

CS492: Operating Systems Notes

Steven DeFalco

Spring 2023

Contents

1	Introduction	5
1.1	Operating System Concepts and Structure	5
1.1.1	Process vs. Program	5
1.1.2	Address Space	5
1.1.3	Files	5
1.1.4	System Calls	6
1.1.5	Kernels	6
1.2	Computer Hardware	6
1.3	Introduction to Linux	7
1.3.1	Create your own kernel	7
1.4	System Calls	8
1.4.1	To add a syscall in Linux	8
2	Processes and Threads	9
2.1	Introduction to Processes	9
2.2	Processes in Linux	11
2.3	Threads	12
2.3.1	PCB	13
2.3.2	Threads Implemented in User Space	13
2.3.3	Threads Implemented in Kernel Space	14
2.3.4	Context Switching with Threads	14
2.3.5	Thread Local Storage	14
2.4	Threads in Linux	15
3	Scheduling	15
3.1	Introduction to Scheduling	15
3.1.1	Clock Interrupts	15
3.1.2	Non-preemptive Scheduling	15
3.1.3	Preemptive Scheduling	15
3.1.4	Goals of All Scheduling Algorithms	15

3.2	Scheduling in Batch Systems	16
3.2.1	Goals for Scheduling in Batch Systems	16
3.2.2	First-come First-served (FCFS)	16
3.2.3	Shortest Job First (SJF)	16
3.3	Scheduling in Interactive Systems	16
3.3.1	Goals for Scheduling in Interactive Systems	16
3.3.2	Round-Robin (RR)	16
3.3.3	Priority (PRIO) Scheduling	17
3.3.4	Multiple Queues (MQ) Scheduling	18
3.3.5	Other Schedulers for Interactive Systems	18
3.4	Scheduling in Real-time Systems	18
3.4.1	Goals for Scheduling in Real-Time Systems	18
3.4.2	Number of Schedulable Processes	19
3.4.3	Rate Monotonic Scheduling	19
3.4.4	Earliest Deadline First Scheduling	19
4	Concurrency	19
4.1	Inter-process Communication	19
4.2	Race Conditions	20
4.2.1	Requirements to Avoid Race Conditions	20
4.2.2	Disabling Interrupts to Avoid Race Conditions	20
4.2.3	Using Lock Variables to Avoid Race Conditions	20
4.2.4	Strict Alternation to Avoid Race Conditions	21
4.2.5	Peterson's Solution	21
4.2.6	Test-and-Set Lock (TSL) Instruction	21
4.2.7	TSL Solution	21
4.2.8	Producer-Consumer Problem with <code>sleep()</code> and <code>wakeup()</code>	21
4.2.9	Semaphores	22
4.2.10	Types of Semaphores	22
4.2.11	Producer-Consumer Problem using Semaphores	23
4.2.12	Producer-Consumer Problem using Semaphores and Mu- texes	23
4.2.13	Monitors	24
4.3	Readers-Writers Problem	24
5	Memory Management	26
5.1	Memory Management	26
5.1.1	Base and Limit Registers	26
5.1.2	Swapping	26
5.1.3	Memory Fragmentation	27
5.1.4	Determine which Memory is Allocated	27
5.1.5	Picking Location of New Process's Memory	27
5.1.6	Overlays	28

6	Virtual Memory	28
6.1	Paging	29
6.1.1	Translating Addresses when Paging	29
6.1.2	Page Table Entry	29
6.1.3	Page Faults	29
6.1.4	Addresses	30
6.1.5	Translation Lookaside Buffer	30
6.1.6	Multilevel Page Tables	31
6.1.7	Inverted Page Tables	32
7	Page Replacement Algorithm	33
7.1	Optimal Algorithm	33
7.2	Not Recently Used Algorithm	33
7.3	First in First Out (FIFO) Algorithm	34
7.4	Second Chance Algorithm	34
7.5	Least Recently Used (LRU) Algorithm	35
7.5.1	Priority List Implementation of LRU	35
7.5.2	Counter Implementation of LRU	35
7.6	Not Frequently Used (NFU) Algorithm	35
7.7	Aging Algorithm	36
7.8	When to Move Pages into Memory	36
7.8.1	Working Set Pages	36
7.9	Working Set Based Algorithm	36
7.10	Working Set Clock (WSClock) Algorithm	37
7.11	Frame Allocation to Processes	37
7.11.1	Fixed Frame Allocation	38
7.11.2	Priority Allocation	38
7.11.3	Local Replacement	38
7.11.4	Global Replacement	38
8	File Systems	39
8.1	Disks	39
8.2	Allocation Methods	39
8.2.1	Contiguous Allocation	40
8.2.2	Linked List Allocation	40
8.2.3	Linked List Allocation with Table	40
8.2.4	Indexed Allocation	41
8.3	Implementation of Directories	41
8.3.1	Symbolic Link	42
8.4	Virtual File System	42
9	Input/Output	43
9.1	I/O Hardware	43
9.2	Communication Mechanisms Between CPU and I/O Devices	43
9.2.1	Register Based I/O	44
9.2.2	Memory Mapped I/O	44

9.2.3	Hybrid I/O Implementation	44
9.2.4	Direct Memory Access (DMA)	45
9.3	Goals of I/O software	46
9.3.1	Programmed I/O	47
9.3.2	Interrupt-driven I/O	47
9.3.3	I/O using DMA	48
10	Disks	48
10.1	RAID	49
10.1.1	RAID 0 (Striping)	49
10.1.2	RAID 1 (Mirroring)	49
10.1.3	RAID 2 (bit-level striping with hamming-code)	50
10.1.4	RAID 3	50
10.1.5	RAID 4	50
10.1.6	RAID 5	50
10.1.7	RAID 6	51
10.2	Layering RAID levels	51
10.2.1	RAID 0 + 1	51
10.2.2	RAID 1 + 0	51
10.2.3	RAID 5 + 0	51
10.2.4	RAID 6 + 0	51
10.3	Error Handling	52

1 Introduction

1.1 Operating System Concepts and Structure

The **operating system** offers functionality through **system calls**. When a system call is made, the kernel runs appropriate code in the **priveleged mode**.

A group of system calls implements **services** such as file system serives and process management services.

1.1.1 Process vs. Program

A **process** is a user-level abstraction to execute a program on behalf of a user. Process is a kernel abstraction in which the program is going to run. Each process has its own **address space**. A **program** is a series of binary CPU instructions stored in a file. A program requires resources to be executed (such as CPU or memory). The operating system needs to represent an instance of a program in execution and that is a process.

1.1.2 Address Space

Address space is the range of valid memory address for a given process. Address space can be read and/or written and/or executed. Each process has its own separate address space which consists of...

Text - binary program code
Stack - function calls data
Heap - dunamically allocated data

1.1.3 Files

A **file** is an abstraction of a (possible) real storage device such as a hard disk. You can read/write data from/to a file by providing postion and an amount of data to transfer. In UNIX, everything is considered a file. Files are maintained in **directories** which keep a name for each file they contain. Directories and files form a hierarchy.

- **Block Special Files** can be written and read from block by block. (e.g. disk)
- **Character Special Files** can be written and read from byte by byte. (e.g. serial port)
- **Pipes** are pseudo files allowing for multiple process to communicate over a FIFO channel.

Device drivers are software inside the kernel that knows how to initialize and communicate with specific hardware devices.

1.1.4 System Calls

System calls are the interface the kernel offers to applications to issue service requests. They are highly specific to the operating system and hardware; therefore, you should use system calls included in the C standard library. If you use C system calls then your code will be portable because C will use the correct internal system call depending on the operating system without requiring the programmer to use the operating system specific call.

1.1.5 Kernels

Monolithic kernel is a very large collection of C functions linked together into a single very large executable binary program that executes in kernel mode.

- every function can call every other function
- a single bug will crash everything
- historically how most kernels were implemented

Microkernel splits the kernel into small well-defined modules.

- Only one kernel runs in kernel mode: the microkernel
- the rest of the modules run as relatively powerless ordinary user processes
- a single bug cannot crash the entire system
- there is more overhead communication between the modules

Hybrid kernel is a compromise between a microkernel and a monolithic kernel

1.2 Computer Hardware

Central Processing Unit (CPU)

- fetches instructions from memory and executes them
- each CPU has a specific set of instructions and thus different instruction sets on different CPUs will not be interchangeable
- has different **registers**
 - **General** register for general data
 - **Program** counter to store the memory address of the next CPU instruction to execute
 - **Stack pointer** to store the address of the top of the stack
 - **Program status word**: readable/writable bits that store the state of the CPU and different conditions that resulted from the previous instruction

Memory

In order of decreasing speed and increasing capacity: registers, cache, main memory, magnetic disk, tape

- the cache temporarily holds a piece of data from a slower memory into a faster memory

Hardware multithreading is when you run multiple processes on a single core. Each core has two (or more) sets of registers. Switches between two different sets of registers on each instruction (where needed) to eliminate waiting for data.

Flash memory stores data even when the electricity is removed (SSD).

An **interrupt** is an electric signal on the bus. Interrupt processing involves taking the interrupt, running the kernel interrupt handler, and returning to the user process. The CPU checks for interrupts in-between instructions. When there is an interrupt, the CPU will jump to a unique function in the kernel (the interrupt handler), the kernel handles the interrupt, then the CPU resumes the original process.

1.3 Introduction to Linux

The C standard library contains the implementation of all the functions required by the definition of the C programming language. The kernel provides an API which is also made of C functions, called system calls. Calling a system call cannot be done directly in a user program written in C because it requires use of a special assembly instruction that switches the CPU to kernel mode. Libc includes C functions (wrappers) which have the same names as the system calls which internally do the real system call for you. If such a wrapper is not available for a given system call, libc also provides the *syscall* function that allows you to call a system call using its system call number.

