INET	implied 0.0.0.0	one sockaddr_in{}	inet_pton(AF_INET)
null hostname string; active	AF_UNSPEC	implied 0::1 implied 127.0.0.1	one sockaddr_in6{}s and one sockaddr_in{}s	inet_pton(AF_INET6) inet_pton(AF_INET)
	AF_INET6	implied 0::1	one sockaddr_in6{}	inet_pton(AF_INET6)
	AF_INET	implied 127.0.0.1	one sockaddr_in{}	inet_pton(AF_INET)



# getnameinfo()

Socket address [ ] [ ] [ ] host [ ] service [ ] [ ] string [ ] [ ] [ ] [ ] [ ] .

[ ] protocol independence [ ] [ ] ... !



sock\_ntop

DNS [ ] [ ] [ ] [ ] [ ] IP address ( IPv4 [ ] [ ] [ ] dotted decimal, IPv6 [ ] [ ] [ ] hex string ) [ ] port number [ ] [ ] [ ] [ ] .

getnameinfo

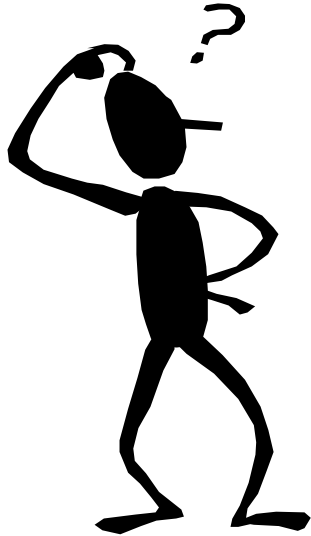
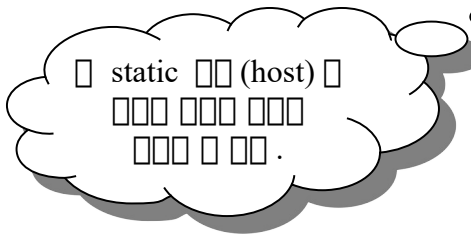
host [ ] service [ ] [ ] name or numeric [ ] [ ] [ ] [ ] [ ] .

Constant	Description
NI_DGRAM	inform datagram service
NI_NAMEREQD	return an error if name cannot be resolved from address
NI_NOFQDN	return only hostname portion of FQDN
NI_NUMERICHOST	return numeric string for hostname
NI_NUMERICSERV	return numeric string for service name

< flags for getnameinfo() >

# Reentrant Functions (1)

## ? Problem



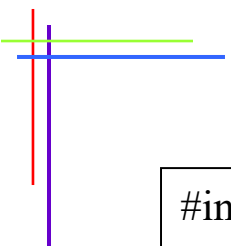
```
static struct hostent host; /* result stored here */
```

```
struct hostent *gethostbyname( const char *hostname )
{
    return( gethostbyname2( hostname, family ) );
}
```

```
struct hostent *gethostbyname2( const char *hostname, int family )
{
    /* call DNS functions for A or AAAA query */
    /* fill in host structure */
    return( &host );
}
```

```
struct hostent *gethostbyaddr( const char *addr, size_t len, int family )
{
    /* call DNS functions for PTR query in in-addr.arpa domain */
    /* fill in host structure */
    return( &host );
}
```



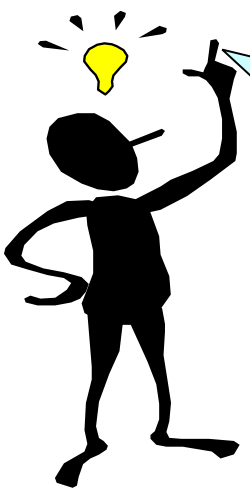


```
#include <netdb.h>

struct hostent *gethostbyname_r( const char *hostname,
                                struct hostent *result, char *buf, int buflen, int *h_errnop );

struct hostent *gethostbyaddr_r( const char *addr, int len, int type,
                                struct hostent *result, char *buf, int buflen, int *h_errnop );

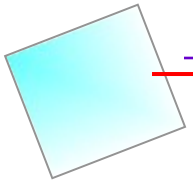
                                both return : nonnull pointer if OK, NULL on error
```

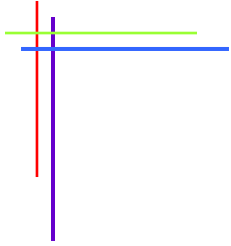


result is a pointer to a struct hostent. The struct hostent is defined in netdb.h. The struct hostent is defined as follows:

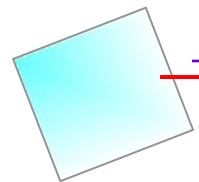
```
struct hostent {
    char **h_name;
    char **h_aliases;
    int h_addrtype;
    char **h_addr_list;
};
```

buf is a buffer for the result. The buffer must be at least 8192 bytes long. The buffer is used to store the result of the lookup. The buffer is used to store the result of the lookup.





# **Chap. 12 *Daemon Processes and inetd Superserver***



# What is Daemon?

- **Daemon is ...**
  - a process that runs in the background
  - independent of control from all terminals
- **Daemon** □ □□□□□ □□
  - □□□□ startup □□ system initialization script □ □□□ □□ (/etc □□ /etc/rc\* □□□□ )
  - inetd superserver □ □□ □□
  - cron daemon □ □□ □□
  - at command □□□ □□□□ □□
  - user terminal □□□ □□

# daemon\_init Function

```

1 #include "unp.h"
2 #include <syslog.h>
3 #define MAXFD 64
4 extern int daemon_proc; /* defined in error.c */
5 void
6 daemon_init(const char *pname, int facility)
7 {
8     int i;
9     pid_t pid;
10    if ( (pid = Fork()) != 0)
11        exit(0); /* parent terminates */
12    /* 1st child continues */
13    setsid(); /* become session leader */
14    Signal(SIGHUP, SIG_IGN);
15    if ( (pid = Fork()) != 0)
16        exit(0); /* 1st child terminates */
17    /* 2nd child continues */
18    daemon_proc = 1; /* for our err_XXX() functions */
19    chdir("/"); /* change working directory */
20    umask(0); /* clear our file mode creation mask */
21    for (i = 0; i < MAXFD; i++)
22        close(i);
23    openlog(pname, LOG_PID, facility);
24 }

```

process is foreground process  
background process

controlling terminal, controlling terminal.

The purpose of this 2<sup>nd</sup> fork is to guarantee that the daemon can't automatically acquire a controlling terminal should it open a terminal device in the future

syslogd logging.





# Example: Daytime Server as a Daemon

```
1 #include "unp.h"
2 #include <time.h>
3 int main(int argc, char **argv)
4 {
5     int     listenfd, connfd;
6     socklen_t  addrlen, len;
7     struct sockaddr *cliaddr;
8     char     buff[MAXLINE];
9     time_t   ticks;
10    daemon_init(argv[0], 0);
11    if (argc == 2)
12        listenfd = Tcp_listen(NULL, argv[1], &addrlen);
13    else if (argc == 3)
14        listenfd = Tcp_listen(argv[1], argv[2], &addrlen);
15    else
16        err_quit("usage: daytimetcpsrv2 [ <host> ] <service or port>");
17    cliaddr = Malloc(addrlen);
18    for ( ; ; ) {
19        len = addrlen;
20        connfd = Accept(listenfd, cliaddr, &len);
21        err_msg("connection from %s", Sock_ntop(cliaddr, len));
22        ticks = time(NULL);
23        snprintf(buff, sizeof(buff), "%.24s\r\n", ctime(&ticks));
24        Write(connfd, buff, strlen(buff));
25        Close(connfd);
26    }
27 }
```

getaddrinfo() + socket() + bind() + listen()

# inetd Daemon

- **inetd daemon = internet superserver**

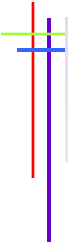
- □□

- inetd □ □□ daemon □□□□ □□□ □□□ □□□□ □□□ □□□□ □□□□
- □□ □□□□□ □□ client request □ □□□ □□□ □□□□ (inetd) □□ □□□□ , □□□□ □□□□ □□ □□□□□ □□

- **/etc/inetd.conf file**

- ftp stream tcp nowait root /usr/sbin/in.ftpd in.ftpd
- telnet stream tcp nowait root /usr/sbin/in.telnetd  
in.telnetd
- login stream tcp nowait root /usr/sbin/in.rlogind  
in.rlogind
- tftp dgram udp wait root /usr/sbin/in.tftpd in.tftpd -s  
/tftpboot

- **/etc/services file**



# inetd



socket()

bind()

listen()  
(if TCP socket)

For each service listed in the /etc/inetd.conf file

select()  
for readability

accept()  
(if TCP socket)

fork()

parent

child

close connected  
socket (if TCP)

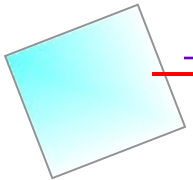
close connected  
socket (if TCP)

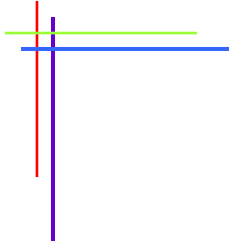
dup socket to  
descriptor 0,1,2;  
close socket

setgid()  
setuid()  
(if user not root)

<http://www.kwang>

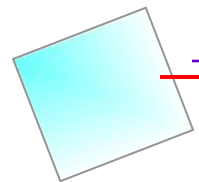
exec() server





# Chap. 13

## *Advanced I/O Functions*



# Socket Timeouts

- **I/O operation timeout**
  - alarm : SIGALRM
  - select : Block waiting for I/O
  - SO\_RCVTIMEO, SO\_SNDTIMEO
- **connect with a Timeout Using SIGALRM**

```
int connect_timeo(int sockfd, const SA *saptr, socklen_t salen, int nsec)
{
    ...
    sigfunc = Signal(SIGALRM, connect_alarm);
    if (alarm(nsec) != 0)
        err_msg("connect_timeo: alarm was already set");

    if ((n = connect(sockfd, (struct sockaddr *) saptr, salen)) < 0) {
        close(sockfd);
        if (errno == EINTR)
            errno = ETIMEDOUT;
    }
    alarm(0); /* turn off the alarm */
    Signal(SIGALRM, sigfunc); /* restore previous signal handler */
    return(n);
}
```

```
static void
connect_alarm(int signo)
{
    return; /* just interrupt the connect() */
}
```

# recv and send Functions

```
#include <sys/socket.h>
ssize_t recv(int sockfd, void *buff, size_t nbyte, int flags);
ssize_t send(int sockfd, const void *buff, size_t nbytes, int flags);
```

Both return: number of bytes read or written if OK, -1 on error

※ flags   read / write

<i>flags</i>	Description	recv	send
<b>MSG_DONTROUTE</b>	<b>bypass routing table lookup</b>		●
<b>MSG_DONTWAIT</b>	<b>only this operation is nonblocking</b>	●	●
<b>MSG_OOB</b>	<b>send or receive out-of-band data</b>	●	●
<b>MSG_PEEK</b>	<b>peek at incoming message</b>	●	
<b>MSG_WAITALL</b>	<b>wait for all the data</b>	●	

# readv and writev / recvmsg and sendmsg Functions

```
#include <sys/socket.h>
ssize_t readv(int filedes, const struct iovec *iov, int iovcnt);
ssize_t writev(int filedes, const struct iovec *iov, int iovcnt);
```

Both return: number of bytes read or written if OK, -1 on error

```
#include <sys/socket.h>
ssize_t recvmsg(int sockfd, struct msghdr *msg, int flags);
ssize_t sendmsg(int sockfd, struct msghdr *msg, int flags);
```

Both return: number of bytes read or written if OK, -1 on error

```
struct iovec {
    void *iov_base; /* starting addr. of buf. */
    size_t iov_len; /* size of buffer */
};
```

```
struct msghdr {
    void *msg_name; /* protocol address */
    socklen_t msg_namelen; /* size of protocol addr. */
    struct iovec *msg_iov; /* scatter/gather array */
    size_t msg_iovlen; /* # elements in msg_iov */
    void *msg_control; /* ancillary data */
    socklen_t msg_controllen; /* len. of ancillary data */
    int msg_flags /* flags returned by recvmsg() */
};
```

# Sockets and Standard I/O

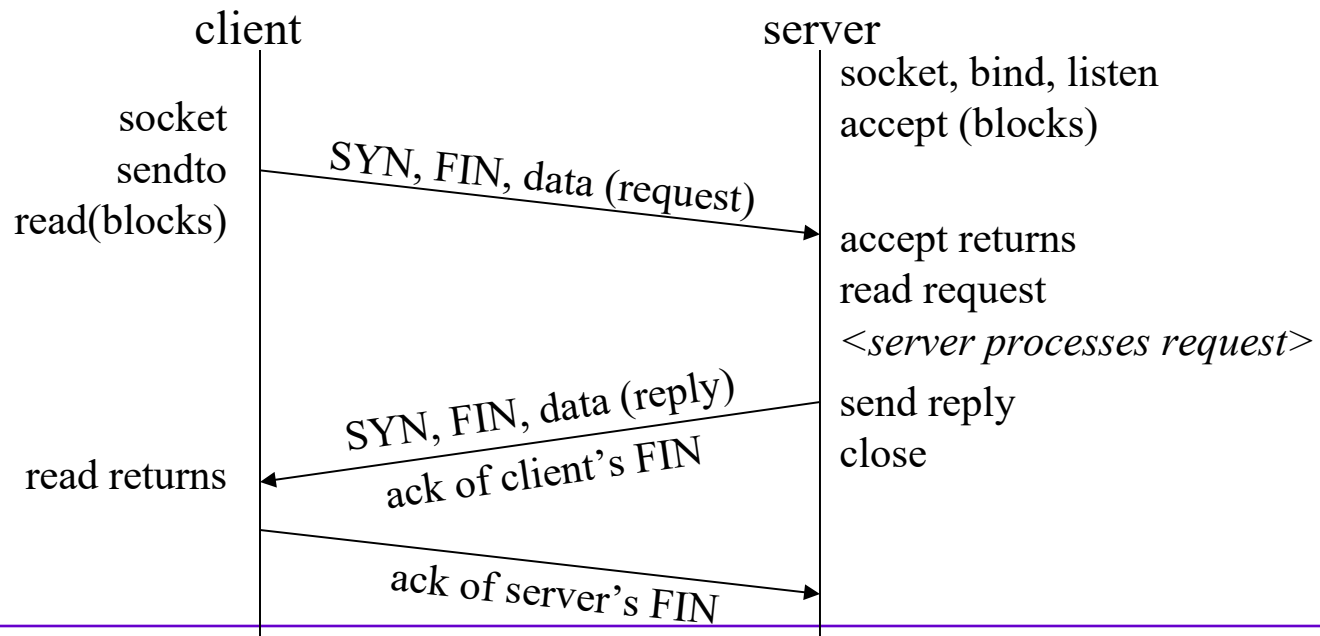
- **Example: str\_echo Function Using Standard I/O**

```
1 #include "unp.h"
2
3 void
4 str_echo(int sockfd)
5 {
6     char    line[MAXLINE];
7     FILE    *fpin, *fpout;
8
9     fpin = Fdopen(sockfd, "r");
10    fpout = Fdopen(sockfd, "w");
11
12    for ( ; ; ) {
13        if (Fgets(line, MAXLINE, fpin) == NULL)
14            return; /* connection closed by other end */
15
16        Fputs(line, fpout);
17    }
18 }
```

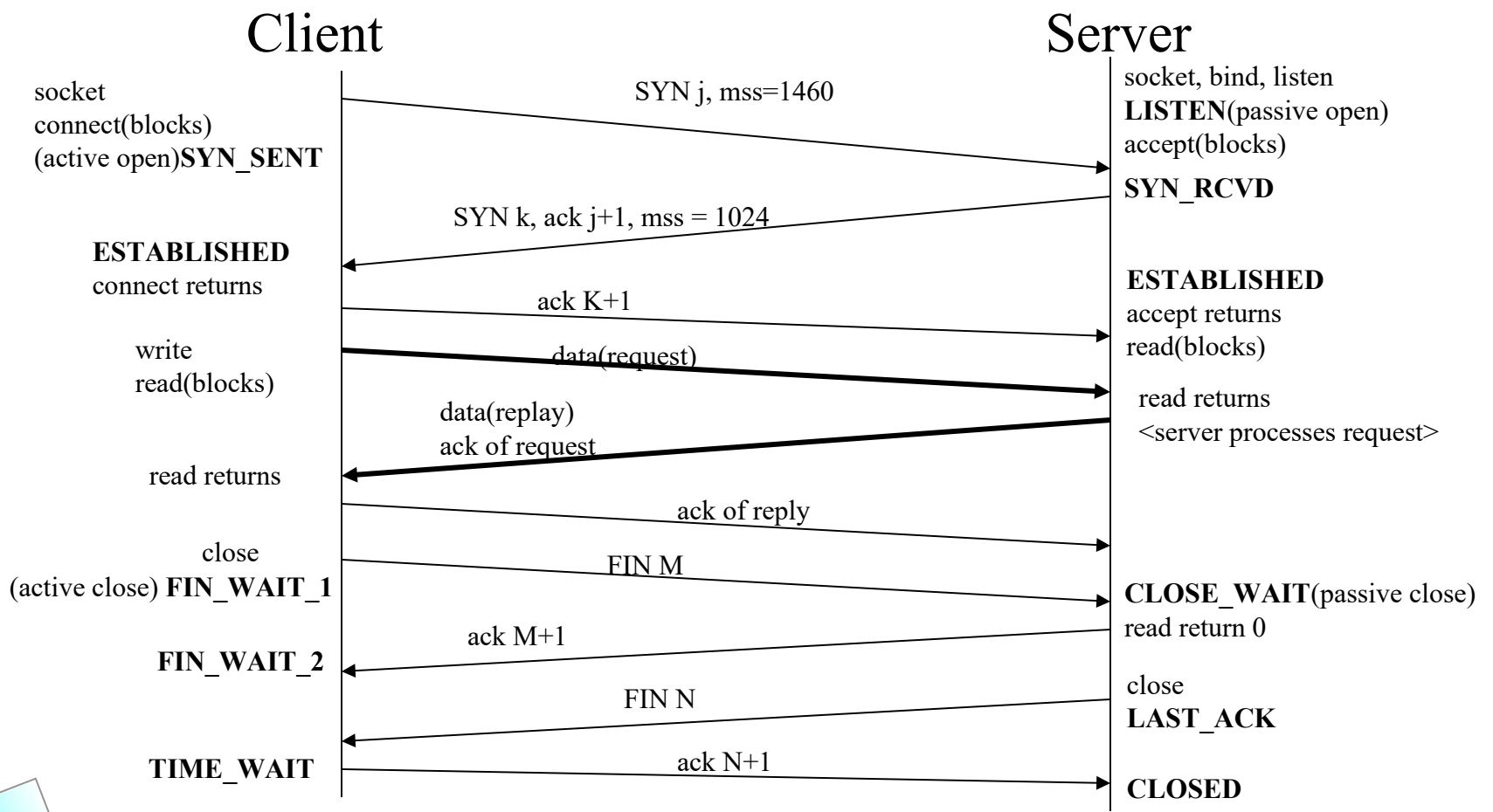


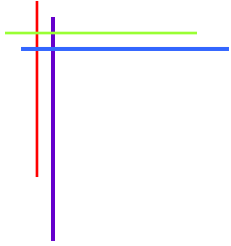
# T/TCP: TCP for Transactions

- can avoid 3-way handshake between hosts that have communicated with each other recently
- can combine the SYN, FIN, and data into a single segment



