# Network Programming

K Hari Babu
Department of Computer Science & Information Systems

**BITS** Pilani
Pilani Campus

**Outline**

# Outline

- Name conversion
  - getaddrinfo()
  - getnameinfo()
  - gethostbyname()
- Advanced I/O
  - recv(), send()
  - readv(), writev()
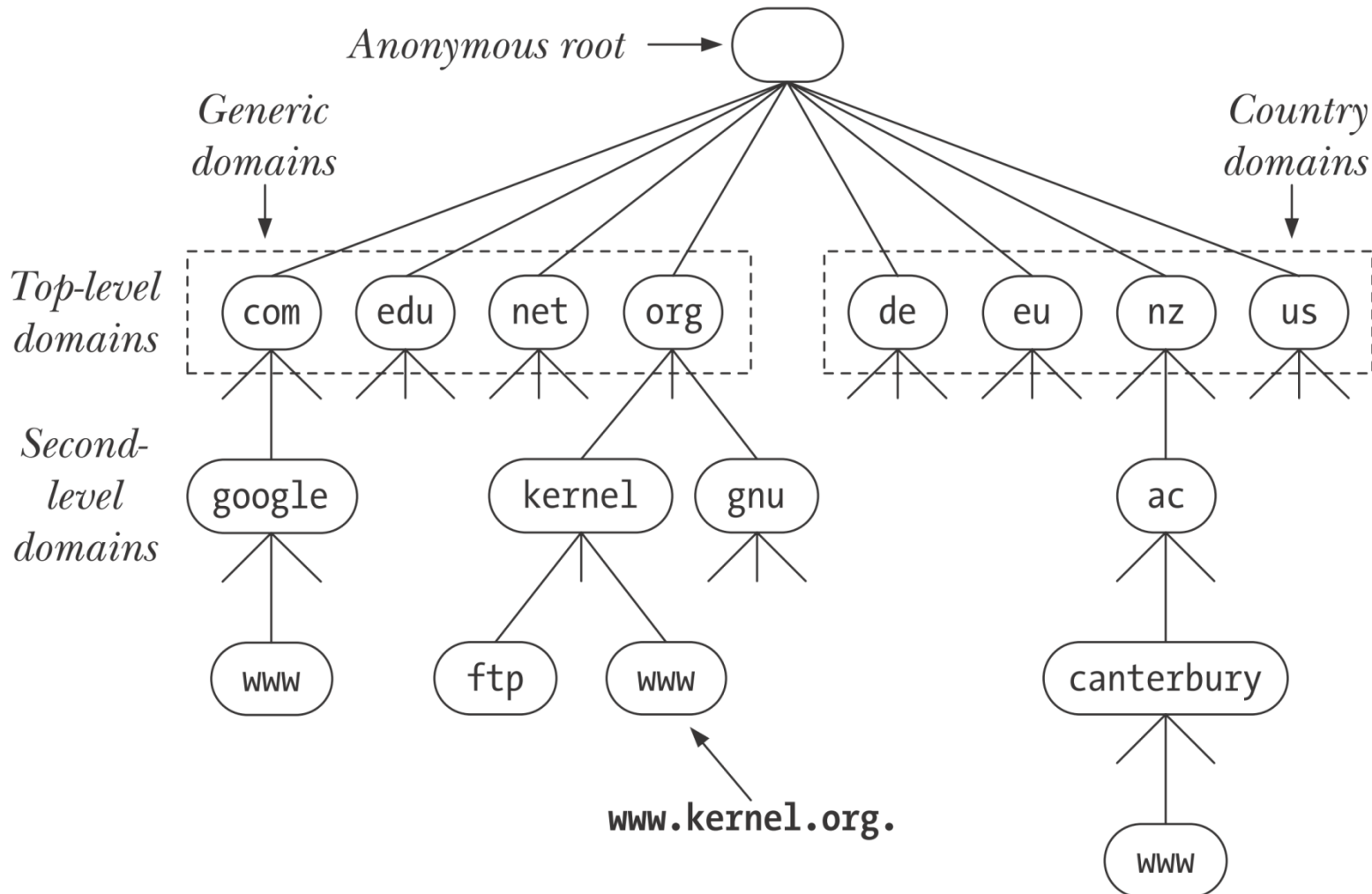  - recvmsg(), sendmsg()

- Framing & Encoding
  - HTTP

# Name and Address Conversions (T1: ch11)

# DNS

- Two functions which we use to convert a domain name to IP address:
  - gethostbyname()  - obsolete
  - getaddrinfo()  - supports both IPv6 and IPv4
- The DNS is used primarily to map between hostnames and IP addresses.
  - A hostname can be either a simple name, such as *solaris* or *freebsd*, or a fully qualified domain name (FQDN), such as *solaris.unpbook.com.*
- Before DNS,  mappings  between  hostnames  and  IP addresses were defined in a manually maintained local file, /etc/hosts

```
1  # IP-address      canonical hostname      [aliases]
2  127.0.0.1         localhost
```

# DNS

- The /etc/hosts scheme scales poorly.

- DNS was devised to address this problem.
  - Hostnames are organized into a hierarchical namespace.
    - Node in the DNS hierarchy has a label (name), which may be up to 63 characters.
    - At the root of the hierarchy is an unnamed node, the "anonymous root."
  - A node's domain name consists of all of the names from that node up to the root concatenated together, with each name separated by a period ( . )
  - No single organization or system manages the entire hierarchy.
    - Instead, there is a hierarchy of DNS servers, each of which manages a branch (a zone) of the tree.
    - For adding a host, admin has to add it to local name server.
  - DNS servers employ caching techniques to avoid unnecessary communication for frequently queried domain names.

# DNS Lookup

- Not every name server knows about every other name server.

- Name server must know the IP address of root servers.

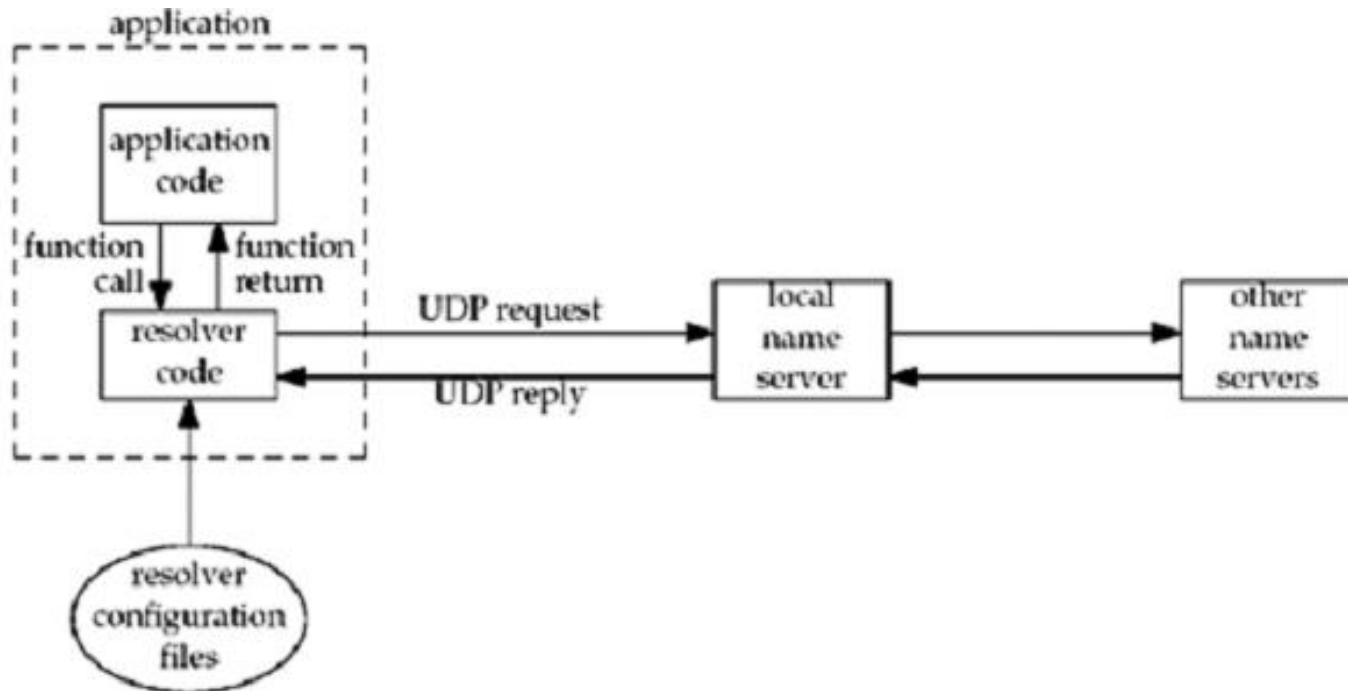- Root servers know the name and location for all second-level domains.