POSIX (Portable Operating System Interface) defines a portable interface including multiple aspects of the OS interface to maintain compatibility between different OSes.

1.3.1 Create your own kernel

1. Configure the kernel by hand (`make menuconfig`)
2. Compile the kernel (`make -h $(nproc) all`)
3. Install the modules (`sudo make INSTALL_MOD_STRIP=1 modules_install`)
4. Install the kernel (`sudo make install`)

Modules are a piece of compiled C code added to the kernel while the kernel is already running, in kernel mode. Modules shorten the development cycle because modules can be loaded and unloaded without having to reboot the system. To load a module that you've made use `sudo insmod ./<module-name>.ko` and to remove the module use `sudo rmmod <module-name>`. Modules have the visibility of all the kernel, so you can use all global variables and exported functions.

1.4 System Calls

System calls are the only way a process can enter the kernel. Used to request OS services and privileged operations such as accessing the hardware, creating other processes, and changing security permissions. System calls switch the execution context of the CPU from user mode to kernel mode.

Syscall trap / Interrupt / Exception

- hardware stacks program counter
- hardware loads new program counter from interrupt vector
- assembly language procedure saves registers
- assembly language procedure sets up new stack
- C interrupt service runs
- scheduler decides which process is to run next
- C procedure returns to the assembly code
- assembly language procedure starts up new current process

1.4.1 To add a syscall in Linux

1. Write your syscall function in either an existing file or a new file. Modify the makefile (`obj-y += my_syscall.o`)
2. Add your syscall in the architecture specific syscall table (`arch/x86/entry/syscall/syscall_64.tbl`)
3. Add your syscall prototype (`include/linux/syscalls.h`)
4. Recompile, reinstall, reboot the kernel
5. Test using the syscall function: invoking the system call from a user-mode program to trigger your code

2 Processes and Threads

2.1 Introduction to Processes

Multiple processes can run the same program. A single CPU can run only one process at a time.

Multiple processes (multiprogramming) increases CPU utilization. This consists of overlapping one process's computation with another's wait and also reduces latency (more interactive). To distinguish multiple processes of the same program, each process has a different **process identifier (PID)**.

A **context switch** pauses execution of one process and continues with the execution of another process. The order in which processes execute is not fixed and may not be reproducible. Thus, programmers cannot make time assumptions and must write their code as if it is the only thing that will be running on the CPU.

Each process has the following key components:

- **stack** stores function arguments and local variables
- **heap** stores dynamically allocated memory
- **executable program** is made up of the initialized data segment and the text segment (CPU instructions from program file)
- **execution context** contains program counter, stack pointer, CPU registers, etc...

Memory layout randomization mixes the location of each part in memory so that the OS is more secure. It will be harder to locate specific processes in memory if it is randomized. There will also be unusable memory in the process that is there to kill the process if anything tries to access the unusable memory.

Stack is built with stack frames and includes local (function) variables). Function parameters are passed on the stack. (LIFO)

Heap is used if it is not known how much data you will need at runtime. Dynamically allocated with random access.

Processes are created during system initialization, from already running processes, due to user interaction, or execution of a scheduled batch job. UNIX process philosophy is to create a new process which is an identical copy of the current process (fork) and then in the new process load the other program.

`fork()` is a syscall used to create a new process

- the process that calls the `fork()` syscall is called the **parent process**
- the newly created process is called the **child process**
- the child is an identical copy of the parent (runs the exact same code)
 - child process returns a zero value
 - parent process returns the PID of the new child
- If C code contains a sequence of n consecutive `fork()` calls, 2^n processes will be created
- most UNIX systems have a limit on the number of processes a given user can have at the same time (`ulimit -a`)

A process **execution state** is what the process is currently doing...

- **Running** - executing instructions on the CPU, it is the process that currently has control of the CPU, cannot have more processes running than number of cores/CPU's
- **Ready** - waiting to be assigned to CPU: ready to execute, but another process is executing on the CPU
- **Waiting/blocked** - waiting for an event: cannot make progress until event happens (most processes are in this state)

Types of **process termination**

- Normal exit (voluntary)
- Error exit (voluntary)
- Fatal error (involuntary) is used when there is a fatal program bug

Orphan processes are children processes whose parents have been terminated. Can be adopted by another process, or can be killed automatically when parent is killed (for security reasons).

Kernel maintains a table called the **process table**. Is stored in the kernel memory and is invisible to processes. Has one entry per process (process's PID is an index into this table). The entry is called the **process control block (PCB)**: contains all information about a process.

When a **context switch** occurs

- hardware saves program counter
- hardware loads new program counter from interrupt vector
- assembly language procedure saves registers

- assembly language procedure sets up new stack
- C interrupt service runs
- scheduler decides which process is to run next
- C procedure returns to the assembly code
- assembly language procedure starts up new current process

2.2 Processes in Linux

Every process has a user and a kernel part in its virtual address space. User space is running the program; kernel space is running kernel code on behalf of the program (kernel space is invisible to the program). There will be both a user-space and kernel-space **execution context** and thus there will be 2 PCs, 2 stack pointer, 2 stacks, etc...

In Linux, a process is a **task** represented by a **task_struct** inside the kernel's memory. This is the **proces control block** for a linux task. A task struct includes:

- scheduling parameters
- memory image
- signals
- CPU registers
- system call state
- file descriptor table
- accounting
- kernel stack

fork() is limited to creating processes only. On Linux, task creation in general can be done by the **clone()** syscall. When **clone()** is called...

- new child starts executing function
- may use a new stack
- has **sharing_flags** to describe the amount of sharing between the caller and the callee
- blurs the distinction between process and thread creation
 - if you share nothing then it's the same as **fork()**
 - if you share everything then it's the same as thread creation

Internally, both **fork()** and **clone()** call the kernel function **kernel_clone()**.

2.3 Threads

A **thread** cannot exist without a process. A process is a "container" for threads. A process may contain one (the default) or more threads. Thread tells you what part of the program is currently being executed. **Thread** is an independent sequential execution stream within a process; programs use one or more threads per process. Items that are **shared per process** include the following:

- address space
- global variables
- open files
- child processes
- pending alarms
- signals and signal handlers
- accounting information

Items that are **shared per thread** (private to each thread) include the following:

- program counter
- registers
- stack
- state

Some of the advantages of threads:

Performance - shared address space and therefore no communication overhead between threads. Thread creation can be 10-100 times faster than process creation.

Efficiency - allows a program to overlap I/O and communication. Allows one process to use multiple cores.

Amdahl's Law identifies performance gains from adding additional cores to an application that has both serial and parallel components. **Serial** components can only be executed by a single thread. **Parallel** components can be executed by multiple threads.

$$speedup \leq \frac{1}{S + \frac{(1-S)}{N}}$$

where S is the serial portion and N is the number of processing cores. Note that the serial portion of an application has disproportionate effect on performance gained by adding additional cores / CPUs.

All threads in a single process share the same address space. Global variables are shared between threads. Multiple threads may access the same global variable concurrently and programmers are responsible to coordinate access (synchronization algorithms). Each thread gets its own local variables (in stack). Each process starts with a single "main" thread; the first thread can create new threads. Any thread can end at any point; if the main thread returns from the `main()` function, the `exit()` system call is called and the whole process terminates.

2.3.1 PCB

When making threads, the kernel will break the PCB into two pieces.

- Information about program execution is stored in **Thread Control Block (TCB)**. This includes the program counter, CPU registers, scheduling information, and pending I/O information.
- Other information is stored in the **Process Control Block (PCB)**. This includes memory management information and accounting information.

2.3.2 Threads Implemented in User Space

Pros:

- **Portable:** can be used on any OS
- a user-level threads library can be implemented on an OS that does not support threads
- thread switching is at least an order of magnitude faster than trapping to the kernel
- thread scheduling is very fast: no process context switching, no kernel trap, no flushing of memory cache
- each process can have its own thread scheduling algorithm

Cons:

- **Doesn't support real parallelism:** kernel will schedule process on single core regardless of number of threads and therefore cannot take advantage of multicore/multiprocessor systems
- if one thread makes a blocking system call then the kernel suspends the whole process
- threads need to voluntarily give up the CPU to each other for multiprogramming so user code is more complex

2.3.3 Threads Implemented in Kernel Space

Pros:

- **supports real parallelism** (multicore/multiprocessor systems)
- no run-time thread system needed in each process
- no thread table in each process
- block system calls are not a problem (the kernel can schedule another thread)

Cons:

- **slower than user-space threads:** if thread operations are common, much more kernel overhead will be incurred because of the many system calls required
- if you fork a multithreaded process, do you copy all threads or just the one that called `fork()`?
- if signals are sent to processes, should the kernel assign it to a specific thread to handle

In practice, user threads are mapped one-to-one with kernel threads: **hybrid implementation**.

2.3.4 Context Switching with Threads

When context switching with threads...

- thread is now the unit of a context switch
- context switch causes CPU state to be copied to/from the TCB
- when context switching two threads of the same process, no need to change address space
- when context switching two threads in different processes, must change address space

2.3.5 Thread Local Storage

Thread local storage (TLS) works like a global variable but each thread has its own copy of it. Use this instead of global variables when making single-threaded code multi-threaded. Visible everywhere in the code, but each thread has its own copy.

2.4 Threads in Linux

Linux, implements threads in the kernel, not in user space. Linux uniformly handles processes and threads: processes and thread tables are unified into a single data structure. Processes and threads are considered **tasks** and are thus represented by a **task_struct**: there is one **task_struct** per thread. Linux assigns PID and TID for each thread; each thread has a unique TID. To check the TID use the `gettid()` syscall.

