

Course topics

- **Networking background**
- **Local storage**
- **Distributed file systems**
- **Storage Architectures**
 - Virtualizing storage, RAID, Storage-area networks
- **Other storage systems**
 - Untrusted storage, scalable systems, peer-to-peer systems
- **Other storage systems**

Class overview

- **Readings & class discussion**
- **Solo labs:**
 - Asynchronous programming: multifinger
 - Network programming: TCP proxy
 - Encrypting file system
- **Group labs:**
 - Basic file system
 - Final project
- **Midterm and final quizzes**

System calls

- **Problem: How to access resources other than CPU**
 - Disk, network, terminal, other processes
 - CPU prohibits instructions that would access devices
 - Only privileged OS “kernel” can access devices
- **Applications request I/O operations from kernel**
- **Kernel supplies well-defined *system call* interface**
 - Applications set up syscall arguments and *trap* to kernel
 - Kernel performs operation and returns result
- **Higher-level functions built on syscall interface**
 - `printf`, `scanf`, `gets`, etc. all user-level code

I/O through the file system

- **Applications “open” files/devices by name**
 - I/O happens through open files
- `int open(char *path, int flags, ...);`
 - flags: `O_RDONLY`, `O_WRONLY`, `O_RDWR`
 - `O_CREAT`: create the file if non-existent
 - `O_EXCL`: (w. `O_CREAT`) create if file exists already
 - `O_TRUNC`: Truncate the file
 - `O_APPEND`: Start writing from end of file
 - mode: final argument with `O_CREAT`
- **Returns file descriptor—used for all I/O to file**

Error returns

- **What if open fails? Returns -1 (invalid fd)**
- **Most system calls return -1 on failure**
 - Specific kind of error in global int errno
- **#include <sys/errno.h> for possible values**
 - 2 = ENOENT “No such file or directory”
 - 13 = EACCES “Permission Denied”
- **perror, strerror print human-readable messages**
 - perror ("initfile");
 - printf ("initfile: %s\n", strerror (errno));
→ “initfile: No such file or directory”

Operations on file descriptors

- `int read (int fd, void *buf, int nbytes);`
 - Returns number of bytes read
 - Returns 0 bytes at end of file, or -1 on error
- `int write (int fd, void *buf, int nbytes);`
 - Returns number of bytes read, -1 on error
- `off_t lseek (int fd, off_t pos, int whence);`
 - whence: 0 – start, 1 – current, 2 – end
 - Returns previous file offset, or -1 on error
- `int close (int fd);`

File descriptor numbers

- **File descriptors are inherited by processes**
 - When one process spawns another, same fds by default
- **Descriptors 0, 1, and 2 have special meaning**
 - 0 – “standard input” (stdin in ANSI C)
 - 1 – “standard output” (stdout, printf in ANSI C)
 - 2 – “standard error” (stderr, perror in ANSI C)
 - Normally all three attached to terminal

Manipulating file descriptors

- `int dup2 (int oldfd, int newfd);`
 - Closes `newfd`, if it was a valid descriptor
 - Makes `newfd` an exact copy of `oldfd`
 - Two file descriptors will share same offset
(`lseek` on one will affect both)
- `int fcntl (int fd, F_SETFD, int val)`
 - Sets *close on exec* flag if `val = 1`, clears if `val = 0`
 - Makes file descriptor non-inheritable by spawned programs

Pipes

- `int pipe (int fds [2]);`
 - Returns two file descriptors in `fds [0]` and `fds [1]`
 - Writes to `fds [1]` will be read on `fds [0]`
 - When last copy of `fds [1]` closed, `fds [0]` will return EOF
 - Returns 0 on success, -1 on error
- **Operations on pipes**
 - `read/write/close` – as with files
 - When `fds [1]` closed, `read (fds [0])` returns 0 bytes
 - When `fds [0]` closed, `write (fds [1])`:
 - Kills process with SIGPIPE, or if blocked
 - Fails with EPIPE

Sockets: Communication between machines

- **Datagram sockets: Unreliable message delivery**
 - On Internet: User Datagram Protocol (UDP)
 - Send atomic messages, which may be reordered or lost
 - Special system calls to read/write: `send/recv`
- **Stream sockets: Bi-directional pipes**
 - On Internet: Transmission Control Protocol (TCP)
 - Bytes written on one end read on the other
 - Reads may not return full amount requested—must re-read

Socket naming

- **Every Internet host has a unique 32-bit *IP address***
 - Often written in “dotted-quad” notation: 204.168.181.201
 - DNS protocol maps names (www.nyu.edu) to IP addresses
 - Network routes packets based on IP address
- **16-bit *port number* demultiplexes TCP traffic**
 - Well-known services “listen” on standard ports: finger—79, HTTP—80, mail—25, ssh—22
 - Clients connect from arbitrary ports to well known ports
 - A connection consists of five components: Protocol (TCP), local IP, local port, remote IP, remote port

