

# CS2106 Intro Op. Systems Notes

AY23/24 Sem 2, github.com/gerteck

## 1. Introduction

**Course objectives:** Introduces basic concepts in operating systems.

Focusing on:

- OS structure and architecture, process management, memory management, file management and OS protection mechanism.
- Identify and understand major functionalities of modern operating systems.
- Extend and apply the knowledge in future courses.

**Supplementary Text:** Modern Operating System (5th Edition), by Andrew S. Tanenbaum, Pearson, 2023.

## Learning Outcomes

- Understand how an **OS manages computational resources for multiple users and applications, and the impact on application performance**
- Appreciate the **abstractions and interfaces provided by OS**
- Write **multi-process / thread programs** and avoid common pitfalls such as **deadlocks, starvation and race conditions**.
- Write system programs that utilizes **POSIX** syscall for process, memory and I/O management.
- Self-learn and explore advanced OS topics.
- Understand important design principles in complex systems.

Areas to focus on: Try to understand how things are running in parallel, since we naturally think sequentially. Secondly, how we can manage memory and how they combine and interact (in strange ways), synchronization.

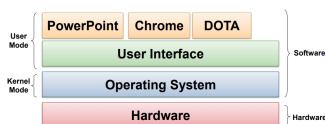
## Operating System OS

An OS is a program that acts as an intermediary between a computer user and the computer hardware. Motivation for OS:

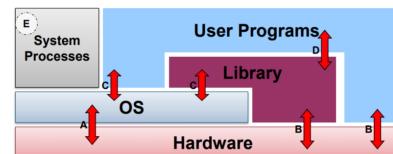
- Manage resources and coordination. (Resource Allocator: Process synchronization, resource sharing)
- Simplify programming (Abstraction of hardware / hardware virtualization, convenient services)
- Enforce usage policies
- Security and protection
- User Program Portability (across different hardware)
- Efficiency (Optimize for particular usage and hardware).

**Kernel Mode:** Complete access to all hardware resources.

**User Mode:** Limited / Controlled access to hardware resources.



## Generic OS Components



- A: OS executing machine instructions
- B: normal machine instructions executed (program/library code)
- C: calling OS using **system call interface**
- D: user program calls library code
- E: system processes
  - Provide high level services, usually part of OS

- OS is known as the **kernel**.  
Program that deals with hardware issues, provide system call interface and special code for interrupt handlers, device drivers.
- Kernel code is different from normal programs:  
No use of system call in kernel code, can't use normal libraries, no normal I/O (must do I/O itself).
- **Implementing OS:** Historically in assembly/machine, now in HLLs (C, C++). Heavily hardware architecture dependent. Challenges include complexity, debugging, codebase size.

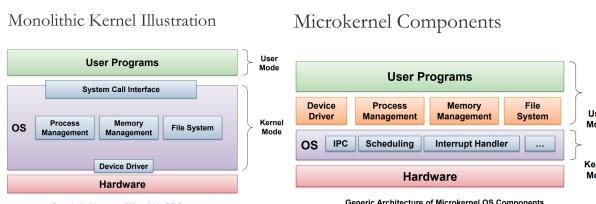
## OS Structures

**Monolithic OS:** One Big program.

- Well understood, good performance, but highly coupled components (everything running in kernel mode) and usually devolved into very complicated internal structure.

**Microkernel OS:**

- Kernel is very small and clean, only providing basic and essential facilities.
- Inter-Process Communication (IPC), Address space management, Thread management etc.
- Higher level services are built on top of basic facilities, run as server process *outside* of OS, use IPC to communicate.
- Kernel is more robust and extendible, better isolation and protection between kernel and high level services. But, lower performance. (Latency)



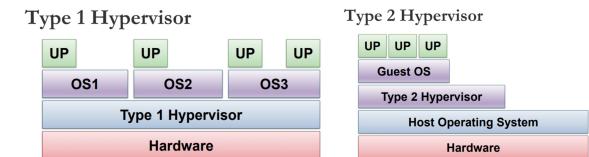
## Other OS Structure

- **Layered Systems:** Generalization of monolithic system, organize components into hierarchy of layers. Lowest is hardware, highest is user interface.
- **Client-Server Model:** Variation of microkernel. Two classes of processes: Client p. request service from server process, server process built on top of microkernel. Client & Server process can be on separate machine.

## Virtual Machines

- **Motivation:** OS assumes total control of hardware, making it hard to run several OS on same hardware at same time. OS is also hard to debug / monitor, hard to observe working of OS, test potentially destructive implementation.
- **Virtual Machine:** Software emulation of hardware.
- **Virtualization of underlying hardware:** Illusion of complete hardware to level above. (Memory, CPU etc.) Normal OS can then run on top of virtual machine. Aka **Hypervisor**.

- **Type 1 Hypervisor:** Provides individual virtual machines to guest OSes (e.g. IBM VM/370)
- **Type 2 Hypervisor:** Runs in host OS, Guest OS runs inside Virtual Machine, (e.g. VMware)



## - Upcoming Topics -

**OS Process Management:** As OS (to maximise efficiency hardware resources), to be able to switch from running one program to the other (share hardware, e.g. CPU), requires information regarding execution of A stored, and A's information replaced with B's information to run. (E.g. the registers in CPU replaced)

- **2. Process Abstraction:** Info describing executing program
- **3. Process Scheduling:** Deciding which process gets to execute
- **4. Inter-Process Communication:** Passing information between processes (tough)
- **5. Threads + Synchronization:** Alternative to Process (Light-weight process aka Thread)
- **6. Memory Management**
- **7. Disjoint Memory Management**
- **8. File System Management**
- **9. File System Implementation**

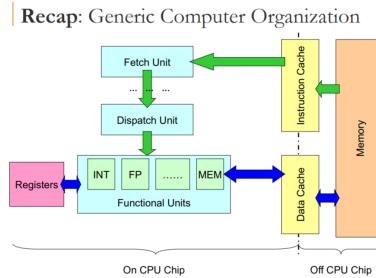
## 2. Process Abstraction

To switch programs, requires information of both programs. Hence, we need abstraction to describe running program, aka **process**.

- (**Process / Task / Job**) is a dynamic abstraction for executing program.
- It is info required to describe a *running program*:
  - Memory Context (Code/Text, Data, Stack, Heap),
  - Hardware Context (Register/PC, Stack/Frame Pointer),
  - OS Context (Process Properties (PID, State), Resources Used).

### Computer Organization (Recap)

- **Components:** Memory, Cache, Fetch Unit (Loads instruction, location indicated by special register **PC**.)
- **Functional Units** (Carry out instruction execution, dedicated to diff. instr. type) (CS2100 looked at INT func. unit)
- Registers (Internal storage, fastest access speed).
  - **GPR:** General Purpose Register, accessible by user program / compiler.
  - **Special Registers:** PC, Stack/Frame Pointer, PSW etc.
- **Binary Executable File:** file in machine language (built by compiler) for specific processor:
  - Executable (binary) consists two major components: Instr. (Text) & Data
  - When under execution, more info: Memory, Hardware, OS context.



