

# OS Condensed Summary Notes For Quick In-Exam Strategic Fact Deployment

Maksymilian Mozolewski

April 27, 2021

## Contents

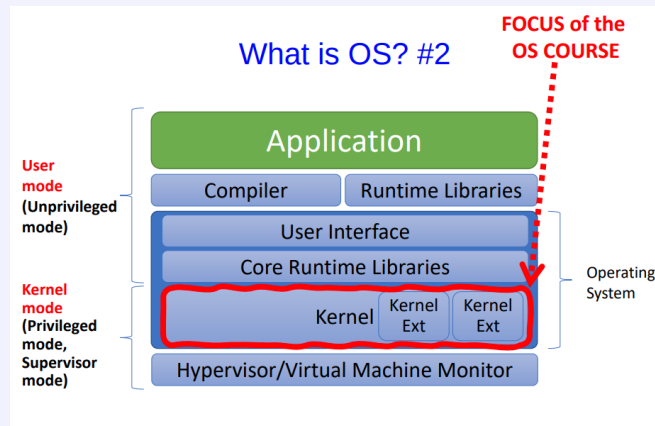
<b>1</b>	<b>OS</b>	<b>2</b>
1.1	Introduction & Structure . . . . .	2
1.2	IO . . . . .	7
1.3	Processes . . . . .	12
1.4	Interprocess communication . . . . .	18
1.5	Threads . . . . .	22
1.6	Scheduling . . . . .	26
1.7	Scheduling algorithms . . . . .	27

# OS

## Introduction & Structure

### OS

An intermediary between the user of a computer and computer hardware. A program itself most intimately connected to the hardware. Everything you don't need to write in order to run your application. Library, all operations on I/O, syscalls. The OS can be an invisible intermediary



Main benefits of OS abstraction:

- Application benefits
  - programming simplicity - see high level abstraction (files) instead of low level hardware details (device registers)
  - abstractions are **reusable** across many platforms
  - **portability** (across machine configurations or architectures) - device independence: 3com or intel card?
- User benefits
  - safety: program sees its own virtual machine, can believe it owns the computer
  - OS **protects** programs from each other
  - OS **fairly multiplexes** resources across programs
  - efficiency - **share** one computer across many users. **Concurrent** execution of multiple programs

## Basic OS Concepts

**Multiprogramming** : ability to keep multiple programs running in parallel (almost). This is done by keeping all possible jobs in the **job pool** on the disk, when a job is ready to take its turn to execute, it's brought to main memory and run for a bit until the OS decides to switch it for another (because of an interrupt) or it completes.

**Multitasking/timesharing** : while multiprogramming makes the OS efficient, it does not facilitate for user interaction! For the user to be able to run multiple tasks simultaneously and have a smooth experience, the CPU needs to actually switch jobs much more frequently so as it appears that all of them are being processed simultaneously. A system which facilitates **timesharing** is **interactive** - it frequently awaits user input and has short response time. A time-shared system enables the system to be used by **multiple** users simultaneously.

**Process** : a program which is loaded in main memory. This may be a full on program, or a printer job. Processes may spawn sub-processes if they wish to using syscalls. This is the **unit of work in a system**.

**Job scheduling** : prioritising which jobs should be in main memory at any time.

**CPU scheduling** : prioritising which jobs in main memory should be executed first.

**Virtual memory** : a technique that allows each program to see the entire memory as theirs and not mess with other processes as well as running programs which require more memory than is **physically available**

**Logical memory** : memory as seen by the programmer, abstracted away from the nitty-gritty mechanical details of the OS.

**Interrupt driven** : if there's no demand for action from the OS, it will IDLE. Events are almost always signalled by the occurrence of an interrupt or a trap

**Trap** : (exception) is a software generated interrupt caused by either an error or a syscall.

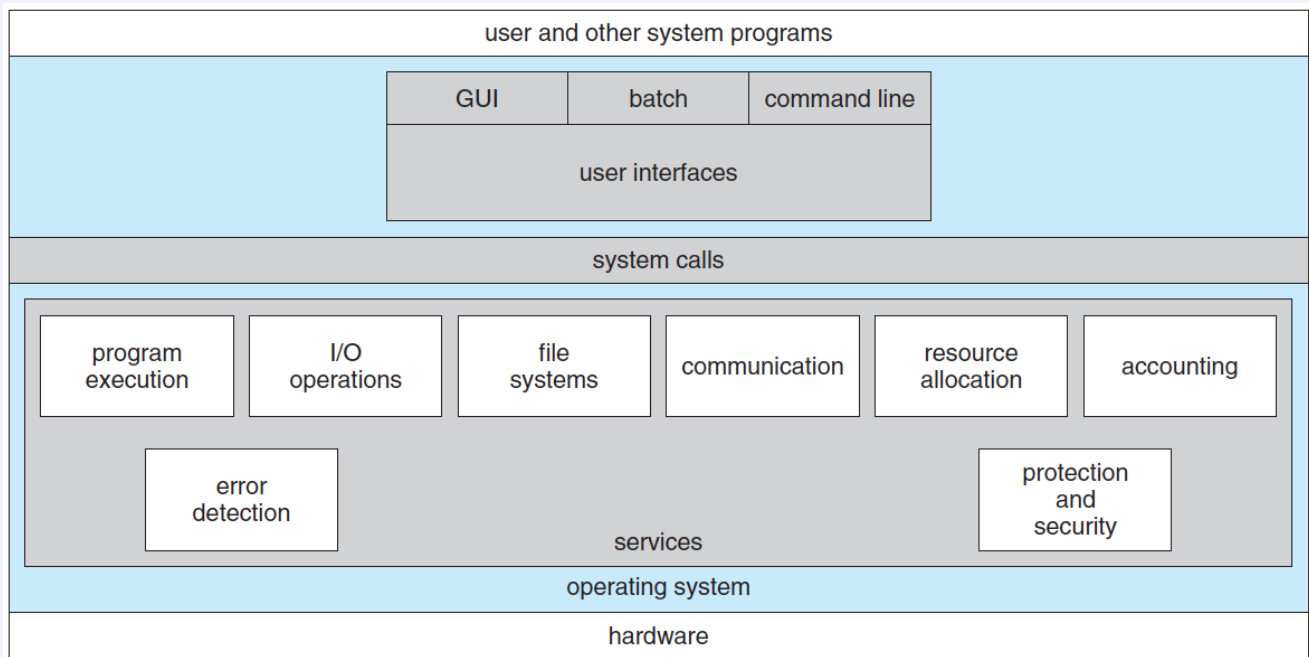
**Dual-mode/multimode operation** : in a multiprogramming system processes are run in parallel, hence protections must be put into place so that processes cannot impede other processes. Most commonly this is done via introduction of **user mode** and **kernel mode**. In user mode certain possibly harmful **privileged instructions** are forbidden by the **hardware itself** which sends a trap signal (switch to kernel mode, timer management, I/O control)

**Virtual machine manager** : virtual machine management software which can be set to run on a third mode (which requires more mode bits), giving it less power than the kernel but more than the user. Virtualisation does not necessarily require its own privilege level

**Mode bits** : reserved bits in the hardware which signify which mode we're in. Typically 0 = user mode and 1 = kernel mode. Always set

**before** passing control to the user program.

**File** : abstract memory concept which is mapped to physical storage space. Each file may contain absolutely any data. The OS is responsible for creating/removing files, creating removing directories to organise files, supporting primitives for manipulating files and directories, mapping files onto secondary storage, backing up files on stable (nonvolatile) storage media.



**User Interface** : can appear in many different forms, CLI, GUI or batch - where commands and directives run directly from files.

**Program execution** : the OS loads programs into memory and runs their instructions, then halts them.

**IO Operations** : interactions with external devices such as the keyboard or mouse.

**File-system manipulation** : search through files and directories, creation and deletion of files, permission management.

**Communications** : facilities for exchanging information between different processes on the same computer or via network. I.e. **Shared memory** or perhaps **message passing**.

**Error detection** : the OS needs to be aware of errors which occur and correct them as they appear. These can happen anywhere in the system, including hardware and software.

**Resource allocation** : distribution of resources available to different jobs and users at the same time efficiently.

**Accounting** : keeping track of who used what resources for either economic purposes or analytics.

**Protection and security** : all data needs to be securely stored and only available to the users who have the correct permissions. Several processes cannot interfere with each other or harm the system.

**Command interpreter** : allows the user to directly interface with the system via commands either in the form of a GUI or CLI or in other forms. When multiple are available these are shown as **shells**. The commands themselves may be stored by the shell, or the shell might simply direct the appropriate loading of file-stored directives which run the appropriate commands (i.e. PATH resolution)

**System calls** : calls to the OS to develop certain specific functions such as opening files or starting sub-processes. Many OS use API's on top of system calls to make portability more achievable.

**System-call interface** : the interface provided by programming languages to interface with the system calls in different OS' which the compiler knows the specifics of.