3 Scheduling

3.1 Introduction to Scheduling

Multiple processes and threads are ready to run; the kernel's scheduler decides which runs next. Scheduling decisions may take place when a process/thread is created, terminates, blocks on an event. Scheduling decisions may also take place when an interrupt occurs; this includes clock interrupts (*running to ready*) or I/O interrupts (*blocked to ready*).

3.1.1 Clock Interrupts

Clock interrupts are a way for a system to keep track of time. Interrupts occur at periodic intervals called a clock tick. Clock interrupts are implemented using a hardware clock interrupt and have high priority in the system.

3.1.2 Non-preemptive Scheduling

In **Non-preemptive Scheduling**, processes execute until completion or until they make a **voluntary process switch** or **process switch on blocking calls**. The scheduler gets involved only at exit or on request (i.e. for every clock interrupt, running processes keep going).

3.1.3 Preemptive Scheduling

In **preemptive scheduling** while a process executes, its execution may be paused and another process resumes its execution. This can take the form of an **involuntary process switch**. This scheduling is used in all OSes today.

3.1.4 Goals of All Scheduling Algorithms

- **fairness** - give each process a fair share of the CPU
- **policy enforcement** - seeing that stated policy is carried out
- **balance** - keeping all parts of the system busy

3.2 Scheduling in Batch Systems

3.2.1 Goals for Scheduling in Batch Systems

- **throughput** - maximize jobs per hour
- **turnaround time** - maximize time between submission and termination
- **CPU utilization** - keep the CPU busy all the time

3.2.2 First-come First-served (FCFS)

Processes are assigned to the CPU in the order they request it (or they arrive). The non-preemptive version will let each job keep running until completion. **Convoy effect** is when a long process delays short processes. **Turnaround time** is the time taken by a job to complete after submission (including the wait time).

$$\text{turnaround time} = \text{time}_{\text{end}} - \text{time}_{\text{submitted}}$$

3.2.3 Shortest Job First (SJF)

Associate the length of its CPU time with each process. Use the CPU time length to schedule the process with the shortest CPU time first. In the **non-preemptive version**, once the CPU is given to the process, it cannot be taken away and then after completion of a process, will consider the next shortest job and begin executing that one until completion. In the **preemptive version**, if a new process arrives with less CPU time than the remaining time of the current executing process: then preempt.

The **preemptive version** is called **shortest remaining time next (SRTN)** and is optimal in terms of average turnaround time (i.e. always gives minimum average turnaround time). This method also prevents the convoy effect. However, there are more context switches and a process may be starved of the processor if short processes keep arriving.

3.3 Scheduling in Interactive Systems

3.3.1 Goals for Scheduling in Interactive Systems

- **response time** - respond to requests quickly
- **proportionality** - match duration expectation of users

3.3.2 Round-Robin (RR)

Each process is allowed to run for a specified time interval called the **time quantum**; in practice, the time quantum is likely not fixed. After this time (the time quantum) has elapsed, the process is preempted and added to the end of

the ready queue and the next process is scheduled. If the process terminates or blocks before the time quantum is entirely used up, then the process loses the rest of its time quantum and the next process is scheduled with a new time quantum.

Advantages of RR

- solution to fairness and starvation
- fair allocation of CPU across jobs
- low average waiting time when job lengths vary
- good for responsiveness if small number of jobs

Disadvantages of RR

- context-switching time may add up for long jobs

Context switching may impact the choice of the time quantum...

- when there are many processes, a long quantum causes a poor response time
- a short time quantum makes things more responsive but gives a higher context switch overhead
- time quantum might be variable depending on CPU load

Cache state must be shared between all jobs which may slow down execution in RR scheduling. There is no cache sharing in FCFS.

3.3.3 Priority (PRIO) Scheduling

PRIO has four priority classes and always executes the highest-priority runnable jobs to completion. Each queue is processed in a round-robin fashion with a time quantum. Note that lower numbers represent higher priority. There are, however, some **problems** with PRIO scheduling...

- **starvation** when lower priority jobs don't get run because higher priority tasks always running
- **priority inversion** happens when a low priority task has resource needed by high priority task, which then must wait

Priorities can be **assigned statically** based on process type, user, how much the user paid (if a cloud server); or **dynamically** based on how much a process runs vs doing I/O. $\text{Priority} = \frac{1}{f}$ where f is the size of the quantum last used. Thus, the longer a process ran, the lower its priority and the process that runs the shortest gets the highest priority to run next.

3.3.4 Multiple Queues (MQ) Scheduling

Same as priority scheduling, but each queue has a different time quantum (shortest quantum for high-priority and longer for low-priority). Processes start at the highest priority. When a process **exceeds** its time quantum, it's moved to the next lowest priority queue. When a process **becomes interactive** it is moved to the higher priority. If the user discovers how to make their tasks more interactive, they can get all of their processes to be very high priority, which is a problem.

3.3.5 Other Schedulers for Interactive Systems

- **Shortest processes next:** SJF but educated guess for how much CPU a process will need next
- **Guaranteed scheduling**
- **Lottery scheduling** (probabilistic scheduling)
- **Fair-share scheduling:** round robin scheduling between users.

3.4 Scheduling in Real-time Systems

3.4.1 Goals for Scheduling in Real-Time Systems

- **meeting deadlines** - avoid losing data
- **predictability** - avoid quality degradation in multimedia systems

Time plays an essential role in real-time system scheduling. One or more physical devices external to the computer generate events and the computer must react appropriately to them within a fixed amount of time. If the computer reacts too late, it is **as bad as not reacting**.

Real-time scheduling can be split into two categories. **Hard real-time** is when there are absolute deadlines that must be met. **Soft real-time** is when missing an occasional deadline is undesirable but tolerable.

Let's assume that process behavior is predictable and known in advance (the software should be written in a way to make this true). Task timing is known and represented through the following values.

Release time (R_i) is the earliest time when a task can start execution

Execution time (C_i) is the expected execution time for a task

Deadline (D_i) is the time by which the process must be completed.

We can also assume that the **periodicity** is known. It can be one of the following options:

- periodic process
- sporadic process (aperiodic, hard deadlines)
- aperiodic process (aperiodic, soft deadline)

3.4.2 Number of Schedulable Processes

$$\sum_{i=1}^m \frac{C_i}{P_i} \leq 1$$

where there are m periodic events. Event i occurs with period P_i , requires C_i time on the CPU. The percentage of CPU usage for event i is $\frac{C_i}{P_i}$.

3.4.3 Rate Monotonic Scheduling

Preemptive algorithm, where the shorter the period, the higher the priority.

3.4.4 Earliest Deadline First Scheduling

Schedules processes according to the shortest remaining time until the next deadline. The shorter the remaining time, the higher the priority.

4 Concurrency

4.1 Inter-process Communication

Processes need to share information or coordinate. Threads are nice for this but do not provide isolation from each other. Processes provide isolation but communication requires system calls. Some **problems with inter-process communication** include...

- how can one process pass information to another in such a way that processes do not get in each other's way?
- how to maintain proper sequencing when dependencies are present between processes?

The same problems apply to inter-thread communication because by design, threads communicate via shared memory. For one process to **pass information** to another, it can...

- pass messages through the kernel
- share memory (which can be setup with system calls)
- share a file
- use asynchronous signals or alerts
- use named pipes (which is just a message queue inside the kernel essentially)

4.2 Race Conditions

A **race condition** occurs when two or more processes/threads are reading or writing shared data and the final result depends on which runs precisely when. The **critical section** is the part of the program where the shared data is accessed. Uncoordinated read/write of the data in critical sections may lead to race conditions. Must find some way to prohibit more than one process to execute in its critical section at the same time. A solution is to have **mutual exclusive** access to critical sections: only one processor or thread in a critical section at any time.

4.2.1 Requirements to Avoid Race Conditions

- No two processes may be simultaneously inside their critical section: **mutual exclusion**
- No processes running outside its critical section may block other processes from entering a critical section: **progress**
- No process should have to wait forever to enter its critical section (be starved): **bounded waiting**

4.2.2 Disabling Interrupts to Avoid Race Conditions

Right before entering the critical section, the process disables all hardware interrupts and when leaving the critical section, the process reenables all interrupts. When interrupts are disabled, the CPU cannot be switched to another process, so the current process can keep the CPU just for itself while it is in the critical section (unless the process makes a system call). The **advantage** of disabling interrupts is that it is very easy to implement in software. The **disadvantages** of disabling interrupts is that it is unwise to give user processes the ability to turn off interrupts and this is not suitable for multicore systems. In practice, disabling interrupts is only used on single core machines for mutual exclusion in kernel space.

4.2.3 Using Lock Variables to Avoid Race Conditions

Use a single shared variable (**lock**) between processes. A value of 0 means that no process is in its critical section. A value of 1 means that some process is in its critical section. This solution creates a critical section with the lock if the lock is read and written to in two separate lines of code (i.e. if the process context switches after reading the lock but before writing/updating the lock, then both processes may enter critical section simultaneously). Continuously testing a variable until some value appears (**busy waiting**) wastes a lot of CPU time. Note that a lock variable that uses busy waiting is called a **spinlock**.

4.2.4 Strict Alternation to Avoid Race Conditions

There is a single turn variable shared between processes. A value of 0 means that it is the turn of process 0. A value of 1 means it is the turn of process 1. the turn variable is initially set with 0 or 1 based on how should start first and the processes strictly alternate. This can easily be generalized to n different processes. An **advantage** of this solution is that there are no race conditions. A **disadvantage** of this solution is that it does not achieve progress because a process can be blocked by a process not in its critical section.

