# Network Programming with Sockets

http://ece297.msrg.utoronto.ca

I added some additional slides that
we wont cover in the lectures; they
are for your information.

Stats :

ῆ 350 submissions (should have been ~380)

ῆ    9 with wrong directory structure (did not

        follow submission instructions)

# Quiz: Week of Feb. 10[th] : Dedicated lab: TBD

- ñ Optiona l & pass/fail

- ñ Entice you to use svn & testsubmit

- ñ Consistent use of svn over a number of days ( start now !)

- ñ See instructions online

- ñ One student on behalf of the whole team

# Midterm: Get an early start!

<span style="color:red">ῆ Read the online instructions now !</span>

ῆ Full lecture dedicated to midterm: Feb 12th

ῆ Evaluates Milestones 1 & 2

   ã Short team presentation

   ã Question-based demo (   by individual      !)

   ã General questions (see online samples)

# Important points for Midterm

- Each student must run demo from his or her own account; <span style="color:red">running the demo from someone else̦s account is not allowed</span>

- It is your responsibility that your account is properly working, i.e., is not revoked, sufficient disk space is available, etc.

- Should there be a problem with the computer you are setting up on, the supervising TA will allow you to try one other computer

- You will have some setup time to verify that all is working, but you will be asked <span style="color:red">to start your demos from scratch</span> (<span style="color:red">note the data loading requirement</span>)!

- Your shells should <span style="color:red">provide meaningful output</span> for erroneous cases, e.g., retrieving tables that do not exist, keys that do not exist, etc.; <span style="color:red">a segmentation fault is not a meaningful output</span>
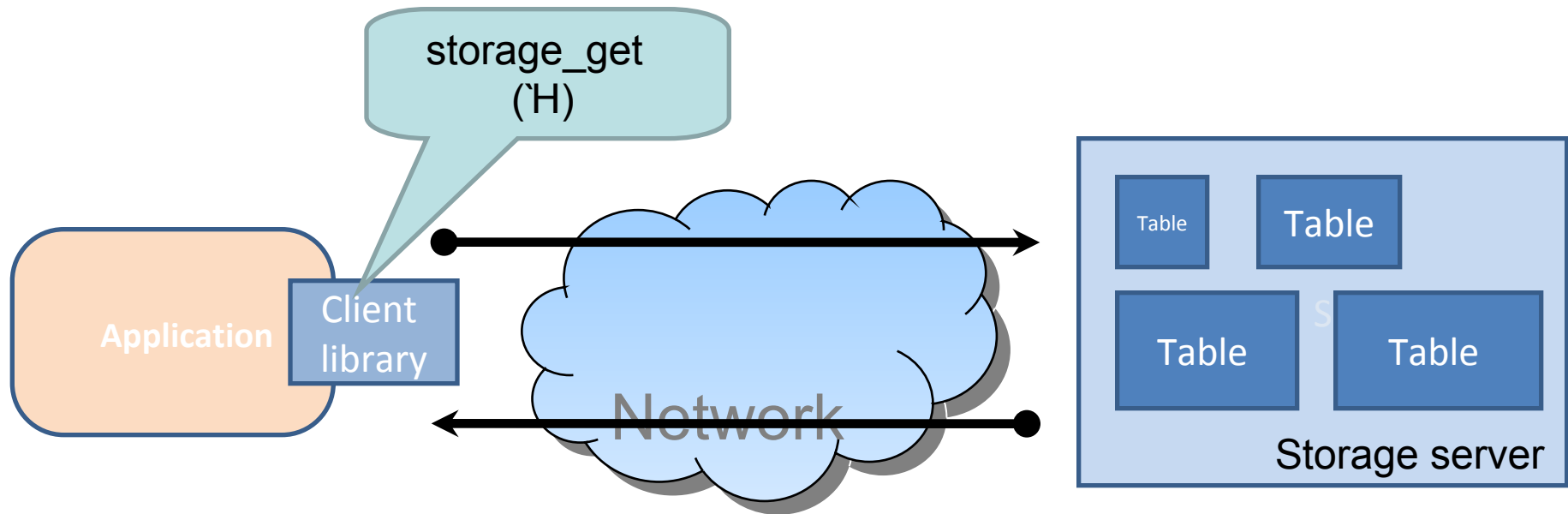
# Agenda

η Network programming

    α Sockets  *et al.*

    α Protocol design

    α Socket API

# Network programming

# How does storage_get(Ẑ) propagate over the network?



storage_get (ˋH)

Application

Client library

Network

Table
Table
Table
Table

Storage server

ECE 297

# On the client side

See client.c
    status =    storage_get    (TABLE, KEY, &r, conn);


See storage.c: in    storage_get
  ...  sendall   (sock, buf, strlen(buf)) `H && `H  recvline(sock, buf,
    sizeof buf) ...

Our string array / buffer

See utils.c: in    sendall
    `H
    while (tosend > 0) {
        ssize_t bytes =    send  (sock, buf, tosend, 0);
        if (bytes <= 0)
                break; // send() was not successful, so stop.
        tosend -= bytes;
        buf += bytes;
    };

ECE 297

# The string buffer

storage_get('MyCourses~, 'ECE344~, &record, conn)

`H GET MyCourses ECE344 `H

Network

# On the server side

See server.c in **main**

`` `H ``

while (wait_for_connections) { ...