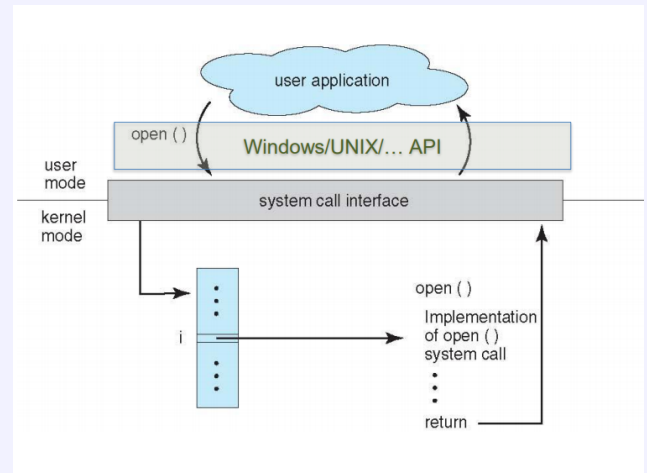
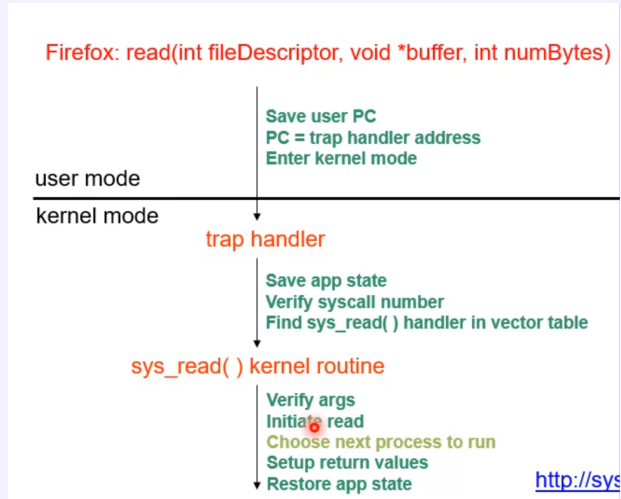
## Syscalls

The user cannot perform IO operations by himself, he must ask the OS to do it for them.

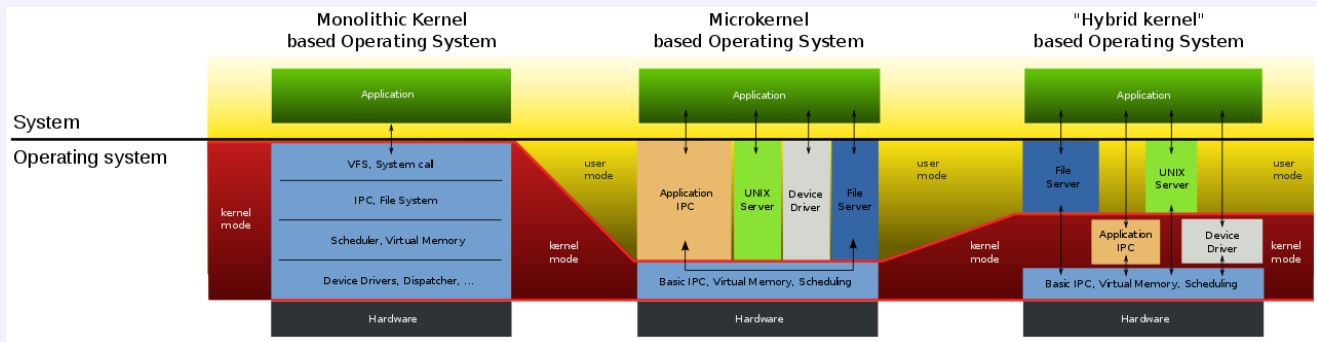
Syscalls define procedures which the OS performs for the user in privileged mode.

Usually implemented as vectors, where each syscall is simply an offset to a base address.

Mechanically just a procedure call (but is not one! in a normal procedure call the caller knows the location of the procedure, in this case a syscall is just an ID), the caller puts arguments in a place the callee expects, then retrieves output from known place. Usually each system call will have wrapper functions provided by each programming language, i.e. the **system call interface**, which the compiler understands.



## OS Structures



### Monolithic Kernel

All major subsystems implemented in kernel.

- + low system interaction cost (procedure call)
- hard to understand and modify
- unreliable (no isolation between system modules)
- hard to maintain

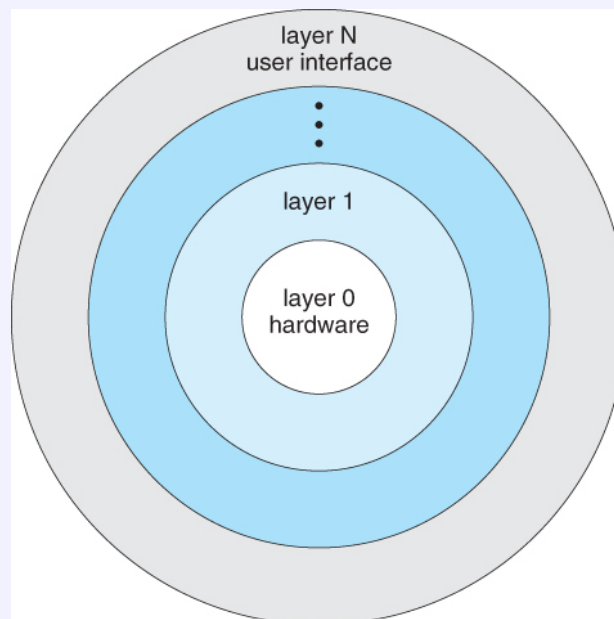
### Microkernel

Minimize what goes into kernel, implement rest of OS as user-level processes

- + better reliability (isolation between components)
- + ease of extension and customization (easy to replace parts)
- poor performance (a lot of user/kernel switches)

### Layered Kernel

Implement OS as a set of layers, each layer interacts only with layer below.



- + more reliable (separation of components)
- strict layering isn't flexible enough - in real life modules might need to communicate with not only nearby layers.
- poor performance, each layer crossing has overhead associated with it (due to API generalization)
- Disjunction between model and reality - system modelled as layers, but not really built that way

## Dynamically loadable kernel modules

Core services in the kernel, others dynamically loaded.

Common in modern implementations:

- **Monolithic**: load the code in kernel space (Solaris, Linux, etc.)
- **Microkernel**: load the code in user space (any)
- + Convenient: no need for rebooting for newly added modules
- + Efficient: no need for message passing unlike microkernel
- + Flexible: any module can call any other module unlike layered model
- Memory fragmentation: fragments OS memory which is normally unfragmented when loaded initially

## Hybrid OS Design

Many different approaches. Key idea: exploit the benefits of monolithic and microkernel designs. Extensibility via kernel modules.

# IO

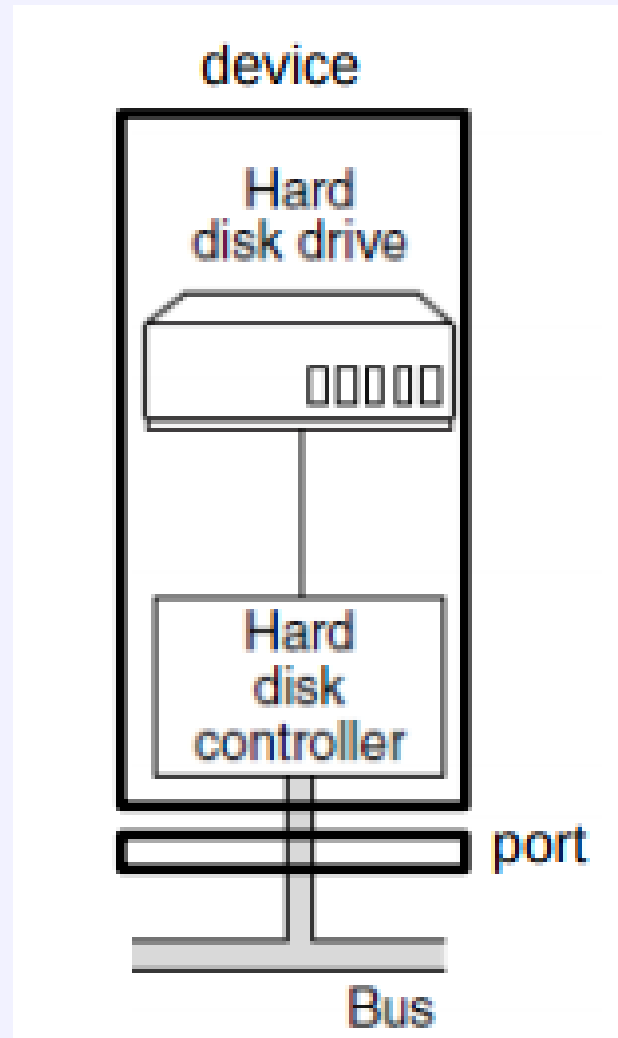
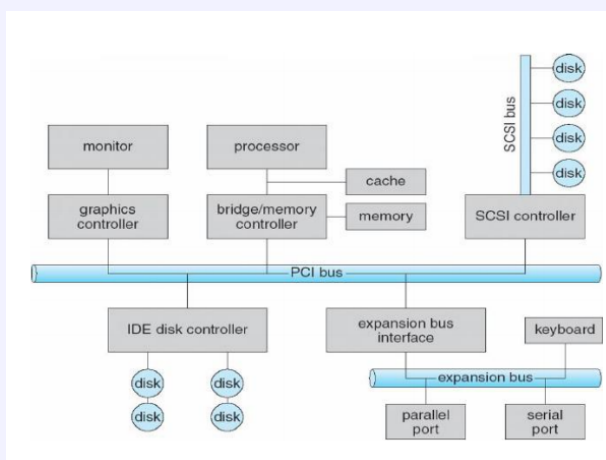
## Variety of I/O Devices