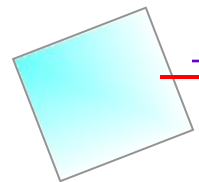
# T/TCP: comparison to general TCP





# Chap. 14

## *Unix Domain Protocols*



# Unix Domain Protocols

- **Unix domain protocol** □ **IPC** □□□□ □□□ □ □□
  - stream socket (similar to TCP)
  - datagram socket (similar to UDP)
- **Unix Domain Socket Address Structure**

```
struct sockaddr_un {  
    uint8_t        sun_len;  
    sa_family_t    sun_family;    /* AF_LOCAL */  
    char           sun_path[104]; /* null-terminated pathname */  
};
```

# Unix Domain Stream Protocol Echo Server

```

1 #include "unp.h"
2 int
3 main(int argc, char **argv)
4 {
5     ... /*   □□ □□□ □□ □□ */
10  listenfd = Socket(AF_LOCAL, SOCK_STREAM, 0);
11  unlink(UNIXSTR_PATH);
12  bzero(&servaddr, sizeof(servaddr));
13  servaddr.sun_family = AF_LOCAL;
14  strcpy(servaddr.sun_path, UNIXSTR_PATH);
15  Bind(listenfd, (SA *) &servaddr, sizeof(servaddr));
16  Listen(listenfd, LISTENQ);
17  Signal(SIGCHLD, sig_chld);
18  for ( ; ) {
19      clen = sizeof(cliaddr);
20      if ( (connfd = accept(listenfd, (SA *) &cliaddr, &clen)) < 0 ) {
21          if (errno == EINTR)
22              continue; /* back to for() */
23          else
24              err_sys("accept error");
25      }
26      if ( (childpid = Fork()) == 0 ) { /* child process */
27          Close(listenfd); /* close listening socket */
28          str_echo(connfd); /* process the request */
29          exit(0);
30      }
31      Close(connfd); /* parent closes connected socket */
32  }
33 }

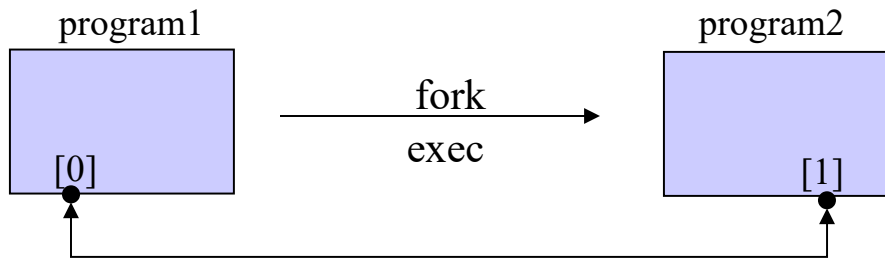
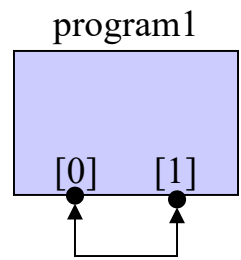
```

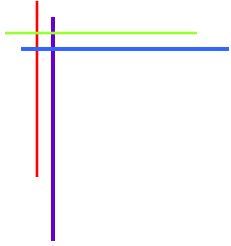
# Unix Domain Stream Protocol Echo Client

```
1 #include "unp.h"
2
3 int
4 main(int argc, char **argv)
5 {
6     int          sockfd;
7     struct sockaddr_un servaddr;
8
9     sockfd = Socket(AF_LOCAL, SOCK_STREAM, 0);
10
11     bzero(&servaddr, sizeof(servaddr));
12     servaddr.sun_family = AF_LOCAL;
13     strcpy(servaddr.sun_path, UNIXSTR_PATH);
14
15     Connect(sockfd, (SA *) &servaddr, sizeof(servaddr));
16
17     str_cli(stdin, sockfd); /* do it all */
18
19     exit(0);
20 }
```

# Unix Domain Protocol Examples

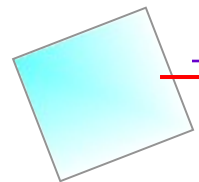
- **Passing Descriptors**





# Unix Network Programming

## Chapter 15 Nonblocking I/O



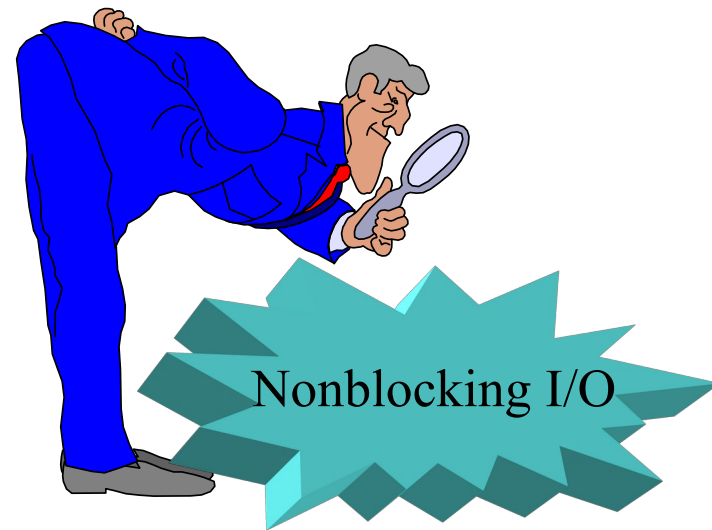


# Contents

- **Introduction**
- **Nonblocking Reads and Writes :**
  - str\_cli function
- **Nonblocking connect**
  - Daytime Client
  - Web Client
  - accept
- **Summary**

# Introduction

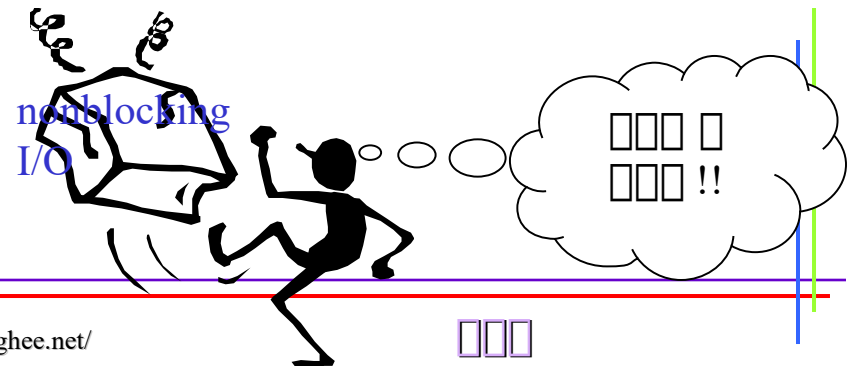
- **By default, sockets are blocking.**
- **blocking**
  - Input Operations
    - † `read()`, `readv()`, `recv()`, `recvfrom()`
  - Output operations
    - † `write()`, `writv()`, `send()`, `sendto()`
  - Accepting incoming connections
    - † `accept()`
  - Initiating outgoing connections
    - † `connect()`
- **Chapter 6 : various I/O model was learned**



# Nonblocking Read and Writes

- **Goal in this section is to develop a version of this function that uses nonblocking I/O**
- **Unfortunately the addition of nonblocking I/O complicates the function's buffer management noticeably.**
- **Is it worth the effort to code an application using nonblocking I/O, given the complexity of the resulting code?**

Writer recommend the simple Approach !!!!

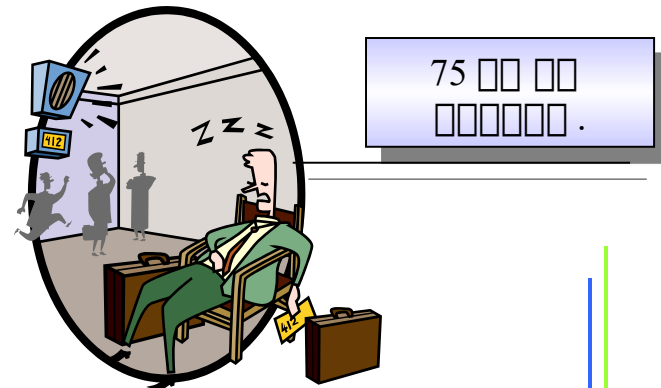


# Nonblocking connect

- **There are 3 uses for nonblocking connect**
  - Overlap other processing with the 3-way handshaking
  - Establish multiple connections at the same time
  - Specify a time limit for *select()*

Must handled !!!

- Even though the socket is nonblocking, if the server to which we are connecting is on the same host, the connection establishment normally takes place immediately when we call `connect()`
- Berkeley-derived implementations have the following 2 rules regarding `select()` and nonblocking connects (1) when the connection completes successfully, the descriptor becomes writable, and (2) when the connection establishment encounters an error, the descriptor becomes both readable and writable



# Daytime Client

```

1. int connect_nonb(int sockfd, const SA *saptr, socklen_t salen, int nsec)
2. { int flags, n, error; socklen_t len; fd_set rset, wset; struct
   timeval tval;
3.   flags = Fcntl(sockfd, F_GETFL, 0);
4.   Fcntl(sockfd, F_SETFL, flags | O_NONBLOCK);
5.   error = 0;
6.   if ( ( n = connect(sockfd, (struct sockaddr *) saptr, salen)) < 0 )
7.     if (errno != EINPROGRESS)
8.       return(-1);
9.   /* Do whatever we want while the connect is taking place. */
10.  if (n == 0)
11.    goto done; /* connect completed immediately */
12.  FD_ZERO(&rset);
13.  FD_SET(sockfd, &rset);
14.  wset = rset;
15.  tval.tv_sec = nsec;
16.  tval.tv_usec = 0;
17.  if ( ( n = Select(sockfd+1, &rset, &wset, NULL, nsec ? &tval : NULL)) == 0 ) {
18.    close(sockfd); /* timeout */
19.    errno = ETIMEDOUT;
20.    return(-1); }
21.  if (FD_ISSET(sockfd, &rset) || FD_ISSET(sockfd, &wset)) {
22.    len = sizeof(error);
23.    if (getsockopt(sockfd, SOL_SOCKET, SO_ERROR, &error, &len) < 0)
24.      return(-1); /* Solaris pending error */ } else
25.    err_quit("select error: sockfd not set");
26. done:
27.  Fcntl(sockfd, F_SETFL, flags); /* restore file status flags */
28.  if (error) {
29.    close(sockfd); /* just in case */
30.    errno = error;
31.    return(-1); }
32.  return(0); }

```

Set socket nonblocking

Overlap processing with connection establishment

Check for immediate completion

Call select

Check for readability or writability

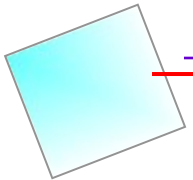
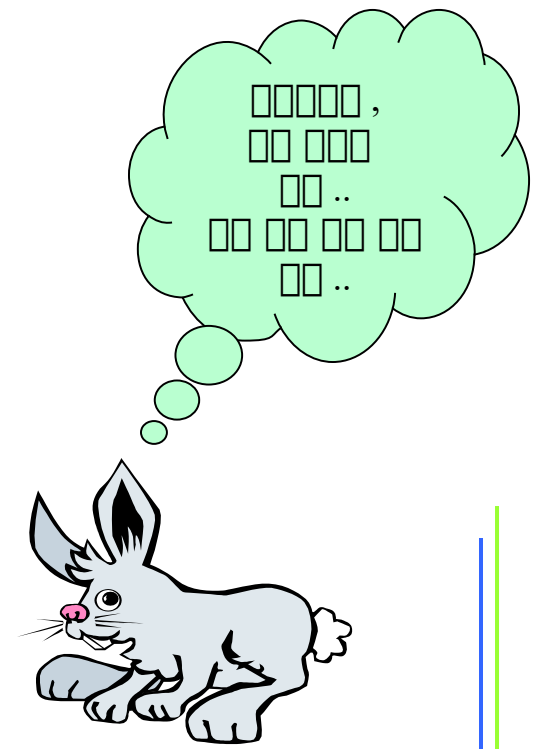
Turn off nonblocking and return



# WebClient

- Performance of simulation Connections

# simultaneous connections	clock time (sec) nonblocking	Clock time threads
1	6.0	6.3
2	4.1	4.2
3	3.0	3.1
4	2.8	3.0
5	2.5	2.7
6	2.4	2.5
7	2.3	2.3
8	2.2	2.3
9	2.0	2.2



# Nonblocking accept

```

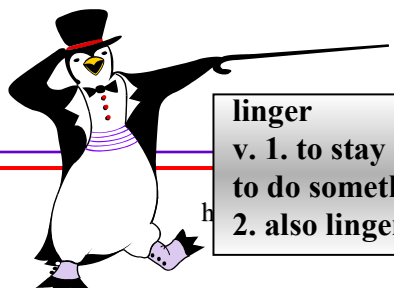
1. #include "unp.h"
2. int
3. main(int argc, char **argv)
4. {
5.     int     sockfd;
6.     struct linger  ling;
7.     struct sockaddr_in servaddr;
8.     if (argc != 2)
9.         err_quit("usage: tcpcli <IPaddress>");
10.    sockfd = Socket(AF_INET, SOCK_STREAM, 0);
11.    bzero(&servaddr, sizeof(servaddr));
12.    servaddr.sin_family = AF_INET;
13.    servaddr.sin_port = htons(SERV_PORT);
14.    Inet_pton(AF_INET, argv[1], &servaddr.sin_addr);
15.    Connect(sockfd, (SA *) &servaddr, sizeof(servaddr));
16.    ling.l_onoff = 1; /* cause RST to be sent on close() */
17.    ling.l_linger = 0;
18.    Setsockopt(sockfd, SOL_SOCKET, SO_LINGER, &ling, sizeof(ling));
19.    Close(sockfd);
20.    exit(0);
21. }

```

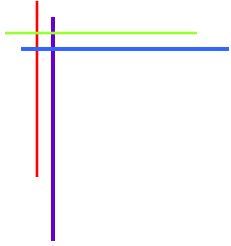


tcp 3way handshake, nonblocking I/O socket option

Set SO\_LINGER socket Option

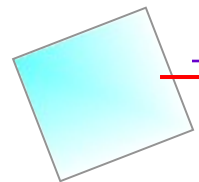


linger  
 v. 1. to stay somewhere longer or to take longer to do something than usual  
 2. also linger on to be slow to disappear



# Unix Network Programming

## Chapter 16 ioctl Operation



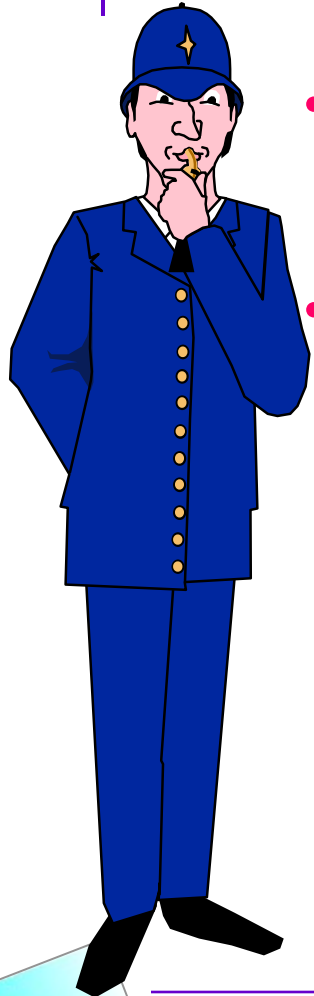


# Contents

- **Introduction**
- **ioctl Function**
- **Socket Operations**
- **File Operations**
- **Interface Configuration**
- **get\_fif\_info Function**
- **Interface Operations**
- **ARP Cache Operations**
- **Routing Table Operations**
- **Summary**



# Introduction



- **ioctl() has traditionally been the system interface used for everything that didn't fit into some other nicely defined category.**
- **This remain for implementation-dependent features related to network programming**
  - obtaining the interface information
  - accessing the the routing table
  - accessing the ARP cache

# ioctl Function

```
#include <unistd.h>
```

```
int ioctl(int fd, int request, .../*void *arg */);
```

Returns : 0 if OK, -1 on error

- **socket operations**
- **file operations**
- **interface operations**
- **ARP cache operations**
- **routing table operations**
- **stream system(Chapter 33)**



# Socket Operation

- **SIOCATMARK**
  - socket's read pointer is currently at the out-of-band mark, or a zero value if the read pointer is not at the out-of-band mark.
- **SIOCGPRP**
  - PID or the process group ID
- **SIOCSPGRP**
  - set either the PID or the process group ID

chapter 7  
fcntl() □□

# File Operations

- **FIONBIO**
  - nonblocking flag for the socket is cleared or turned on, depending whether the third argument to `ioctl()` points to a zero or nonzero value
- **FIOASYNC**
  - governs the receipt of asynchronous I/O signals(SIGIO)
- **FIONREAD**
  - the number of bytes in the socket receive buffer
- **FIOSETOWN**
  - same as `SIOCSPGRP`
- **FIOGETOWN**
  - same as `SIOCGPGRP`

# Interface Operations

- **SIOCGIFADDR**
  - Return the unicast address
- **SIOCSIFADDR**
  - Sets the interface address
- **SIOCGIFFLAGS**
  - Return the interface flags
- **SIOCSIFFLAGS**
  - Set the interface flags
- **SIGIFDSTTADDR**
  - Return the point-to-point address
- **SICSIFDSTADDR**
  - Set the point-to-point address
- **SIOCSIFBRDADDR**
  - Set the broadcast address
- **SIOCGIFNETMASK**
  - Return the subnet mask
- **SIOCSIFNETMASK**
  - Set the subnet mask
- **SIOCGIFMETRIC**
  - Return the interface metric
- **SIOCIFMETRIC**
  - Set the interface routing metric



# ARP Cache Operations

```

1. #include <net/if_arp.h>
2. int main(int argc, char **argv)
3. {
4.     int    family, sockfd;
5.     char   str[INET6_ADDRSTRLEN];
6.     char   **pptr;
7.     unsigned char *ptr;
8.     struct arpreq arpreq;
9.     struct sockaddr_in *sin;
10.    pptr = my_addrs(&family);
11.    for ( ; *pptr != NULL; pptr++) {
12.        printf("%s: ", Inet_ntop(family, *pptr, str, sizeof(str)));
13.        switch (family) {
14.            case AF_INET:
15.                sockfd = Socket(AF_INET, SOCK_DGRAM, 0);
16.                sin = (struct sockaddr_in *) &arpreq.arp_pa;
17.                bzero(sin, sizeof(struct sockaddr_in));
18.                sin->sin_family = AF_INET;
19.                memcpy(&sin->sin_addr, *pptr, sizeof(struct in_addr));
20.                ioctl(sockfd, SIOCGARP, &arpreq);
21.                ptr = &arpreq.arp_ha.sa_data[0];
22.                printf("%x:%x:%x:%x:%x:%x\n", *ptr, *(ptr+1),
23.                    *(ptr+2), *(ptr+3), *(ptr+4), *(ptr+5));
24.                break;
25.            default:
26.                err_quit("unsupported address family: %d", family);
27.        } } exit(0);}

