

# 1 Linking

Linking is the process of collecting and combining various pieces of code and data into a single file that can be *loaded* (copied) into memory and executed. Programs are translated and linked using a compiler driver: `gcc -Og -o out main.c sum.c`.

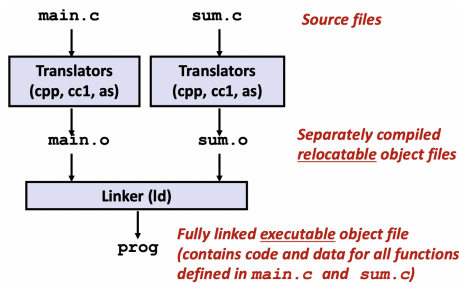


Figure 1. The static linking.

## 1.1 Why linkers?

- **Modularity.** Program can be written as a collection of smaller source files. Can build libraries of common functions.
- **Time efficiency.** Separate compilation. Change one source file, compile, and then relink.
- **Space efficiency.** Running memory images contain only code for the functions we actually use, rather than a whole library.

## 1.2 What do linkers do?

### 1.2.1 Symbol resolution

Programs define and reference *symbols* (global variables and functions). Symbol definitions are stored in object file by assembler in *symbol table*, array of struct which includes name, size, and location of each symbol. During symbol resolution, the linker associates each symbol reference with exactly one symbol definition.

### 1.2.2 Relocation

- Merge separate code and data sections into single section.
- Relocate symbols from their relative locations in the .O files to their final absolute memory locations in the executable.
- Update all references to symbols to reflect their new positions.

## 1.3 Object files (Modules)

### 1.3.1 Three kinds of object files

- **Relocatable object file** (.O file): Produced from one source (.C) file. Contains code and data in a form that can be combined with other relocatable object files to form executable object file.
- **Executable object file** (.out file): Contains code and data in a form that can be copied directly into memory and then executed.
- **Shared object file** (.so file): Special type of relocatable object file that can be loaded into memory and linked dynamically, at either load time or run-time. Called *dynamic link libraries* (DLLs) in Windows.

### 1.3.2 Executable and linkable format (ELF)

ELF binary is the standard binary format for object files. It is one unified format for .O, .out, and .so files.

- Elf header: word size, byte ordering, file type (.O, exec, .so), machine type
- Segment header table: page size, virtual addresses memory segments (sections), segment sizes.
- .text: code
- .rodata: read only data, e.g., jump tables
- .data: initialized global variables
- .bss (better save space): uninitialized global variables. Has sec-

ELF header
Segment header table (required for executables)
.text section
.rodata section
.data section
.bss section
.symtab section
.rel .text section
.rel .data section
.debug section
Section header table

Figure 2. ELF object file format.

tion header but occupies no space

- .symtab: symbol table, procedure and static variable names, section names and locations
- .rel .text: relocation info for code (.text)
- .rel .data: relocation info for .data section. Addresses of pointer data that will need to be modified in the merged executable
- .debug: info for symbolic debugging `gcc -g`
- Section header table: offsets and sizes of each section

## 1.4 Linker symbols

- **Global symbols:** Defined by module that can be referenced by other modules, e.g., non-static C functions and global variables
- **External symbols:** Referenced by module but defined by some other module
- **Local symbols:** Defined and referenced exclusively by module, e.g., static C functions and global variables

Local linker symbols are *not* local program variables.

### 1.4.1 Local symbols

- Local non-static C variables: stored on the stack
- Local static C variables: stored in either .bss or .data

Compiler allocates space in .data for each definition of x. Creates local symbols in the symbol table with unique names: x.1 and x.2.

```

1 int f(){
2     static int x = 0;
3     return x; }
4 int g(){
5     static int x = 1;
6     return x; }
  
```

### 1.4.2 Global variables

Avoid if you can. Otherwise, use static if you can or initialize. Use extern if you reference an external global variable.

## 1.5 Step 1: Symbol resolution

The input to the linker is relocatable object modules. What happens if multiple modules define global symbols with same name?

### 1.5.1 Global symbols are either *strong* or *weak*

- **Strong:** procedures and initialized globals
- **Weak:** uninitialized globals

### 1.5.2 Linker's symbol rules

1. **Multiple strong symbols are not allowed.** Each item can be defined only once. Otherwise: linker error
2. **Given a strong symbol and multiple weak symbols, choose the strong symbol.** References to the weak symbol resolve to the strong symbol.
3. **If there are multiple weak symbols, pick an arbitrary one.** Can override this with `gcc -fno-common`.

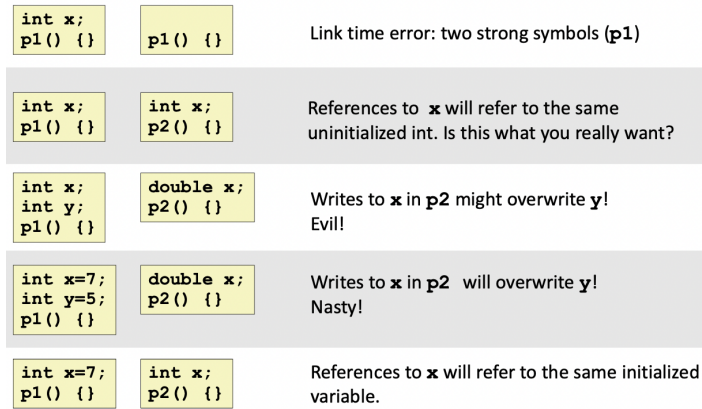


Figure 3. Linker puzzles.

## 1.6 Step 2: Relocation

Once the linker has completed the symbol resolution, it knows the exact sizes of the code and data sections in its input object modules. Then linker begin relocation, where it merges the input modules and assigns run-time addresses to each symbol. Relocation consist of two steps:

- Relocating sections and symbol definitions.** The linker merges all sections of the same type into a new aggregate section. The linker then assigns run-time memory addresses to the new aggregate section, to each section and to each symbol defined by input modules, so that each instruction and global variable in the program has a unique run-time memory address.
- Relocating symbol references within sections.** The linker modifies every symbol reference in the bodies of the code and data sections so that they point the correct run-time addresses.

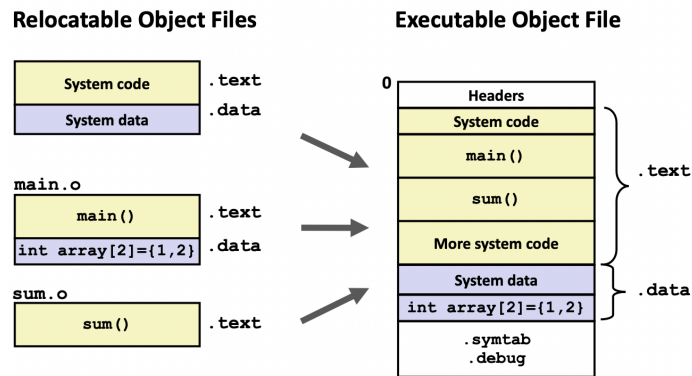


Figure 4. Relocation.

### 1.6.1 Relocation entries

When an assembler generates an object module, it does not know where the code, data, and externally defined functions or global variables that are referenced be stored in memory. So whenever assembler encounters a reference to an object whose location is unknown, it generates a *relocation entry* that tells the linker how to modify the reference when it merges the object file into an executable. Relocation entries for code are placed in `.rel.text` and data are placed in `.rel.data`.

#### Format of an ELF relocation entry

```
1 typedef struct {
2     long offset; // Offset of the reference to relocate
3     long type:32, // Relocation type
4     long symbol:32; // Symbol table index
5     long addend; // Constant part of relocation
6     expression
7 } Elf64_Rela;
```

ELF defines 32 different relocation types, but we are concerned with only the two most basic types:

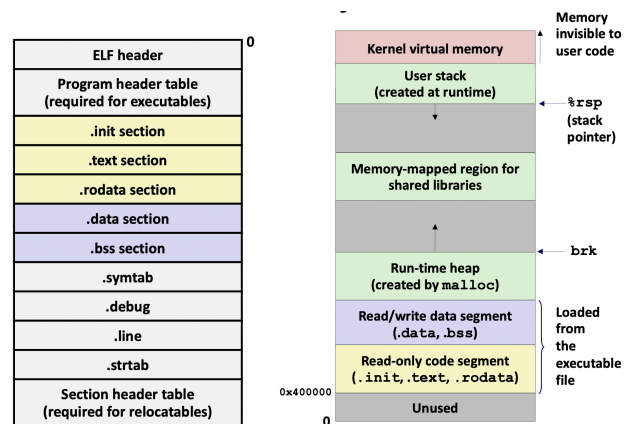
- R\_X86\_64\_PC32.** Relocate a reference that uses a 32-bit PC-relative address, *i.e.*, offset from current run-time value of PC.
- R\_X86\_64\_32.** Relocate a reference that uses a 32-bit absolute address. CPU directly uses 32-bit value encoded in the instruction as *effective address*, *e.g.*, the target of `call` instruction

### 1.6.2 Relocating symbol references

#### Pseudocode for the linker's relocation algorithm

```
1 foreach section s {
2     foreach relocation entry r {
3         rp = s + r.offset; // ptr to ref to be relocated
4         if (r.type == R_X86_64_PC32) {
5             ra = ADDR(s) + r.offset; // run-time address of ref
6             *rp = (unsigned) (ADDR(r.symbol) + r.addend - ra); }
7         if (r.type == R_X86_64_32){
8             *rp = (unsigned) (ADDR(r.symbol) + r.addend);}
9     }
10 }
```

## 1.7 Executable object files



(a) ELF executable object file (b) Runtime memory image

- `.text`, `.rodata`, `.data`: similar to object file
- `.init`: defines a small function called `_init` that will be called by the program's initialization code
- No `.rel` sections since the executable is *fully linked*
- Program header table: describes mapping from contiguous chunks of the executable file to contiguous memory segments

The loader copies the code and data in the executable object file from disk into memory and then runs the program by jumping to its first instruction, or *entry point*.

## 1.8 Static libraries (.a archive files)

- Concatenate related relocatable object files into a single file with an index, called an *archive*
- Enhance linker so that it tries to resolve unresolved external references by looking one or more archives
- If archive member file resolves reference, link it

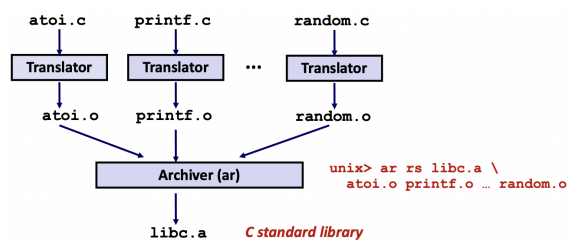


Figure 5. Creating static libraries

Archiver allows incremental updates. It recompiles function that changes and replace `.o` file in archive.

### 1.8.1 Commonly used libraries

- **libc.a** (C standard library): I/O, memory allocation, ...
- **libm.a** (C math library): sin, cos, log, exp, ...

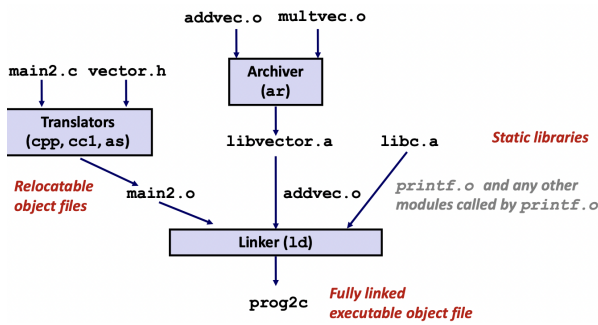


Figure 6. Linking with static libraries

### 1.8.2 Solving external references

1. Scan .o files and .a files in the command line order
2. For each scan, try to resolve each unresolved reference
3. If any entry is unresolved at the end of scan, error

Command line order matters.

## 1.9 Shared libraries (.so files)

Static libraries have the following disadvantages:

- Duplication in the stored executables (every function needs libc)
- Duplication in running executable
- Minor bug fixes of system libraries require each application to explicitly relink

Solution is using *shared libraries*: object files that contain code and data that are loaded and linked into an application *dynamically*, at either load-time and run-time. It is also called dynamic link libraries (DLLs) and performed by a *dynamic linker*. A single copy of the .text section of a shared library in memory can be shared by different running processes.

### 1.9.1 Load-time linking

gcc -shared -o libvector.so addvec.c multvec.c

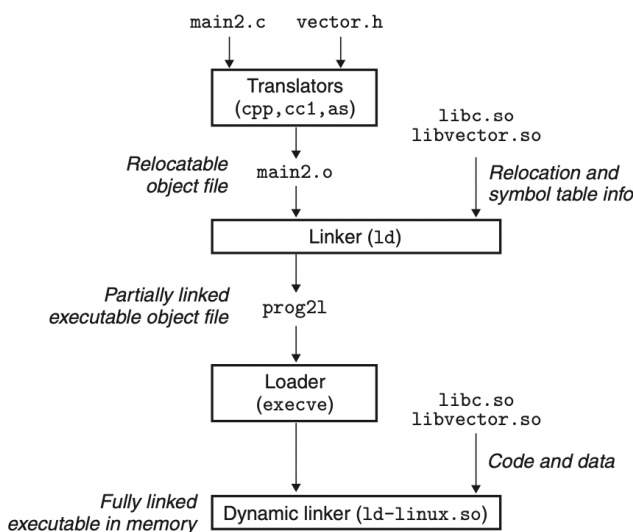


Figure 7. Dynamic linking in load time

### 1.9.2 Run-time linking

In Linux, this is done by calls to the dlopen().

Code 1. Dynamic linking at run-time

```

1 int x[2] = {1, 2};
2 int y[2] = {3, 4};
3 int z[2];
4
5 int main() {
6     void *handle;
7     void (*addvec)(int *, int *, int *, int);
8     char *error;
9     /* Dynamically load the shared library that contains
10    addvec() */
11    handle = dlopen("./libvector.so", RTLD_LAZY);
12    if (!handle) {
13        fprintf(stderr, "%s\n", dlerror());
14        exit(1); }
15    /* Get a pointer to the addvec() function we just
16    loaded */
17    addvec = dlsym(handle, "addvec");
18    if ((error = dlerror()) != NULL) {
19        fprintf(stderr, "%s\n", error);
20        exit(1); }
21    /* Now we can call addvec() just like any other
22    function */
23    addvec(x,y, z, 2);
24    printf("z = [%d %d]\n", z[0], z[1]);
25    /* Unload the shared library */
26    if (dlclose(handle) < 0) {
27        fprintf(stderr, "%s\n", dlerror());
28        exit(1); }
29 }

```

## 1.10 Library interpositioning

*Library interpositioning* is a powerful linking technique that allows programmers to intercept calls to arbitrary functions. Given target function, create a wrapper function whose prototype is identical to the target function. Then trick the system into calling the wrapper function instead of the target function. Interpositioning can occur at compile time, link time, or load/run time.

### 1.10.1 Applications

- Security
- Debugging
- Monitoring and profiling
  - Count number of calls to functions
  - Characterize call sites and arguments to functions
- Malloc tracing
  - Detecting memory leaks
  - Generating address traces

## 2 Exceptional control flow

Processors do only one thing: read and execute a sequence of instructions, one at a time. This sequence is the CPU's *control flow*. The abrupt changes to the flow is caused by instructions such as jumps, calls, and returns. Such instructions allow programs to react to changes in internal program state.

But systems must also be able to react to changes in system state that are not captured by internal program variables and are not necessarily related to the execution of the program, for example:

- User hits Ctrl-C at the keyboard
- Instruction divides by zero
- Data arrives from a disk or a network adapter

Modern systems react to these situations by making abrupt changes in the control flow, which is called *exceptional control flow* (ECF).