**Port** : Connection point for a device (e.g., USB, parallel, serial, ethernet)

**Bus** : Peripheral buses (e.g. PCI/PCIe), Expansion bus connects relatively slow devices

### Device

**Controller(host adapter)** : electronics that operate port, bus, device (sometimes integrated). Contains processor, microcode, private memory, bus controller etc.



Controllers have **registers** for data and control as well as **buffers**, mostly for data. Communication Methods:

- **IO Ports**
- **Memory-mapped IO**
- **Hybrid**



## I/O Ports

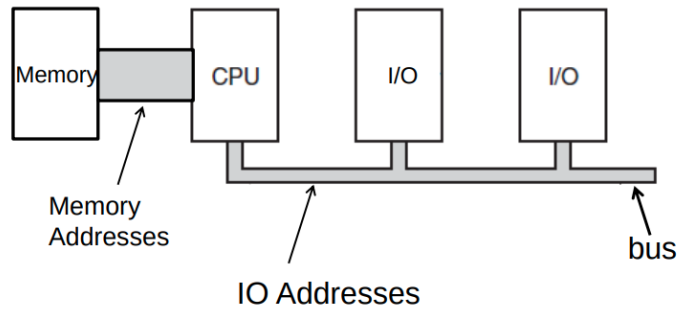
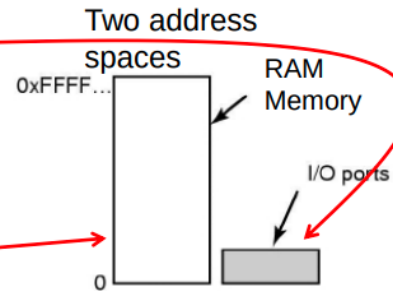
- Each control register has an I/O port number
- Special instructions exist to access the I/O port space
- CPU reads in from device I/O PORT to CPU register (IN REG, PORT)
- CPU writes to device I/O PORT from CPU register (OUT PORT, REG)
- Instructions are privileged (OS kernel only)
- Separate I/O Port space and memory space:

- I/O instructions

- *IN R0, 4*
- *OUT 4, R0*

- Similar memory access instruction

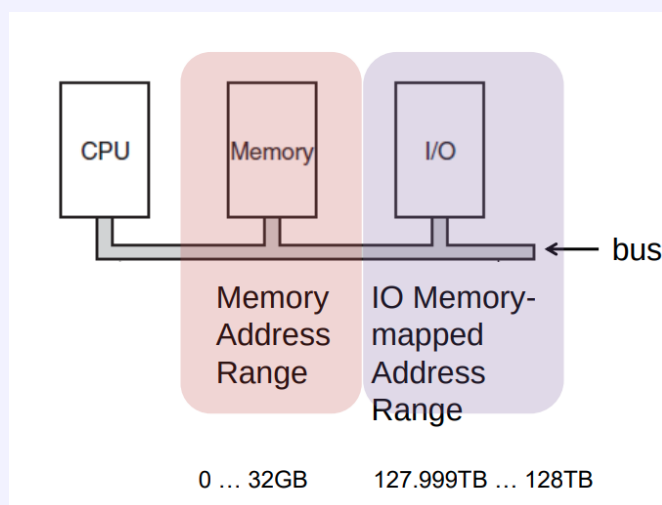
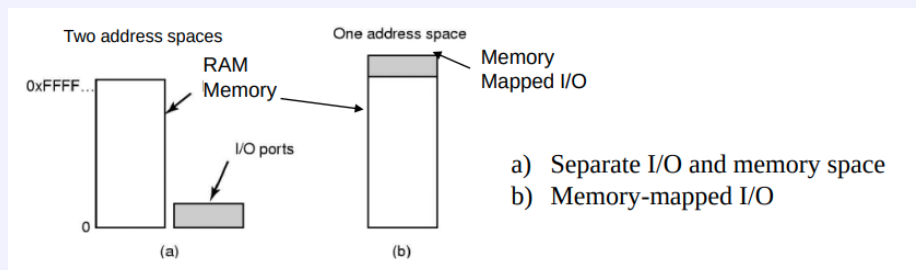
- *MOV R0, 4*
- *MOV 4, R0*



I/O address range (hexadecimal)	device
000-00F	DMA controller
020-021	interrupt controller
040-043	timer
200-20F	game controller
2F8-2FF	serial port (secondary)
320-32F	hard-disk controller

## Memory Mapped I/O

- All control registers and buffers mapped into the memory space
- Each control register is assigned a unique memory address
- There is no actual RAM memory for that address
- Such addresses may be at the top of the physical address space



## Hybrid I/O

Simply do both, use memory mapped I/O for the **data buffers**, and keep separate I/O ports for **control registers**.

## Offloaded Communication

The CPU can request data from an I/O controller one byte at a time (Programmed IO). This wastes CPU time for large data transfers, small data transfers are ok.

CPU can instead offload data transfers using DMA:

**DMA (Direct Memory Access) controller** transfers data from the CPU, either from/to IO or between IO devices.

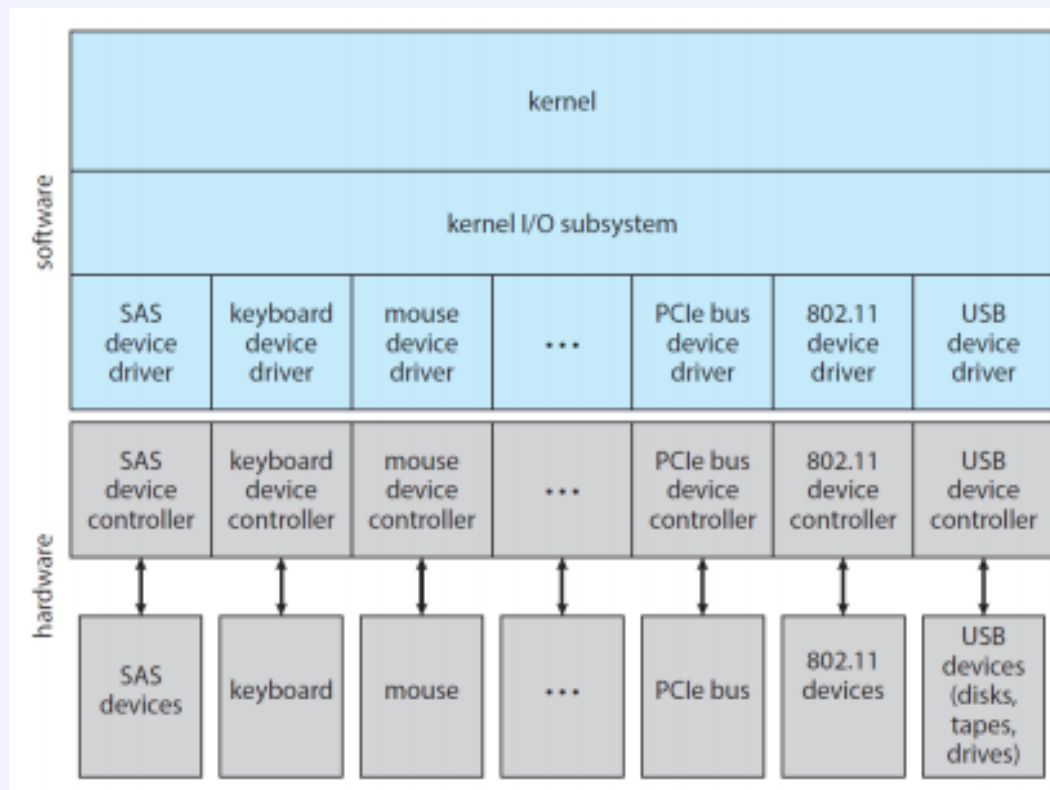
This requires a DMA controller either on the device's host controller, and the motherboard. This controller contains registers to be read/written by the software:

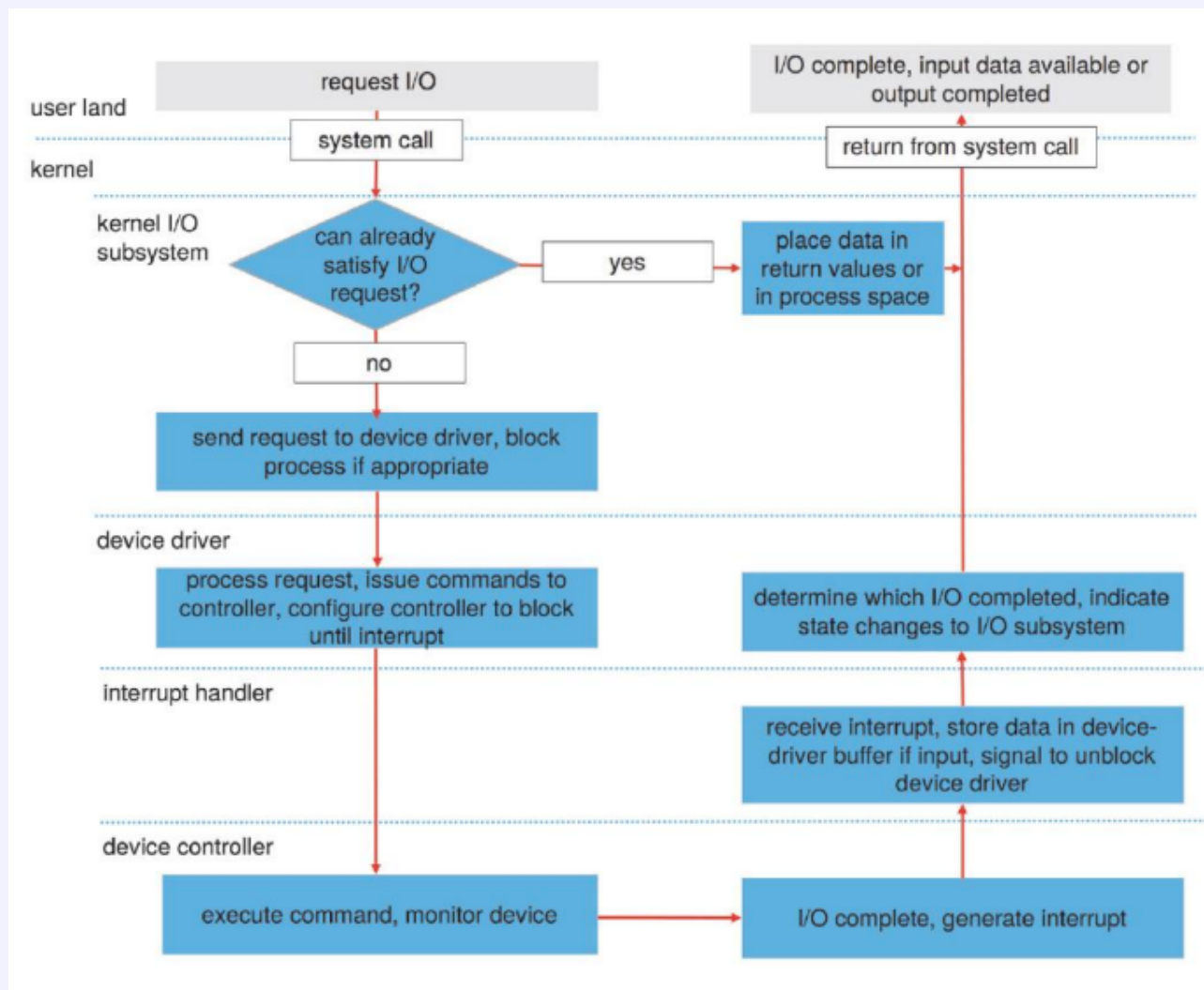
- Memory address register
- byte count register
- Control registers: direction, unit, byte burst size etc..

## OS Device drivers

Great variety of devices, each one has very different specs.

OS Deals with IO devices in a standard and uniform way. Each type of driver is an interface which the vendor can implement as a class. Each OS has its own standard.





## Processes

## Process

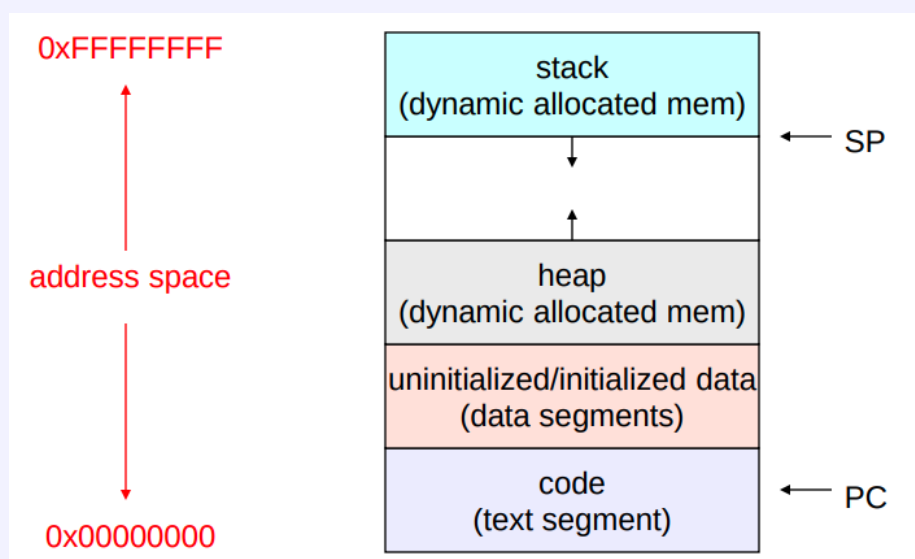
Process is the OS's abstraction for execution.

- Program is the list of instructions, initialized data, etc
- A process is a **program in execution**
- A single flow/sequence of instruction in execution
- An address space (an abstraction of the CPU)

Only one process can be running on a processor core at any instant  
Contents:

- An address space, containing:
  - Code (instructions) for running program
  - Data for the running program (static data, heap data, stack)
- A CPU state, consisting of
  - Program counter, indicating the next instruction
  - Stack pointer, current stack position
  - Other general-purpose register values
- A set of OS resources
  - Open files
  - network connections
  - sound channels

I.e. everything needed to run the program



Each process is identified by a process ID (**PID**). PID's are unique and global. With certain exceptions (cgroups)

Operations that create processes return a PID, and those which operate on processes accept PID's as arguments

## Process representation

Each process is represented internally by the OS with a **Process control block (PCB)** or process/task descriptor, identified by the PID.

OS keeps all of a process's execution state in (or linked from) the PCB when the process isn't running:

- PC, SP, registers etc.
- when a process execution is stopped, its state is transferred out of the hardware into the PCB

When the process is running its state is spread between the PCB and the hardware (CPU regs)

PCB's contain:

- Process ID (PID)
- Parent process ID
- Execution state
- PC, SP, registers
- Address space info
- UNIX user ID, group ID
- Scheduling priority
- Accounting info
- Pointers for state queues

and likely many more.

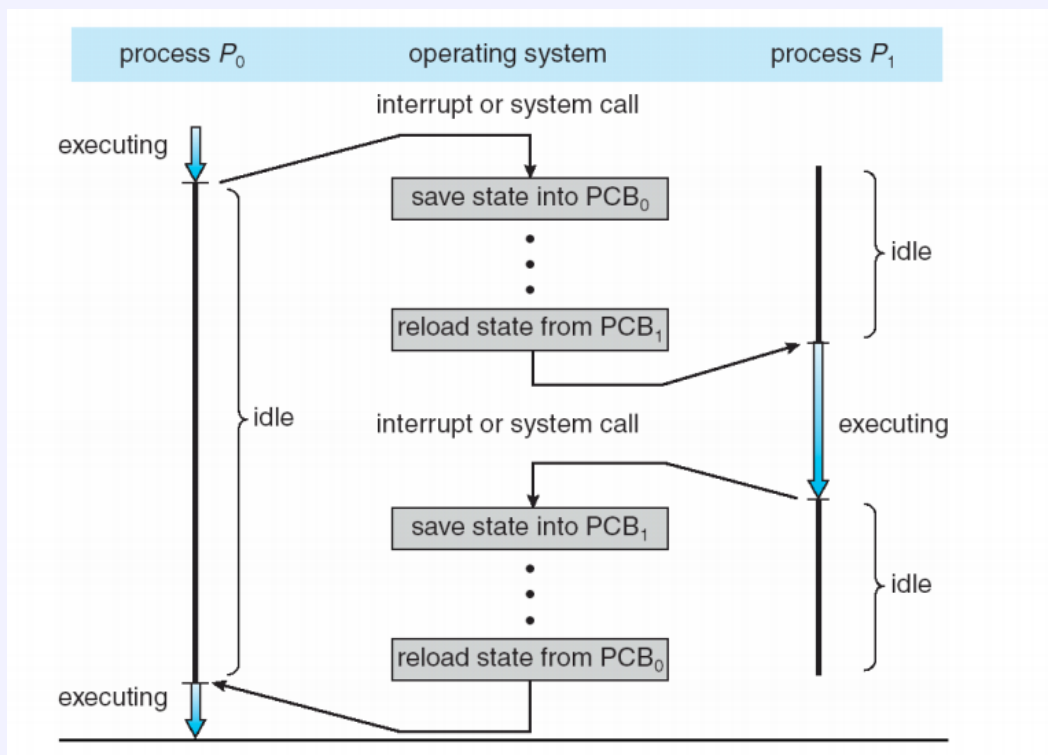
Whenever the OS gets control because of :

- Syscalls
- Exceptions
- Interrupts

The OS then saves the CPU state into the PCB. Whenever the process is resumed into execution again, its PCB is loaded onto the machine registers SP, PC etc.