```

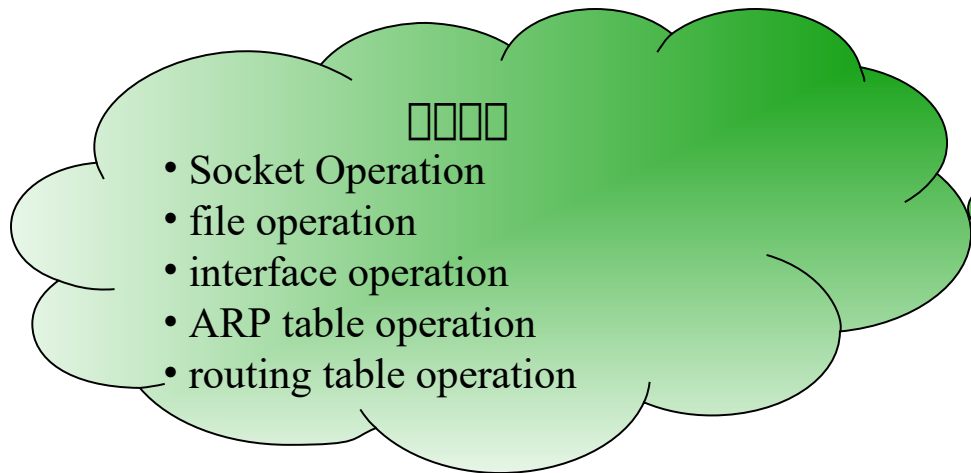
get list of address and loop through each one

Print IP address

Issue ioctl and print hardware address

# Routing Table Operations

- **TWO ioctl() requests provided**
  - SIOCADDRT
    - † Add an entry to the routing table
  - SIOCDELRT
    - † Delete an entry from the routing table







# UNIX Network Programming

( chap 17 - chap 18 )

## Chap.17 - Routing Socket

- **Routing table** □□ □□ □□
- **Message Types**
- **Structure Types**
- **Constants in routing messages**
- **RMT\_GET example**
- **sysctl Operations**

# We must study ... ?



📖 Network      Routing table   
  (2  )

1. Routing socket  (superuser)

2. *sysctl*

## Routing table □□ □□ □□

1. Use ioctl command
2. Use netstat program
3. Use routing daemon ( only superuser )
  - monitor ICMP redirection message by creating a raw ICMP socket
  - routing domain □□ □□□□ □□□ socket type □ raw socket □□ .
4. Use sysctl command

# Message Types

☞ Routing table □ □□ / □□ □□□ Message type - <net/route.h>

Message type	To kernel	From kernel	Description	Structure type
RTM_ADD	•	•	add route	rt_msghdr
RTM_CHANGE	•	•	change gateway, metrics, or flags	rt_msghdr
RTM_DELADDR		•	address being removed from interface	ifa_msghdr
RTM_DELETE	•	•	delete route	rt_msghdr
RTM_GET	•	•	report metrics and other route information	rt_msghdr
RTM_IFINFO		•	interface going up, down etc.	if_msghdr
RTM_LOCK	•	•	lock specified metrics	rt_msghdr
RTM_LOSING		•	kernel suspects route is failing	rt_msghdr
RTM_MISS		•	lookup failed on this address	rt_msghdr
RTM_NEWADDR		•	address being added to interface	ifa_msghdr
RTM_REDIRECT		•	kernel told to use different route	rt_msghdr
RTM_RESOLVE		•	request to resolve destination to link layer address	rt_msghdr

# Structure Types

```

typedef struct rt_msghdr {
    ushort_t rtm_msglen; /* to skip over non-understood messages */
    uchar_t rtm_version; /* future binary compatibility */
    uchar_t rtm_type; /* message type */
    ushort_t rtm_index; /* index for associated ifp */
    int rtm_flags; /* flags, incl. kern & message, e.g. DONE */
    int rtm_addr; /* bitmask identifying sockaddrs in msg */
    pid_t rtm_pid; /* identify sender */
    int rtm_seq; /* for sender to identify action */
    int rtm_errno; /* why failed */
    int rtm_use; /* from rtenry */
    uint_t rtm_inits; /* which metrics we are initializing */
    struct rt_metrics rtm_rmx; /* metrics themselves */
} rt_msghdr_t;

```

```

typedef struct if_msghdr {
    ushort_t ifm_msglen; /* to skip over non-understood messages */
    uchar_t ifm_version; /* future binary compatibility */
    uchar_t ifm_type; /* message type */
    int ifm_addr; /* like rtm_addr */
    int ifm_flags; /* value of if_flags */
    ushort_t ifm_index; /* index for associated ifp */
    struct if_data ifm_data; /* statistics and other data about if */
} if_msghdr_t;

```

```

typedef struct ifa_msghdr {
    ushort_t ifam_msglen; /* to skip over non-understood messages */
    uchar_t ifam_version; /* future binary compatibility */
    uchar_t ifam_type; /* message type */
    int ifam_addr; /* like rtm_addr */
    int ifam_flags; /* route flags */
    ushort_t ifam_index; /* index for associated ifp */
    int ifam_metric; /* value of ipif_metric */
} ifa_msghdr_t;

```

# Constants in routing messages

Bitmask		Array index		Socket address structure containing
constant	value	constant	value	
<i>RTA_DST</i>	0x01	<i>RTAX_DST</i>	0	destination address
<i>RTA_GATEWAY</i>	0x02	<i>RTAX_GATEWAY</i>	1	gateway address
<i>RTA_NETMASK</i>	0x04	<i>RTAX_NETMASK</i>	2	network mask
<i>RTA_GENMASK</i>	0x08	<i>RTAX_GENMASK</i>	3	cloning mask
<i>RTA_IFP</i>	0x10	<i>RTAX_IFP</i>	4	interface name
<i>RTA_IFA</i>	0x20	<i>RTAX_IFA</i>	5	interface address
<i>RTA_AUTHOR</i>	0x40	<i>RTAX_AUTHOR</i>	6	author of redirect
<i>RTA_BRD</i>	0x80	<i>RTAX_BRD</i>	7	broadcast or point-to-point destination address
		<i>RTAX_MAX</i>	8	Max #elements



# RMT\_GET example (1)

```

#include "unproute.h"
#define BUFLen (sizeof(struct rt_msghdr) + 512)
/* 8 * sizeof(struct sockaddr_in6) = 192 */
#define SEQ 9999

int
main(int argc, char **argv)
{
    int sockfd;
    char *buf;
    pid_t pid;
    ssize_t n;
    struct rt_msghdr *rtm;
    struct sockaddr *sa, *rti_info[RTAX_MAX];
    struct sockaddr_in *sin;

    if (argc != 2)
        err_quit("usage: getrt <IPaddress>");

sockfd = Socket(AF_ROUTE, SOCK_RAW, 0); /* need superuser privileges */

    buf = Calloc(1, BUFLen); /* and initialized to 0 */

    rtm = (struct rt_msghdr *) buf;
    rtm->rtm_msglen = sizeof(struct rt_msghdr) + sizeof(struct sockaddr_in);
    rtm->rtm_version = RTM_VERSION;
    rtm->rtm_type = RTM_GET;
    rtm->rtm_addrs = RTA_DST;
    rtm->rtm_pid = pid = getpid();
    rtm->rtm_seq = SEQ;

    sin = (struct sockaddr_in *) (rtm + 1);
    sin->sin_family = AF_INET;
    Inet_pton(AF_INET, argv[1], &sin->sin_addr);

    Write(sockfd, rtm, rtm->rtm_msglen);

    do {
        n = Read(sockfd, rtm, BUFLen);
    } while (rtm->rtm_type != RTM_GET || rtm->rtm_seq != SEQ ||
             rtm->rtm_pid != pid);

    rtm = (struct rt_msghdr *) buf;
    sa = (struct sockaddr *) (rtm + 1);
    get_rtaddrs(rtm->rtm_addrs, sa, rti_info);
    if (sa == rti_info[RTAX_DST]) != NULL)
        printf("dest: %s\n", Sock_ntop_host(sa, sa->sa_len)); /* sa_len = 16 or 24 */

    if ((sa == rti_info[RTAX_GATEWAY]) != NULL)
        printf("gateway: %s\n", Sock_ntop_host(sa, sa->sa_len)); /* sa_len = 16 or 24 */

    if ((sa == rti_info[RTAX_NETMASK]) != NULL)
        printf("netmask: %s\n", Sock_masktop(sa, sa->sa_len)); /* sa_len = 0,5,6,7,8 */

    if ((sa == rti_info[RTAX_GENMASK]) != NULL)
        printf("genmask: %s\n", Sock_masktop(sa, sa->sa_len)); /* sa_len = 0,5,6,7,8 */

    exit(0);
}

```

Page 454 □□

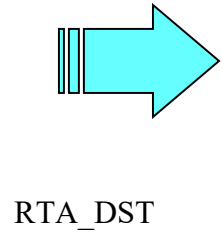


145 **RMT\_GET example (2)**

```
Example >
bsdi # getrt 206.62.226.32
dest : 206.62.226.32
gateway : AF_LINK, index=2
netmask : 255.255.255.224
```

buffer sent  
to kernel

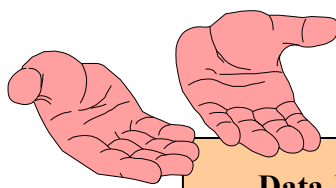
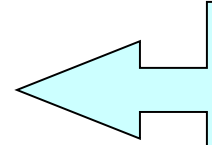
```
rt_msghdr{
  rtm_type =
    RTM_GET
  destination
  socket address
  structure
```



buffer returned  
by kernel

```
rt_msghdr{
  rtm_type =
    RTM_GET
  destination
  socket address
  structure
  gateway
  socket address
  structure
  netmask
  socket address
  structure
  genmask
  socket address
  structure
```

RTA\_DST  
RTA\_GATEWAY  
RTA\_NETMASK  
RTA\_GENMASK



**Data-link socket address structure**

```
struct sockaddr_dl {
  uint8_t    sdl_len;
  sa_family_t sdl_family; /* AF_LINK */
  uint16_t   sdl_index;
  uint8_t    sdl_type;
  uint8_t    sdl_nlen;
  uint8_t    sdl_alen;
  uint8_t    sdl_slen;
  char       sdl_data[12];
}
```

Port

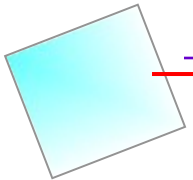
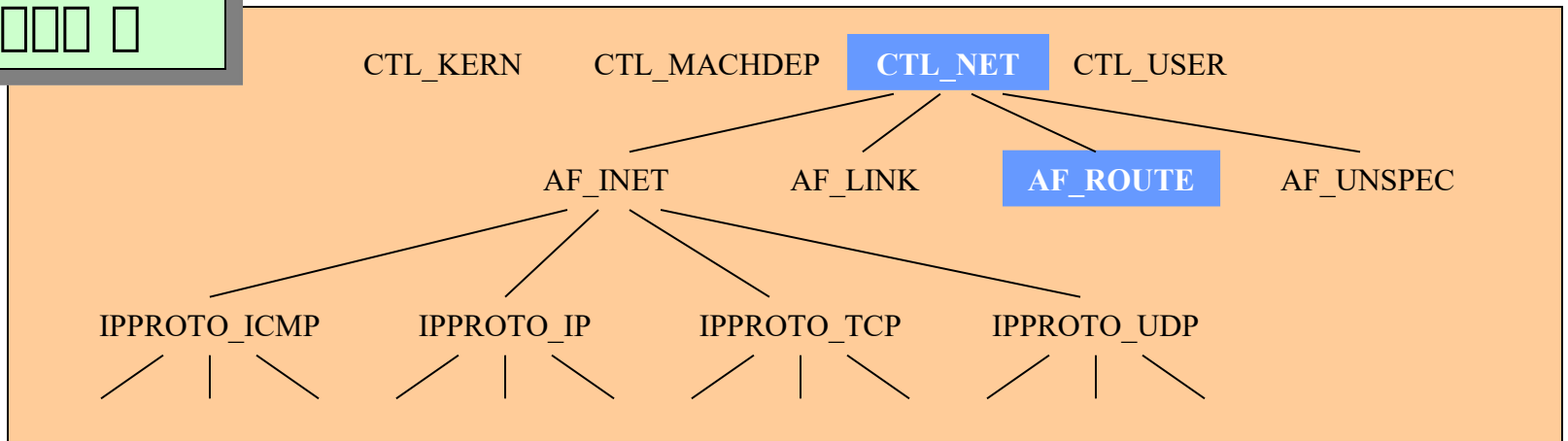
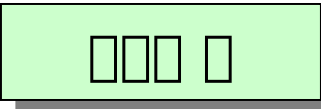
# sysctl Operations (1)

- ▶ routing table □ interface list □□□ □□□□ .
- ▶ □□ process □ □□ □□□□ .  
 ↔ cf) routing socket □□□ superuser □□ □□□□ .

```
#include <sys/param.h>
#include <sys/sysctl.h>

int sysctl( int *name, u_int namelen, void *oldp, size_t *oldlenp,
           void *newp, size_t newlen );

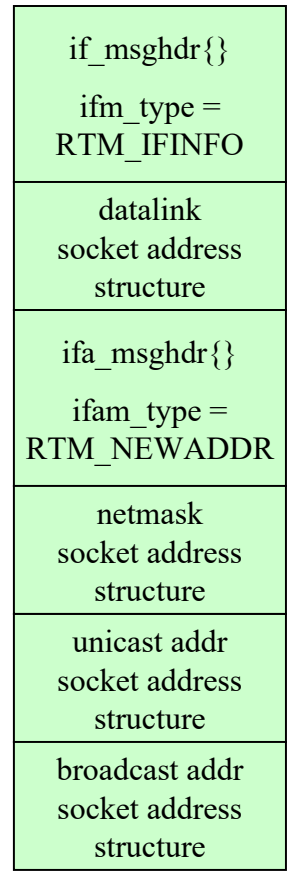
return : 0 if OK, -1 on error
```



# sysctl Operations (2)

name[]	Return IPv4 routing table	Return IPv4 ARP cache	Return IPv4 interface list
0	CTL_NET	CTL_NET	CTL_NET
1	AF_ROUTE	AF_ROUTE	AF_ROUTE
2	0	0	0
3	AF_INET	AF_INET	AF_INET
4	NET_RT_DUMP	NET_RT_FLAGS	NET_RT_IFLIST
5	0	RTF_LLINFO	0

buffer returned by kernel



One per interface : interface name, index and hardware address

One per address configured for the interface










```

mib[0] = CTL_NET;
mib[1] = AF_ROUTE;
mib[2] = 0;
mib[3] = family;
mib[4] = NET_RT_IFLIST;
mib[5] = flag /* interface index or 0*/
sysctl( mib, 6, buff, lenp, NULL, 0 );

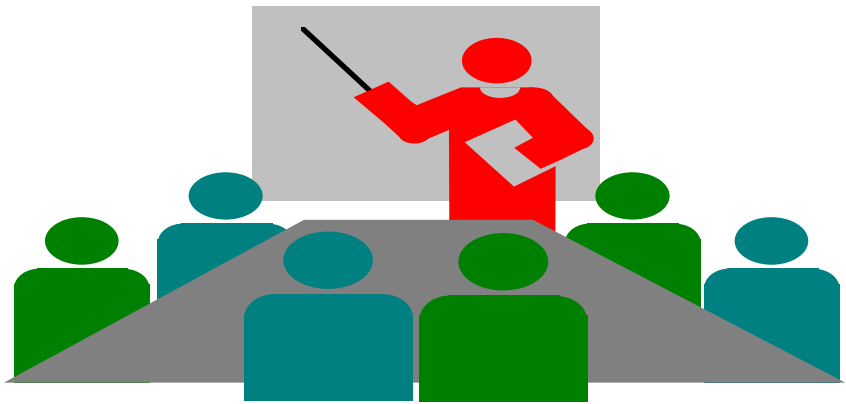
```

# Contents

## Chap.18 - Broadcast

-  **Broadcast support and address**
-  **Broadcast datagram flow example**
-  **Race Conditions**
  -  **blocking and unblocking the signal**
  -  **blocking and unblocking the signal with  
pselect**
  -  **using sigsetjmp and siglongjmp**
  -  **using IPC(Interprocess communication)  
from signal handler to function**

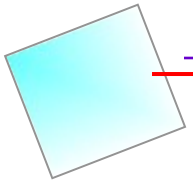
# We must study ... ?



📖 Unicast □ Broadcast □ □□

📖 Broadcast □ □□□ □□

📖 Race condition □ □□□ □ □□□ (4 □□ )



# Broadcast support and address



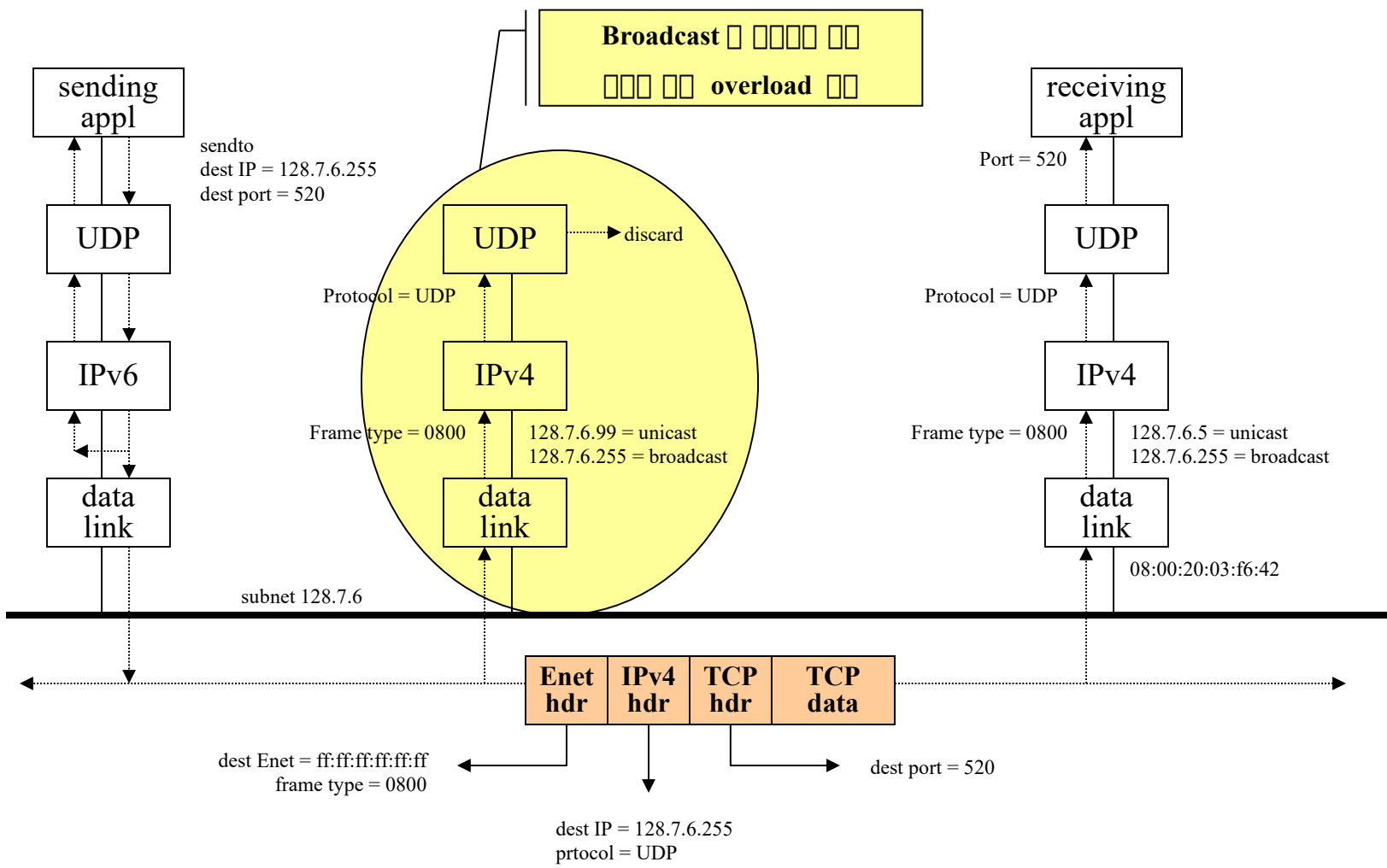
- 1. Multicast □□□ IPv4 □□□ optional □□□ IPv6 □□□ □□□□□ .
- 2. Broadcast □□□ IPv6 □□□ □□□□ □□□ . IPv6 □ □□ broadcast □ multicast □ □□ □□□ □□□ □□ .
- 3. Broadcast □ multicast □ UDP □ □□□□ TCP □□□ □ □□□ □□□ .

