

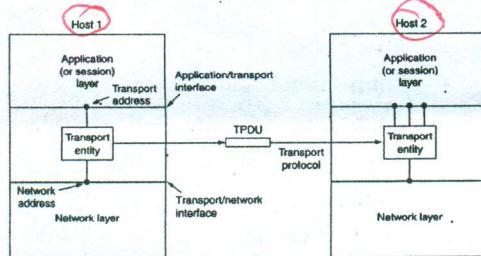
Chapter 6

The Transport Layer

✓ The Transport Service

- Services Provided to the Upper Layers
- Transport Service Primitives
- Berkeley Sockets
- An Example of Socket Programming:
 - An Internet File Server

Services Provided to the Upper Layers



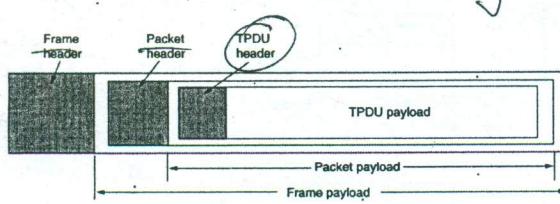
The network, transport, and application layers.

✓ Transport Service Primitives

Primitive	Packet sent	Meaning
LISTEN	(none)	Block until some process tries to connect
CONNECT	CONNECTION REQ.	Actively attempt to establish a connection
SEND	DATA	Send information
RECEIVE	(none)	Block until a DATA packet arrives
DISCONNECT	DISCONNECTION REQ.	This side wants to release the connection

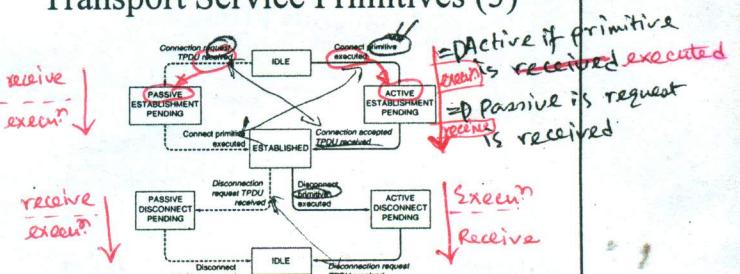
The primitives for a simple transport service.

Transport Service Primitives (2)



The nesting of TPDUs, packets, and frames.

Transport Service Primitives (3)



A state diagram for a simple connection management scheme. Transitions labeled in italics are caused by packet arrivals. The solid lines show the client's state sequence. The dashed lines show the server's state sequence.

Berkeley Sockets

Primitive	Meaning
SOCKET	Create a new communication end point
BIND	Attach a local address to a socket
LISTEN	Announce willingness to accept connections; give queue size
ACCEPT	Block the caller until a connection attempt arrives
CONNECT	Actively attempt to establish a connection
SEND	Send some data over the connection
RECEIVE	Receive some data from the connection
CLOSE	Release the connection

The socket primitives for TCP.

Socket Programming Example: Internet File Server

Client code using sockets.

```

/* This page contains a client program that does request file from the server program
 * on the next page. The server responds by sending the whole file.
 */
#include <sys/types.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <sys/conf.h>
#include <sys/buf.h>
#define SERVER_PORT 12345
#define BUF_SIZE 4096
#define QUEUE_SIZE 10
int main(int argc, char *argv)
{
    int c, bytes;
    char buf[BUF_SIZE];
    struct hostent *h;
    struct sockaddr_in channel;
    int channellen, hlen;
    if (argc != 3) fatal("Usage: client server-name file-name");
    h = gethostbyname(argv[1]);
    if (!h) fatal("gethostbyname failed");
    s = socket(AF_INET, SOCK_STREAM, IPPROTO_TCP);
    if (s < 0) fatal("socket failed");
    /* Build address structure to bind to socket. */
    memzero(&channel, sizeof(channel));
    channel.sin_family = AF_INET;
    channel.sin_addr.s_addr = htonl(INADDR_ANY);
    channel.sin_port = htons(SERVER_PORT);
    /* Create open. Wait for connection. */
    if (connect(s, (struct sockaddr *)&channel, sizeof(channel)) < 0)
        fatal("connect failed");
    /* Bind, (which socketd *) channel, sizeof(channel));
    if (bind(s, (struct sockaddr *)&channel, sizeof(channel)) < 0)
        fatal("bind failed");
    /* Listen, QUEUE_SIZE */
    if (listen(s, 10) < 0) fatal("listen failed");
    /* Socket is now set up and bound. Wait for connection and process it. */
    while (1) {
        if (select(0, 0, 0, 0) < 0) /* block for connection request */
            fatal("select failed");
        if (bytes = read(s, buf, BUF_SIZE) < 0)
            fatal("read from socket failed");
        /* read file name from socket */
        if (bytes == 0) break; /* check for end of file */
        write(1, buf, bytes); /* write bytes to socket */
        close(s);
        close(s);
        /* close file */
    }
    /* Go get the file and write it to standard output. */
    while (1) {
        bytes = read(s, buf, BUF_SIZE);
        if (bytes == 0) exit(0); /* check for end of file */
        write(1, buf, bytes);
        /* write to standard output */
    }
    fatal("client");
}

```

Socket Programming Example: Internet File Server (2)