4.2.5 Peterson's Solution

Use two shared variables (`turn`, `flags[]`) between processes. `int turn` imposes an access order: which process can enter its critical section next. `bool flags[]` (array) specifies for each process whether it is currently interested in entering its critical section. This array is initialized to be false. Some **advantages** of this solution are that there are no race conditions and it satisfies all three conditions for critical regions. **Disadvantages** are that busy waiting uses a lot of CPU and generalization to more than two processes is complex.

4.2.6 Test-and-Set Lock (TSL) Instruction

Reads and modifies the content of a memory word atomically. Takes arguments `TSL register`, `memory_address`. Returns in a register the current value of the memory word at `memory_address` and then sets the value of the memory word at `memory_address` to true (usually abstracted as a non-zero value). Operations cannot be interrupted during execution because this locks the memory bus. Essentially, there can only be one read and one write in one atomic step.

4.2.7 TSL Solution

There is a single lock variable in memory at `LOCK_ADDR`, shared between processes. Uses TSL instead of while loops and assignment. **Advantages** of this solution are that there are no race conditions, mutual exclusion is achieved, and progress can be made. **Disadvantages** of this solution are that you must have a TSL instruction implemented in the CPU, there is no bounded waiting, and it has busy-waiting.

4.2.8 Producer-Consumer Problem with `sleep()` and `wakeup()`

Instead of busy-waiting, let the process sleep (requires kernel syscalls). The `sleep()` syscall causes the caller to give up the CPU for some duration of time or until some other process wakes it up. The `wakeup()` syscall causes the caller to wake up some sleeping process. Note that this solution doesn't solve the priority inversion problem because process B might still end up having to wait a long time for process A to exit the critical section, even if process B is now sleeping.

```

1  Producer
2  while(1){
3      produce an item A;
4      if (count == N) //full buffer
5          sleep();
6      insert item;
7      count++;
8      if (count == 1) //was buffer empty
9          wakeup(consumer);
10 }
11
12 Consumer
13 while(1){
14     if (count == 0) //empty buffer
15         sleep();
16     remove item;
17     count--;
18     if (count == N-1) //was buffer full?
19         wakeup(producer);
20     consume an item;
21 }

```

In this solution both the producer and consumer may sleep forever. An **advantage** of this solution is that there is no busy-waiting. **Disadvantages** of this solution are that there is a race condition. A wakeup sent to a process that is not (yet) sleeping is lost. Decision to go to sleep and calling sleep() must be paired. So that the test of count must be done at the same time as the sleep/wakeup in some way. To achieve this we use **semaphores**...

4.2.9 Semaphores

The purpose of **semaphores** is to solve the lost wakeup problem (explained just above) and avoid race conditions. Semaphores will try to store the number of wakeups. To make a semaphore, define a count variable (the semaphore): **Down(<semaphore>)** is equivalent to consume (decrease) or sleep(), **Up(<semaphore>)** is equivalent to produce (increase) and (potentially) wakeup. **down()** and **up()** are **atomic** (i.e. once the operation has started, no other process can access the semaphore until completed or blocked).

The value of the semaphore... if **positive**, represents the number of resources available and that there are no pending wakeups. If **zero** represents that there are no resources available and no pending wakeups. If **negative** represents that there are no resources available and signifies the number of pending wakeups (number of processes waiting).

4.2.10 Types of Semaphores

Counting semaphore is a semaphore based on the number of controlled resources. The value is an integer.

Binary Semaphore is a semaphore with a value of 0 or 1 which implements mutual exclusion; this is called a **mutex**. The value is either a 0 or 1. A value of 0 represents **locked** and that the critical region is not available. A value of 1 represents **unlocked** and that the critical region is available.

4.2.11 Producer-Consumer Problem using Semaphores

```
1  Shared variables
2  const int N=100;
3  semaphore empty=N, full=0;
4
5  Producer
6  while(1) {
7      produce an item A;
8      down(empty);
9      insert item;
10     up(full);
11 }
12
13 Consumer
14 while(1) {
15     down(full);
16     remove item;
17     up(empty);
18     consume an item;
19 }
```

The sum of the semaphores must be equal to N . There could be a race condition with the insert and remove actions: if the buffer is not empty or full, then they may both try to insert/remove at the same time.

4.2.12 Producer-Consumer Problem using Semaphores and Mutexes

```
1  Shared variables
2  const int N=100;
3  semaphore empty=N, full=0;
4  mutex mux=1;
5
6  Producer
7  while(1) {
8      produce an item A;
9      down(empty);
10     lock(mux); // down(mux)
11     insert item;
12     unlock(mux); // up(mux)
13     up(full);
14 }
15
16 Consumer
17 while(1) {
18     down(full);
19     lock(mux); // down(mux)
```

```

20     remove item;
21     unlock(mux); // up(mux)
22     up(empty);
23     consume an item;
24 }

```

The order of the `up()` and `down()` calls is important and can't be changed.

```

Producer = down(empty), down(mux), up(mux), up(full)
Consumer = down(full), down(mux), up(mux), up(empty)

```

4.2.13 Monitors

Monitors are a higher-level primitive. Monitors are a programming language construct: a package or module or class. Some **rules for monitors** include:

- only one process/thread can be active in a monitor at any time
- processes cannot access internal data of monitor (fields are private)
- put all critical sections inside monitor methods

A **mutex** is used to achieve mutual exclusion: locked/unlocked automatically at function/method invocation and return. They also have a **condition variable** which is used for processes to wait/block when attempting to enter a monitor which is already currently in use and to wake up the process.

`wait(condition)`

- adds process to wait queue for condition variable
- causes the calling process to block
- releases the mutex to allow other processes to enter monitor
- re-acquire mutex once blocked process is woken up

`signal(control)`

- wakes up one blocked process on the condition, if any
- exits the monitor, leaving critical section

4.3 Readers-Writers Problem

Problem is that multiple processes / threads need to access a shared resource concurrently (but some only need to read it). Mutual exclusion allows one process to access the resource at a time, but can we allow more than one reader at a time? We want multiple readers to access the resource concurrently, but writers must access it with mutual exclusion.


```

1  typedef int semaphore;
2  semaphore mutex = 1;
3  semaphore db = 1;
4  int rc = 0;
5
6  void reader(void) {
7      while(TRUE) {
8          down(&mutex);
9          rc = rc + 1;
10         if (rc == 1) down(&db)
11         up(&mutex)
12         read\_data\_base();
13         down(&mutex);
14         rc = rc - 1;
15         if (rc == 0) up(&db);
16         up(&mutex);
17         use\_data\_read();
18     }
19 }
20
21 void writer(void) {
22     while(TRUE) {
23         think\_up\_data();
24         down(&db);
25         write\_date\_base();
26         up(&db);
27     }
28 }

```

5 Memory Management

Physical memory (RAM) is a hardware array consisting of words. **words** are a fixed-size unit of data (64 bits). When there is **one program** with **no memory abstraction**, the program lives in and directly works with physical memory; the program can even share physical memory with the OS. When there are **multiple programs** with **no memory abstraction** every program operates in physical memory, there is no isolation and thus programs can access each other's memory and OS memory.

Problems with having multiple programs and no memory abstraction include:

- No protection
- expensive relocation process
- security
- **stability**: user program could write over something in OS

Additionally, note that programs must be **relocatable**: written to be placed and run at any physical memory address. The **loader** is the part of the OS that loads the program from disk to RAM and is responsible for deciding where to place programs based on the available physical memory.

5.1 Memory Management

The **virtual address space** abstracts physical memory space so that each process has its own memory address space.

5.1.1 Base and Limit Registers

One such tool for implementing virtual memory is the use of **base and limit registers**. These are hardware support to ease relocation (store the base register and limit register). Thus, the program does not need to be relocatable, so there is a faster load time. Each process also gets its own private address space. In this implementation, the memory management unit will convert virtual address in program (offset) to the physical memory address, then retrieve the data from the memory and return it to the user program.

5.1.2 Swapping

1. Save address space of idle processes to disk and reclaim memory
2. When a swapped process needs to run, bring its address space from disk back to memory

3. Go back to step 1

Know that idle processes are mostly stored on the disk, so they do not take up any memory when they are not running. Swapped processes are kicked out of memory and brought back into memory in full when it is needed again.

5.1.3 Memory Fragmentation

Memory fragmentation happens during process creation, exit, and swapping when memory holes are created. This brings about the problem that new processes may not fit in any available memory hole, even if all the holes together are big enough for the new process. A solution to this is **compaction**. In **compaction** all the processes are moved tightly together, which results in one single hole which is big enough for the new process. Compaction has to freeze all the programs to perform the operation, so it is not good to use on interactive systems where the freeze will be noticeable.

5.1.4 Determine which Memory is Allocated

With all implementations for memory management, recall that we need to allow extra room for process growth. To figure out which part of memory is allocated, divide the memory into **blocks**.

$$1 \text{ block} = 4\text{KB}$$

5.1.4.1 Bitmaps

Allocation units (blocks) are as small as a few words and as large as many KB. Each allocation unit has a bit in the bitmap which is 1 if that block is occupied and 0 if that block is free. An issue with this implementation is the long search time.

5.1.4.2 Linked Lists

Each node in the list contains a process or a hole. Each node has a start address and length. Since the list is sorted by increasing addresses, it's easy to merge adjacent nodes when a process exits.

5.1.5 Picking Location of New Process's Memory

When picking where to allocate the memory of a new process, there are many options. The following are some of the common options:

- **First fit**: take the first fitting hole. This is the simplest option. Tends to allocate at the beginning of memory
- **Next fit**: take the next fitting hold. Faster than first fit in practice (allocates from all over the memory) but is not as memory efficient.

