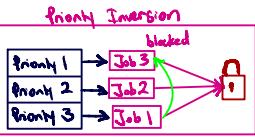


**Scheduling**

Scheduling decides which threads are given access to resources moment to moment  
 Goals: 1) ↓ Response Time 2) ↑ Throughput (operations/sec), ↓ overhead 3) ↑ Fairness  
 Waiting Time: time before it got scheduled  
 Completion Time: Waiting time + running time  
 Priority handles differences in importance, watch Starving  
 Priority Inversion: high priority task blocked waiting on low priority thread

**Choose**

- FCFS / FIFO
- SRTF Approx.
- Linux CFS
- Round Robin
- EDF
- Priority

**For**

- CPU Throughput
- Avg Response Time
- I/O Throughput
- Fairness (CPU Time)
- Fairness (CPU Wait Time)
- Meeting Deadlines
- Favor important

**Policies****1) First Come, First Served (FCFS/FIFO)**

Idea: One program scheduled until done

Pro: Least overhead, simple

Con: Convo effect (short processes stuck behind large ones)

**3) Shortest Job First (SJF), Shortest Remaining Time First (SRTF)**

Idea: Run job with least amount of computation to do

Pro: Optimal!!

Con: Need to be able to see future, know process length

**5) Multiple-Level Feedback Scheduling**

Idea: Multiple queues, adjusts queue as process is run  
 have queues w/ fixed priority scheduling, time slice

Pro: Approximates SJF

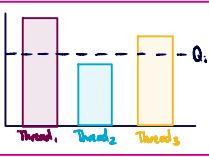
Con: Can counter by regaining I/O and staying in highest

**7) Stride Scheduling, Linux Completely Fair Scheduler (CFS)**

Idea: Track virtual CPU time per thread, gets equal share, choose thread w/ least CPU time

Basic equal share:  $Q_i = \text{Target Latency} \cdot \frac{1}{N}$  threads Target Latency: period of time where every service gets service

Weighed share:  $Q_i = \left( \frac{w_i}{\sum w_p} \right) \text{Target Latency}$  Add min granularity to ensure each process gets to run, min time slice

**Deadlocks**

**Deadlock:** cyclic waiting for resources, deadlock  $\rightarrow$  starvation, starvation  $\not\rightarrow$  deadlock

**Requirements for Deadlock**

- 1) Mutual Exclusion and bounded resources  
     -one thread at a time use resources
- 2) Hold and wait  
     -thread holding resource waits to acquire more
- 3) No preemption  
     -resource released voluntarily
- 4) Circular Wait  
     -set of waiting threads waiting on each other

**Deadlock Prevention**

- 1) Provide sufficient resources, VM unlimited
- 2) Abort requests, acquire atomically
- 3) Fail if waiting too long, force give up
- 4) Order resources usage in same order

**Deadlock Avoidance**

- Prevent system from entering unsafe state  
 ↳ Use Banker's Algorithm
- Safe space: can prevent by delaying acquisition  
 unsafe space: can unavoidably lead to deadlock, with certain acquisition
- Deadlocked state: exists a deadlock

**Deadlock recovery****Deadlock detection****Banker's Algorithm**

- Check if resource request leads to unsafe state
- State max resource needs in advance
- Allow thread to continue if available resources - # requested  $\geq$  max

**State**

Idea:  
 Allocate resources dynamically

- Evaluate each request & grant if some ordering is deadlock free
- Pretend request granted, run deadlock detection algo

**Deadlock detection algo [Banker]**

Add to unfinished  
 for each thread unfinished  
 if (request  $[max - Alloc] \leq Avail]$ )  
 remove from unfinished  
 Avail = Avail + Alloc

**Memory**

Virtual Memory to multiplex memory, protection, controlled overlap

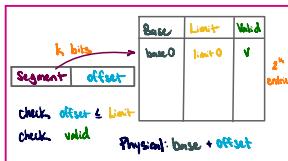
Pages: small fixed size physical memory chunks

Page Table: one per process, has physical page #, permission (R/W/Valid)

Memory Management Unit (MMU): Translation box converts between virtual & physical address; kernel handles evicting, invalidating, disk

## 1) Base & Bound / Segment Mapping

- Idea: set registers w/ base and limit
- Pro: Simple
- Con: Internal and External Fragmentation, no sparse address space support or interprocess sharing



## 2) Inverted Page Table

- Idea: use a hash table to map VPN to PPN

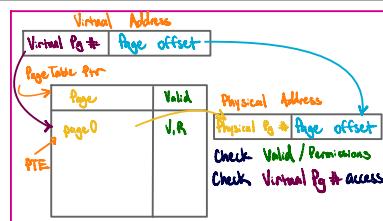
Pro: Efficient lookup

Con: complexity of hash chains, poor cache locality

## 3) Simple Paging

- Idea: Translations in Page Table

- Pro: Able to share memory, point to same physical page #, easy to reallocate memory
- Cons: Page Table too big, Internal fragmentation