# Recursive and Iterative Lookups

- DNS resolution requests fall into two categories: recursive and iterative.

  o In a recursive request, the requester asks the server to handle the entire task of resolution.

- When an application on the local host calls *getaddrinfo*(), that function makes a recursive request to the local DNS server.

- If the local DNS server does not itself have the information to perform the resolution, it resolves the domain name iteratively.

Resolver is part of the application

# The /etc/services File

- Well-known port numbers are centrally registered by IANA.
  - Each of these ports has a corresponding service name.
    - Because service numbers are centrally managed and are less volatile than IP addresses, an equivalent of the DNS server is usually not necessary. Instead, the port numbers and service names are recorded in the file /etc/services .
  - The getaddrinfo() and getnameinfo() functions use the information in this file to convert service names to port numbers and vice versa.

# The /etc/services File

```
 1   # Service name    port/protocol    [aliases]
 2   echo              7/tcp            Echo      # echo service
 3   echo              7/udp            Echo
 4   ssh               22/tcp                     # Secure Shell
 5   ssh               22/udp
 6   telnet            23/tcp                     # Telnet
 7   telnet            23/udp
 8   smtp              25/tcp                     # Simple Mail Transfer Protocol
 9   smtp              25/udp
10   domain            53/tcp                     # Domain Name Server
11   domain            53/udp
12   http              80/tcp                     # Hypertext Transfer Protocol
13   http              80/udp
14   ntp               123/tcp                    # Network Time Protocol
15   ntp               123/udp
16   login             513/tcp                    # rlogin(1)
17   who               513/udp                    # rwho(1)
18   shell             514/tcp                    # rsh(1)
19   syslog            514/udp                    # syslog
```

# Host and Service Conversion

- The getaddrinfo() function converts host and service names to IP addresses and port numbers.
    - successor to the obsolete gethostbyname() and getservbyname() functions

- Given a host name and a service name, getaddrinfo() returns a list of socket address structures, each of which contains an IP address and port number.

```
1   #include <sys/socket.h>
2   #include <netdb.h>
3   int getaddrinfo(const char * host , const char * service ,
4           const struct addrinfo * hints , struct addrinfo ** result );
5   //Returns 0 on success, or nonzero on error
```

# The getaddrinfo() Function

- As input, getaddrinfo() takes the arguments host, service, and hints.
  - Host:
    - It can be hostname or numeric address string 172.24.2.19
  - Service:
    - This contains either service name or port number.
  - Hints:
    - The hints argument points to an addrinfo structure that specifies further criteria for selecting the socket address structures returned via result.

```
1  #include <sys/socket.h>
2  #include <netdb.h>
3  int getaddrinfo(const char * host , const char * service ,
4         const struct addrinfo * hints , struct addrinfo ** result );
5  //Returns 0 on success, or nonzero on error
```
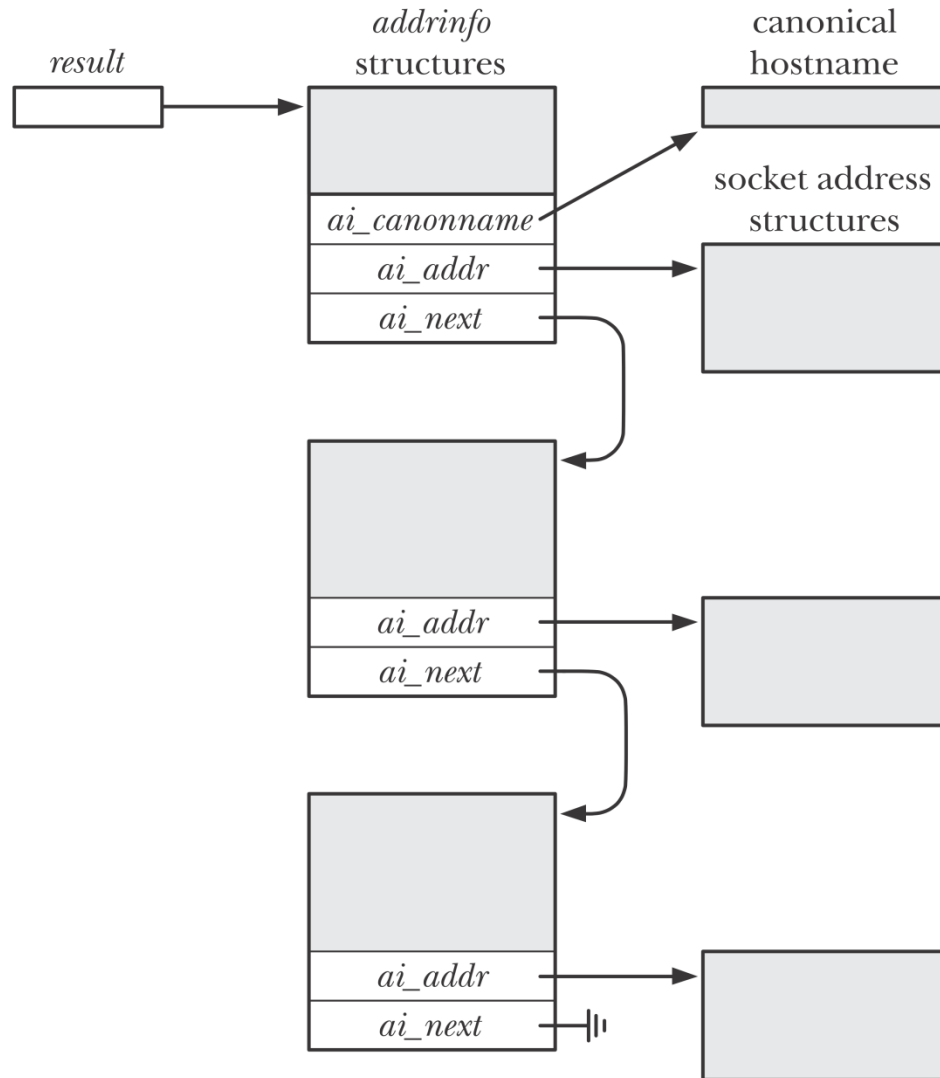
# Addrinfo structure

- As output, getaddrinfo() dynamically allocates a linked list of addrinfo structures and sets result pointing to the beginning of this list.

- Each of these addrinfo structures includes a pointer to a socket address structure corresponding to host and service.

```
1   struct addrinfo {
2       int     ai_flags;          /* Input flags (AI_* constants) */
3       int     ai_family;         /* Address family */
4       int     ai_socktype;       /* Type: SOCK_STREAM, SOCK_DGRAM */
5       int     ai_protocol;       /* Socket protocol */
6       size_t  ai_addrlen;        /* Size of structure pointed to by ai_addr */
7       char    *ai_canonname;     /* Canonical name of host */
8       struct sockaddr *ai_addr;  /* Pointer to socket address structure */
9       struct addrinfo *ai_next;  /* Next structure in linked list */
10  };
```

# Result

# The *hints* Argument

- Hints is either a null pointer or a pointer to an addrinfo structure.
  - the caller fills in this structure with hints about the types of information the caller wants returned.
- The members of the hints structure that can be set by the caller are:
  - ai_flags (zero or more AI_XXX values OR'ed together)
  - ai_family (an AF_xxx value)
  - ai_socktype (a SOCK_xxx value)
  - ai_protocol
- For example,
  - if the specified service is provided for both TCP and UDP, set ai_socktype member of the hints structure to SOCK_DGRAM. Then only information returned will be for datagram sockets.