  // Get commands from client.

| Our string array / buffer |
| --- |

  do { `` `H ``

    // Read a line from the client.

    int status = **recvline** (clientsock, cmd, MAX_CMD_LEN);

    `` `H ``

    // Handle the command from the client.

    int status = **handle_command** (clientsock, cmd);

# On the server side

```
int  handle_command      (int sock, char *cmd){ `H
    LOG(("Processing command '%s'\n", cmd));


    // For now, just send back the command to the
    client.
    sendall   (sock, cmd, strlen(cmd));
    `H
}
```
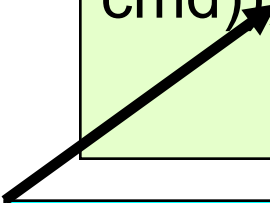
ECE 297

# Using the 'network˜

Application
(communication end-points)

Networking API

Protocols (TCP, `H)

Network

Application
(communication end-points)

Networking API

Protocols (TCP, `H)

**Sockets**  are network programming abstractions.

# Local call vs. sending data

```
int handle_command(int s, char *cmd){
        LOG (("Processing command '%s'\n",
cmd));
        sendall   (sock, cmd, strlen(cmd) );
```

Local call:    **Function call**      in the same address space (process)

Network

```
int sendall(int s, char *buf, size_t len){
        ssize_t tosend = len;
        while (tosend > 0) {
                size_t bytes =      send  (sock, buf, tosend, 0);
```

Sending data: A    **system call**      that initiates    **transfer of data across the net**

# Communication end-points & sockets

- ñ Identify the *application* that sends and receives data and the *host* the application runs on
  - α Multiple applications running on the same machine
- ñ End-points are programmatically represented by sockets ( #include <sys/socket.h> )
- ñ Create the end-points
- ñ Connect the end-points
- ñ Listen and accept connections on end-points

ECE 297

- ñ Transmit messages via sending and

# Sockets are file descriptors

ǐ Socket is a **programming abstraction** for communicating among processes (applications) based on (Unix) *file descriptors*

ǐ A **file descriptor** is an integer representing an open file managed by the OS

  ã In Unix any I/O is done by reading/writing from/to file descriptors

  ã For sockets, also called **socket descriptor**

# Our focus in ECE 297

ñ Our focus is **Internet sockets**
   α there are others

ñ Our focus is sockets of type **stream**, i.e., **SOCK_STREAM** (i.e., based on TCP)
   α there are other types

ñ Our focus is **IPv4**

ñ All this is provided for you in the skeleton code.
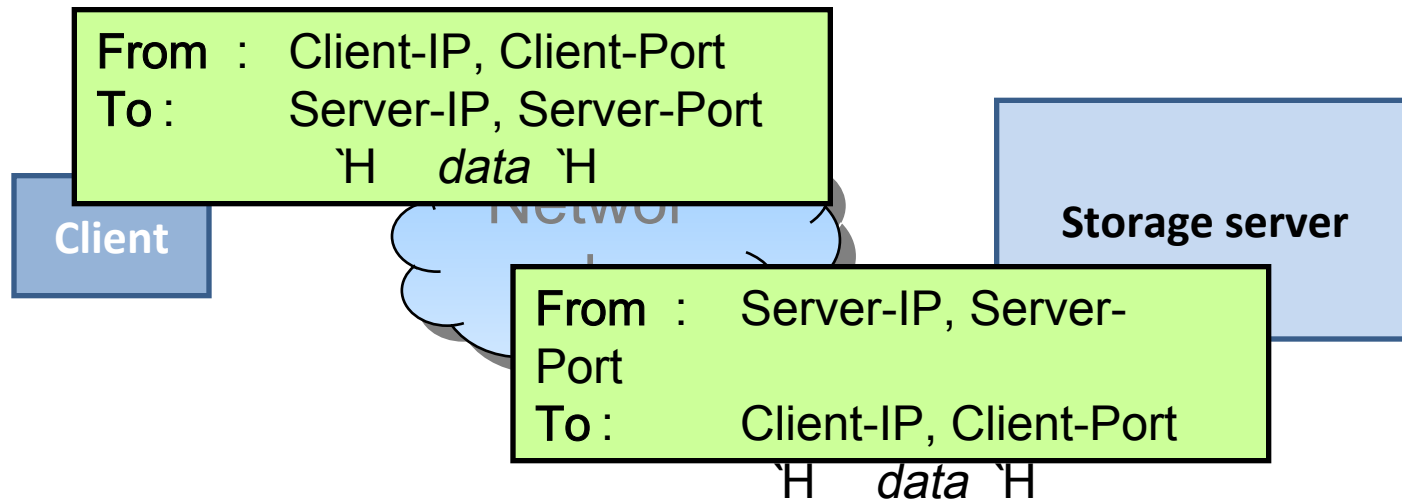
# SOCK_STREAM Sockets

- Are **reliable**, **two-way connected** communication streams

  - Output two items into the socket in the order "1, 2", they arrive in the order "1, 2" at the opposite end (*'`H if they arrive at all* ☹)

  - Transmission is **error-free** (i.e., messages are not corrupted)

# Two end-points determine a connection

ɧ End-point is determined by

   ã **IP address**     (host address); e.g., 192.168.100.100

   ã **Port number**    , e.g., 8888

| | | |
|---|---|---|
| From : | Client-IP, Client-Port | |
| To : | Server-IP, Server-Port | |
| | `H *data* `H | |

**Client**

Netwo~~~

**Storage server**

| | |
|---|---|
| From : | Server-IP, Server-Port |
| To : | Client-IP, Client-Port |

`H *data* `H

Try `netstat` to see open sockets on your machine

ECE 297

man netstat for help

# The configuration file & server IP, port

- ñ The configuration file sets *properties* for the storage serv
- ñ It is read at startup to initialize the storage server.
- ñ It is **not available to client** (i.e., it resides across network

File: `default.conf` (given)

Server IP

```
server_host localhost
server_port 1111
table marks
```

Port the server listens on

Most software has
configuration
files
and ways to
overwrite
their settings.

