# 1 Overview

## 1.1 What Operating Systems Do

$$\left.\begin{array}{r}\text{hardware}\\\text{software}\\\text{data}\end{array}\right\}\text{Computer System}\left\{\begin{array}{l}\text{Hardware}\left\{\begin{array}{l}\text{CPU (Central Processing Unit)}\\\text{Memory}\\\text{I/O (Input/Output) Devices}\end{array}\right.\\\text{Application Programs}\quad\text{defines how resources are used}\\\textbf{Operating System}\quad\text{controls and coordinates hardware}\end{array}\right.$$
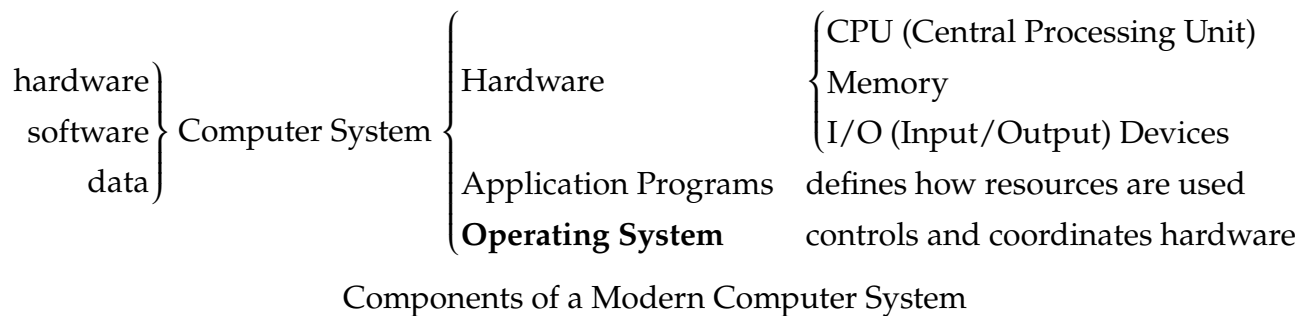
<div align="center">Components of a Modern Computer System</div>
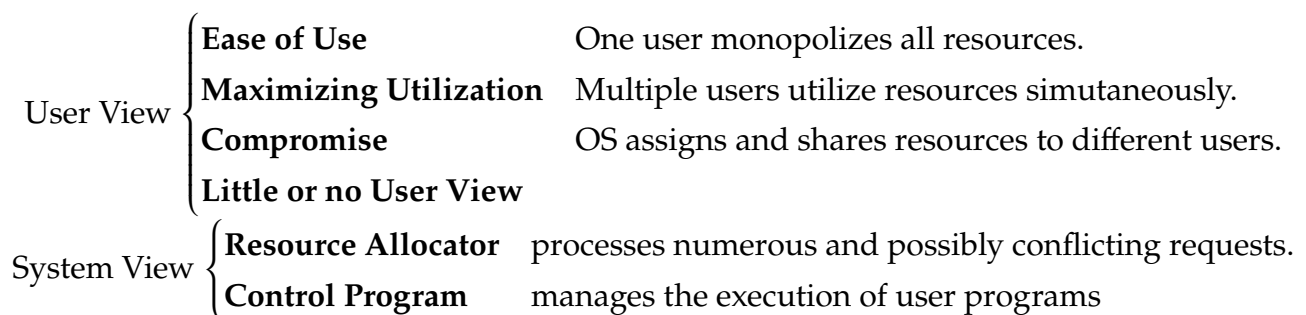
An **operating system** is a software that
- manages and controls a computer's hardware;
- coordinates and optimizes utilization of hardware;
- provides a basis for application programs.