- **Best fit:** take the best fitting hole. This is slower as it must search the whole list. Prone to creating lots of small holes.
- **Worst fit:** take the worst fitting hole. Poor performance in practice
- **Quick fit:** keep hole lists for holes of different sizes. Poor coalescing performance.

5.1.6 Overlays

Overlays are a strategy where you only load part of the program. The programmer breaks the address space into pieces so that each piece can fit into physical memory. The pieces are called **overlays** and are loaded and overwritten by the program when necessary. This requires an overlay manager which likely takes the form of some software library. In this library, programmers call functions of the overlay manager to load an overlay segment into memory when it is not in RAM. This overwrites a segment previously in RAM. The **issues** of this implementation are that it is implemented in the program itself and therefore there is a high complexity for the developer. However, this is still used in some embedded systems where the hardware does not support paging.

6 Virtual Memory

Virtual memory is a powerful concept for the address-space abstraction that decouples process address space from physical memory. It gives the illusion of large and private address spaces. It also allows for programs to execute even if only partially loaded in memory. It enables different memory protection policies as well. In this, each process has its own virtual address space from 0x00 to some maximum. The program and CPU operate on virtual addresses which are translated to the physical addresses: a process which is transparent to the program.

The virtual address space is broken into **chunks**. Each chunk is a contiguous range of addresses mapped onto a contiguous range of physical memory. The process always sees the entire virtual address space. Not all chunks need be in physical memory to run the program. Chunk switching is hidden from the process. When a program references a chunk loaded in physical memory, hardware performs necessary virtual-to-physical address translation on the fly. When a program references a chunk not in physical memory, the kernel gets it from the disk and re-executes the memory instruction a second time. A chunk is either a **segment** or a **page**.

6.1 Paging

A **virtual address space** consists of fixed-size chunks called **pages** which each correspond to chunks in the **physical memory** called **frames**. Pages and frames are *always* the same size (which is usually 4KB). The hardware **memory management unit (MMU)** translates pages to frames; this translation is invisible to the process and the programmer does not have to worry about it.

The **application** sees a single, private, contiguous address space while in reality application code and data are scattered across physical memory. Note that not all pages need to be loaded into physical memory.

6.1.1 Translating Addresses when Paging

To translate the address when paging, the virtual address is broken into **page number** and **page offset**. **Page number** is used to get the **frame number**. This uses a mapping function and is implemented with a page table. The page number is used as the index into the page table. Each page table entry stores a frame number plus some additional information. The number of entries in the page table is equal to the size of the virtual address space divided by the page size.

6.1.2 Page Table Entry

The *exact* content of a **page table entry** is highly OS and MMU dependent; however, all must contain at least the frame number. The entry will also have a **modified/dirty** bit (referenced bit) which is set by the MMU when the page is used to allow the kernel to keep track of page usage. The entry also has **protection bits** (i.e. what kind of access is permitted: read, write, execute). The entry also has a **present/absent bit**: 1 if the entry is valid and page can be used, 0 if the page is not currently in memory. Accessing a page with absent/invalid entry causes a **page fault**.

6.1.3 Page Faults

If a process makes a reference to a page which is not currently in RAM, the reference to that page will result in the MMU detecting that the page is marked as invalid in the process' page table, and the MMU will trap to the kernel: **page fault**.

1. Kernel looks at another table in the process' PCB to decide
 - If invalid reference, then abort
 - If just not in memory, but exists somewhere else
2. Kernel finds a free frame in RAM
3. Page-in the needed page into the free frame in RAM using a disk operation

4. Kernel changes the page table to indicate that the page is now in memory: set the valid bit to true
5. Kernel restarts the process at the same CPU instruction that caused the page fault
6. The second time the instruction will succeed without a page fault.

During a process context switch, the kernel tells the MMU where in RAM to find the page table for the process executing next on the CPU. Everytime the process reads/writes a virtual address...

- MMU uses the page number bits of the virtual address to lookup in the page table the corresponding frame number
- If a page table entry is found, the MMU uses that frame number and the offset bits to generate the physical memory request
- If a page table entry is not found then a page fault is generated.
 - Either the kernel kills the process if it was trying to access nonexistent memory
 - Or the kernel loads the requested page from disk and updates the page table before resuming the process's execution

6.1.4 Addresses

In an n -bit virtual address space, a virtual address consists of n bits and there are 2^n virtual addresses.

In an n -bit virtual address space with k bits for the page offset. Each page consists of 2^k virtual addresses. The virtual page number takes $(n - k)$ bits. There are $2^{(n-k)}$ virtual pages and $2^{(n-k)}$ entries in the table.

In an m -bit physical address space, a physical address consists of m bits. There are 2^m physical addresses.

In an m -bit physical address space with k bits for the frame offset. Each page consists of 2^k physical addresses. The frame number takes $(m - k)$ bits. There are at most $2^{(m-k)}$ physical frames.

6.1.5 Translation Lookaside Buffer

The use of a **translation lookaside buffer** helps reduce accesses to the page table by caching recently used page table entries. TLB is a hardware cache for page table entries inside the MMU, used to speed up address translation. This works because most programs often make a large number of references to

a small number of pages: locality of memory references.

The **translation lookaside buffer** translates virtual addresses into physical addresses without going through the page table. It is usually implemented as an **associative memory** that caches most recent virtual to physical address translations. All TLB entries are searched in parallel.

- If the page number is found in the TLB: **hit**. The matching TLB entry contains the right frame number and thus the virtual address can easily be translated to the physical memory.
- If the page number is not found in the TLB: **miss**. The MMU accesses the page table in the main memory to get the right frame number. MMU translates the virtual address into a physical address. The MMU updates the TLB so accessing the same page later will then result in a hit.

6.1.5.1 Metrics

Hit ratio (h) is the percentage of time that a page number is found in the TLB.

Effective access time (EAT) = $a + (2 - h)m$. Where a is TLB lookup time, m is memory access time. a is negligible; if h is close to 1 (100%) then EAT is close to m .

6.1.6 Multilevel Page Tables

There is a high memory overhead with a **flat page table**. If the virtual address space is large then the page table will be large as well. Even if the program uses a small fraction of its virtual address space, the entire page table is still needed.

In **multilevel page tables**, the virtual address bits are split into three parts: high page bits, low page bits, and offset bits. The basic idea is to *page* the page table. This allows for only portions of a page table to be kept in memory at one time (the rest can be on the disk); portions of a page table can simply not be present at all. When there is a TLB miss, a process might now require three memory accesses:

- primary page table needed by the MMU
- Secondary page table needed by the MMU
- Data requested by the process

In practice, the TLB hit ratio (h) is very close to 100% so requiring 3 memory accesses only happens very rarely.

6.1.6.1 Example

Assume the 32-bit address is allocated as following: 10 bits to the primary page, 10 bits to the secondary page, 12bits to the page offset.

1st-level page size = 1024×4 bytes = 4KB

2nd-level page size = $2^{10} = 1024 \times 4$ bytes = 4KB. The secondary page gets 10 bits, so 1024 unique addresses and each one is 4 bytes.

Total size of the page table = $4KB \times 1024 + 4KB$. All 1024 entries in the primary table point to one secondary page.

Thus, 8KB of memory is needed for one virtual address translation.

6.1.7 Inverted Page Tables

Another approach, instead of paging, is using **inverted paged tables**. To do this, you must have a frame table where the frame number is the index and the page number is the context. The table size is proportional to RAM size, which is much smaller than the size of the virtual address space of a process. Only a single frame is required for all processes. The problem with this implementation is that a TLB miss requires searching the frame table entry by entry, so is very slow on a miss. This takes **less space** than paging but **more time**.

To implement an inverted page table using hashed maps, hash the page number to a bucket. The table size only depends on the hash function used. Hashing is slow but only needs to be done on a TLB miss. Multiple page numbers might hash to the same bucket, so then need to search a (short) list to find the right frame number.

7 Page Replacement Algorithm

If there is no free frame in physical memory, then must identify some *victim* pages in memory and move them to disk to free the corresponding frames.

Question: Which page to evict when physical memory is full and a free frame is required?

Goal: Achieve the lowest number of (future) page faults

Input of Algorithm: A particular string of memory references

Output of Algorithm: The number of page faults on that string

7.1 Optimal Algorithm

The **optimal algorithm** will replace the page in memory that will be used the furthest in the future. This is the page that will create a page fault the furthest in the future, thereby minimizing the number of page faults. This is not achievable in real systems because this only works if we know the whole sequence of page references in advance.

In the following algorithms, we will use the following notation. In **page table entries**, **R** is set whenever the page is referenced, **M** is set when the page is written to.

7.2 Not Recently Used Algorithm

The general idea of this algorithm is to use virtual memory hardware tracking bits to determine what page was not recently accessed.

1. When a process starts, the R and M bits for all its pages are set to 0 by the kernel
2. Periodically (on the clock interrupt), the R bit is cleared (the M bit remains the same).
 - To distinguish pages that have not been referenced recently from those that have been
3. When a page fault occurs, the kernel inspects all the pages and divides them into four categories
 - **Class 0:** not referenced but not modified
 - **Class 1:** not referenced but modified
 - **Class 2:** referenced but not modified
 - **Class 3:** referenced and modified

4. The algorithm removes a page at **random** from the lowest-numbered non-empty class

Disadvantages of NRU algorithm:

- only looks at referenced/modified bits changed since the last clock tick
- pages split among only four categories
- performance is, at best, adequate

Note that the page removed is from the same process that triggered the page fault in the first place.