# The *hints* Argument

- AI_PASSIVE
  - Returnsocket address structures suitable for a passive open.
  - If *host* is null, then IP will be INADDR_ANY or IN6ADDR_ANY_INIT.
- AI_CANONNAME
  - Tells the function to return the canonical name of the host.
- AI_NUMERICHOST
  - the hostname argument must be a numeric address string.
  - Prevents name resolution.
- AI_NUMERICSERV
  - the service argument must be a decimal port number string.
  - Prevents any kind of name-to-service  resolution;
- AI_V4MAPPED
  - If specified along with an ai_family of AF_INET6, then returns IPv4-mapped IPv6 addresses corresponding to A records if there are no available AAAA records.

# The *hints* Argument

- Returns linked list of addrinfo structures, linked through the ai_next pointer.

- There are two ways that multiple structures can be returned:
  - Multiple ips per hostname; one sockaddr structure for each ip
  - Service is provided for multiple socket types; SOCK_STREAM or SOCK_DGRAM

- For example, if no hints are provided and if the *domain* service is looked up for a host with two IP addresses, four addrinfo structures are returned:
  - One for the first IP address and a socket type of SOCK_STREAM
  - One for the first IP address and a socket type of SOCK_DGRAM
  - One for the second IP address and a socket type of SOCK_STREAM
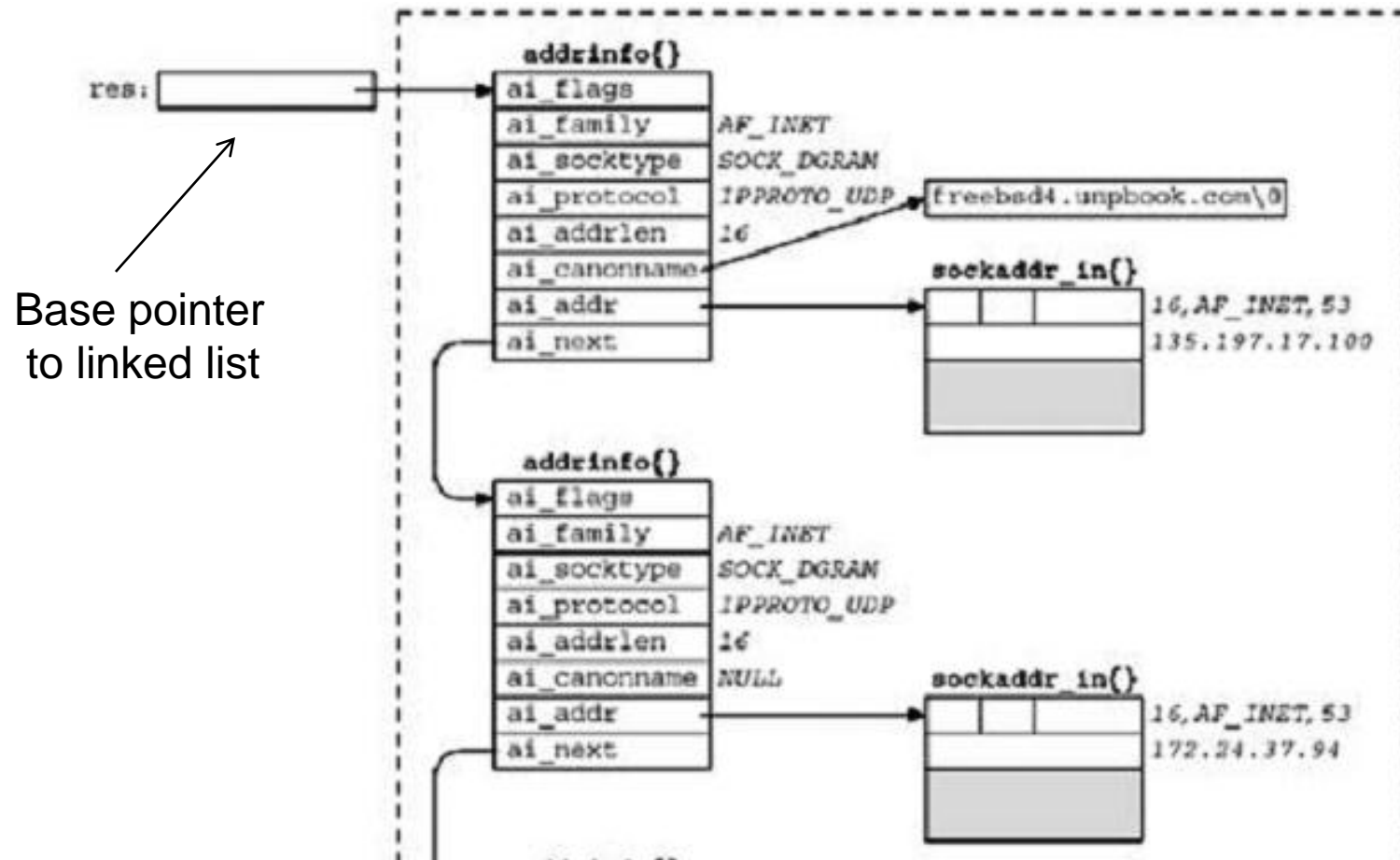  - One for the second IP address and a socket type of SOCK_DGRAM

- Consider the following code

```
1   struct addrinfo hints, *res;
2   bzero(&hints, sizeof(hints) ) ;
3   hints.ai_flags = AI_CANONNAME;
4   hints.ai_family = AF_INET;
5   getaddrinfo("freebsd4", "domain", &hints, &res);
```
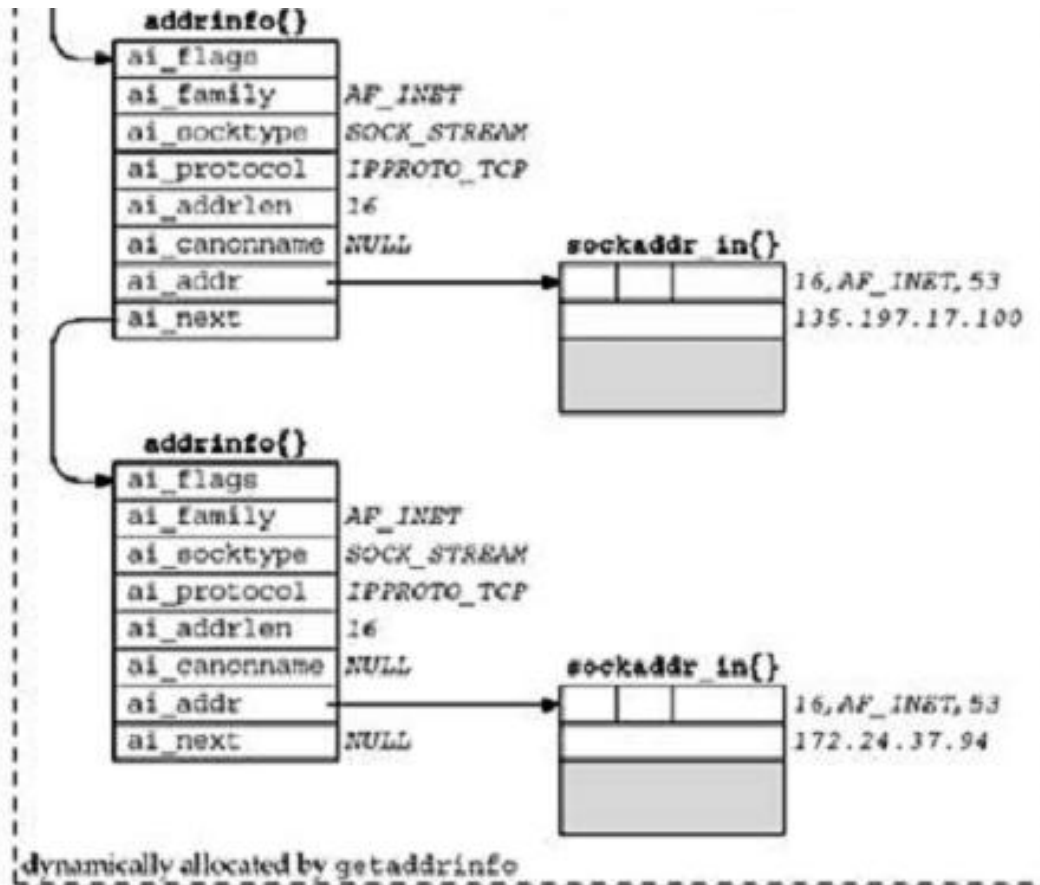
# Result