## 4) Multilevel Page Table

- Idea: Tree of Page Tables w/ fixed size, Save Page Table Ptr (ATB)

Virtual Pg # Virtual Pg # Page offset (10b-10b-12b)

- Pros: allocate just needed PTE, easy memory allocation, sharing
- Cons: one pointer per page, ≥ 2 lookups per reference

**Page Table Entry (PTE):** pointer to next level page table or actual page, permission bits

## Caching

Translation Lookaside Buffer (TLB): cache for translations

records just end result, recent VPN to PPN, include ProcessID, hardware

**Temporal Locality:** Time locality, recently accessed closer

**Spatial Locality:** Space locality, contiguous blocks

## Sources of Cache Misses

1) Compulsory: first access to block

↳ Clustering, Working set tracking

2) Capacity: cache cannot contain all blocks accessed

↳ increase cache size

3) Conflict: multiple mem location mapped to same cache location

↳ increase cache sizes, increase associativity

4) Coherence: other process updates memory

## Types of Caches

1) Direct Mapped Cache: single block for set, index

Valid	Tag	Data
-------	-----	------

Cache typically physically indexed

Can lookup TLB and Cache simultaneously



## Cache



Block: minimum quantum of caching

Index: lookup candidates in cache, identify set

Tag: identify actual copy

Write Through: info written to both block and lower level memory

Write Back: info written only to block, write when evict

Zif distribution: increasing size of cache has diminishing returns

## Average Memory Access Time (AMAT)

$$AMAT = \text{Hit Rate} \times \text{Hit Time} + \text{Miss Rate} \times \text{Miss Time}$$



## Demand Paging

Demand Paging: only keep active pages in memory, as cache: fully associative, LRU, write back, 1 pg block

If invalid PTE:

- 1) MMU traps to OS w/ Page Fault
- 2) Find & replace page w/ page from disk
- 3) Reset Page Table & restart instruction

## Page Replacement Policies

- 1) FIFO: evict oldest page  
Con: evicts heavily used pages
- 2) RANDOM: choose random page for replacement  
Con: unpredictable
- 3) MIN: replaces page not used for longest time, optimal  
Con: don't know future
- 4) LRU: replace page not used for longest time  
Con: too much overhead

Freelist: keep set of free pages by Clock Algorithm

Working Set: group of pages accessed by process recently

Swapping: Some or all of previous process moved to disk to make room  
- Can share code segment, setting read-only

## 5) Second Chance

Split into Active and Second Chance List

Pro: few disk accesses

Con: increased overhead trapping

Approx LRU:

6) Clock Algorithm: Replace an old page, partition into old and young

7) Nth Chance: N chances to stay in memory

## Allocation of Memory for Processes:

Equal allocation, proportional allocation, priority allocation

- can set lower and upper bound for memory

Thrashing: busy swapping pages in and out w/ little progress w/o enough pages

## Operating System Overview

- Purpose: Special layer of software that provides application software access to hardware resources
- 1) Illusionist: Provide simple abstractions of physical resources (infinite memory, virtualization)
  - 2) Referee: Manage protection, isolation, and sharing of resources (resource allocation, communication)
  - 3) Glue: Common Services (Storage, Networking, sharing, look and feel)

## Four Fundamental OS Concepts

### 1) Threads

- Single unique execution context
- has own Program Counter, Registers, Execution Flag, Stack, Memory State
- When executing and resident on processor: running
- When not loaded in: suspended
- In order to execute multiple processes, multiplex in time, virtual cores (TCB)
- Store other threads in Thread Control Block

### 2) Addresses

- Address Space: the set of accessible addresses + state associated w/ them

- OS must protect user programs from one another & protect itself from other programs

- 1) Base & Bound  
have base register and bound register to check address
- 2) Address Space Translation (Page Table)  
Program operates in virtual (pages) translated to memory add

### Abstraction

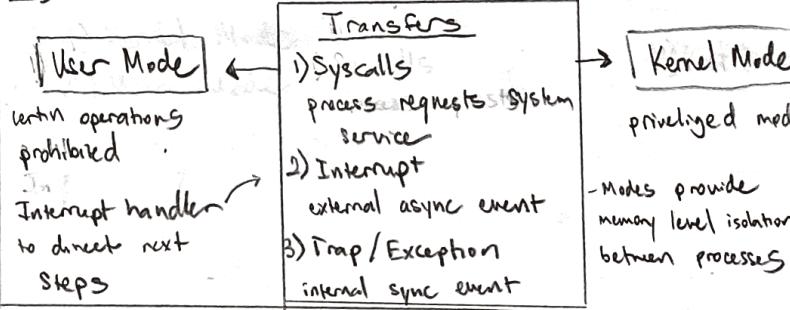
- Processor → Threads
- Memory → Address Space
- Disk, SSD → Files
- Networks → Sockets
- Machines → Processes

### 3) Process

- Execution environment with restricted rights
- Protected Address Space w/ dt threads
- Running program w/ memory protection
- Processes provide protection, isolation, thread provides concurrency