An operating system is similar to a *government*, who performs no useful function, but provides an environment within which other programs can do useful work.

$$\text{User View}\left\{\begin{array}{ll}\textbf{Ease of Use} & \text{One user monopolizes all resources.}\\\textbf{Maximizing Utilization} & \text{Multiple users utilize resources simutaneously.}\\\textbf{Compromise} & \text{OS assigns and shares resources to different users.}\\\textbf{Little or no User View}\end{array}\right.$$

$$\text{System View}\left\{\begin{array}{ll}\textbf{Resource Allocator} & \text{processes numerous and possibly conflicting requests.}\\\textbf{Control Program} & \text{manages the execution of user programs}\end{array}\right.$$

## 1.2 Computer-System Organization

### 1.2.1 Computer-System Operation

Rquires Memory Controllers — Memory Cycles

Execute in Parallel

CPU(s)  Device Controllers $\left\{\begin{array}{l}\text{disk drives}\\\text{video displays}\\\dots\end{array}\right.$

Common Bus

Memory
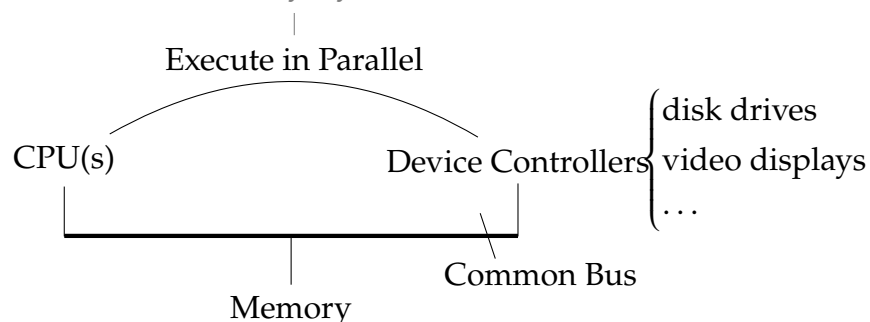
<div align="center">Figure: Components of Modern General-Purpose Computer</div>
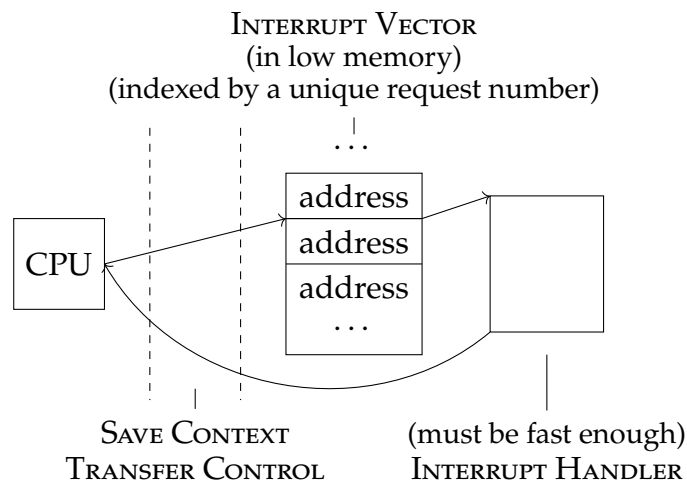
For a computer to start running, it
1. runs **bootstamp program**, which
   - tends to be simple.
   - is stored in **read-only memory** (ROM), or Electrically Erasable Programmable ROM
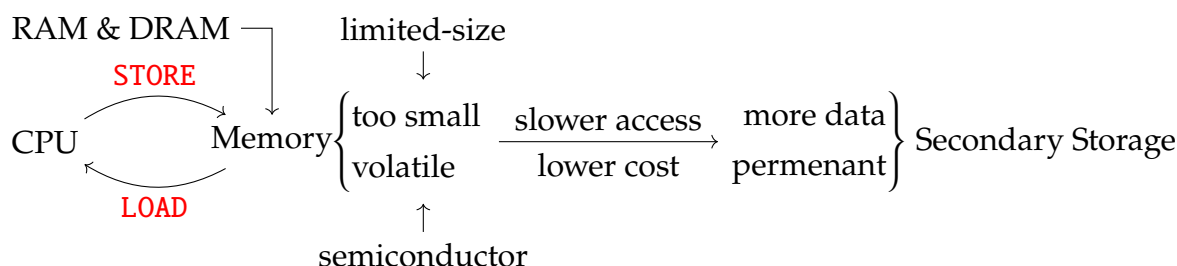
- initializes all aspects of the OS, from CPU registers to device controllers to memory.
- locates the operating system and loads it to memory ($\Leftarrow$ know how to load and start)
2. loads service programs (**system daemons**: outside kernel, loaded at boot, runs entire time)

The event is signaled by an **interrupt** from either hardware or software.