Base pointer to linked list

# Usage

- Sockaddr structure in addrinfo structures is ready for
  - a call to socket
  - then either a call to connect or sendto (for a client), or bind (for a server).
- The arguments to socket are the members *ai_family, ai_socktype, and ai_protocol.*
- The second and third arguments to either connect or bind are *ai_addr*, and *ai_addrlen*

# On client side

```c
int tcp_connect (const char *host, const char *serv)
{
  int sockfd, n;
  struct addrinfo hints, *res, *ressave;
 bzero(&hints, sizeof (struct addrinfo));
 hints.ai_family = AF_UNSPEC;
 hints.ai_socktype = SOCK_STREAM;
 if ( (n = getaddrinfo (host, serv, &hints, &res)) != 0)
     err_quit("tcp_connect error for %s, %s: %s",
             host, serv, gai_strerror (n));
 ressave = res;
 do {
     sockfd = socket (res->ai_family, res->ai_socktype, res->ai_protocol);
     if (sockfd < 0)
         continue;            /*ignore this one */
     if (connect (sockfd, res->ai_addr, res->ai_addrlen) == 0)
         break;               /* success */
     Close(sockfd);           /* ignore this one */
 } while ( (res = res->ai_next) != NULL);
  if (res == NULL)            /* errno set from final connect() */
 err_sys ("tcp_connect error for %s, %s", host, serv);
 freeaddrinfo (ressave);
 return (sockfd);
}
```

# Server side usage

```
1    int tcp_listen(const char *host, const char *serv, socklen_t *addrlenp)
2    {   struct addrinfo hints, *res, *ressave;
3        bzero(&hints, sizeof (struct addrinfo)) ;
4        hints.ai_flags = AI_PASSIVE;
5        hints.ai_family = AF_UNSPEC;
6        hints.ai_socktype = SOCK_STREAM;
7        if ( (n = getaddrinfo (host, serv, &hints, &res)) != 0)
8            err_quit("tcp_listen error for %s, %s: %s",
9                    host, serv, gai_strerror(n)) ;
10       ressave = res;
11   do {listenfd =socket(res->ai_family, res->ai_socktype, res->ai_protocol);
12       if (listenfd < 0)continue;            /* error, try next one */
13       setsockopt(listenfd, SOL_SOCKET, SO_REUSEADDR, &on, sizeof (on) ) ;
14       if (bind(listenfd, res->ai_addr, res->ai_addrlen) == 0)
15           break;                    /* success */
16       close (listenfd);        /* bind error, close and try next one */
17       } while ( (res = res->ai_next) != NULL);
18       if (res == NULL)            /* errno from final socket () or bind () */
19           err_sys ("tcp_listen error for %s, %s", host, serv);
20       listen (listenfd, LISTENQ);
21       if (addrlenp)
22           *addrlenp = res->ai_addrlen;    /* return size of protocol address */
23       freeaddrinfo (ressave);
24       return (listenfd);
25   }
```

```
1    listenfd = tcp_listen (NULL, argv[1], NULL);
2    |
```

# Freeing *addrinfo* Lists

- The *getaddrinfo()* function dynamically allocates memory for all of the structures referred to by *result*.

```
1  #include <sys/socket.h>
2  #include <netdb.h>
3  void freeaddrinfo(struct addrinfo * result );
```

- Diagnosing errors
  - Errors returned by getaddrinfo() are not stored in errno.
  - They have to be looked up using

```
1  #include <netdb.h>
2  const char *gai_strerror(int  errcode );
3  //Returns pointer to string containing error message
```

# Error Codes

| Error constant | Description |
|---|---|
| EAI_ADDRFAMILY | No addresses for *host* exist in *hints.ai_family* (not in SUSv3, but defined on most implementations; *getaddrinfo()* only) |
| EAI_AGAIN | Temporary failure in name resolution (try again later) |
| EAI_BADFLAGS | An invalid flag was specified in *hints.ai_flags* |
| EAI_FAIL | Unrecoverable failure while accessing name server |
| EAI_FAMILY | Address family specified in *hints.ai_family* is not supported |
| EAI_MEMORY | Memory allocation failure |
| EAI_NODATA | No address associated with *host* (not in SUSv3, but defined on most implementations; *getaddrinfo()* only) |
| EAI_NONAME | Unknown *host* or *service*, or both *host* and *service* were NULL, or AI_NUMERICSERV specified and *service* didn't point to numeric string |
| EAI_OVERFLOW | Argument buffer overflow |
| EAI_SERVICE | Specified *service* not supported for *hints.ai_socktype* (*getaddrinfo()* only) |
| EAI_SOCKTYPE | Specified *hints.ai_socktype* is not supported (*getaddrinfo()* only) |
| EAI_SYSTEM | System error returned in *errno* |

# getnameinfo()

- This function is the converse of *getaddrinfo*(). Takes socket address structure and returns host, and service.

```
1  #include <sys/socket.h>
2  #include <netdb.h>
3  int getnameinfo(const struct sockaddr * addr , socklen_t  addrlen,
4  char * host , size_t  hostlen , char * service ,
5  size_t  servlen , int  flags );
6  //Returns 0 on success, or nonzero on error
```

- NI_DGRAM
  - Default is stream socket service. This make datagram service.
- NI_NAMEREQD
  - Default: if no name found, return numeric ip and port. If this is on, error is returned.

# gethostbyname()

- Returns only A type (IPv4) records

```
1   #include <netdb.h>
2    struct hostent *gethostbyname (const char *hostname);
3    //Returns: non-null pointer if OK,NULL on error with h_errno set
4    struct hostent {
5      char  *h_name;        /* official (canonical) name of host */
6      char **h_aliases;     /* pointer to array of pointers to alias names */
7      int    h_addrtype;    /* host address type: AF_INET */
8      int    h_length;      /* length of address: 4 */
9      char **h_addr_list;   /* ptr to array of ptrs with IPv4 addrs */
10   };
```

- can return multiple IP addresses pointed by h_addr_list.

# gethostbyname() example

```c
void main(int argc, char* argv){
char    *ptr, **pptr;
char      str [INET_ADDRSTRLEN];
struct hostent *hptr;
   ptr = argv[1];
   if ( (hptr = gethostbyname (ptr) ) == NULL) {
        err_msg ("gethostbyname error for host: %s: %s",
        ptr, hstrerror (h_errno) );
printf ("official hostname: %s\n", hptr->h_name);
for (pptr = hptr->h_aliases; *pptr ! = NULL; pptr++)
printf ("\talias: %s\n", *pptr);
switch (hptr->h_addrtype) {
case AF_INET:
  pptr = hptr->h_addr_list;
   for ( ; *pptr != NULL; pptr++)
     printf ("\taddress: %s\n",
        inet_ntop (hptr->h_addrtype, *pptr, str, sizeof (str)));
          break;
default:
err_ret ("unknown address type");
          break;
    }}
}
```

# Advanced I/O Functions (T1: ch 14)

# recv() and send()

- The recv() and send() system calls perform I/O on connected sockets (TCP or connected UDP sockets).