### Memory Context for Function Call (Stack Memory)

#### Memory Context Challenges of Functional Calls:

- Control Flow Issues: Need to jump to function body, resume after, need to store PC of caller.
- Data Storage Issues: Need to pass params to function, capture return result, may need declare local variables.
- Require region of memory dynamically used by function invocations.

Hence, portion of memory space used as **stack memory** that stores executing function using **stack frame**, which includes usage of *Stack Pointer, Frame Pointer*.

#### Stack Memory Region

- **Memory region to store information function invocation.**
- **Stack Frame:** Describes information of function invocation.
- Stack frame added on top when function is invoked, stack "grows", removed from top when function call ends, stack "shrinks".
- Stack Frame contains return PC address of caller, arguments for function, storage for local variable, etc.
- **Stack Pointer:** Indicates top of stack region (first unused memory location). Usually indicated in specialized register.

### Function Call Convention: Stack Frame Setup / Teardown

There are different ways to setup stack frame, known as function call convention, differences about (info stored in frame, which portion of stack frame prepared & cleared by caller / callee etc). Dependent on hardware & programming language.

#### Example Scheme:

##### Stack Frame Setup

- Prepare to make a function call:
  - Caller: Pass parameters with registers and/or stack
  - Caller: Save Return PC on stack
  - Transfer Control from Caller to Callee
  - Callee: Save the old Stack Pointer (SP)
  - Callee: Allocate space for local variables of callee on stack
  - Callee: Adjust SP to point to new stack top

##### Stack Frame Teardown

- On returning from function call:
  - Caller: Place return result on stack (if applicable)
  - Caller: Restore saved Stack Pointer
  - Transfer control back to caller using saved PC
  - Caller: Utilize return result (if applicable)
  - Caller: Continues execution in caller

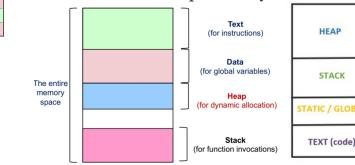
##### Stack Frame Setup / Teardown [Updated example]

- On executing function call:
  - Caller: Pass arguments with registers and/or stack
  - Caller: Save Return PC on stack
  - Transfer control from caller to callee
  - Caller: Save registers used by callee. Save old FP, SP
  - Callee: Allocate space for local variables of callee on stack
  - Callee: Adjust SP to point to new stack top
- On returning from function call:
  - Callee: Restore saved registers, FP, SP
  - Transfer control from callee to caller using saved PC
  - Caller: Continues execution in caller

### Heap Memory

- Managing heap memory trickier due to variable size, variable allocation / deallocation timing.
- Common situation where heap memory alloc/dealloc creating "holes" in memory. Free memory block squeezed between occupied memory blocks.
- Covered in memory management.

#### Illustration for Heap Memory



#### Summary so Far

- It is the dynamic memory that can grow or shrink as per our need. Also called the free storage. The size of all other segments is decided at compile time but heap can grow during runtime. We can control memory allocation and de-allocation in heap space.
- The entire memory space
- It is the space where the local variables gets space. The variables which are declared within functions live in stack.
- A space to store all the global variables. Variables that are declared outside functions and are accessible to all functions.
- To store all the instructions in the program. The instructions are compiled instructions in machine language.

Todd Sauer

### Stack Frame: Other Information

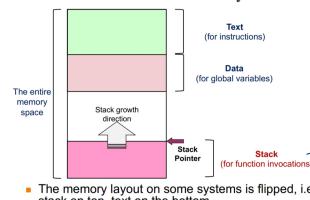
#### Frame Pointer

- To facilitate access of various stack frame items. As stack pointer hard to use as it can change, some processors provide dedicated register Frame Pointer.
- Frame Pointer points to fixed location in stack frame, other items accessed as displacement from frame pointer, usage of FP is platform dependent.

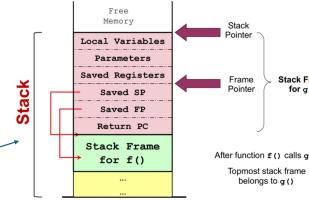
#### Saved Registers

- Since number of GPR limited, when GPR exhausted, use memory to temp. hold GPR values for reuse.
- **Known as Register Spilling.** Function can spill registers it intends to use before function starts, then restore registers at end of function.

#### Illustration: Stack Memory



#### Illustration: Stack Frame v2.0



### OS Context: Process ID, Process State

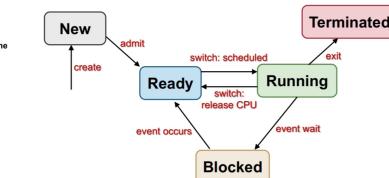
#### Process Identification:

- **Process ID:** To distinguish processes from each other (Just a number, unique among processes)
- PIDs are OS dependent as well, including if PIDs reused, if limits maximum no. of processes or any PIDs reserved.

#### Process State:

- Processes require a process state as indication of execution status. (Running / Not Running / Ready to Run etc.)
- **Process Model:** Set of states and transitions, describes behaviors of a process.
- **Global View of Process States:** Given  $n$  processes,
  - With 1 CPU,  $\leq 1$  process in running state, 1 transition at a time.
  - With  $m$  CPUs,  $\leq m$  processes running state, possibly parallel transitions.
- Different processes may be in different states, each process may be in different part of its state diagram.
- **5-State Process Model:**

#### Generic 5-State Process Model



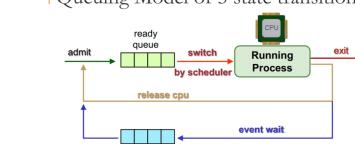
#### Process States for 5-Stage Model

- **New:**
  - New process created
  - May still be under initialization → not yet ready
- **Ready:**
  - process is waiting to run
- **Running:**
  - Process is being executed on CPU
- **Blocked:**
  - Process waiting (sleeping) for event
  - Cannot execute until event is available
- **Terminated:**
  - Process has finished execution, may require OS cleanup

#### Process State Transitions in 5-Stage Model

- **Create** (nil → New):
  - New process is created
- **Admit** (New → Ready):
  - Process ready to be scheduled for running
- **Switch** (Ready → Running):
  - Process selected to run
- **Switch** (Running → Ready):
  - Process gives up CPU voluntarily or *preempted* by scheduler
- **Event wait** (Running → Blocked):
  - Process requests event/resource/service which is not available/in progress
  - Example events:
    - System call, waiting for I/O, (*more later*)
- **Event occurs** (Blocked → Ready):
  - Event occurs → process can continue