Mechanism	Level	Implemented by
Exceptions	Low	Hardware and OS software
Process context switch	High	Os software and hardware timer
Signals	High	OS software
Nonlocal jumps	High	C runtime library

Table 2. ECF exists at all levels of a computer system.

### 2.1 Exceptions

An *exception* is a transfer of control to the OS *kernel* in response to some event, e.g., divide by 0, arithmetic overflow, page fault, I/O request completes, typing Ctrl-C. As shown in Fig. 8, a change in

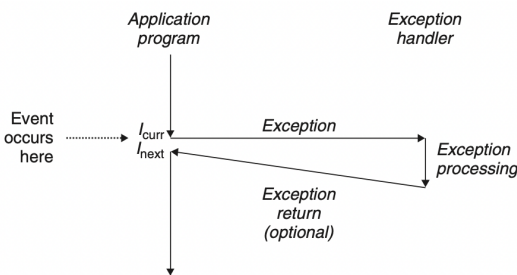


Figure 8. Anatomy of an exception.

the processor's state (an event) triggers an abrupt control transfer (an exception) from the application program to an exception handler. After finishing processing exception, one of below happens:

- Handler returns control to current instruction  $I_{curr}$ .
- Handler returns control to next instruction  $I_{next}$ , i.e., the instruction that would have executed after  $I_{curr}$  if event not occurred
- Handler aborts the interrupted program

#### 2.1.1 Exception handling

Each type of possible exception in a system is assigned a unique nonnegative integer *exception number*. At the system boot time, OS allocates and initializes a jump table called *exception table*, so that entry  $k$  contains the address of the handler for exception  $k$ . The starting address of exception table is contained in special CPU register called the *exception table base register*.

Class	Cause	Example	Sync	Return behavior
Interrupt	Signal from I/O device	Ctrl-C, arrival of packet or data	Async	Next
Trap	Intentional exception	<b>System calls</b> , breakpoints	Sync	Next
Fault	Potentially recoverable error	Page and protection faults, floating point	Sync	Current or abort
Abort	Nonrecoverable error	Illegal instruction, parity error	Sync	Abort

Table 1. Classes of exceptions.

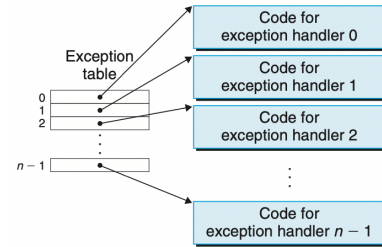


Figure 9. Exception table.

At run time, the processor detects that an event has occurred and determines the corresponding exception number  $k$ . The processor then triggers the exception by making an *indirect procedure call* through entry  $k$ .

1. Processor pushes a return address on the stack, either current or next instruction depending on the class of exception
2. Processor pushes an additional processor state on the stack that will be necessary to restart the interrupted program when the handler returns, e.g., current condition codes
3. When control is transferred from a user program to the kernel, all of these items are pushed onto the kernel's stack
4. Exception handlers run in *kernel mode*, which means they have complete access to all system resources

#### 2.1.2 Classes of exceptions

- **Asynchronous:** caused by events external to the processor
- **Synchronous:** caused by events that occur as a result of executing an instruction

### 2.2 System calls

*System call* is user process calling a system (kernel).

Number	Name	Description
0	read	Read file
1	write	Write file
2	open	Open file
3	close	Close file
4	stat	Get info about file
57	fork	Create process
59	execve	Execute a program
60	_exit	Terminate process
62	kill	Send signal to process

Table 3. List of system calls

#### 2.2.1 Example: Opening file

User calls `open(filename, options)`. It calls `__open` function, which invokes system call instruction `syscall`.

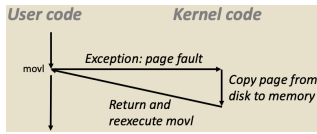
```

0000000000e5d70 <__open>:
...
e5d79: b8 02 00 00 00  mov $0x2,%eax # open is syscall #2
e5d7e: 0f 05          syscall # Return value in %rax
e5d80: 48 3d 01 f0 ff ff  cmp $0xffffffffffff001,%rax
...
e5dfa: c3            retq
    
```

Table 1. Classes of exceptions.

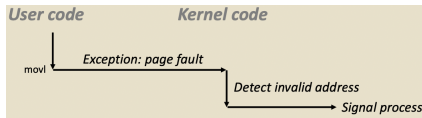
### 2.2.2 Example: Page fault

User writes to memory location, but that page of user's memory is currently on disk.



### 2.2.3 Example: Invalid memory reference

Send SIGSEGV signal to user process. User process exit with segmentation fault.



## 2.3 Processes

A *process* is an instance of a running program. Process provides each program with two key abstractions:

- **Logical control flow:** Seems to have exclusive use of the CPU. Provided by kernel mechanism called *context switching*
- **Private address space:** Seems to have exclusive use of main memory. Provided by kernel mechanism called *virtual memory*.

### 2.3.1 Multiprocessing

Traditionally, single processor executes multiple processes concurrently. Register values for processes saved in memory and loaded. Recently, we use mutlicore processors. There are multiple CPUs on a single chip, sharing main memory and some of the caches. So, each can execute a separate process. Scheduling of processors onto cores done by kernel.

### 2.3.2 Concurrent processes

Each process is a logical control flow. Two processes run *concurrently* if their flows overlap in time. Otherwise, they're *sequential*.

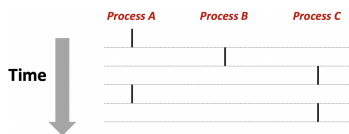


Figure 10. A & B, A & C are concurrent. B & C is sequential.

Control flows for concurrent processes are physically disjoint in time. However, user views they are running in parallel.

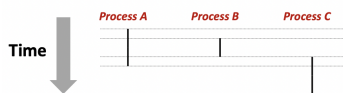


Figure 11. User view of concurrent processes.

### 2.3.3 Context switching

Processes are managed by a shared chunk of memory-resident OS code called *kernel*. The kernel is not a separate process, but rather runs as part of some existing process. Control flow passes from one process to another via a *context switch*.

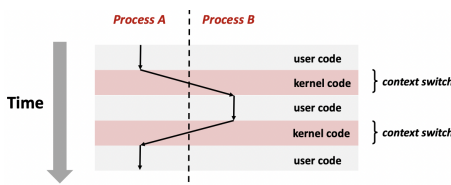


Figure 12. Context switch.

## 2.4 Process control

### 2.4.1 System call error handling

On error, Linux system-level functions typically return -1 and set global variable `errno` to indicate cause. So check return status of every system-level function, except a few that return `void`.

#### Code 6. Error handling

```

1 /* Error-reporting function */
2 void unix_error(char *msg) {
3     fprintf(stderr, "%s: %s\n", msg, strerror(errno));
4     exit(0);
5 }
6 /* Error-handling wrappers */
7 pid_t Fork(void) {
8     pid_t pid;
9     if ((pid = fork()) < 0) unix_error("Fork error");
10    return pid;
11 }
12 pid = Fork();
    
```

### 2.4.2 Obtaining process IDs

- `pid_t getpid(void)` returns PID of current process
- `pid_t getppid(void)` returns PID of parent process

### 2.4.3 State of processes

- **Running:** either executing, or waiting to be executed and will eventually be scheduled (*i.e.*, chosen to execute) by the kernel
- **Stopped:** execution is suspended and will not be scheduled until further notice
- **Terminated:** stopped permanently

### 2.4.4 Terminating processes

Process becomes terminated for one of three reasons:

- Receiving a signal whose default action is to terminate
- Returning from the `main` routine
- Calling the `exit` function. `void exit(int status)` terminates with an exit status of `status`. Normal returns status is 0 and nonzero indicates error. `exit` is called once and never returns.

### 2.4.5 Creating processes

Parent process creates a new running child process by calling `fork`. `int fork(void)`

- Returns 0 to the child process
- Returns child's PID to parent process
- Child gets an identical copy of the parent's virtual address space
- Child gets identical copies of the parent's open file descriptors
- Child has a different PID with the parent

### 2.4.6 Process graphs

A *process graph* captures the partial ordering of statements in a concurrent program.

- Each vertex is the execution of a statement
- $a \rightarrow b$  means  $a$  happens before  $b$
- Edges labeled with current value of variables
- Graph begins with a vertex with no incoming edges

Any topological sort of the graph corresponds to a feasible total ordering. It is when all edges point from left to right.

### 2.4.7 Process groups

Every process belongs to exactly one process group.

- `getpgrp()`: return process group of current process
- `setpgid()`: change process group of a process

### 2.4.8 Fork examples

#### Code 7. Fork example 1

```

1 int main() {
2     pid_t pid;
3     int x = 1;
4     pid = Fork();
5     if (pid == 0) { /* Child */
6         printf("child: x=%d\n", ++x); exit(0); }
7     /* Parent */
8     printf("parent: x=%d\n", --x); exit(0);
9 }
10 /* Output: parent: x=0  child: x=2 */
    
```

- Call one, return twice
- Concurrent execution: cannot predict execution order
- Duplicate but separate address spaces
  - x has a value of 1 when fork returns in parent and child
  - Changes in x are independent
- Shared open files: stdout is the same in both parent and child

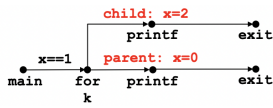


Figure 13. Process graph of fork example.

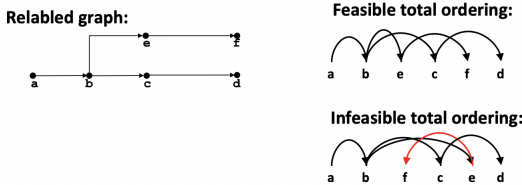


Figure 14. Feasible ordering of graph.

```

void fork2()
{
    printf("L0\n");
    fork();
    printf("L1\n");
    fork();
    printf("Bye\n");
}
        
```

Feasible output:

```

L0
L1
Bye
Bye
L1
Bye
Bye
        
```

Infeasible output:

```

L0
Bye
L1
Bye
L1
Bye
Bye
        
```

```

void fork5()
{
    printf("L0\n");
    if (fork() == 0) {
        printf("L1\n");
        if (fork() == 0) {
            printf("L2\n");
        }
    }
    printf("Bye\n");
}
        
```

Feasible output:

```

L0
Bye
L1
L2
Bye
Bye
        
```

Infeasible output:

```

L0
L0
Bye
L1
Bye
Bye
L2
        
```

```

void fork4()
{
    printf("L0\n");
    if (fork() != 0) {
        printf("L1\n");
        if (fork() != 0) {
            printf("L2\n");
        }
    }
    printf("Bye\n");
}
        
```

Feasible output:

```

L0
L1
Bye
Bye
L2
Bye
        
```

Infeasible output:

```

L0
L0
Bye
L1
Bye
Bye
L2
        
```

Figure 15. Fork examples.

### 2.4.9 Reaping child processes

When process terminates, it still consumes system resources, called a *zombie*. *Reaping* is performed by parent on terminated child. Parent is given exit status information. Kernel then deletes zombie child process. If any parent terminates without reaping a child, then the orphaned child will be reaped by init process (pid==1). So, we only need explicit reaping in long-running processes, e.g., shells and servers.

**Zombie Example**

```

void fork7() {
    if (fork() == 0) {
        /* Child */
        printf("Terminating Child, PID = %d\n", getpid());
        exit(0);
    } else {
        printf("Running Parent, PID = %d\n", getpid());
        while (1)
            ; /* Infinite loop */
    }
}
        
```

```

linux> ./forks 7 &
[1] 6639
Running Parent, PID = 6639
Terminating Child, PID = 6640
linux> ps
PID TTY          TIME CMD
6585 tty9          00:00:00 tcsh
6639 tty9          00:00:03 forks
6640 tty9          00:00:00 forks <defunct>
6641 tty9          00:00:00 ps
linux> kill 6639
[1] Terminated
linux> ps
PID TTY          TIME CMD
6585 tty9          00:00:00 tcsh
6642 tty9          00:00:00 ps
        
```

ps shows child process as "defunct" (i.e., a zombie)

Killing parent allows child to be reaped by init

forks.c

**Non-terminating Child Example**

```

void fork8()
{
    if (fork() == 0) {
        /* Child */
        printf("Running Child, PID = %d\n",
            getpid());
        while (1)
            ; /* Infinite loop */
    } else {
        printf("Terminating Parent, PID = %d\n",
            getpid());
        exit(0);
    }
}
        
```

```

linux> ./forks 8
Terminating Parent, PID = 6675
Running Child, PID = 6676
linux> ps
PID TTY          TIME CMD
6585 tty9          00:00:00 tcsh
6676 tty9          00:00:06 forks
6677 tty9          00:00:00 ps
linux> kill 6676
linux> ps
PID TTY          TIME CMD
6585 tty9          00:00:00 tcsh
6678 tty9          00:00:00 ps
        
```

Child process still active even though parent has terminated

Must kill child explicitly, or else will keep running indefinitely

forks.c

Figure 16. Zombie and non-terminating child examples

#### 2.4.10 wait: Synchronizing with children

int wait(int \*child\_status).

- Suspends current process until one of its children terminates
- Return value is the pid of the child process that terminated
- If child\_status!=NULL, then the integer it points to will be set to a value that indicates the reason the child terminated and the exit status

```

void fork9() {
    int child_status;
    if (fork() == 0) {
        printf("HC: hello from child\n");
        exit(0);
    } else {
        printf("HP: hello from parent\n");
        wait(&child_status);
        printf("CT: child has terminated\n");
    }
    printf("Bye\n");
}
        
```

Feasible output:

```

HC
HP
CT
Bye
        
```

Infeasible output:

```

HP
CT
Bye
HC
        
```

```

void fork10() {
    pid_t pid[N];
    int i, child_status;
    for (i = 0; i < N; i++)
        if ((pid[i] = fork()) == 0) {
            exit(100+i); /* Child */
        }
    for (i = 0; i < N; i++) { /* Parent */
        pid_t wpid = wait(&child_status);
        if (WIFEXITED(child_status))
            printf("Child %d terminated with exit status %d\n",
                wpid, WEXITSTATUS(child_status));
        else
            printf("Child %d terminate abnormally\n", wpid);
    }
}
        
```

Feasible output:

```

Child 0 terminate abnormally
Child 1 terminate abnormally
Child 2 terminate abnormally
        
```

Infeasible output:

```

Child 0 terminate abnormally
Child 1 terminate abnormally
Child 2 terminate abnormally
        
```

Figure 17. Wait example

If multiple children completed, it will take in arbitrary order. We can use macros WIFEXITED and WEXITSTATUS to get information about exit status.

### 2.4.11 waitpid: Waiting for a specific process

pid\_t waitpid(pid\_t pid, int &status, int options) suspends current process until specific process terminates.

```
pid_t wpid = waitpid(pid[i], &child_status, 0);
if (WIFEXITED(child_status))
    printf("Child %d terminated with exit status %d\n",
        wpid, WEXITSTATUS(child_status));
else
    printf("Child %d terminate abnormally\n", wpid);
```

Figure 18. Waitpid example

### 2.4.12 execve: Loading and running programs

int execve(char \*filename, char \*argv[], char \*envp[]) loads and runs executable file filename in the current process with argument list argv and environment variable list envp. By convention, argv[0] == filename. envp is in format of name=value.

It overwrites code, data, and stack. IT retains PID, open files, and signal context. It is called once and never returns except error.