INTERRUPT VECTOR
(in low memory)
(indexed by a unique request number)

CPU

address
address
address
. . .

SAVE CONTEXT
TRANSFER CONTROL

(must be fast enough)
INTERRUPT HANDLER

### 1.2.2 Storage Structure

All forms of memory provide **an array of bytes**. Each byte has its own address.



RAM & DRAM — limited-size

STORE

CPU — Memory $\begin{cases} \text{too small} \\ \text{volatile} \end{cases}$ $\xrightarrow[\text{lower cost}]{\text{slower access}}$ $\begin{rcases} \text{more data} \\ \text{permenant} \end{rcases}$ Secondary Storage

LOAD

semiconductor

Other types of memory:
- Cache: stores data to reduce time cost of further request for that data.
- ROM: cannot be changed $\Rightarrow$ ONLY static programs (e.g., bootstamp program).
- EEPROM: change is slow $\Rightarrow$ mostly static programs (e.g., factory-installed programs).

| Hierarchy | Magnitude | Volatility | Implementation |
|---|---|---|---|
| Registers | bytes | $\checkmark$ | MOSFET |
| Cache | 16KB $\sim$ 50MB | $\checkmark$ | MOSFET |
| Main Memory | 8GB $\sim$ 64GB | $\checkmark$ | MOSFET |
| **Solid-state Disks** | $\geq$ 100 GB | $\bigcirc$ / $\times$ | Flash Memory |
| **Magnetic Disks** | $\geq$ 500 GB | $\times$ | Magnetic Poles |
| **Optical Disks** | | $\times$ | Pits & Lands |
| **Magnetic Tapes** | TB | $\times$ | Magnetic Memory |

Table: Information and Hierarchy of Storage
(higher in hierarchy $\implies$ larger capacity, more expensive, and faster)
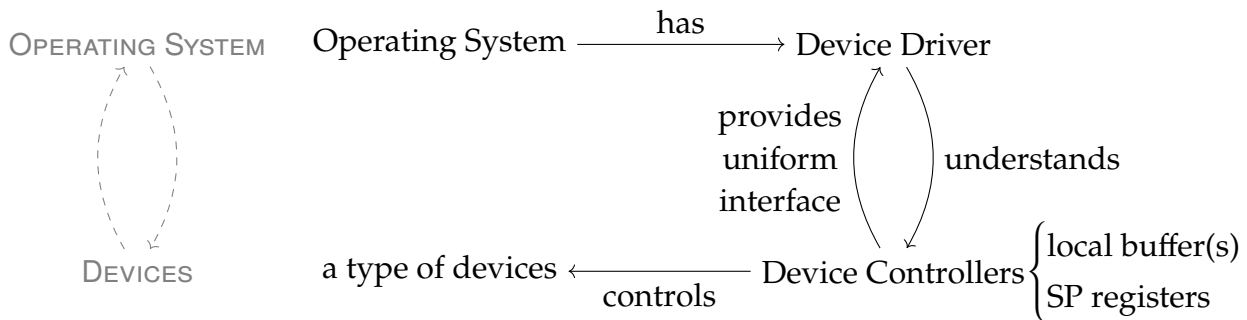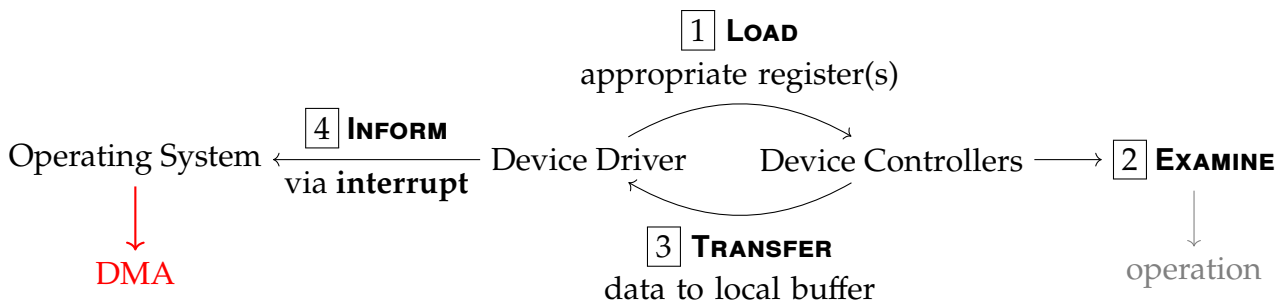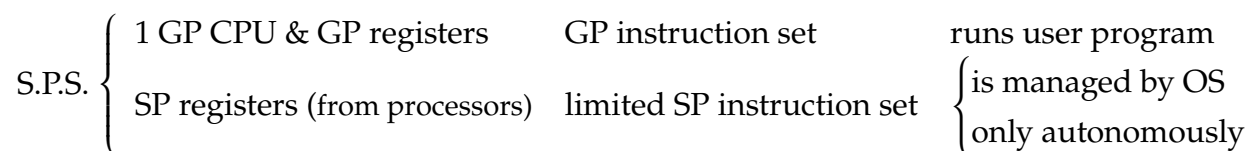
### 1.2.3 I/O Structure