7.3 First in First Out (FIFO) Algorithm

In FIFO, the kernel maintains a linked list of all pages in the order that they came into physical memory. The page at the beginning of the list is removed from the list and from physical memory when a free frame is needed; this is also the page that has been in physical memory for the longest time, so may very well be a frequently used page.

An **advantage** of FIFO algorithm is that it is easy to implement.

Disadvantages of FIFO algorithm:

- The page which is in memory the longest is replaced. This does not consider page usage and thus may throw out pages that are in use right now.
- Suffers from "**Belady's Anomaly**": the page fault rate might increase when there are more frames

The linked list can contain pages from many different processes, so the page removed might not belong to the process that triggered the page fault.

7.4 Second Chance Algorithm

Second chance is a FIFO variant that adds the concept of page usage (check reference bit). This examines pages in FIFO order starting from the beginning of the linked list but considers the reference bit, R of the page at the front of the list:

- If $R=0$, the oldest page has not been used since it was added at the end of the list, so remove page (go to third bullet)
- If $R=1$, set $R=0$ and place the oldest page back at the end of FIFO list (second chance) because it was used recently (go to first bullet)
- Add new page at the end of FIFO

If no victim page is found on the first pass, the algorithm automatically reverts to pure FIFO on the second pass. The problem with the second chance algorithm is that moving pages around the list is not efficient. The **clock algorithm** is a more performant implementation that uses a circular list and simply moves a pointer around the list.

7.5 Least Recently Used (LRU) Algorithm

LRU is a good approximation of the optimal algorithm that is based on the observations that pages that have been heavily used in the last few instructions will probably be heavily used again soon. Additionally, we can infer that pages that have not been used for a long time will probably remain unused for a long time. When a page fault occurs, remove the page that has been unused for the longest time in the past.

7.5.1 Priority List Implementation of LRU

Keep a doubly linked list of page numbers with the most recently used at the front and the least recently at the rear. A page that gets referenced is moved to the front, and on a page fault, you always remove from memory the page at the rear of the linked list. In this implementation, it is easy to find the least recently used page, but a list update is required on every memory reference; this is time consuming and not practical.

7.5.2 Counter Implementation of LRU

Equip the hardware with an N-bit counter and on every instruction that references memory, increment the counter. Every page table entry must have a field to save the counter. When a page is referenced, copy the counter to that field. The more recently that the page is used, the larger the counter value is in its table entry. For replacement during a page fault: evict the page with the lowest counter value (i.e. the oldest page). This operation requires searching the page table for the entry with the lowest counter, so it is linear time in the size of the page table. This also means that the evicted page always belongs to the process that triggered the page fault.

7.6 Not Frequently Used (NFU) Algorithm

1. At each clock interrupt, for each page, the R bit is added to a software counter for each page which keeps track (roughly) of how often each page has been referenced in the past
2. The page with the lowest counter is evicted during page fault

This algorithm never forgets: looks at overall usage of each page rather than recently used usage.

7.7 Aging Algorithm

This algorithm is a variant of NFU that forgets the past after a while. Uses an N-bit counter per page and periodically...

- Shift counter to the right to forget the oldest bit: reduces counter values over time by dividing it by two
- Add R as the new leftmost bit to represent recent usage. This recent R bit is added as the most significant bit and, therefore, more weight is given to recent references

When a page fault occurs, we replace the page with the lowest counter. Because this algorithm uses 8 bit integers, it only remembers what has happened within the last 8 clock ticks (anything older than that would be shifted to zero). If two pages have a value of 0 then we know they have not been used for n clock ticks, but we do not know which of the two was used more recently.

7.8 When to Move Pages into Memory

There are many options of how to decide when to move pages into memory. Here are some of the options.

Demand paging is when pages are loaded on demand, not in advance

Prepaging loads a group of pages at once in the hope of minimizing the number of page faults later. The group of pages to get is a guess but could include all of a given program's pages, the first/last N pages: the **working set pages**

7.8.1 Working Set Pages

The **working set** is the set of pages that a process is currently using. This set may change over time in the regard that there will be different pages and that there will be a different total size due to the pages themselves changing. The working set of pages are those used by the k most recent memory references. $w(k, t)$ is the size of the working set at time t .

If the entire working set is in physical memory, there are no page faults. If there is insufficient physical memory for the working set, **thrashing** where a process busies itself copying pages in and out of memory. The goal is to keep the working set in physical memory to minimize the number of page faults.

7.9 Working Set Based Algorithm

This page replacement algorithm has a threshold T (working set window size), and for every page table entry keeps track of the time of last use, reference bit

R and modified bit M. During a page fault, the algorithm will scan the page table, and for each entry:

- If **R=1**, set time of last use to current time, set R to 0. Page has been referenced since last clock tick, so in working set, so remember when it was last used.
- If **R=0**, the page is a candidate for removal. First calculate the age: age = (current time – time of last use). If the age is greater than threshold T, then the page is replaced. If the age is less than threshold T, the page is still in the working set, but may be removed if it is the oldest page in the working set.
- If not page has R=0, choose the oldest and when they are all the same age, pick randomly

This algorithm works fine in terms of page faults, but requires scanning the whole page table which is expensive.

7.10 Working Set Clock (WSClock) Algorithm

This algorithm uses a circular list of **frames** (including R and M bits and time of last use). When there is a page fault:

- If **R=1**, set R=0 and advance the clock hand around the list. This page was used recently and is part of the working set.
- If **R=0** and **age is greater than threshold T**:
 - Page is clean (M=0), replace this page
 - Page is dirty (M=1), schedule write but advance hand to check other pages
- At the end of the first pass, if no page has been replaced...
 - If the write has been scheduled, keep moving the hand until the write is done and the page is clean. Evict the first clean page.
 - If no write scheduled, claim any clean page even though it is in the working set.

This algorithm is simple and has good performance so is used a lot in practice.

7.11 Frame Allocation to Processes

Each process requires a minimum number of frames. The following are some of the options for choosing how many frames to allocate.

7.11.1 Fixed Frame Allocation

In this implementation there could be **equal allocation**; so, if-for example- there are 100 frames and 5 processes, then you give each process 20 frames. There could also be **proportional allocation** which means that we allocate frames according to the size of the process.

$$S \sum s_i$$

S_i = virtual size used by process p_i

M = total number of frames

$$a_i = \text{allocation for } p_i = \frac{s_i}{S \times m}$$

7.11.2 Priority Allocation

This implementation uses a proportional allocation scheme that uses priorities rather than size.

7.11.3 Local Replacement

In this implementation, for each process, the kernel only selects from the process's own set of allocated frames; therefore a process has a fixed number of frames. The good thing about this is that when a process thrashes, it does not cause other processes to thrash. The bad thing about this is that it can lead to wasted memory if a process's working set decreases in size.

7.11.4 Global Replacement

For each process, the kernel selects a replacement frame from the set of all frames. This is better for dynamic working set which may grow or shrink. However, there can be domino-style thrashing. Most modern operating systems use global replacement to avoid wasting memory.

7.11.4.1 Avoiding Thrashing

To prevent thrashing, adopt a local allocation policy (wastes memory), do not schedule a process unless its working set of pages is in memory (implies pre-paging), and use a page fault frequency approach.

7.11.4.2 Page Fault Frequency Scheme

Measure the page fault rate of process and if its too low then take pages from that process and if too high, then allocate additional frames to the process. If many processes have high page fault frequency, swap out one or more of these processes (the process to swap out is determined by the scheduling algorithm). This may use process priority to affect the range of acceptable page fault frequencies.

8 File Systems

File systems are used for long term information storage, i.e.

- It must be possible to store a very large amount of information
- Information must survive termination of a process using it
- Multiple processes must be able to access information concurrently

A **file** is an abstraction that is a logical unit of information created by processes managed by the operating system. The **file system** is the subsystem that deals with files (part of the kernel). In UNIX, extensions are just conventions that sometimes determine how a file should be opened, but are mostly information for the user: the kernel does not care.

To read a byte using **sequential access**, all previous bytes must be read first. To read a byte using **random access**, you can read any byte at any time.

Directories are persistent data structures for organizing and maintaining information about files (they are also files themselves). The **file system** consists of the kernel's file system subsystem organized into layers and the file system data is stored in storage.

8.1 Disks

A **disk** can be divided up into one or more partitions each with an independent file system. The **master boot record (MBR)** is a 0.5 KB partition that looks at the partition table and decides which to use. The **partition boot record** finds the boot loader for the kernel: there is a PBR in each partition. To boot...

- The computer's BIOS runs first, finds boot storage device, and loads the MBR
- MBR contains partition table and boot code that finds and loads PBR
- Boot block = PBR = boot code to find kernel loader inside the file system of the partition (which itself then finds the kernel)
- Superblock stores key parameters about the file system

8.2 Allocation Methods

A **disk block** is the minimal unit of allocation. Physical addresses are some kind of tuple that identifies a single block on a disk. The kernel drive translates physical addresses to block addresses. The logical file address is relative to the beginning of the file and is measured in bytes.

8.2.1 Contiguous Allocation

In contiguous allocation, each file occupies a set of contiguous disk blocks. Entire blocks are used independent from the file size (i.e. file size is not a consideration).

Pros of contiguous allocation:

- **Simple:** only need to maintain starting location and length
- **Performance excellent:** the entire file system can be read in a single operation

Cons of contiguous allocation:

- **Wasteful of space:** external fragmentation, may require compaction
- **Files cannot grow:** need to know the final size upfront

8.2.2 Linked List Allocation

In linked list allocation, each file is a linked list of disk blocks. The first bytes of the block contain a pointer to the next block. Blocks may be scattered anywhere on the disk.