Client code using sockets.

```

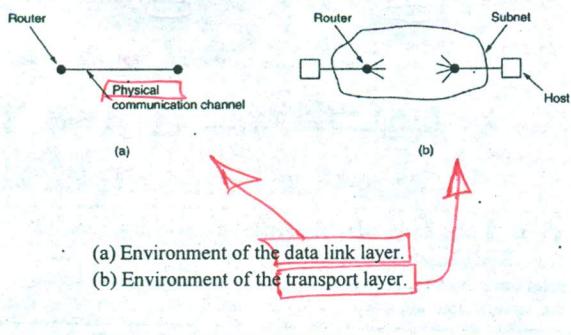
/* This is the server code */
#include <sys/types.h>
#include <sys/conf.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <sys/conf.h>
#define SERVER_PORT 12345
#define BUF_SIZE 4096
#define QUEUE_SIZE 10
int main(int argc, char *argv)
{
    int c, bytes, on = 1;
    char buf[BUF_SIZE];
    struct sockaddr_in channel;
    /* Build address structure to bind to socket. */
    memzero(&channel, sizeof(channel));
    channel.sin_family = AF_INET;
    channel.sin_addr.s_addr = htonl(INADDR_ANY);
    channel.sin_port = htons(SERVER_PORT);
    /* Create open. Wait for connection. */
    s = socket(AF_INET, SOCK_STREAM, IPPROTO_TCP); /* create socket */
    if (s < 0) fatal("socket failed");
    if (bind(s, (struct sockaddr *)&channel, sizeof(channel)) < 0)
        fatal("bind failed");
    if (listen(s, 10) < 0) fatal("listen failed");
    /* Bind, (which socketd *) channel, sizeof(channel));
    if (listen(s, 10) < 0) fatal("listen failed");
    /* Socket is now set up and bound. Wait for connection and process it. */
    while (1) {
        if (select(0, 0, 0, 0) < 0) /* block for connection request */
            fatal("select failed");
        if (bytes = read(s, buf, BUF_SIZE) < 0)
            fatal("read from socket failed");
        /* read file name from socket */
        if (bytes == 0) break; /* check for end of file */
        write(1, buf, bytes); /* write bytes to socket */
        close(s);
        close(s);
        /* close file */
    }
    /* Go get the file and write it to standard output. */
    while (1) {
        bytes = read(s, buf, BUF_SIZE); /* read from file */
        if (bytes == 0) break; /* check for end of file */
        write(1, buf, bytes); /* write bytes to socket */
        close(s);
        close(s);
        /* close file */
    }
    /* Go get the file and write it to standard output. */
    while (1) {
        bytes = read(s, buf, BUF_SIZE);
        if (bytes == 0) exit(0); /* check for end of file */
        write(1, buf, bytes);
        /* write to standard output */
    }
    fatal("server");
}

```

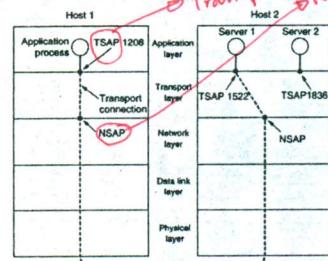
Elements of Transport Protocols

- Addressing
- Connection Establishment
- Connection Release
- Flow Control and Buffering
- Multiplexing
- Crash Recovery

Transport Protocol

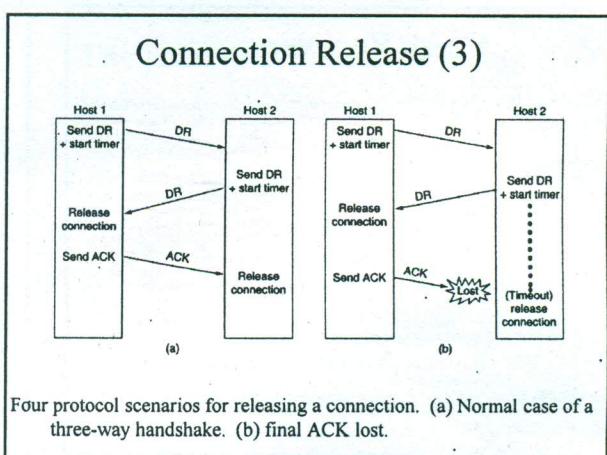
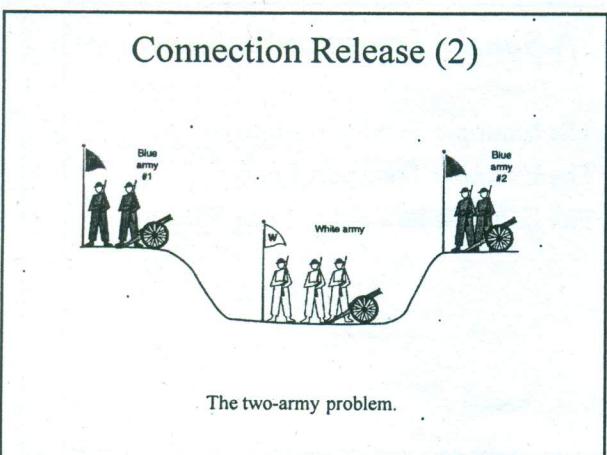
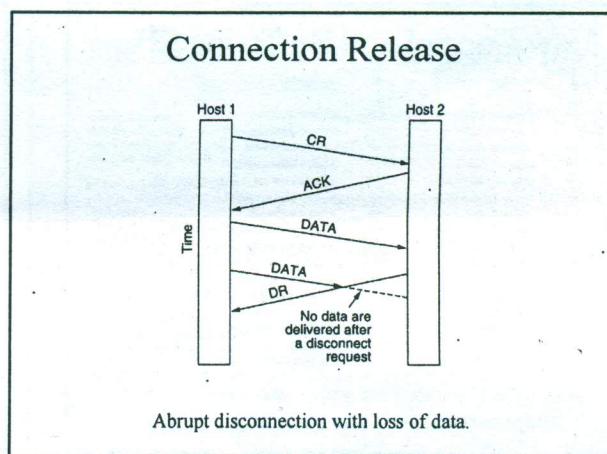
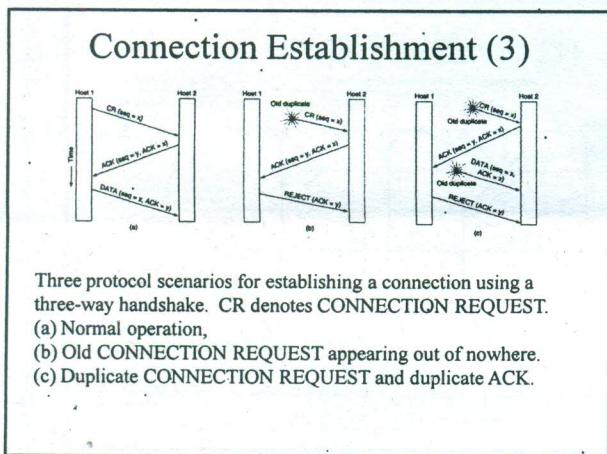
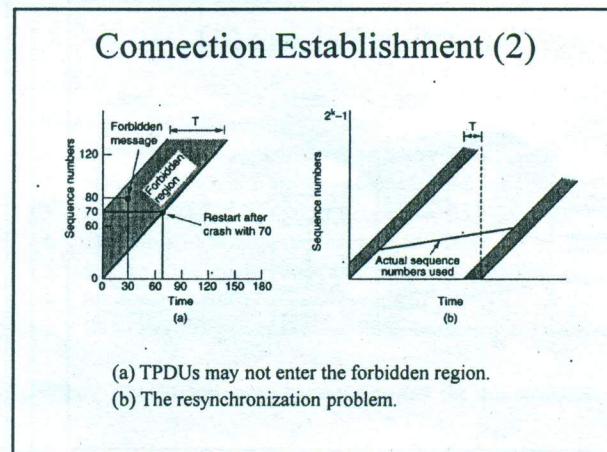
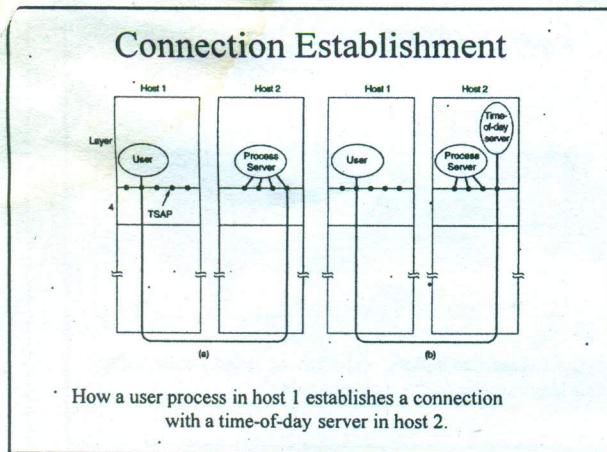