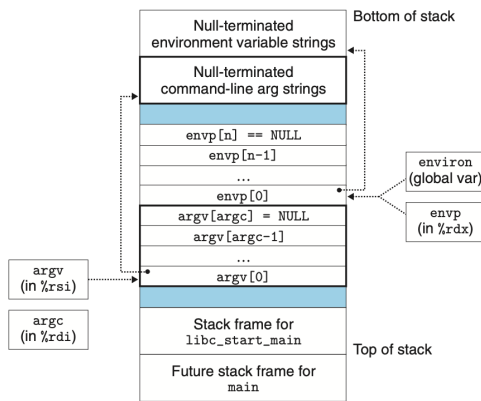
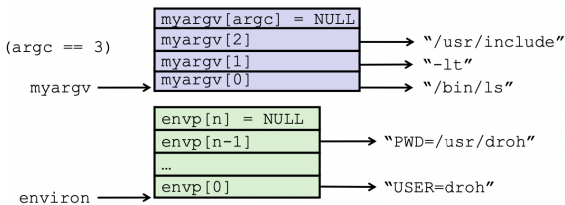


Figure 19. Organization of user stack when new program starts.

Executes "/bin/ls -lt /usr/include" in child process using current environment:



```
if ((pid = Fork()) == 0) { /* Child runs program */
    if (execve(myargv[0], myargv, environ) < 0) {
        printf("%s: Command not found.\n", myargv[0]);
        exit(1);
    }
}
```

Figure 20. execve example

## 2.5 Shell

Shell is application program that runs programs on behalf of user.

```
void eval(char *cmdline)
{
    char *argv[MAXARGS]; /* Argument list execve() */
    char buf[MAXLINE]; /* Holds modified command line */
    int bg; /* Should the job run in bg or fg? */
    pid_t pid; /* Process id */

    strcpy(buf, cmdline);
    bg = parse_line(buf, argv);
    if (argv[0] == NULL)
        return; /* Ignore empty lines */

    if (!builtin_command(argv)) {
        if ((pid = Fork()) == 0) { /* Child runs user job */
            if (execve(argv[0], argv, environ) < 0) {
                printf("%s: Command not found.\n", argv[0]);
                exit(0);
            }
        }

        /* Parent waits for foreground job to terminate */
        if (!bg) {
            int status;
            if (waitpid(pid, &status, 0) < 0)
                unix_error("waitfg: waitpid error");
        }
        else
            printf("%d %s", pid, cmdline);
    }
    return;
}
```

Figure 21. Simple shell eval function

The example shell correctly waits for and reaps foreground jobs. But what about background jobs?

- Will become zombie when they terminate
- Will never be reaped because shell will not terminate
- Will create a memory leak that could run kernel out of memory

### 2.5.1 Process states

ps w shows the state of the processes.

No.	Letter	Meaning
First	S	sleeping
	T	stopped
	R	running
Second	s	session leader
	+	foreground group

Table 4. STAT legend

## 2.6 Signals

A signal is a message that notifies a process that an event of some type has occurred in the system. It is sent from kernel to process.

ID	Name	Default action	Event
2	SIGINT	Terminate	User typed Ctrl-C
9	SIGKILL	Terminate	Kill program
11	SIGSEGV	Terminate	Segmentation violation
14	SIGALRM	Terminate	Timer signal
17	SIGCHLD	Ignore	Child stopped or terminated

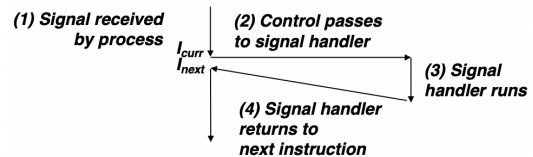
Table 5. List of signals

Kernel sends a signal to a destination process by updating some state in the context of the process. Kernel sends a signal because:

- Kernel has detected a system event such as divide-by-zero or the termination of a child process
- Another process has invoked the kill call to request the kernel to send a signal to the destination process

A destination process receives a signal when it is forced by the signal to react in some way to the delivery of the signal.

- Ignore signal: do nothing
- Terminate to process
- Catch signal by executing a user-level function: signal handler



A signal is pending if sent but not yet received. Signals are not queued. There can be at most one pending signal of certain type. A process can block receipt of certain signals. Blocked signals can be delivered, but will not be received until signal is unblocked.

Kernel maintains pending and blocked bit vectors in the context of each process.

- Kernel sets/clears bit k in pending when a signal of type k is delivered/received
- blocked can be set and cleared using the sigprocmask function

### 2.6.1 Sending signals

/bin/kill sends arbitrary signal to a process or process group.

- /bin/kill -9 24818: send SIGKILL to process 24818
- /bin/kill -9 -24817: send SIGKILL to processes in group 24817

Typing Ctrl-C/Ctrl-Z causes the kernel to send a SIGINT/SIGTSTP to every job in the foreground process group.

### 2.6.2 Receiving signals

Suppose kernel is returning from an exception handler and is ready to pass control to process `pnb`. Kernel computes `pnb = pending & ~ blocked`, the set of pending nonblocked signals of process `p`.

- If `pnb==0`, pass control to next instruction in logical flow for `p`
- Else, repeat for all nonzero `k` in `pnb`:
  - Choose least nonzero bit `k` in `pnb` and force processes `p` to receive signal `k`
  - The receipt of the signal triggers some action by `p`
  - Pass control to next instruction in logical flow for `p`

Each signal type has a predefined default action, which is one of:

- The process terminates
- The process stops until restarted by a `SIGCONT` signal
- The process ignores the signal

### 2.6.3 Installing signal handlers

The signal function modifies the default action associated with the receipt of signal `sigum`:

`handler_t *signal(int sigum, handler_t *handler).`

- `SIG_IGN`: ignore signals of type `sigum`
- `SIG_DFL`: revert to the default action of type `sigum`
- Otherwise, handler is the address of a user-level signal handler.

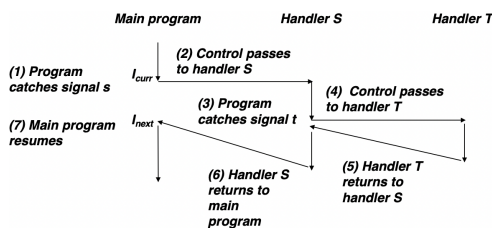
#### Code 8. Signal handling example

```
1 /* Install the SIGINT handler */
2 if (signal(SIGINT, sigint_handler) == SIG_ERR)
3     unix_error("signal error");
4 pause(); // Wait for the receipt of a signal
```

Signal handler is a separate logical flow, not process, that runs concurrently with the main program.

### 2.6.4 Nested signal handlers

Handlers can be interrupted by other handlers.



### 2.6.5 Blocking and unblocking signals

To avoid nested signal handling, block signals. Kernel **implicitly** blocks any pending signals of type currently being handled. **Explicit blocking** is done by `sigprocmask` function. Supporting functions are:

<code>sigemptyset</code>	Create empty set
<code>sigfillset</code>	Add every signal number to set
<code>sigaddset</code>	Add signal number to set
<code>sigdelset</code>	Delete signal number from set

#### Code 9. Temporarily blocking signals

```
1 sigset_t mask, prev_mask;
2 Sigemptyset(&mask);
3 Sigaddset(&mask, SIGINT);
4 /* Block SIGINT and save previous blocked set */
5 Sigprocmask(SIG_BLOCK, &mask, &prev_mask);
6 /* Code region that will not be interrupted by SIGINT */
7 /* Restore previous blocked set, unblocking SIGINT */
8 Sigprocmask(SIG_SETMASK, &prev_mask, NULL);
```

### 2.6.6 Async-Signal-Safety

Function is `async-signal-safe` if either *reentrant*, i.e., all variables stored on stack frame, or non-interruptible by signals.

- Safe: `_exit`, `write`, `wait`, `watipid`, `sleep`, `kill`
- Not safe: `printf`, `malloc`, `exit`

`write` is the only `async-signal-safe` output function.

### 2.6.7 Example: Wait all child process

Put `wait` in a loop to reap all terminated children. **This is wrong.** Parent process don't have to wait for all child to terminate. It needs to do its own work. Since we install handler to `SIGCHLD`, just handle processes that gave parent `SIGCHLD`.

```
1 void child_handler2(int sig){
2     int olderrno = errno;
3     pid_t pid;
4     while ((pid = wait(NULL)) > 0) {
5         ccount--;
6         Sio_puts("Handler reaped child ");
7         Sio_putl((long)pid);
8     }
9     errno = olderrno;
10 }
```

### 2.6.8 Example: Synchronizing flows

`mask_one` corrects the synchronization error that assumes parent runs before child. Without it, parent would not receive `SIGCHLD` since it did not added child yet.

```
1 void handler(int sig) {
2     int olderrno = errno;
3     sigset_t mask_all, prev_all;
4     pid_t pid;
5     Sigfillset(&mask_all);
6     while ((pid = waitpid(-1, NULL, 0)) > 0) { /* Reap
7         child */
8         Sigprocmask(SIG_BLOCK, &mask_all, &prev_all);
9         deletejob(pid);
10        Sigprocmask(SIG_SETMASK, &prev_all, NULL);
11    }
12    errno = olderrno;
13 }
14 int main(int argc, char **argv){
15     int pid;
16     sigset_t mask_all, mask_one, prev_one;
17     Sigfillset(&mask_all);
18     Sigemptyset(&mask_one);
19     Sigaddset(&mask_one, SIGCHLD);
20     Signal(SIGCHLD, handler);
21     initjobs(); /* Initialize the job list */
22     while (1) {
23         Sigprocmask(SIG_BLOCK, &mask_one, &prev_one); /*
24         Block SIGCHLD */
25         if ((pid = Fork()) == 0) { /* Child process */
26             Sigprocmask(SIG_SETMASK, &prev_one, NULL); /*
27             Unblock SIGCHLD. Child may make its children! */
28             Execl("/bin/date", argv, NULL);
29         }
30         Sigprocmask(SIG_BLOCK, &mask_all, NULL); /*
31         Parent process */
32         addjob(pid);
33         Sigprocmask(SIG_SETMASK, &prev_one, NULL); /*
34         Unblock SIGCHLD */
35     }
36     exit(0);
37 }
```



### 2.6.9 Example: Explicitly waiting for signal

Handlers for program explicitly waiting for SIGCHLD to arrive. Code below is a still **nonsense** (but correct) since it immediately returns child process after fork, and it waits for SIGCHLD instead of adding child process to the job list. The parent process explicitly wait for the global variable pid to change by the sigchld\_handler. This is still inaccurate, but similar to a shell waiting for a foreground job.

```

1 volatile sig_atomic_t pid;
2 void sigchld_handler(int s) {
3     int olderrno = errno;
4     pid = Waitpid(-1, NULL, 0); /* Main is waiting for
5     nonzero pid */
6     errno = olderrno;
7 }
8 void sigint_handler(int s) { }
9 int main(int argc, char **argv) {
10    sigset_t mask, prev;
11    Signal(SIGCHLD, sigchld_handler);
12    Signal(SIGINT, sigint_handler);
13    Sigemptyset(&mask);
14    Sigaddset(&mask, SIGCHLD);
15    while (1) {
16        Sigprocmask(SIG_BLOCK, &mask, &prev); /* Block
17        SIGCHLD */
18        if (Fork() == 0) /* Child */
19            exit(0);
20        /* Parent */
21        pid = 0;
22        Sigprocmask(SIG_SETMASK, &prev, NULL); /* Unblock
23        SIGCHLD */
24        /* Wait for SIGCHLD to be received (wasteful!) */
25        while (!pid);
26        /* Do some work after receiving SIGCHLD */
27        printf(".");
28    }

```

Let's think of other options than waiting. First approach is using **pause**, which waits until any signal comes. However it could go around (*race*) when the process get the SIGCHLD signal just after entering the while but just before **pause**. Then **pause** will wait forever, because the process already got signal!

```

1 while(!pid) pause();

```

Second approach is using **wait**, which never going to stuck you. However, it is too *slow*. The child may end just after 1ms. But process should still wait for 1s.

```

1 while(!pid) sleep(1);

```

Solution is using: `int sigsuspend(const sigset_t *mask).`

```

1 while(!pid) sigsuspend(&prev);

```

It is equivalent to uninterreuptable (*atomic*) version of:

```

1 sigprocmask(SIG_BLOCK, &mask, &prev);
2 pause();
3 sigprocmask(SIG_SETMASK, &prev, NULL);

```

Since we blocked **mask** and the operations are atomic, **pid** cannot change before **pause**.

## 3 Virtual Memory (VM)

### 3.1 Physical addressing versus virtual addressing

Physical addressing	Virtual addressing
Simple systems like camera Directly access physical	Most of the devices MMU translates virtual to physical

### 3.2 Why VM?

- Use main memory efficiently: Use DRAM as a cache for parts of a virtual address space
- Simplifies memory management: Each process gets the same uniform linear address space
- Isolates address spaces
  - One process can't interfere with another's memory
  - User program cannot access privileged kernel info and code

#### 3.2.1 Address spaces

Linear	Ordered set of contiguous integer	{0, 1, ...}
Virtual	Set of $N = 2^n$ virtual addresses	{0, 1, ..., $N - 1$ }
Physical	Set of $M = 2^m$ physical addresses	{0, 1, ..., $M - 1$ }

Table 6. List of address spaces

### 3.3 VM as a tool for caching

Conceptually, VM is an array of contiguous bytes stored on disk. Since it is typically larger than physical memory, VM provides mechanism for using the DRAM as cache.

#### 3.3.1 Pages

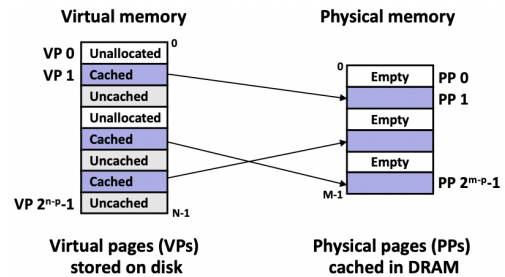


Figure 22. Virtual pages

The memory is splitted into *page*, a chunk of caches. Size of page is  $P = 2^p$  bytes. In Linux system,  $p = 12$ , which means that we need 12 bits to access particular bytes in page. The number of virtual pages is  $2^{n-p} - 1$ , excluding the first VP which is unallocated.

#### 3.3.2 DRAM cache organization

DRAM is about 10x slower than SRAM, and disk is about 10,000x slower than DRAM. DRAM cache organization is driven by the enormous miss penalty, *i.e.*, the cost of accessing something that is not close to CPU. As a result,

- Large page size, typically 4KB: increase the hit rate
- Fully associative: Any VP can be placed in any PP
- Highly sophisticated, expensive replacement algorithms
- Write-back rather than write-through

#### 3.3.3 Enabling data structure: Page table

**Replacement.** If there's a miss, a system determines where the VP is stored on disk, select a victim page in physical memory, and copy the VP from disk to DRAM, replacing the victim page. This is provided by a combination of OS software, MMU, and data structure stored in physical memory known as a *page table* that maps virtual pages to physical pages.

