

## Introduction to Operating Systems

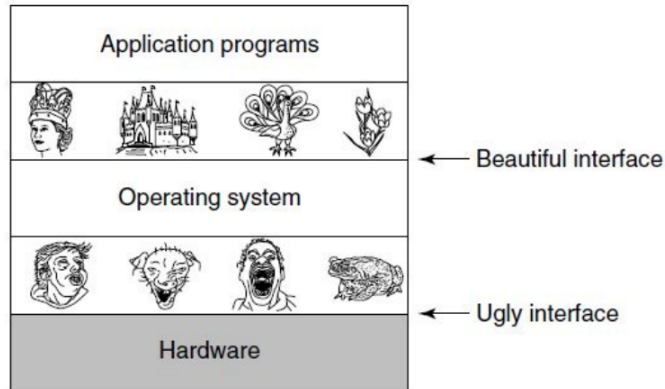
## Content:

- Basic Principles of Operating Systems
- Computer Hardware Review

## Basic Principles of Operating Systems

**Operating System (OS)** as an extended machine

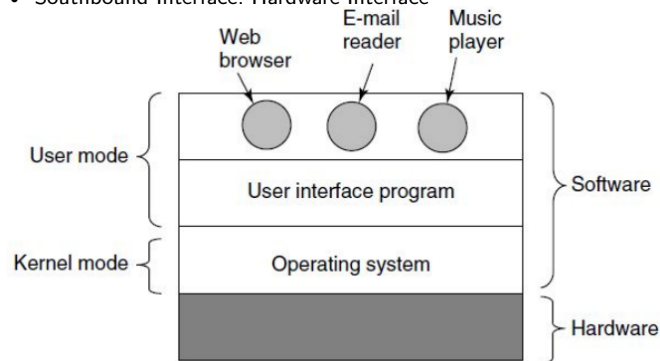
- Software that manages computer hardware
- Provides services for programs
- Acts as intermediary between user and hardware



Hardware is very complicated: The job of the operating system is to create good abstractions and manage the abstract objects thus created.

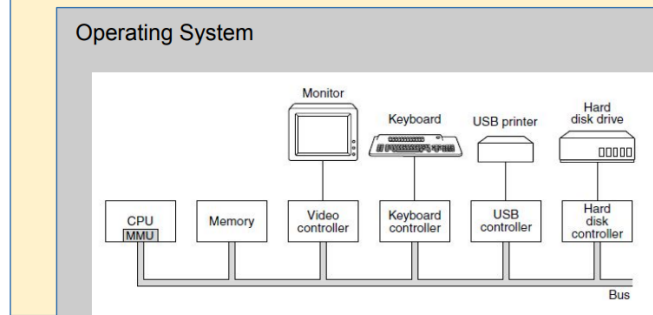
**Operating System (OS)**

- User Mode vs Kernel Mode
- Northbound Interface: User Interface
- Southbound Interface: Hardware Interface

**OS as a Resource Manager**

- Process Management
- Memory Management
- File System Management
- Device Management
- Security

## Applications



The OS controls resource usage, grants resource request, accounts for usage and mediates conflicting requests.

## Computer Hardware Review

**CPU** Central Processing Unit

- Basic cycle: Fetch, Decode, Execute
- CPUs feature some registers to hold key variables and temporary results
- Special registers for internal use: Program Counter (PC), Stack Pointer (SP), Program Status Word (PSW), etc.

CPUs and their Instruction Sets are architecture-specific:

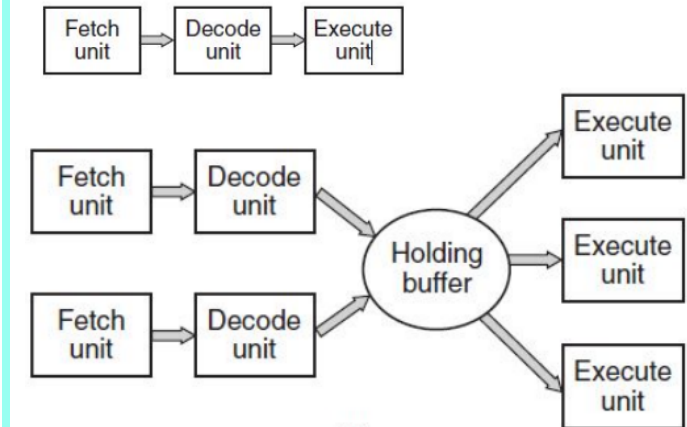
- ARM, RISC, X86, etc.
- Instructions are classified along Execution Privileges, enforced by CPU:
  - Intel: User Mode (limited set of instructions) vs Kernel Mode (full set of instructions)
  - ARM: UnPrivileged Mode vs Privileged Mode

**CPU cycles**

Basic cycle of every CPU:

- Fetches the first instructions from memory into registers
- Decodes instructions (determining type and operands)
- Executes instructions
- fetch, decode, execute, ...

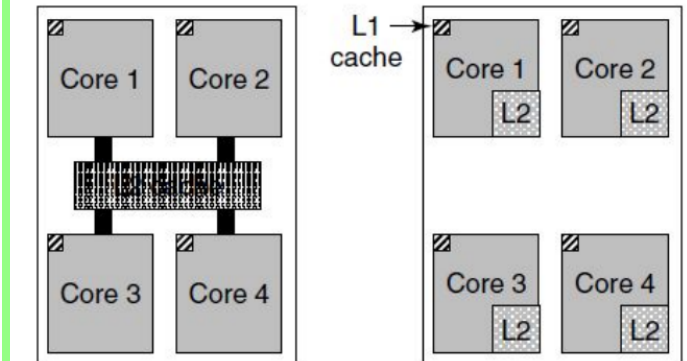
CPUs have multiple cores, each having multiple execution units and parallel pipelines:

**CPU Caches**

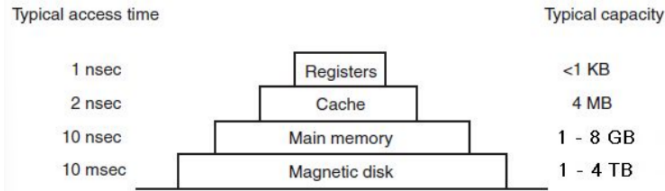
CPUS may have multiple levels of caches:

- L1: Small, fast, close to CPU
- L2: Larger, slower, further away
- L3: Even larger, even slower, even further away
- Caches are used to store frequently accessed data and instructions

Example: Quad-core chip with a shared L2 cache and a quad-core chip with separate L2 caches for each core.



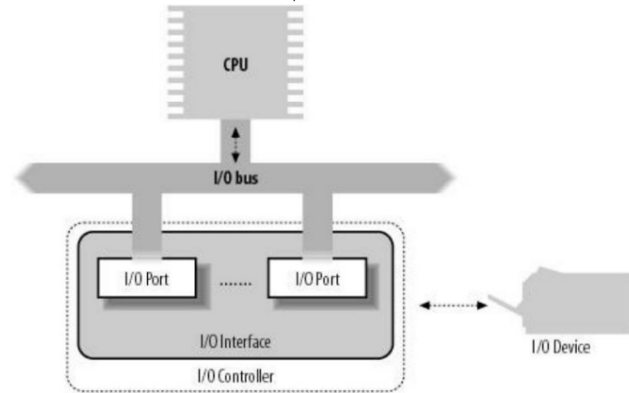
**Memory** The memory system is constructed as a hierarchy of layers, according to access time:



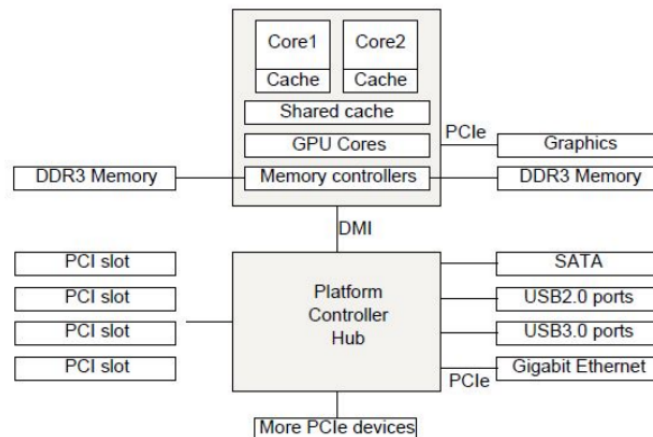
**Input/Output** Huge amount of I/O Devices, Hard Drives, Network Interfaces, Serial Ports, Keyboard, Mouse, Graphics, Cameras, etc

Devices connected with CPU via a Bus System (SW), I/O Interfaces (SW) and I/O Controllers (HW)

Software that talks to the I/O Controller (Commands/Responses), called Device Driver Runs in Kernel or User Mode. Built-into Kernel or modular and loadable at boot-/run-time



**IO and Bus System** X86 System with different Bus standards:



## Operating System Variants

- Mainframe OS: IBM z/OS, IBM z/VM
- Server OS: Windows Server, Linux, Solaris
- Multiprocessor OS: Windows, Linux, Solaris
- Personal Computer OS: Windows, MacOS, Linux
- Handheld Computer OS: Android, iOS
- Embedded OS: VxWorks, QNX
- Real-Time OS: VxWorks, QNX
- Sensor Node OS: TinyOS, Contiki
- Cloud OS: OpenStack, OpenNebula
- Smart Card OS: JavaCard, MULTOS

## Operating Systems vs. Distributions

- Operating System: Kernel, System Libraries, System Utilities
  - Distribution: OS + Applications, Tools, Documentation, etc.
- Example: Linux Kernel + GNU Tools + X11 + Gnome + Firefox + LibreOffice = Ubuntu **evtl.** add more info from slides

## Operation System Concepts

### Basics

Interacting with the OS:

- 

## Introduction to Operating Systems

## Defining Operating Systems

### The Operating System as an Extended Machine

The operating system serves as an abstraction layer between hardware and applications:

- Hardware is complex and difficult to program directly
- Operating systems create good abstractions for hardware resources
- These abstractions are implemented and managed by the OS
- Applications interact with these abstractions rather than directly with hardware

### User Mode vs Kernel Mode

Operating systems operate in two fundamental modes:

- **Kernel Mode:** Privileged execution mode with full access to all hardware resources and instructions
- **User Mode:** Restricted execution mode where applications run with limited privileges

The OS provides two key interfaces:

- **Northbound Interface:** Towards user applications
- **Southbound Interface:** Towards hardware

### The Operating System as a Resource Manager

The operating system:

- Controls resource usage
- Grants resource requests
- Accounts for resource usage
- Mediates conflicting requests
- Ensures fair and efficient resource allocation

## Computer Hardware Review

### Central Processing Unit (CPU)

The CPU follows a basic operation cycle:

- Fetch instructions from memory into registers
- Decode the instruction to determine type and operands
- Execute the instruction
- Repeat for subsequent instructions

Key CPU components include:

- Registers to hold variables and temporary results
- Special registers for internal use (Program Counter, Stack Pointer, Program Status Word)
- Execution units and pipelines for processing instructions
- Cache memory organized in hierarchical levels (L1, L2, L3)

CPUs enforce execution privileges through:

- Intel: Priority Rings (User Mode vs Kernel Mode)
- ARM: Privileged vs Unprivileged Mode

### Memory System

Memory is organized as a hierarchy of layers based on access time:

- Registers (fastest, smallest)
- Cache memory (L1, L2, L3)
- Main memory (RAM)
- Secondary storage (disks, SSDs)
- Tertiary storage (backup systems)

### Input/Output (I/O) System

The I/O system comprises:

- Various devices (hard drives, network interfaces, serial ports, keyboards, etc.)
- Bus systems to connect devices (PCI, PCIe, SATA, USB, etc.)
- I/O controllers (hardware)
- Device drivers (software)