This is called a **Context switch**

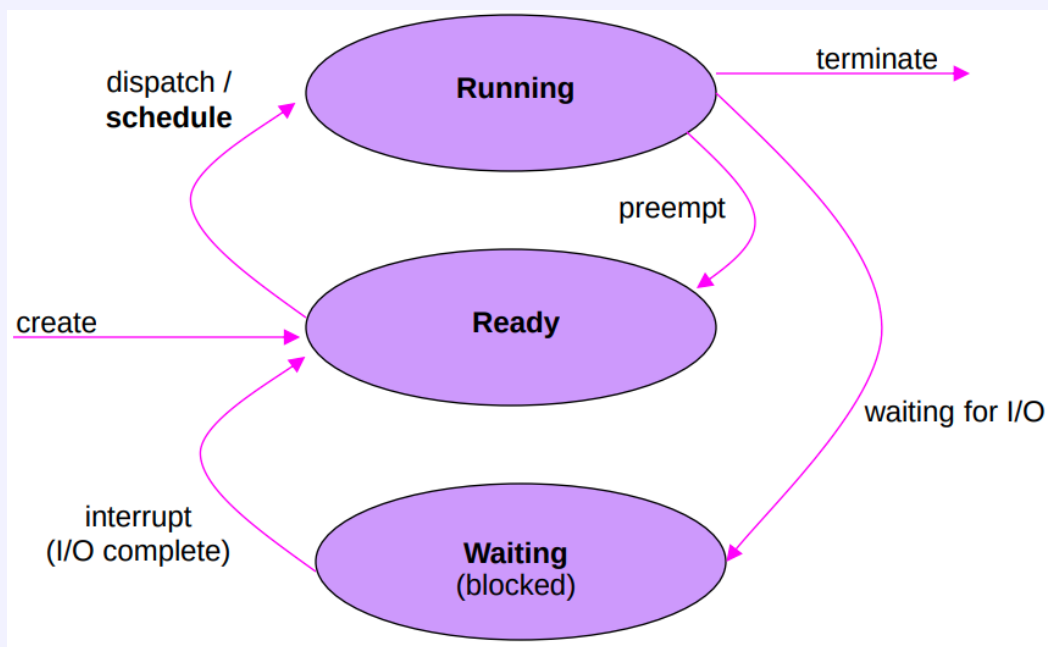
## Context Switch



## Execution states

Each process has an **Execution state**, which indicates what it's currently doing.

- **Ready:** waiting to be assigned to a CPU
- **Running:** executing on a CPU
- **Waiting:** Waiting for an event, e.g. IO completion, or a message from another process.



## State queues

The OS maintains a collection of queues, that represent the state of all processes in the system. Typically one queue for each state (executing, waiting etc..) Each PCB is queued onto a state queue according to the current state of the process it represents. As a process changes state, its PCB is unlinked from one queue, and then linked onto another.

There may be many wait queues, one for each type of wait (specific device, timer, message) etc.

## Process creation

New processes are created by existing processes (parent-child)

The first process is started by the OS, everything else stems from it (init in linux)

Depending on OS, child processes inherit certain attributes of parent, (i.e. open file table: implies stdin/stdout/stderr) Some systems divide resources of parent between children.

## UNIX - fork()

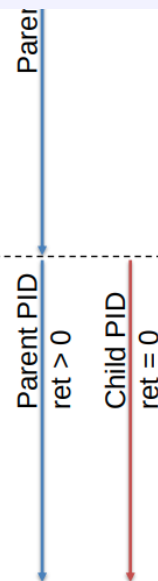
On UNIX systems, process creation is done through the fork() system call:

- Creates and initializes a new PCB
- Initializes kernel resources of new process with resources of parent (e.g. open files)
- Initializes PC,SP to be same as parent
- Creates a new address space, which is an identical copy of this of the parent's (by value)

The fork call returns "twice" once in the parent, and once in the child. In the child the PID returned is 0 and in the parent, the child's PID is returned.

```
#include <sys/types.h>
#include <unistd.h>
#include <stdio.h>

int main(int argc, char **argv)
{
    char *name = argv[0];
    int ret = fork();
    if (ret < 0) { /* error */
        printf("Error\n");
        return 1;
    } else if (ret > 0) { /* parent */
        printf("Child of %s is %d\n", name, ret);
        return 0;
    } else { /* child */
        printf("My child is %d\n", ret);
        return 0;
    }
}
```



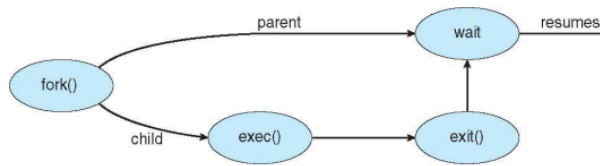


## UNIX - exec()

In order to start a new program instead of just copying the old one, we must use `exec()`. Which is the call which stops the current process, loads a new program into the address space, initializes the hardware context and args for the new program and finally places the PCB onto the ready queue

```
#include <sys/types.h>
#include <stdio.h>
#include <unistd.h>
```

```
int main() {
    pid_t pid;
    pid = fork(); /* fork a child process */
    if (pid < 0) { /* error occurred */
        fprintf(stderr, "Fork Failed");
        return 1;
    }
    else if (pid == 0) { /* child process */
        execlp("/bin/ls", "ls", NULL);
    }
    else { /* parent process */
        wait(NULL); /* parent waits for the child to complete */
        printf("Child Complete");
    }
    return 0;
}
```



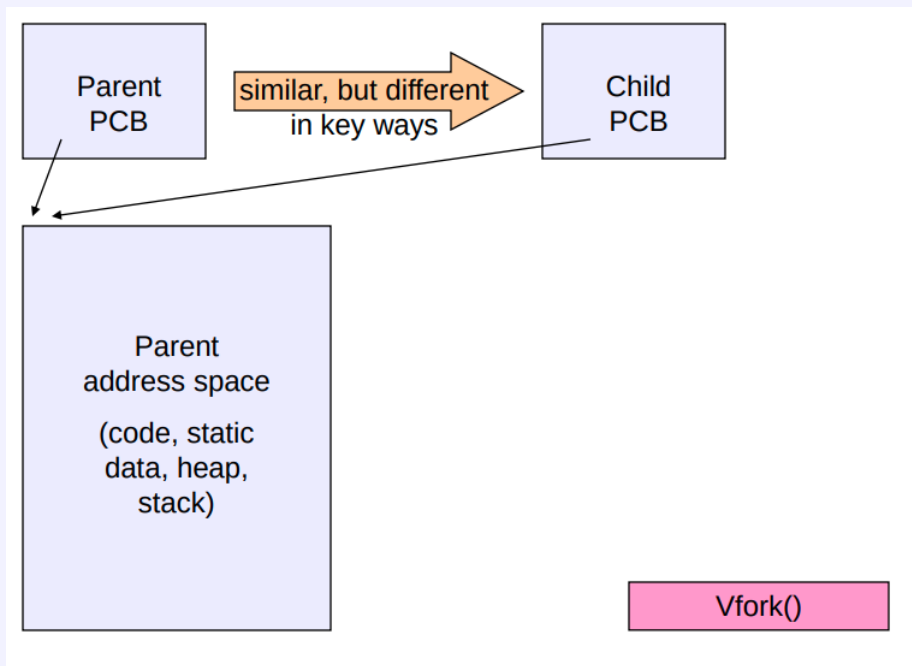
28

Alternatively use `vfork()` which is faster but less safe

## UNIX - vfork()

Same as `fork`, but the child's address space **IS** by address the same space as the parents.

Usage relies on the child not modifying the address space before doing an `execve()` call, otherwise bad things can happen.