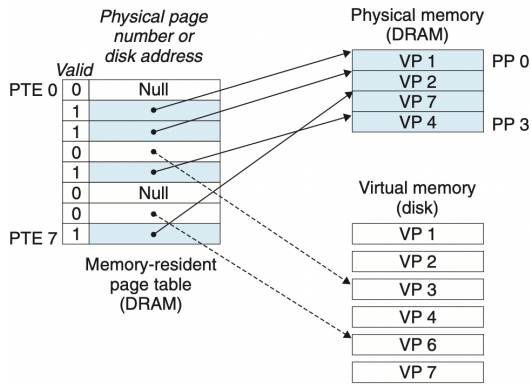


Figure 23. Page table

A *page table* is an array of *page table entries* (PTEs). Each page in the virtual address space has a PTE at a fixed offset in the page table. Valid bit indicates whether the VP is cached in DRAM:

- 1 (cached): the address indicates the start of the corresponding PP in DRAM where the VP is cached
- 0 (unallocated): null address indicates that VP not allocated
- 0 (uncached): the address points to the start of the VP on disk

### 3.3.4 Handling page faults

- Page hit: reference to VM word that is in physical memory
- Page fault: reference to VM word that is not in physical memory

Page fault is also an exception and handled by:

1. Page fault handler selects a victim to be evicted.
2. If victim is modified, the kernel copies it back to disk.
3. Update victim PTE: point virtual memory, change valid bit.
4. Kernel copies referenced page from disk to physical memory, and update PTE: point physical memory, change valid bit
5. Handler returns. It restarts the faulting instruction, which re-sends the faulting virtual address to the address translation hardware. But now it is cached in main memory, so page hit.

Waiting until the miss to copy the page to DRAM is known as *demand paging*.

### 3.3.5 Allocating pages

When OS allocates new page of VM, *e.g.*, as a result of `malloc`, VP is allocated by creating room on disk and updating PTE to point the created page on disk.

### 3.3.6 Locality to the rescue!

Will the large miss penalties destroy the program performance? Although the total number of pages that program references during an entire run might exceed the physical memory, the program will tend to work on smaller set of *active pages* known as the *working set* or *resident set*. After an initial overhead where the working set is paged to the memory, subsequent references to the working set will result in hits. However, if the working set size exceeds the size of physical memory, then the program will cause *thrashing*, where pages are swapped in and out continuously.

## 3.4 VM as a tool for memory management

OS provide a separate page table, and thus a separate virtual address space, for each process. Note that multiple VPs can be mapped to the same shared PP. The combination of demand paging and separate virtual address spaces has a profound impact on memory management:

- **Simplifying linking.** Each program has similar virtual address space. Code, data, and heap always start at the same addresses
- **Simplify loading**

## 3.5 VM as a tool for memory protection

PTE can be extended with additional permission bits, so that MMU can check these bits on each access.

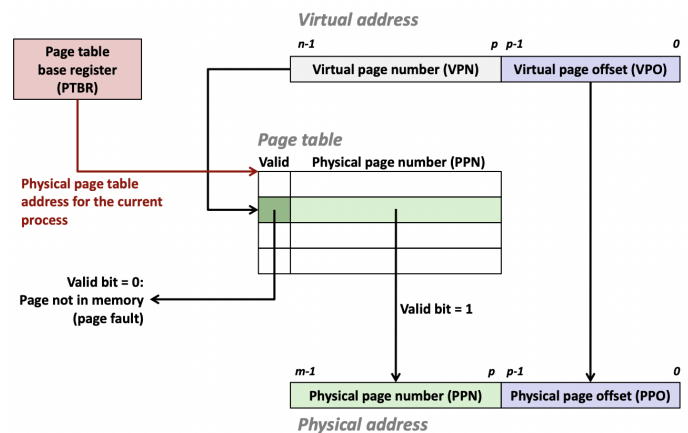
## 3.6 Address translation

Say virtual address space  $V = \{0, \dots, N - 1\}$ , physical address space  $P = \{0, \dots, M - 1\}$ , address translation is the mapping  $MAP : V \rightarrow P \cup \{\emptyset\}$ .

### 3.6.1 Address translation with a page table

The  $n$ -bit virtual address has two components: a  $p$ -bit VP offset (VPO) and an  $(n - p)$ -bit VP number (VPN). A control register in CPU, page table base register (PTBR) points to the current page table. The MMU adds VPN to PTBR to get PTE address (PTEA).

- If page hit, the corresponding physical address (PA) is the concatenation of the PP number (PPN) from PTE and VPO.
- If page fault, *i.e.*, valid bit is zero, MMU triggers page fault exception. See 3.3.4.



### 3.6.2 Integrating VM and cache

Consider L1 cache in CPU chip. If there's PTEA or PA hit, get it from L1 cache. If miss, get it from memory, and store in L1 cache.

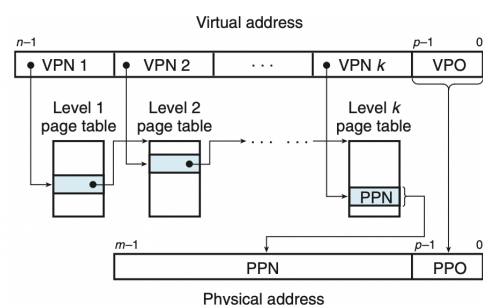
### 3.6.3 Speeding up translation with a TLB

Let's eliminate even the cost looking L1 cache by including a small cache of PTEs in MMU called a translation lookaside buffer (TLB). TLB maps VPN to PPN instead of storing a whole PTE. VPN is consisted of  $t$ -bits TLB index (TLBI) and TLB tag (TLBT). TLB is consisted of multiple sets containing the tag and PTE. TLBI decides the set and compare TLBT with the tags in the set.

### 3.6.4 Multi-level page tables

Suppose a 4KB ( $2^{12}$ ) page size, 48-bit address space, and 8-byte PTE. Then we will need  $2^{48-12} \times 2^3 = 2^{39}$  bytes!

→ Multi-level page table. Consequently, multiple VPNs. For example, level-1 table points to a level-2 table. Level-2 table points to a page. Since most of the level-2 table would be empty, most of the level-1 table would be null. So, most of the level-2 tables are not allocated. Usually the number of levels is 4.



## 4 Malloc

### 4.1 Dynamic memory allocation

Dynamic memory allocators manage an area of process VM known as heap. Allocator maintains heap as collection of variable sized blocks, which are either *allocated* or *free*.

- **Explicit allocator:** application allocates and frees space
- **Implicit allocator:** application allocates, but does not free space, e.g., garbage collection in Java

#### 4.1.1 The malloc package

void \*malloc(size\_t size)

- Successful: Returns a pointer to a memory block of at least size bytes aligned to 16-byte (x86-64). If size==0, returns NULL
- Unsuccessful: Returns NULL and sets errno

void free(void \*p)

- Returns the block pointed at by p
- p must come from a previous call to malloc or realloc

Other functions are:

- realloc: changes the size of a previously allocated block
- sbrk: used internally by allocators to grow or shrink the heap

### 4.2 Performance goals

- Handling arbitrary request sequences
- Making immediate responses to requests
- Using only the heap
- Aligning blocks (alignment requirement)
- Not modifying allocated blocks

Within the constraints, an allocator attempts to meet the often conflicting performance goals:

- **Throughput:** number of request it completes per unit time
- **Memory utilization:** peak utilization  $U_k = \frac{\max_{i \leq k} P_i}{H_k}$  where  $H_k$  denote the current size of the heap and  $P_i$  denote the sum of payloads of currently allocated blocks

#### 4.2.1 Fragmentation

Poor memory utilization is caused by *fragmentation*.

#### 4.2.2 Internal fragmentation

Occurs if payload is smaller than block size. It is caused by

- Overhead of maintaining heap data structures
- Padding for alignment
- Explicit policy decisions

It depends on pattern of previous requests, so is easy to measure.

#### 4.2.3 External fragmentation

Occurs when there is enough heap memory, but no single free block is large enough. It depends on pattern of future requests, so is hard to measure.

### 4.3 Implementation

#### 4.3.1 Knowing how much to free

Use *header/footer* to keep the length of the block.

#### 4.3.2 Keeping track of free blocks

- Implicit list using length: links all blocks
- Explicit list: links free blocks using pointers
- Segregated free list: different free lists for different size classes
- Block sorted by size: balanced tree

### 4.4 Implicit list

Header and footer (boundary tags) indicates size and tag. Footer allows a bidirectional coalescing.

#### 4.4.1 Finding a free block (placement policy)

- **First fit:** Search list from beginning, choose first that fits
  - Can take linear time in total number of blocks
  - Cause splinters at beginning of list
- **Next fit:** First fit, search starts where previous search finished
  - Faster than first fit: avoid re-scanning
  - Fragmentation is worse
- **Best fit:** fits with the fewest bytes left over
  - Keep fragments small: higher memory utilization
  - Typically slower than first fit

#### 4.4.2 Allocating in free block

Split the remaining space when allocating free block.

#### 4.4.3 Freeing a block

Only clearing tag leads to false fragmentation → Coalescing!

- Immediate coalescing: coalesce each time freeing
- Deferred coalescing: coalesce when the amount of external fragmentation reaches some threshold

### 4.5 Explicit list

Include next, prev pointers after header tag. Maintain the list of free blocks. Only free blocks, so we can use payload area.

#### 4.5.1 Insertion policy

Where in the free list do you put a newly freed block?

- LIFO: beginning of the free list
  - Pro: simple and constant time
  - Con: fragmentation
- Address-ordered policy: blocks are always in address order
  - Pro: lower fragmentation
  - Con: require search

### 4.6 Segregated free list

Each size class of blocks has its own free list. Separate classes for each small size, two-power size classes for larger sizes.

- Higher throughput
- Better memory utilization: First-fit search of segregated free list approximates a best-fit search of entire heap.

#### 4.6.1 Allocating

1. Search appropriate free list
2. If an appropriate block is found: split block and place fragmentation on appropriate list
3. If no block is found, try next larger class

If no block is found, request heap memory from OS using sbrk(). Extend chunk size, and put remaining free block to appropriate list.

#### 4.6.2 Freeing

Coalesce and place on appropriate list.

Cost	Implicit	Explicit	Segregated
Allocate	Linear to all	Linear to free	Log time
Free	Constant	Constant	Constant
Memory	Depends	Better	Best

Table 7. Comparison on allocators

## 5 System-level I/O

### 5.1 Unix I/O

A Linux file is a sequence of bytes. All I/O devices, even the kernel, are represented as files. Mapping of files to devices allows kernel to export simple interface called Unix I/O:

- Opening and closing files: `open()`, `close()`
- Reading and writing a file: `read()`, `write()`
- Changing the current file position (`seek`): `lseek()`

#### 5.1.1 File types

Each file has a type indicating its role in the system.

- **Regular file:** Contains arbitrary data
  - Text files are with only ASCII or Unicode character. Text file is a sequence of text lines, terminated by newline char (0xa).
  - Binary files are everything else.
  - Kernel doesn't know the difference!
- **Directory:** Index for a related group of files
  - Consists of an array of links mapping a filename to a file
  - Contains at least two entries: link to itself (`.`), and link to the parent directory (`..`)
- **Socket:** For communicating with a process on another machine

#### 5.1.2 Directory hierarchy

All files are organized as a hierarchy anchored by root directory named `/` (slash). Kernel maintains *current working directory* (`cwd`) for each process. Locations of files in hierarchy are denoted by pathnames.

- Absolute: path from root and starts with `/`
- Relative: path from `cwd`

#### 5.1.3 Opening and closing files

Opening/closing a file informs the kernel that you are ready/finished accessing that file. `open` returns a identifying integer *file descriptor*. `fd==-1` indicates that an error occurred.

```
1 int fd; /* file descriptor */
2 if ((fd = open("/etc/hosts", O_RDONLY)) < 0) exit(1);
3 int retval; /* return value */
4 if ((retval = close(fd)) < 0) exit(1);
```

Each process created by a Linux shell begins life with three open files associated with a terminal:

- 0: standard input (`stdin`)
- 1: standard output (`stdout`)
- 2: standard error (`stderr`)

#### 5.1.4 Reading and writing files

Reading/writing a file copies bytes CFP→memory/memory→CFP, and then updates file position. Writing a file copies bytes from memory to the current file position, and then updates file position. Returns a number of bytes read/written. `nbytes<0` indicates that an error occurred.

```
1 char buf[512];
2 int fd;
3 int nbytes;
4 if ((nbytes = read(fd, buf, sizeof(buf))) < 0) exit(1);
5 if ((nbytes = write(fd, buf, sizeof(buf))) < 0) exit(1);
```

*Short counts* (`nbytes < sizeof(buf)`) are possible and are not errors. Short count can occur when encountering EOF on reads, reading from terminal, or reading and writing network sockets. But never occur when reading and writing from disk files.

#### 5.1.5 Example

```
1 int main(void) {
2     char c;
3     while(Read(STDIN_FILENO, &c, 1) != 0)
4         Write(STDOUT_FILENO, &c, 1);
5     exit(0);
6 }
```

### 5.2 Metadata

*Metadata* is data about data. Per-file metadata maintained by kernel. Users can access with the `stat` and `fstat` functions.

### 5.3 File sharing

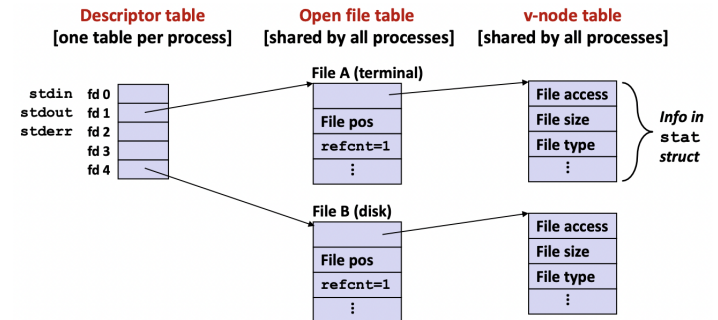


Figure 24. How the kernel represents open files

File sharing is two distinct descriptors sharing the same disk file through two **distinct** open file table entries, e.g., by `open` twice.

Processes share files using **same** open file table entries. A child process inherits parent's open files. After fork, child process's descriptor table is same as parent's, and `refcnts` are incremented.

#### 5.3.1 I/O redirection

Linux shells provide I/O redirection operators that allow users to associate standard input and output with disk files. For example, `ls > foo.txt` causes the shell to load and execute the `ls` program, with standard output redirected to disk file `foo.txt`.

Shell implement it by calling `dup2(oldfd, newfd)` function. It copies per-process descriptor table entry `oldfd` to entry `newfd`.

1. Open file to which `stdout` should be redirected. Happens in child executing shell code, before `exec`.
2. Call `dup2`. Change `refcnt` accordingly.

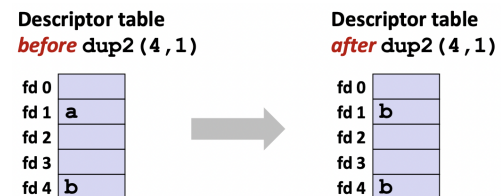


Figure 25. I/O redirection

### 5.4 Standard I/O

C standard library `libc.so` contains a collection of higher-level standard I/O functions.

- Opening and closing files: `fopen`, `fclose`
- Reading and writing bytes: `fread`, `fwrite`
- Reading and writing text lines: `fgets`, `fputs`
- Formatted reading and writing: `fscanf`, `fprintf`

#### 5.4.1 Streams

Standard I/O models open files as *streams*, abstraction for a file descriptor and a buffer in memory. C programs begin life with three open streams: `stdin`, `stdout`, and `stderr`.

	Unix I/O	Standard I/O
Pros	Most general and lowest overhead Provides functions for accessing metadata Async-signal-safe and can be used safely in signal handlers	Buffering increases efficiency by reducing system calls Short counts automatically handled
Cons	Dealing with short counts is tricky No buffering → Lower efficiency	No functions for accessing metadata Not async-signal-safe, not appropriate for signal handlers Not appropriate for I/O in network sockets

Table 8. Comparison on I/O

```

1 #include <stdio.h>
2 extern FILE *stdin;
3 extern FILE *stdout;
4 extern FILE *stderr;
5 int main() {
6     fprintf(stdout, "Hello, world\n");
7 }
    
```

### 5.4.2 Buffered I/O

Applications often read/write one character at a time. Implementing as Unix I/O calls is expensive. Solution is *buffered read*.