For a further example see onetable.
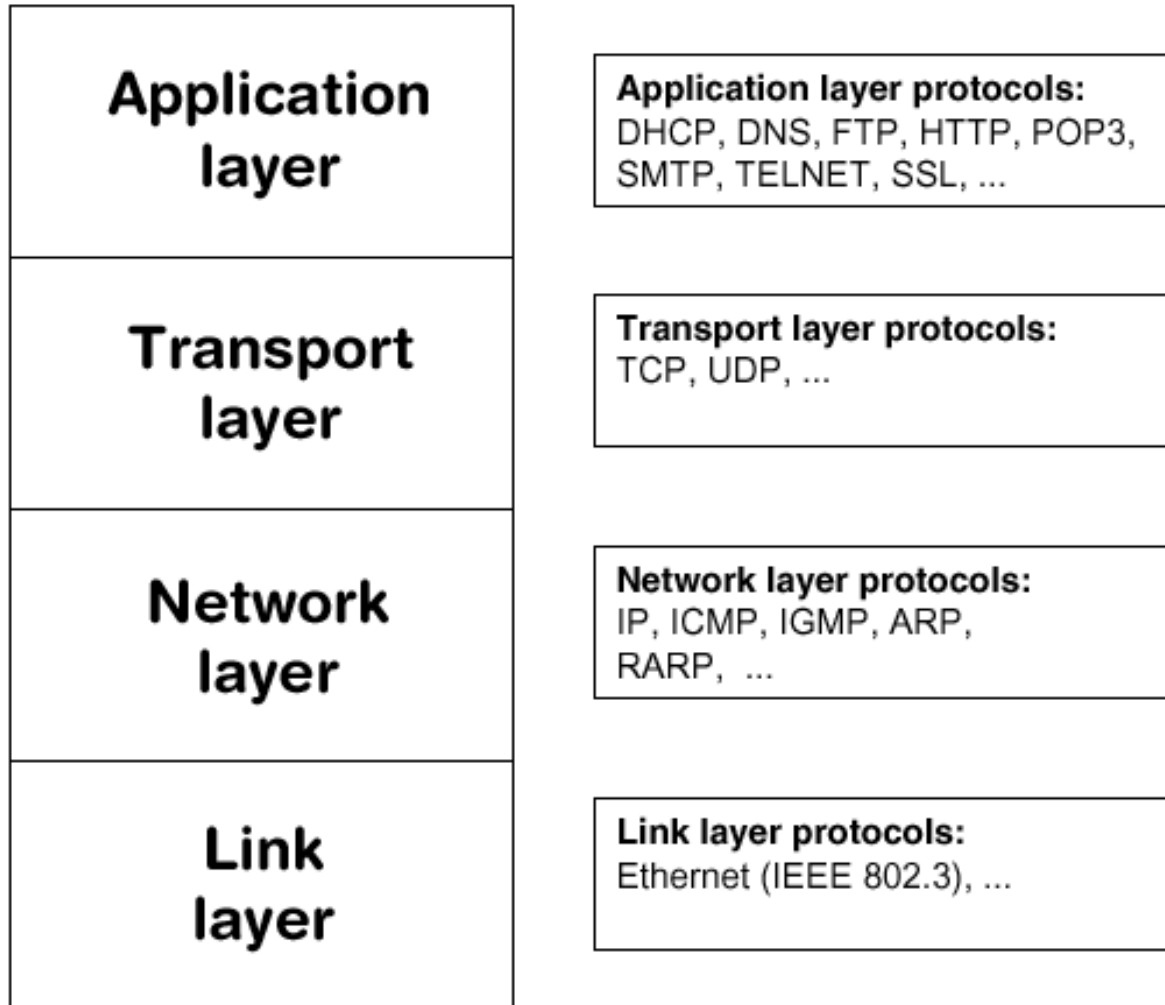conf in the skeleton code distribution

19

# IP address
## (Only what we need to know at this point)

- ñ Internet Protocol (IP) address of a host (32bit number)
- ñ Split up into four eight-bit numbers
- ñ At home you may see 192.168.0.0 `H
  - ã Reserved block of IP address for private networks
  - ã Your modem / router is assigned a **public IP** by Bell/Rogers
- ñ For development we can use **localhost**
  - ã Meaning **this computer** / host
  - ã Translated by OS to loopback IP address **127.0.0.1**

- ñ **Domain Name System** (DNS ) maps names to IP addresses

ECE 297

  - ã It is easier to remember www.example.com than
  192.168.212.4

# The network stack

| | |
|---|---|
| **Application layer** | **Application layer protocols:** DHCP, DNS, FTP, HTTP, POP3, SMTP, TELNET, SSL, ... |
| **Transport layer** | **Transport layer protocols:** TCP, UDP, ... |
| **Network layer** | **Network layer protocols:** IP, ICMP, IGMP, ARP, RARP, ... |
| **Link layer** | **Link layer protocols:** Ethernet (IEEE 802.3), ... |

# Ports

- Example: Computer with a given IP address handles Mail, IM, Torent, and Web browsing
  - *How should the OS differentiate among these networked applications, if packets arrive?*
- A port is a **16 bit number** to locally (per host) **identify** the **connection**
- Numbers (vary by OS)
  - 0-1023 'reserved˜, must be *root* to use
  - **1024 – 65535 are available to regular user**
- Well-known, reserved services
  - http **80/tcp**
  - ftp 21/tcp
  - ssh 22/tcp
  - Others: telnet 23/tcp; finger 79/tcp; snmp 161/udp
- See `/etc/services` in Unix for all available services

# 'Error binding socket"

ñ Port is used by another process

```
–ps aux | grep server
```

ñ Kill it, if it is your own

ñ Use a different port number, if it is
   someone elseĄs who might be logged
   in to the machine via ssh