All bits '1'



- 1. Subnet-directed broadcast address : { netid, subnetid, -1 }
- 2. All-subnets-directed broadcast address : { netid, -1, -1 }
- 3. Network-directed broadcast address : { netid, -1 } → Subnet □ □□ □□
- 4. Limited broadcast address : { -1, -1, -1 } → Router □□ □ □□ 1 □□□ □□□□ □□□□ □□□ □□ .

# Broadcast datagram flow example



## Example continue ...

[54 feel:kjschaos ~/mywork/network/unpv12e/bcast ] udpcli03 150.150.55.127

hi

from 150.150.55.52: Mon Jan 10 18:01:23 2000

from 150.150.55.53: Mon Jan 10 18:04:38 2000

from 150.150.55.55: Mon Jan 10 18:05:48 2000

from 150.150.55.99: Mon Jan 10 17:50:41 2000

from 150.150.55.10: ☐☐ 6:02:48 2000-01-10

from 150.150.55.54: Mon Jan 10 18:06:23 2000

[55 feel:kjschaos ~/mywork/network/unpv12e/bcast ] udpcli03 255.255.255.255

hi

from 150.150.55.52: Mon Jan 10 18:01:42 2000

from 150.150.55.53: Mon Jan 10 18:04:57 2000

from 150.150.55.55: Mon Jan 10 18:06:07 2000

from 150.150.55.99: Mon Jan 10 17:51:00 2000

from 150.150.55.10: ☐☐ 6:03:07 2000-01-10

from 150.150.55.100: Mon Jan 10 06:20:44 2000

from 150.150.55.54: Mon Jan 10 18:06:42 2000

[56 feel:kjschaos ~/mywork/network/unpv12e/bcast ] udpcli03 150.150.55.52

hi

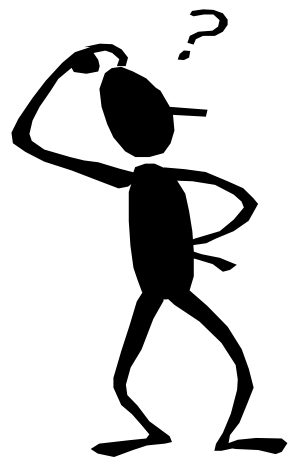
from 150.150.55.52: Mon Jan 10 18:02:01 2000



# Race Conditions (1)

## When happen ?

- ☛ process 访问 data 时
- ☛ Signal 被阻塞



## Why concern ?

? Broadcast 消息 主机 信号

阻塞 阻塞 阻塞 blocking 阻塞

阻塞 阻塞 阻塞 .

# Race Conditions (2)

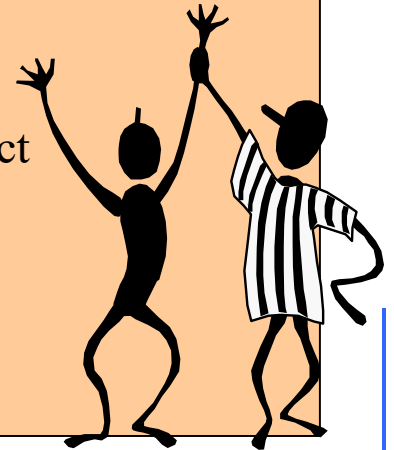
## Solution

🚗 `process` `data`

- 1) mutual exclusion variables( `mutex` )
- 2) condition variables

🚗 Signal `signal`

- 1) blocking and unblocking the signal
- 2) blocking and unblocking the signal with `pselect`
- 3) using `sigsetjmp` and `siglongjmp`
- 4) using IPC(Interprocess communication)  
from signal handler to function



# Race Conditions - Signal case (1)

```

void dg_cli( ... )
{
  Setsockopt( sockfd, SOL_SOCKET, SO_BROADCAST, &on, sizeof(on) );
  ...
  Sigemptyset( &sigset_alm );
  Sigaddset( &sigset_alm, AIGALRM );

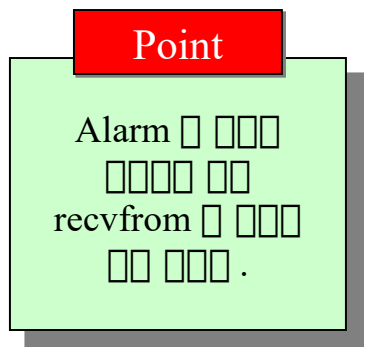
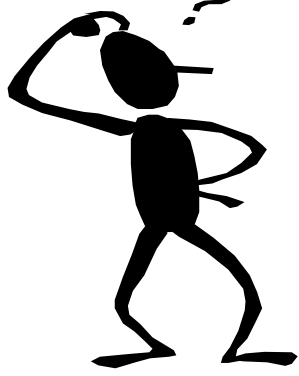
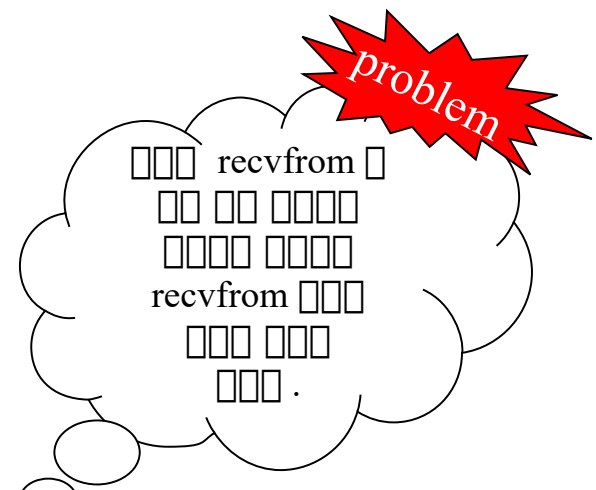
  Signal( SIGALRM, recvfrom_alm );

  while ( ... ) {
    Sendto( sockfd, ... )

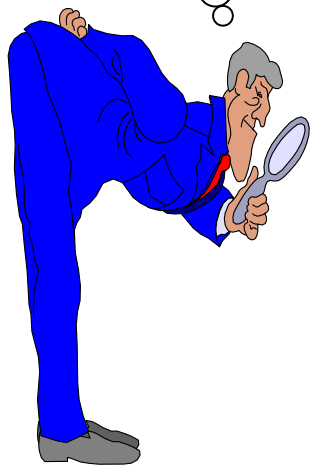
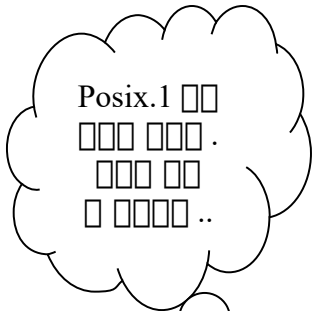
    alarm(5);
    for ( ; ; ) {
      Sigprocmask( SIG_UNBLOCK, &sigset_alm, NULL );
      n = recvfrom( sockfd, recvline, ... );
      Sigprocmask( SIG_BLOCK, &sigset_alm, NULL );
      ...
    }
  }
}

static void recvfsfrom_alm( int signo )
{
  return;
}

```



# Race Conditions - Signal case (2)



```

void dg_cli( ... )
{
    Setsockopt( sockfd, SOL_SOCKET, SO_BROADCAST, &on, sizeof(on) );
    ...
    Sigemptyset( &sigset_alarm ); Sigemptyset( &sigset_alarm );
    Sigaddset( &sigset_alarm, AIGALRM );

    Signal( SIGALRM, recvfrom_alarm );

    while ( ... ) {
        Sendto( sockfd, ... )
        Sigprocmask( SIG_BLOCK, &sigset_alarm, NULL );

        alarm(5);
        for ( ; ; ) {
            FD_SET( sockfd, &rset );
            n = pselect( sockfd+1, &rset, NULL, NULL, NULL, &sigset_empty );
            ...
        }
    }
}

int pselect( ... )
{
    sigprocmask( SIG_SETMASK, signask, &savemask ); /* caller's mask */
    n = select( ... );
    sigprocmask( SIG_SETMASK, &savemask, NULL ); /* restore maks */
}

```

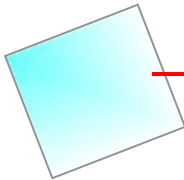
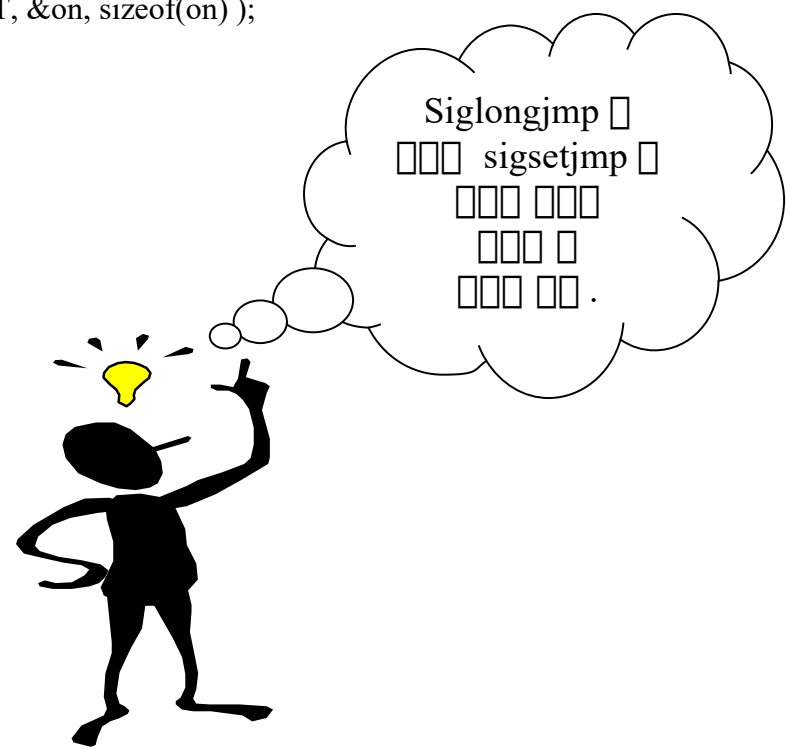
# Race Conditions - Signal case (3)

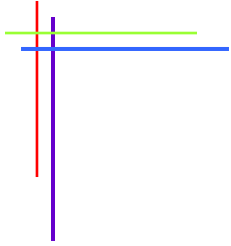
```

void dg_cli( ... )
{
  Setsockopt( sockfd, SOL_SOCKET, SO_BROADCAST, &on, sizeof(on) );
  Signal( SIGALRM, recvfrom_alarm );
  while ( ... ) {
    Sendto( sockfd, ... )
    alarm(5);
    for ( ; ; ) {
      if ( sigsetjmp( jmpbuf, 1 ) != 0 );
        break;
      n = Recvfrom( sockfd, recvline, ... );
      ...
    }
  }
}

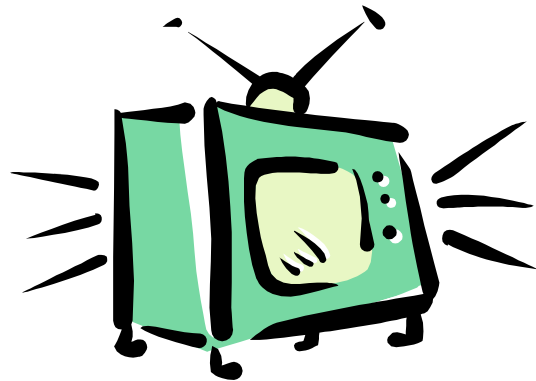
static void recvfrom_alarm( int signo )
{
  siglongjmp( jmpbuf, 1 );
}

```

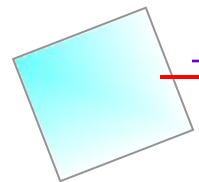




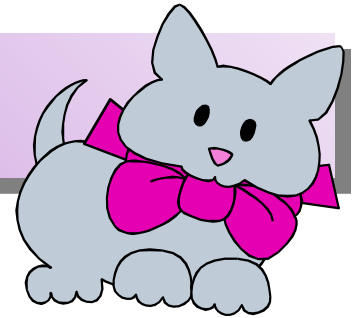
# Unix Network Programming



## Chapter 19. Multicasting



# Contents



- **Introduction**
- ***Multicast Address***
- ***Multicasting vs Broadcasting on A LAN***
- **Multicast Socket Options**
- *mcast\_join()* and Related Functions
- *dg\_cli()* Functions Using Multicasting
- **Receiving MBone Session Announcements**
- **Sending and Receiving**
- **SNTP : Simple Network Time Protocol**
- **SNTP(Continued)**
- **Summary**



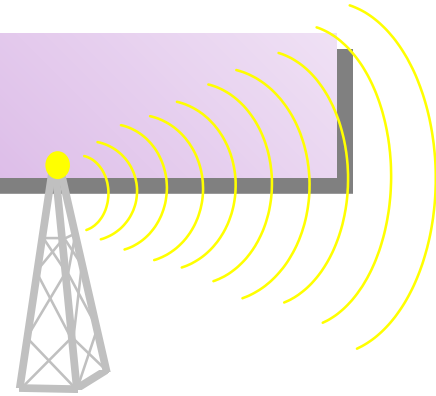
# Introduction



- **Broadcasting is normally limited to a LAN**
- **Multicasting can be used on a LAN or across a WAN.**
- **Five socket options**
  - 3 that affect the sending of UDP datagrams to a multicast address, and
  - 2 that affect the host's receptions of multicast datagrams



# Multicast Address

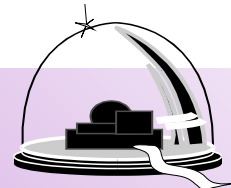


- **IPv4 Class D Address**

- range : 224.0.0.0 ~ 239.255.255.255
- low order bit 28 bits : multicast group ID
- 32-bit address : group address
- IPv4 Ethernet multicast address : *01:00:5e*

- **IPv6 Multicast Address**

- high-order byte of an IPv6 multicast address : *ff*
- special IPv6 Multicast address
  - † ff02::1 is all-node group
  - † ff02::2 is the all-routers group

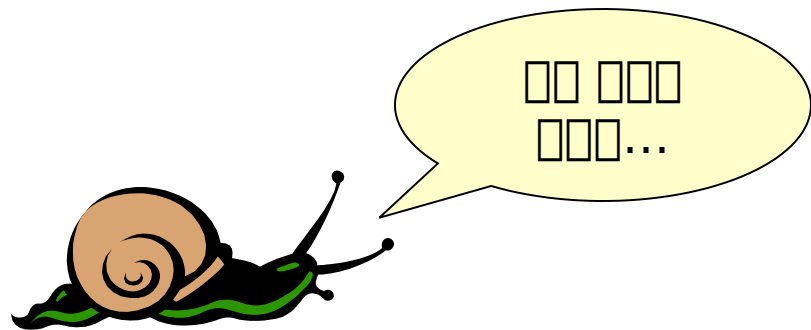


# Multicast Socket Options

Command	Datatype	Description
<b>IP_ADD_MEMBERSHIP</b>	<b>struct ip_mreq</b>	join a multicast group
<b>IP_DROP_MEMBERSHIP</b>	<b>struct ip_mreq</b>	leave a multicast group
<b>IP_MULTICAST_IF</b>	<b>struct in_addr</b>	specify default interface for outgoing multicasts
<b>IP_MULTICAST_TTL</b>	<b>u_char</b>	specify TTL for outgoing multicasts
<b>IP_MULTICAST_LOOP</b>	<b>u_char</b>	enable or disable loopback of outgoing multicasts
<b>IPV6_ADD_MEMBERSHIP</b>	<b>struct ipv6_mreq</b>	join a multicast group
<b>IPV6_DROP_MEMBERSHIP</b>	<b>struct ipv6_mreq</b>	leave a multicast group
<b>IPV6_MULTICAST_IF</b>	<b>u_int</b>	specify default interface for outgoing multicasts
<b>IPV6_MULTICAST_TTL</b>	<b>int</b>	specify TTL for outgoing multicasts
<b>IPV6_MULTICAST_LOOP</b>	<b>u_int</b>	enable or disable loopback of outgoing multicasts

# Receiving MBone Session Announcements

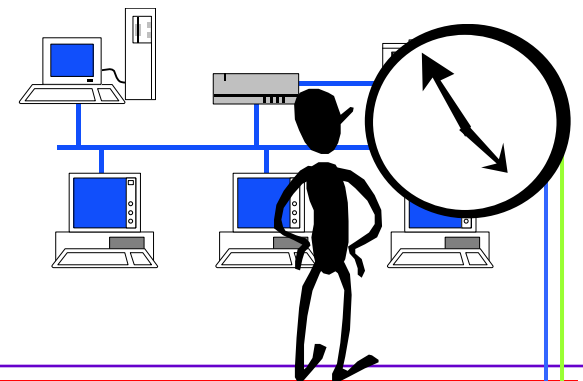
- **MBone**
  - multimedia conference
- **SAP**
  - Session Announcement Protocol
- **SDP**
  - Session Description Protocol



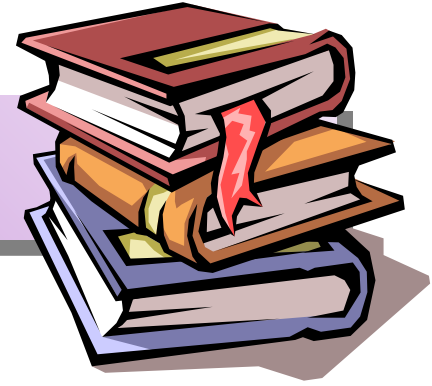


# SNTP (Simple Network Protocol)

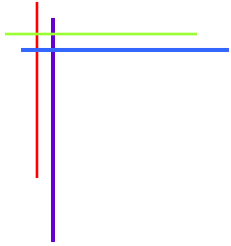
- **NTP(Network Time Protocol)**
  - sophisticated protocol for synchronizing clocks across a WAN or LAN, and
  - can often millisecond accuracy
- **SNTP**
  - common for a few hosts on a LAN to synchronize their clocks across the Internet to other NTP hosts, and
  - then redistribute this time on the LAN using either broadcasting or multicasting