- Use Unix read to grab block of bytes
- User input functions take one byte at a time from buffer
- Refill buffer when empty

Buffer is flushed to output fd on newline char, call to fflush, exit, or return from main.

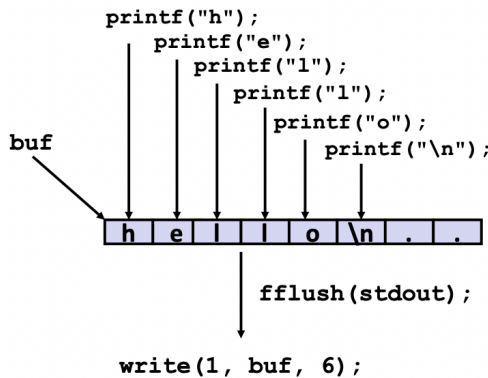
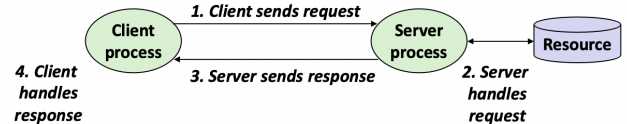


Figure 26. Buffered write

## 6 Network programming

### 6.1 A client-server transaction



### 6.2 Computer networks

A *network* is a hierarchical system of boxes and wires organized by geographical proximity.

- LAN (local area network): spans a building or campus
- WAN (wide area network): spans country or world

An *internetwork* (internet) is an interconnected set of networks. The Global IP Internet (uppercase “I”) is the unique and united form of an internet (lowercase “i”).

#### 6.2.1 Internet protocol

Protocol is set of rules of how hosts and routers should cooperate when they transfer data from network to network. Internet protocol defines:

- **Naming scheme.** Defines a uniform format of *host address*.
- **Delivery mechanism.** Defines a standard transfer unit (*packet*) consisting of *header* and *payload*.
  - Header: contains info such as packet size, source and destination address
  - Payload: contains data bits sent from source host

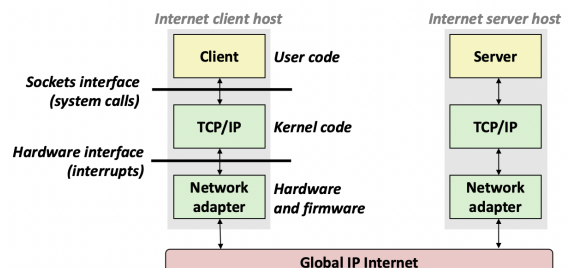
#### 6.2.2 Global IP Internet

Based on TCP/IP protocol family.

- IP (internet protocol): Provides basic naming scheme and unreliable delivery of packets (datagrams) from host-to-host
- UDP (unreliable datagram protocol): Uses IP to provide unreliable datagram delivery from process-to-process
- TCP (transmission control protocol): Uses IP to provide reliable byte streams from process-to-process over connections

Via a mix of Unix file I/O and functions from sockets interface.

#### 6.2.3 Hardware and software organization of an internet



### 6.3 Programmer's view of the Internet

#### 6.3.1 IP address

Hosts are mapped to a set of 32-bit IP addresses

- Stored in an IP address struct
- Always stored in memory in network byte order (big-endian)
- Dotted decimal notation: each byte in address is represented by its decimal value and separated by a period

#### 6.3.2 Domain naming system (DNS)

The set of IP addresses is mapped to a set of identifiers called Internet domain names in DNS.

- Each host has a locally defined domain name *localhost* which always maps to the *loopback address* 127.0.0.1
- Use *hostname* to determine real domain name of local host
- Mapping is many-to-many

#### 6.3.3 Connection

Clients and servers communicate by sending streams of bytes over *connections*. Each connection is:

- Point-to-point: connects a pair of processes
- Full-duplex: data can flow in both directions at the same time
- Reliable: stream of bytes sent by the source is eventually received by the destination in the same order it was sent

A *socket* is endpoint of a connection. Socket address is an IP address:port pair. A *port* is a 16-bit integer that identifies a process.

- Ephemeral port: Assigned automatically by client kernel when client makes a connection request.
- Well-known port: Associated with service provided by a server

Popular services have permanently assigned well-known ports and corresponding well-known service names:

- echo server: 7/echo
- ssh servers: 22/ssh
- email server: 25/smtp
- Web servers: 80/http

A connection is uniquely identified by the socket addresses of its endpoints (*socket pair*).

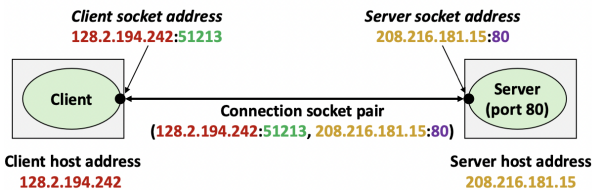


Figure 27. 51213 is an ephemeral port allocated by the kernel.

### 6.4 Socket

To the kernel, a socket is an endpoint of communication. To an application, a socket is a file descriptor that lets the application read/write from/to the network (networks are modeled as files!). Clients and servers communicate with each other by reading from and writing to socket descriptors. The main distinction between regular file I/O and socket I/O is how the application *opens* the socket descriptors.

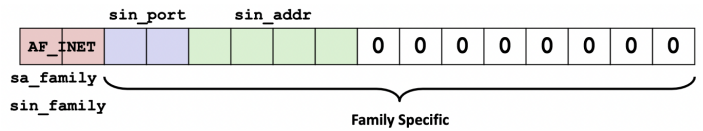
#### 6.4.1 Socket address structures

Generic socket address:

```
1 struct sockaddr {
2     uint16_t sa_family; /* Protocol family */
3     char sa_data[14]; /* Address data. */
4 };
```

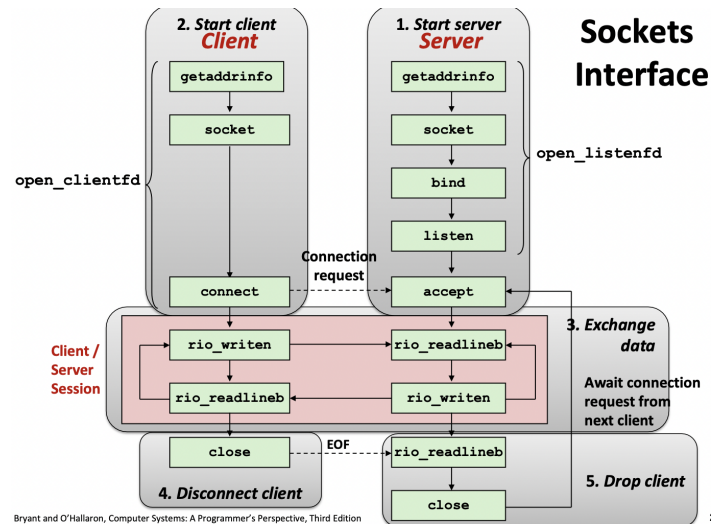
Internet-specific socket address:

```
1 struct sockaddr_in {
2     uint16_t sin_family; /* Protocol family (always AF_INET) */
3     uint16_t sin_port; /* Port num in network byte order */
4     struct in_addr sin_addr; /* IP addr in network byte order */
5     unsigned char sin_zero[8]; /* Pad to sizeof(struct sockaddr) */
6 };
```



Must cast to (struct sockaddr \*) for functions that take socket address arguments.

### 6.5 Sockets interface



#### 6.5.1 getaddrinfo

Convert strings of hostnames, host addresses, ports, and service names to socket address structures.

```
1 int getaddrinfo(const char *host, // Hostname or address
2               const char *service, // Port or service name
3               const struct addrinfo *hints, // Input params
4               struct addrinfo **result); // Output linked list
```

Given host and service, *getaddrinfo* returns result that points to a linked list of *addrinfo* structs, each of which points to a corresponding socket address struct, and which contains arguments for the sockets interface functions.

Each *addrinfo* struct returned by *getaddrinfo* contains arguments that can be passed directly to *socket* function. Also points to a socket address struct that can be passed directly to *connect* and *bind* functions.

Clients walk the list, trying each socket address in turn, until the calls to *socket* and *connect* succeed. Servers walk the list until calls to *socket* and *bind* succeed.

#### 6.5.2 getnameinfo

*getnameinfo* is the inverse of *getaddrinfo*, converting a socket address to the corresponding host and service.

### 6.5.3 socket

```
int socket(int domain, int type, int protocol)
```

Clients and servers use the socket function to create a socket descriptor. It is protocol specific. Use `getaddrinfo` to generate the parameters automatically, so that code is protocol independent.

### 6.5.4 bind

```
int bind(int sockfd, SA *addr, socklen_t addrlen);
```

A server uses `bind` to ask the kernel to associate the server's socket address with a socket descriptor. Use `getaddrinfo` to supply parameters.

### 6.5.5 listen

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

By default, kernel assumes that descriptor from socket function is an *active socket* that will be on the client end of a connection. A server calls the `listen` function to tell the kernel that a descriptor will be used by a server rather than a client. It converts `sockfd` from an active socket to a listening socket that can accept connection requests from clients.

### 6.5.6 accept

```
int accept(int listenfd, SA *addr, int *addrlen);
```

Servers wait for connection requests from clients by calling `accept`. Returns a *connected descriptor* that can be used to communicate with the client via Unix I/O routines.

### 6.5.7 connect

```
int connect(int sockfd, SA *addr, socklen_t addrlen);
```

A client establishes a connection with a server by calling `connect`. If successful, then `clientfd` is now ready for reading and writing.

### 6.5.8 Listening vs. connected descriptors

- Listening: Created once and exists for lifetime of the server. End point for client connection requests.
- Connected: End point of the connection between client and server
- New descriptor is created each time the server accepts a connection request from a client

### 6.5.9 Example

#### Helpers

```
1 int open_clientfd(char *hostname, char *port) {
2     int clientfd;
3     struct addrinfo hints, *listp, *p;
4     // Get a list of potential server addresses
5     memset(&hints, 0, sizeof(struct addrinfo));
6     hints.ai_socktype = SOCK_STREAM; // Open a connection
7     hints.ai_flags = AI_NUMERICSERV; // using numeric port
8     hints.ai_flags |= AI_ADDRCONFIG;
9     Getaddrinfo(hostname, port, &hints, &listp);
10    // Walk the list for one that we can successfully connect to
11    for (p = listp; p; p = p->ai_next) {
12        /* Create a socket descriptor */
13        if ((clientfd = socket(p->ai_family, p->ai_socktype, p->
14        ai_protocol)) < 0) continue; // failed, try next
15        /* Connect to the server */
16        if (connect(clientfd, p->ai_addr, p->ai_addrlen) != -1)
17            break; // Success
18        Close(clientfd); // Connect failed, try another
19    }
20    Freeaddrinfo(listp); // Clean up
21    if (!p) return -1;
22    else return clientfd;
23 }
24 int open_listenfd(char *port) {
25     struct addrinfo hints, *listp, *p;
26     int listenfd, optval=1;
```

```
26     /* Get a list of potential server addresses */
27     memset(&hints, 0, sizeof(struct addrinfo));
28     hints.ai_socktype = SOCK_STREAM; // Accept connect.
29     hints.ai_flags = AI_PASSIVE | AI_ADDRCONFIG; // any addr
30     hints.ai_flags |= AI_NUMERICSERV; // using port no.
31     Getaddrinfo(NULL, port, &hints, &listp);
32     /* Walk the list for one that we can bind to */
33     for (p = listp; p; p = p->ai_next) {
34         /* Create a socket descriptor */
35         if ((listenfd = socket(p->ai_family, p->ai_socktype, p->
36         ai_protocol)) < 0) continue; // failed, try next
37         // Eliminates "Address already in use" error from bind
38         Setsockopt(listenfd, SOL_SOCKET, SO_REUSEADDR, (const
39         void *)&optval, sizeof(int));
40         /* Bind the descriptor to the address */
41         if (bind(listenfd, p->ai_addr, p->ai_addrlen) == 0)
42             break; // Success */
43         Close(listenfd); // Bind failed, try the next */
44     }
45     Freeaddrinfo(listp); // Clean up
46     if (!p) return -1;
47     // Make it a listening socket ready to accept conn. requests
48     if (listen(listenfd, LISTENQ) < 0) {
49         Close(listenfd); return -1; }
50     return listenfd;
51 }
```

#### echoclient.c

```
1 int main(int argc, char **argv) {
2     int clientfd;
3     char *host, *port, buf[MAXLINE];
4     rio_t rio;
5     host = argv[1];
6     port = argv[2];
7     clientfd = Open_clientfd(host, port);
8     Rio_readinitb(&rio, clientfd);
9     while (Fgets(buf, MAXLINE, stdin) != NULL) {
10        Rio_writen(clientfd, buf, strlen(buf));
11        Rio_readlineb(&rio, buf, MAXLINE);
12        Fputs(buf, stdout);
13    }
14    Close(clientfd);
15    exit(0);
16 }
```

#### echoserver.c

```
1 void echo(int connfd);
2 int main(int argc, char **argv) {
3     int listenfd, connfd;
4     socklen_t clientlen;
5     struct sockaddr_storage clientaddr; // Enough room for any
6     char client_hostname[MAXLINE], client_port[MAXLINE];
7     listenfd = Open_listenfd(argv[1]);
8     while (1) {
9         clientlen = sizeof(struct sockaddr_storage);
10        connfd = Accept(listenfd, (SA*)&clientaddr, &clientlen);
11        Getnameinfo((SA*)&clientaddr, clientlen,
12        client_hostname, MAXLINE, client_port, MAXLINE, 0);
13        printf("Connected to (%s, %s)\n", client_hostname,
14        client_port);
15        echo(connfd);
16        Close(connfd);
17    }
18    exit(0);
19 }
```

#### echo.c

```
1 void echo(int connfd) {
2     size_t n;
3     char buf[MAXLINE];
4     rio_t rio;
5     Rio_readinitb(&rio, connfd);
6     while((n = Rio_readlineb(&rio, buf, MAXLINE)) != 0) {
7         printf("server received %d bytes\n", (int)n);
8         Rio_writen(connfd, buf, n);
9     }
10 }
```

### 6.5.10 Testing servers using telnet

```
linux> telnet <host> <portnumber>
```

The telnet program is invaluable for testing servers that transmit ASCII strings over Internet connections

## 6.6 Web servers

Clients and servers communicate using the HyperText Transfer Protocol (HTTP).

1. Client and server establish TCP connection
2. Client requests content
3. Server responds with requested content
4. Client and server close connection (eventually)

### 6.6.1 Web content

Web servers return content to clients, a sequence of bytes with an associated MIME (Multipurpose Internet Mail Extensions) type, e.g., HTML, and PNG. The content returned in HTTP responses can be either static or dynamic.

- Static: stored in files and retrieved in response to HTTP request
- Dynamic: produced on-the-fly in response to HTTP request

### 6.6.2 URLs

Unique name for a file: URL (Universal Resource Locator).

```
http://www.cmu.edu:80/index.html
```

- Clients use prefix (`http://www.cmu.edu:80`) to infer protocol, where the server is, and port
- Servers use suffix (`/index.html`) to determine if request is for static or dynamic content

### 6.6.3 HTTP requests

HTTP request is a request line followed by zero or more request headers.

Request line: `<method> <uri> <version>`

- `<method>` is one of GET, POST, OPTIONS, HEAD, PUT, DELETE, or TRACE
- `<uri>` is typically URL for proxies, URL suffix for servers
- `<version>` is HTTP version of request

Request headers: `<header name>: <header data>`

### 6.6.4 HTTP responses

HTTP response is a response line followed by zero or more response headers, possibly followed by content, with blank line separating headers from content.