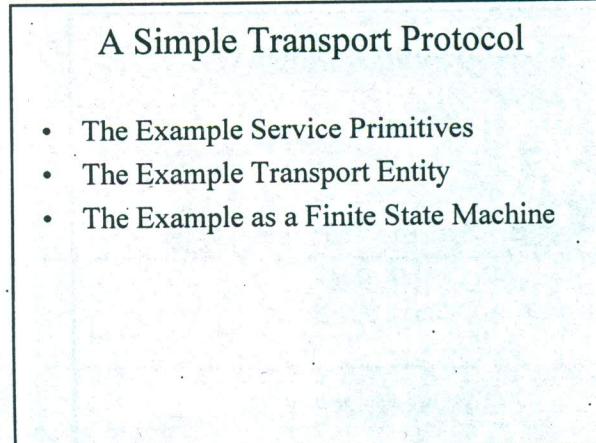
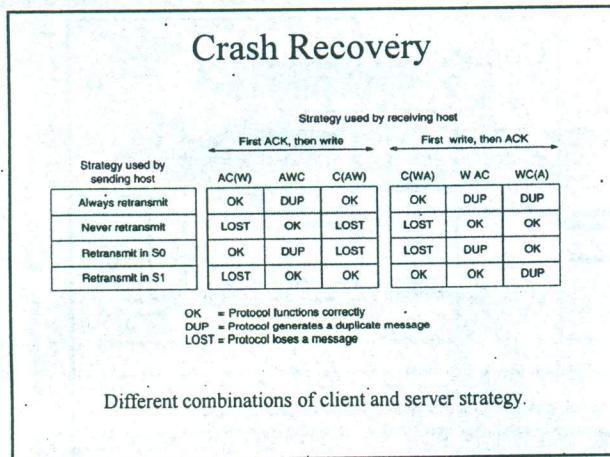
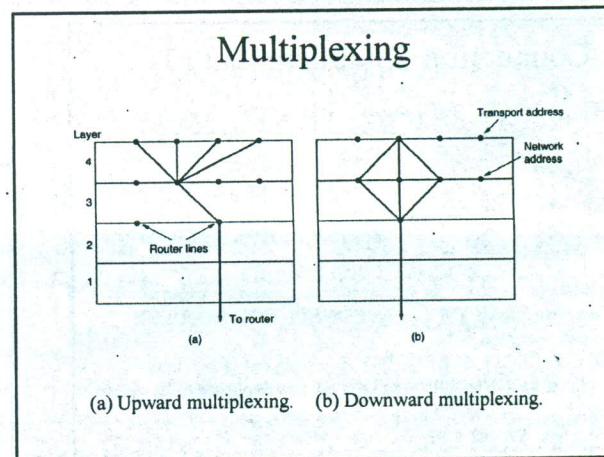
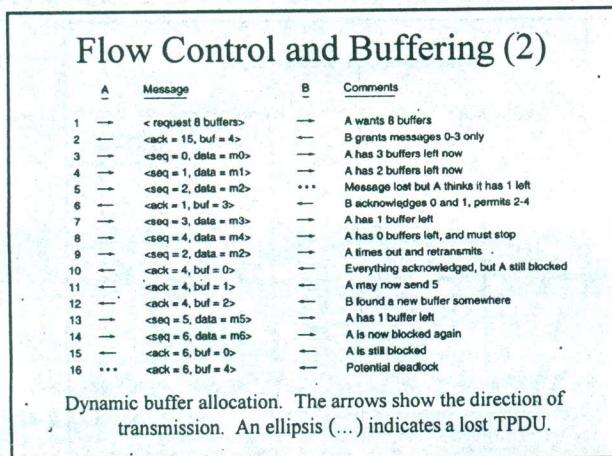
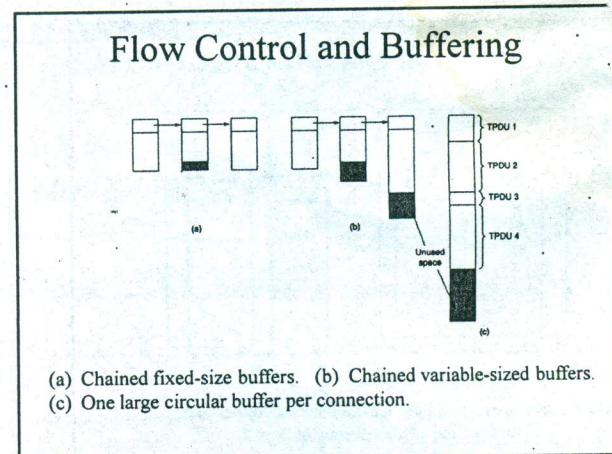
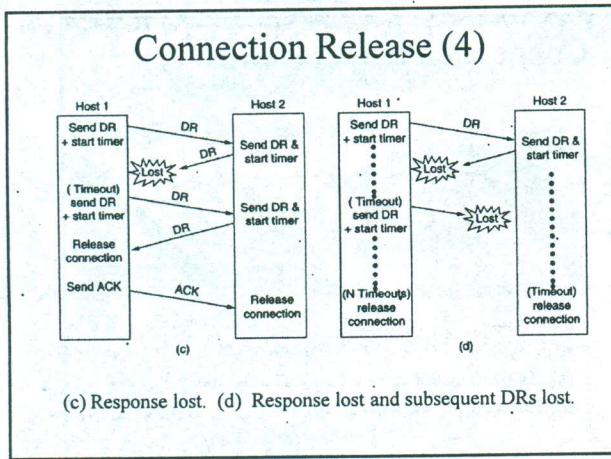
Addressing