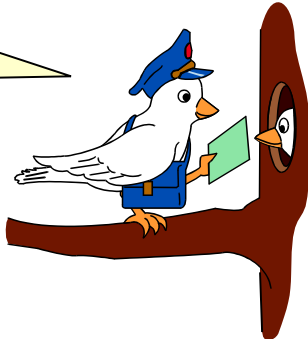
# Summary



- **Multicast application starts by joining the multicast group assigned to the application**
    - tell the IP layer to join the group
  - **Using hardware filtering reduces the load on all the other hosts that are not participating in the application**
  - **5 Five socket option**
    - join a multicast group on an interface
    - leave a multicast group,
    - set the default interface for outgoing multicasts,
    - set the TTL or hop limit for outgoing multicasts,
    - enable or disable loopback of multicasts
- } for receiving
- } for sending

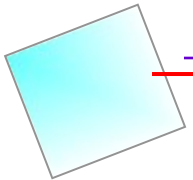


☐☐☐ reliable  
UDP ☐☐☐ ☐☐☐  
☐☐☐☐☐ ~



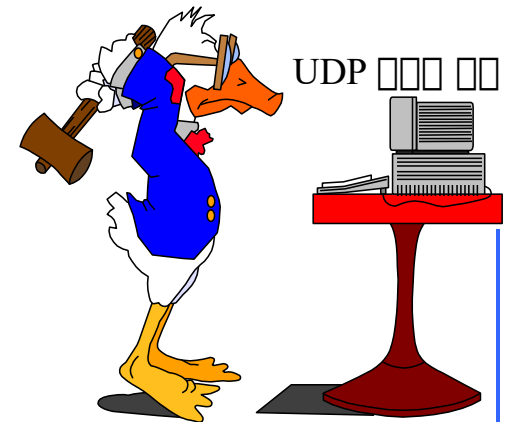
# Unix Network Programming

## Chapter 20 Advanced UDP Sockets

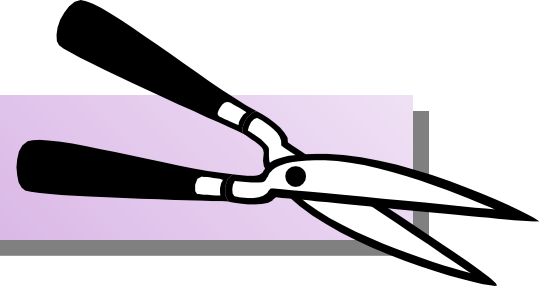


# Contents

- Introduction
- Receiving Flags, Destination IP address, and Interface Index
- **Datagram Truncation**
- **When to Use UDP Instead of TCP**
- *Adding Reliability to a UDP Application*
- Biding Interface Addresses
- **Concurrent UDP servers**
- **IPV6 Packet Information**
- **Summary**



# Datagram Truncation



- **When a UDP datagram arrives that is larger than the application's buffer, `recvmsg` sets the `MSG_TRUNC` flag.**
- **3** □□ □□□ □□□□
  - Discard the excess bytes and return the `MSG_TRUNC` flag to the application. This requires that the application call `recvmsg` to receive the flag
  - Discard the excess bytes but do not tell the application
  - Keep the excess bytes and return them in subsequent read operations on the socket

BSD/OS

solaris 2.5

early ver. of SVR4



# When to Use UDP Instead of TCP



- **Advantages of UDP**

- support broadcasting & multicasting
- has no connection setup or teardown
  - † UDP :  $RTT + SPT$  (T/TCP □ UDP □ □ □ )
  - † TCP :  $2 \times RTT + SPT$

- **only TCP Feature**

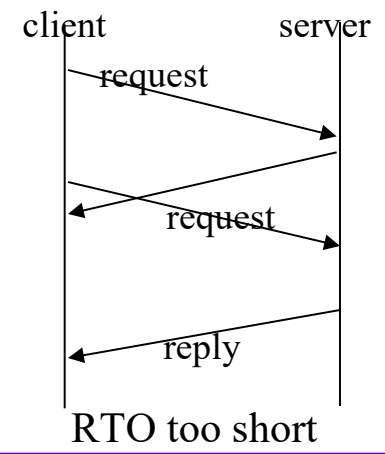
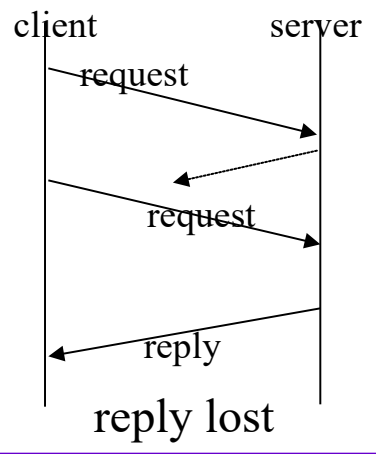
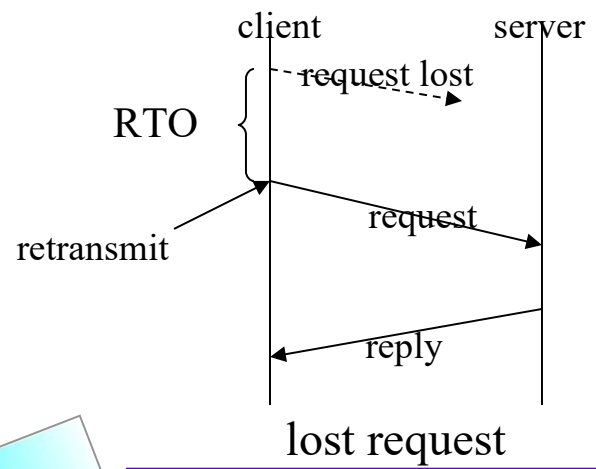
- Positive acknowledgments, retransmission of lost packets, duplication detection, and sequencing of packets reordered by the network
- Windowed flow control
- Slow start and congestion avoidance

# When to Use UDP

- **Recommendations**
  - UDP must be used for **broadcast** or **multicast** application
  - UDP can be used for ***simple request-reply applications*** but error detection must be built into the application
  - UDP should **not** be used for bulk data transfer

# Adding Reliability to a UDP Application

- **2 features to our client**
  - **Timeout / Retransmission** to handle datagrams
  - †
  - **Sequence numbers** so the client can verify reply
  - † ( ) DNS, SNMP, TFPT, RPC
- **Retransmission ambiguity Problem**



# Adding Reliability to a UDP Application

```

static sigjmp_buf jmpbuf;
{
  .....
  form request
  signal(SIGALRM, sig_alarm);
  rtt_newpack();
sendagain:
  sendto();
  alarm(rtt_start());
  if (sigsetjmp(jmpbuf, 1) !=0) {
    if (rtt_timeout() ) give up
    goto sendagain; }
  do {
    recvfrom();
  } while(wrong sequence#);
  alarm(0);
  rtt_stop();
  process_reply();
  .....
}
void sig_alarm(int signo) {
  siglongjmp(jmpbuf, 1);
}

```

establish signal handler

initialize rexmt counter to 0

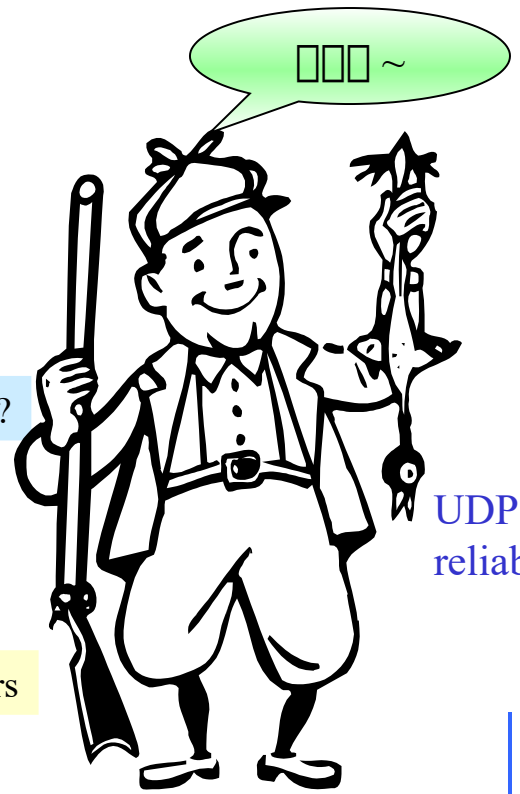
set alarm for RTO seconds

double RTO, retransmitted enough?

Retransmit

turn off alarm

calculate RTT and update estimators



UDP reliable

# Concurrent UDP Server(1)

Process involved in stand-alone concurrent UDP server

client

server (parent)  
port 69

creates socket, binds *recvfrom*, blocks until client request, *fork*, another *recvfrom*

server (child)  
port 2134

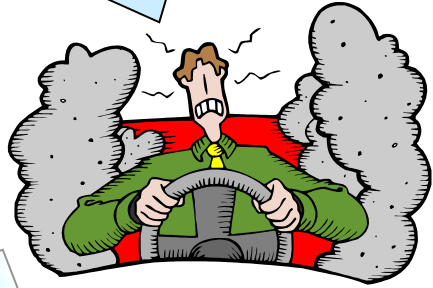
create new socket, bind ephemeral port(2134), process clients request, exchange additional datagrams with client

1<sup>st</sup> datagram from client

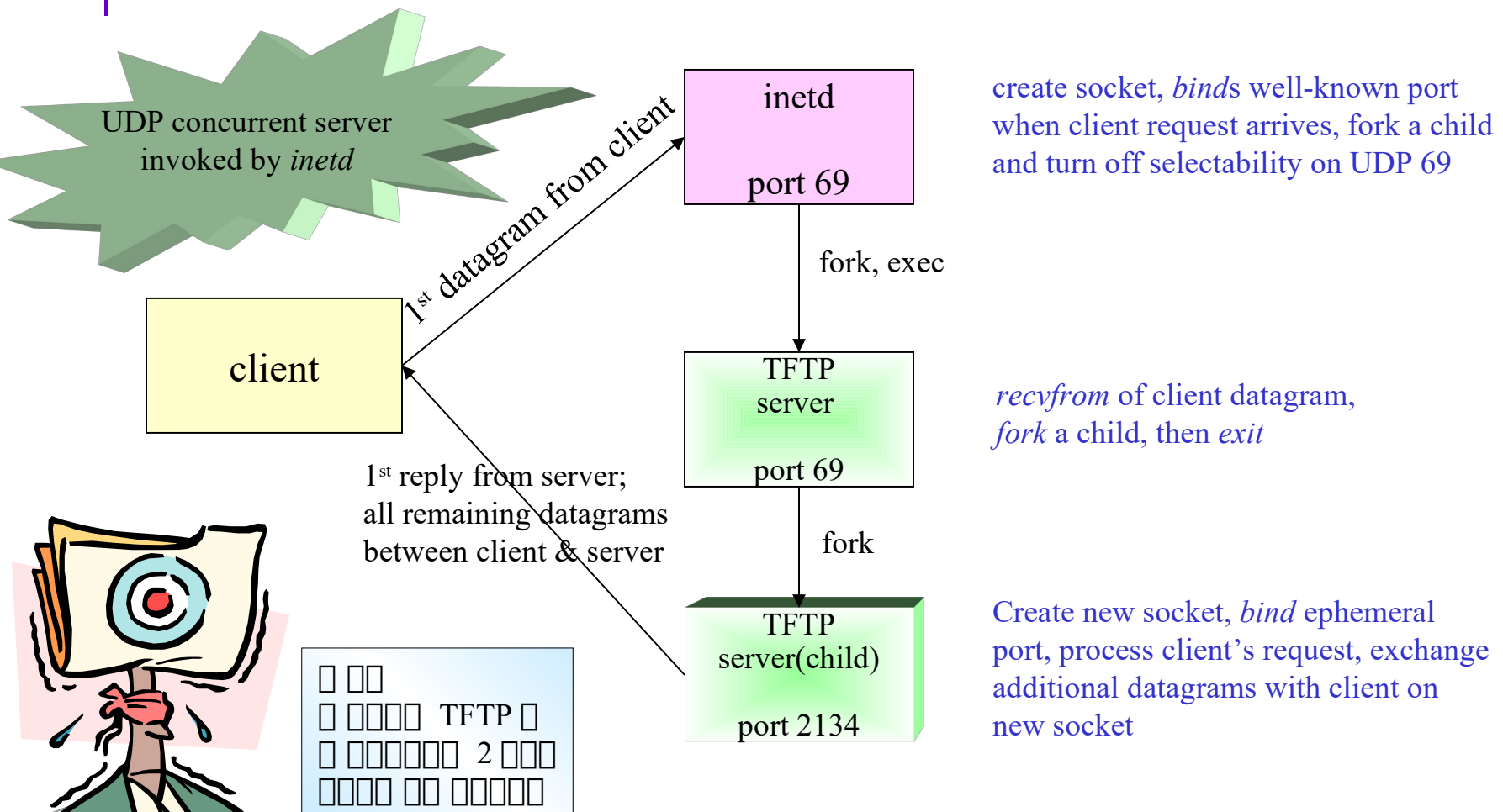
fork

1<sup>st</sup> reply from server; all remaining datagrams between client & server

request



# Concurrent UDP server(2)



# Summary

- **IP\_RECVSTADDR / IP\_RECVIF socket options can be enabled to return this information as ancillary data with each datagrams**
- **UDP**
  - broadcasting or multicasting
  - simple request-reply scenarios
  - Not used for bulk data transfer
- **Added reliability**
  - by detecting lost packets using a timeout and retransmission
  - RTT / timestamp

# Race Conditions - Signal case (4)

```

void dg_cli( ... )
{
  Setsockopt( sockfd, SOL_SOCKET, SO_BROADCAST, &on,
  sizeof(on) );

  Pipe( pipefd );

  Signal( SIGALRM, recvfrom_alarm );

  while ( ... ) {
    Sendto( sockfd, ... )

    alarm(5);
    for ( ; ; ) {
      FD_SET( sockfd, &rset );
      FD_SET( pipefd[0], &rset );
      if ( ( n=select( ... ) < 0 ) {
        if ( errno == EINTR )
          continue;
        else
          err_sys( "select error" );
      }
      if ( FD_ISSET( socfd, &rset ) ) {
        n = Recvfrom( sockfd, .. );
      }
    }
  }
}

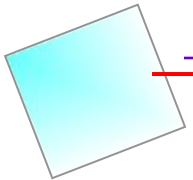
```

..... Continue

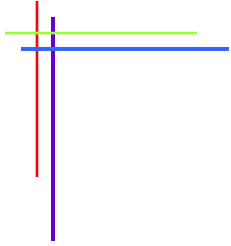
```

Continue .....
if ( FD_ISSET( pipefd[0], &rset ) ) {
  Read( pipefd[0], &n, 1 ); /* timer expired */
  break;
}
}
}
}
static void recvfrom_alarm( int signo )
{
  Write( pipefd[1], "", 1 ); return;
}

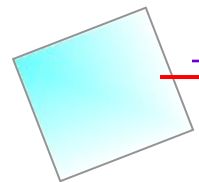
```





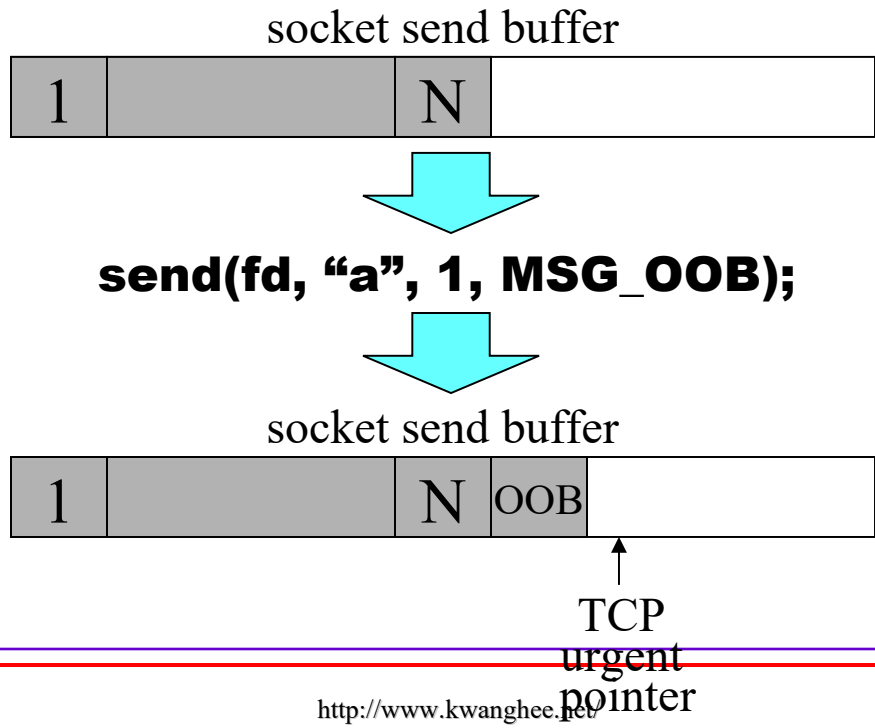


# **Chap. 21 *Out-of-Band Data***

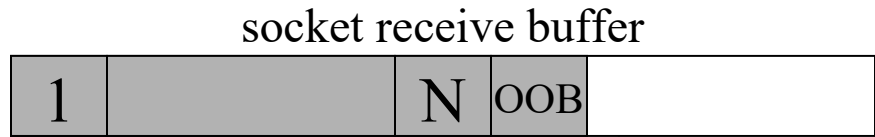


# Out-of-Band Data ... ? (1)

- **Out-of-band data** ≠ **normal(“inband”) data** priority
- **TCP** urgent mode ...