Response line: `<version> <status code> <status msg>`

- `<version>` is HTTP version of the response
- `<status code>` is numeric status
- `<status msg>` is text: OK, Moved, or Not found

Response headers: `<header name>: <header data>`

## 6.7 Tiny web server

### 6.7.1 Tiny operations

- Accept connection from client
- Read request from client (via connected socket)
- Split into `<method> <uri> <version>`
- If `<method>` is not GET, then return error
- If URI contains `cgi-bin` then serve dynamic content
- Otherwise serve static content

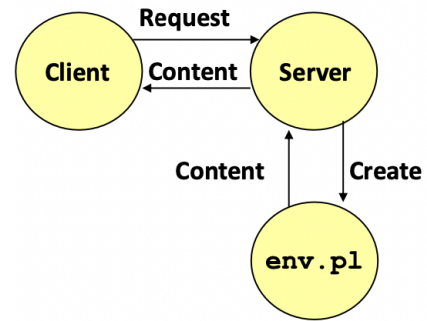


Figure 28. Serving dynamic content

### 6.7.2 Serving static content

```

1 void serve_static(int fd, char *filename, int filesize)
2 {
3     int srcfd;
4     char *srcp, filetype[MAXLINE], buf[MAXBUF];
5     /* Send response headers to client */
6     get_filetype(filename, filetype);
7     sprintf(buf, "HTTP/1.0 200 OK\r\n");
8     sprintf(buf, "%sServer: Tiny Web Server\r\n", buf);
9     sprintf(buf, "%sConnection: close\r\n", buf);
10    sprintf(buf, "%sContent-length: %d\r\n", buf,
11            filesize);
12    sprintf(buf, "%sContent-type: %s\r\n\r\n", buf,
13            filetype);
14    Rio_writen(fd, buf, strlen(buf));
15    /* Send response body to client */
16    srcfd = Open(filename, O_RDONLY, 0);
17    srcp = Mmap(0, filesize, PROT_READ, MAP_PRIVATE,
18               srcfd, 0);
19    Close(srcfd);
20    Rio_writen(fd, srcp, filesize);
21    Munmap(srcp, filesize);
22 }
  
```

### 6.7.3 Serving dynamic content

1. The server creates a child process and runs the program identified by the URI in that process
2. The child runs and generates the dynamic content
3. The server captures the content of the child and forwards it without modification to the client

Common Gateway Interface (CGI) defines a simple standard for transferring information between the client (browser), the server, and the child process. Because the children are written according to the CGI spec, they are often called CGI programs.

```
http://add.com/cgi-bin/adder?15213&18213
```

- `adder` is the CGI program on the server that will do the addition
- argument list starts with `?`
- arguments separated by `&`
- spaces represented by `+` or `%20`

Server pass the arguments to child in environment variable `QUERY_STRING`. For example, `QUERY_STRING = "15213&18213"`

Child generates output on `stdout`. Server uses `dup2` to redirect it to its connected socket.

```

1 void serve_dynamic(int fd, char *filename, char *cgiargs)
2 {
3     char buf[MAXLINE], *emptylist[] = { NULL };
4     /* Return first part of HTTP response */
5     sprintf(buf, "HTTP/1.0 200 OK\r\n");
6     Rio_writen(fd, buf, strlen(buf));
7     sprintf(buf, "Server: Tiny Web Server\r\n");
  
```



```

7  Rio_writen(fd, buf, strlen(buf));
8  if (Fork() == 0) { /* Child */
9      /* Real server would set all CGI vars here */
10     setenv("QUERY_STRING", cgiargs, 1);
11     Dup2(fd, STDOUT_FILENO); /* Redirect stdout to
client */
12     Execve(filename, emptylist, environ); /* Run CGI
program */
13 }
14 Wait(NULL); /* Parent waits for and reaps child */
15 }
    
```

```

bash:makoshark> telnet whaleshark.ics.cs.cmu.edu 15213
Trying 128.2.210.175...
Connected to whaleshark.ics.cs.cmu.edu (128.2.210.175).
Escape character is '^]'.
GET /cgi-bin/adder?15213&18213 HTTP/1.0
-----
HTTP/1.0 200 OK
Server: Tiny Web Server
Connection: close
Content-length: 117
Content-type: text/html
-----
Welcome to add.com: THE Internet addition portal.
<p>The answer is: 15213 + 18213 = 33426
<p>Thanks for visiting!
-----
Connection closed by foreign host.
bash:makoshark>
    
```

Figure 29. Dynamic content result

```

1  int main(int argc, char **argv) {
2      int listenfd, connfd;
3      int port = atoi(argv[1]);
4      struct sockaddr_in clientaddr;
5      int clientlen=sizeof(clientaddr);
6      signal(SIGCHLD, sigchld_handler);
7      listenfd = open_listenfd(port);
8      while (1) {
9          connfd = accept(listenfd, (SA *) &clientaddr, &
clientlen);
10         if (fork() == 0) {
11             close(listenfd); // Child closes its listening
socket
12         }
13         echo(connfd); // Child services client
14         close(connfd); // Child closes connection with
client
15         exit(0); // Child exits
16         Close(connfd); // Parent closes connected socket
(important!)
17     }
18 }
19 /* Listening server process must reap zombie children to
avoid fatal memory leak */
20 void sigchld_handler(int sig) {
21     while (waitpid(-1, 0, WNOHANG) > 0) ;
22     return;
23 }
    
```

## 7 Concurrent programming

### 7.1 Iterative servers

Iterative servers process one request at a time.

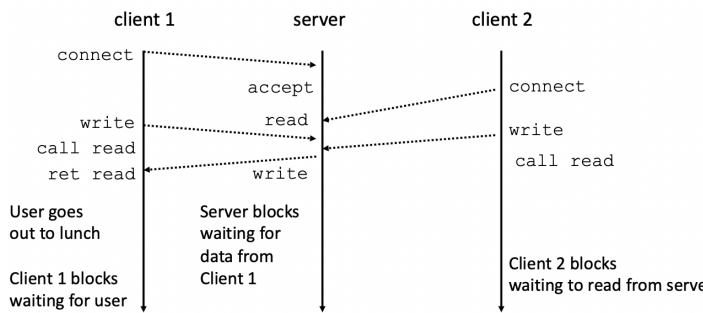
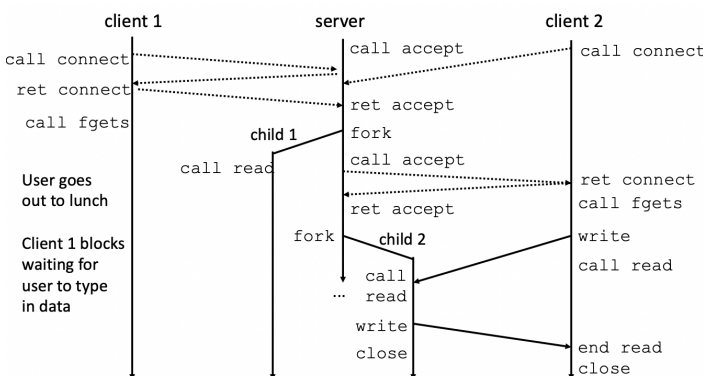


Figure 30. Fundamental flaw of iterative servers. Server waits until client 1 is closed.

Solution: use concurrent servers instead! Allow server to handle multiple clients simultaneously.

### 7.2 Process-based concurrent server

Spawn separate process for each client.



- Kernel automatically interleaves multiple logical flows
- Each flow has its own private address space

### 7.2.1 Pros and cons

- + Handle multiple connections concurrently
- + Clean sharing model
- + Simple and straightforward
- – Additional overhead for process control
- – Nontrivial to share data between processes

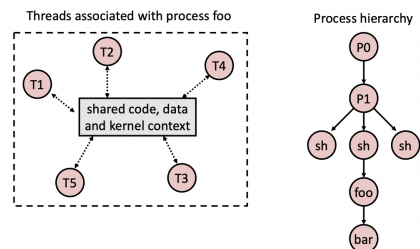
### 7.3 Thread-based concurrent server

Multiple threads can be associated with a process.

- Each thread has its own logical control flow
- Each thread shares the same code, data, and kernel context
- Share common virtual address space
- Each thread has its own thread id (TID)

Two threads are (logically) concurrent if their flows overlap in time. Otherwise, they are sequential. True concurrency is only possible in multi-core processor.

#### 7.3.1 Threads vs. Processes



Threads share code and some data. Processes (typically) do not. Threads are less expensive than processes.

#### 7.3.2 Posix threads (Pthreads) interface

Standard interface for functions that manipulate threads from C.

- Creating and reaping threads: pthread\_create(), pthread\_join()
- Determining your thread ID: pthread\_self()
- Terminating threads: pthread\_cancel(), pthread\_exit(), exit() (terminates all threads), RET (terminates current thread)

```

1 void *thread(void *vargp) {
2     printf("Hello, world!\n");
3     return NULL;
4 }
5 int main() {
6     pthread_t tid;
7     Pthread_create(&tid, NULL, thread, NULL);
8     Pthread_join(tid, NULL);
9     exit(0);
10 }
    
```

## 8 Synchronization

### 8.1 Sharing

Say variable  $x$  is shared if and only if multiple threads reference instance of  $x$ . Which variables in threaded C program are shared?

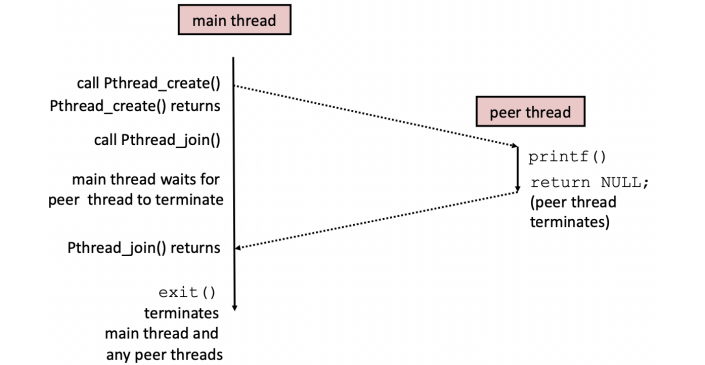
#### 8.1.1 Threads memory model: Conceptual model

- Multiple threads run within the same context of a single process
- Each thread has its own separate thread context

This model is not strictly enforced. Any thread can read and write the stack of any other thread.

#### 8.1.2 Mapping variable instances to memory

- Global/Local static variable: VM contains exactly one instance
- Local variables: Each thread stack contains one instance



Spawn new thread for each client.

```

1 int main(int argc, char **argv) {
2     int port = atoi(argv[1]);
3     struct sockaddr_in clientaddr;
4     int clientlen = sizeof(clientaddr); pthread_t tid;
5     int listenfd = open_listenfd(port);
6     while (1) {
7         int *conncfd = malloc(sizeof(int));
8         *conncfd = accept(listenfd, (SA *) &clientaddr, &clientlen);
9         pthread_create(&tid, NULL, echo_thread, conncfd);
10    }
11 }
12 void *echo_thread(void *vargp) {
13     int conncfd = *((int *)vargp); pthread_detach(pthread_self());
14     free(vargp);
15     echo(conncfd);
16     close(conncfd);
17     return NULL;
18 }
    
```

Global var: 1 instance (ptr [data])

Local vars: 1 instance (i.m, msgs.m)

Local var: 2 instances (myid.p0 [peer thread 0's stack], myid.p1 [peer thread 1's stack])

Local static var: 1 instance (cnt [data])

```

char **ptr; /* global */

int main()
{
    int i;
    pthread_t tid;
    char *msgs[2] = {
        "Hello from foo",
        "Hello from bar"
    };
    ptr = msgs;

    for (i = 0; i < 2; i++)
        pthread_create(&tid,
            NULL,
            thread,
            (void *)i);
    pthread_exit(NULL);
}

/* thread routine */
void *thread(void *vargp)
{
    int myid = (int)vargp;
    static int cnt = 0;

    printf("[%d]: %s (svar=%d)\n",
        myid, ptr[myid], ++cnt);
}
    
```

Variable instance	Referenced by main thread?	Referenced by peer thread 0?	Referenced by peer thread 1?
ptr	yes	yes	yes
cnt	yes	yes	yes
i.m	yes	no	no
msgs.m	yes	yes	yes
myid.p0	no	yes	no
myid.p1	no	no	yes

ptr, cnt, msgs are shared but i, myid are not shared.

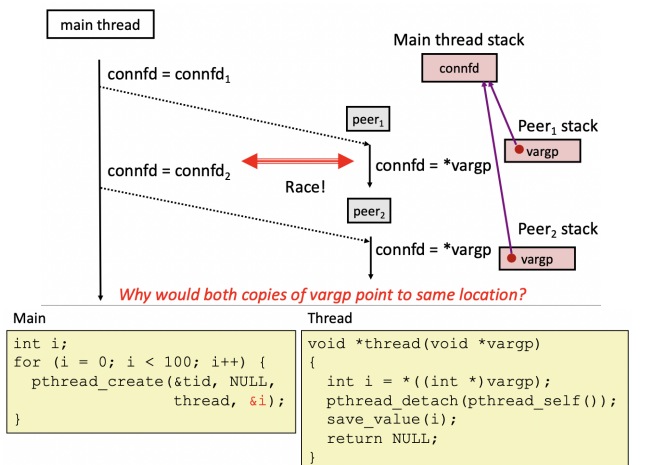
### 8.2 Mutual exclusion

#### 8.2.1 Improper synchronization example

### 7.3.3 Pros and cons

- + Easy to share data structures between threads
- + Threads are more efficient than processes
- – Unintentional sharing

### 7.3.4 Unintended sharing



```

volatile int cnt = 0; /* global */

int main(int argc, char **argv)
{
    int niters = atoi(argv[1]);
    pthread_t tid1, tid2;

    pthread_create(&tid1, NULL,
        thread, &niters);
    pthread_create(&tid2, NULL,
        thread, &niters);
    pthread_join(tid1, NULL);
    pthread_join(tid2, NULL);

    /* Check result */
    if (cnt != (2 * niters))
        printf("BOOM! cnt=%d\n", cnt);
    else
        printf("OK cnt=%d\n", cnt);
    exit(0);
}

/* Thread routine */
void *thread(void *vargp)
{
    int i, niters = *((int *)vargp);

    for (i = 0; i < niters; i++)
        cnt++;

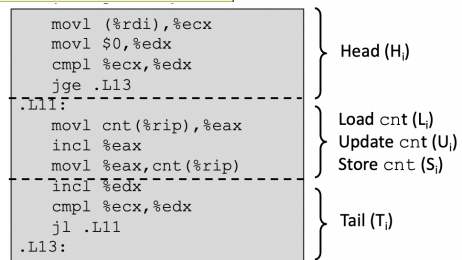
    return NULL;
}
    
```

```

linux> ./badcnt 10000
OK cnt=20000
linux> ./badcnt 10000
BOOM! cnt=13051
linux>
    
```

cnt should equal 20,000.

What went wrong?



Incorrect ordering: two threads increment the counter, but the result is 1 instead of 2.

i (thread)	instr <sub>i</sub>	%eax <sub>1</sub>	%eax <sub>2</sub>	cnt
1	H <sub>1</sub>	-	-	0
1	L <sub>1</sub>	0	-	0
1	U <sub>1</sub>	1	-	0
2	H <sub>2</sub>	-	-	0
2	L <sub>2</sub>	-	0	0
1	S <sub>1</sub>	1	-	1
1	T <sub>1</sub>	1	-	1
2	U <sub>2</sub>	-	1	1
2	S <sub>2</sub>	-	1	1
2	T <sub>2</sub>	-	1	1

### 8.2.2 Process graphs

A *progress graph* depicts the discrete execution state space of concurrent threads. Each axis corresponds to the sequential order of instructions in a thread. Each point corresponds to a possible *execution state*. A *trajectory* is a sequence of legal state transitions that describes one possible concurrent execution of the threads.