## Copy-on-write (COW) fork()

Retains original semantics, but copies "only what is necessary" rather than the entire address space.  
On fork():

- create a new address space
- initialize page tables with same mappings as parents (identical)
- Set both parent and child page tables to make all pages read-only
- if either parent or child writes to memory, an exception occurs
- On exception, OS copies the page, adjusts page tables, etc..

# Interprocess communication

## Shared memory

Allow processes to communicate and synchronize:

- Sharing part of address space
- OS doesn't mediate communication (no overhead)
- Usually OS prevents processes from accessing each other's memory
- Processes should agree to void this restriction

Data:

- Format decided by application
- Direct access (not mediated by the OS) - very fast
- Application programmer fully manages the data transfer - not trivial

Possible use cases:

- Passing of large (single) objects (image)
- Notification variable

```
#include <stdio.h>
#include <stdlib.h>
#include <string.h>
#include <fcntl.h>
#include <sys/shm.h>
#include <sys/stat.h>
#include <sys/mman.h>
```

```
int main() {
    const int SIZE = 4096; /* the size (bytes) of shared memory object */
    const char *name = "OS"; /* name of the shared memory object */
    const char *message 0 = "Hello"; /* written to shared memory */
    const char *message 1 = "World!"; /* written to shared memory */
    int fd; /* shared memory file descriptor */
    char *ptr; /* pointer to shared memory object */

    /* create the shared memory object */
    fd = shm_open(name, O_CREAT | O_RDWR, 0666);
    /* configure the size of the shared memory object */
    ftruncate(fd, SIZE);
    /* memory map the shared memory object */
    ptr = (char *)
        mmap(0, SIZE, PROT_READ | PROT_WRITE, MAP_SHARED, fd, 0);
    /* write to the shared memory object */
    sprintf(ptr, "%s", message 0);
    ptr += strlen(message 0);
    sprintf(ptr, "%s", message 1);
    ptr += strlen(message 1);

    return 0;
}
```

```
#include <stdio.h>
#include <stdlib.h>
#include <fcntl.h>
#include <sys/shm.h>
#include <sys/stat.h>
#include <sys/mman.h>
```

```
int main() {
    const int SIZE = 4096; /* the size (bytes) of shared memory object */
    const char *name = "OS"; /* name of the shared memory object */

    int fd; /* shared memory file descriptor */
    char *ptr; /* pointer to shared memory object */

    /* open the shared memory object */
    fd = shm_open(name, O_RDONLY, 0666);
    /* memory map the shared memory object */
    ptr = (char *)
        mmap(0, SIZE, PROT_READ | PROT_WRITE, MAP_SHARED, fd, 0);
    /* read from the shared memory object */
    printf("%s", (char *)ptr);
    /* remove the shared memory object */
    shm_unlink(name);

    return 0;
}
```

## Message Passing

Allow processes to communicate and synchronize:

- Without sharing part of address space
- OS mediates communication (overhead likely)

Works with processes on the same machine and also those on different inter-networked machines! This is not possible with shared memory.

Message passing facility provides at least two operations:

- Send(message)
- Receive(message)

Communication link:

- Several implementation tradeoffs, .e.g. messages size
- Fixed
- Variable

Communicating processes must refer to each other:

- Direct communication:
  - Symmetric: explicit name of sender and receiver
    - \* send(P, message) - send message to P
    - \* receive(P, message) - receive message from P
  - Asymmetric: Explicit at least on one end:
    - \* send(P, message) - send to p
    - \* receive(id, message) - receive from any process, sender saved in id
- Indirect communication:
  - No need to know/explicitly in advance sender and/or receiver
  - Mailboxes (e.g., POSIX mailbox)
    - \* send(A, message) - send message to mailbox A
    - \* receive(A, message) - receive message from mailbox A
  - A mailbox can be accessed by more than two processes
  - Multiple mailboxes might exist between processes

send() and receive() calls might be implemented as **blocking** or **synchronous** as well as **nonblocking** or **asynchronous**. different combinations of these might be offered:

- Blocking send
- Nonblocking send
- Blocking receive
- Nonblocking receive

**Rendezvous** - When both send and receive are blocking

## Buffering

Messages exchanged reside in temporary buffers/queues with either:

- Zero capacity (no buffering) - no message waiting, sender must block until recipient receives message
- Bounded capacity - n messages may reside. If the queue is not full, then nonblocking, otherwise blocking
- Unbounded - never blocks

## Example implementation: Pipes

A pipe acts as a conduit allowing two processes to communicate **one-way** only and either **Annonymously** or **Named**.

```
#include <sys/types.h>
#include <stdio.h>
#include <string.h>
#include <unistd.h>
#define READ_END 0
#define WRITE_END 1
int main(void) {
    char write_msg[256] = "Greetings";
    char read_msg[256];
    int fd[2];
    pid_t pid;

    if (pipe(fd) == -1) { /* create the pipe */
        fprintf(stderr, "Pipe failed");
        return 1;
    }
    pid = fork(); /* fork a child process */
    if (pid < 0) { /* error occurred */
        fprintf(stderr, "Fork Failed");
        return 1;
    }
    if (pid > 0) { /* parent process */
        close(fd[READ_END]); /* close the unused end of the pipe */
        write(fd[WRITE_END], write_msg, strlen(write_msg)+1); /* write to the pipe */
        close(fd[WRITE_END]); /* close the write end of the pipe */
    }
    else { /* child process */
        close(fd[WRITE_END]); /* close the unused end of the pipe */
        read(fd[READ_END], read_msg, 256); /* read from the pipe */
        printf("read %s", read_msg);
        close(fd[READ_END]); /* close the read end of the pipe */
    }
    return 0;
}
```

Mechanically act like file descriptors (streams)

## Client-Server communication

- **Sockets** abstraction
  - endpoint for communication
  - identified by an IP address concatenated with a port number
  - servers implementing specific services (SSH, FTP, HTTP) listen to well-known ports
  - an SSH server listens to port 22, an FTP server listens to port 21, a web or http server listens to port 80
- **Remove procedure call (RPC)**
  - Abstract the procedure-call mechanism
  - for use between systems with network connections
  - similar in many respects to the IPC
  - uses message-based communication to provide remote service

## Signals

- OS mechanism to notify a process (one way)
- From the OS POV can be thought as a software-generated interruption/exception (synchronous or asynchronous)
- From other processes POV, it is only a notification, no data is transferred. A communication method for: management, synchronization etc.

UNIX: signal handlers must be registered, and provide code for handling the signal.

```
#include <unistd.h>
#include <signal.h>
#include <stdio.h>
#include <sys/types.h>
#include <sys/wait.h>
#define SIG_STOP_CHILD SIGUSR1

main() {
    pid_t pid;
    sigset_t newmask, oldmask;

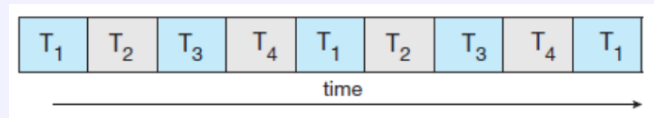
    if ((pid = fork()) == 0) { /* Child */
        struct sigaction action;
        void catchit();
        sigemptyset(&newmask);
        sigaddset(&newmask, SIG_STOP_CHILD);
        sigprocmask(SIG_BLOCK, &newmask, &oldmask);
        action.sa_flags = 0;
        action.sa_handler = catchit;
        if (sigaction(SIG_STOP_CHILD, &action, NULL) == -1) {
            perror("sigusr: sigaction");
            _exit(1);
        }
        sigsuspend(&oldmask);
    }
    else { /* Parent */
        int stat;
        sleep(10);
        kill(pid, SIG_STOP_CHILD);
        pid = wait(&stat);
        printf("Child exit status = %d\n", WEXITSTATUS(stat));
        _exit(0);
    }
}

void catchit(int signo) { /* Signal Handler */
    printf("Signal %d received from parent\n", signo);
    _exit(0);
}
```

# Threads

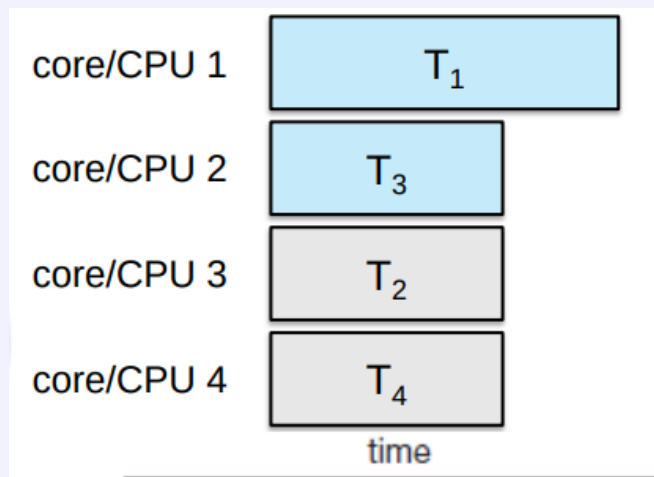
## Concurrency

Concurrency is carrying out multiple tasks in parallel but only one at the same time instance:



## Parallelism

Parallelism is carrying out multiple tasks in parallel at the same time instance:



## Communication problems

Multiple processes are required for both concurrency and parallelism, this requires communication. But the methods discussed thus far have limited usability.