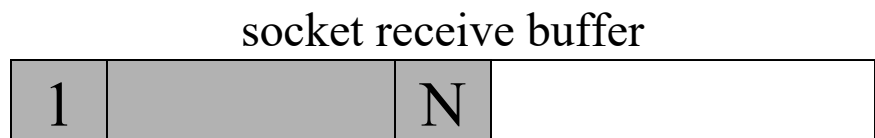
# Out-of-Band Data [...] ? (2)



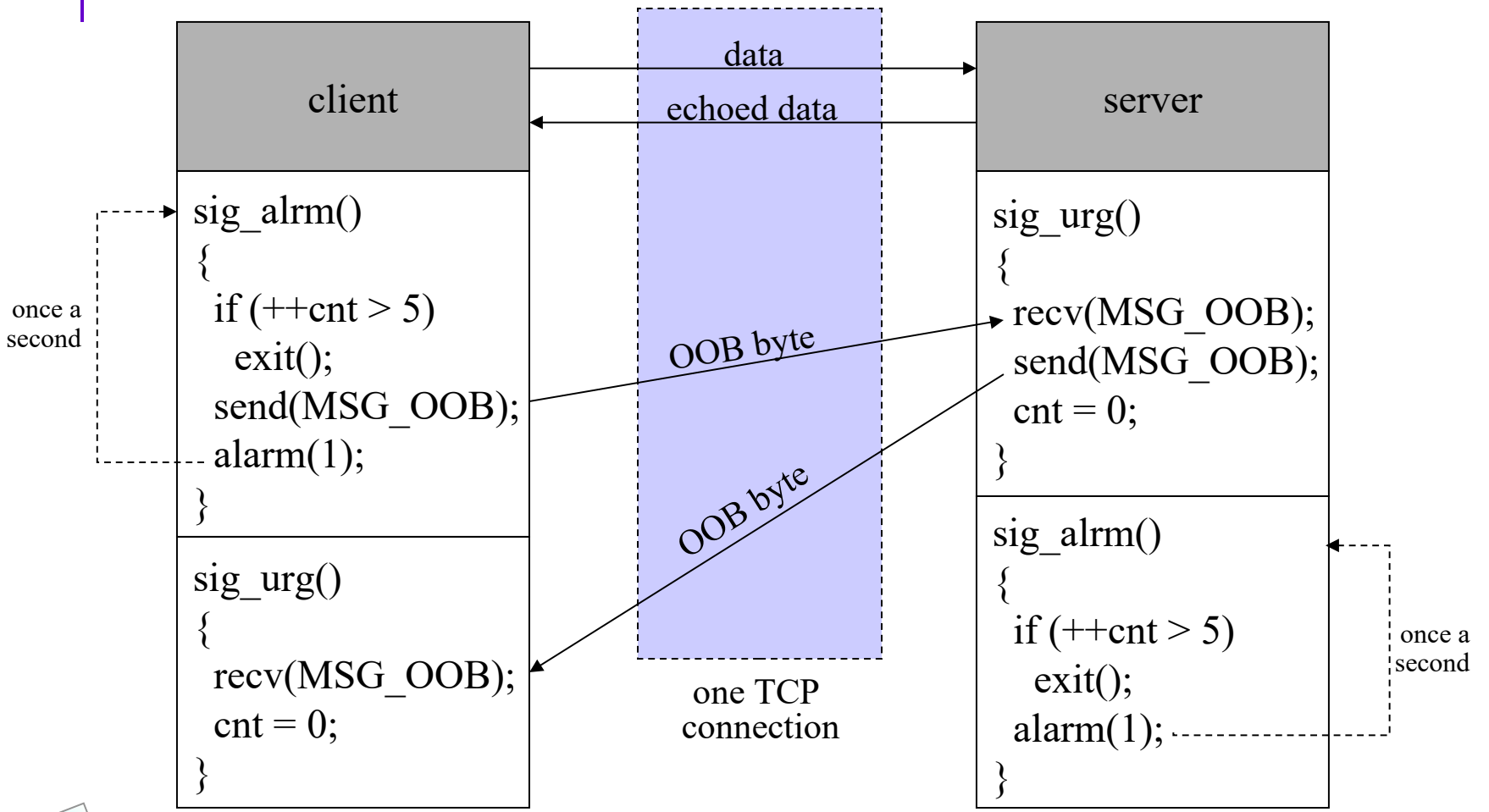
TCP  
urgent  
pointer

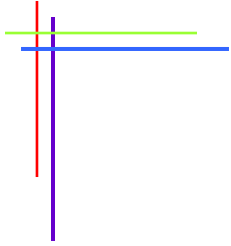
1. **SIGURG** []
2. **select** [] [] [] [] **except\_set ready**

```
recv(fd, ..., MSG_OOB);
```

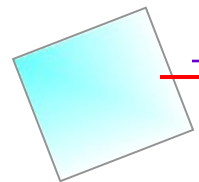


# Example: Client-Server Heartbeat Functions





# Chap. 22 *Signal-Driven I/O*

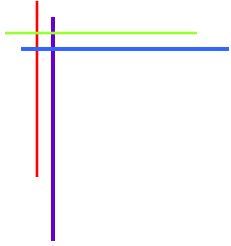


# Signal-Driven I/O

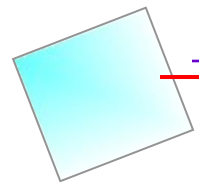
- **3 Steps to use signal-driven I/O with a socket (SIGIO)**
  1. A signal handler must be established for the SIGIO signal
    - ※ `signal(SIGIO, sigio_handler);`
  2. The socket owner must be set
    - ※ `fcntl(fd, F_SETOWN, getpid());`
  3. Signal-driven I/O must be enabled for the socket
    - ※ `fcntl(fd, O_ASYNC, on);`

# When to drive SIGIO

- **SIGIO with UDP Sockets**
  - a datagram arrives for the socket
  - an asynchronous error occurs on the socket
- **SIGIO with TCP Sockets**
  - a connection request has completed on a listening socket
  - a disconnect request has been initiated
  - a disconnect request has completed
  - half of a connection has been shut down
  - data has arrived on a socket
  - data has been sent from a socket (i.e., the output buffer has free space)
  - an asynchronous error occurred



# Chap. 23 *Threads*





# What is Threads?

- **Comparison to fork...**
  - fork is expensive (e.g., memory is copied, all descriptors are duplicated, and so on...)
  - IPC is required to pass info. between the parent and child
- **Threads**
  - all threads within a process share...
    - † process instructions
    - † most data
    - † open files (e.g., descriptors)
    - † signal handlers and signal dispositions
    - † current working dir
    - † user ID and group ID
  - each thread has its own...
    - † thread ID
    - † set of registers, including PC and SP
    - † stack
    - † errno
    - † signal mask
    - † priority

# Basic Thread Functions

- **int pthread\_create(pthread\_t \*tid, const pthread\_attr\_t \*attr, void \* (\*func) (void \*), void \*arg);**
- **int pthread\_join(pthread\_t tid, void \*\*status);**
- **pthread\_t pthread\_self(void);**
- **int pthread\_detach(pthread\_t tid);**
- **void pthread\_exit(void \*status)**

# TCP Echo Server Using Threads

```

#include "unpthread.h"

static void *doit(void *); /* each thread executes
                           this function */

int
main(int argc, char **argv)
{
    int      listenfd, *iptr;
    socklen_t  addrlen, len;
    struct sockaddr *cliaddr;

    if (argc == 2)
        listenfd = Tcp_listen(NULL, argv[1], &addrlen);
    else if (argc == 3)
        listenfd = Tcp_listen(argv[1], argv[2], &addrlen);
    else
        err_quit("usage: tcpserv01 [ <host> ] <service
or port>");

    cliaddr = Malloc(addrlen);

    for ( ;; ) {
        len = addrlen;
        iptr = Malloc(sizeof(int));
        *iptr = Accept(listenfd, cliaddr, &len);

        Pthread_create(NULL, NULL, &doit, iptr);
    }
}

```

```

static void *
doit(void *arg)
{
    int  connfd;

    connfd = *((int *) arg);
    free(arg);

    Pthread_detach(pthread_self());
    str_echo(connfd); /* same function as before */
    Close(connfd); /* we are done with connected socket */
    return(NULL);
}

```

# More Else about Threads

- **Thread-Specific Data**


- `int pthread_once(pthread_once_t *onceptr, void (*init)(void));`
- `int pthread_key_create(pthread_key_t *keyptr, void (*destructor)(void *value)`
- `void *pthread_getspecific(pthread_key_t key);`
- `int pthread_setspecific(pthread_key_t key, const void *value);`

- **Mutexes: Mutual Exclusion**

- `int pthread_mutex_lock(pthread_mutex_t *mptr);`
- `int pthread_mutex_unlock(pthread_mutex_t *mptr);`

- **Condition Variables**

- `int pthread_cond_wait(pthread_cond_t *cptr, pthread_mutex_t *mptr);`
- `int pthread_cond_signal(pthread_cond_t *cptr);`

 **Chap.24 - IP Options**

☞ **IP Options** □ □□

☞ **IPv4 Source Route Options**

☞ **IPv6 Extension Headers**

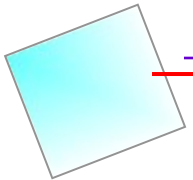
☞ **IPv6 Routing Header**

# We must study ... ?



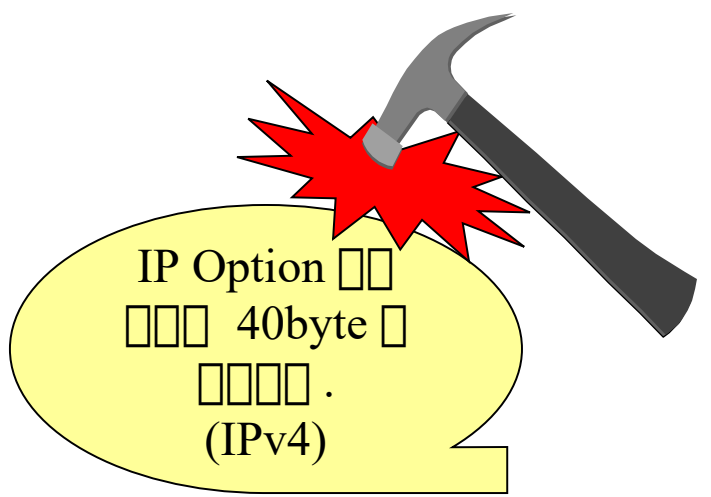
📖 IP option □ □□

📖 program □□ □□□ □□□□ □□□□ .



# IP Option ☐ ☐☐

- 1. NOP : no-operation, padding 1byte
- 2. EOL : end-of-list, indicate option end
- 3. LSRR : loose source and record route
- 4. SSRR : strict source and record route
- 5. Time stamp
- 6. Record route
- 7. Basic security
- 8. Extended security
- 9. Stream identifier
- 10. Router alert

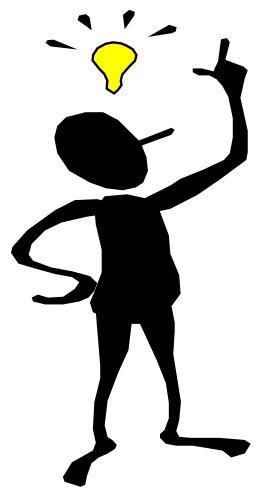


```
getsockopt ☐ setsockopt ☐☐☐☐
level = IPPROTO_IP, optname = IP_OPTIONS ☐
☐☐ ☐☐ ☐☐ ☐☐ ☐☐☐☐ .
```

# IPv4 source route option

☞ 0000 00 0000 IP 000000 00

☞ 0000 0 000 00 0 00 00 , 000000 0000 00  
00 0000 0000 0000 0 0 00 00 .



↪ SSRR

datagram 0 000 000 000 node 00 0000  
00 , 00 0000 node 00 0000 0 00 .

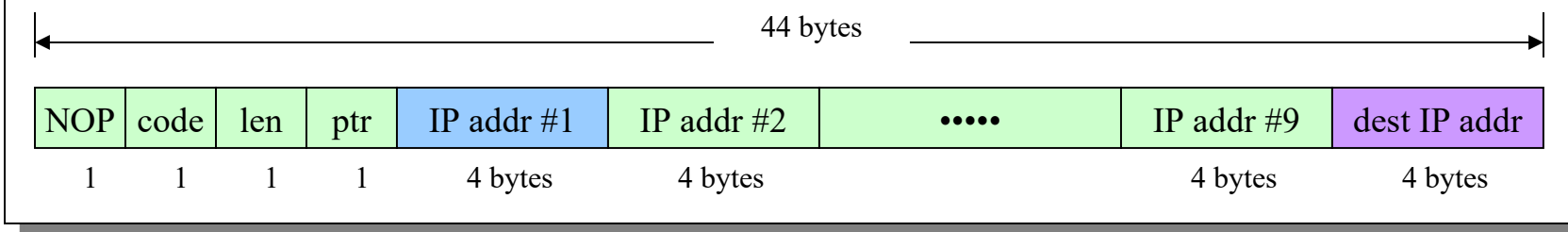
↪ LSRR

datagram 0 000 000 node 00 0000 0000  
0000 00 00 node 00 0000 00 00 .

※ Kernel 0 000 000 source route 000 000 IP address 0 0000 application 00  
0000 .

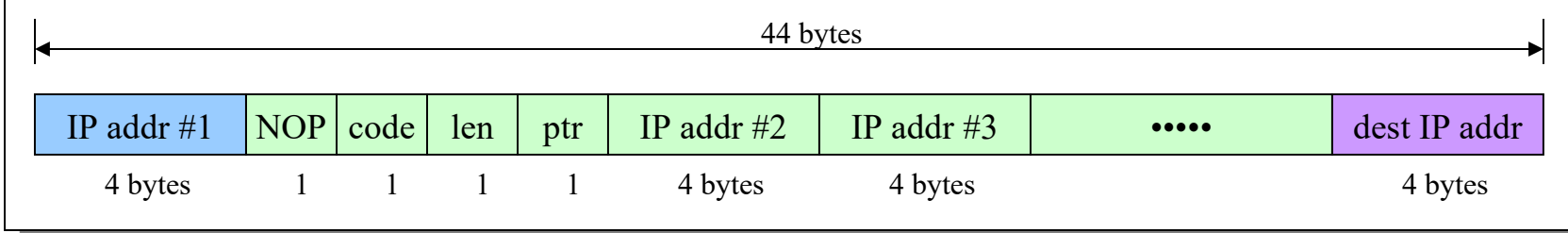


Send : to kernel by setsockopt



code : LSRR - 0x83, SSRR - 0x89  
 len : code ~ dest IP addr □□□ byte □  
 ptr : □□□□ □□□ □□ □ IP □□□ offset

Receive : to application by getsockopt



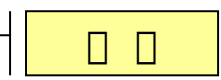
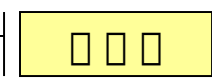
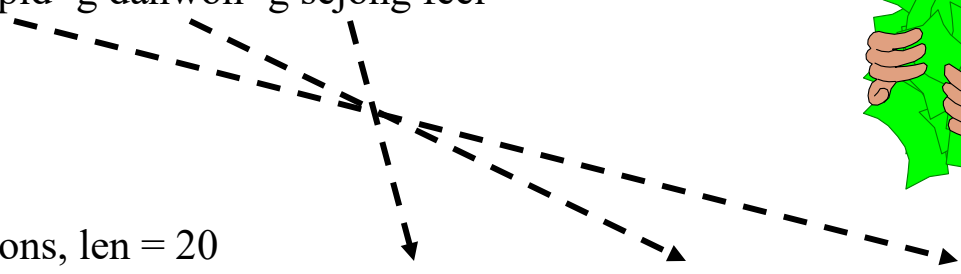
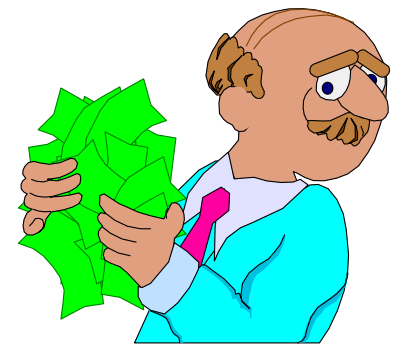
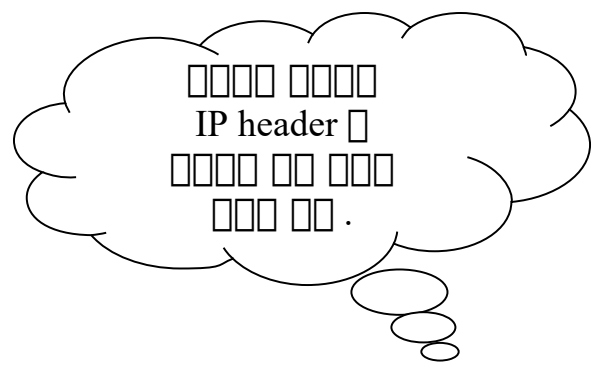
# Example

```
> tcpcli01 -g sejong -g cupid feel
hi
hi
```

```
# tcpserv01
received IP options, len = 16
received LSRR: 150.150.55.52 150.150.55.53 150.150.55.55
```

```
> tcpcli01 -g cupid -g danwon -g sejong feel
hi
hi
```

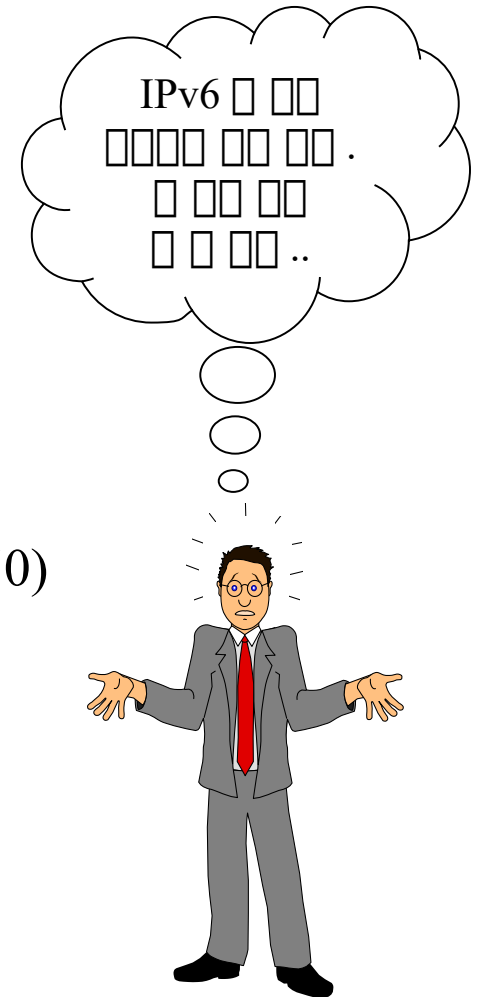
```
# tcpserv01
received IP options, len = 20
received LSRR: 150.150.55.52 150.150.55.55 150.150.55.54 150.150.55.53
child 21180 terminated
```



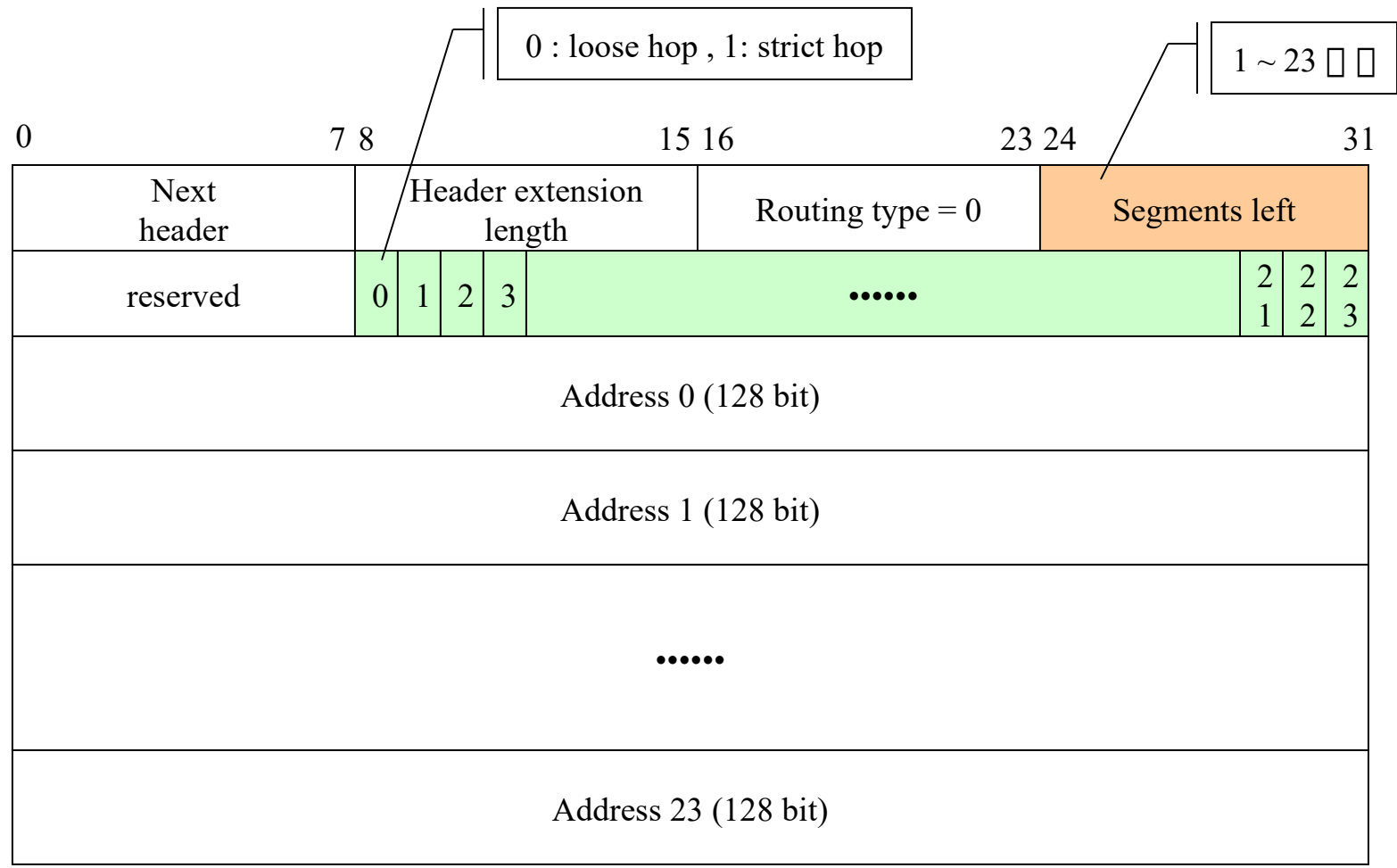
# IPv6 extension headers

1. Hop-by-hop option(0) : currently non-defined
2. Destination option(60) : currently non-defined
3. Routing header(43) : source routing option
4. Fragmentation header(44)
5. Authentication header(51)
6. ESP(Encapsulating security payload) header(50)