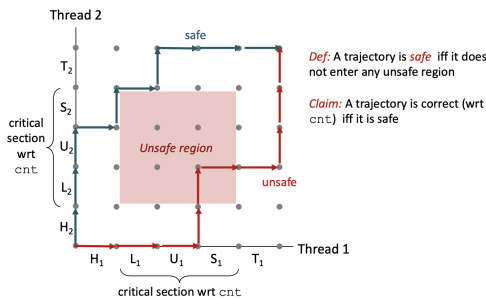


Figure 31. L, U, and S form a critical section with respect to the shared variable cnt.

To enforce safe trajectory, we must need to guarantee mutually exclusive access to critical regions.

### 8.3 Semaphores

*Semaphore* is non-negative global integer synchronization variable. It is manipulated by:

- P(s): [ while (s == 0) wait(); s--; ] (test)
- V(s): [ s++; ] (increment)

OS kernel guarantees that operations between [] are executed indivisibly: Only one P or V operation at a time can modify s.

#### 8.3.1 Proper synchronization

- Binary semaphore: semaphore whose value is always 0 or 1
- Mutex: binary semaphore used for mutual exclusion
  - P operation: locking the mutex
  - V operation: unlocking the mutex

```

1 int sem_init(sem_t *sem, 0, unsigned int val);
2 int sem_wait(sem_t *s); // P(s)
3 int sem_post(sem_t *s); // V(s)
4 void P(sem_t *s); // Wrapper func for sem_wait void V(
    sem_t *s); // Wrapper func for sem_post
5
6 volatile int cnt = 0; // Counter
7 sem_t mutex; // Semaphore that protects cnt
8 sem_init(&mutex, 0, 1); // mutex = 1
9
10 /* Surround critical section with P and V */
11 for (i = 0; i < niters; i++) {
12     P(&mutex);
13     cnt++;
14     V(&mutex);
15 }
    
```

Semaphore creates a forbidden region that encloses unsafe region that cannot be entered by any trajectory.

## 9 Attack Lab

### 9.1 Buffer overflow

When exceeding the memory size allocated for an array. It can cause security vulnerabilities.

#### Vulnerable buffer example

```

1 void echo(){
2     char buf[4]; // Way too small!
3     gets(buf); // Runs until EOF so no way to specify
                // limit on number of characters
4     puts(buf); }
5 void call_echo() {
6     echo(); }
    
```

Before call to gets	After call to gets	With canary
Stack Frame for call_echo	Stack Frame for call_echo	Stack Frame for call_echo
Return Address (8 bytes)	00 00 00 00 00 40 00 34	Return Address (8 bytes)
20 bytes unused	33 32 31 30 39 38 37 36 35 34 33 32 31 30 29 28 37 36 35 34 33 32 31 30	Canary (8 bytes)
[3][2][1][0]	[3][2][1][0]	[3][2][1][0]

Figure 32. Buffer overflowed and corrupted the return pointer.

Buffer overflow can corrupt the return pointer and may corrupt state (cause segmentation fault) or cause undesired action.

#### 9.1.1 Code injection attack

Return address can be corrupted to point exploit code. To prevent,

- **Avoid vulnerabilities in code.** Use methods that limits string lengths, e.g., fgets, strncpy instead of gets, strcpy.
- **System-level protections.** At the start of the program, allocate random amount of space on stack, making difficult for hacker to predict the beginning of the inserted code.
- **Non-executable code segments.** Mark stack as non-executable.
- **Stack canaries.** Place special value beyond buffer and check for corruption before exiting function. Add -fstack-protector.

```

400733: mov    %fs:0x28,%rax
40073c: mov    %rax,0x8(%rsp)
400741: xor    %eax,%eax
400743: mov    %rsp,%rdi
400746: callq 4006e0 <gets>
40074b: mov    %rsp,%rdi
40074e: callq 400570 <puts@plt>
400753: mov    0x8(%rsp),%rax
400758: xor    %fs:0x28,%rax
400761: je     400768 <echo+0x39>
400763: callq 400580 <__stack_chk_fail@plt>
    
```

Figure 33. Disassembly of canary.

### 9.2 Code injection attack

Set return address to the address of touch1.

1. disas getbuf to see how much rsp is decreased
2. Get the start address of function by disas touch1
3. Put bytes to fill the size, and then put the start address of function in little endian (gh ef cd ab 00 00 00 00)

Set the assembly operations in bytes in the return address.

```

1 pushq $touch2 address$
2 movq $COOKIE$, %rdi
3 retq
    
```

## 10 Shell Lab

```

1  /* Misc manifest constants */
2  #define MAXLINE 1024 /* max line size */
3  #define MAXARGS 128 /* max args on a command line */
4  #define MAXJOBS 16 /* max jobs at any point in time */
5  #define MAXJID 1 << 16 /* max job ID */
6
7  /* Job states */
8  #define UNDEF 0 /* undefined */
9  #define FG 1 /* running in foreground */
10 #define BG 2 /* running in background */
11 #define ST 3 /* stopped */
12
13 /*
14  * Jobs states: FG (foreground), BG (background), ST (stopped)
15  * Job state transitions and enabling actions:
16  *   FG -> ST : ctrl-z
17  *   ST -> FG : fg command
18  *   ST -> BG : bg command
19  *   BG -> FG : fg command
20  * At most 1 job can be in the FG state.
21  */
22
23 /* Global variables */
24 extern char **environ; /* defined in libc */
25 char prompt[] = "tsh> "; /* command line prompt */
26 int nextjid = 1; /* next job ID to allocate */
27 char sbuf[MAXLINE]; /* for composing sprintf messages */
28
29 struct job_t
30 {
31     /* The job struct */
32     pid_t pid; /* job PID */
33     int jid; /* job ID [1, 2, ...] */
34     int state; /* UNDEF, BG, FG, or ST */
35     char cmdline[MAXLINE]; /* command line */
36 };
37 struct job_t jobs[MAXJOBS]; /* The job list */
38 /* End global variables */
39
40 /* Function prototypes */
41 void eval(char *cmdline);
42 int builtin_cmd(char **argv);
43 void do_bgfgkl(char **argv);
44 void do_export(char **argv);
45 void waitfg(pid_t pid);
46
47 void sigchld_handler(int sig);
48 void sigtstp_handler(int sig);
49 void sigint_handler(int sig);
50
51 /* Safe API print functions for signal handler */
52 ssize_t sio_puts(char s[]);
53 ssize_t sio_putl(long v);
54 void sio_error(char s[]);
55
56 /* Here are helper routines */
57 int parseline(const char *cmdline, char **argv);
58 void sigquit_handler(int sig);
59
60 void clearjob(struct job_t *job);
61 void initjobs(struct job_t *jobs);
62 int maxjid(struct job_t *jobs);
63 int addjob(struct job_t *jobs, pid_t pid, int state, char *
64     cmdline);
65 int deletejob(struct job_t *jobs, pid_t pid);
66 pid_t fgpj(struct job_t *jobs);
67 struct job_t *getjobpid(struct job_t *jobs, pid_t pid);
68 struct job_t *getjobjid(struct job_t *jobs, int jid);
69 int pid2jid(pid_t pid);
70 void listjobs(struct job_t *jobs);
71
72 void unix_error(char *msg);
73 void app_error(char *msg);
74 typedef void handler_t(int);
75 handler_t *Signal(int signum, handler_t *handler);
76
77
78
79 /* The shell's main routine */
80 int main(int argc, char **argv)
81 {
82     char c;
83     char cmdline[MAXLINE];
84     int emit_prompt = 1; /* emit prompt (default) */
85
86     /* Install the signal handlers */
87     Signal(SIGINT, sigint_handler); /* ctrl-c */
88     Signal(SIGTSTP, sigtstp_handler); /* ctrl-z */
89     Signal(SIGCHLD, sigchld_handler); /* Terminated or stopped
90     child */
91     Signal(SIGQUIT, sigquit_handler);
92
93     /* Initialize the job list */
94     initjobs(jobs);
95
96     /* Execute the shell's read/eval loop */
97     while (1)
98     {
99         /* Read command line */
100        if (emit_prompt) {
101            printf("%s", prompt);
102            fflush(stdout);
103        }
104        if ((fgets(cmdline, MAXLINE, stdin) == NULL) && ferror(
105            stdin))
106            app_error("fgets error");
107        if (feof(stdin)) { /* End of file (ctrl-d) */
108            fflush(stdout);
109            exit(0);
110        }
111        /* Evaluate the command line */
112        eval(cmdline);
113        fflush(stdout);
114        fflush(stdout);
115    }
116    exit(0); /* control never reaches here */
117 }
118
119 /* eval - Evaluate the command line that the user has just typed
120 in. If the user has requested a built-in command (quit,
121 jobs, bg or fg) then execute it immediately. Otherwise,
122 fork a child process and run the job in the context of the
123 child. If the job is running in the foreground, wait for it
124 to terminate and then return. Note: each child process
125 must have a unique process group ID so that our background
126 children don't receive SIGINT (SIGTSTP) from the kernel
127 when we type ctrl-c (ctrl-z) at the keyboard. */
128 void eval(char *cmdline)
129 {
130     char *argv[MAXARGS];
131     char buf[MAXLINE];
132     int bg;
133     pid_t pid;
134
135     strcpy(buf, cmdline);
136     bg = parseline(buf, argv);
137     if (!builtin_cmd(argv)) {
138         sigset_t mask, prev;
139         sigemptyset(&mask);
140         sigaddset(&mask, SIGCHLD);
141         sigprocmask(SIG_BLOCK, &mask, &prev); // block SIGCHLD
142         if ((pid = fork()) == 0) { // child process
143             setpgid(0, 0);
144             sigprocmask(SIG_SETMASK, &mask, NULL); // unblock
145             SIGCHLD
146             if (execve(argv[0], argv, environ) < 0) {
147                 printf("%s: Command not found\n", argv[0]);
148                 exit(1);
149             }
150         }
151         // parent process
152         addjob(jobs, pid, bg ? BG : FG, cmdline);
153         sigprocmask(SIG_SETMASK, &prev, NULL); // unblock
154         SIGCHLD
155         if (!bg) waitfg(pid);
156         else printf("[%d] (%d) %s", pid2jid(pid), pid, cmdline);
157     }
158     return;
159 }

```

```

148 /* parseline - Parse the command line and build the argv array. 222 /*
    Characters enclosed in single quotes are treated as a 223 * do_bgfgkl - Execute the builtin bg, fg and kill commands
    single argument. Return true if the user has requested a BG 224 */
    job, false if the user has requested a FG job. */ 225
149 int parseline(const char *cmdline, char **argv) 226 void do_bgfgkl(char **argv)
150 { 227 {
151     static char array[MAXLINE]; /* holds local copy of command 228     struct job_t * job;
    line */ 229     char * stop;
152     char *buf = array; /* ptr that traverses command line */ 230     if (argv[1] == 0) {
153     char *delim; /* points to first space delimiter */ 231         printf("%s command requires PID or %%jobid argument\n",
154     int argc; /* number of args */ 232         argv[0]);
155     int bg; /* background job? */ 233     return;
156 234 }
157 strcpy(buf, cmdline); 235 if(argv[1][0] == '%'){
158     buf[strlen(buf) - 1] = ' '; /* replace trailing '\n' with 236     int jid = strtol(argv[1]+1, &stop, 10);
    space */ 237     if (stop == argv[1]+1) {
159     while (*buf && (*buf == ' ')) /* ignore leading spaces */ 238         printf("%s: argument must be a PID or %%jobid\n",
    buf++; 239         argv[0]);
160 240     return;
161 241 }
162 /* Build the argv list */ 242     job = getjobjid(jobs, jid);
163     argc = 0; 243     if (job == NULL) {
164     if (*buf == '\') { 244         printf("%s: No such job\n", argv[1]);
165         buf++; 245     return;
166         delim = strchr(buf, '\'); 246 }
167     } else delim = strchr(buf, ' '); 247 } else {
168 248     int pid = strtol(argv[1], &stop, 10);
169     while (delim) { 249     if (stop == argv[1]) {
170         if (*buf == '$') { 250         printf("%s: argument must be a PID or %%jobid\n",
171             buf++; 251         argv[0]);
172             char * copy = malloc(strlen(buf)-1); 252         return;
173             strncpy(copy, buf, strlen(buf)-1); 253     }
174             if(getenv(copy)!=NULL) argv[argc++] = getenv(copy); 254     job = getjobpid(jobs, pid);
175             *delim = '\0'; 255     if (job == NULL) {
176             buf = delim + 1; 256         printf("%s): No such process\n", argv[1]);
177             while (*buf && (*buf == ' ')) /* ignore spaces */ 257         return;
178             buf++; 258     }
179             if (*buf == '\') { 259     }
180             buf++; 260     kill(-job->pid, SIGCONT);
181             delim = strchr(buf, '\'); 261     if(!strcmp(argv[0], "bg")){
182         } else delim = strchr(buf, ' '); 262         job->state = BG;
183         free(copy); 263         printf("[%d] (%d) %s", job->jid, job->pid, job->cmdline)
184     } else { 264     ;
185         argv[argc++] = buf; 265     }
186         *delim = '\0'; 266     if(!strcmp(argv[0], "fg")){
187         buf = delim + 1; 267         job->state = FG;
188         while (*buf && (*buf == ' ')) /* ignore spaces */ 268         waitfg(job->pid);
189         buf++; 269     }
190         if (*buf == '\') { 270     if(!strcmp(argv[0], "kill")) kill(-job->pid, SIGKILL);
191             buf++; 271     return;
192             delim = strchr(buf, '\'); 272 }
193         } else delim = strchr(buf, ' '); 273 }
194     } 274 }
195 } 275 }
196 276 }
197 argv[argc] = NULL; 277 /* do_export - Execute the builtin export commands */
198 if (argc == 0) /* ignore blank line */ 278 void do_export(char **argv)
199     return 1; 279 {
200 280     char * name = strtok(argv[1], "=");
201 281     char * value = strtok(NULL, "=");
202     /* should the job run in the background? */ 282     setenv(name, value, 1);
203     if ((bg = (*argv[argc - 1] == '&')) != 0) { 283     return;
204         argv[--argc] = NULL; 284 }
205     } 285 }
206     return bg; 286 }
207 } 287 }
208 /* builtin_cmd - If the user has typed a built-in command then 288 /* sigint_handler - The kernel sends a SIGINT to the shell
    execute it immediately. */ 289 whenever the user types ctrl-c at the keyboard. Catch it
209 int builtin_cmd(char **argv) 290 and send it along to the foreground job. */
210 { 291 void sigint_handler(int sig)
211     if(!strcmp(argv[0], "quit")) exit(0); 292 {
212     if(!strcmp(argv[0], "jobs")){ listjobs(jobs); return 1; } 293     pid_t pid;
213     if(!strcmp(argv[0], "bg")) { do_bgfgkl(argv); return 1; } 294     if ((pid = fgpid(jobs)) == 0) return;
214     if(!strcmp(argv[0], "fg")) { do_bgfgkl(argv); return 1; } 295     kill(-pid, sig);
215     if(!strcmp(argv[0], "kill")) { do_bgfgkl(argv); return 1; } 296     return;
216     if(!strcmp(argv[0], "export")) { do_export(argv); return 1; 297 }
    } 298 }
217     return 0; /* not a builtin command */ 299 }
218 } 300 }
219 301 }
220 302 }
221 303 }