#### Queuing Model of 5 state transition



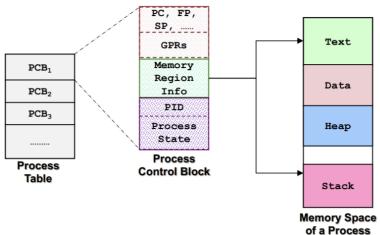
Notes:

- More than 1 process can be in ready + blocked queues
- May have separate event queues
- Queuing model gives global view of the processes, i.e. how the OS views them

## Process Table & Process Control Block

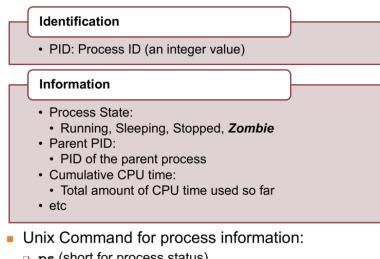
- Since the OS is just a program as well, need to make use of data structures to track these processes.
- Process Control Block (PCB) or Process Table Entry:** Entire Execution Context for a process.
- Kernel maintains PCB for all processes. (Conceptually stored as one table representing all processes.)
- Factors to consider:**
  - Scalability (how many concurrent processes at once).
  - Efficiency (should provide efficient access with minimum space wastage).

Illustration of a Process Table



## Process Abstraction in Unix

### • Process Identification, Information



### • Process Creation, Termination, Parent-Child Synchronization

#### Note: Command Line Argument in C

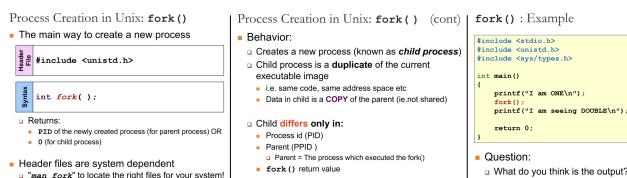
- We can pass arguments to a program in C.
- argc**: Number of CL arguments, including program name.
- argv**: A char strings array, each element in **argv[]** is a C character string.

```
int main( int argc , char* argv[] )
{
    int i;
    for ( i = 0; i < argc; i++ ){
        printf("Arg %i: %s , ", i, argv[ i ] );
    }
    return 0;
}
```

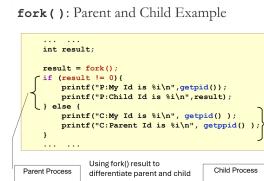
Example Run: "a.out 123 hello world"

Output: "Arg 0: a.out, Arg 1: 123, Arg 2: hello, Arg 3: world"

## Process Creation: fork()



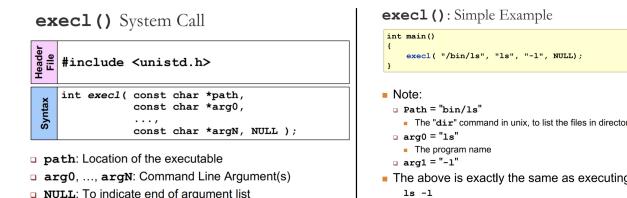
- fork()**: Create exact copy of the parent, including any variables.
- Output**: Both parent and child resume execution after the point **fork()**.
- Note clone()**: **fork()** not versatile, for scenarios where partial duplication preferred, **clone()**, which supersedes **fork()**.
- Both parent and child processes continue executing, common usage is to use the parent/child process differently. (Parent spawns off child to carry out some work, parent ready to take another order.)
- Use return value of **fork()** to distinguish parent and child.



## Process Replacement: exec()

- Function replaces current executing process image with a new process image specified by path. No return is made because the calling process image is replaced by the new process image.

### Code Replacement, but PID and other information still intact.



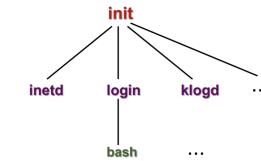
- By combining **fork()** and **exec()**, we can spawn off a child process (let it perform task through **exec()**), while parent process around to accept another request.
- This combination of mechanisms is main way in Unix to get new process for running new program!

## The Master Process: init

- Every process has parent process, consider special initial process.
- init process**: Created in kernel at boot up time, usually PID = 1.
- Purpose**: Watches and respawns other (critical) processes where needed.
- fork()** creates the process tree, where **init** is the root process.

## Simplified Process Tree Ex.

### Process Tree Example (simplified)

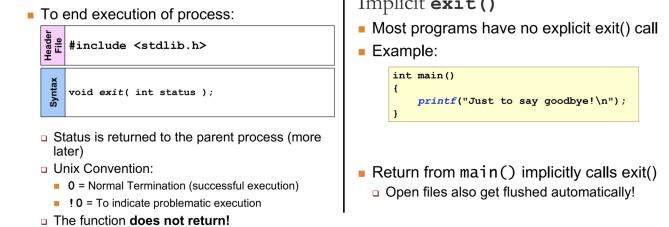


Note: just a simple example, actual process tree varies according to Unix setup

- d**, (e.g. **klogd**) at end of process name usually means server process (Daemon, background process.)

## Process Termination in Unix: exit

### Process Termination in Unix

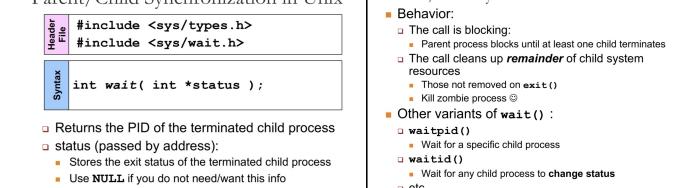


- Process finished execution**: Most system resources used by process are released on exit. (e.g. file descriptors).
- Certain basic process resources not Releasable**: PID, status needed. For parent-children synchronization, for parent to check status of child, For process accounting info (e.g. cpu time).
- Process table entry may still be needed after termination.

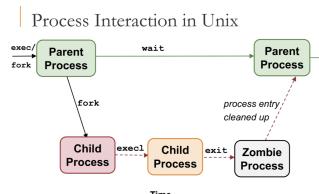
## Parent/Child Synchronization in Unix

- Parent process can wait for child process to terminate.
- Argument is a pointer to a variable that will store the return value (**\*status**).

### Parent/Child Synchronization in Unix



- Kills zombie processes! With enough zombies, process table finite size, run out of space.



Note: example uses one ordering of execution, others are possible!

## Zombie Processes (2 Cases)

- `wait()` "creates" the zombies (and later cleans it up) as on process exit, process becomes zombie.
- Since it cannot delete all process info (if parent asks for info in `wait()` call, remainder of process data structure can be cleaned up only when `wait()` happens.)
- We cannot kill `PID` zombie process, is already dead. Until restart system or modern OS look through table and remove them.