Device drivers:

- Software that communicates with I/O controllers
- Can run in Kernel or User Mode
- May be built into the kernel or loaded as modules

## Operating System Variants

### Types of Operating Systems

Operating systems are categorized based on their purpose and environment:

- Mainframe Operating Systems
- Server Operating Systems
- Multiprocessor Operating Systems
- Personal Computer Operating Systems
- Handheld Computer Operating Systems
- Embedded Operating Systems
- Sensor Node Operating Systems
- Real-Time Operating Systems
- Smart Card Operating Systems

Operating System vs Distributions

- **Operating System (Kernel):** The core component that directly manages hardware resources
- **Distribution:** A complete package including kernel, utilities, and applications

Examples:

- Linux: The kernel available at <https://www.kernel.org/>
- Linux Distributions: OpenSUSE, RedHat, Ubuntu (kernel plus utilities)
- Microsoft: OS plus proprietary extensions and bundled software
- Apple: OS X kernel (Darwin) plus proprietary GUI components
- BSD: FreeBSD, NetBSD, OpenBSD (Unix-like systems)

Basic OS Concepts

Fundamental OS Concepts

- **Interacting with OS:** Through terminal or graphical user interface
- **Users:** Regular vs. privileged users
- **Data organization:** Files, directories, filesystems
- **Programs vs. Processes:** A program is a compiled executable, while a process is an instance of a running program
- **Multi-tasking:** Running multiple processes concurrently
- **Context-switching:** Switching execution from one process to another
- **System calls:** Interface for processes to request OS services
- **Inter-Process Communication:** Methods for processes to communicate
- **Signals:** Mechanism for notifying processes of events

OS Structures

Operating systems can be organized in different ways:

- **Monolithic Systems:** Single executable containing all OS functionality
- **Modular Systems:** Core kernel with loadable modules
- **Microkernels:** Minimal kernel with most services running in user space

To check CPU information in Linux:

```
1  # Display kernel version information
2  uname -a
3
4  # View CPU details
5  less /proc/cpuinfo
6
7  # Display CPU architecture information
8  lscpu
```

Working with Processes in Linux

Viewing running processes

- Use `ps aux` to display all processes
- Use `top` for an interactive process viewer
- Check environment variables with `env`

Process control

- Background a process with `&` or `bg`
- Bring to foreground with `fg`
- List background jobs with `jobs`
- Terminate a process with `kill PID`

System information

- View process hierarchy with `pstree`
- Check system call activity with `strace`

## Linux Primer

### Terminal Basics

#### Terminal

A terminal is a text-based user interface program that:

- Accepts text input via a prompt
- Interprets text as commands
- Executes operations based on user input
- Returns output to the user

Terminal commands can be:

- Binary programs loaded from disk (e.g., `mkdir`)
- Built-in functions of the terminal (e.g., `cd`)
- Operations like regular expressions for complex command preparation

#### Essential Terminal Commands

##### System information

- `whoami` - Display current user
- `uname -a` - Display kernel information
- `lsmod` - List loaded kernel modules
- `dmesg` - Display kernel messages

##### Command output manipulation

- `command > file.txt` - Redirect output to a file
- `command | grep pattern` - Filter output through `grep`
- `clear` - Clear terminal screen

##### System control

- `env` - Display environment variables
- `exit` - Exit the terminal
- `shutdown` - Shutdown the system

### User Management

#### User Management in Linux

Linux is a multi-user system with user management tools:

- `users` - List users currently logged in
- `who` - Show who is logged in
- `id` - Display user and group IDs
- `passwd` - Change user password
- `usermod` - Modify user account
- `groupmod` - Modify a group

Administrative access:

- `root` - Superuser account with full system privileges
- `sudo` - Execute a command as another user (typically `root`)
- `su` - Switch user
- `su -` - Switch user and load their environment

### Files, Directories, and Filesystems

#### File System Operations

Basic file and directory operations:

- `pwd` - Print working directory
- `cd` - Change directory
- `ls -al` - List all files with details
- `touch` - Create empty file or update timestamp
- `mkdir` - Create directory
- `tree` - Display directory structure as a tree
- `cp` - Copy files/directories
- `mv` - Move/rename files/directories
- `rm` - Remove files/directories

File permissions and ownership:

- `chown` - Change file owner
- `chmod` - Change file permissions

Disk and filesystem management:

- `fdisk -l` - List disk partitions
- `mount` - Mount a filesystem

#### Working with Files in Linux

##### Creating and viewing files

- `touch filename` - Create empty file
- `cat filename` - Display entire file contents
- `less filename` - View file with pagination
- `tail filename` - Display last 10 lines of file
- `tail -f filename` - Continuously monitor file for changes

##### Text editing with vim

- `vim filename` - Open file in vim
- Press `i` for insert mode
- Press `Esc` to exit insert mode
- Type `:w` to save
- Type `:q` to quit
- Type `:wq` to save and quit
- Type `:q!` to quit without saving

### Process Management

#### Process Management Commands

Commands for monitoring and controlling processes:

- `ps` - Display current processes
- `top` - Interactive process viewer
- `pstree` - Display process tree
- `pidof` - Find process ID of a program
- `kill` - Send signal to process
- `kill -9` - Force terminate process
- `killall` - Kill processes by name

### System Information and Configuration

#### System Information Commands

Commands for system information and configuration:

- `hwinfo` - Hardware information
- `lshw` - List hardware
- `/proc` - Virtual filesystem for kernel information
- `/etc` - Configuration files directory
- `sysctl` - Read/modify kernel parameters
- `systemctl` - Control the systemd system and service manager

Network-related commands:

- `ifconfig` - Configure network interfaces
- `ip` - Show/manipulate routing, devices, policy routing
- `dig` - DNS lookup utility
- `hostname` - Show or set system hostname

Package management:

- `apt-get` - Package handling utility

Working with the Linux terminal:

```
1  # Check current user
2  whoami
3
4  # View hardware information
5  lshw -short
6
7  # Monitor system messages
8  dmesg | grep usb
9
10 # Check disk usage
11 df -h
12
13 # Find a process ID
14 pidof firefox
15
16 # Install a package
17 sudo apt-get install http
```

#### OpenStack Lab Environment Setup

##### Initial setup

- Connect to ZHAW VPN if not at ZHAW facilities
- Navigate to OpenStack Horizon dashboard:  
`https://ned.cloudlab.zhaw.ch`
- Log in with provided credentials
- Change password if using default credentials

##### Creating SSH key pair

- Go to Compute → Key Pairs → Create Key Pair
- Download the private key file (\*.pem)
- Set appropriate permissions: `chmod 600 key.pem`

##### Creating a VM

- Go to Compute → Instances → Launch Instance
- Select Ubuntu image
- Attach to 'internal' network
- Set SSH key and security groups
- Launch VM and associate a Floating IP

##### Connecting to VM

- SSH to the VM: `ssh -i key.pem ubuntu@IP_ADDRESS`

## Linux Lab Basic Tasks

### Creating and managing files:

```
1  # Create a test file
2  touch delta.txt
3
4  # Add content to the file
5  echo "Hello, this is a test" > delta.txt
6
7  # View the file content
8  cat delta.txt
9
10 # Stop and restart the VM from OpenStack dashboard
11 # Then check if the file persists
12 cat delta.txt
```

## Booting an OS

### Boot Process Overview

#### Boot Process Stages

The typical boot process follows these common steps:

- Hardware initialization (BIOS/UEFI)
- Bootloader execution
- Kernel loading and initialization
- System initialization (services, environment)
- User interface presentation

### BIOS and Hardware Initialization

#### Basic Input Output System (BIOS)

BIOS is the first program loaded on boot:

- Stored in ROM on the motherboard
- Contains low-level I/O software for interfacing with devices
- Performs Power-On Self-Test (POST)
- Discovers devices by scanning PCI buses
- Initializes hardware devices
- Selects a boot device from the list in CMOS
- Reads the first sector from the boot device (Master Boot Record)
- Loads MBR into memory and executes it

### Bootloader

#### Bootloader Function

The bootloader (loaded from MBR):

- Accesses the boot partition containing the OS
- Loads the OS kernel into memory
- Sets up registers (Program Counter, Processor Status Word)
- Transfers control to the OS by jumping to the first instruction

#### GRUB (GRand Unified Bootloader)

GRUB is a common bootloader for Linux systems:

- Requires filesystem access to read the OS
- Contains device drivers and filesystem modules
- Provides a menu to select the OS to boot
- Loads the selected kernel and optional initial RAM disk (initrd)
- Passes control to the kernel

## OS Initialization

#### Kernel Initialization

The OS kernel initialization process:

- Queries hardware information
- Loads and initializes device drivers
- Initializes internal management structures (e.g., process table)
- Sets up virtual memory and memory management
- Creates system services
- Launches the first user process (init)

Linux kernel initialization phases:

- Architecture-specific assembly code
  - Sets up OS memory map
  - Identifies CPU type
  - Calculates total RAM
  - Disables interrupts
  - Enables MMU and caches
- C main() procedure
  - Initializes process tables, interrupt/system-call tables
  - Sets up virtual memory and page cache
  - Configures resource control
  - Loads drivers and initializes OS services

### BIOS vs UEFI

#### BIOS Limitations

Traditional BIOS has several limitations:

- Operates in 16-bit mode
- Relies on MBR (Master Boot Record) from 1983
- Limited partition number and size (2 TB max)
- Not designed for extendability
- Vulnerable to rootkit and bootkit attacks

#### Unified Extensible Firmware Interface (UEFI)

UEFI is the modern replacement for BIOS:

- Replaces MBR with GPT (GUID Partitioning Table)
- Supports arbitrary number of partitions
- Addresses disk space up to 2<sup>64</sup> bytes
- Uses unique UUIDs for partitions to avoid collisions
- Features modular design for extendability
- Includes Secure Boot to restrict which binaries can be executed
- Uses cryptographic signatures and X.509 certificates for verification

#### UEFI Functioning

UEFI uses an architecture-independent virtual machine:

- Executes special binary files compiled for UEFI (\*.efi)
- These files can be device drivers, bootloaders, or extensions
- Files are stored in the EFI System Partition (ESP)
- ESP uses FAT filesystem and can be reused in multi-boot systems
- EFI Boot Manager configures which EFI binary to execute

Examining UEFI boot configuration:

```
1 # Display current boot entries
2 sudo efibootmgr -v
3
4 # List contents of EFI system partition
5 sudo ls -la /boot/efi/EFI/
6
7 # View disk partition table format (MBR or GPT)
8 sudo fdisk -l
```

### Initial RAM Disk (initrd/initramfs)

#### Initial RAM Disk

The initial RAM disk addresses the "chicken-and-egg" problem:

- OS needs drivers to access hard drive and its filesystem
- These drivers are stored on the hard drive itself
- Solution: initrd/initramfs provides temporary root filesystem
- Contains kernel modules and basic device files
- Bootloader uncompresses both kernel and initrd into RAM
- Kernel mounts initrd as the initial filesystem
- Kernel uses tools found in initrd to find and mount the real filesystem

Examining GRUB configuration for kernel and initrd:

```
1 # View GRUB configuration
2 cat /boot/grub/grub.cfg
3
4 # Examine a specific menu entry showing kernel and
5 grep -A 10 "menuentry 'Ubuntu'" /boot/grub/grub.cfg
```

### System Services Initialization

#### System V Init

Traditional System V initialization:

- Init process runs initialization shell scripts from /etc/rc# directories
- Uses predefined runlevels (0-6) to determine system state
- Each runlevel has a set of services defined in scripts
- Dependencies among services are coded in the scripts themselves
- Results in complex initialization process

System V runlevels:

- 0: Halt - Shuts down the system
- 1: Single-user mode - Administrative tasks
- 2: Multi-user mode without networking
- 3: Multi-user mode with networking
- 4: Not used/user-definable
- 5: Multi-user mode with GUI
- 6: Reboot - Restarts the system

#### systemd

Modern initialization system (systemd):

- Replacement for SysV Init
- Provides coordinated and parallel service startup
- Features on-demand activation and runtime management
- Uses dependency-based service control logic
- Takes a holistic management approach for the entire system



systemd Units

- systemd organizes system components as units:
- Units encapsulate system objects (services, mounts, devices, etc.)
  - Units have states (active, inactive, activating, deactivating, failed)
  - Units can depend on other units
  - Most units are configured in unit configuration files
  - Some units can be created automatically or programmatically

- Common unit types:
- .service - A system service
  - .target - A group of systemd units (similar to runlevels)
  - .mount - A filesystem mount point
  - .device - A device file recognized by the kernel
  - .socket - An inter-process communication socket
  - .timer - A systemd timer

systemd Dependencies

- systemd supports various dependency types:
- **Requires=**: Units that must be started when this unit is started
  - **Wants=**: Units that should be started (but not required) when this unit is started
  - **Conflicts=**: Units that must be stopped when this unit is started
  - **After=**: This unit should be started after the listed units
  - **Before=**: This unit should be started before the listed units

systemd Service Unit File

A systemd service unit file consists of three sections:

```
1 [Unit]
2 Description=Example Service
3 Documentation=https://example.com/docs
4 After=network.target
5 Wants=network-online.target
6 Requires=example-dependency.service
7
8 [Service]
9 Type=simple
10 ExecStart=/usr/bin/example-service
11 Restart=on-failure
12 User=exampleuser
13
14 [Install]
15 WantedBy=multi-user.target
```

- The sections have specific purposes:
- Unit : Metadata and dependencies
  - Service : Execution configuration for services
  - Install : Configuration for enabling the unit

Working with systemd

Viewing unit status

- `systemctl status -all` - Show status of all units
- `systemctl status SERVICE` - Show status of specific service
- `systemctl list-units -t service` - List only service units
- `systemctl list-unit-files -type service` - List service unit files
- `systemctl cat SERVICE` - View the content of a unit file

Managing units

- `systemctl start SERVICE` - Start a service
- `systemctl stop SERVICE` - Stop a service
- `systemctl restart SERVICE` - Restart a service
- `systemctl enable SERVICE` - Enable service autostart
- `systemctl disable SERVICE` - Disable service autostart

Working with targets

- `systemctl list-units -t target` - List available targets
- `systemctl get-default` - Show default target
- `systemctl set-default TARGET` - Set default target

Creating a Simple systemd Service

Creating a custom service to write to a file when started:

```
1 # Create a service unit file
2 sudo nano /etc/systemd/system/myservice.service
3
4 # Add the following content
5 [Unit]
6 Description=My Simple Service
7 After=network.target
8
9 [Service]
10 Type=simple
11 ExecStart=/bin/bash -c 'echo "Service started at
12   $(date)" > /home/ubuntu/service_started'
13 Restart=on-failure
14
15 [Install]
16 WantedBy=multi-user.target
17
18 # Reload systemd to recognize the new service
19 sudo systemctl daemon-reload
20
21 # Start the service
22 sudo systemctl start myservice
23
24 # Verify it worked
25 cat /home/ubuntu/service_started
26
27 # Enable auto-start on boot
28 sudo systemctl enable myservice
```

systemd User Services

- systemd supports user instance services:
- Services managed by individual users without requiring root privileges
  - User service units are stored in:
    - `/usr/lib/systemd/user/`: Units provided by packages
    - `~/.local/share/systemd/user/`: User-installed package units
    - `/etc/systemd/user/`: System-wide user units
    - `~/.config/systemd/user/`: User's own units
  - By default, user services start on login and stop on logout
  - Lingering allows user services to start at boot without login

# Processes and Threads

## Process Model

### Programs vs. Processes

A program is fundamentally different from a process:

- **Program:** A compiled executable (binary)
  - Set of CPU instructions and related data
  - Targets a specific platform (OS + hardware)
  - Static entity stored on disk
- **Process:** An active instance of a program
  - Has a program, input, output, and state
  - Dynamic entity in memory
  - Multiple instances of the same program can run as separate processes

### Process Characteristics

Processes have several important characteristics:

- Can run sequentially or in parallel (multi-tasking)
- Selected for execution by the OS scheduler
- Associated with an owner (defines access privileges)
- Run within an environment (environment variables)
- Can run in foreground (interactive) or background (non-interactive)
- Can be user processes or system processes

Viewing process information in Linux:

```
1 # List all processes with details
2 ps aux
3
4 # Interactive process viewer
5 top
6
7 # Show environment variables
8 env
9
10 # Run a process in the background
11 wc /dev/zero &
12
13 # List background jobs
14 jobs
15
16 # Bring a background job to foreground
17 fg %job_number
```

## Process Creation

### Process Creation by the OS

When creating a process, the OS performs several operations:

- Loads the executable into RAM
- Sets up the memory map
- Updates scheduler and process table entries
- Sets program counter and stack pointer
- Switches from system mode to user mode

### Memory Layout

Process memory is typically divided into several segments:

- **Text:** Contains CPU instructions and constant data
- **Data:** Global and static variables
  - Initialized data segment: Pre-initialized variables
  - Uninitialized data segment (BSS): Zero-initialized variables
- **Stack:** Local variables and function call information (LIFO structure)
- **Heap:** Dynamically allocated memory (controlled by the programmer)

### Process Creation Events

Processes can be created in several ways:

- System boot (initial processes)
- User request (launching an application)
- Process creating another process (fork/exec)
- Scheduled creation (cron jobs)
- System request (responding to interrupts)

### Parent-Child Process Relationship

When a process creates another process:

- The creating process becomes the parent
- The new process becomes the child
- Child's memory map is initially a copy of the parent's memory map
- Two memory handling approaches:
  - Distinct address space: Separate memory regions for parent and child
  - Copy-on-write: Memory shared until a change is made by either process
- Forms a process hierarchy (starting with PID 1)

### Linux Process Creation

Linux terminal command execution involves process creation:

```
1 # When you run a command in the terminal:
2 ls -la
3
4 # The shell:
5 # 1. Forks itself (duplicate the process)
6 # 2. Child process executes the 'ls' binary
7 # 3. Parent waits for child to complete
8 # 4. Shell continues after child terminates
```

This process can be visualized with pstree showing the parent-child relationships.

## Process Termination

### Process Termination

A process can terminate in several ways:

- Voluntary normal exit (job completed, return from main())
- Voluntary error exit (required resource unavailable)
- Involuntary error exit (segmentation fault, division by zero)
- Termination by another process (kill signal)

Process termination in Linux:

```
1 # Start a process
2 wc /dev/zero &
3
4 # Get its process ID
5 pidof wc
6
7 # Terminate the process
8 kill -9 <PID>
```

## Process States

### Process States

Each process has a life-cycle with specific states:

- **Running:** The process is currently executing on the CPU
- **Ready:** The process is ready to run but waiting for CPU allocation
- **Blocked:** The process is unable to run (waiting for an event or resource)
  - Dependencies not met for running
  - Waiting for external resource
  - Waiting for I/O completion
  - Sleeping
  - Under job control or debugger

### Process State Changes

Processes change state for various reasons:

- **Timer expiration:** CPU allocation time ended
- **Interrupt:** Hardware/resource calls for service
- **Page fault:** Data not in memory, requires disk access
- **System call:** Explicit OS service request

State changes involve a switch between user mode and system (kernel) mode.

### User Mode vs. System Mode

CPU supports two execution modes:

- **User Mode:** Limited privileges
  - Application logic execution
  - Application data manipulation
  - Restricted access to hardware and system resources
- **System Mode (Kernel Mode):** Full privileges
  - System management operations
  - Hardware access
  - Interrupt handling
  - Device management

### Context Switch vs. Mode Switch

Two types of switches occur during system operation:

- **Mode Switch:** Transition between user mode and kernel mode
  - Occurs during system calls, interrupts, exceptions
  - Same process continues execution in different mode
  - No scheduler involvement
- **Context Switch:** Changing from one process to another
  - Involves saving the state of the current process
  - Loading the state of another process
  - Typically involves mode switches (user → kernel → user)
  - Requires scheduler involvement



## Process Management

### Process Control Block

The OS maintains a Process Control Block (PCB) for each process:

- Process identification (PID, UID, GID)
- Process state information
- Program counter and CPU registers
- CPU scheduling information (priority)
- Memory management information
- I/O status information
- Accounting information

In Linux, PCBs are implemented as `task_struct` entries in the process table.

## Threads

### Threads

Threads are execution entities within a process:

- Created and owned by a process
- Share the address space and all data of the owning process
- Allow multiple executions to take place within the same process environment
- Enable a process to continue even when some operations would block
- Cooperate toward the objective of the owning process

Each thread has:

- Program counter (tracks next instruction)
- Registers (hold working variables)
- Stack (with frames for procedure calls)

### Thread Implementation Approaches

Threads can be implemented in different ways:

- **User Space Threading (M:1):**
  - All threads in user space appear as a single process to the OS
  - Thread functionality provided by a library
  - Scheduling handled by the process (no mode switch required)
  - Based on cooperation (threads voluntarily yield CPU)
  - Disadvantages: Single thread can monopolize CPU time, blocked threads block the entire process, no SMP advantage
- **Kernel-supported Threading (1:1):**
  - One kernel structure per user thread
  - Kernel schedules all threads
  - Advantages: Better handling of blocking I/O, full SMP exploitation
  - Disadvantages: Higher overhead, slower creation/removal, mode switching for scheduling

### Thread Implementation in Linux

Linux doesn't have a dedicated concept of threads:

- All threads are standard processes (tasks)
- A thread is a process that shares resources with other processes
- Shared resources can include: address space, file descriptors, sockets, signal handlers, etc.
- Implemented using system calls: `fork()` and `clone()`
- Two implementation frameworks:
  - LinuxThreads (legacy)
  - NPTL (Native POSIX Thread Library, current standard)

### Linux Process vs. Thread Creation

Linux offers different system calls for process and thread creation:

- `fork()`: Creates a child process by duplicating the parent
  - Child and parent run in separate memory spaces
- `clone()`: Provides precise control over what execution context is shared
  - Allows sharing of address space, file descriptors, signal handlers, etc.
  - Used to implement threads in Linux

### Kernel Threads

Kernel threads are special processes that:

- Run exclusively in system (kernel) mode
- Have a PID and state like any process
- Are listed in the process table but flagged as "kernel thread"
- Are scheduled and dispatched like regular processes
- Are uninterruptible (need to voluntarily yield CPU)
- Listen to kernel signals

Kernel threads are organized around working queues and thread pools:

- Working queues ordered by priority
- Mapped onto a pool of reusable kernel threads
- Number of threads dynamically managed based on workload
- Managed by the `ktthread` daemon

## Process and Thread Management in Linux

### Viewing process information

- `ps aux` - List all processes
- `ps -ef` - List processes in full format
- `ps -eLF` - List processes and threads
- `top` - Interactive process viewer
- `top -H` - Show threads in top
- `pstree` - Display process tree

### Creating processes

- Write a C program using `fork()` to create a child process
- Use `getpid()` to retrieve the process ID
- Use `sleep()` to pause execution

### Creating threads

- Use POSIX threads library (`pthread`)
- Include `<pthread.h>`
- Create threads with `pthread_create()`
- Join threads with `pthread_join()`

### Identifying zombie processes

- Create a parent process that doesn't wait for its child
- Child terminates while parent continues running
- Check process state with `ps` (state `SZ` indicates zombie)