Stream socket system calls

Client

socket – make socket

bind – assign address

connect – connect to listening socket

Server

socket – make socket

bind – assign address

listen – listen for clients

accept – accept connection

Example client

```
struct sockaddr_in {
    short    sin_family; /* = AF_INET */
    u_short  sin_port;   /* = htons (PORT) */
    struct   in_addr  sin_addr;
    char     sin_zero[8];
} sin;
```

```
int s = socket (AF_INET, SOCK_STREAM, 0);
bzero (&sin, sizeof (sin));
sin.sin_family = AF_INET;
sin.sin_port = htons (13); /* daytime port */
sin.sin_addr.s_addr = htonl (IP_ADDRESS);
connect (s, (sockaddr *) &sin, sizeof (sin));
while ((n = read (s, buf, sizeof (buf))) > 0)
    write (1, buf, n);
```

Example server

```
struct sockaddr_in sin;
int s = socket (AF_INET, SOCK_STREAM, 0);
bzero (&sin, sizeof (sin));
sin.sin_family = AF_INET;
sin.sin_port = htons (9999);
sin.sin_addr.s_addr = htonl (INADDR_ANY);
bind (s, (sockaddr *) &sin, sizeof (sin));
listen (s, 5);

for (;;) {
    socklen_t len = sizeof (sin);
    int cfd = accept (s, (sockaddr *) &sin, &len);
    /* do something with cfd */
    close (cfd);
}
```

Concurrent connections

- **Servers must handle multiple clients concurrently**
 - Read or write of a socket connected to slow client can block
 - Overlap network latency with CPU, transmission, disk I/O
 - Keep disk queues full when server accesses disk
- **Can use one process per client: accept, fork, close**
 - High overhead, cannot share state between clients
- **Can use threads for concurrency**
 - Data races and deadlock make programming tricky
 - Must allocate one stack per request
- **Use non-blocking read/write calls**
 - Unusual programming model

Non-blocking I/O

- `fcntl` sets `O_NONBLOCK` flag on descriptor
- **Non-blocking semantics of system calls:**
 - `read` immediately returns `-1` with `errno` `EAGAIN` if no data
 - `write` may not write all data, or may return `EAGAIN`
 - `connect` may fail with `EINPROGRESS`
 - `accept` may fail with `EAGAIN` if no pending connections

How do you know when to read/write?

```
struct timeval {
    long    tv_sec;           /* seconds */
    long    tv_usec;        /* and microseconds */
};

int select (int nfd, fd_set *readfds, fd_set *writefds,
            fd_set *exceptfds, struct timeval *timeout);

FD_SET(fd, &fdset);
FD_CLR(fd, &fdset);
FD_ISSET(fd, &fdset);
FD_ZERO(&fdset);
```

Asynchronous programming model

- **Many non-blocking file descriptors in one process**
 - Wait for pending I/O events on file many descriptors
 - Each event triggers some *callback* function
- **Lab: libasync – supports event-driven model**
 - Register callbacks on file descriptors
 - Call `amain()` – main select loop
 - Add/delete callbacks from within callbacks

callback.h

- **Problem: Need state from one callback to next**
- **wrap bundles a function with its arguments**

```
callback<void, int>::ref errwrite = wrap (write, 2);  
(*errwrite) ("hello", 5); // writes "hello" to stderr
```

- `void fdcb(int fd, selop op, cb_t cb);`
registers callbacks on file descriptor fd
 - op is selread or selwrite
 - cb is void callback (no arguments), or NULL to clear

libasync example server

```
void doaccept (int lfd) {
    sockaddr_in sin;
    bzero (&sin, sizeof (sin));
    socklen_t sinlen = sizeof (sin);
    int cfd = accept (lfd, (sockaddr *) &sin, &sinlen);
    if (cfd >= 0) { /* ... */ }
}

int main (int argc, char **argv) {
    // ...
    int lfd = inetsocket (SOCK_STREAM, your_port, INADDR_ANY);
    if (lfd < 0) fatal << "socket: " << strerror (errno) << "\n";
    if (listen (lfd, 5) < 0) fatal ("listen: %m\n");
    fdcb (lfd, selread, wrap (doaccept, lfd));
    amain ();
}
```

Remote procedure call

- **Abstract away network in distributed programs**
 - Idea: Distributed programming looks like function call
 - Reality: Can't abstract away everything (e.g., failure)
- **Next class: Sun RPC**
 - XDR defines structures that can be transmitted on the wire
 - RPC is simple layer that sends arg. structs and gets returns