### 4) Dual Mode Operation



## Threads & Processes: programmer POV

- Allow parallel programs to be run
- Multiprocessing: Multiple CPU (cores)
- Multiprogramming: multiple jobs / processes
- Multi threading: multiple threads / processes, same CPU
- Threads have non-determinism: can run in any order, leads to race cond.
- Process fork: copy current process: page table
  - 1) copy, new process has pid 0, parent pid > 0
- Each process/thread has kernel segment with PCB/TCB, kernel stack
  - KSP stores kernel stack pointer in order to reduce I/O blocking in kernel

### Thread States

- Running - currently in CPU
- Ready - eligible, not running
- Blocked - ineligible to run

Process	Thread
Creation	pthread-create()
Page Table	Distinct
Registers, IP	Distinct
Stack	Separate & inaccessible
Heap, static var	Separate
File descriptors	Separate
Synchronization	Wait(), Waitpid() pthread-join(), semaph locks
Overhead	Higher
Protection	Higher

parallel ⇒ concurrent

Concurrent ≠ parallel

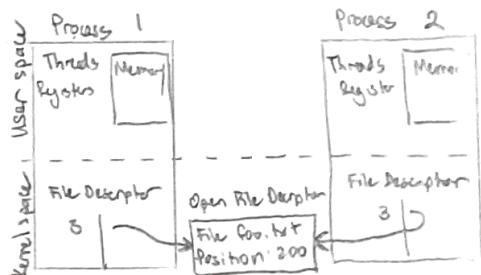
## File I/O, Devices

Everything is a file: open, read, write, close

High level FILE: buffered, has fd, buffer, lock

Low level file: returns fd (not buffered)

Drivers: device specific code in kernel that interacts directly w/ device hardware



### Streams

High Level I/O

Low Level I/O

Syscall

File System

I/O Direct

Disk, Flash

can have same file diff descriptors

## IPC, Pipes, Sockets

### Interprocess Communication (IPC)

- communication between protected environments (processes)

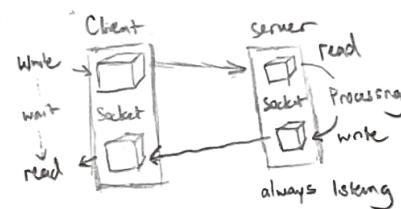
Pipes: act as single queue between processes

- write only on one side, read only on other

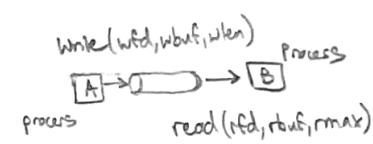
Sockets: allow two queues, communication between

- communication between multiple processes on different machines, socket / bind / connect / listen

### Sockets



### Pipes



## Synchronization

- Many different solutions to fix synchronization issues, wait at least busy waiting
- Atomic Operation: operation that always runs to completion or not at all
- Mutual Exclusion: ensuring only one thread does particular thing at a time, exclude the other
- Critical Section: piece of code only one thread can execute at once
- Locks: synchronization mechanism for enforcing mutual exclusion on critical sections to support atomic operations
- Hardware atomic primitives: disabling interrupts, test & set, swap, compare & swap, load-linked & store conditional
- Separate lock variable, use hardware mechanisms to protect modifications of that var

### Semaphores

- Down() or p(): waits for semaphore to become positive then decrements by 1, like (wait)
- Up() or v(): atomic operation increments by 1, wake up waiting P (signal)

Monitors: a lock and zero or more condition variables managing concurrent access to shared data

- locks for mutual exclusion and condition var scheduling constraint

- condition variable: queue of threads waiting for inside critical section

### Hoare Monitor

```
if (!IsEmpty(&queue))
    cond-wait(&buf_cv, &buf_lock)
```

- Wait(block): Atomically release lock and go to sleep

- Signal(): Wake up one waiter

- Broadcast(): Wake up all waiters

### Mesa Monitor

```
while (!IsEmpty(&queue))
    cond-wait(&buf_cv, buf_lock)
```

Lock: prevent others from changing critical section

acquire(&lock): wait till lock is free, then grab, run critical section  
release(&lock): unlock, wake up anyone waiting

### Implementation:

```
require():
    disable interrupts
    if (value == BUSY)
        put thread wait queue
        go to sleep
    else: value = BUSY
    enable interrupts
```

```
released():
    disable interrupts
    if anyone on waitqueue
        take thread off queue
        place on ready queue
    else: value == 0
    enable interrupts
```

### futex

Kernelspace wait queue attached to user space atomic integer  
faster, no syscalls, FUTEX\_WAIT, FUTEX\_WAKE  
when(fd, count, &len) + 1  
word-count \* NC  
ref-count?

thread-current(): get current thread  
strncpy(dst, src, len): copy from src to dst strnlen() + 1

(x) list-entry(e, word-count-t, elem) list-init(&cv->writers)  
lock\_cv(&cv)

### Reader's/Writers Problem

```
while (...) { test & set (guard))
```