### 1. Parent process terminates before child process

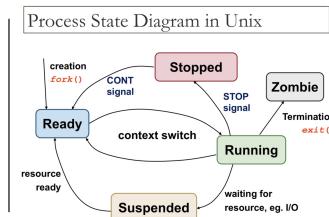
- `init` process becomes "pseudo" parent of child processes.
- Child termination sends signal to `init`, which utilizes `wait()` to cleanup

### 2. Child process terminates before parent but parent did not call wait

- Child process become a zombie process
- Can fill up / hog process table. May need a reboot to clear the table on older Unix implementations

#### Summary of Unix Process System calls

- `fork()`:
  - Process creation
- `exec()` family:
  - Change executing image/program
  - `exec1, execv, execve, execle, execvp`
- `exit()`:
  - Process termination
- `wait()` family:
  - Get status, synchronize with child
  - `wait, waitpid, waiatid, etc`
- `getpid()` family:
  - Get process information
  - `getpid, getppid, etc`



## Implementation Issues

### Implementing `fork()`

- Implementing `fork()`
  - Behavior of `fork()`:
    - Makes an almost exact copy of parent process
- Simplified implementation:
  1. Create address space of child process
  2. Allocate `p' = new PID`
  3. Create kernel process data structures
    - E.g. Entry in Process Table
  4. Copy kernel environment of parent process
    - E.g. Priority (for process scheduling)
  5. Initialize child process context:
    - `PPID=p', PPID=parent id, zero CPU time`

Implementing `fork()` (cont)

- Copy memory regions from parent
  - Program, Data, Stack
  - Very expensive operation that can be optimized (more later)
- Acquires shared resources:
  - Open files, current working directory etc
  - Inherited hardware context for child process:
    - Copy memory, etc, from parent process
    - Child process is now ready to run
      - add to scheduler queue

Memory copy is very expensive:
 

- Potentially need to copy the whole memory space

- **Copying entire memory space is wasteful** and not always needed! (E.g. copy entire 200mb program image of Zoom etc). Mostly, only need contents, and PC, register values.

- **Give Rise to COW**. (copy on write)

### Memory Copy Operation

- If child just read from location, unchanged, just use a shared version.
- **Only when write is perform on a location, then two independent copies needed.**
- **Copy on Write** is possible optimization, only duplicate "memory location" when it is written to, otherwise parent and child share same "memory location".
- Note: memory organized into memory pages (consec range of mem locations), memory managed on page level instead of individual location.

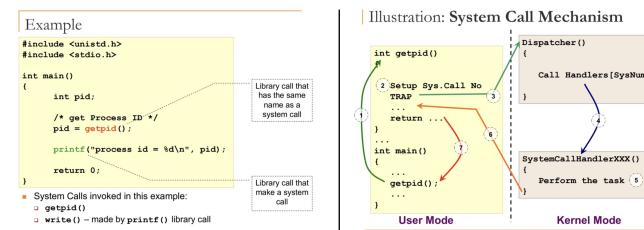
## System Calls (Process Interaction with OS)

### API to OS: Application Program Interface to OS

- OS API provides way of calling facilities/services in kernel.
- **Not same as normal function call**: Change from *user mode to kernel mode*.
- Different OS have different APIs: Unix Variants most follow POSIX standards, small no. of calls 100. Windows family uses Win API across diff. windows, huge no. of calls 1000.

### Unix System Calls in C/C++ program

- In C/C++ program, **system call can be invoked almost directly**, as library version very closely reflects these calls.
- Majority of system calls have library version with **same name** and parameters, library version acts as **function wrapper**.
- A few library functions present more user friendly version, e.g. less no./more flexible parameters). Library version acts as **function adapter**.



### General System Call Mechanism

#### General System Call Mechanism

1. User program invokes the library call
  - Using the normal function call mechanism as discussed
2. Library call (usually in assembly code) places the **system call number** in a designated location
  - E.g. Register
3. Library call executes a special instruction to switch from user mode to kernel mode
  - That instruction is commonly known as **TRAP**
4. Now in kernel mode, the appropriate system call handler is determined:
  - Using the system call number as index
  - This step is usually handled by a **dispatcher**
5. System call handler is executed:
  - Carry out the actual request
6. System call handler ended:
  - Control return to the library call
  - Switch from kernel mode to user mode
7. Library call return to the user program:
  - via normal function return mechanism

## Exception and Interrupt

### Exception

- **Executing machine level instruction** can cause **exception**. For example:

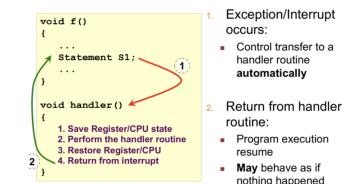
- Arithmetic Errors: Overflow, Underflow, Division by Zero
- Memory Accessing Errors: (accessing memory not belonging to program), Illegal memory address, mis-aligned memory access.

- Exception is **Synchronous**: Determinate, occurs due to program execution at exact points of time.
- **Effect of Exception**: Have to execute **exception handler**, which is similar to a "**forced function call**"!
- Exception ≠ Interrupt!

### Interrupt

- **External events** can interrupt the execution of a program.
- Interrupt request lines connected to CPU, lines are checked during instruction execution cycle.
- Usually hardware related, e.g.: Timer, Mouse move, Keyboard press etc
- Interrupt is **asynchronous**: Events occurs independent of program execution.
- **Effect of interrupt**: Program execution, suspended, execute an interrupt handler.

#### Exception/Interrupt Handler: Illustration



### Summary

- Using **process as an abstraction** of running program.
- Includes necessary information (environment) of execution, Memory, Hardware and OS contexts.
- **Process from OS perspective**: PCB and process table
- **How OS & Process interact**: System calls, Exception / Interrupt

### REFER TO TEXTBOOK:

#### • Modern Operating System (3rd Edition)

- Section 2.1: Processes
- Section 2.4: Process Scheduling
- Section 2.2: Threads

#### • Operating System Concepts (8th Edition)

##### • Section 3.1

### 3. Process Scheduling

A multiprogrammed computer frequently has multiple processes/threads computing for CPU at the same time. Occurs whenever  $\geq 2$  simultaneously in ready state.

- **Scheduler:** Part of OS to decide which process to run next.
- **Scheduling Algorithm:** Algo used.
- **Scheduling Problem:** Choosing, ready process  $>$  available CPUs.
- In addition to picking right process, need **efficient use of CPU** as process switching is expensive. (Switch user to kernel mode, save state of process, memory map, memory cache may need to reload, etc.)
- **I/O Input/Output** (disk or network): When process enters blocked state waiting for external device to complete work.

#### Process Behavior

- Process' unique **requirement of CPU time**.
- Process goes through phases of CPU-activity & IO-activity.
- **Compute/CPU-Bound Process** (computation, e.g. number crunching) vs. **IO-bound Process** (e.g. read/write to file, print to screen)

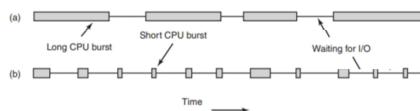


Figure 2-39. Bursts of CPU usage alternate with periods of waiting for I/O.  
(a) A CPU-bound process. (b) An I/O-bound process.