※ ( ) : IP next header field value. if value is 59, no next header



# IPv6 routing header



# Contents

## Chap.25 - Raw Socket

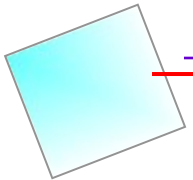
- ☞ **Raw socket's features**
- ☞ **Raw socket creation, output, input**
- ☞ **Ping program**
- ☞ **Traceroute program**
- ☞ **ICMP message daemon**

# We must study ... ?



📖 Raw socket □ □□

📖 Raw socket □ □□□ programming



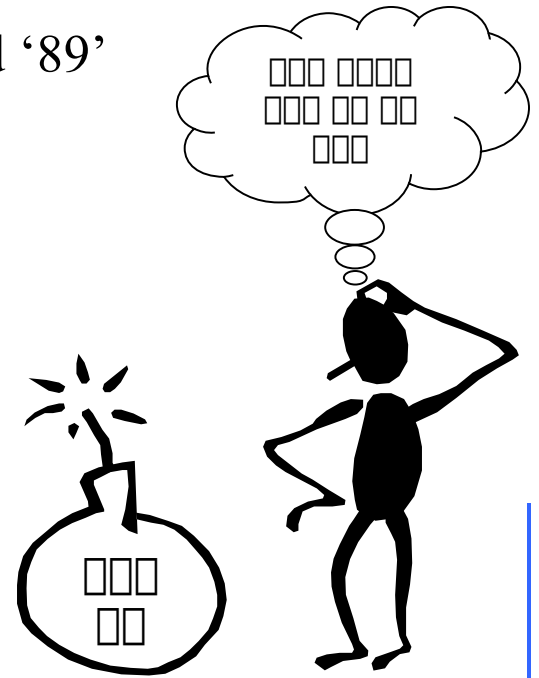
# Raw socket's features

1. ICMPv4, IGMPv4, ICMPv6 packet [ ] [ ] [ ] [ ] [ ] [ ] .

Ex) ping program

2. Kernel [ ] [ ] [ ] [ ] [ ] IPv4 protocol field [ ] [ ] IPv4 datagram [ ] [ ] [ ] [ ] [ ] .

Ex) OSPF routing protocol - protocol field '89'  
cf) ICMP-1, IGMP-4, TCP-6, UDP-17



# Raw socket creation, output, input

## Creation

How to create raw socket in C.

```
sockfd = socket( AF_INET, SOCK_RAW, protocol )
```

### ▶ IP\_HDRINCL socket option

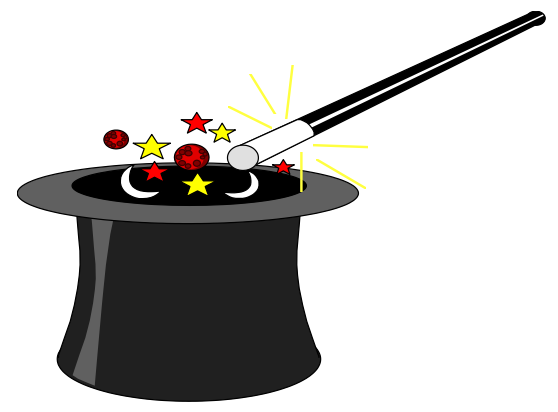
```
const int on = 1;
if ( setsockopt( sockfd, IPPROTO_IP, IP_HDRINCL, &on, sizeof(on) ) < 0 )
    error
```

### ▶ Bind

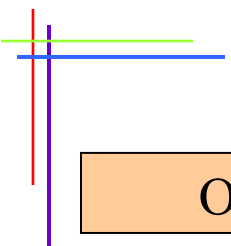
local address

### ▶ Connect

remote address







# Output

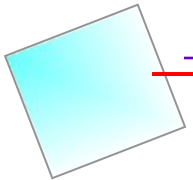
- ▶ sendto() sendmsg() destination IP address

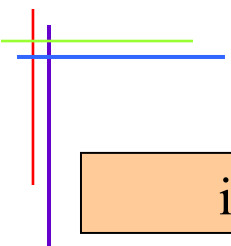
```
Ex] sendto( sockfd, sendbuf, len, 0, pr->sasend, pr->salen );
```

- ▶ IP\_HDRINCL kernel IP header process data IP header byte

- ▶ IP\_HDRINC kernel IP header byte

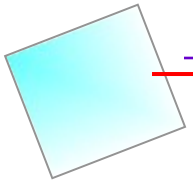
- ▶ kernel interface MTU packet





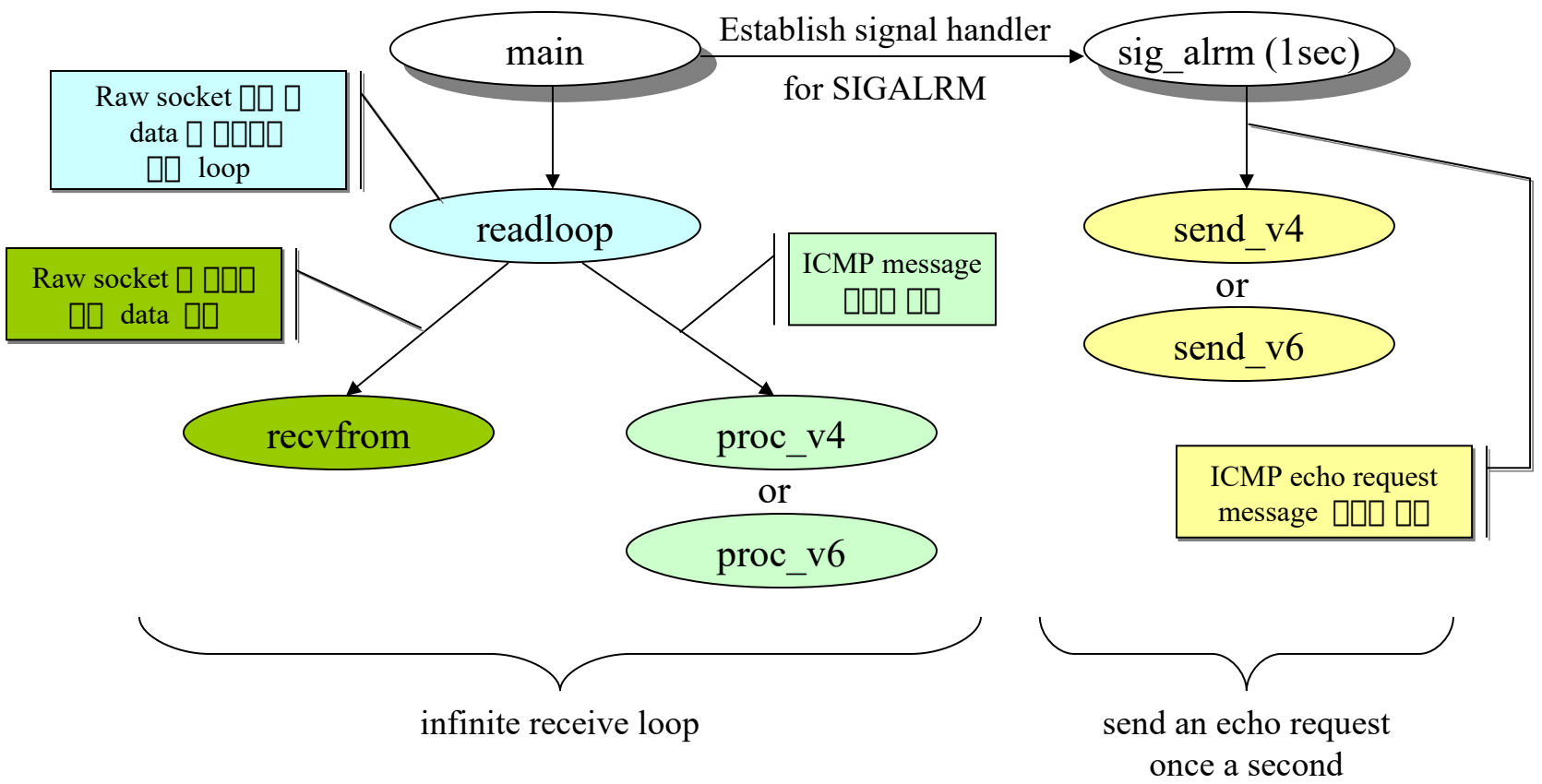
input

- ▶ UDP packets □ TCP packets □ raw socket □ □□□□ □□□□ . Process □ □□□□ data link □□□□ □□□□ □□ .
- ▶ kernel □ ICMP message □ □□□ □□ □□□□ **ICMP packet** □□ raw socket □□ □□□□ . Echo, timestamp, □□ □□ □□□□ □□□□ kernel □ □□□ □□□□ .
- ▶ kernel □ IGMP message □ □□□ □□ □□ **IGMP packet** □□ raw socket □□ □□□□ .
- ▶ kernel □ □□□□ □□□□ protocol field □ □□ □□ IP datagram □ raw socket □□ □□□□ .
- ▶ datagram □ fragment □□□□ □□□□ □□ □□□□ □□□□ □□□□ □□□□ □□□□ raw socket □□ □□□□ □□□□ .



# Ping program

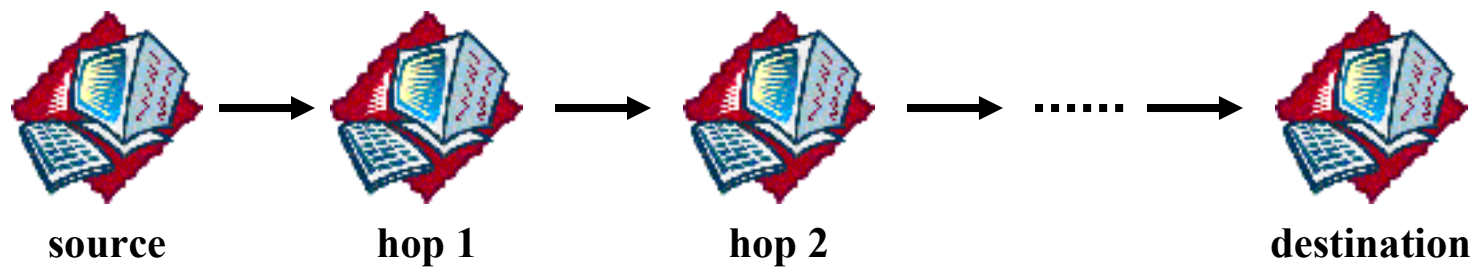
## Function diagram



# Traceroute program

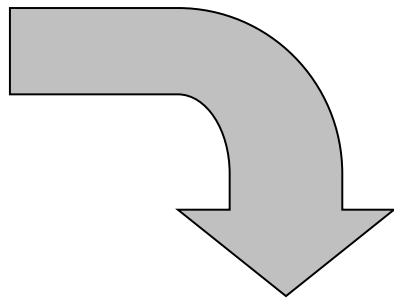
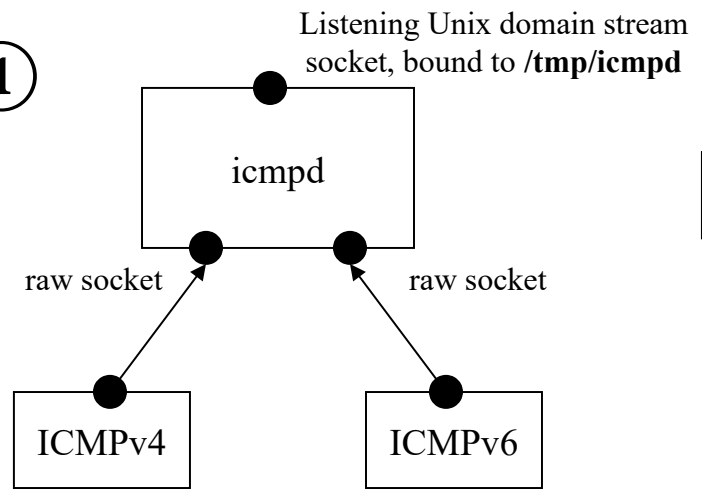
☞ Host → destination host → datagram flow → trace → .

- ① TTL( hop count) = 1 → → UDP datagram → destination → .
- ② First hop router → ICMP 'time exceeded in transit' →
- ③ TTL = → → UDP datagram → destination → .
- ④ ③ → → hop host → ICMP 'time exceeded in transit' → .
- ⑤ Destination host → ICMP 'port unreachable' → .