Message passing:

- slow, OS mediates

Shared Memory:

- limited shareability, not all pointers work (both processes have different virtual memory layouts)
- OS resources not shared by default - cumbersome

Possible solution:

1. Fork several processes
2. cause each of them to map to the same shared memory (`shmget()`)
3. make them open the same OS resources

Couple problems:

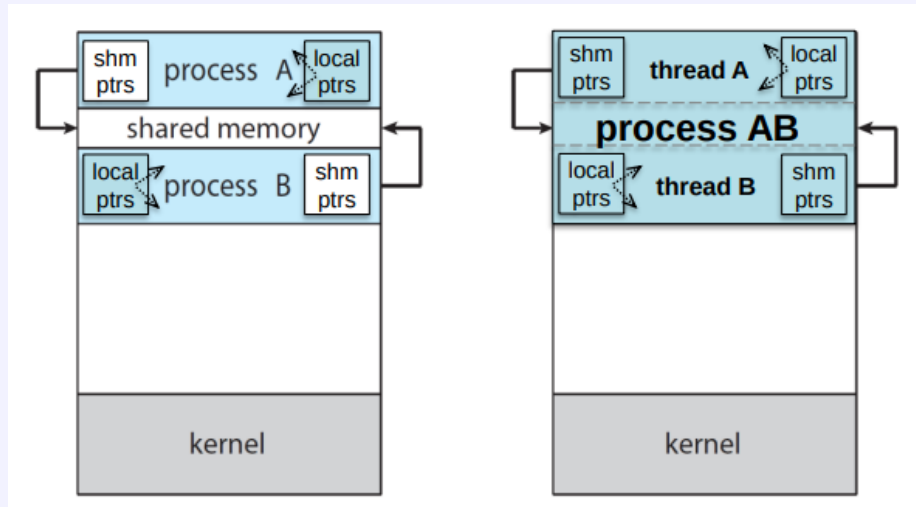
- cumbersome
- has limited shareability again
- inefficient - takes a long time to create all this, requires a PCB per process etc..

## Threads

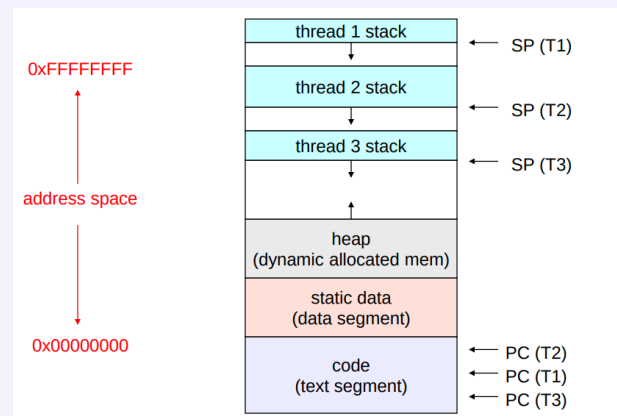
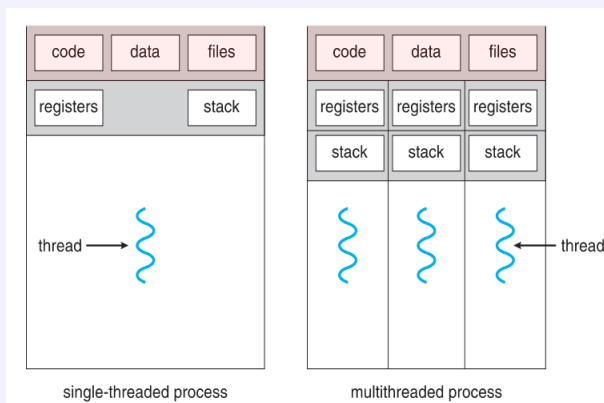
Threads are a great solution to the communication problem: Instead of spawning a new process per task, use **threads**.

Each thread is part of the same spawning process, shares **address space & OS resources**.

Threads only differ in their **execution state (instruction flow)** i.e. private stack, and CPU state



A thread abstracts the execution state away from a process, and now a process represents the static parts of a task (address space, OS resources etc)



Threads become the **unit of scheduling** (depending on implementation of course).  
Think processes are boxes for threads in which they execute.

## Thread Control Block - TCB

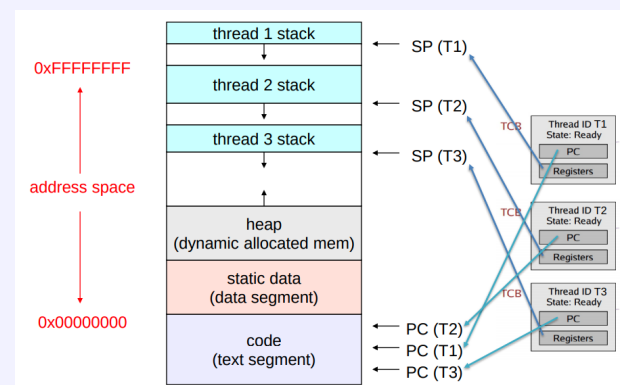
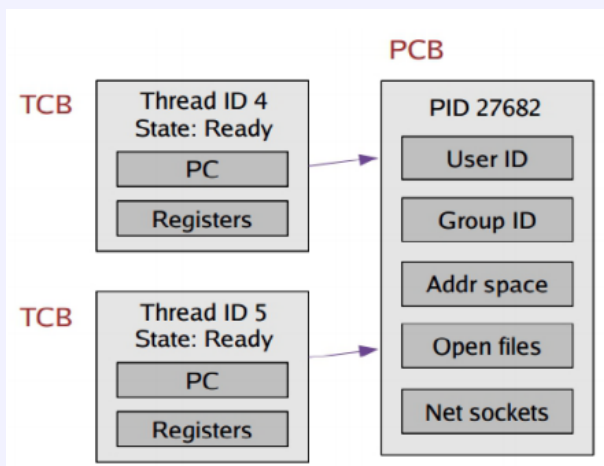
On the OS side, the PCB need to be adjusted to accomodate threads, easiest way is to create sub blocks for each thread representing execution state:

TCB contains:

- Program counter
- CPU registers
- Scheduling information
- Pending I/O information

PCB stores:

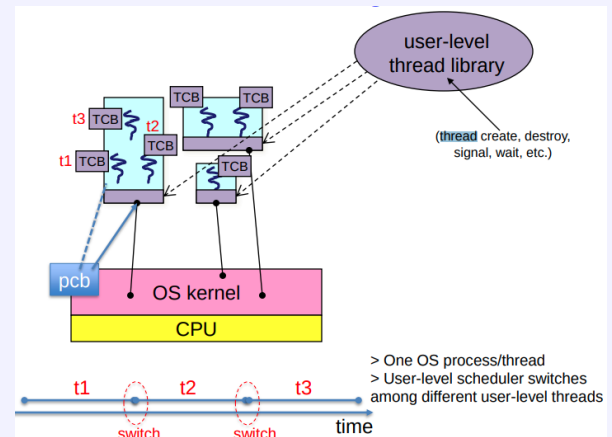
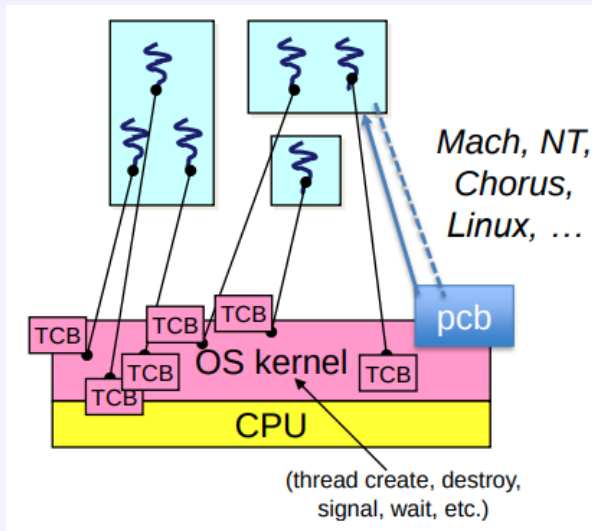
- Memory management information
- Accounting information





## User vs Kernel level threading

Threading can be either implemented as part of the OS, or as a library in the user space



### Kernel level threading 1:1

- OS allocates and manages threads
- TID's are used to identify threads
- + if one thread blocks, the OS can run other threads within the same process
- + possibility to efficiently overlap IO and CPU time within a process
- + Threads are cheaper than processes - less state to manage
- pretty expensive for fine-grained use
  - Orders of magnitude more expensive than procedure calls
  - thread operations are syscalls (context switches + argument checks)
  - Must maintain kernel state for each thread

### User Level threading 1:N

- Threads managed at the user level, within the process
- A library in the program manages the threads
- the thread manager doesn't need to manipulate address spaces (Only OS can)
- Threads differ only in hardware contexts (PC, SP, registers), which can be manipulated by user-level code
- The thread package multiplexes user-level threads in a process
- TID's are now unique per process not globally
- No context switching between thread operations, these are done via procedure calls now (10-100x times faster than kernel threads)
- if one thread tries to do IO, the whole process is blocked!