- Socket-Specific I/O System Calls:
  - They provide socket-specific functionality not available with read(0 and write().

```
1  #include <sys/socket.h>
2  ssize_t recv(int  sockfd , void * buffer , size_t  length , int  flags );
3  //Returns number of bytes received, 0 on EOF, or -1 on error
4  ssize_t send(int  sockfd , const void * buffer , size_t  length , int  flags );
5  //Returns number of bytes sent, or -1 on error
```

  - Same as read() and write() except for *flags*.
  - Return values are same as read() and write().

# recv() flags

- MSG_DONTWAIT:
  - perform a non-blocking recv().
  - can be done using fcntl() call but that will make sock fd non-blocking. Here only this operation is non-blocking.

- MSG_OOB:
  - receive out-of-band data on the socket.

- MSG_PEEK:
  - retrieve a copy of the requested bytes from the socket buffer.
  - Data is not removed from the socket buffer.
  - Used for knowing the no of bytes available on the buffer.

- MSG_WAITALL
  - Blocks until *length* bytes are read from socket buffer.
  - May get interrupted by signals.

# send() flags
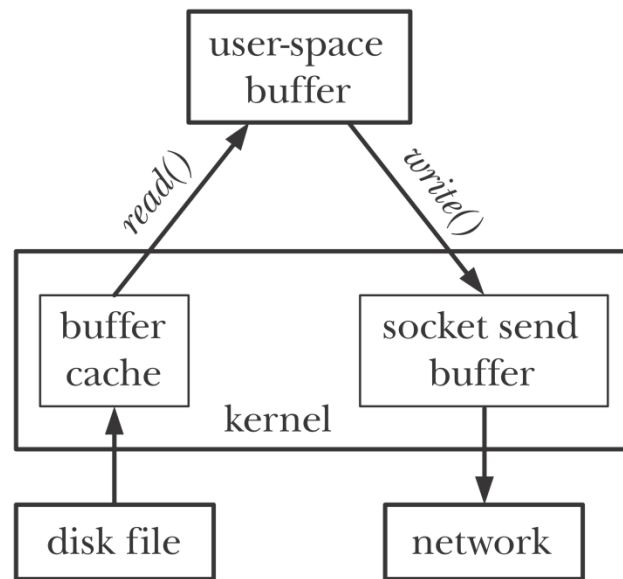
- MSG_DONTWAIT
  - Perform a non-blocking send.

- MSG_MORE
  - Data written using send() or sendto() calls with this flag is packaged into a single datagram until a send() without this flag.

- MSG_NOSIGNAL
  - Do not generate SIGPIPE signal. Return only EPIPE error.

- MSG_OOB
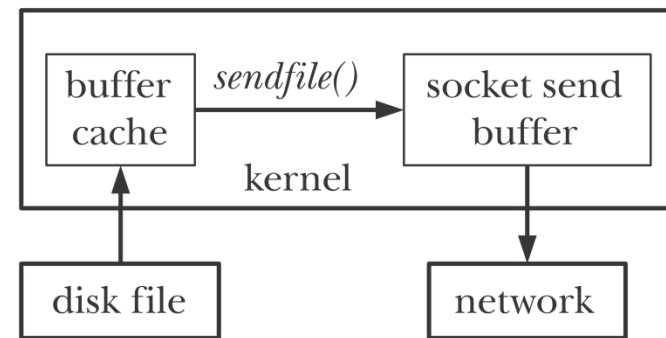  - Write out of band data on TCP.

# sendfile() sys call (R1: 61.4)

- Transferring large file in web servers requires repeated calls to read() and write().
  - This is inefficient.

```
1  while ((n = read(diskfilefd, buf, BUZ_SIZE)) > 0)
2      write(sockfd, buf, n);
```



a) *read() + write()*    b) *sendfile()*

# sendfile() sys call (R1: 61.4)

- The sendfile() sys call is designed to eliminate copying file data into user space.
  - File contents are directly transferred to the socket without going through user space.
  - This is referred as a *zero-copy transfer.*

```
1  #include <sys/sendfile.h>
2  ssize_t sendfile(int out_fd, int  in_fd, off_t * offset, size_t count );
3  //Returns number of bytes transferred, or −1 on error
```

  - *out_fd:* is the socket fd.
  - *In_fd:* is regular file fd.
  - *off_t:* is the offset. This is a value-result argument.
  - *count* is the number of bytes to be transferred.
  - *sendfile* doesn't change the file offset for *in_fd*.

# TCP_CORK socket option

- When a file is requested, web server would first send the HTTP headers and followed by the file contents.
  - o Normally this will result in 2 TCP segments.
  - o Leads to inefficient use of network bandwidth.
- TCP_CORK option if enabled buffers data in TCP until either
  - o upper limit on the size of a segment is reached,
  - o the TCP_CORK option is disabled,
  - o the socket is closed, or
  - o a maximum of 200 milliseconds passes from the time that the first corked byte is written.
- This can be achieved by putting both into single buffer or using *writev()*. But these can't be used with *sendfile*().

# TCP_CORK

```
1   int optval;
2   /* Enable TCP_CORK option on 'sockfd' - subsequent TCP output is corked
3       until this option is disabled. */
4   optval = 1;
5   setsockopt(sockfd, IPPROTO_TCP, TCP_CORK, sizeof(optval));
6   write(sockfd, ...);                    /* Write HTTP headers */
7   sendfile(sockfd, ...);                 /* Send page data */
8   /* Disable TCP_CORK option on 'sockfd' - corked output is now transmitted
9       in a single TCP segment. */
10  optval = 0
11  setsockopt(sockfd, IPPROTO_TCP, TCP_CORK, sizeof(optval));
```

- There is also UDP_CORK option that buffers multiple data outputs into a single datagram.

# readv() and writev()

- The readv() and writev() system calls perform scatter-gather I/O.
- *iov* points to an array of buffers, each in *iovec* structure.

```
1  #include <sys/uio.h>
2  ssize_t readv(int  fd , const struct iovec * iov , int  iovcnt );
3  //Returns number of bytes read, 0 on EOF, or -1 on error
4  ssize_t writev(int  fd , const struct iovec * iov , int  iovcnt );
5  //Returns number of bytes written, or -1 on error
```

```
8   struct iovec {
9       void  *iov_base; /* Start address of buffer */
10      size_t iov_len;  /* Number of bytes to transfer to/from buffer */
11  };
```

# readv() and writev()

- The readv() system call performs scatter input:
  - o Reads from the file and puts the data into the buffer starting at iov[0]. Once the first buffer is full, it goes to another.
- readv() completes atomically:
  - o Kernel performs a single data transfer.
  - o Assured that all the bytes read are contiguous in the file. File offset can't be changed by other process.
- The writev() call performs gather output:
  - o Starting from the first buffer, writes the data contiguously into the file.
  - o Partial write is possible.
- writev() completes atomically.
- readv() and writev() are used for convenience and speed.
  - o Reduce number of sys calls.

# sendmsg() & recvmsg() sys calls

- The sendmsg() and recvmsg() system calls are the most general purpose of the socket I/O system calls.
  - The sendmsg() system call can do everything that is done by write(), send(), and sendto();
  - the recvmsg() system call can do everything that is done by read(), recv(), and recvfrom().

```
1  #include <sys/socket.h>
2  ssize_t recvmsg(int sockfd, struct msghdr *msg, int flags);
3  ssize_t sendmsg(int sockfd, struct msghdr *msg, int flags);
4  //Both return: number of bytes read or written if OK, -1 on error
```

```
1  struct msghdr {
2    void            *msg_name;        /* protocol address */
3    socklen_t        msg_namelen;     /* size of protocol address */
4    struct iovec    *msg_iov;         /* scatter/gather array */
5    int              msg_iovlen;      /* # elements in msg_iov */
6    void            *msg_control;     /* ancillary data (cmsghdr struct) */
7    socklen_t        msg_controllen;  /* length of ancillary data */
8    int              msg_flags;       /* flags returned by recvmsg() */
9  };
```

# sendmsg() & recvmsg() sys calls

- Can be used to send or receive ancillary data (control information).
- Flags are