# Byte order
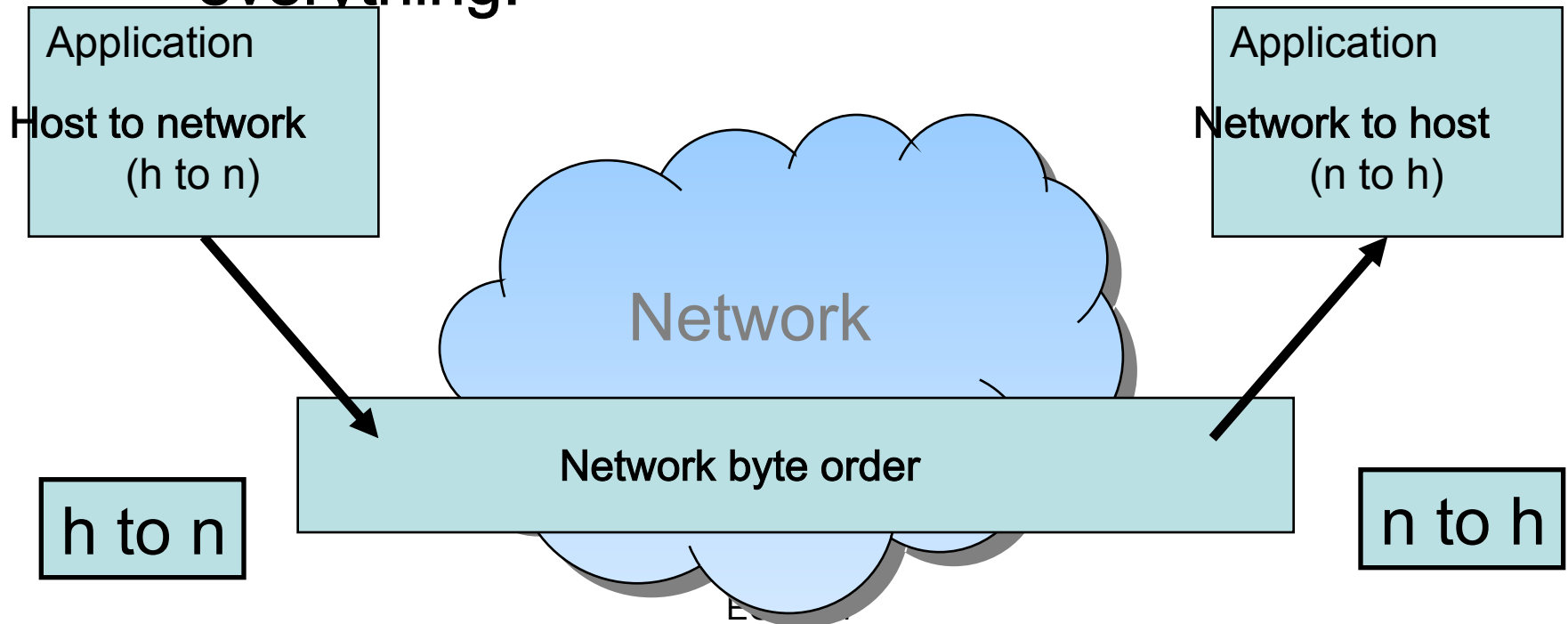
- The order bytes are represented by a computer architecture in memory, on the wire, `H

- E.g., 0x *A3 F1* as two bytes // 0x HEX number
  - Stored as A3 F1 // Big-Endian (' *Big end first* ˜)
    - Motorola 6800, 68k, PowerPC (Intel), `H
    - Network byte order
  - Stored as F1 A3 // Little-Endian (' *Little end first* ˜)
    - Intel or Intel-compatible processors
    - x86, 6502, Z80, VAX, PDP-11
- Host byte order : the order used by the host vs. the network byte order , the order used by the network

- Important to honor network byte order when building messages and filling out structures!

# How do we know the host byte order?

Example ( wrong ): (Motorola) 0xA3F1 `H (Network) 0xA3F1 `H    (Intel) 0xF1A3

        41969                                               61859

1010 0011   1111 0001           1111 0001   1010 0011

But, donẠt worry;    just transform everything!

| | | |
|---|---|---|
| Application | | Application |
| Host to network (h to n) | Network | Network to host (n to h) |
| | Network byte order | |
| h to n | | n to h |

# h-to-n    and    n-to-h

htons( `H)   host to network short

htonl( `H)   host to network long

ntohs( `H)   network to host short

ntohl( `H)   network to host long

*Great, but how do we handle floats?*

Left as an exercise for the reader.

# Protocol design
*What̃s relevant for us for the assignments?*

ñ Simplicity
- ã Simple encoding of *client requests*
- ã Simple encoding of *server response*
- ã Need to worry about *error conditions*
- ã Don̞t worry about saving a bit of bandwidth or computation
- ã Consider debugging your communication

ñ Interoperability
- ã We **can not assume** client and server run on the same hardware / OS