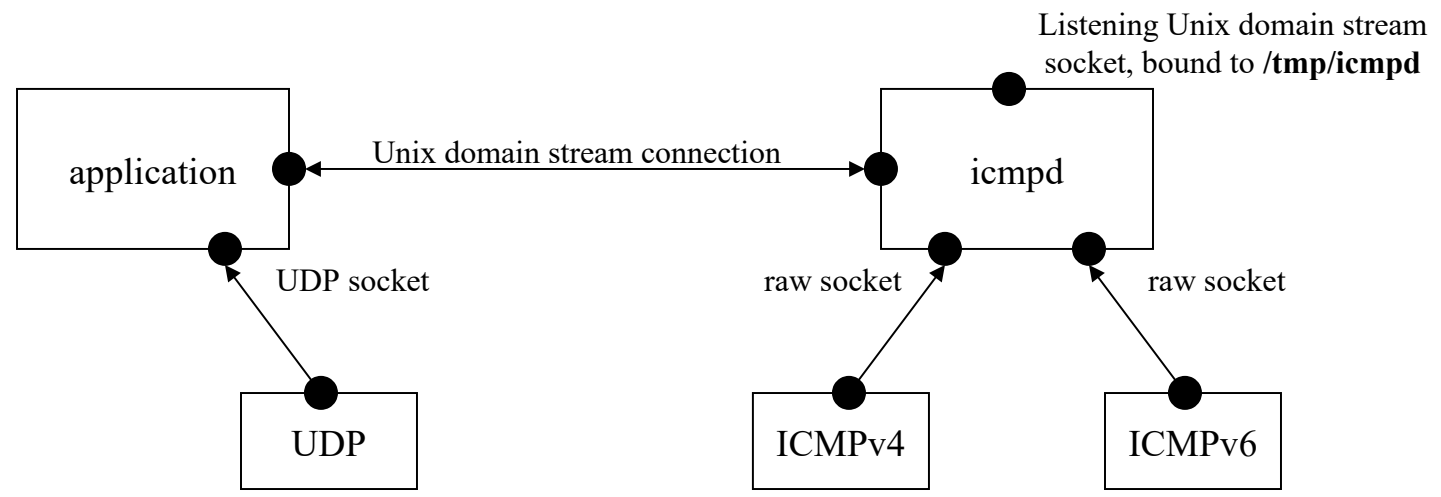


# ICMP Message daemon

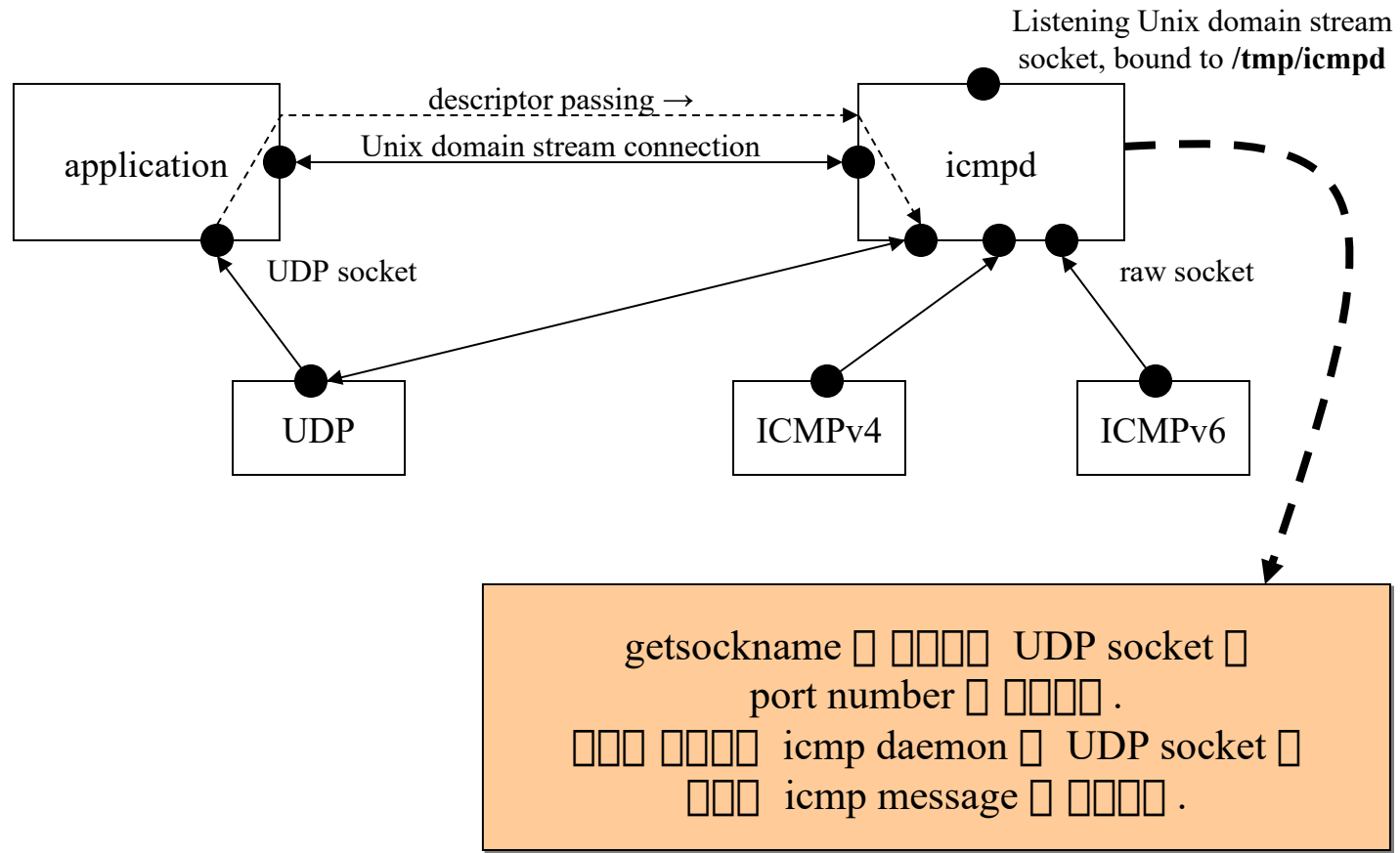
①



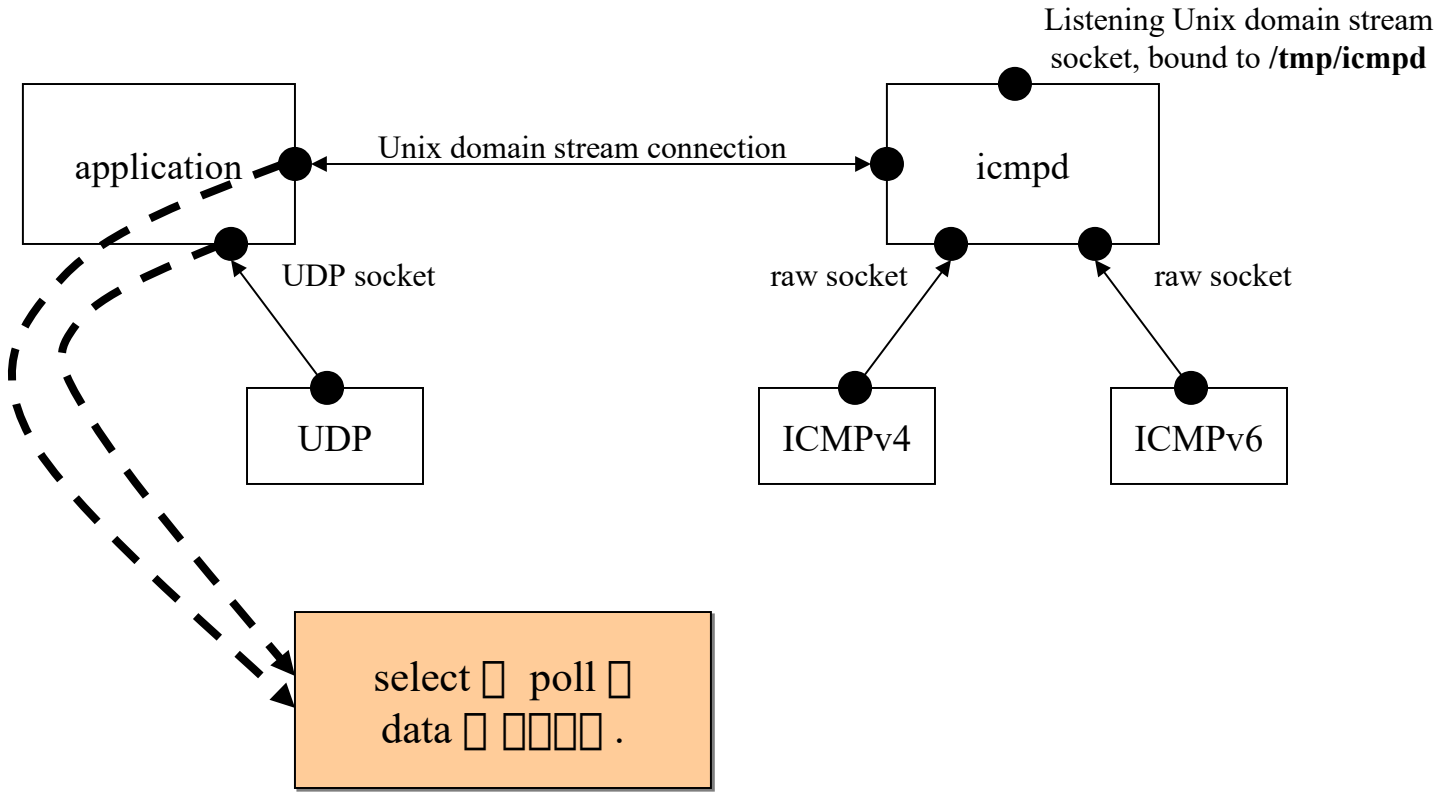
②

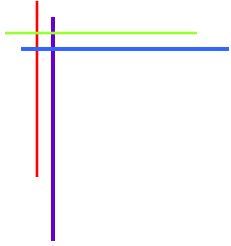


3

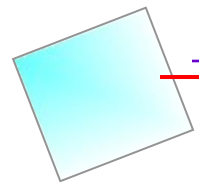


4





# Chap. 26 *Datalink Access*

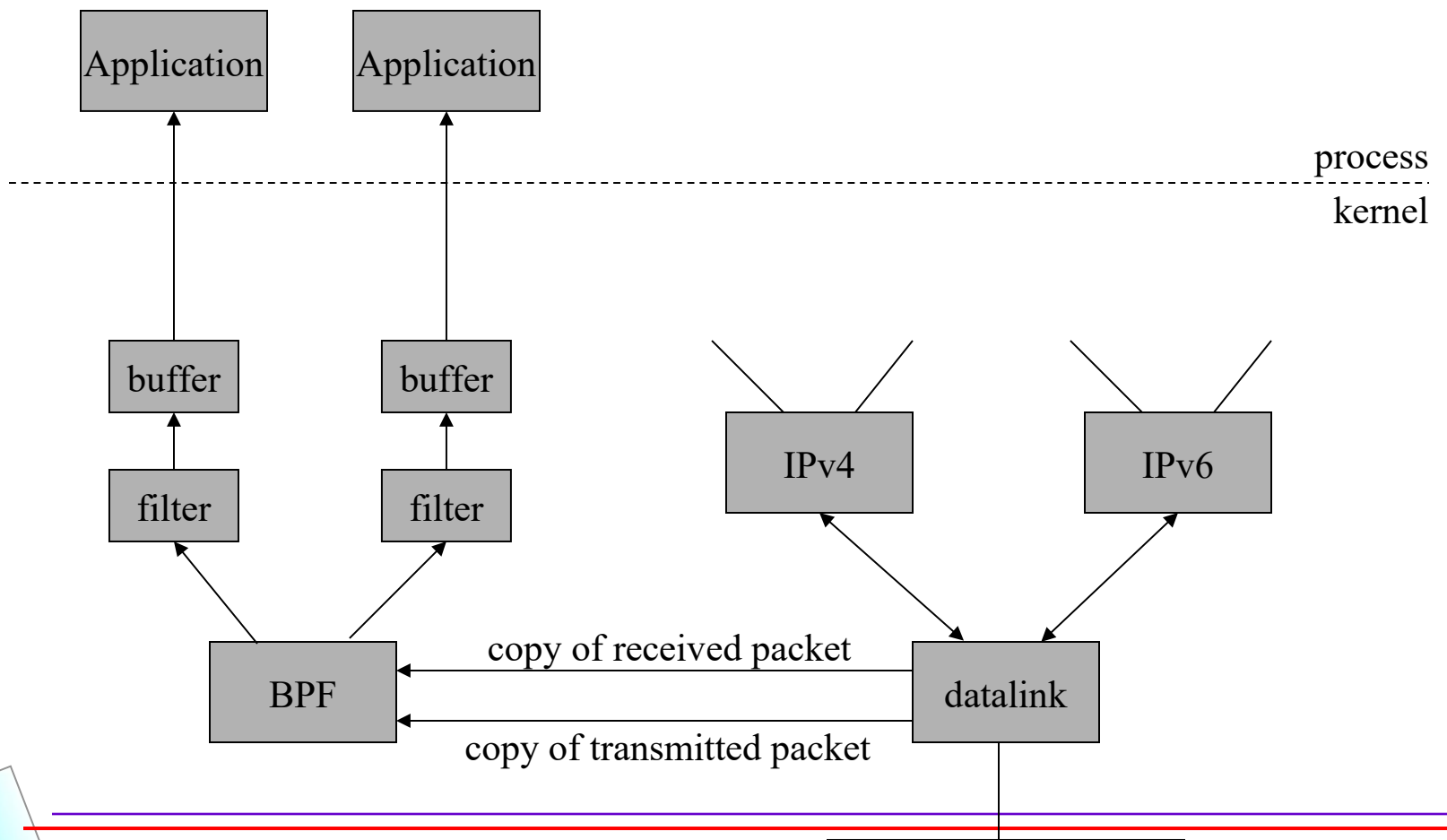




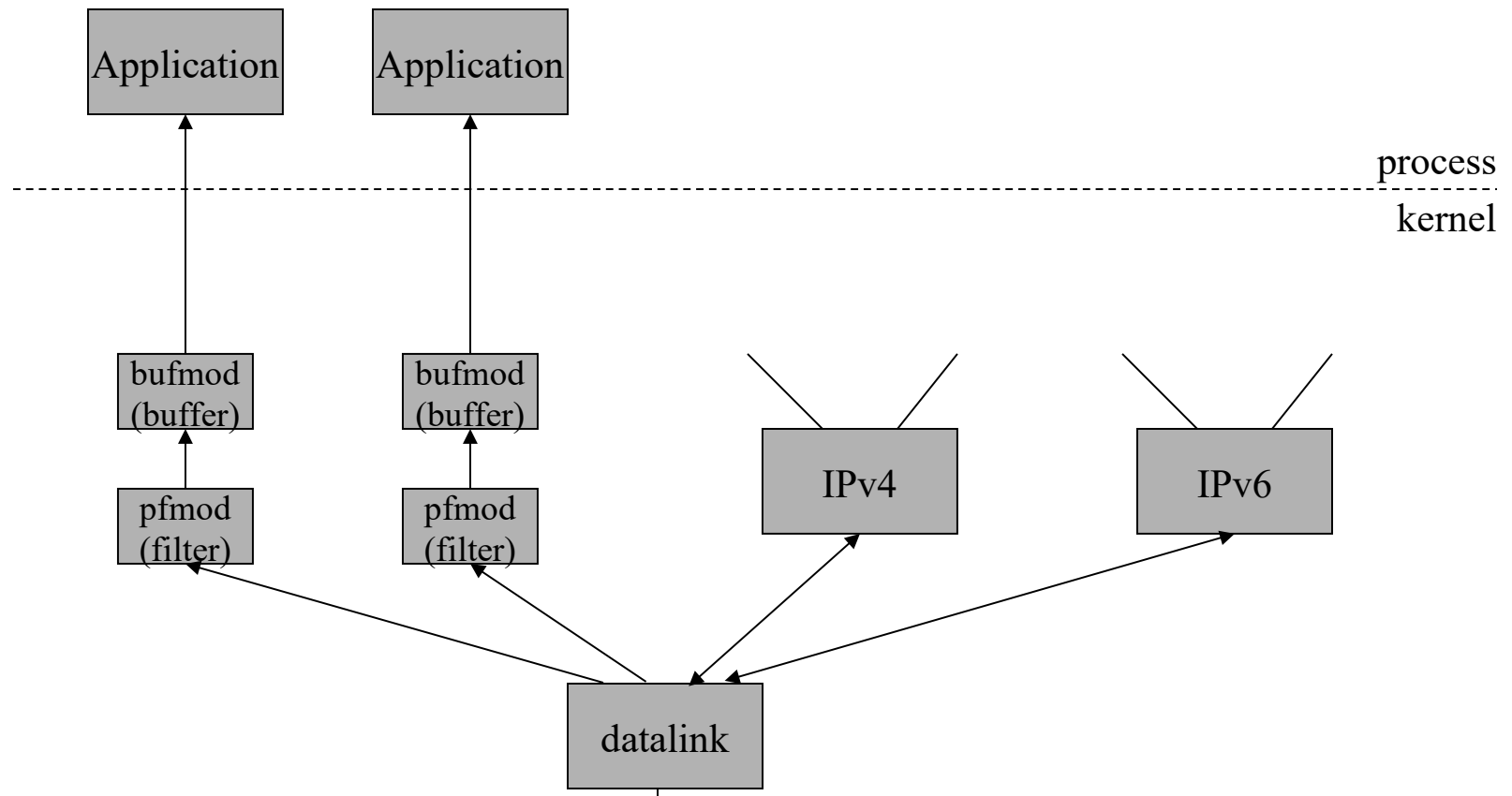
# Introduction

- **Datalink access provides...**
  - The ability to watch the packets received by the datalink layer
  - The ability to run certain programs as normal applications instead of as part of the kernel(e.g., RARP server)
- **Common methods to access the datalink layer**
  - BSD Packet Filter(BPF)
  - SVR4 Data Link Provider Interface (DLPI)
  - Linux SOCK\_PACKET interface
  - **libpcap, the publicly available packet capture library**

# Packet Capture Using BPF

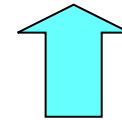
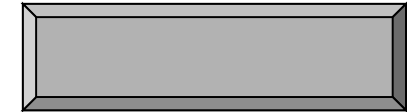


# Packet Capture Using DLPI



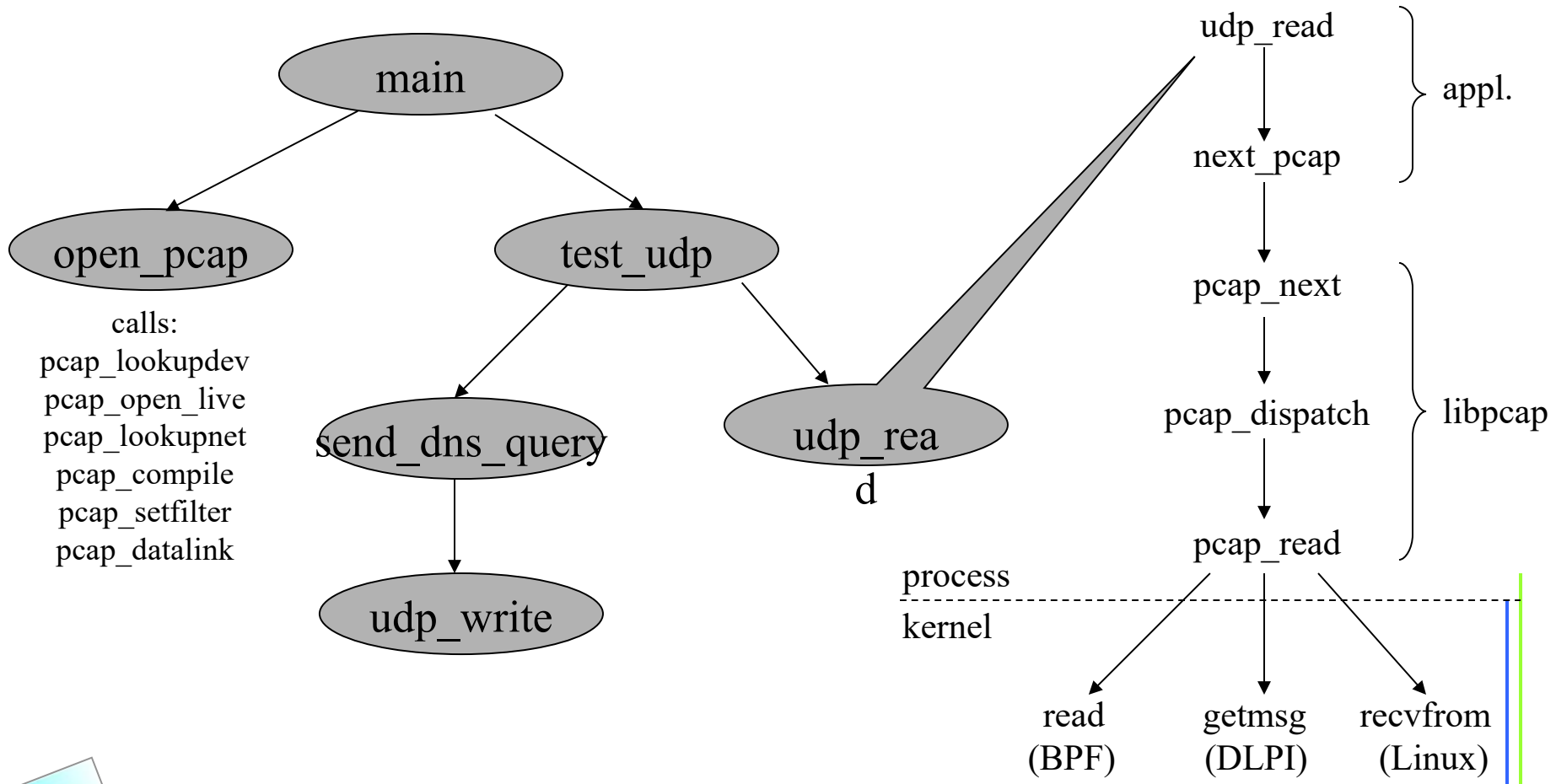
# Linux: SOCK\_PACKET

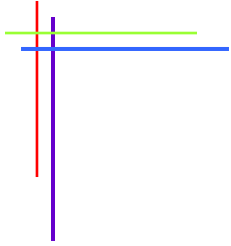
- **fd = socket(AF\_INET, SOCK\_PACKET, htons( ETH\_P\_ALL ));**



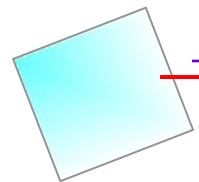
ETH\_P\_ALL  
ETH\_P\_IP  
ETH\_P\_ARP  
ETH\_P\_IPV6

# libpcap Example: Examining the UDP Checksum Field





# **Chap. 27 *Client-server Design Alternatives***

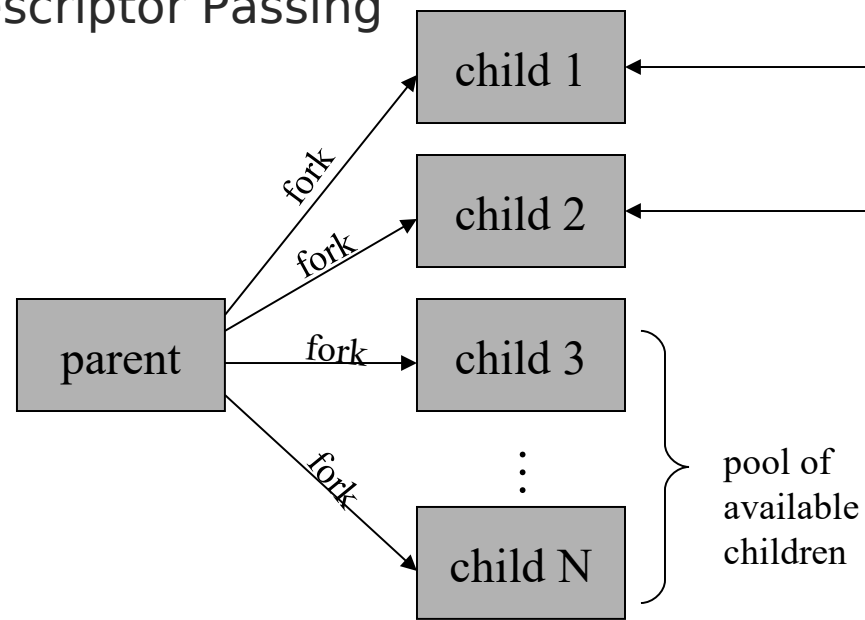




# Alternative TCP Server Models

## • Preforking server

- No Locking around accept (BSD only)
- File Locking around accept
- Thread Locking around accept
- Descriptor Passing



## ■ Prethreading server

- per-Thread accept
- Main Thread accept