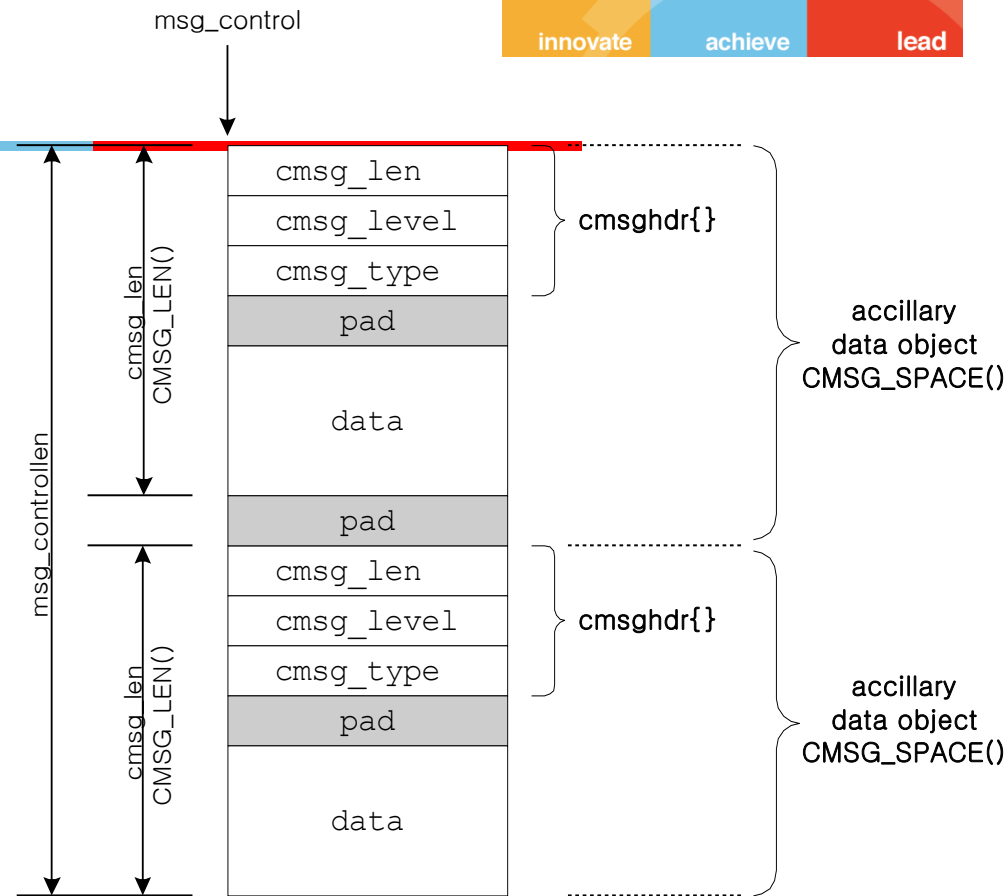| Flag | Examined by:<br>Send flags<br>Sendto flags<br>Sendmsg flags | Examined by:<br>recv flags<br>recvfrom flags<br>recvmsg flags | Returned by:<br><br>Recvmsg<br>msg_flags |
|---|:---:|:---:|:---:|
| MSG_DONTROUTE<br>MSG_DONTWAIT<br>MSG_PEEK<br>MSG_WAITALL | ●<br>● | <br>●<br>●<br>● | |
| MSG_EOR<br>MSG_OOB | ●<br>● | <br>● | ●<br>● |
| MSG_BCAST<br>MSG_MCAST<br>MSG_TRUNC<br>MSG_CTRUNC | | | ●<br>●<br>●<br>● |

# Flags returned by rcvmsg()

- MSG_BCAST
  - is returned if the datagram was received as as a broadcast.
- MSG_MCAST
  - is returned if the datagram was received as a link-layer multicast.
- MSG_TRUNC
  - is returned if the datagram was truncated
- MSG_CTRUNC
  - is returned if the ancillary data was truncated
- MSG_EOR
  - is turned on if the returned data ends a logical record.
- MSG_OOB
  - This flag is never returned for TCP out-of-band data. This flag is returned by other protocol suites (e.g., the OSI protocols).
- MSG_NOTIFICATON
  - This flag is returned for SCTP receivers to indicate that the message read is an event notification, not a data message.

# Ancillary Data

- Ancillary data can be sent and received using the msg_control and msg_controllen members of the msghdr structure.

  o Another term for ancillary data is control information.



```
1 ▾ struct cmsghdr {
2     socklen_t  cmsg_len;   /* length in bytes, including this structure */
3     int        cmsg_level; /* originating protocol */
4     int        cmsg_type;  /* protocol-specific type */
5 ▾        /* followed by unsigned char cmsg_data[] */
6    };
```

# Ancillary Data

- Ancillary data is domain specific.

| Protocol | cmsg_level | Cmsg_type | Description |
|---|---|---|---|
| IPv4 | IPPROTO_IP | IP_RECVDSTADDR | receive destination address with UDP datagram |
| | | IP_RECVIF | receive interface index with UDP datagram |
| IPv6 | IPPROTO_IPV6 | IPV6_DSTOPTS | specify / receive destination options |
| | | IPV6_HOPLIMIT | specify / receive hop limit |
| | | IPV6_HOPOPTS | specify / receive hop-by-hop options |
| | | IPV6_NEXTHOP | specify next-hop address |
| | | IPV6_PKTINFO | specify / receive packet information |
| | | IPV6_RTHDR | specify / receive routing header |
| Unix domain | SOL_SOCKET | SCM_RIGHTS | send / receive descriptors |
| | | SCM_CREDS | send / receive user credentials |

# Ancillary Data

- File descriptors and process credentials can be passed between unrelated processes using ancillary data.



| cmsghdr{} | |
|---|---|
| cmsg_len | *16* |
| cmsg_level | *SOL_SOCKET* |
| cmsg_type | *SCM_RIGHTS* |
| discriptor | |

| cmsghdr{} | |
|---|---|
| cmsg_len | *16* |
| cmsg_level | *SOL_SOCKET* |
| cmsg_type | *SCM_CREDS* |
| **fcred{}** | |

# Framing, Encoding & Decoding

# Encoding & Framing

- Wire format: format of message in the network.
  - Application protocol specifies the wire format for the messages to be exchanged between sender and receiver.
  - Can be: binary or text represented.
- Encoding:
  - Sender has to fill in the fields of the message ( bits or byes or multiple bytes) considering the network byte order.
- Framing
  - How to identify boundaries between messages

# Wrapping TCP Sockets in Streams

- Wrapping TCP Sockets in Streams
  - Using fdopen
  - fread() and fwrite() functions read/write a number of objects from/to the stream.
  - They return no of objects written or read.

```
2  size_t fwrite(const void * ptr, size_t size, size_t nmemb, FILE * stream)
3  size_t fread(void * ptr, size_t size, size_t nmemb, FILE * stream)
```

- FILE-streams can only be used with TCP sockets

```
1  sock = socket(/*...*/);
2  /* ... connect socket ...*/
3  // wrap the socket in an output stream
4  FILE *outstream = fdopen(sock, "w");
5  // send message, converting each object to network byte order before sending
6  if (fwrite(&val8, sizeof(val8), 1, outstream) != 1)
7  //do
```

# Wrapping TCP Sockets in Streams

- If we have to exchange data

```
1   struct addressInfo {
2       uint16_t streetAddress;
3       int16_t aptNumber;
4       uint32_t postalCode;
5   } addrInfo;
```

```
1   // ... put values in addrInfo ...
2   // convert to network byte order
3   addrInfo.streetAddress = htons(addrInfo.streetAddress);
4   addrInfo.aptNumber = htons(addrInfo.aptNumber);
5   addrInfo.postalCode = htonl(addrInfo.postalCode);
6   if (send(sock, &addrInfo, sizeof(addrInfo), 0) != sizeof(addrInfo)) ...
```

```
1   struct addressInfo addrInfo;
2   // ... sock is a connected socket descriptor ...
3   FILE *instream = fdopen(sock, "r");
4   if (fread(&addrInfo, sizeof(struct addressInfo), 1, instream) != 1) {
5   // ... handle error
6   }
7   // convert to host byte order
8   addrInfo.streetAddress = ntohs(addrInfo.streetAddress);
9   addrInfo.aptNumber = ntohs(addrInfo.aptNumber);
10  addrInfo.postalCode = ntohl(addrInfo.postalCode);
11  // use information from message...
```