TSAPs, NSAPs and transport connections.

Transport Service Access Point
Network Service Access Point
TSAP → like port
NSAP → like IP address
ATM → AAL - SAP (like port)





The Example Transport Entity

Network packet	Meaning
CALL REQUEST	Sent to establish a connection
CALL ACCEPTED	Response to CALL REQUEST
CLEAR REQUEST	Sent to release a connection
CLEAR CONFIRMATION	Response to CLEAR REQUEST
DATA	Used to transport data
CREDIT	Control packet for managing the window

The network layer packets used in our example.

The Example Transport Entity (2)

Each connection is in one of seven states:

1. Idle – Connection not established yet.
2. Waiting – CONNECT has been executed, CALL REQUEST sent.
3. Queued – A CALL REQUEST has arrived; no LISTEN yet.
4. Established – The connection has been established.
5. Sending – The user is waiting for permission to send a packet.
6. Receiving – A RECEIVE has been done.
7. DISCONNECTING – a DISCONNECT has been done locally.

The Example Transport Entity (3)

```
#define MAX_CONN 32
#define MAX_MSG_SIZE 8192
#define MAX_PKT_SIZE 512
#define TIMEOUT 20
#define CRED 1
#define OK 0

#define ERR_FULL -1
#define ERR_REJECT -2
#define ERR_CLOSED -3
#define LOW_ERR -3

typedef int transport_address;
typedef enum {CALL_REQ,CALL_ACC,CLEAR_REQ,CLEAR_CONF,DATA_PKT,CREDIT} pkt_type;
typedef enum {IDLE,WAITING,QUEUED,ESTABLISHED,SENDING,RECEIVING,DISCONN} state;

/* Global variables */
transport_address listen_address; /* local address being listened to */
int listen_conn; /* connection identifier for listen */
unsigned char data[MAX_PKT_SIZE]; /* scratch area for packet data */

struct conn {
    transport_address local_address, remote_address;
    state state; /* state of this connection */
    unsigned char *user_buf_addr; /* pointer to receive buffer */
    int byte_count; /* send/receive count */
    int clr_req_received; /* set when CLEAR_REQ packet received */
    int timer; /* used to time out CALL_REQ packets */
    int credits; /* number of messages that may be sent */
} conn[MAX_CONN + 1];
```

The Example Transport Entity (4)

```
void sleep(void); /* prototypes */
void wakeup(void);
void to_net(int cid, int q, int m, pkt_type pt, unsigned char *p, int bytes);
void from_net(int *cid, int *q, int *m, pkt_type *pt, unsigned char *p, int *bytes);

int listen(transport_address t) /* User wants to listen for a connection. See if CALL_REQ has already arrived. */
{
    int i, found = 0;

    for (i = 1; i <= MAX_CONN; i++) /* search the table for CALL_REQ */
        if (conn[i].state == QUEUED && conn[i].local_address == t) {
            found = i;
            break;
        }

    if (found == 0) /* No CALL_REQ is waiting. Go to sleep until arrival or timeout. */
        listen_address = t; sleep(); i = listen_conn;
}

conn[i].state = ESTABLISHED; /* connection is ESTABLISHED */
conn[i].timer = 0; /* timer is not used */
```

The Example Transport Entity (5)

```
listen_conn = 0;
to_net(0, 0, CALL_ACC, data, 0); /* tell net to accept connection */
return();
}

int connect(transport_address l, transport_address r)
{ /* User wants to connect to a remote process; send CALL_REQ packet. */
int i;
struct conn *cptr;

data[0] = r; data[1] = l; /* CALL_REQ packet needs these */
i = MAX_CONN; /* search table backward */
while (conn[i].state != IDLE && i > 1) i--;
if (conn[i].state == IDLE) {
    /* Make a table entry that CALL_REQ has been sent. */
    cptr = &conn[i];
    cptr->local_address = l; cptr->remote_address = r;
    cptr->state = WAITING; cptr->clr_req_received = 0;
    cptr->credits = 0; cptr->timer = 0;
    to_net(i, 0, 0, CALL_REQ, data, 2);
    sleep(); /* wait for CALL_ACC or CLEAR_REQ */
    if (cptr->state == ESTABLISHED) return();
    if (cptr->clr_req_received) {
        /* Other side initiated call. */
        cptr->state = IDLE; /* back to IDLE state */
        to_net(i, 0, 0, CLEAR_CONF, data, 0);
        return(ERR_REJECT);
    }
} else return(ERR_FULL); /* reject CONNECT: no table space */
}
```