#### Process Environment

- **Batch Processing:** No user, no interaction / responsiveness required.
- **Interactive / Multiprogramming:** Active user interacting with system. Need responsive, consistent in response time.
- **Real Time Processing:** Deadline to meet, usually periodic process.

#### Scheduler Evaluation Criteria

- Many criteria to evaluating algo, largely influenced by p. environment. May be conflicting.
- **All Systems:**
  - **Fairness:** CPU time (per process basis / per user basis). **No starvation.**
  - **Balance:** All parts of computing system utilized.
- **Batch Systems:**
  - **Throughput:** Maximize jobs per hour.
  - **Turnaround time:** Minimize time btwn. submission & termination.
  - **++(Waiting Time):** Related to turnaround, time waiting for CPU
  - **CPU utilization:** keep CPU busy all the time.
- **Interactive Systems:**
  - **Response time:** respond to requests quickly.
  - **Proportionality:** meet users' expectations.
- **Real-time Systems:**
  - **Meeting deadlines:** avoid losing data. (e.g. livestream)
  - **Predictability:** avoid quality degradation in multimedia.

#### Concurrent Execution

- **Concurrent Processes:** Logical concept to cover multitasked processes.
- **Virtual parallelism:** Illusion of parallelism (pseudo-parallelism)
- **Physical parallelism:** Multiple CPUs/Cores, multiple parallel exec

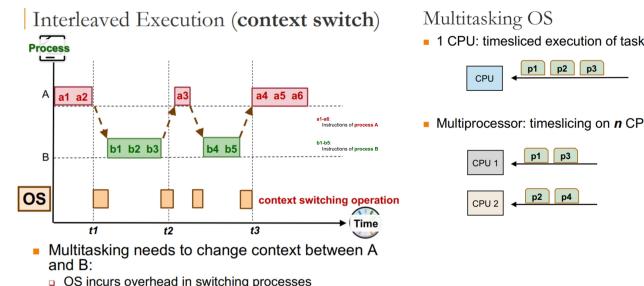
#### When to Schedule

Key issue, when to make schedule decisions.  
E.g. When new process created, run parent or child. When process exits, when process blocks on I/O, or I/O interrupt.

- **Non-preemptive (cooperative):** Process stays scheduled (running state) until it blocks, or gives up CPU voluntarily.
- **Preemptive (Fixed Quota):** Process given fixed time quota to run (possible to block / give up early). At end of quota, running process suspended, another picked if available.

#### Interleaved Execution (Timeslicing)

- **Concurrent Execution on 1 CPU:** Interleave instructions from both processes.
- **OS overhead:** Multitasking needs to change context between programs, incurs overhead.



#### Scheduling a Process:

1. Scheduler triggered (OS takes over)
2. If Context Switch needed, save context, place on blocked/ready queue.
3. Pick suitable process  $P$  to run base on scheduling algo.
4. Setup context for  $P$ .
5. Let process  $P$  run.

#### Scheduling in Batch Systems

**Environment:** No user interaction, non-preemptive scheduling predominant.

Scheduling algorithms generally easier to understand and implement, with variants and improvements for other type of system.

**Criteria:** Turnaround time (related to waiting time, time spent waiting for CPU). Throughput. CPU utilization.

#### Batch Systems Scheduling Algorithms

##### FCFS: First-Come First-Served

- Tasks stored on **FIFO queue based on arrival time**.
- Pick first task in queue to run until done / blocked. Blocked task removed from FIFO queue, when ready, place at back of queue ("newly arrived").
- **Evaluation:**
  - **No starvation.** (Every task eventually processed)
  - **Covoy Effect.** (CPU-Bound followed by IO-Bound tasks heavily inefficient.) Simple reordering can reduce average waiting time.

##### SJF: Shortest Job First (Nonpreemptive)

- Select task with **smallest total CPU time**.
- Need to know total CPU time for task in advance.
- Possible to guess future CPU time by previous CPU-bound phases.

Common approach (Exponential Average):

$$\text{Predicted}_{n+1} = \alpha \text{Actual}_n + (1-\alpha) \text{Predicted}_n$$

- **Actual<sub>n</sub>** = The most recent CPU time consumed
- **Predicted<sub>n</sub>** = The past history of CPU Time consumed
- $\alpha$  = Weight placed on recent event or past history
- **Predicted<sub>n+1</sub>** = Latest prediction

##### Evaluation:

- **Starvation Possible** (Biased towards short jobs) - **Minimize average waiting time.** - Optimal only when all jobs available simultaneously.

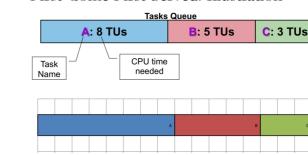
##### SRT: Shortest Remaining Time Next (Preemptive)

##### Preemptive ver of SJF.

- Scheduler chooses process whose remaining run time is shortest.
- When new job arrives, total time compared to current process' remaining (or expected) time.

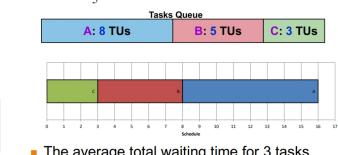
#### Batch System Scheduling Examples

##### First-Come First-Served: Illustration



- The average total waiting time for 3 tasks
- $(0 + 3 + 8)/3 = 3.66$  Time Units

##### Shortest Job First: Illustration



- The average total waiting time for 3 tasks
- $(0 + 8 + 13)/3 = 7$  Time Units
- Can be shown that SJF guarantees smallest average waiting time

## Scheduling in Interactive Systems

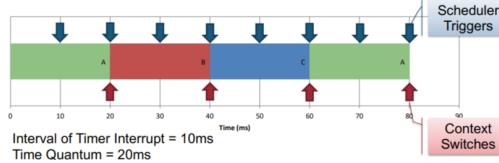
**Criteria:** Response time (time btwn request & response). Predictability (Less variation in response time).

**Environment:** User interaction. Preemptive scheduling used to ensure good response time. Scheduler needs to run periodically.

**Ensuring Periodic Scheduler:** Use timer **interrupts** (based on hardware clock). OS ensures timer interrupt cannot be intercepted by any other program. Interrupt handle **invokes scheduler**.

- **ITI: Interval of Timer Interrupt:** Timing Interval that interrupt happens, and OS scheduler triggered. Typical values (1ms - 10ms).
- **Time Quantum:** Execution duration given to each process, constant / variable. Must be multiples of ITI. Typical values (5ms - 100ms).

Illustration: ITI vs Time Quantum



## Interactive Systems Scheduling Algorithms

### RR: Round Robin

- Preemptive ver. of FCFS.
- **Tasks stored in FIFO queue.** Each process assigned time quantum.
- If task still running at end of quantum / blocked / gives up CPU voluntarily, CPU preempted, given to another process.
- Task placed at end of queue to wait for another turn.
- **Response time guarantee:**  $n$  tasks,  $q$  quantum. Time before task gets CPU bounded by  $(n - 1)q$ .
- **Evaluation:**
  - Too short quantum = many process switches, lower CPU efficiency
  - Too long quantum causes poor response to short interactive requests.