### Creating Processes in C

Simple process creation with `fork()`:

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <unistd.h>
4
5 int main() {
6     pid_t pid = fork();
7
8     if (pid < 0) {
9         // Error
10        fprintf(stderr, "Fork failed\n");
11        return 1;
12    } else if (pid == 0) {
13        // Child process
14        printf("Child process: PID = %d\n", getpid());
15        sleep(5);
16    } else {
17        // Parent process
18        printf("Parent process: PID = %d, Child PID = %d\n", getpid(), pid);
19        sleep(10);
20    }
21
22    return 0;
23 }
```

### Creating Threads in C

Simple thread creation with POSIX threads:

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <pthread.h>
4 #include <unistd.h>
5
6 void* thread_function(void* arg) {
7     int thread_id = *(int*)arg;
8     printf("Thread %d: running\n", thread_id);
9     sleep(3);
10    printf("Thread %d: exiting\n", thread_id);
11    return NULL;
12 }
13
14 int main() {
15     pthread_t threads[2];
16     int thread_args[2];
17
18     for (int i = 0; i < 2; i++) {
19         thread_args[i] = i;
20         pthread_create(&threads[i], NULL,
21                       thread_function, &thread_args[i]);
22         printf("Main: created thread %d\n", i);
23     }
24
25     // Wait for threads to finish
26     for (int i = 0; i < 2; i++) {
27         pthread_join(threads[i], NULL);
28     }
29
30     printf("Main: all threads completed\n");
31     return 0;
32 }
```

## Scheduling

### Scheduling Fundamentals

#### Scheduling Problem Domain

Scheduling is a resource-time management activity:

- Focuses on managing CPU time allocation among processes
- Requirements vary by platform (mobile vs. supercomputer)
- Requirements vary by application type (batch processing vs. real-time)

Processes can be categorized by behavior:

- **CPU-bound:** Computationally intensive tasks (e.g., image processing)
- **I/O-bound:** Tasks that frequently wait for I/O operations (e.g., interactive applications)

#### Queueing Theory

Scheduling is based on queueing theory principles:

- Deals with waiting for and dispatching resources
- Aims to provide sufficient resources to avoid under-capacity
- Ensures urgent tasks are not kept waiting

#### Scheduling Queues

The scheduler manages several types of queues:

- **Ready queue:** Processes waiting to be executed
- **Device queues:** Processes waiting for specific devices
- **Job queue:** All processes in the system

The scheduler sorts the ready queue according to policy, and the dispatcher moves the process from the head of the queue to the CPU.

#### Scheduling Metrics

Scheduling policies can be evaluated using various metrics:

- **CPU utilization:** Keeping the CPU as busy as possible
- **Throughput:** Number of processes completed per time unit
- **Turnaround time:** Time from submission to completion
- **Waiting time:** Time spent in the ready queue
- **Response time:** Time from request to first response
- **Fairness:** Equal distribution of CPU time

#### When to Schedule

Scheduling decisions are made in several situations:

- When a new process is created
- When a process exits
- When a process blocks on I/O
- When an I/O interrupt occurs
- Regularly on a timer

Based on when scheduling occurs, schedulers can be:

- **Non-preemptive:** Only schedules when a process blocks or terminates
- **Preemptive:** Can interrupt a running process and schedule another

### Scheduling Algorithms

#### Simple Schedulers

Simple scheduling approaches:

- **Uniprocessing:** One machine, one task - no scheduler required
- **Multiprocessing:** Single task running at one time with job control
- **Multitasking:** Multiple tasks running with scheduler

#### First-In-First-Out (FIFO) / First-Come-First-Served (FCFS)

FIFO is a basic scheduling algorithm:

- Processes are executed in the order they arrive
- Single queue, no preemption
- Processes run to completion
- Simple to implement
- Non-preemptive
- No awareness of process type (interactive vs. batch)

#### Round Robin (RR)

Round Robin is a time-sharing algorithm:

- Each process runs for a fixed time slice (quantum)
- After quantum expires, process is moved to the end of the queue
- Preemptive, as running tasks are interrupted after quantum
- Arriving tasks have priority over adjourned tasks (by convention)
- No starvation, as all processes eventually get CPU time
- Performance depends on quantum size
  - Too large: approaches FIFO
  - Too small: too much context switching overhead

#### Multi-level Scheduling

Multi-level schedulers address priority needs:

- Multiple queues with different priorities
- Round Robin within each queue
- Challenge: High-priority tasks can cause starvation of low-priority tasks

#### Multi-level Feedback Queue

Addresses starvation in multi-level queues:

- Task priority sinks after each execution interval
- Low-priority queues may get larger time quantum
- Balances responsiveness with throughput

#### Fair Share Scheduling

Fair Share approaches aim for equal CPU time:

- Maintains a running clock of CPU time used per process
- Uses a Virtual Clock (VC) to track usage
- Ensures average run-time is roughly equal for all tasks
- Better balance between CPU-bound and I/O-bound tasks

### Real-Time Schedulers

#### Real-Time Systems

Real-time systems have specific scheduling needs:

- Need responsiveness to I/O
- Must fulfill deadlines (hard or soft)
- Hard deadlines MUST be met, soft deadlines can occasionally be missed
- Deadlines may be in milliseconds, seconds, or hours

A Real-Time Operating System (RTOS):

- Completes system calls in deterministic time
- Schedules user tasks to meet deadlines
- Facilitates hard real-time requirements

#### Rate Monotonic Scheduling

A static-priority scheduling algorithm for real-time systems:

- Assigns higher priority to tasks with higher repetition rates
- Guaranteed schedule if utilization meets criteria:
  - $U = \sum_{i=1}^n \frac{C_e}{T_r} \leq n(2^{1/n} - 1)$
  - Where  $C_e$  is execution time and  $T_r$  is period
- Maximum guaranteed utilization converges to 69%
- Schedule determined at compile-time, not run-time

#### Earliest Deadline First (EDF)

A dynamic scheduling algorithm for real-time systems:

- Scheduler determines the task with the next deadline
- Task with the earliest deadline gets highest priority
- Schedule is achievable if utilization does not exceed 100%
- Can achieve full CPU utilization
- Schedule determined at run-time, not compile-time

### Linux Scheduling

#### Linux Scheduling Policies

Linux supports two general scheduling policies:

- **Real-Time Policy:** For time-critical tasks
  - SCHED\_DEADLINE: Earliest Deadline First + Constant Bandwidth Server
  - SCHED\_FIFO: First-In-First-Out real-time scheduling
  - SCHED\_RR: Round Robin real-time scheduling
- **Normal Policy:** For regular tasks
  - SCHED\_OTHER: Default Linux scheduling policy (Completely Fair Scheduler)
  - SCHED\_BATCH: For batch processing tasks
  - SCHED\_IDLE: For extremely low priority background tasks

#### Linux Priority System

Linux priority representation varies between kernel and user space:

- **Kernel space:** Priorities from high to low
  - RT: 0-99
  - Normal: 100-139
- **User space:** "nice" values from -20 (high) to +19 (low)
  - Maps to kernel priorities: nice + 20 = kernel priority - 100
  - Not used in RT policies

#### SCHED\_DEADLINE

Highest priority scheduler in Linux:

- Implements Earliest Deadline First + Constant Bandwidth Server
- Takes parameters: runtime, period, and deadline
- Tasks scheduled with SCHED\_DEADLINE cannot fork
- Tasks may yield CPU time when not needed

#### SCHED\_FIFO

First-In-First-Out real-time scheduler:

- Uses one queue per priority level (1-99)
- Higher priority than SCHED\_RR at same priority level
- Immediately preempts any Normal policy thread
- Runs to completion unless:
  - Preempted by a higher priority RT thread
  - Blocked by I/O call
  - Voluntarily yields the CPU

SCHED\_RR

Round Robin real-time scheduler:

- Similar to SCHED\_FIFO but with time quanta
- When quantum expires, thread moved to end of its priority queue
- Quantum size configurable (default 100ms)
- RT bandwidth limiting prevents RT tasks from monopolizing CPU

Completely Fair Scheduler (CFS)

Default scheduler in Linux (SCHED\_OTHER):

- Uses a red-black tree sorted by execution time (O(log N) operations)
- Tracks virtual runtime to achieve fairness
- Considers "nice" values to adjust CPU share
- Tasks can be grouped and scheduled together
- Aims to model an ideal, precise multi-tasking CPU"
- Time accounting managed according to configurable granularity

SCHED\_BATCH and SCHED\_IDLE

Low-priority schedulers:

- **SCHED\_BATCH:**
  - Same static priority as SCHED\_OTHER
  - Designed for batch-type, CPU-intensive tasks
  - Applies penalty due to CPU usage
  - SCHED\_OTHER has precedence over SCHED\_BATCH at same nice value
- **SCHED\_IDLE:**
  - Extremely low priority
  - Lower than static priority 0 and nice 19
  - Used for background tasks that should only run when system is idle

Multi-core Scheduling

Load Balancing

On multicore systems, Linux implements load balancing:

- Dynamic distribution of tasks across CPU cores
- Applied based on scheduling policy
- Balances competing concerns:
  - Moving tasks incurs management overhead
  - Moving tasks incurs cache penalty (lost cache advantage)

Cache Affinity

Cache affinity affects scheduling decisions:

- Task data may remain in CPU cache after context switch
- Rerunning on the same core avoids cache misses
- Linux considers estimated cache live-time when migrating tasks
- Default cache live-time: /proc/sys/kernel/sched\_migration\_cost\_ns

Analyzing Scheduling in Linux

Viewing scheduler information

- Check process priorities: ps -el (PRI and NI columns)
- View real-time processes: ps -eo pid,cls,pri,rtprio,nice,cmd | grep -E 'CLS|FIFO|RR'
- Check scheduler statistics: cat /proc/schedstat
- See current I/O scheduler: cat /sys/block/sda/queue/scheduler

Modifying process priorities

- Start process with nice value: nice -n [value] command
- Change nice value: renice [value] -p [pid]
- Set real-time priority: chrt -f [priority] command (SCHED\_FIFO)
- Set round-robin priority: chrt -r [priority] command (SCHED\_RR)

Analyzing schedule

- Trace scheduling events: trace-cmd record -e sched
- View trace results: trace-cmd report
- Check CPU affinity: taskset -p [pid]
- Set CPU affinity: taskset -c [cpu\_list] -p [pid]

FIFO and RR Scheduling Comparison

Given the following task list:

| Task | Arrival Time | Burst Time |
|------|--------------|------------|
| T1   | 0            | 10         |
| T2   | 3            | 6          |
| T3   | 7            | 1          |
| T4   | 8            | 3          |

FIFO Schedule:

T1 runs from 0-10, T2 runs from 10-16, T3 runs from 16-17, T4 runs from 17-20

Round Robin Schedule (quantum = 2):

T1(0-2), T1(2-4), T2(4-6), T1(6-8), T2(8-10), T3(10-11), T4(11-13), T1(13-15), T2(15-17), T1(17-19), T1(19-20)

RR provides better response time but potentially longer turnaround time due to context switches.

Rate Monotonic Scheduling

Given tasks with the following characteristics:

| Task | WCET (C) | Period (T) |
|------|----------|------------|
| T1   | 10       | 20         |
| T2   | 10       | 50         |
| T3   | 5        | 30         |

Analysis:

1. T1 has the highest priority (shortest period) 2. Utilization:  $U = \frac{10}{20} + \frac{10}{50} + \frac{5}{30} = 0.5 + 0.2 + 0.167 = 0.867$  3. Maximum utilization for 3 tasks:  $3(2^{1/3} - 1) \approx 0.779$  4. Since  $0.867 > 0.779$ , there is no guaranteed schedule  
However, a feasible schedule might still exist. Testing would be required to confirm.