The Example Transport Entity (6)

```
int send(int cid, unsigned char bufptr[], int bytes)
{ /* User wants to send a message. */
int i, count, m;
struct conn *cptr = &conn[cid];

/* Enter SENDING state. */
cptr->state = SENDING;
cptr->byte_count = 0; /* bytes sent so far this message */
if (cptr->clr_req_received == 0 && cptr->credits == 0) sleep();
if (cptr->clr_req_received == 0) {
    /* Credit available; split message into packets if need be. */
    do {
        if (bytes - cptr->byte_count > MAX_PKT_SIZE) /* multipacket message */
            count = MAX_PKT_SIZE; m = 1; /* more packets later */
        else /* single packet message */
            count = bytes - cptr->byte_count; m = 0; /* last pkt of this message */
    } while (cptr->byte_count < bytes); /* loop until whole message sent */
}

for (i = 0; i < count; i++) data[i] = bufptr[cptr->byte_count + i];
to_net(cid, 0, m, DATA_PKT, data, count); /* send 1 packet */
cptr->byte_count = cptr->byte_count + count; /* increment bytes sent so far */
}
```

The Example Transport Entity (7)

```

    // each message uses up one credit */

    if (cptr->credits == 0)
        cptr->state = ESTABLISHED;
    return(OK);
}

else {
    cptr->state = ESTABLISHED;
    return(ERR_CLOSED);                                /* send failed: peer wants to disconnect */
}
}

int receive(int cid, unsigned char bufptr[], int +bytes)
/* User is prepared to receive a message. */
struct conn *cptr = &conn[cid];

if (cptr->clr_req_received == 0) {
    /* Connection still established: try to receive. */
    cptr->state = RECEIVING;
    cptr->user_buf_addr = bufptr;
    cptr->byte_count = 0;
    data[0] = CRED;
    data[1] = 1;
    to_ner(cid, 1, 0, CREDIT, data, 2);           /* send credit */
    sleep();                                     /* block awaiting data */
    +bytes = cptr->byte_count;
}

cptr->state = ESTABLISHED;
return(cptr->clr_req_received ? ERR_CLOSED : OK);
}

```

The Example Transport Entity (9)

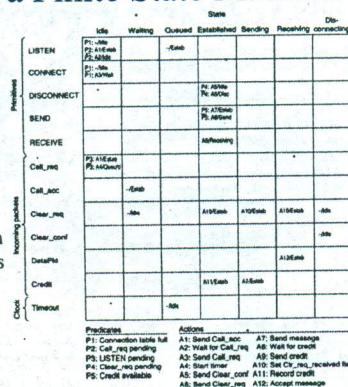
```

switch (phytype) {
    case CALL_REQ:           /* remote user wants to establish connection */
        cpt->remote_address = data[0]; cpt->remote_address = data[1];
        if (cpt->local_address == listen_address) {
            listen_conn = cpt; cpt->state = ESTABLISHED; wakeup();
        } else {
            cpt->state = QUEUED; cpt->timer = TIMEOUT;
        }
        cpt->clr_req_received = 0; cpt->credits = 0;
        break;
    case CALL_ACCEPT:         /* remote user has accepted our CALL_REQ */
        cpt->state = ESTABLISHED;
        wakeup();
        break;
    case CLEAR_REQ:           /* remote user wants to disconnect or reject call */
        cpt->clr_req_received = 1;
        if (cpt->state == DISCONN) cpt->state = IDLE; /* clear collision */
        if (cpt->state == WAITING || cpt->state == RECEIVING) cpt->state = SENDING; wakeup();
        break;
    case CLEAR_CONF:          /* remote user agrees to disconnect */
        cpt->state = IDLE;
        break;
    case CREDIT:               /* remote user is waiting for data */
        cpt->credits += data[1];
        if (cpt->state == SENDING) wakeup();
        break;
    case DATA_PKT:             /* remote user has sent data */
        for (i = 0; i < count; i++) cpt->user_buf_addr[cpt->byte_count + i] = data[i];
        cpt->byte_count += count;
        if (m == 0) wakeup();
}

```

The Example as a Finite State Machine

The example protocol as a finite state machine. Each entry has an optional predicate, an optional action, and the new state. The tilde indicates that no major action is taken. An overbar above a predicate indicate the negation of the predicate. Blank entries correspond to impossible or invalid events.



The Example Transport Entity (8)

```

int disconnect(int cid)
{ /* User wants to release a connection. */
  struct conn *cptr = &conn[cid];

  if (cptr->clr_req_received) {           /* other side initiated termination */
    cptr->state = IDLE;                  /* connection is now released */
    to_net(cid, 0, 0, CLEAR_CONF, data, 0);
  } else {                                /* we initiated termination */
    cptr->state = DISCONN;               /* not released until other side agrees */
    to_net(cid, 0, 0, CLEAR_REQ, data, 0);
  }
  return(OK);
}

void packet_arrival(void)
/* A packet has arrived, get and process it. */
{
  int cid;
  int count, i, q, m;
  pkt_type type;
  unsigned char data[MAX_PKT_SIZE];        /* data portion of the incoming packet */
  struct conn *cptr;

  from_net(&cid, &q, &m, &type, data, &count); /* go get it */
  cptr = &conn[cid];
}