### Priority Scheduling: Priority Based

- Each process assigned a priority, runnable process with highest priority allowed to run.
- **Preemptive ver.:** Higher p. process can preempt running low p. process.
- **Non-preemptive ver.:** Late high p. process wait for next scheduling.
- **Evaluation:**
  - **Possible Starvation:** High p. process may hog CPU. To prevent this, scheduler may **decrease priority** of currently running process at each clock tick (clock interrupt). Or, give some max time quantum, when used up, allow next in line to run.
  - **Priority Inversion:** When lower priority task preempts higher priority task. (If lower priority task locks some resource, e.g. file, gets switched, higher priority task cannot run)

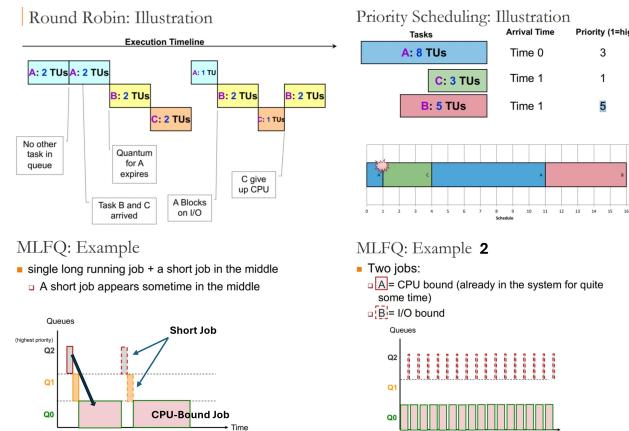
### MLFQ: Multi-Level Feedback Queue

- More efficient to give CPU-bound process large quantum once in a while, rather than small quanta frequently to reduce swapping. Also, giving all processes large quantum means poor response time. Hence, solution to set **multiple priority queues**.
- As (CPU-bound) process sinks deeper into priority queues, run less frequently, saving CPU for short, interactive processes.
- Hence, **adaptively learn process behavior**, minimize both:
  - Min. Response time for IO bound processes
  - Min. Turnaround time for CPU bound processes.
- **MLFQ Rules:**
  - **Basic Rule:** Higher priority process runs. If same p, run in RR.
  - **Priority Setting:** New job given highest priority. If job fully utilize time slice, priority reduced. If job gives up / blocks before time slice finished, priority retained.
- **Evaluation:** - **Can be gamed:** User typing carriage returns at random every few seconds doing wonders for his response time.

### Lottery Scheduling

- Give processes lottery tickets for system resources. When scheduling, chose ticket at random.
- "All processes equal, some processes more equal." More important processes given extra tickets.
- **Evaluation:**
  - **Responsive:** New process can participate in next lottery.
  - **Good Control:** Each process can distribute to child processes proportionally w.r.t. need.

### Interactive System Scheduling Examples



### Others

- **Shortest Process Next:** Estimate (using calculated aging).
- **Guaranteed / Fair-Share Scheduling:** Track CPU usage, run accordingly. Each user gets agreed allocation of CPU.

## 4. Threads

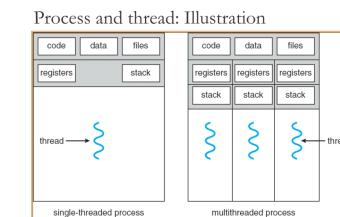
In trad. OS, each process has an address space and single thread of control. Desirable instead to have multiple threads of control in the same address space running in quasi-parallel, as though (almost) separate processes (except for shared address space).

### Motivation for Thread

- **Process is Expensive:** Process creation under `fork()` model duplicate memory space, most of process context.
- Context switch also requires overhead, save/restore process info.
- **Communication:** Hard for independent processes to communicate with each other. Independent memory space, no shared variable, requires IPC.
- **Thread:** "quick hack" into popular mechanism.
- **Basic Idea:** Traditional process only single thread of control. (Only one instruction executing at a time).
- Add more threads (of control) to same program, multiply parts of programs executing at same time conceptually.

### Process and Thread

- **Multi-threaded Process:** Single process can have multiple threads.
- Threads in same process share same: *Memory Context* (text, data heap), *OS Context* (Process id, files etc.)
- **Unique info per thread:** Id (thread id), Registers (GPR and special), "Stack".
- **Process Context Switch vs. Thread Switch:**
  - **Process Context Switch:** Switch OS, Hardware, Memory Context.
  - **Thread Switch (same process):** Just hardware context (registers, stack).
  - Thread is **lightweight process**.



### Threads: Benefits

- **Resource Sharing:** Ability for parallel entities to share address space and all of data. No need add. mechanism for passing info.
- **Economy:** Threads are lighter weight than process, faster to create and destroy. Less resources to manage.
- **Responsiveness:** No performance gain when all CPU bound, but when substantial computing and substantial I/O, threads allows activities to overlap, speeding up application.
- **Scalability:** Multithreaded program can take advantage of multiple CPUs.

### Threads: Problems

- **System Call Concurrency:** Parallel execution of multiple threads: parallel system call possible. Need to guarantee correctness and determine correct behavior.
- **Process Behavior (OS dependent):** Impact on process operations, (e.g. If one thread executes exec() / exit(), what about other threads / whole process).

## Thread Models (ways to support threads)

- User Thread:** Thread implemented as a **user library**. (Runtime system (in the process) will handle thread related operation.)
- Kernel not aware of threads in process.
- Kernel Thread:** Thread implemented in the OS. Operation handled as system calls.
- Kernel thread-level scheduling possible, where kernel schedule by threads instead of by process.
- Kernel may make use of threads for own execution.
- Threads on Modern Processor:** Threads started as software mechanism, now exists hardware support on modern processors (multiple register sets on same core, *simultaneous multi-threading (SMT)*, "Hyperthreading" on Intel).

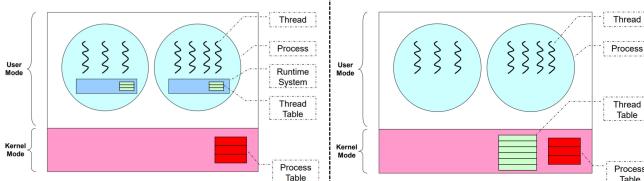
## User Thread

- Advantages:** Can have any multithreaded program on any OS, thread operations are just library calls, generally more configurable and flexible (e.g. customized thread scheduling policy.)
- Disadvantages:** OS not aware of threads, scheduling is performed at process level. (One thread blocked, process blocked, all threads blocked). Cannot exploit multiple CPUs.