## Resource Control

### Problem Statement and Motivation

#### Resource Competition

Operating systems must manage competition for resources:

- **Physical resources:** Devices, bus systems, memory, etc.
- **Virtual resources:** Timers, locks, etc.

Resource control requires:

- Visibility: Knowing what resources are available
- Access control: Determining who can use resources
- Usage monitoring: Tracking resource consumption
- Unified approach: Coordinated resource management

#### Nice Process Concept

Traditional Unix approach to resource control:

- Users can be "nice" by voluntarily releasing resources
- Modified by the `nice` command
- Limitations:
  - No global definition, only relative values
  - No strict and precise enforcement
  - Users can only modify niceness of their own processes
  - Limited applicability to specific resources

### Linux Control Groups (cgroups)

#### Control Groups (cgroups)

Linux kernel feature for organizing tasks into hierarchical groups:

- Allows processes to be organized and controlled collectively
- Provides strict enforcement of access and usage criteria
- Applies to sets of processes and all their future children
- Controls various types of resources (CPU, memory, I/O, etc.)

#### cgroups Terminology

Key concepts in cgroups:

- **cgroup:** Collection of processes bound to limits/parameters
- **Subsystem** (Controller): Kernel component related to a resource type
- **Hierarchy:** Arrangement of cgroups in a tree structure
- **Resource Controller:** Implements behavior for a specific resource type

Various controllers have been implemented:

- CPU controller: Limits CPU time
- Memory controller: Limits memory usage
- I/O controller: Controls disk I/O
- Network controller: Manages network bandwidth
- Devices controller: Controls device access

#### cgroups Implementation Approach

cgroups uses a filesystem-based approach:

- Communicates with Linux kernel via filesystem
- Virtual filesystem, stored in RAM
- Provides structured, standardized operations via I/O system calls
- Similar approach to `/proc` filesystem
- Implementation steps:
  - Create tmpfs mount
  - Create directory
  - Mount resource control interfaces as files in directory
  - Configure controls by editing files
  - Associate processes (PIDs) with control configuration

#### cgroups Hierarchy

cgroups are organized hierarchically:

- Created by making subdirectories in the cgroup filesystem
- Limits defined at each level of the hierarchy
- Limits apply throughout the sub-hierarchy
- Descendant cgroups cannot exceed limits of ancestor cgroups

### cgroups Versions

#### cgroups v1 vs v2

Linux has two implementations of cgroups:

- **cgroups v1:**
  - Initial release in Linux 2.6.24
  - Rapid adoption but uncoordinated development
  - Some inconsistencies between controllers
- **cgroups v2:**
  - Added in Linux 4.5
  - Intended to replace v1 (but v1 continues to exist for compatibility)
  - Currently implements a subset of v1 controllers
  - Both versions can coexist, but a controller can't be used in both simultaneously

### cgroups Version 1

#### cgroups v1 Concepts

cgroups v1 has two approaches to resource control:

- **Individual Resource Control:**
  - Each controller mounted against a separate filesystem
  - One controller associated with one hierarchy
- **Collective Resource Control:**
  - Multiple controllers co-mounted against a single filesystem
  - Co-mounted controllers manage the same hierarchy

Each cgroup is represented by a directory:

- Child cgroups represented as subdirectories
- Configuration files under each directory
- Files reflect resource limits and properties

#### Tasks vs Processes in cgroups v1

cgroups v1 distinguishes between tasks and processes:

- Process can consist of multiple threads (tasks)
- cgroups v1 allows independent manipulation of thread cgroup membership
- This can cause problems for controllers like memory (all threads share an address space)
- This ability has been limited in cgroups v2

#### Mounting cgroups v1

Mounting cgroups controllers:

```
1 # Create a tmpfs mount for cgroups
2 sudo mount -t tmpfs -o size=10M tmpfs /sys/fs/cgroup
3
4 # Mount a single controller (individual resource
   control)
5 sudo mount -t cgroup -o cpu none /sys/fs/cgroup/cpu
6
7 # Co-mount multiple controllers (collective resource
   control)
8 sudo mount -t cgroup -o cpu,cpuacct none
   /sys/fs/cgroup/cpu,cpuacct
```

It's not possible to mount the same controller against multiple hierarchies.

#### Creating and Managing cgroups v1

Basic cgroup management:

```
1 # Create a new cgroup
2 mkdir /sys/fs/cgroup/cpu/cg1
3
4 # Add the current process to the cgroup
5 echo $$ > /sys/fs/cgroup/cpu/cg1/cgroup.procs
6
7 # Add a specific process to the cgroup
8 echo <PID> > /sys/fs/cgroup/cpu/cg1/cgroup.procs
9
10 # Remove a cgroup (must be empty of processes and
    child cgroups)
11 rmdir /sys/fs/cgroup/cpu/cg1
```

When adding a process to a cgroup, all threads in the process are moved together.

### cgroups Controllers

#### CPU Controller

The CPU controller manages CPU usage:

- Controls upper and lower limits of CPU shares
- Lower limit guarantees minimum CPU time when system is busy
- Upper limit restricts CPU time in each scheduling period
- Does not limit CPU usage if CPUs are not busy

#### Memory Controller

The Memory controller manages memory usage:

- Reports and limits process memory, kernel memory, and swap usage
- Sets hard and soft limits on memory consumption
- Can enforce OOM (Out Of Memory) killing within cgroups

#### Blkio Controller

The Blkio controller manages block device I/O:

- Limits access to block devices through throttling
- Sets upper I/O rate limits on devices
- Implements proportional-weight time-based division of disk I/O
- Can limit both read and write operations

## Cpuset Controller

The Cpuset controller manages CPU assignment:

- Binds processes to specific CPUs
- Allows process isolation on multi-core systems
- Can be used to improve cache utilization

## CPU Controller Example

Limiting CPU usage for a group of processes:

```
1 # Create a cgroup for CPU control
2 mkdir /sys/fs/cgroup/cpu/limited_group
3
4 # Set maximum CPU usage to 10% (100ms in a 1000ms
5   period)
6 echo 100000 >
7   /sys/fs/cgroup/cpu/limited_group/cpu.cfs_period_us
8 echo 10000 >
9   /sys/fs/cgroup/cpu/limited_group/cpu.cfs_quota_us
10
11 # Add a process to the cgroup
12 echo $PID >
13   /sys/fs/cgroup/cpu/limited_group/cgroup.procs
14
15 # Run a CPU-intensive task and observe it being limited
16 # Even on an idle system, it won't exceed 10% of one
17   CPU
```

## Device Controller Example

Restricting access to a device:

```
1 # Create a cgroup for device control
2 mkdir /sys/fs/cgroup/devices/group0
3
4 # By default, full permissions exist
5 cat /sys/fs/cgroup/devices/group0/devices.list
6 # Output: a *:rwm (all devices, all permissions)
7
8 # Deny access to /dev/null (major 1, minor 3)
9 echo 'c 1:3 rwm' >
10   /sys/fs/cgroup/devices/group0/devices.deny
11
12 # Add the current shell to the group
13 echo $$ > /sys/fs/cgroup/devices/group0/tasks
14
15 # Try to use /dev/null - it will fail
16 echo "test" > /dev/null
17 # Output: bash: /dev/null: Operation not permitted
```

## cgroups Version 2

### cgroups v2 Overview

cgroups v2 has several key differences from v1:

- All mounted controllers reside in a single unified hierarchy
- Simplified controller set
  - io (replaces blkio)
  - memory
  - pids
  - perf\_event
  - rdma
  - cpu
  - freezer
- Supports delegation (non-privileged users can manage subtrees)
- Supports thread mode (thread-level control)

### Working with cgroups

#### Exploring cgroups

- List available subsystems: `cat /proc/cgroups`
- View cgroup hierarchy: `tree /sys/fs/cgroup`
- Check cgroups created by systemd: `systemd-cgls`
- Monitor cgroup usage: `systemd-cgtop`

#### Creating and managing cgroups

- Create a cgroup: `mkdir /sys/fs/cgroup/[controller]/[name]`
- Add process to cgroup: `echo [PID] > /sys/fs/cgroup/[controller]/[name]/cgroup.procs`
- Set CPU limit: `echo [value] > /sys/fs/cgroup/cpu/[name]/cpu.cfs_quota_us`
- Set memory limit: `echo [value] > /sys/fs/cgroup/memory/[name]/memory.limit_in_bytes`
- Remove cgroup: `rmdir /sys/fs/cgroup/[controller]/[name]`

#### Working with systemd cgroups

- Create transient cgroup: `systemd-run -unit=[name] -slice=[slice] [command]`
- Set resource limits: `systemctl set-property [unit] [property]=[value]`
- Example: `systemctl set-property user.slice CPUQuota=20%`

## Using systemd Resource Control

Controlling resource limits with systemd:

```
1 # View current system slices
2 systemctl list-units --type=slice
3
4 # Check user slice settings
5 systemctl show user.slice
6
7 # Limit CPU usage for all user processes to 20%
8 sudo systemctl set-property user.slice CPUQuota=20%
9
10 # Create a resource-limited service
11 cat << EOF >
12   /etc/systemd/system/limited-service.service
13 [Unit]
14   Description=Resource Limited Service
15
16 [Service]
17   ExecStart=/usr/bin/sha1sum /dev/zero
18   CPUQuota=10%
19   MemoryLimit=100M
20
21 [Install]
22   WantedBy=multi-user.target
23 EOF
24
25 # Reload systemd and start the service
26 sudo systemctl daemon-reload
27 sudo systemctl start limited-service
28
29 # Monitor resource usage
30 systemd-cgtop
```

## Creating a Custom CPU Control

Creating a CPU affinity control from scratch:

```
1 # Create a tmpfs mount
2 sudo mkdir -p /mnt/cgroups
3 sudo mount -t tmpfs none /mnt/cgroups
4
5 # Mount cgroup controller
6 sudo mkdir -p /mnt/cgroups/cpuset
7 sudo mount -t cgroup -o cpuset none /mnt/cgroups/cpuset
8
9 # Create a cgroup
10 sudo mkdir /mnt/cgroups/cpuset/group1
11
12 # Configure required memory settings first
13 echo 0 > /mnt/cgroups/cpuset/group1/cpuset.mems
14
15 # Restrict to CPUs 0 and 1 only (assuming 4 CPUs)
16 echo "0-1" > /mnt/cgroups/cpuset/group1/cpuset.cpus
17
18 # Run CPU-intensive processes
19 for i in {1..4}; do
20     # Create processes
21     shasum /dev/zero &
22     # Get PID and add to cgroup
23     PID=$!
24     echo $PID > /mnt/cgroups/cpuset/group1/cgroup.procs
25     echo "Added process $PID to CPU-restricted group"
26 done
27
28 # Check CPU assignment with taskset
29 for pid in $(cat
30     /mnt/cgroups/cpuset/group1/cgroup.procs); do
31     taskset -p $pid
```



## Memory Management

### Memory Management Basics

#### Memory Management Introduction

Memory is a critical system resource managed by the operating system:

- The Memory Manager controls main memory allocation and usage
- Secondary memory can be used as a buffer zone (swap)

Memory is organized in a hierarchy:

- Fast cache in 1-3 layers (L1, L2, L3)
- Main memory (RAM) - slower but larger
- Secondary memory (hard disks, SSDs) - for programs and files
- Tertiary memory (backup storage, tapes)

#### Memory Management Tasks

Operating systems handle several memory management tasks:

- Determining how much memory processes require
- Deciding where in memory processes are located (position of residency)
- Managing how long processes remain in memory (length of residency)
- Subdividing memory for co-existence of multiple processes
- Handling memory fragmentation