ñ Portability
- ã We restrict ourselves to IPv4

# The Protocol

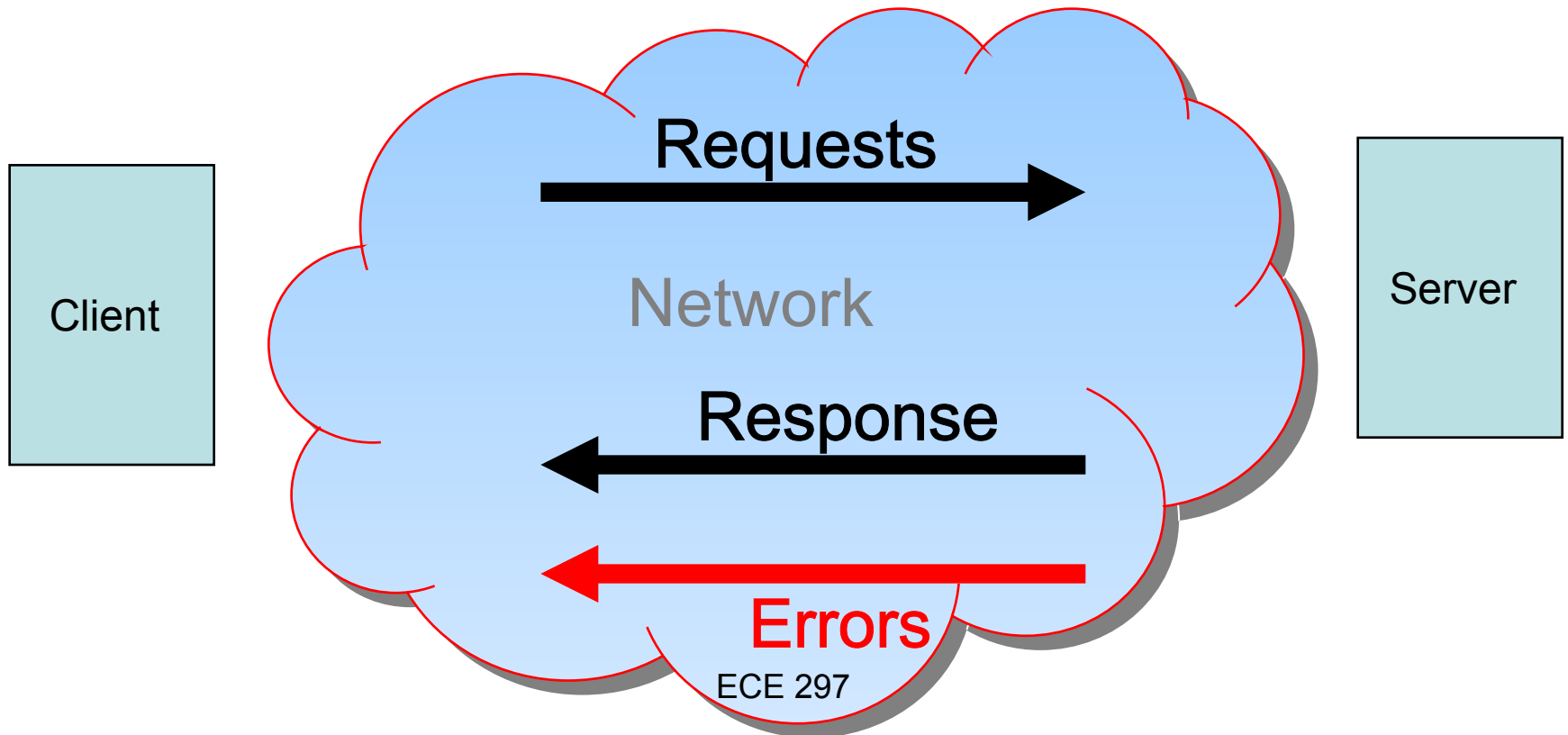**Example**

Request:    GET <path>/index.html HTTP/1.0

Response:            HTTP/1.0 200 OK

Messages

Error response:    HTTP/1.0 404 Not Found

Requests →

Client

Network

Server

← Response

← Errors

# Setting up & using sockets
(given in the skeleton code)

**Server**

`socket()`

`bind()`

`listen`

`accept`

*(Block) until connection*

**Client**

`socket`

*WHandshakew*

`connec`

Data (request)

`send()`

`recv()`

Data (reply)

`send()`

`recv()`

End-of-File

`recv()`

`close(`

`close(`

`)`

# Sending data

```
int send( int sockfd,
          const void *msg,
          int len, int flags);
```

-1 is returned on error, and `errno` is set to the error number.

- ñ **msg** is a pointer to the data to send
- ñ **len** is the length of the data to send
- ñ **flag** controls specifics, set to 0 for us
- ñ send(`H) returns the number of bytes sent, which might be less than what was requested

```
char *msg = '   ECE297 is great fun, we learn useful stuff
int len, bytes_sent;
. . .                    serĄs responsibility to send the rest            !
len = strlen(msg);                                                          !
bytes_sent     = send(sockfd, msg, len, 0);
```

# On the client side

See utils.c: in  **sendall**

  `H

  while (tosend > 0) {

        ssize_t bytes =    **send**  (sock, buf, tosend, 0);

        if (bytes <= 0)

                    break; // send() was not successful, so stop.

        tosend -= bytes;

        buf += bytes;

  };

# ADDITIONAL MATERIAL

# Protocol design considerations

- How are API calls (e.g., get / set) sent from the client to the server?
  - *May they have given us any hints in the code?*
- What parts of a call need to be conveyed to the server?
  - E.g., int storage_get(*table, *key, *record, *conn);
- What's the simplest possible way of sending this information?
- What else needs to be represented in the protocol between the client library and the storage server?

# Associate socket with port on host

```
int bind( int sockfd,
          struct sockaddr *myaddr,
          socklen_t addrlen);
```

Assigns a socket an address.

- ñ **sockfd** is socket descriptor from `socket(...)`
- ñ **myaddr** is a pointer to address struct with:
  - ã port number and IP address
  - ã if port is 0, then host picks port (>1023)
  - ã IP address = INADDR_ANY
- ñ **addrlen** is length of structure
- ñ returns 0 on success, -1 on failure and sets errno
  - ã Common is EADDRINUSE (' *Address already in use* ˜)

Especially when rerunning the server in development

# Client connecting to server

```
int connect(    int sockfd,
                struct sockaddr *serv_addr,
                int addrlen);
```
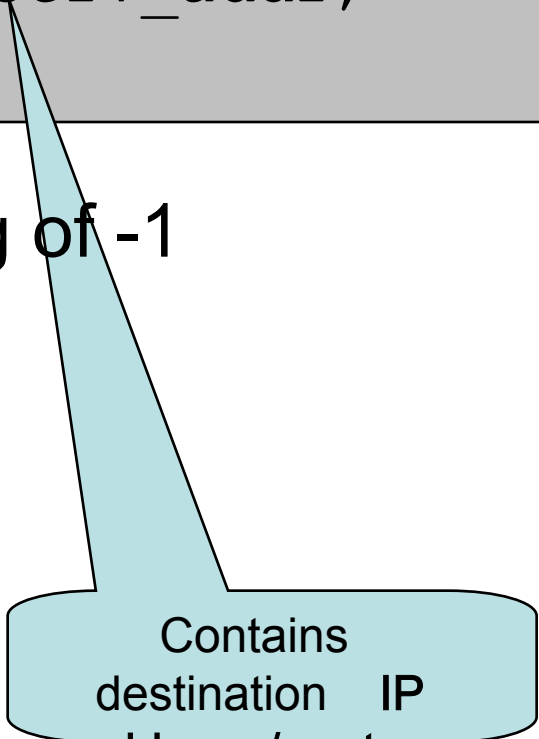
- ñ Failures indicated by the returning of -1 and the setting of    errno
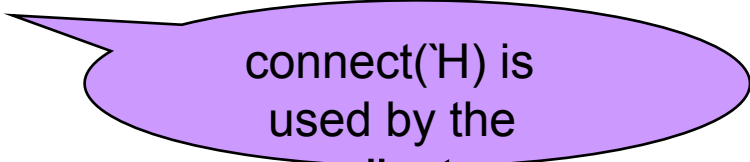- ñ E.g., bind(`H) was not called!
- ñ Setting local port is not important
- ñ The destination IP/port matters
- ñ OS selects local port & conveys to remote site