# Wrapping TCP Sockets in Streams

- Using recv() will not help because there is no guarantee that all the bytes of the message will be returned in a single call.
  - We can use MSG_WAITALL

# Case

- Client sends a request message to the server:
  - inquiry request: how many votes are polled for a candidate.
  - Vote request: vote for a candidate.
- For both the requests server responds with
  - Response flag
  - Candidate id
  - Vote count



Vote Request
Candidate = 775

Vote Response
Candidate = 775
Vote Count = 21527

```
1   struct VoteInfo {
2   uint64_t count; // invariant: !isResponse => count==0
3   int candidate; // invariant: 0 <= candidate <= MAX_CANDIDATE
4   bool isInquiry;
5   bool isResponse;
6   };
```

# Encoding & decoding

- Application uses VoteInfo structure for its internal bookkeeping.
- Encoding: convert VoteInfo structure to wire format.
- Decoding: convert back to VoteInfo structure

```
1 ▾  /* Routines for Text encoding of vote messages.
2    * Wire Format:
3    * "Voting <v|i> [R] <candidate ID> <count>"
4    */
```

```
1 ▾  /* Routines for binary encoding of vote messages
2    * Wire Format:
3    *                              1 1 1 1 1 1
4    *  0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5
5    * +--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
6    * |     Magic        |Flags|        ZERO          |
7    * +--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
8    * |                  Candidate ID                 |
9    * +--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
10   * |                                               |
11   * |          Vote Count (only in response)        |
12   * |                                               |
13   * |                                               |
14   * +--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
15   *
16   */
```

# Encoding

```
1  /* Routines for Text encoding of vote messages.
2   * Wire Format:
3   * "Voting <v|i> [R] <candidate ID> <count>"
4   */
```

lead

```c
1   static const char *MAGIC = "Voting";
2   static const char *VOTESTR = "v";
3   static const char *INQSTR = "i";
4   static const char *RESPONSESTR = "R";
5   static const char *DELIMSTR = " ";
6   enum {
7     BASE = 10};
8   /* Encode voting message info as a text string.
9    * WARNING: Message will be silently truncated if buffer is too small!
10   * Invariants (e.g. 0 <= candidate <= 1000) not checked.
11   */
12  size_t Encode(const VoteInfo *v, uint8_t *outBuf, const size_t bufSize) {
13    uint8_t *bufPtr = outBuf;
14    long size = (size_t) bufSize;
15    int rv = snprintf((char *) bufPtr, size, "%s %c %s %d", MAGIC,
16        (v->isInquiry ? 'i' : 'v'), (v->isResponse ? "R" : ""), v->candidate);
17    bufPtr += rv;
18    size -= rv;
19    if (v->isResponse) {
20      rv = snprintf((char *) bufPtr, size, " %llu", v->count);
21      bufPtr += rv;
22    }
23    return (size_t) (bufPtr - outBuf);
24  }
```

```c
1   struct VoteInfo {
2     uint64_t count;
3     int candidate;
4     bool isInquiry;
5     bool isResponse;
6   };
```

# Decoding

```
1 ▾  /* Routines for Text encoding of vote messages.
2      * Wire Format:
3      * "Voting <v|i> [R] <candidate ID> <count>"
4      */
```

```
1 ▾ bool Decode(uint8_t *inBuf, const size_t mSize, VoteInfo *v) {
2     char *token;
3     token = strtok((char *) inBuf, DELIMSTR);
4     // Check for magic
5     if (token == NULL || strcmp(token, MAGIC) != 0)
6       return false;
7     // Get vote/inquiry indicator
8     token = strtok(NULL, DELIMSTR);
9     if (token == NULL)
10      return false;
11    if (strcmp(token, VOTESTR) == 0)
12      v->isInquiry = false;
13    else if (strcmp(token, INQSTR) == 0)
14      v->isInquiry = true;
15    else return false;
16    // Next token is either Response flag or candidate ID
17    token = strtok(NULL, DELIMSTR);
18    if (token == NULL)     return false; // Message too short
19 ▾  if (strcmp(token, RESPONSESTR) == 0) { // Response flag present
20      v->isResponse = true;
21      token = strtok(NULL, DELIMSTR); // Get candidate ID
22      if (token == NULL)
23        return false;
24 ▾  } else { // No response flag; token is candidate ID;
25      v->isResponse = false;  }
```

# Decoding

```c
1  /* Routines for Text encoding of vote messages.
2   * Wire Format:
3   * "Voting <v|i> [R] <candidate ID> <count>"
4   */
```

```c
26     // Get candidate #
27   v->candidate = atoi(token);
28   if (v->isResponse) {
29   // Response message hould contain a count value
30     token = strtok(NULL, DELIMSTR);
31     if (token == NULL)
32       return false;
33     v->count = strtoll(token, NULL, BASE);
34   } else {
35     v->count = 0L;
36   }
37   return true;
38 }
```

# Encoding a Binary Wire Format

- Text based formats vary in length.
- Binary formats always have fixed size.
- Magic value 010101 helps in ensuring that we are receiving the right message.

```
1  ▾  /* Routines for binary encoding of vote messages
2      * Wire Format:
3      *                                  1  1  1  1  1  1
4      *  0  1  2  3  4  5  6  7  8  9  0  1  2  3  4  5
5      * +--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
6      * |      Magic       |Flags|         ZERO         |
7      * +--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
8      * |                  Candidate ID                 |
9      * +--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
10     * |                                               |
11     * |          Vote Count (only in response)        |
12     * |                                               |
13     * |                                               |
14     * +--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
15     *
16     */
```

# Encoding a Binary Wire Format

- Generally declare an equivalent structure for the wire format.

```
1    enum {
2        REQUEST_SIZE = 4,
3        RESPONSE_SIZE = 12,
4        COUNT_SHIFT = 32,
5        INQUIRE_FLAG = 0x0100,
6        RESPONSE_FLAG = 0x0200,
7        MAGIC = 0x5400,
8        MAGIC_MASK = 0xfc00
9    };
10
11   typedef struct voteMsgBin voteMsgBin;
12
13   struct voteMsgBin {
14       uint16_t header;
15       uint16_t candidateID;
16       uint32_t countHigh;
17       uint32_t countLow;
18   };
```

# Encoding a Binary Wire Format

```
20   size_t Encode(VoteInfo *v, uint8_t *outBuf, size_t bufSize) {
21     if ((v->isResponse && bufSize < sizeof(voteMsgBin)) || bufSize < 2
22         * sizeof(uint16_t))
23       DieWithUserMessage("Output buffer too small", "");
24     voteMsgBin *vm = (voteMsgBin *) outBuf;
25     memset(outBuf, 0, sizeof(voteMsgBin)); // Be sure
26     vm->header = MAGIC;
27     if (v->isInquiry)
28       vm->header |= INQUIRE_FLAG;
29     if (v->isResponse)
30       vm->header |= RESPONSE_FLAG;
31     vm->header = htons(vm->header); // Byte order
32     vm->candidateID = htons(v->candidate); // Know it will fit, by invariants
33     if (v->isResponse) {
34       vm->countHigh = htonl(v->count >> COUNT_SHIFT);
35       vm->countLow = htonl((uint32_t) v->count);
36       return RESPONSE_SIZE;
37     } else {
38       return REQUEST_SIZE;
39     }
40   }
```

# Decoding a Binary Wire Format

```c
1    /* Extract message info from given buffer.
2     * Leave input unchanged.
3     */
4    bool Decode(uint8_t *inBuf, size_t mSize, VoteInfo *v) {
5      voteMsgBin *vm = (voteMsgBin *) inBuf;
6      // Attend to byte order; leave input unchanged
7      uint16_t header = ntohs(vm->header);
8      if ((mSize < REQUEST_SIZE) || ((header & MAGIC_MASK) != MAGIC))
9        return false;
10     /* message is big enough and includes correct magic number */
11     v->isResponse = ((header & RESPONSE_FLAG) != 0);
12     v->isInquiry = ((header & INQUIRE_FLAG) != 0);
13     v->candidate = ntohs(vm->candidateID);
14     if (v->isResponse && mSize >= RESPONSE_SIZE) {
15       v->count = ((uint64_t) ntohl(vm->countHigh) << COUNT_SHIFT)
16               | (uint64_t) ntohl(vm->countLow);
17     }
18     return true;
19   }
```

# Framing

- Framing refers to the general problem of enabling the receiver to locate the boundaries of a message.
    - the application protocol must specify how the receiver of a message can determine when it has received all of the message.
        - Whether information is encoded as text, as multibyte binary numbers, or as some combination of the two
    - This is trivial in UDP.
    - But TCP doesn't preserve the boundaries.
        - If the message length is fixed, then we can wait until we read all of the bytes. But if the length is varying then?