### Memory Allocation Approaches

#### Memory Division and Fragmentation

Two basic approaches to memory division:

- **Static memory division:**
  - Memory divided into fixed-size segments
  - Problem: Internal fragmentation (wasted space within allocated segments)
- **Dynamic memory division:**
  - Memory divided according to process needs
  - Problem: External fragmentation (free space becomes fragmented)
  - Solution: Compaction (expensive operation)

#### Buddy Algorithm

The Buddy Algorithm addresses fragmentation:

- Memory divided into blocks of power-of-2 sizes
- When a request arrives, the system:
  - Finds the smallest block that fits the request
  - If no suitable block exists, splits a larger block into two "buddies"
  - Allocates one buddy and keeps the other free
- When a block is freed, the system:
  - Checks if its buddy is also free
  - If both are free, merges them into a larger block
  - Continues merging recursively if possible
- Simpler to implement than other dynamic allocation schemes
- Still experiences some internal fragmentation

#### Free Space Management

Finding free space during allocation requires efficient algorithms:

- **Bitmap approach:**
  - Space-efficient representation
  - One bit per allocation unit
  - Fast free-space finding
- **Linked list approach:**
  - List of free blocks
  - Supports various placement algorithms (first/next/best fit)

#### Swapping

When main memory becomes scarce, swapping is used:

- Secondary memory serves as temporary storage for processes
- When processes are suspended or exit, their memory is freed
- When processes restart, they are reloaded into memory
- Allows more processes to run concurrently than physical memory would permit

### Virtual Memory

#### Virtual Memory Concept

Virtual memory leverages program behavior characteristics:

- Programs exhibit spatial locality (tend to use a limited area of code at any time)
- Entire processes don't need to be fully resident in memory
- Non-required code/data can be swapped out when not immediately needed
- Enables more processes to run concurrently
- Must be transparent to programmer/process

#### Paging

Paging is the mechanism that enables virtual memory:

- Process memory is divided into fixed-size pages
- Physical memory is divided into frames of the same size
- Pages are loaded into frames as needed
- Memory Management Unit (MMU) manages mapping between pages and frames
- Typically uses on-demand paging (lazy loading)
  - Only loads pages when they are accessed
  - Sometimes prefetches additional pages based on locality
- Process then has a set of resident pages in memory
  - Resident set: All process pages in memory
  - Working set: Pages currently being used

#### Address Translation

Virtual memory requires translation between logical and physical addresses:

- Logical address consists of page number and page offset
- Page number used to look up frame number in page table
- Physical address formed by combining frame number with page offset
- Translation process:
  - Extract page number and offset from logical address
  - Use page number to index page table
  - Retrieve frame number from page table
  - Combine frame number with offset to form physical address

#### Page Table

Page tables maintain the mapping between pages and frames:

- Each entry contains a frame number
- Entries also include status bits:
  - Valid bit: Indicates if page holds valid data
  - Present bit: Indicates if page is in memory
  - Modified bit (dirty): Indicates if page has been written to
  - Referenced bit: Indicates recent usage
  - Protection bits: Control read/write/execute permissions
- Page tables are stored in main memory
- Can be very large for large address spaces

#### Translation Lookaside Buffer (TLB)

The TLB addresses the performance overhead of page table lookups:

- Cache for recently accessed page table entries
- Small (typically 64 entries)
- Uses content-addressable memory (CAM) for fast lookups
- Memory access process with TLB:
  - Check TLB for page number
  - If found (TLB hit), use frame number directly
  - If not found (TLB miss), search page table
  - If not in page table, trigger page fault
  - Add entry to TLB for future accesses

#### Page Replacement

When memory is full and a new page is needed, page replacement occurs:

- System must decide which resident page to evict
- Can use global strategy (any process's pages) or local strategy (only faulting process's pages)
- Common page replacement algorithms:
  - Optimal: Replace page used furthest in future (theoretical only)
  - Least Recently Used (LRU): Replace page unused for longest time
  - First-In-First-Out (FIFO): Replace oldest page

### Linux Memory Management

#### Linux Page Table Organization

Linux uses a hierarchical page table structure:

- Multi-level page directory to reduce size and improve lookup speed
- Typically 4-level structure:
  - Page Global Directory (PGD)
  - Page Upper Directory (PUD)
  - Page Middle Directory (PMD)
  - Page Table Entry (PTE)
- Allows efficient handling of sparse address spaces
- Only allocates page tables for used parts of address space

#### Huge Pages

Linux supports huge pages to improve performance:

- Standard page size is 4KB
- Huge pages can be 2MB or 1GB (architecture-dependent)
- Advantages:
  - Reduces TLB pressure (fewer entries needed to cover same memory)
  - Improves performance for memory-intensive applications
  - More efficient for large memory allocations
- Uses higher-level page table entries (PUD/PMD)
- Requires explicit configuration

Memory Zones

- Linux divides physical memory into zones to handle hardware limitations:
- **ZONE\_DMA**: Memory addressable by DMA controllers (typically below 16MB)
  - **ZONE\_NORMAL**: Regularly mapped memory in kernel space
  - **ZONE\_HIGHMEM**: Memory beyond what the kernel can directly address
  - Zones are managed separately to accommodate different hardware constraints
  - Each zone has its own free lists and allocation policies

Buddy Allocator

- Linux uses the buddy system for frame allocation:
- Fundamental allocation unit is the page frame
  - Maintains lists of free blocks of various sizes (powers of 2)
  - When a process requests memory:
    - System finds the smallest block size that fits the request
    - If necessary, splits larger blocks into "buddies"
    - Allocates memory from appropriate free list
  - When memory is freed:
    - System checks if buddy is also free
    - If so, merges buddies to form larger block
    - Continues merging recursively if possible
  - Maintains free lists up to MAX\_ORDER-1 (typically 10, so up to 512 contiguous pages)

Slab Allocator

- Linux uses the slab allocator for kernel objects:
- Kernel often requires small allocations for data structures
  - Pages (4KB) are too large for many kernel objects
  - Slab allocator:
    - Gets pages from buddy allocator
    - Divides them into smaller objects of specific types
    - Maintains caches of frequently used object types
    - Reuses recently freed objects (helps prevent fragmentation)
    - Preserves object state between uses
  - Improves memory utilization and allocation speed
  - Minimizes internal fragmentation

Memory Compaction

- Linux performs memory compaction to address fragmentation:
- Problem: Buddy allocator may not find large contiguous blocks
  - Solution: kcompactd daemon performs compaction
  - Process:
    - Balances memory zones by swapping out non-working-set pages
    - Moves movable pages toward the top of the memory zone
    - Leaves bottom of memory free for new allocations
    - Creates larger contiguous free blocks
  - Performed on-demand or periodically
  - Enables allocation of huge pages and other large memory blocks

Shared Libraries

- Linux uses shared libraries to reduce memory usage:
- Multiple processes can use the same library code
  - Libraries compiled with -fPIC (Position Independent Code)
  - Dynamically linked with processes using ld.so
  - Read-only code pages memory-mapped into processes
  - Benefits:
    - Reduces memory footprint
    - Shared code/text pages only loaded once
    - Only data pages need to be process-specific
  - Challenge: Version compatibility ("DLL hell")

Page Reclamation

- Linux uses page reclamation to recover memory:
- Working set: Pages actively in use by processes
  - Resident set: All pages in memory
  - A page is in the working set if:
    - Accessed via process address space
    - Accessed via system call
    - Accessed via device driver
  - Linux identifies non-working set pages using a bitmap
  - Pages marked idle can be reclaimed when memory is needed

Least Recently Used in Linux

- Linux implements a two-stage LRU algorithm:
- Maintains two lists of page frames:
    - Active list: Recently accessed pages
    - Inactive list: Less recently accessed pages
  - Pages move between lists based on access patterns
  - Inactive pages are candidates for reclamation
  - Recently accessed inactive pages can be promoted to active list
  - Linux uses a global strategy for page reclamation
  - Recent development: Multi-Generational LRU (MGLRU)
    - Assigns generation numbers to page frames based on recent access
    - Older generations are reclaimed first
    - Improves performance and responsiveness

Out-of-Memory (OOM) Killer

- Linux has an OOM Killer to handle critical memory shortages:
- Activates when system is critically low on memory
  - Linux tends to be optimistic in memory allocation
    - Processes typically request more memory than needed
    - System may over-allocate (more than physical memory)
  - OOM Killer selects processes to terminate based on heuristics
    - Considers memory usage, runtime, nice value, etc.
    - Each process has an oom\_score in /proc/\$PID/oom\_score
    - Can be adjusted through oom\_score\_adj
  - Prioritizes system stability over individual process survival

Memory Management Analysis in Linux

Basic memory information

- Display system memory usage: free -h
- View memory details: cat /proc/meminfo
- Check process memory usage: ps -eo pid,ppid,cmd,vsz,rss
- Interactive memory monitor: top or htop

Process memory analysis

- Check process memory maps: cat /proc/\$PID/maps
- View process memory status: cat /proc/\$PID/status
- Analyze memory usage in detail: pmap \$PID
- Track memory over time: smem

Memory page information

- Get page size: getconf PAGE\_SIZE
- Check huge pages: cat /proc/meminfo | grep Huge
- View page stats: cat /proc/pagetypeinfo
- Check page faults: ps -o min\_flt,maj\_flt \$PID

Memory limits and control

- Set memory limits: ulimit -m [size]
- Control cgroup memory: echo [value] > /sys/fs/cgroup/memory/[group]/memory.limit\_in\_bytes
- Check swappiness: cat /proc/sys/vm/swappiness
- Adjust swappiness: sysctl vm.swappiness=[value]

Analyzing Process Memory

Check memory usage details for a process:

```
1 # Get PID of a process
2 PID=$(pidof firefox)
3
4 # Check virtual and resident memory size
5 ps -o pid,comm,vsz,rss -p $PID
6
7 # Analyze memory map segments
8 cat /proc/$PID/maps | head -10
9
10 # View detailed memory status
11 grep -E 'VmSize|VmRSS|VmData|VmStk|VmExe'
   /proc/$PID/status
12
13 # Map process address space in detail
14 pmap -x $PID | head -20
15
16 # Check page faults
17 ps -o min_flt,maj_flt -p $PID
```

Explanation of memory terms:

- **VSZ (Virtual Size)**: Total virtual memory allocated to process
- **RSS (Resident Set Size)**: Actual physical memory used
- **VmData**: Size of data segment
- **VmStk**: Size of stack
- **VmExe**: Size of text segment
- **min\_flt**: Minor page faults (page in memory but not in process's page table)
- **maj\_flt**: Major page faults (page had to be loaded from disk)

## Working with Page Size and Memory Allocation

Analyzing page size and memory allocation:

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <unistd.h>
4 #include <sys/mman.h>
5
6 int main() {
7     // Get process ID
8     pid_t pid = getpid();
9     printf("Process ID: %d\n", pid);
10
11     // Get page size
12     int page_size = getpagesize();
13     printf("Page size: %d bytes\n", page_size);
14
15     // Allocate memory for 10 pages
16     size_t size = 10 * page_size;
17     char *buffer = malloc(size);
18     printf("Allocated %zu bytes (%zu pages)\n",
19           size, size / page_size);
20
21     // Check page faults before access
22     printf("Check page faults with: ps -o
23           min_flt,maj_flt %d\n", pid);
24     printf("Press Enter to continue...\n");
25     getchar();
26
27     // Access the memory (causes page faults)
28     for (int i = 0; i < size; i += page_size) {
29         buffer[i] = 1; // Touch one byte per page
30     }
31
32     // Check page faults after access
33     printf("Check page faults again with: ps -o
34           min_flt,maj_flt %d\n", pid);
35     printf("Press Enter to continue...\n");
36     getchar();
37
38     // Allocate page-aligned memory
39     void *aligned_buf = NULL;
40     int result = posix_memalign(&aligned_buf,
41                               page_size, size);
42     if (result == 0) {
43         printf("Allocated %zu bytes aligned on page
44               boundary\n", size);
45     }
46
47     // Free memory
48     free(buffer);
49     free(aligned_buf);
50
51     return 0;
52 }
```

This example demonstrates:

- Getting the system page size
- Allocating memory
- Observing page faults due to lazy allocation
- Creating page-aligned memory allocations

## Input/Output Basics

## I/O Challenges

I/O management presents unique challenges for operating systems:

- Huge variety of I/O devices (keyboards, mice, drives, sensors, etc.)
- Diverse interfaces (ATA, SATA, USB, PCI, etc.)
- Wide range of speeds and capacities
- Different data transfer characteristics
- Need for a unified approach to device interaction

The operating system must abstract this heterogeneity by providing consistent interfaces for applications to access diverse devices.

## Device Categories

From a user perspective, devices fall into two main categories:

- **Block Devices:**
  - Operate on fixed-size blocks of data
  - Support random access to any block
  - Example: Hard drives, SSDs
  - Block operations are independent from each other
  - Block size defined when formatting (logical)
  - Sector size is the physical unit on the device
- **Character Devices:**
  - Operate on streams of characters
  - Generally sequential access
  - Example: Keyboards, mice, printers
  - Characters are interpretations of bit patterns according to specifications (ASCII, Unicode)
  - Character devices ultimately operate on bit-level too

## I/O Hardware Architecture

## I/O Hardware Components

Key components in I/O hardware architecture:

- **I/O Controller:** Electronic interface to the device
  - Typically one per device category (SCSI, IDE, USB, Ethernet)
  - Contains registers that control device operation
  - Maintains buffers for read/write operations
  - Communicates with CPU and main memory
- **I/O Ports:** Addresses that point to controller registers
- **Buffers:** Memory areas for data transfer
- **Bus System:** Connects CPU, memory, and I/O devices

## I/O Address Space

X86 architecture has two approaches for I/O address allocation:

- **Port Mapped I/O (PMIO):**
  - Peripheral I/O registers assigned to a dedicated address range
  - Distinct from system memory
  - Uses dedicated I/O instructions (IN and OUT)
  - Limited to 64K ports (16-bit addressing)
  - Common for older peripherals (e.g., ISA cards)
  - Listed in `/proc/ioports`
- **Memory Mapped I/O (MMIO):**
  - Peripheral registers mapped into main memory address space
  - Uses regular memory instructions (e.g., MOV)
  - Default in modern systems (e.g., PCIe, PCI)
  - Listed in `/proc/iomem`

## Direct Memory Access (DMA)

DMA improves I/O performance:

- Allows devices to transfer data directly to/from memory
- Bypasses CPU for bulk data transfer
- CPU sets up transfer parameters:
  - Memory address
  - Count of bytes to transfer
  - Direction (read/write)
  - Device-specific parameters
- DMA controller handles the actual transfer
- Interrupts CPU when transfer completes
- Reduces CPU overhead for data-intensive I/O operations

## I/O Principles

## I/O Access Methods

Two primary methods for CPU to interact with I/O devices:

- **Programmed I/O:**
  - CPU does all the work
  - CPU executes instructions to transfer data
  - CPU repeatedly checks device status (polling/busy-waiting)
  - Synchronous operation (CPU waits for I/O completion)
  - Simple but inefficient for slow devices
- **Interrupt-driven I/O:**
  - Devices signal CPU when they need attention
  - CPU initiates operation then continues other work
  - Device interrupts when operation completes
  - Asynchronous operation
  - More efficient, especially for slow devices

## Interrupt Handling

Interrupts are fundamental to efficient I/O:

- **Interrupt:** Event that defers the normal flow of CPU execution
- When an interrupt occurs:
  - Current execution is suspended
  - CPU state is saved
  - Special routine (interrupt handler) is executed
  - After completion, previous execution resumes
- **Interrupt types:**
  - **Synchronous:** Generated by executing an instruction (e.g., divide by zero)
  - **Asynchronous:** Generated by external events (e.g., I/O completion)
- **Interrupt classifications:**
  - **Maskable:** Can be ignored by setting interrupt mask
  - **Non-maskable:** Cannot be ignored (critical events)

## Interrupt Flow

Hardware interrupt flow:

- Device raises interrupt on its IRQ line
- Programmable Interrupt Controller (PIC) converts IRQ to vector number
- PIC signals CPU via INTR pin
- CPU acknowledges interrupt
- CPU executes the appropriate interrupt handler
- In Linux, interrupt handling occurs in three phases:
  - Critical: Minimal processing, acknowledge interrupt
  - Immediate: Essential processing, can't be deferred
  - Deferred: Non-critical processing, scheduled for later

## I/O Access Patterns

Different access patterns affect I/O performance:

- **Exclusive vs. Shared Access:**
  - Exclusive: Device dedicated to one process
  - Shared: Multiple processes access same device
  - Shared access requires scheduling (e.g., disk I/O scheduling)
- **Sequential vs. Random Access:**
  - Sequential: Data accessed in order (e.g., tape)
  - Random: Data accessed in any order (e.g., disk)
- **Blocking vs. Non-Blocking:**
  - Blocking: Process waits until I/O completes
  - Non-Blocking: Process continues, checks completion later
- **Synchronous vs. Asynchronous:**
  - Synchronous: Process execution synchronized with I/O completion
  - Asynchronous: Process continues, notified of I/O completion

## Buffered vs. Direct I/O

Buffering improves I/O performance:

- **Buffered I/O:**
  - Data passes through intermediate buffer
  - Decouples data access from data generation
  - Handles different speeds between devices
  - Enables rate control
  - Allows data manipulation before final transfer
  - Supports data verification
- **Direct I/O:**
  - Data transferred directly to/from device
  - Bypasses system caches
  - Reduces memory usage and CPU overhead
  - Useful for applications with their own caching
  - May be slower for some workloads

## Error Handling

I/O systems must handle errors effectively:

- Errors can occur at various levels:
  - Bit-level, byte-level, packet-level, block-level
  - Hardware vs. software detection
  - User space vs. kernel space handling
- Error handling approaches:
  - Error Detection: Identify errors (e.g., checksums, parity)
  - Error Correction: Fix errors without retransmission
  - Error Recovery: Return to consistent state after error
- Trade-offs between performance and reliability

## Linux I/O Subsystem

## Linux Device Model

The Linux Device Model maintains the state and structure of the system:

- Maintains information about devices, drivers, buses, etc.
- Key entities:
  - **Device:** Physical device attached to a bus
  - **Driver:** Software entity that operates devices
  - **Bus:** Device to which other devices can be attached
  - **Class:** Type of device with similar behavior
  - **Subsystem:** View on the system structure
- Represented in user space via `sysfs` (mounted at `/sys`)

## sysfs

sysfs is a virtual filesystem that exposes the Linux Device Model:

- Located at /sys
- Key directories:
  - /sys/block: Block devices
  - /sys/bus: Bus types
  - /sys/class: Device classes
  - /sys/devices: Hierarchical device structure
  - /sys/firmware: Firmware information
  - /sys/fs: Filesystem information
  - /sys/kernel: Kernel status
  - /sys/module: Loaded modules
  - /sys/power: Power management
- Contains attributes in files for configuration and status

## udev

udev is the userspace device manager for Linux:

- Part of systemd
- systemd-udev listens to kernel events
- Executes rules based on device information from sysfs
- Creates or removes device nodes in /dev
- Provides consistent device naming
- Enables automatic device setup
- Supports user-defined rules

## /dev Directory

The /dev directory contains special files representing devices:

- Each file represents an I/O device
- Allows standard file operations (open, read, write, close) on devices
- When accessed, kernel routes operations to appropriate device drivers
- Types of device files:
  - Character device files: For character devices
  - Block device files: For block devices
- Naming conventions:
  - /dev/sdX: SCSI/SATA disk devices
  - /dev/ttyX: Terminal devices
  - /dev/nullX: Special device files

## Device Access in Linux

Linux provides a unified interface for device access:

- Applications use standard file operations:
  - open(): Open device
  - read(): Read from device
  - write(): Write to device
  - ioctl(): Device-specific operations
  - close(): Close device
- Kernel translates these operations to device-specific commands
- Virtual File System (VFS) provides abstraction layer
- Device drivers implement the specific operations

## Working with I/O in Linux

### Exploring device information

- List I/O ports: cat /proc/ioports
- List I/O memory: cat /proc/iomem
- View interrupts: cat /proc/interrupts
- List block devices: lsblk
- Show hardware: lshw
- Display PCI devices: lspci
- List USB devices: lsusb

### Working with devices

- Check device information: udevadm info -query=all -name=/dev/sda
- Monitor device events: udevadm monitor
- Show device properties: udevadm info -attribute-walk -name=/dev/sda
- Check I/O performance: iostat -x
- Monitor I/O activity: iotop

### Device file operations

- Create device file: mknod /dev/example c 1 3
- Read from device: cat /dev/input/mouse0 | hexdump
- Write to device: echo "test" > /dev/tty1
- Control device: ioctl system call in C programs

## I/O Performance Testing

Testing disk write performance with different options:

```
1 # Basic write test (with caching)
2 dd if=/dev/zero of=speedtest bs=10M count=100
3 rm speedtest
4
5 # Write test with synchronous I/O (forces data to disk)
6 dd if=/dev/zero of=speedtest bs=10M count=100
7   conv=fdatasync
8 rm speedtest
9
10 # Write test with direct I/O (bypasses the page cache)
11 dd if=/dev/zero of=speedtest bs=10M count=100
12   oflag=direct
13 rm speedtest
14
15 # Monitor disk I/O activity during test
16 iostat -x 1
17
18 # Check disk utilization statistics
19 iostat -p sda
20
21 # Explanation:
22 # - Standard dd shows high performance but may not be
23   on disk yet
24 # - fdatasync ensures data is physically written to
25   disk
26 # - direct bypasses the OS cache, showing raw device
27   performance
```

## Device Information Analysis

Exploring device information in Linux:

```
1 # List block devices with details
2 lsblk -f
3
4 # Get detailed information about a specific device
5 udevadm info --query=all --name=/dev/sda1
6
7 # Examine sysfs entries for the device
8 ls -l /sys/block/sda/sda1/
9
10 # Check device attributes
11 cat /sys/block/sda/queue/scheduler
12 cat /sys/block/sda/queue/read_ahead_kb
13
14 # Get PCI information for a disk controller
15 lspci | grep -i sata
16
17 # Check all properties of a device
18 udevadm info --attribute-walk --name=/dev/sda1 | less
19
20 # Monitor udev events when plugging in a USB device
21 udevadm monitor --property
22 # (Now plug in a USB device to see events)
```



## Building and Using a Custom Linux Kernel

### Custom Kernel Motivation

There are various reasons to build a custom kernel:

- Create a minimalist kernel (disable unused features, load needed ones as modules)
- Add custom OS features for specific requirements
- Support special hardware that may not be in the standard kernel
- Optimize for specific workloads or hardware platforms
- Learn about kernel internals and development processes

### Linux Kernel Development Process

Linux follows a time-based release process:

- New major kernel releases every 2-3 months (e.g., 5.11, 5.12)
- Development cycle phases:
  - Merge window (first two weeks): New features merged into mainline
  - Release candidates (weekly): Bug fixes only, no new features
  - Final release: After several release candidates
- Long-term support (LTS) kernels receive updates for extended periods
- Community development model with maintainers for various subsystems

### Choosing a Kernel Version

When building a custom kernel, version selection matters:

- Mainline: Latest development version (might be unstable)
- Stable: Recent release with proven stability
- Long-term: Longer support period (good for production systems)
- Distribution-specific: Modified by Linux distributions

Version numbers indicate:

- First number: Major version (rarely changes)
- Second number: Minor version (major features)
- Third number: Patch level (bug fixes and minor improvements)

## Building a Custom Linux Kernel

### Getting the source code

- Download from kernel.org: `wget https://cdn.kernel.org/pub/linux/kernel/v5.x/linux-5.xx.tar.xz`
- Extract: `tar xf linux-5.xx.tar.xz`
- Change directory: `cd linux-5.xx`

### Installing build dependencies

- On Debian/Ubuntu: `sudo apt-get install build-essential gcc bc bison flex libssl-dev libncurses5-dev libelf-dev`

### Configuring the kernel

- Start from existing config (recommended):
  - Current running kernel: `cp /boot/config-$(uname -r) ./.config`
  - Distribution default: `make defconfig`
- Update config for new options: `make oldconfig`
- Interactive configuration tools:
  - Text-based menu: `make menuconfig`
  - GUI-based: `make xconfig` or `make gconfig`
- Important configuration options:
  - Custom version name: `CONFIG_LOCALVERSION`
  - Module support: `CONFIG_MODULES`
  - CPU and architecture options
  - Device drivers
  - File systems

### Building the kernel

- Compile: `make -j$(nproc)`
- Build modules: `make modules`
- Create Debian/Ubuntu packages: `make -j$(nproc) deb-pkg`
- For RPM-based systems: `make -j$(nproc) rpm-pkg`

### Installing the kernel

- Debian packages: `sudo dpkg -i ../linux-*.deb`
- RPM packages: `sudo rpm -ivh ../linux-*.rpm`
- Manual installation:
  - Install modules: `sudo make modules_install`
  - Install kernel: `sudo make install`
- Update bootloader: `sudo update-grub` (GRUB)

### Booting the new kernel

- Reboot: `sudo reboot`
- Select the new kernel at the bootloader menu
- Verify running kernel: `uname -a`

## Linux Kernel Modules

### Kernel Modules Concept

Kernel modules extend kernel functionality without rebuilding:

- Loadable code that can be added to or removed from a running kernel
- Allow dynamic extension of kernel capabilities
- Provide device drivers, filesystem drivers, system calls, etc.
- Reduce the size of the base kernel
- Enable support for hardware that is hot-pluggable
- Support for special features used in specific applications

## Module Structure

A kernel module follows a specific structure:

- Must include necessary kernel headers
- Initialization function (`init_module()` or custom named)
  - Called when the module is loaded
  - Sets up resources, registers with kernel subsystems
  - Returns success (0) or error code
- Cleanup function (`cleanup_module()` or custom named)
  - Called when the module is unloaded
  - Releases resources, unregisters from kernel subsystems
- Uses `MODULE_LICENSE()` macro to specify license
- Can include other macros for author, description, etc.

### Hello World Kernel Module

Minimal kernel module example:

```
1 #include <linux/module.h>    /* Needed by all modules */
2 #include <linux/kernel.h>    /* Needed for KERN_INFO */
3 #include <linux/init.h>      /* Needed for macros */
4
5 MODULE_LICENSE("GPL");
6 MODULE_AUTHOR("Your Name");
7 MODULE_DESCRIPTION("A simple Hello World module");
8
9 static int __init hello_init(void)
10 {
11     printk(KERN_INFO "Hello, World!\n");
12     return 0;    /* Success */
13 }
14
15 static void __exit hello_exit(void)
16 {
17     printk(KERN_INFO "Goodbye, World!\n");
18 }
19
20 module_init(hello_init);
21 module_exit(hello_exit);
```

Key components:

- `module_init` and `module_exit` macros register functions
- `printk` is the kernel's version of `printf`
- `KERN_INFO` sets the message priority level
- `MODULE_LICENSE` declares the license (important for symbol exports)

### Module Building Process

Building a kernel module requires:

- Pre-built kernel with module support
- Module source code
- Makefile defining the build process
- Build tools and headers

The build process:

- `Kbuild` system builds `<module_name>.o` from source
- Links to create `<module_name>.ko` (kernel object)
- Kernel modules must be built against the same kernel version they will run on
- Distribution-specific kernel headers needed for module compatibility



## Module Makefile

Example Makefile for a kernel module:

```
1 # Define the module name
2 obj-m := hello.o
3
4 # For standalone module build
5 all:
6     make -C /lib/modules/$(shell uname -r)/build M=$(PWD) modules
7
8 # Clean up
9 clean:
10     make -C /lib/modules/$(shell uname -r)/build M=$(PWD) clean
```

For multiple source files:

```
1 # Module with multiple source files
2 hello-y := hello_main.o hello_func.o
3 obj-m := hello.o
4
5 # For standalone module build
6 all:
7     make -C /lib/modules/$(shell uname -r)/build M=$(PWD) modules
8
9 # Clean up
10 clean:
11     make -C /lib/modules/$(shell uname -r)/build M=$(PWD) clean
```

## Working with Kernel Modules

### Building a module

- Create module source file
- Create Makefile
- Build module: make
- Result: <module\_name>.ko file

### Loading and unloading modules

- List loaded modules: lsmod
- Insert module: sudo insmod <module\_name>.ko
- Remove module: sudo rmmod <module\_name>
- Load module with dependencies: sudo modprobe <module\_name>
- Unload module and dependencies: sudo modprobe -r <module\_name>
- View kernel messages: dmesg

### Module information

- Show module info: modinfo <module\_name>.ko
- Check if module is loaded: lsmod | grep <module\_name>
- Display module parameters: systool -vm <module\_name>
- View module details in sysfs: ls -l /sys/module/<module\_name>/

### Module autoloading

- Install module: sudo make modules\_install
- Update module dependencies: sudo depmod -a
- Configure autoloading: echo <module\_name> sudo tee /etc/modules-load.d/<module\_name>.conf

## Character Device Drivers

### Character Device Drivers

Character device drivers are a common type of Linux device driver:

- Handle devices that transfer data as a stream of bytes
- Support operations like read, write, open, release, ioctl
- Examples: serial ports, keyboards, mice, sensors
- Appear as files in /dev with major and minor numbers
- Major number identifies the driver
- Minor number identifies the specific device

## Basic Character Device Driver

Structure of a simple character device driver:

```
1 #include <linux/kernel.h>
2 #include <linux/module.h>
3 #include <linux/fs.h>
4 #include <linux/uaccess.h>
5
6 MODULE_LICENSE("GPL");
7
8 /* Prototypes */
9 static int device_open(struct inode *, struct file *);
10 static int device_release(struct inode *, struct file *);
11 static ssize_t device_read(struct file *, char *,
12                             size_t, loff_t *);
13 static ssize_t device_write(struct file *, const char
14                             *, size_t, loff_t *);
15
16 #define DEVICE_NAME "chardev"
17 #define BUF_LEN 80
18
19 /* Global variables */
20 static int Major; /* Major number assigned */
21 static int Device_Open = 0; /* Is device open? */
22 static char msg[BUF_LEN]; /* Message for the device */
23 static char *msg_Ptr;
24
25 static struct file_operations fops = {
26     .read = device_read,
27     .write = device_write,
28     .open = device_open,
29     .release = device_release
30 };
31
32 /* Initialization function */
33 int init_module(void)
34 {
35     Major = register_chrdev(0, DEVICE_NAME, &fops);
36     if (Major < 0) {
37         printk(KERN_ALERT "Failed to register with %d\n", Major);
38         return Major;
39     }
40
41     printk(KERN_INFO "Major number assigned: %d\n", Major);
42     printk(KERN_INFO "Create device with: 'mknod /dev/%s c %d 0'\n",
43            DEVICE_NAME, Major);
44     return 0;
45 }
46
47 /* Cleanup function */
48 void cleanup_module(void)
49 {
50     unregister_chrdev(Major, DEVICE_NAME);
51 }
52
53 /* Device methods */
54 static int device_open(struct inode *inode, struct file *file)
55 {
56     static int counter = 0;
57
58     if (Device_Open)
59         return -EBUSY;
60
61     Device_Open++;
62     sprintf(msg, "Called device_open %d times\n",
```

## Building Installing and Using a Character Device Module

Complete workflow for a character device driver:

```
1 # 1. Create source file (chardev.c) and Makefile as
   shown above
2
3 # 2. Build the module
4 make
5
6 # 3. Load the module
7 sudo insmod chardev.ko
8
9 # 4. Check kernel messages for major number
10 dmesg | tail
11
12 # 5. Create device node (assuming major number is 250)
13 sudo mknod /dev/chardev c 250 0
14 sudo chmod 666 /dev/chardev
15
16 # 6. Test reading from the device
17 cat /dev/chardev
18
19 # 7. Check module information in sysfs
20 ls -l /sys/module/chardev/
21 cat /sys/module/chardev/parameters/* 2>/dev/null
22
23 # 8. Unload the module when done
24 sudo rmmod chardev
25
26 # 9. Clean up the device node
27 sudo rm /dev/chardev
```

This demonstrates:

- Building and loading a custom character device driver
- Creating a device node with appropriate permissions
- Interacting with the device through standard file operations
- Examining module information through sysfs
- Proper cleanup when the module is no longer needed

## Linux Kernel Module Debugging

### Using printk

- Add debug messages: `printk(KERN_DEBUG "Debug: %d\n", value);`
- Set console log level: `echo 7 > /proc/sys/kernel/printk`
- View kernel messages: `dmesg`
- Follow kernel messages: `dmesg -w`
- Clear buffer: `dmesg -c`

### Debug filesystem

- Mount debugfs: `mount -t debugfs none /sys/kernel/debug`
- Create debug files in your module:
  - Include `<linux/debugfs.h>`
  - Create directory: `debugfs_create_dir()`
  - Create files: `debugfs_create_file()`
- Access from user space: `cat /sys/kernel/debug/mymodule/myfile`

### Module parameters

- Define parameters: `module_param(name, type, permissions);`
- Load with parameters: `insmod mymodule.ko debug=1`
- Change at runtime: `echo 1 > /sys/module/mymodule/parameters/debug`

### Kernel/system crashes

- Configure kdump: Install and configure kdump-tools
- Analyze crash dumps with crash utility
- Check system logs after reboot: `journalctl -b -1`

## I/O Performance Testing with Custom I/O Scheduler

Testing I/O performance with different schedulers:

```
1 # 1. Check available I/O schedulers
2 cat /sys/block/sda/queue/scheduler
3
4 # 2. Test performance with default scheduler
5 echo 3 > /proc/sys/vm/drop_caches # Clear cache
6 dd if=/dev/zero of=testfile bs=1M count=1000
   oflag=direct
7 rm testfile
8
9 # 3. Change to a different scheduler (e.g., BFQ)
10 echo bfq > /sys/block/sda/queue/scheduler
11
12 # 4. Test performance with new scheduler
13 echo 3 > /proc/sys/vm/drop_caches # Clear cache
14 dd if=/dev/zero of=testfile bs=1M count=1000
   oflag=direct
15 rm testfile
16
17 # 5. Test with different I/O priorities
18 # Start a background write process
19 dd if=/dev/zero of=bg_file bs=1M count=2000
   oflag=direct &
20 BG_PID=$!
21
22 # Run a foreground process with higher priority
23 ionice -c2 -n0 dd if=/dev/zero of=fg_file bs=1M
   count=500 oflag=direct
24
25 # Clean up
26 kill $BG_PID
27 rm bg_file fg_file
28
29 # 6. Return to the default scheduler
30 echo deadline > /sys/block/sda/queue/scheduler
```