## Kernel Thread

- Advantages:** Kernel can schedule on thread levels: > 1 thread in same process can run simultaneously on multiple CPUs.
- Disadvantages:** Thread operations now system calls, slower and more resource intensive. Generally less flexible, used by all multithreaded programs, so too many / few features, overkill / not flexible enough for different programs.

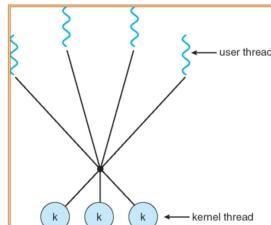
### User Thread: Illustration



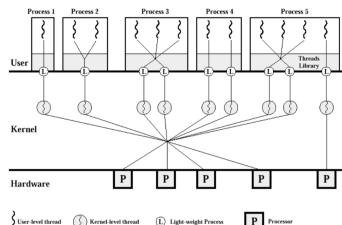
## Hybrid Thread Model

- Have both Kernel and User threads, where OS schedules on kernel threads only, and user thread binds to a kernel thread.
- Offers flexibility, limit concurrency of any process / user.

### Hybrid Thread Model



### Hybrid Model Example: Solaris



## POSIX Threads

IEEE defined standards (1003.1c) for threads, called *Pthreads*, most UNIX systems support it. Defines over 60 function calls.

### Pthread

- POSIX:** Portable Operating System Interface. (family of standards)
- Defines the API as well as the behavior, but not implementation.
- pthread can be implemented as user / kernel thread.

#### Basics of pthread

- Header File:** #include <pthread.h>
- Compilation (flag is system dependent):** gcc XXXXX.c -lpthread
- Useful datatypes:**
  - #pthread\_t: Data type to represent a thread ID(TID)
  - #pthread\_attr\_t: Data type to represents attributes of a thread

#### pthread Creation Syntax

```
int pthread_create(pthread_t *tidCreated,
                  const pthread_attr_t *threadAttributes,
                  void (*startRoutine)(void *),
                  void *argStartRoutine);
```

▪ Returns (0 = success; 10 = errors)

- Parameters:**
  - #tidCreated: Thread id for the created thread
  - #threadAttributes: Control the behavior of the new thread
  - #startRoutine: Function pointer to the function to be executed by thread
  - #argStartRoutine: Arguments for the startRoutine function

#### pthread Termination Syntax

```
int pthread_exit(void *exitValue);
```

- Returns:**
  - #exitValue: Value to be returned to whoever synchronizes with this thread (more later)

- If pthread\_exit() is not used, a thread will terminate automatically at the end of the startRoutine.**
  - If a 'return XYZ;' statement is used, then 'XYZ' is captured as the exitValue
  - Otherwise, the exitValue is not well defined

#### pthread Creation & Termination: Example

```
//header files not shown
void* sayHello( void* arg ) {
    // Function to be executed
    // by a pthread
    printf("Just to say hello!\n");
    pthread_exit( NULL );
}

int main()
{
    pthread_t tid;
    // Pthread Creation
    pthread_create( &tid, NULL, sayHello, NULL );
    printf("Thread created with tid %i\n", tid);
    return 0;
}
```

#### pthread: Sharing of memory space

```
//header files not shown
int globalVar;
// Variable shared between pthreads
void* doSum( void* arg )
{
    int i;
    for ( i = 0; i < 1000; i++ )
        globalVar++;
}

int main()
{
    pthread_t tid[5]; //5 threads id
    int i;
    for ( i = 0; i < 5; i++ )
        pthread_create( &tid[i], NULL, doSum, NULL );
    printf("Global variable is %i\n", globalVar);
    return 0;
}
```

#### Bc. Race Condition, Sum unpredictable

```
#include <stdio.h>
#include <pthread.h>
int globalVar;

void* doSum( void* arg){
    int i;
    for ( i = 0; i < 1000; i++ ) {
        globalVar++;
    }
}

int main() {
    pthread_t tid[5]; //5 threads id
    int i;

    for ( i = 0; i < 5; i++ ) {
        pthread_create( &tid[i], NULL, doSum, NULL );
    }

    // Wait for all threads to finish
    for ( i = 0; i < 5; i++ ) {
        pthread_join( tid[i], NULL );
    }

    printf("Global variable is %i\n", globalVar);
    return 0;
}
```

#### pthread Simple Synchronization - Join

```
int pthread_join( pthread_t threadID,
                  void **status );
```

#### Pthrcd: A lot more!

- There are more interesting stuff about **pthread**.
  - Yielding (giving up CPU voluntarily)
  - Advanced synchronization
  - Scheduling policies
  - Binding to kernel threads
  - Etc
- As we cover new topics, you can explore the **pthread** library to see the application!

## Summary of Pthread

- All Pthread threads have certain properties:** each one has identifier, set of registers (including program counter), set of attributes which are stored in a structure.
- Attributes include stack size, scheduling params etc.
- pthread\_create call:** New thread created, returns thread identifier of newly created thread as function value.
- pthread\_exit:** Thread terminate by calling, stops thread and releases stack.
- pthread\_join:** Thread needing to wait for another thread to finish work and exit before continuing can call to wait for specific other thread to terminate. (TID) of thread to wait for given as parameter.

Thread call	Description
Pthread_create	Create a new thread
Pthread_exit	Terminate the calling thread
Pthread_join	Wait for a specific thread to exit
Pthread_yield	Release the CPU to let another thread run
Pthread_attr_init	Create and initialize a thread's attribute structure
Pthread_attr_destroy	Remove a thread's attribute structure

Figure 2-14. Some of the Pthreads function calls.

## 5. Inter-Process Communication (IPC)

Given that processes frequently need to communicate, some need for communication, preferably in some well-structured way not using interrupts. Consider that cooperating processes have independent memory space, making IPC necessary.

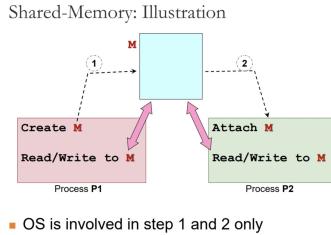
### Common Communication Mechanisms

#### Shared Memory

- General Idea:** Process 1 creates shared memory region  $M$ , process 2 attaches  $M$  to its own memory space. Both processes can now communicate using memory region  $M$ . Applicable to multiple processes.
- $M$  behaves similar to normal memory region, any writes to region can be seen by all other parties.

POSIX Shared Memory in \*nix

- Basic steps of usage:
  - Create/locate a shared memory region  $M$
  - Attach  $M$  to process memory space
  - Read from/Write to  $M$ 
    - Values written visible to all process that share  $M$
  - Detach  $M$  from memory space after use
  - Destroy  $M$ 
    - Only one process need to do this
    - Can only destroy if  $M$  is not attached to any process



#### Advantages:

- Efficient, as only initial steps of create & attach  $M$  involves OS.
- Ease of use, as  $M$  behaves as normal memory space, info of any size/type writable easily.

#### Disadvantages:

- Synchronization: Shared resource, need to sync access still.
- Implementation is harder.