Contains destination   IP address / port

connect(`H) is used by the client

# Accepting connections

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

- ñ Return value & arguments are similar to the above (cf. bind(`H))
- ñ Returns brand new descriptor      created by OS
- ñ Former descriptor continues to listen for new connections
- ñ New descriptor is used to send and receive data

accept(`H) is used
by the   server

# Sending data

```
int send( int sockfd,
          const void *msg,
          int len, int flags);
```

ñ **msg**   is a pointer to the data to send

ñ **len**   is the length of the data to send

ñ **flag**   controls specifics, set to 0 for us

ñ send(`H) returns the number of bytes sent, which might be less than what was requested

-1 is returned on error, and `errno` is set to the error number.

```
char *msg = '     ECE297 is great fun, we learn useful stuff
int len, bytes_sent;
. . .                     serÃs responsibility to send the rest
len = strlen(msg);
bytes_sent     = send(sockfd, msg, len, 0);
```

!

# Receiving data & close

```
int recv( int sockfd,
          void *buf,
          int len, int flags);
```

- ñ  **buf**   is a buffer to read data into
- ñ  **len**   is the maximum length of the buffer
- ñ  **flag**   controls specifics, set to 0 for us
- ñ  recv(`H) returns the number of bytes read
- ñ  recv(`H) may return 0;      remoteˌs side  closed connection on us

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

- ñ  attempts to send any unsent data
- ñ  closes socket for sending and receiving
- ñ  returns -1 if error

-1 is returned on error, and `errno` is set to the error number.

# Read the code !

ñ A lot of what we ask you to do is already in the code

  ã or at least hinted in the code and the comments

ñ Scared of code

  ã look at the doxygen output and then read the code

ñ Read the handout

  ã a lot of what we ask you to do is explained in English

# ADDITIONAL MATERIAL

# Actually `send(...)` is a system call

- Program control transitions to the OS
- Switch from user space to kernel space
- Implementation requires 100s of assembly instructions
- Much more expensive than a pure **function call**, like `sendall(...)`
- A function call is a few assembly instructions
- **Latency cost** : Remote call > > Function call

ECE 297

For far more details, see ECE344: Operating Systems

# The *ẄLatency is Zerȯw* Fallacy

ñ In the context of remote calls a fatal misconception

ñ Latency is the time it takes to move data from client to storage server vs amount of data transferred (bandwidth)

ñ Typical WAN round trip times are 30ms+

ñ *Transfer as much data with as few calls as possible vs. many*

# "The Network is Reliable"
## ' Fallacies of Distributed Computing Explained~

- ñ Consider: Power supply failures, disk failures, CPU / power fan failures, A/C failures, accidents (cables get disconnected etc.), flash crowds, bugs, `H

- ñ Dependency on third party services (e.g., credit card validation, ad serving, `H)

- ñ Requires hardware and software redundancy (cost-benefit trade-off )

- ñ Loose messages (acknowledgements, time-outs & re-tries), detect duplicates (idempotent operations), re-order messages (donẠt depend on order), `H

ECE 297

# Stream sockets in practice

ñ Telnet uses stream sockets

ñ Web browsers & HTTP use stream sockets to fetch pages

ñ Stream sockets are based on the Transmission Control Protocol (TCP)

ñ TCP ensures that data arrives sequentially and error-free (given no failures)

ñ *When is this kind of reliability not required?*

ECE 297

# /etc/services

```
# This file contains port numbers for well-known services defined by
IANA
# Format:
# <service name>  <port number>/<protocol>  [aliases...]
[#<comment>]
echo              7/tcp
echo              7/udp
discard           9/tcp     sink null
discard           9/udp     sink null
systat            11/tcp    users                          #Active users
systat            11/udp    users                          #Active users
daytime           13/tcp
daytime           13/udp
qotd              17/tcp    quote                          #Quote of the day
qotd              17/udp    quote                          #Quote of the day
chargen           19/tcp    ttytst source                  #Character
generator
chargen           19/udp    ttytst source                  #Character
generator
ftp-data          20/tcp                    #FTP, data
ftp               21/tcp          ECE 297  #FTP. control
ssh               22/tcp                    #SSH Remote Login Protocol
telnet            23/tcp
```

ECE 297

# Client/Server

- Client
  - Issues requests to server (e.g., send & receive)
- Server
  - Starts up and listens for connections, requests, and sends/receives
- Client/Server examples
  - telnet/telnetd
  - ftp/ftpd (sftp/sftpd)
  - Firefox/Apache
- Client and server are **roles**, our following discussion also applies in **peer-to-peer** contexts
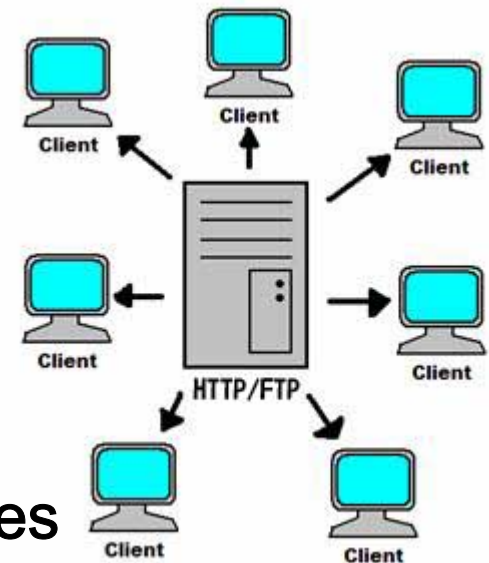- Same application may play **both roles**

Client/Server



HTTP/FTP

Client
Client
Client
Client
Client
Client
Client

Image adopted from http://simonlin.info/

ECE 297

# Creating a socket

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

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

Creates a socket Descriptor.