```

```

296
297 /* sigtstp_handler - The kernel sends a SIGTSTP to the shell
    whenever the user types ctrl-z at the keyboard. Catch it
    and suspend the foreground job by sending it a SIGTSTP. */
298 void sigtstp_handler(int sig)
299 {
300     pid_t pid;
301     if ((pid = fgpuid(jobs)) == 0) return;
302     kill(-pid, sig);
303     return;
304 }
305
306
307 /* sigchld_handler - The kernel sends a SIGCHLD to the shell
    whenever a child job terminates, or stops because it
    received a SIGSTOP or SIGTSTP signal. The handler reaps all
    available zombie children, but doesn't wait for any other
    currently running children to terminate. */
308 void sigchld_handler(int sig)
309 {
310     int olderrno = errno;
311     int status;
312     pid_t pid;
313     while ((pid = waitpid(-1, &status, WNOHANG | WUNTRACED)) >
314            0) {
315         struct job_t *job = getjobpid(jobs, pid);
316         if (WIFEXITED(status)) {
317             deletejob(jobs, pid);
318         } else if (WIFSIGNALED(status)) {
319             if (WTERMSIG(status) == SIGINT){
320                 sio_puts("Job [");
321                 sio_putl(job->jid);
322                 sio_puts("] (");
323                 sio_putl(pid);
324                 sio_puts(") terminated by signal ");
325                 sio_putl(WTERMSIG(status));
326                 sio_puts("\n");
327             }
328             deletejob(jobs, pid);
329         } else if (WIFSTOPPED(status)) {
330             sio_puts("Job [");
331             sio_putl(job->jid);
332             sio_puts("] (");
333             sio_putl(pid);
334             sio_puts(") stopped by signal ");
335             sio_putl(WSTOPSIG(status));
336             sio_puts("\n");
337             job->state = ST;
338         }
339     }
340     /* other handlers can't overwrite the value of errno */
341     errno = olderrno;
342     return;
343 }
344
345 /* sigquit_handler - The driver program can gracefully terminate
    the child shell by sending it a SIGQUIT signal. */
346 void sigquit_handler(int sig)
347 {
348     printf("Terminating after receipt of SIGQUIT signal\n");
349     exit(1);
350 }
351
352 /* Signal - wrapper for the sigaction function. Different
    versions of Unix can have different signal handling
    semantics. Some older systems restore the action to default
    after catching signal, and some systems don't block
    signals of the type being handled. So use this. */
353 handler_t *Signal(int signum, handler_t *handler)
354 {
355     struct sigaction action, old_action;
356     action.sa_handler = handler;
357     sigemptyset(&action.sa_mask); /* block sigs of type being
        handled */
358     action.sa_flags = SA_RESTART; /* restart syscalls if
        possible */
359     if (sigaction(signum, &action, &old_action) < 0)
360         unix_error("Signal error");
361     return (old_action.sa_handler);
362 }

```

## 11 Malloc lab

```

1 #define ALIGNMENT 8
2 #define ALIGN(size) (((size) + (ALIGNMENT - 1)) & ~0x7)
3
4 #define WSIZE 4 /* Word and header/footer size (bytes) */
5 #define DSIZE 8 /* Double word size (bytes) */
6 #define INITCHUNK (1<<6)
7 #define CHUNKSIZE (1<<12) /* Extend heap by this amount (bytes)
    */
8 #define LISTSIZE 20
9
10 #define MAX(x, y) ((x) > (y) ? (x) : (y))
11 /* Pack a size and allocated bit into a word */
12 #define PACK(size, alloc) ((size) | (alloc))
13
14 /* Read and write a word at address p */
15 #define GET(p) (*(unsigned int *) (p))
16 #define PUT(p, val) (*(unsigned int *) (p) = (val))
17 #define PUT_ADD(p, bp) (*(unsigned int *) (p) = (unsigned int) (bp
    ))
18 /* Read the size and allocated fields from address p */
19 #define GET_SIZE(p) (GET(p) & ~0x7)
20 #define GET_ALLOC(p) (GET(p) & 0x1)
21 /* Given block ptr bp, compute address of its header and footer
    */
22 #define HDRP(bp) ((char *) (bp) - WSIZE)
23 #define FTRP(bp) ((char *) (bp) + GET_SIZE(HDRP(bp)) - DSIZE)
24 /* Given block ptr bp, compute address of next and previous
    blocks */
25 #define NEXT_BLKP(bp) ((char *) (bp) + GET_SIZE(((char *) (bp) -
    WSIZE)))
26 #define PREV_BLKP(bp) ((char *) (bp) - GET_SIZE(((char *) (bp) -
    DSIZE)))
27
28 #define PREV_ADD(bp) ((char *) (bp))
29 #define NEXT_ADD(bp) ((char *) (bp) + WSIZE)
30 #define PREV(bp) (*(char **) (bp))
31 #define NEXT(bp) (*(char **) (NEXT_ADD(bp)))
32
33 static void* extend_heap(size_t words);
34 static void* coalesce(void* bp);
35 static void append(void* bp, size_t asize);
36 static void delete(void* bp);
37 static void* place(void* bp, size_t asize);
38
39 static char* heap_listp;
40 char* segfree_list[LISTSIZE];
41
42 /* mm_init - initialize the malloc package. */
43 int mm_init(void) {
44     for (int i = 0; i < LISTSIZE; i++) {
45         segfree_list[i] = NULL;
46     }
47     if ((heap_listp = mem_sbrk(4 * WSIZE)) == (void*) -1)
48         return -1;
49     PUT(heap_listp, 0); /* Alignment padding */
50     PUT(heap_listp + (1 * WSIZE), PACK(DSIZE, 1)); /* Prologue
        header */
51     PUT(heap_listp + (2 * WSIZE), PACK(DSIZE, 1)); /* Prologue
        footer */
52     PUT(heap_listp + (3 * WSIZE), PACK(0, 1)); /* Epilogue
        header */
53     /* Extend the empty heap with a free block of CHUNK bytes */
54     if (extend_heap(INITCHUNK) == NULL)
55         return -1;
56     return 0;
57 }
58
59 static void* extend_heap(size_t words)
60 {
61     char* bp;
62     size_t size;
63     size = ALIGN(words);
64
65     if ((long) (bp = mem_sbrk(size)) == -1)
66         return NULL;
67
68     /* Initialize free block header/footer and the epilogue
        header */
69     PUT(HDRP(bp), PACK(size, 0)); /* Free block header */
70     PUT(FTRP(bp), PACK(size, 0)); /* Free block footer */

```

```

70  PUT(HDRP(NEXT_BLKp(bp)), PACK(0, 1)); /* New epilogue header 149
    */
71
72  append(bp, size);
73
74  /* Coalesce if the previous block was free */
75  return coalesce(bp);
76 }
77
78 static void* coalesce(void* bp)
79 {
80     size_t prev_alloc = GET_ALLOC(FTRP(PREV_BLKp(bp)));
81     size_t next_alloc = GET_ALLOC(HDRP(NEXT_BLKp(bp)));
82     size_t size = GET_SIZE(HDRP(bp));
83
84     if (prev_alloc && next_alloc) { /* Case 1 */
85         return bp;
86     }
87     else if (prev_alloc && !next_alloc) { /* Case 2 */
88         delete(bp);
89         delete(NEXT_BLKp(bp));
90         size += GET_SIZE(HDRP(NEXT_BLKp(bp)));
91         PUT(HDRP(bp), PACK(size, 0));
92         PUT(FTRP(bp), PACK(size, 0));
93     }
94     else if (!prev_alloc && next_alloc) { /* Case 3 */
95         delete(bp);
96         delete(PREV_BLKp(bp));
97         size += GET_SIZE(HDRP(PREV_BLKp(bp)));
98         PUT(FTRP(bp), PACK(size, 0));
99         PUT(HDRP(PREV_BLKp(bp)), PACK(size, 0));
100        bp = PREV_BLKp(bp);
101    }
102    else { /* Case 4 */
103        delete(bp);
104        delete(NEXT_BLKp(bp));
105        delete(PREV_BLKp(bp));
106        size += GET_SIZE(HDRP(PREV_BLKp(bp))) + GET_SIZE(HDRP(
107            NEXT_BLKp(bp)));
108        PUT(HDRP(PREV_BLKp(bp)), PACK(size, 0));
109        PUT(FTRP(NEXT_BLKp(bp)), PACK(size, 0));
110        bp = PREV_BLKp(bp);
111    }
112    append(bp, size);
113    return bp;
114 }
115
116 static void append(void* bp, size_t asize) {
117     void* prev_p = NULL;
118     void* curr_p = NULL;
119     int lidc = 0;
120
121     while (lidc < LISTSIZE - 1) {
122         if (asize <= 1)
123             break;
124         asize >>= 1;
125         lidc++;
126     }
127     prev_p = segfree_list[lidc];
128     while (prev_p != NULL) {
129         if (asize <= GET_SIZE(HDRP(prev_p)))
130             break;
131         curr_p = prev_p;
132         prev_p = PREV(prev_p);
133     }
134
135     if (prev_p == NULL && curr_p == NULL) {
136         PUT_ADD(PREV_ADD(bp), NULL);
137         PUT_ADD(NEXT_ADD(bp), NULL);
138         segfree_list[lidc] = bp;
139     }
140     else if (prev_p == NULL && curr_p != NULL) {
141         PUT_ADD(PREV_ADD(bp), NULL);
142         PUT_ADD(NEXT_ADD(bp), curr_p);
143         PUT_ADD(PREV_ADD(curr_p), bp);
144     }
145     else if (curr_p == NULL) {
146         PUT_ADD(PREV_ADD(bp), prev_p);
147         PUT_ADD(NEXT_ADD(prev_p), bp);
148         PUT_ADD(NEXT_ADD(bp), NULL);
149         segfree_list[lidc] = bp;
150     }
151     else {
152         PUT_ADD(PREV_ADD(bp), prev_p);
153         PUT_ADD(NEXT_ADD(prev_p), bp);
154         PUT_ADD(NEXT_ADD(bp), curr_p);
155         PUT_ADD(PREV_ADD(curr_p), bp);
156     }
157     return;
158 }
159
160 static void delete(void* bp) {
161     size_t size = GET_SIZE(HDRP(bp));
162     int lidc = 0;
163
164     while (lidc < LISTSIZE - 1) {
165         if (size <= 1)
166             break;
167         size >>= 1;
168         lidc++;
169     }
170
171     if (PREV(bp) == NULL && NEXT(bp) == NULL) {
172         segfree_list[lidc] = NULL;
173     }
174     else if ((PREV(bp) == NULL && NEXT(bp) != NULL)) {
175         PUT_ADD(PREV_ADD(NEXT(bp)), NULL);
176     }
177     else if ((NEXT(bp) == NULL)) {
178         PUT_ADD(NEXT_ADD(PREV(bp)), NULL);
179         segfree_list[lidc] = PREV(bp);
180     }
181     else {
182         PUT_ADD(NEXT_ADD(PREV(bp)), NEXT(bp));
183         PUT_ADD(PREV_ADD(NEXT(bp)), PREV(bp));
184     }
185     return;
186 }
187
188 static void* place(void* bp, size_t asize) {
189     size_t csize = GET_SIZE(HDRP(bp));
190
191     delete(bp);
192
193     if ((csize - asize) > (2 * DSIZE)) {
194         if (asize < 100) {
195             PUT(HDRP(bp), PACK(asize, 1));
196             PUT(FTRP(bp), PACK(asize, 1));
197             PUT(HDRP(NEXT_BLKp(bp)), PACK(csize - asize, 0));
198             PUT(FTRP(NEXT_BLKp(bp)), PACK(csize - asize, 0));
199             append(NEXT_BLKp(bp), csize - asize);
200         }
201         else {
202             PUT(HDRP(bp), PACK(csize - asize, 0));
203             PUT(FTRP(bp), PACK(csize - asize, 0));
204             PUT(HDRP(NEXT_BLKp(bp)), PACK(asize, 1));
205             PUT(FTRP(NEXT_BLKp(bp)), PACK(asize, 1));
206             append(bp, csize - asize);
207             return NEXT_BLKp(bp);
208         }
209     }
210     else {
211         PUT(HDRP(bp), PACK(csize, 1));
212         PUT(FTRP(bp), PACK(csize, 1));
213     }
214     return bp;
215 }
216
217 /* mm_malloc - Allocate a block by incrementing the brk pointer.
218    Always allocate a block whose size is a multiple of the
219    alignment. */
220 void* mm_malloc(size_t size) {
221     size_t asize; /* Adjusted block size */
222     size_t extendsize; /* Amount to extend heap if no fit */
223     size_t tsize;
224     int lidc = 0;
225     char* bp = NULL;
226     /* Ignore spurious requests */
227     if (size == 0)
228         return NULL;
229
230     /* Adjust block size to include overhead and alignment reqs.
231     */
232     if (size <= DSIZE) {
233         asize = 2 * DSIZE;

```

```

227 }
228 else {
229     asize = ALIGN(size + DSIZE);
230 }
231 tsize = asize;
232 /* Search the free list for a fit */
233 while (lidx < LISTSIZE && bp == NULL) {
234     if (((tsize <= 1) && (segfree_list[lidx] != NULL)) || (
235         lidx == LISTSIZE - 1)) {
236         bp = segfree_list[lidx];
237         while (bp != NULL) {
238             if (asize > GET_SIZE(HDRP(bp)))
239                 bp = PREV(bp);
240             else
241                 break;
242         }
243         tsize >>= 1;
244         lidx++;
245     }
246 }
247 /* No fit found. Get more memory and place the block */
248 if (bp == NULL) {
249     extendsize = MAX(asize, CHUNKSIZE);
250     if ((bp = extend_heap(extendsize)) == NULL)
251         return NULL;
252 }
253 bp = place(bp, asize);
254 return bp;
255 }
256
257 /* mm_free - Freeing a block does nothing. */
258 void mm_free(void* ptr) {
259     size_t size = GET_SIZE(HDRP(ptr));
260     PUT(HDRP(ptr), PACK(size, 0));
261     PUT(FTRP(ptr), PACK(size, 0));
262     append(ptr, size);
263     coalesce(ptr);
264 }
265
266 /* mm_realloc - Implemented simply in terms of mm_malloc and
267    mm_free */
268 void* mm_realloc(void* ptr, size_t size) {
269     void* new_ptr = ptr;
270     int extendsize = 0;
271     int sizesum = 0;
272     size_t copysize = size;
273
274     if (size == 0) return NULL;
275     if (copysize <= DSIZE) {
276         copysize = 2 * DSIZE;
277     } else {
278         copysize = ALIGN(size + DSIZE);
279     }
280
281     if (GET_SIZE(HDRP(ptr)) >= copysize)
282         return ptr;
283
284     if (GET_ALLOC(HDRP(NEXT_BLKPTR(ptr))) && GET_SIZE(HDRP(
285         NEXT_BLKPTR(ptr)))) {
286         new_ptr = mm_malloc(copysize - DSIZE);
287         memcpy(new_ptr, ptr, copysize);
288         mm_free(ptr);
289     } else {
290         sizesum = GET_SIZE(HDRP(ptr)) + GET_SIZE(HDRP(NEXT_BLKPTR(
291             ptr)));
292         if (sizesum < (int)copysize) {
293             extendsize = MAX((int)copysize - sizesum, CHUNKSIZE)
294             ;
295             if (extend_heap(extendsize) == NULL)
296                 return NULL;
297         }
298         delete(NEXT_BLKPTR(ptr));
299         PUT(HDRP(ptr), PACK(sizesum + extendsize, 1));
300         PUT(FTRP(ptr), PACK(sizesum + extendsize, 1));
301     }
302     return new_ptr;
303 }

```