Pros of linked list allocation:

- **Simple:** only need starting address
- **Free-space management system:** No waste of space

Cons of linked list allocation: **No efficient random access:** must go from block to block to find the data

8.2.3 Linked List Allocation with Table

This implementation is a variant of linked list allocation that uses a **file-allocation table (FAT)**. We keep the pointer in a separate table rather than as part of data blocks. In the table, each entry corresponds to the disk block number and contains a pointer to the next block. The table is on disk but copied into the memory for speed.

Pros: Faster random access

Cons: Entire table must be in memory (huge problem for very large disks)

8.2.4 Indexed Allocation

In indexed allocation, we associate a block-sized data structure called index-node (i-node) to each file which lists the attributes and lists the disk addresses of the file's blocks. The last address(es) is for another block of addresses.

Pros of indexed allocation:

- efficient random access
- dynamic access without external fragmentation
- lower overhead than FAT

Con of indexed allocation is that it can be hard to support large files (that do not fit the i-node)

An i-node can contain a pointer to direct block (data blocks), single indirect, double indirect, triple indirect block, etc...

Indirection allows files to grow and be incredibly large. Indirection refers to when i-node points to indirection blocks which can point to other indirection blocks or data blocks.

8.3 Implementation of Directories

To implement directories we first start by mapping the file name to some attributes and data. Directories are stored as special files: the file contents are the directory entries. **Directory blocks** include an entry for each file within that directory; directories with many files may span multiple blocks. Every file and directory is represented by an **i-node** (index node). This index node contains two kinds of information:

- Metadata describing the file or directory's owner, access rights, etc...
- Location of the file or directory's data blocks on disk

For a directory, the data blocks simply store the directory's file entries: file name + corresponding i-node number. The beginning part of the file system contains all the inodes.

Because it can be convenient for a shared file to appear in different directories, the directory system is implemented as a **directed acyclic graph**. We can thus create multiple names for the same inode in the same directory or even in different directories. All the names must be on the file system since each file system has its own inode. In UNIX and Windows NTFS this is called a **"hard link"**, which refers to the specific location of physical data. Use the Linux command `ln <source> <destination>` to create a link.

8.3.1 Symbolic Link

A **symbolic link** is a special kind of file, the contents of which indicate the location of another file. These files have a special status that is indicated by the file attribute "sym link" bit being set. During name interpretation, symbolic link is read and used to find the read file. You can use the Linux command `ln -s <source> <destination>` to create a symbolic link. To read the contents of a symbolic link file itself, you can use the `readlink` command. Only the original file actually links to the i-node; symbolic links only link to path names, not i-nodes. When the owner removes the real file, the link now points to a non-existent path name.

Disadvantage: extra overhead to parse symbolic links

Advantage: can link across file systems.

8.4 Virtual File System

Many different file systems are in use even on the same computer and under the same operating system. From the user's point of view, there is a single file-system hierarchy which encompasses multiple (incompatible) file systems. The UNIX kernel uses a **virtual system layer** (VFS) which integrates multiple file systems into a single structure and abstracts over the details of each file system. All system calls related to file are directed to the VFS layer using standard POSIX calls. The lower VFS interface calls into the file system code. Internally, the VFS is object oriented and contains a superblock, v-node, and directory.

9 Input/Output

The operating system...

- **Controls** all the computer's I/O devices
 - Issues commands to the devices
 - Catches and manages interrupts
 - Handles errors
- **Provides an interface** between the devices and the user-space software that is simple, easy to use, and potentially the same for all devices

One type of device that may perform I/O with the system are **block devices** which store information in fixed-size blocks. In block devices, each block has its own address and transfers are in terms of entire blocks.

Another type is a **character device** which delivers or accepts a stream of characters without any block structure. These are not strictly addressable and do not have any *seek* operations.

9.1 I/O Hardware

A **port** is a connection point for a device. There are **buses** such as peripheral buses such as PCI and PCIe and **expansion buses** which connect relatively slow devices. There will be a **controller (host adapter)** which operates the port, bus, and devices. Sometimes this controller will be integrated on the motherboard while other times it may be on a separate circuit board. The controller contains the processor, microcode, private memory, bus controller, etc... **Buses** are used to handle traffic between I/O devices and the processor. **PCI/PCIe** connects with speed graphics, networking, etc and connects to low speed buses.

Each device has both a mechanical and electronic component. The **mechanical component** (the device itself) will have its own electronic controller too, but this is invisible to the kernel. The **electronic component** (chip on the motherboard) consists of the controller (host adapter) which may be able to handle multiple or only a single device at once. The controller's tasks are to convert bit streams to bytes, perform error correction when necessary, and move data to/from host memory.

9.2 Communication Mechanisms Between CPU and I/O Devices

Communication between the CPU and I/O devices can happen through registers or a buffer. The kernel device driver reads/writes to **registers** for configuration/initialization, reading/writing data, and to check the status of a device. The kernel device driver reads/writes to **buffers** to move data.

9.2.1 Register Based I/O

Each controller register has an I/O port number. There are special instructions to access the I/O port space: CPU reads in from device I/O port to CPU registers (IN REG, PORT) and the CPU writes to device I/O port from CPU register (OUT PORT, REG). These instructions are privileged (kernel-use only) and must be written in assembly. There are separate I/O port space and memory space.

9.2.2 Memory Mapped I/O

All controller registers and buffers are mapped into the normal RAM address space. Each controller register is assigned a unique memory address: there is not actual RAM memory for that address. Such addresses may be at the top of the physical address space.

Advantages of memory mapped I/O:

- Software can read/write to I/O devices as to any other memory location. All normal load/store memory instructions apply to I/O devices as well. Read/write easy to program in C using a normal pointer.
- Protection of I/O device registers and buffers is the same as memory protection. Can be accessed from kernel mode and user mode (if you have the right permissions).

Disadvantages of memory mapped I/O:

- Need to disable caching to memory-mapped I/O
- For one address, RAM controller and I/O devices will examine the memory reference.

9.2.3 Hybrid I/O Implementation

A **hybrid I/O implementation** uses both I/O ports and memory-mapped I/O. To **send data to I/O device**, the CPU wants to read a word either from memory or I/O port:

1. Puts the address it needs on the bus address line
2. Asserts a read signal on a bus control line
3. A second signal line to tell I/O space or memory space
 - If it is memory space, memory responds to request
 - If it is I/O space, I/O device responds to the request
 - If there is only memory space (no I/O ports)

- RAM controller and every I/O device compares the address lines to the range of addresses that it services
- If the address falls in its range, it responds to the request
- Since no address is ever assigned to both memory and I/O device, there is no ambiguity and no conflict

The CPU can request data from an I/O controller one byte/word at a time; this wastes the CPU's time for big transfers. In **offloaded communication**, the CPU can offload data transfers and do something else in the meantime.

9.2.4 Direct Memory Access (DMA)

The **direct memory access** (DMA) controller transfers data for the CPU from/to the I/O devices and between I/O devices. This implementation requires a DMA controller on the device host controller or on the motherboard. A **DMA controller** contains registers to be read/written by the kernel include the memory address register, byte counter register, and control registers to indicate the direction of the transfer and the transfer unit. The CPU commands the DMA controller by writing to control registers of DMA device. DMA is no faster than CPU but frees the CPU from having to wait for everything to finish and thus is used to take care of all big memory transfers. To do **multiple data transfers** at once (multiple channels), we need multiple sets of registers (one for each transfer channel).

9.2.4.1 Data Transfer Bus Mode

In **word-at-a-time** mode, the DMA controller requests for transfer of one word. If the CPU wants the bus, it has to wait: this is called **cycle stealing**.

In **block** mode, the DMA controller acquires the bus, issues a series of transfers, and releases the bus: this is also called **burst** mode. An advantage of block mode is that it takes less time for acquiring the bus. A disadvantage of block mode is that it can block the CPU and other devices for a long time if there is a long burst.

In **fly-by** mode, the DMA controller demands the device controller to transfer data directly to the destination.

In **flow-through** mode, data transferred passes through the DMA controller. The DMA controller first reads data into an internal buffer and then writes it to the destination. This scheme requires an extra bus cycle per word transferred, but is more flexible because it can perform both device-to-device and memory-to-memory copies.

9.2.4.2 Interrupts

DMA controllers are able to read and write **directly from/to memory** and/or I/O devices with no kernel intervention in the meantime. DMA hardware enables data transfer to be accomplished without using the CPU except for configuration. Interrupts are sent to the interrupt controller from the device; the controller issues the interrupt to the CPU and the CPU acknowledges the interrupt.

There are two types of interrupts: external interrupt and internal interrupts (referring to whether the interrupt comes from within or outside of the CPU). **Processing** any type of **interrupt** is almost identical:

1. Kernel (with help from CPU) is responsible to save the state of the interrupted process
2. Kernel processes the interrupt
3. The interrupted process is resumed after the interrupt is handled by the kernel.

External interrupts deal with interactions with the I/O devices.

Internal interrupts (traps, exceptions) occur entirely within the CPU. They are used to handle exceptional conditions that occur during the execution of programs such as divide by zero.

9.2.4.3 Processing Interrupts

The CPU checks for interrupts after executing every instruction. In the case that there is an interrupt, the hardware looks up the **interrupt vector** (configured by the kernel at boot time) to fetch a new program counter. The new program counter points at the code of the right kernel **interrupt handler** which is then executed.

9.3 Goals of I/O software

Some of the goals of I/O software include the following:

- **Device independence:** programs can easily access any I/O device without having to specify the device in advance.
- **Uniform naming:** name of a file or device should be a string independently of device
- **Error handling:** errors should be handled as close to the hardware as possible

- **Synchronous** (blocking) vs. **asynchronous** (interrupt-driven)
- **Buffering** involves considerable copying which may negatively impact I/O performance
- **Shareable** vs. **Dedicated devices**: Some I/O devices, like disks, can be used by many users at the same time while other devices, such as printers, have to be dedicated to a single user.