- family
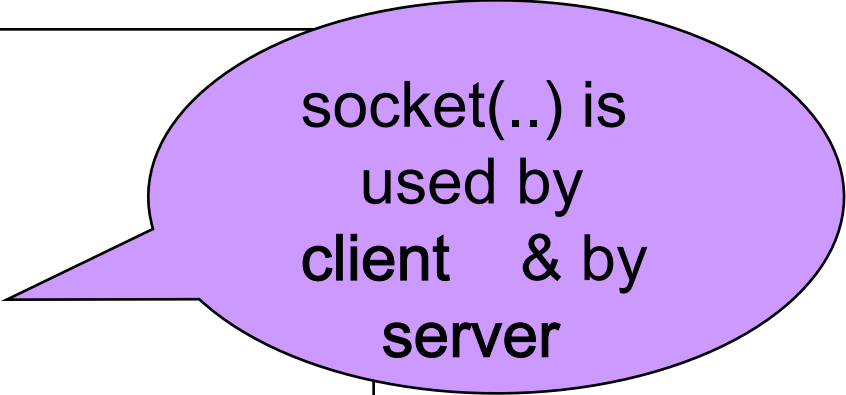  - AF_INET (IPv4),     AF_INET6 (IPv6), AF_LOCAL (local Unix), AF_ROUTE, AF_KEY
- type
  - SOCK_STREAM (TCP)    , SOCK_DGRAM (UDP), SOCK_RAW
- protocol    set to 0 (retrieve the default protocol for the given family and type); other values see <netinet/in.h>
- Upon success returns socket descriptor, otherwise  -1  indicating failure    and  sets errno

# socket(`H) often used as follows

```
int s;
struct addrinfo *    res ;

// do the lookup
getaddrinfo("www.example.com",
            "http",
            &hints,
            & res );


s = socket(res->ai_family,
            res->ai_socktype,
            res->ai_protocol);
```

socket(..) is used by client & by server

# Associate socket with port on host

```
int bind( int sockfd,
          struct sockaddr *myaddr,
          socklen_t addrlen);
```

Assigns a socket an address.

- ñ **sockfd** is socket descriptor from `socket(...)`
- ñ **myaddr** is a pointer to address struct with:
  - ã port number and IP address
  - ã if port is 0, then host picks port (>1023)
  - ã IP address = INADDR_ANY

Especially when rerunning the server in development

- ñ **addrlen** is length of structure
- ñ returns 0 on success, -1 on failure and sets errno
  - ã Common is EADDRINUSE ('*Address already in use* ˜)

# Example

struct addrinfo hints, *res;
int sockfd;

memset(&hints, 0, sizeof hints);
hints.ai_family = AF_UNSPEC;
hints.ai_socktype = SOCK_STREAM;
hints.ai_flags =    AI_PASSIVE   ;

getaddrinfo(NULL, "    3490  ", &hints, &res);

sockfd = socket(res->ai_family, res->ai_socktype, res->ai_protocol);

bind (sockfd, res->ai_addr, res->ai_addrlen);

> Load up address structs with getaddrinfo():

> Fill in my IP for me

> Bind socket to port we passed to `getaddrinfo`

> bind(`H) is used by the **server**

> We'd better do error checking
>
> **error = bind(`H)**

# Client connecting to server

```
int connect(    int sockfd,
                struct sockaddr *serv_addr,
                int addrlen);
```

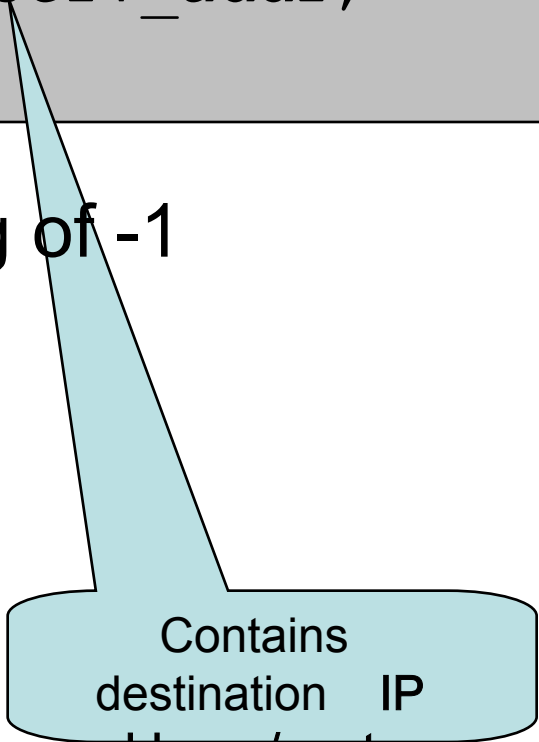- ñ Failures indicated by the returning of -1 and the setting of    errno
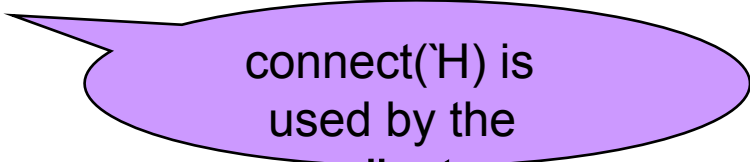- ñ E.g., bind(`H) was not called!
- ñ Setting local port is not important
- ñ The destination IP/port matters
- ñ OS selects local port & conveys to remote site