```

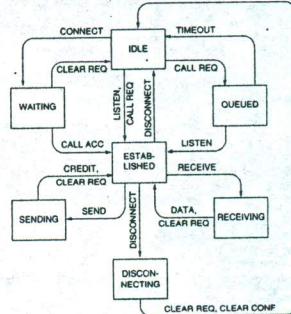
The Example Transport Entity (10)

```

void clock(void)
{ /* The clock has ticked, check for timeouts of queued connect requests. */
    int i;
    struct conn *cptr;
    for (i = 1; i <= MAX_CONN; i++) {
        cptr = &conn[i];
        if (cptr->timer != 0) {                                /* timer was running */
            cptr->timer--;
            if (cptr->timer == 0) {                            /* timer has now expired */
                cptr->state = IDLE;
                to_net(i, 0, 0, CLEAR_REQ, data, 0);
            }
        }
    }
}

```

The Example as a Finite State Machine (2)



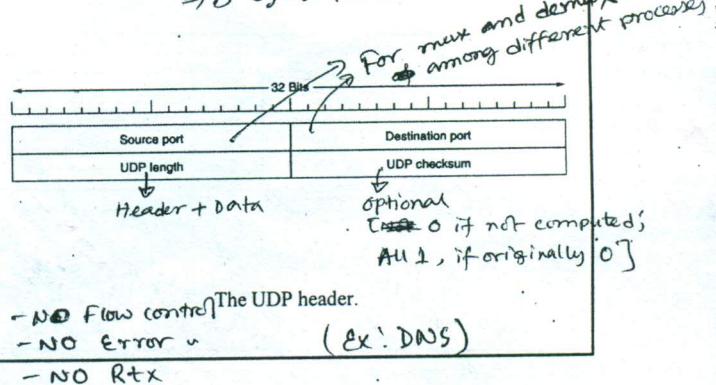
The example protocol in graphical form. Transitions that leave the connection state unchanged have been omitted for simplicity.

The Internet Transport Protocols: UDP

- Introduction to UDP
- Remote Procedure Call
- The Real-Time Transport Protocol

Introduction to UDP

→ 8 bytes header



- NO Flow control

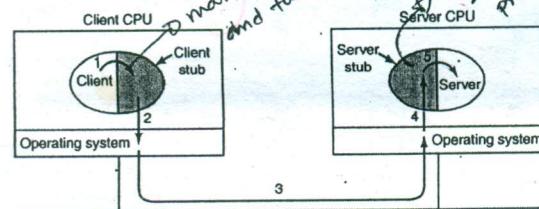
- NO Error

(Ex: DNS)

- NO Rtx

(RPC)
⇒ Call a procedure from
a distant m/c
(caller)

Remote Procedure Call



- Difficulties!
- ① No pointer passing as addr space is different
- ② No global var
- ③ In weakly-typed lang (ex: c), where array can be terminated with own special character, the length of array cannot be determined
- ④ Deducing parameter types from formal spec (ex: printf) may not always possible [int → might be 2B or 4B]

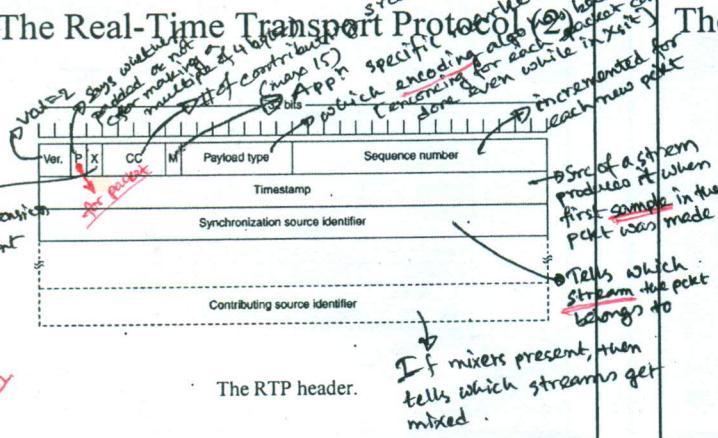
⇒ RPC generally uses UDP, But it is NOT mandatory. The stubs are shaded.

① If pair or result is larger than UDP pkt

② If the requested op is not idempotent (can NOT be repeated safely)

TCP in place of UDP

The Real-Time Transport Protocol

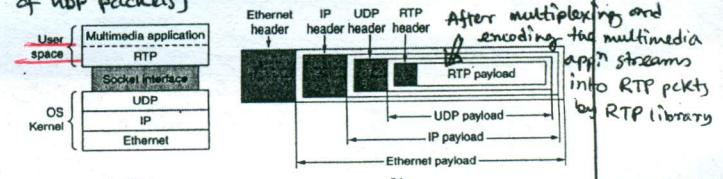


RTCP =>

- i) feedback, synch!, user interface hand
- ii) No data xport
- iii) Name the sources

The Real-Time Transport Protocol

→ Transport protocol implemented in App layer
[multiplex several real-time data streams into a single stream of UDP packets]



(a) The position of RTP in the protocol stack. (b) Packet nesting.

(b) UDP pkts generated at one end of the socket

Multimedia app streams ⇒ fed to RTP library ⇒ stuff to socket

Multiplex the streams and encode them into RTP packets by RTP library

⇒ RTP packets are sent to the other end of the socket

The Internet Transport Protocols: TCP

- Introduction to TCP
- The TCP Service Model
- The TCP Protocol
- The TCP Segment Header
- TCP Connection Establishment
- TCP Connection Release
- TCP Connection Management Modeling
- TCP Transmission Policy
- TCP Congestion Control
- TCP Timer Management
- Wireless TCP and UDP
- Transactional TCP

TCP \Rightarrow full duplex + P2P

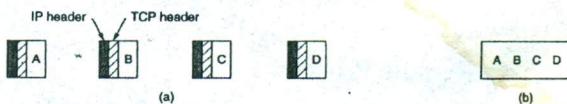
- ② Byte Stream (not msg stream)
- ③ Does NOT support Broadcasting & multicasting

④ Byte stream (Net msg stream)