### N:M Threading

Best of both worlds, can start OS level threads for threads which will use IO.

# Scheduling

## Dispatcher

Mechanism used to switch between tasks (save and restore state)

## Scheduler

Decides on policy (implemented by scheduling algorithms) for ordering execution of tasks (threads/processes)

## CPU Bursts

Bursts of CPU processing done by a task. **Application dependent**

## IO Bursts

Similar to CPU Burst but for IO operations

## Performance goals

- CPU Utilization
- Throughput (processes completed per unit time)
- Turnaround time (time from submission of task to completion)
- Waiting time (all periods spent waiting in the ready queue from submission)
- Response time (time from submission of request to when response produced)
- Energy (joules per instruction) subject to some constraint (fps)

In most cases we optimize the **average metric**. Goals may be conflicting

## Fairness

No single compelling definition of fair for process resource allocation.

Sometimes goal is to be unfair and prioritize some classes of requests higher.

We want to avoid starvation - everyone needs at least some service

## Classes of schedulers

- Batch - throughput / utilization oriented
- Interactive - response time oriented
- Real time - deadline driven

## Preemptive scheduling

**Non-preemptive** scheduling:

- Processes/threads execute until completion or until they want
- The scheduler gets involved only at exit or on request

**Preemptive** scheduling:

- While a process/thread executes, its execution may be paused, and another process/thread resumes its execution
- Involuntary process switch

# Scheduling algorithms

## First-come First-served (FCFS)

Processes/tasks served in the order they arrive:

Process	CPU time	Turnaround time
P1	24	24
P2	3	$24 + 3 = 27$
P3	3	$24 + 3 + 3 = 30$

Execution order: P1,P2,P3

Avg. Turnaround time:  $\frac{24 + 27 + 30}{3} = 27$

- Non pre-emptive
- Poor average response time
- poor utilisation of **other resources** - a CPU-intensive job prevents I/O- intensive job from tiny bit of computation on the CPU before returning to IO and keeping disk busy

## Shortest Job First (SJF)

Associate with each process the length of its CPU time Sort jobs, shortest CPU time goes first Can be preemptive (simply re-sort including the current process - Shortest Remaining Time Next, SRTN)

### Non-Preemptive

Process	Arrival time	CPU time	Turnaround time
P1	0	7	7
P2	2	4	$12 - 2 = 10$
P3	4	1	$8 - 4 = 4$
P4	5	4	$16 - 5 = 11$

Execution order: P1,P3,P2,P4

Avg. Turnaround time:  $\frac{7 + 10 + 4 + 11}{4} = 8$

### Preemptive

Process	Arrival time	CPU time	Turnaround time
P1	0	7	16
P2	2	4	$7 - 2 = 5$
P3	4	1	$5 - 4 = 1$
P4	5	4	$11 - 5 = 6$

Execution order: P1 (2s),P2 (2s),P3(1s),P2(2s),P4(4s),P1(5s)

Avg. Turnaround time:  $\frac{16 + 5 + 1 + 6}{4} = 7$

- + Preemptive is optimal
- Too complex, to be implemented in practice
- not always possible to determine the CPU/IO burst lengths

## Round-robin (RR)

Processes run in discrete time slots, after each time slot a new process/task is chosen to be run

Time quantum = 20	Process	CPU time	Turnaround time
	P1	53	$125 - 0 = 125$
	P2	8	$28 - 0 = 28$
	P3	68	$153 - 0 = 153$
	P4	24	$112 - 0 = 112$

Execution order: P1(20s),P2(8s),P3(20s),P4(20s),P1(20s),P3(20s),P4(4s),P1(13s),P3(20s),P3(8s)

Avg. Turnaround time:  $\frac{125 + 28 + 153 + 112}{4} = 104.5$

Long time quanta cause poor response times with a lot of processes, a too low time quanta, causes a lot of context switching loss.

- + Solves fairness and starvation
- + Fair allocation of CPU across jobs
- + Low average waiting time when job lengths vary
- + Good for responsiveness (interactivity) if small number of jobs
- Context switching time may add up for long jobs

## Priority (PRIO)

Always execute highest-priority runnable jobs to completion

Process	CPU time	Priority	Turnaround time
P1	10	3	16
P2	1	1	1
P3	2	4	18
P4	1	5	19
P5	5	2	6

Execution order: P2,P5,P1,P3,P4

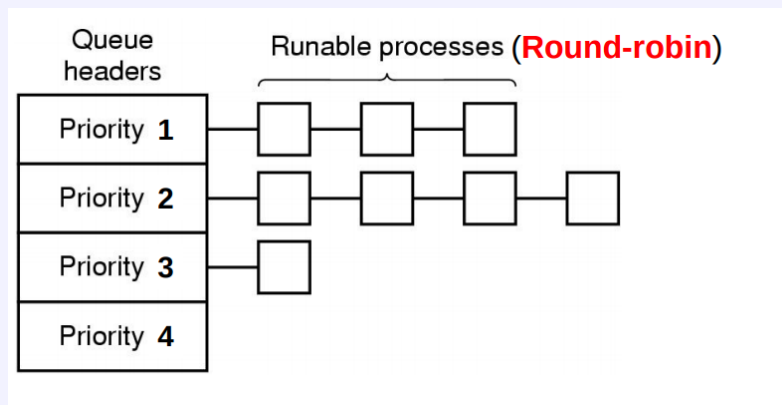
Avg. Turnaround time:  $\frac{16 + 1 + 18 + 19 + 6}{5} = 12$

How to assign priorities ? Based on process type, User, price paid etc. or dynamically, based on how long the process ran etc..

- Starvation - lower priority jobs dont get to run because higher priority always running
- deadlock - priority inversion - happens when a low priority task has lock needed by high priority task (busy waiting)

## Multiple Queues(MQ)

Multiple round-robin scheduled queues, with queues of higher priority always scheduled first



Time quantum = 2

Process	CPU time	Priority	Turnaround time
P1	10	3	19
P2	1	1	1
P3	2	3	10
P4	1	3	11
P5	5	2	6

Execution order: P2(1s),P5(5s),P1(2s),P3(2s),P4(1s),P1(8s)

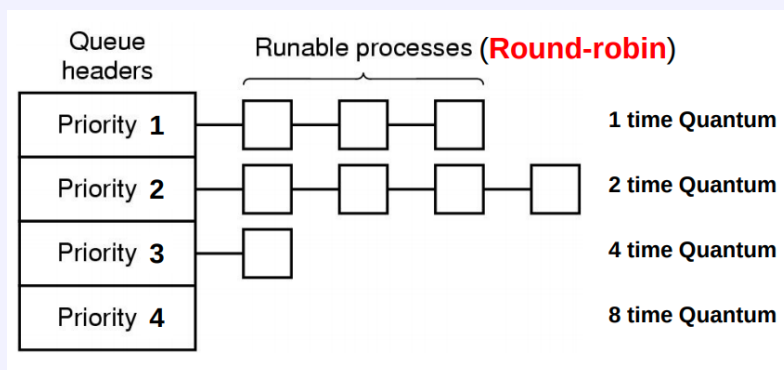
Avg. Turnaround time:  $\frac{19 + 1 + 10 + 11 + 6}{5} = 9.4$

How to assign priorities ? Based on process type, User, price paid etc. or dynamically, based on how long the process ran etc..

- Starvation - lower priority jobs dont get to run because higher priority always running
- deadlock - priority inversion - happens when a low priority task has lock needed by high priority task (busy waiting)

## Multilevel Feedback Queue(MLFQ)

Same as MQ but each queue has a different time quanta. Time quanta are increasing inversely with priority of queue (higher priority lower time quanta). Each process starts in queue 1 - but when it exceeds its quanta it's pushed lower in the queues. When a process becomes inactive it is moved to a higher priority. This can be gamed by making a process interactive.



Time quanta = (1,2,4,8,16,32,64)

Process	CPU time	Turnaround time
P1	100	102
P2	2	3

Execution order: P1(1s),P2(1s),P1(2s),P2(1s),P1(4s),P1(16s),P1(37s)

Avg. Turnaround time:  $\frac{102 + 3}{2} = 52.5$

8 context switches vs 101 with fixed quanta

- Starvation - lower priority jobs dont get to run because higher priority always running
- deadlock - priority inversion - happens when a low priority task has lock needed by high priority task (busy waiting)