#### Example: Master program

```
#include <sys/types.h>
#include <sys/ipc.h>
#include <sys/shm.h>
#include <sys/conf.h>
#include <sys/conf.h>

int main()
{
    int shmid, i, *shm;
    Step 1. Create Shared Memory region.
    shmid = shmget( IPC_PRIVATE, 40, IPC_CREAT | 0600 );
    if (shmid == -1) {
        perror("Cannot create shared memory!\n");
        exit(1);
    } else
        printf("Shared Memory Id = %d\n", shmid);

    shm = (int*) shmat(shmid, NULL, 0); Step 2. Attach Shared Memory region.
    if (shm == (int*) -1) {
        perror("Cannot attach shared memory!\n");
        exit(1);
    }

    shm[0] = 0;
    while (shm[0] == 0) { Step 3. Values produced by the slave program.
        sleep(3);
    }

    for (i = 0; i < 3; i++) {
        printf("Read %d from shared memory.\n", shm[i+1]);
    }

    shmdt((char*) shm); Step 4+5. Detach and destroy Shared Memory region.
    shmctl(shmid, IPC_RMID, 0);
    return 0;
}
```

#### Example: Slave program

```
/*similar header files
int main()
{
    int shmid, i, input, *shm;
    Step 1. By using the shared memory region id directly, we skip shmat() in this case.
    printf("Shared memory id for attachment: ");
    scanf("%d", &shmid);

    shm = (int*) shmat(shmid, NULL, 0); Step 2. Attach to shared memory region.
    if (shm == (int*) -1)
        perror("Error: Cannot attach!\n");
        exit(1);

    for (i = 0; i < 3; i++) {
        scanf("%d", &input); Step 3. Values into shm[i+1]
        shm[i+1] = input;
    }

    shm[0] = 1; Let master program know we are done!
    shmdt((char*)shm); Step 4. Detach Shared Memory region.
    return 0;
}
```

### Message Passing

- General Idea:** Process 1 prepares, send message  $M$ . Process 2 receives it.
- Message sending / receiving usually provided as system calls.
- Additional properties:** Naming (identify other party), synchronization (behavior of send/rec ops).
- Msg has to be stored in kernel memory space. Each send/rec op needs to go through OS (system call).

#### Direct Communication (Naming Scheme)

- Sender/Receiver of message **explicitly names other party**.
- E.g. Send(P2, Msg) & Receive(P1, Msg).
- One link per pair of process, need know identity of other party.

#### Indirect Communication (Naming Scheme)

- Messages sent / received from message storage, aka **mailbox / port**.
- E.g. Send(MB, Msg) & Receive(MB, Msg).
- One mailbox can be shared among a number of processes

#### Synchronization Behaviors

- Blocking Primitives (Synchronous):** Sender blocked until message received, receiver blocked until message has arrived.
- Non-Block Primitives (Asynchronous):** Sender resume operation immediately, receiver either receive, or indicate message not ready yet.
- Advantages:**
  - Portable, easily implemented on diff. processing env.
  - Easier synchronization: Esp. when synchronous primitive used.
- Disadvantages:**
  - Inefficient, require OS intervention
  - Harder to use, less flexi, messages limited in size / format.

### Pipe (Unix Specific)

- Unix process has 3 default communication channels:
  - **stdin** (standard in): Commonly linked to keyboard input.
  - **stderr** (standard error): Only used print out error msg.
  - **stdout** (standard out): Commonly linked to screen.

#### Piping in Shell

- For example ("A | B"):
 
 A diagram showing two boxes labeled 'A' and 'B'. Box A has an arrow labeled 'stdout' pointing to a pipe. Box B has an arrow labeled 'stdin' pointing from the pipe. The pipe is labeled 'Unix Pipe'.
  - The first element in the shared memory region is used as "control" value in this example (0: values not ready, 1: values ready).
  - The next 3 elements are values produced by the slave program.
- General Idea:
  - A communication channel is created with 2 ends:
    - 1 end for reading, the other for writing
    - Just like a water pipe in the real world
  - The piping "|" in shell is achieved using this mechanism internally

- Unix Shell Piping:** Unix shell provides | symbol to link input/output channels of one process to another, aka *piping*.

#### Unix Pipes: as a IPC Mechanism

- Process P Write → [d c b a] → Read → Process Q
- A pipe can be shared between two processes
- A form of Producer-Consumer relationship
  - P produces (writes)  $n$  bytes
  - Q consumes (reads)  $m$  bytes
- Behavior:
  - Like an anonymous file
  - FIFO → must access data in order

#### Unix Pipes: Semantic

- Pipe functions as **circular bounded byte buffer** with **implicit synchronization**:
- Writers **wait** when buffer is **full**
- Readers **wait** when buffer is **empty**
- Variants:
  - Can have multiple readers/writers
    - The normal shell pipe has 1 writer and 1 reader
  - Depends on Unix version, pipes may be **half-duplex**
    - unidirectional: with one write end and one read end
  - Or **full-duplex**
    - bidirectional: any end for read/write

### Unix Pipe: System Calls

```
Header File #include <unistd.h> Syntax int pipe( int fd[] )
```

#### Returns:

- 0 to indicate success; !0 for errors
- An array of file descriptors is returned:
  - fd[0] == reading end fd[1]
  - fd[1] == writing end Write → data → Read

#### Unix Pipes: Example Code

```
#define READ_END 0
#define WRITE_END 1

int main()
{
    int pipeFd[2], pid, len;
    char buf[100], *str = "Hello There!";
    pipe( pipeFd );
    if ((pid = fork()) > 0) { /* parent */
        close(pipeFd[READ_END]);
        write(pipeFd[WRITE_END], str, strlen(str)+1);
        close(pipeFd[WRITE_END]); /* child */
        len = read(pipeFd[READ_END], buf, sizeof buf);
        printf("Proc %d read: $s\n", pid, buf);
        close(pipeFd[READ_END]);
    }
}
```

### Signal (Unix Specific)

- Form of IPC. **Asynchronous notification** regarding an event. Sent to process / thread.
- Recipient of signal must handle signal by: Default set of handlers or user supplied handler.
- Common UNIX signals:** Kill, Stop, Continue, Arithm. error etc.

#### Example: Custom Signal Handler

```
#include <stdio.h>
#include <signal.h>
#include <unistd.h>

void myOwnHandler( int signo )
{
    if (signo == SIGSEGV) {
        printf("Memory access blows up!\n");
        exit(1);
    }
}

int main(){
    int *ip = NULL;

    if (signal(SIGSEGV, myOwnHandler) == SIG_ERR)
        printf("Failed to register handler\n");

    *ip = 123; This statement will cause
    return 0; a segmentation fault.
}
```

User defined function to handle signal. In this example, we handle the "SIGSEGV" signal, i.e. the memory segmentation fault signal.

Register our own code to replace the default handler.

End.