9.3.1 Programmed I/O

In programmed I/O, the CPU does all the work. The CPU writes/reads a byte/-word at a time from/to main memory to/from the device. The CPU makes a request, then waits for the device to become read. Buses are only byte/word wide, so the last few steps are repeated for large transfers.

CPU time is wasted; if the device is slow, the CPU may wait a long time. The CPU continuously polls the device to see if it is ready to accept another byte/word: **polling** or **busy waiting**.

Listing 1: Writing a string to the printer using programmed I/O

```

1 copy_from_user(buffer,p,count);          /* p is the kernel buffer */
2 for (i=0;i>count;i++) {                  /* loop on every character */
3     while(*printer_status_reg != READY);  /* loop until ready */
4     *print_data_register = p[i];          /* output one character */
5 }
6 return_to_user();

```

Pros: simple to implement/understand

Cons: CPU cannot be used for anything else until I/O ends, polling is inefficient

9.3.2 Interrupt-driven I/O

In interrupt-driven I/O, the CPU does the work but interrupts tell when to do that work. We use interrupts because I/O devices are slower than memory or CPU and there is uncertainty of when the device will be ready. The OS needs to know when the I/O device has completed an operation and when the I/O operation has encountered an error. **Instead of waiting**, the CPU continues with other computations. The device, then, interrupts the processor when operation completes or there is an error.

Listing 2: Writing a string to the printer using interrupt driven I/O

```

1 /* Executes at the time the print system call is made */
2 copy_from_user(buffer,p,count);
3 enable_interrupts();
4 while(*printer_status_reg != READY);
5 *printer_data_register = p[0];
6 scheduler();
7
8 /* Interrupt service procedure for the printer */

```

```

9  if (count == 0) {
10     unblock_user();
11 } else {
12     *printer_data_register = p[i];
13     count = count - 1;
14     i = i + 1;
15 }
16 acknowledge_interrupt();
17 return_from_interrupt();

```

Pros: Better than programmed I/O because CPU can execute another process in the meantime. Also avoid inefficient polling

Cons: An interrupt occurs on every character. Interrupts take time, so it wastes CPU cycles.

9.3.3 I/O using DMA

In I/O using DMA, the DMA controller does all the work, it uses interrupts for notification, but CPU needs to program the DMA controller. The DMA controller reads/writes directly from/to memory. It has access to the system bus independent of the CPU.

1. The CPU sets up the operation
2. DMAC does the transfer
3. When the transfer is complete, the DMAC notifies the CPU with a single interrupt

Listing 3: Printing a string using DMA

```

1  /* Syscall code. Executed when the print system call is made */
2  copy_from_user(buffer,p,count);
3  set_up_DMA_controller();
4  scheduler();
5
6  /* Interrupt handler code */
7  acknowledge_interrupt();
8  unblock_user();
9  return_from_interrupt();

```

10 Disks

There are different types of persistent storage technologies. **Magnetic disks** such as HDD and floppy disks have moving parts and on physical media, the reads and writes are equally fast. **Solid-state drives** such as SSD and USB thumb drives do not have any moving parts (fully electronic) and on physical media, reads are faster than writes. **Optical disks** such as CDs, DVDs, and blu-ray are mechanical (have moving parts) and are considered write-once, read

multiple times implying that if re-writing, then writing is way slower than reading. There are, however, some **problems** with using a single storage device such as any of the ones previously mentioned:

- **Slow**: disk is slower than other components of a computer
- **Modest capacity** compared to organizational needs **Unreliable**: a disk may fail totally (where it does not read/write data) or may have bad sectors which would imply that you cannot read what you have written.

10.1 RAID

RAID stands for **redundant array of inexpensive disks** which consists of building a storage system that has an array of drives where data is distributed over multiple drives, data can be read in parallel, and data may be replicated. This system will appear like a single disk and can be applied to any storage technology.

To implement a RAID system, one may use a **RAID controller card** (hardware). In this implementation, each drive is connected to the RAID controller and the RAID controller can issue read in parallel among all disks. RAID can also be implemented at the operating system level through **software**. This implementation is cheaper but slower, and typically fewer RAID levels/configurations are supported in OS than in specialized hardware.

RAID is a layer of hardware/software that does more than just store the data. The following are some of the different RAID levels:

10.1.1 RAID 0 (Striping)

In RAID 0, data is striped (split) across disks. The total storage space across all disks is divided into **strips** where every strip includes k sectors. Strips are mapped round-robin to consecutive disks, and multiple steps can be fetched in parallel. There is **no redundancy** in RAID 0.

RAID 0 is **fast** because it can do reads and writes in parallel.

If one disk fails then all the files will be messed up because every file will have a piece of it missing). For smaller requests (small files), there are no performance gains; if software asks for one sector at a time, result is correct, but there is no parallelism. RAID 0 has **worse reliability** than a single disk: adding more disks means that failures happen more often. To increase reliability, we can replicate data across multiple disks...

10.1.2 RAID 1 (Mirroring)

RAID 1 achieves redundancy by **duplicating** all the data. Every disk has a mirror disk, storing the same data. A **read** can be serviced by either of the two

disks which contains the requested data (up to twice as fast as RAID 0). A **write** must be executed on both disks but can be done in parallel (same speed as RAID 0).

RAID 1 requires double the amount of disks which can be **expensive**. When installing a new drive to replace a broken one, have to copy an entire disk which will take a non-trivial amount of time.

10.1.3 RAID 2 (bit-level striping with hamming-code)

RAID 2 stripes (splits) each word/byte into bits, one bit per disk. It also adds **hamming code bits**, one bit per extra disk. The hamming code is used to reconstruct the original data. The all drives **write in sync** the same word/byte.

RAID 2 is too complicated to do in practice. Additionally, all drives need to be rotationally in sync. RAID 2 requires very performant disk controllers: many disks for a good ratio of data to hamming bits.

10.1.4 RAID 3

RAID 3 is the same as RAID 2 but uses **parity** bits instead of hamming codes. Given N bits b_1, b_2, \dots, b_N , the parity bit P is the bit 0,1 that yields an even number of "1" bits in the set b_1, b_2, \dots, b_N, P . If any bit in b_1, b_2, \dots, b_N is lost, it can cause the remaining P bits to recover it.

To implement RAID 3, we need to add an extra disk to store the parity bits. This handles one failing drive only. Additionally, RAID 3 is complex to implement because all drives need to be rotationally in sync.

10.1.5 RAID 4

RAID 4 uses **large strips** with a **parity strip** (XOR). This is similar to RAID 3's parity bit but it is done at block level rather than bit level. If a drive crashes, lost bytes can be recomputed with parity. You must recompute parity whenever you change something on the disks.

RAID 4 can handle only failing drive. Additionally, RAID 4 performs poorly for small updates. If one sector is changed, you need to read all drives in order to recalculate the parity which must then be rewritten on parity drive. This implementation results in a **heavy load** on the parity drive.

10.1.6 RAID 5

RAID 5 uses interleaved data and parity blocks. You rotate the assignment of data blocks and parity blocks among the drives. This avoids the bottleneck of a single disk for a parity disk and also allows multiple reads/writes to occur in

parallel.

RAID 5 handles one failing drive only. Reconstructing the contents of a failed drive is expensive.

10.1.7 RAID 6

RAID 6 uses an additional parity block. The data is striped across the disks with **two parity blocks instead of one**: two different parities. In this implementation, writes are a bit more expensive because of the parity calculations; reads incur no performance penalty. RAID 6 protects against two simultaneous disk failures and is thus good for very large RAID arrays.

10.2 Layering RAID levels

Note that layering one RAID level on top of the other provides both performance and redundancy...

10.2.1 RAID 0 + 1

Striped disks (RAID 0) are then mirrored (RAID 1).

After one disk fails, one stripe is down, and the whole system will fail if any other disks fails on the other side of the mirror.

10.2.2 RAID 1 + 0

Mirrored disk pairs (RAID 1) are then striped (RAID 0).

After one disk fails, no pair is down, and the whole system will fail only if another disk fails in the same pair. This layering is more reliable than RAID 0 + 1.

10.2.3 RAID 5 + 0

Stripes with distributed parity (RAID 5) and then striped (RAID 0).

This layering increases reliability. RAID 5 can handle a maximum of one failed drive before data loss, but striping increases that to two failed drives.

10.2.4 RAID 6 + 0

Stripes with double distributed parity (RAID 6) and then striped (RAID 0).

RAID 6 can handle a maximum of two failed drives before data loss, but striping increases that to four failed drives.

10.3 Error Handling

Storage devices must be capable of dealing with **errors**. **Magnetic disks** are difficult to make perfect and without defects. There are manufacturing defects which introduce bad sectors, what you read may not be what your wrote, and there may be too few bits in which case you can use error correcting codes (if implemented). **Solid-state drives** are the same as magnetic disks but the unit is the block, not a sector. Errors in SSDs can be derived from bad block management in the SSD firmware.

To deal with bad sectors: **sector errors**.

- **Internal hardware controller.** Before the disk is shipped it is tested. For each bad error, use spare sectors to substitute in. If the number attempting to be accessed corresponds to a bad sector, then the controller will direct the request to the designated spare sector that replaces the erroneous one.
- **Operating system.** Do the same as with the hardware controller except in software. Acquire a list of bad sectors; nbuid remapping tables; use the remapping tables when accessing the disk.

ECC (error correcting code) protects against small sector errors while RAID protects against whole-disk failures. However, neither of these are a substitute for a backup system.