The TCP Service Model

Port	Protocol	Use
21	FTP	File transfer
23	Telnet	Remote login
25	SMTP	E-mail
69	TFTP	Trivial File Transfer Protocol
79	Finger	Lookup info about a user
80	HTTP	World Wide Web
110	POP-3	Remote e-mail access
119	NNTP	USENET news

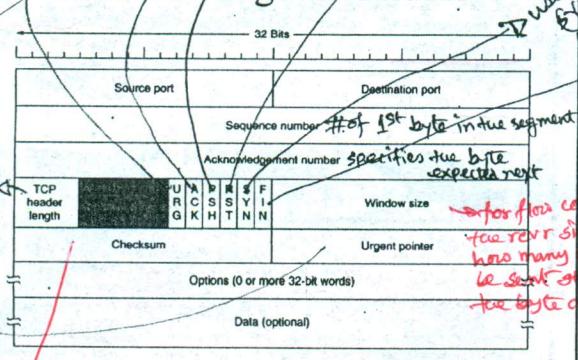
Some assigned ports



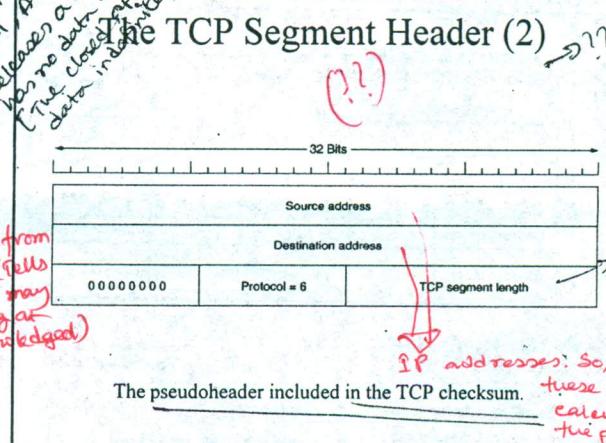
(a) Four 512-byte segments sent as separate IP datagrams.

The 2048 bytes of data delivered to the application in a single READ CALL.

The TCP Segment Header



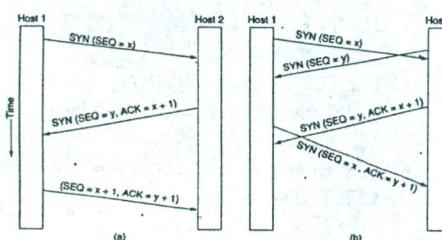
TCP Header



The pseudoheader included in the TCP checksum

These is checksum
calculated violates
the protocol
hierarchy.

TCP Connection Establishment



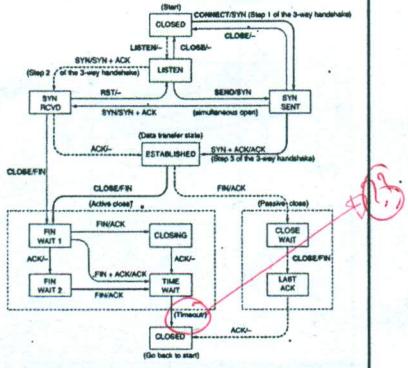
- (a) TCP connection establishment in the normal case.
- (b) Call collision.

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

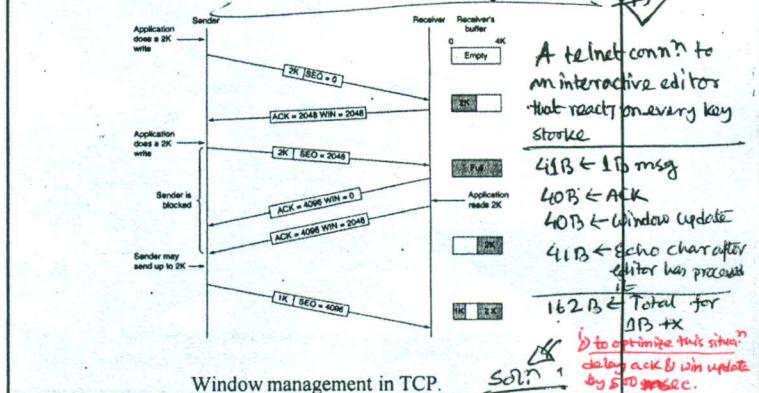
The states used in the TCP connection management finite state machine.

TCP Connection Management Modeling (2)

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



TCP Transmission Policy



problem for
slow sender

A telnet conn'n to
an interactive editor
that reacts on every key
stroke

$41B \leftarrow 1B$ msg
 $40B \leftarrow ACK$
 $40B \leftarrow$ Window update
 $41B \leftarrow$ echo char after
editor has processed
1B

$\Delta B + x$
to optimize this situation
delay ack B win update
by 500 msec.

If sender sends 1B at a time, then the sender waits after sending the 1st B. It buffers all generated pkts until ACK received. All buffered pkts are sent in TCP segment(s) after getting ACK

Nagle's
Alg O

can work
together TCP Conge

Problem for
Slow receiver

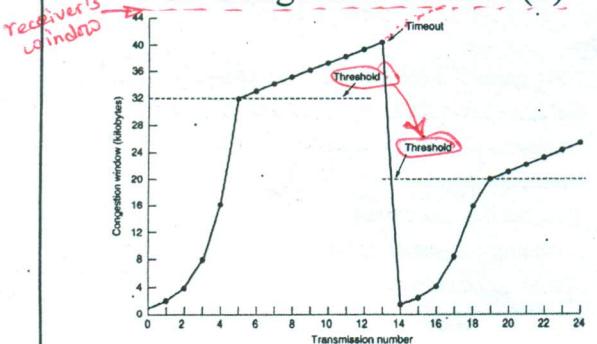
44

Don't send window update just after getting 1 B space. → Header
Send window update Header
only after getting 1 Byte
decent amount of space available
(max segment size, half buffer)