Contains destination    IP address / port

connect(`H) is used by the client

# Example

struct addrinfo hints, *res; int sockfd;

memset(&hints, 0, sizeof hints);
hints.ai_family = AF_UNSPEC;
hints.ai_socktype = SOCK_STREAM;
getaddrinfo("www.ex ... .com", "3490", &hints,
   &res);

sockfd = socket(          res->ai_family, res-
   >ai_socktype,
                    res->ai_protocol);

connect   (sockfd, res->ai_addr, res-
   >ai_addrlen);

# listen(…)

int listen(int sockfd, int backlog);

ñ backlog is maximum number of queued up connections (often about 20)

ñ Pattern of us

```
getaddrinfo(`H);
socket(`H);
bind(`H);
listen(`H);
/* accept(`H) goes here */
```

ñ DonȦt ignore any failures

listen(`H) is used by the **server**

# Accepting connections

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

- ñ Return value & arguments are similar to the above
- ñ Returns brand new descriptor           created by OS
- ñ Former descriptor continues to listen for new connections
- ñ New descriptor is used to send and receive data

accept(`H) is used by the   server

# Sending data

```
int send( int sockfd,
          const void *msg,
          int len, int flags);
```

-1 is returned on error, and `errno` is set to the error number.

- ñ **msg** is a pointer to the data to send
- ñ **len** is the length of the data to send
- ñ **flag** controls specifics, set to 0 for us
- ñ send(ˋH) returns the number of bytes sent, which might be less than what was requested

```
char *msg = ˋ    ECE297 is great fun, we learn useful stuff
int len, bytes_sent;
. . .
len = strlen(msg);
bytes_sent    = send(sockfd, msg, len, 0);
```

serˌAs responsibility to send the rest

!

# Receiving data

```
int recv( int sockfd,
          void *buf,
          int len, int flags);
```

- ñ **buf** is a buffer to read data into
- ñ **len** is the maximum length of the buffer
- ñ **flag** controls specifics, set to 0 for us
- ñ recv(`H) returns the number of bytes read
- ñ recv(`H) may return 0;    remoteẠs side closed connection on us

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

- ñ attempts to send any unsent data
- ñ closes socket for sending and receiving
- ñ returns -1 if error

-1 is returned on error, and `errno` is set to the error number.

ECE 297

# Protocols

ñ Set of rules that specifies data transfer between computing end-points, often including

  ã Connection establishment & tear-down

  ã Communication

  ã Data representation

ñ Often based on standards, de facto standards, open specifications, `H

# Common protocols

- IP (Internet Protocol
- UDP (User Datagram Protocol)
- TCP (Transmission Control Protocol)
- DHCP (Dynamic Host Configuration Protocol)
- HTTP (Hypertext Transfer Protocol)
- FTP (File Transfer Protocol)
- Telnet (Telnet Remote Protocol)
- SSH (Secure Shell Remote Protocol)
- SIP (Session Initiation Protocol)
- POP3 (Post Office Protocol 3)
- SMTP (Simple Mail Transfer Protocol)
- IMAP (Internet Message Access Protocol)

# Configuration files

ñThe configuration file sets      *properties*      for the storage serve
ñIt is read at startup to initialize the storage server.
ñIs not available to client            , which may reside across the net

File:  `default.conf` (given)

Server
location

Most software has
configuration
files
and ways to
overwrite
its settings.

```
server_host localhost
server_port 1111
table marks

# Data directory.
data_directory ./data
```

Port the
server
listesnts

File:`twotables.conf` (given)

```
server_host localhost
server_port 8888
table foo
table bar

# Data directory
data_directory ./data
```

Location
of data
files

Table
names

ECE 297

For a further example see onetable.

How the customer explained it

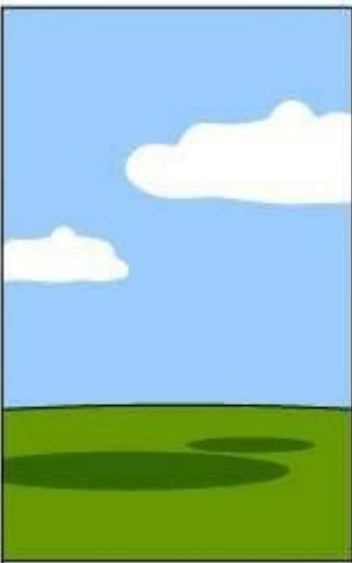How the Project Leader understood it

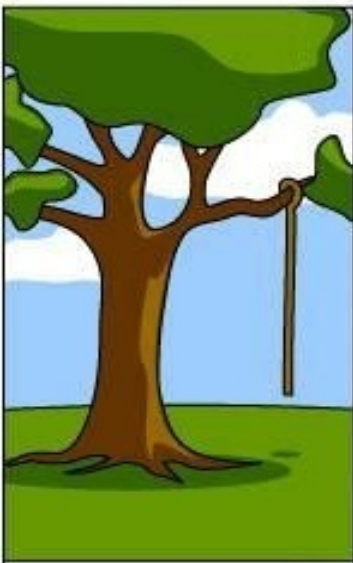How the Analyst designed it

How the Programmer wrote it

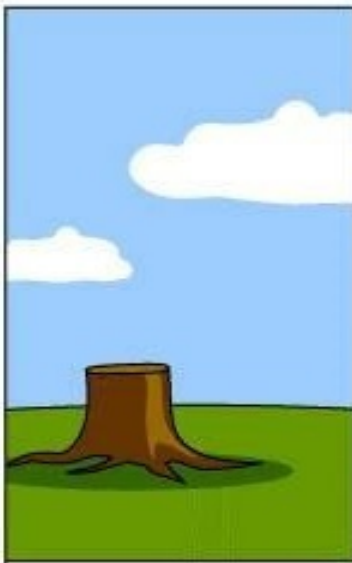How the Business Consultant described it

How the project was documented

What operations installed

How the customer was billed

How it was supported

What the customer really needed

Source: http://onproductmanagement.files.wordpress.com/2007/07/treecomicbig.jpg