- Two techniques:
    - Delimiter based: The end of the message is indicated by a unique marker.
    - Explicit length: The variable-length field or message is preceded by a length field that tells how many bytes it contains.

# Framing in HTTP

# HTTP Request

- HTTP is a request-response stateless protocol. HTTP/1.1 supports persistent connections.

- HTTP Request:
  o GET request doesn't include anything in the body.

```
1   GET /about.html HTTP/1.1
2   Host: www.bits-pilani.ac.in      //must in HTTP 1.1
3   Connection: Keep-Alive
4   User-Agent: Mozilla/4.06 [en] (X11; U; Linux 2.1.121 i686)
5   Accept: image/gif, image/x-xbitmap, image/jpeg, image/pjpeg, image/png
6   Accept-Encoding: gzip
7   Accept-Language: en
8   Accept-Charset: iso-8859-1,utf-8
9       <blank line>
```

  o Each line is terminated by \r\n.
  o Headers are terminated by an \n.

- Framing: a blank line

# HTTP Request

- POST request includes content in the body.

  o Length of the content is specified in the header.

```
1   POST /about.html HTTP/1.1
2   Host: www.bits-pilani.ac.in      //must in HTTP 1.1
3   Connection: Keep-Alive
4   User-Agent: Mozilla/4.06 [en] (X11; U; Linux 2.1.121 i686)
5   Accept: image/gif, image/x-xbitmap, image/jpeg, image/pjpeg, image/png
6   Accept-Encoding: gzip
7   Accept-Language: en
8   Accept-Charset: iso-8859-1,utf-8
9   Content Length:35
10
11  idno=2007A1PS001&item=test1&name=Krishna
```

  o Each line is terminated by \r\n.

  o Headers are terminated by an \n.

- Framing: a blank line is end of headers.

  o Body's boundary is identified by explicit-length method
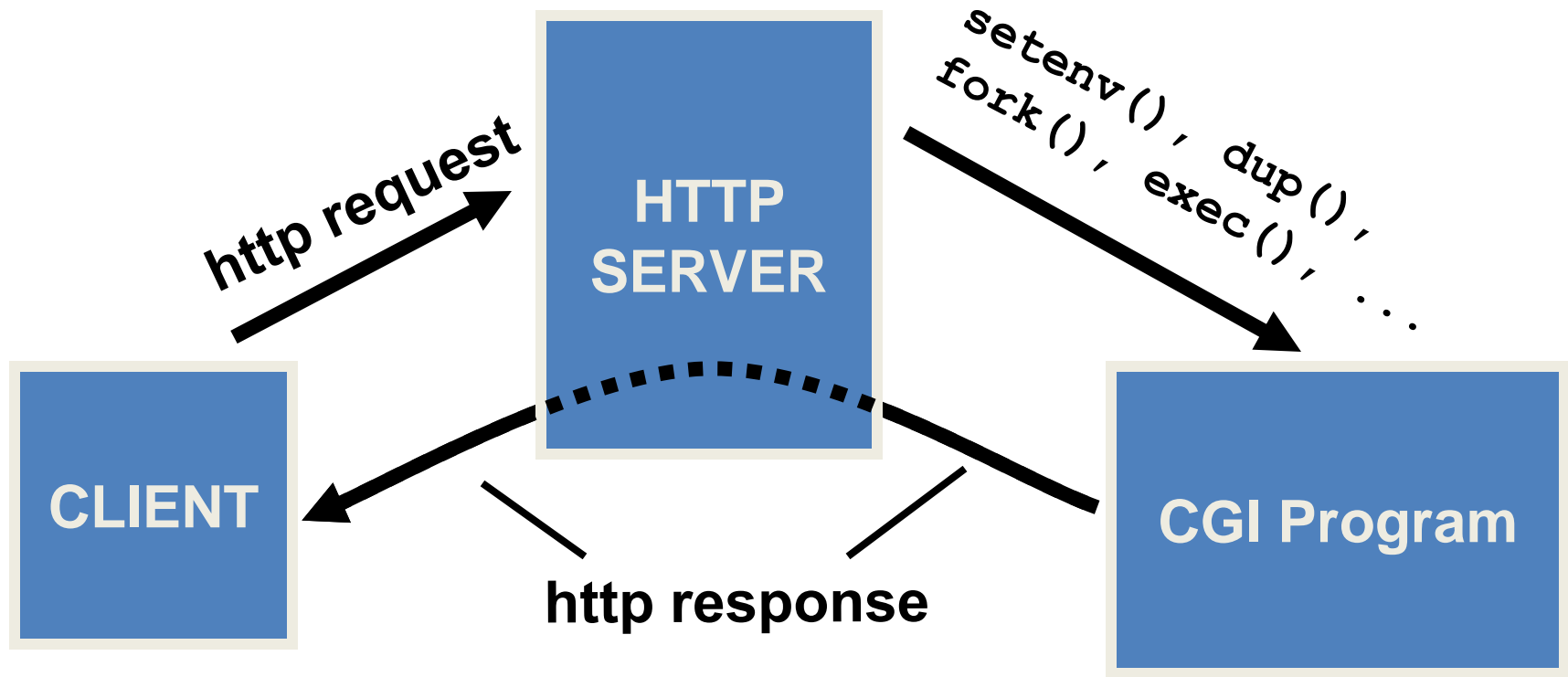
# HTTP Response

- Response header contains content-length field.

```
1   HTTP/1.1 200 OK
2   Server: Microsoft-IIS/5.0
3   Date: Fri, 08 Oct 2010 05:08:14 GMT
4   Connection: open
5   Content-length: 175
6   Content-Type: text/html
7
8   content is very big. ...................
```

- In dynamic output cases

```
1   HTTP/1.1 200 OK
2   Server: Microsoft-IIS/5.0
3   Date: Fri, 08 Oct 2010 05:08:14 GMT
4   Connection: close
5   Content-Type: text/html
```

# CGI Programming

# HTTP & Dynamic Outputs

- It may  not be convenient or even possible for  a  server to know  the length of  an item before sending.
- Servers use the Common Gateway Interface  (CGI) mechanism to create dynamic documents.
- To provide for dynamic Web pages, the HTTP standard specifies that if  the server does not know the length of  an item a priori,  the server can inform the browser that it will close the connection after transmitting the item

```
1   HTTP/1.1 200 OK
2   Server: Microsoft-IIS/5.0
3   Date: Fri, 08 Oct 2010 05:08:14 GMT
4   Connection: close
5   Content-Type: text/html
```

# HTTP & Dynamic Outputs

- Closing connection results in poor performance.

- If server doesn't know the output length a priori, it can also use chunked transfer encoding.

- The sender breaks the message body into chunks of arbitrary length, and each chunk is sent with its length prepended;

- It marks the end of the message with a zero-length chunk.

- The sender uses the *Transfer-Encoding: chunked* header to signal the use of chunking.

- This mechanism allows the sender to buffer small pieces of the message, instead of the entire message, without adding much complexity or overhead.

# Acknowledgements

# Q&A

# Thank You