↑
Silly window
send win update after
getting this amount of space available

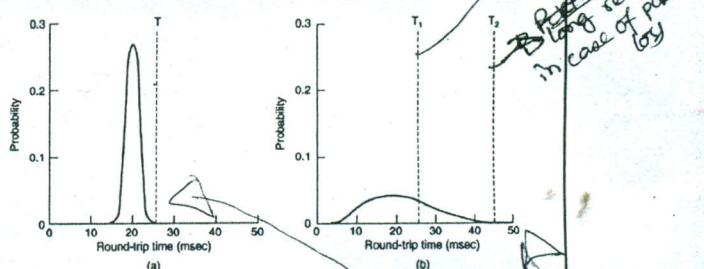
Δ
Clark's Algo

TCP Congestion Control (2)



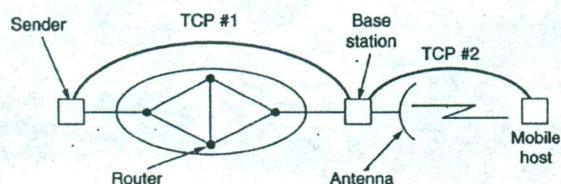
An example of the Internet congestion algorithm.

TCP Timer Management



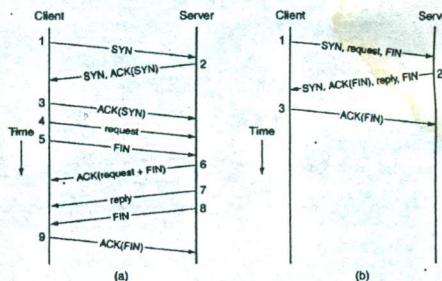
(a) Probability density of ACK arrival times in the data link layer.
 (b) Probability density of ACK arrival times for TCP.

Wireless TCP and UDP



Splitting a TCP connection into two connections.

Transitional TCP



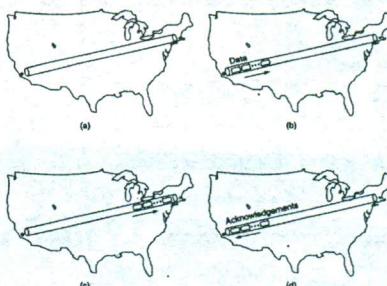
(a) RPC using normal TPC.

(b) RPC using T/TCP.

Performance Issues

- Performance Problems in Computer Networks
- Network Performance Measurement
- System Design for Better Performance
- Fast TPDU Processing
- Protocols for Gigabit Networks

Performance Problems in Computer Networks



The state of transmitting one megabit from San Diego to Boston
 (a) At $t = 0$, (b) After 500 μ sec, (c) After 20 msec, (d) after 40 msec.

Network Performance Measurement

The basic loop for improving network performance.

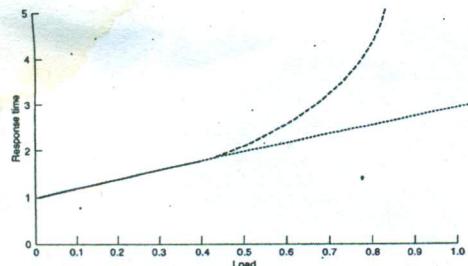
1. Measure relevant network parameters, performance.
2. Try to understand what is going on.
3. Change one parameter.

System Design for Better Performance

Rules:

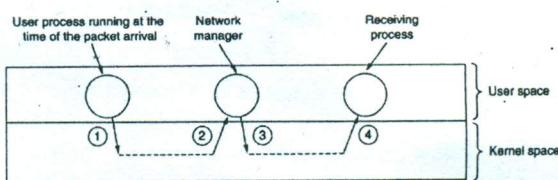
1. CPU speed is more important than network speed.
2. Reduce packet count to reduce software overhead.
3. Minimize context switches.
4. Minimize copying.
5. You can buy more bandwidth but not lower delay.
6. Avoiding congestion is better than recovering from it.
7. Avoid timeouts.

System Design for Better Performance (2)



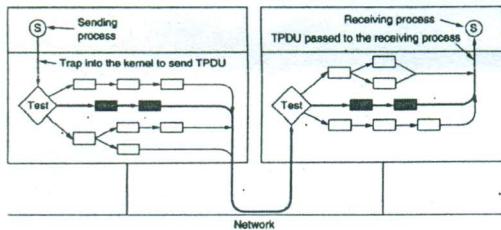
Response as a function of load.

System Design for Better Performance (3)



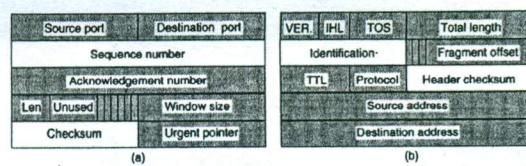
Four context switches to handle one packet with a user-space network manager.

Fast TPDU Processing

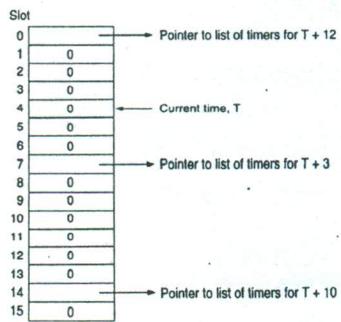


The fast path from sender to receiver is shown with a heavy line.
The processing steps on this path are shaded.

Fast TPDU Processing (2)

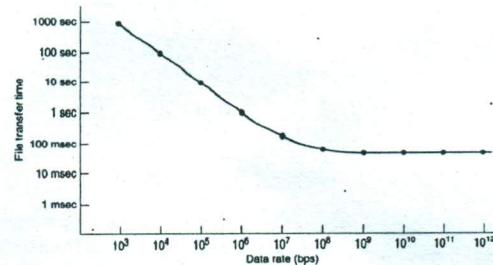


Fast TPDU Processing (3)



A timing wheel.

Protocols for Gigabit Networks



Time to transfer and acknowledge a 1-megabit file over a 4000-km line.