Figure: I/O Structure



Figure: I/O Operation

This form creates overhead when bulk and/or frequent data movement, like disk and keyboard. By **direct memory access** (DMA), the driver fires only one interrupt and transfers a block of data from its local buffer the main memory, without CPU's intervention.

## 1.3 Computer-System Architecture
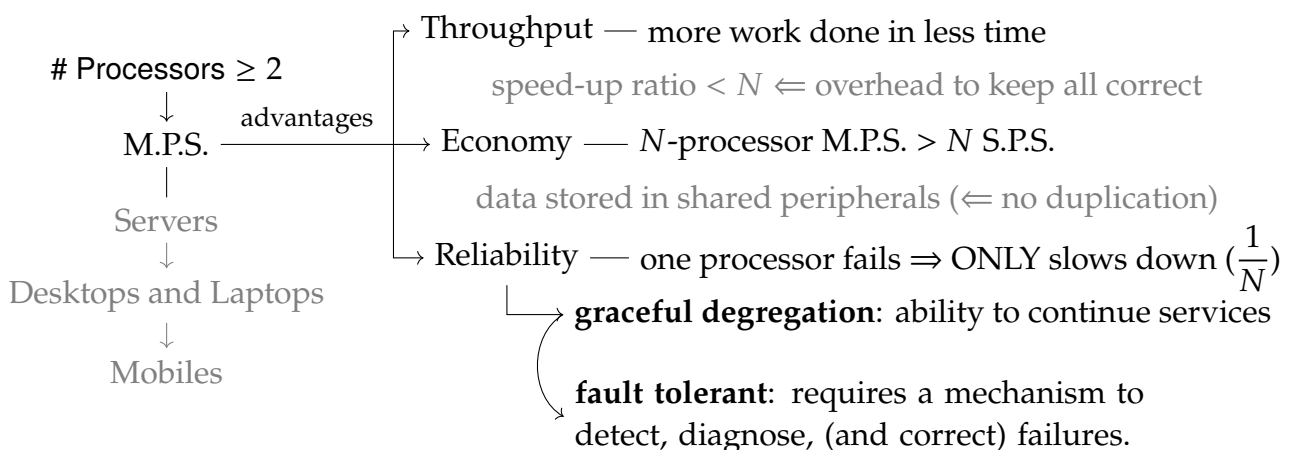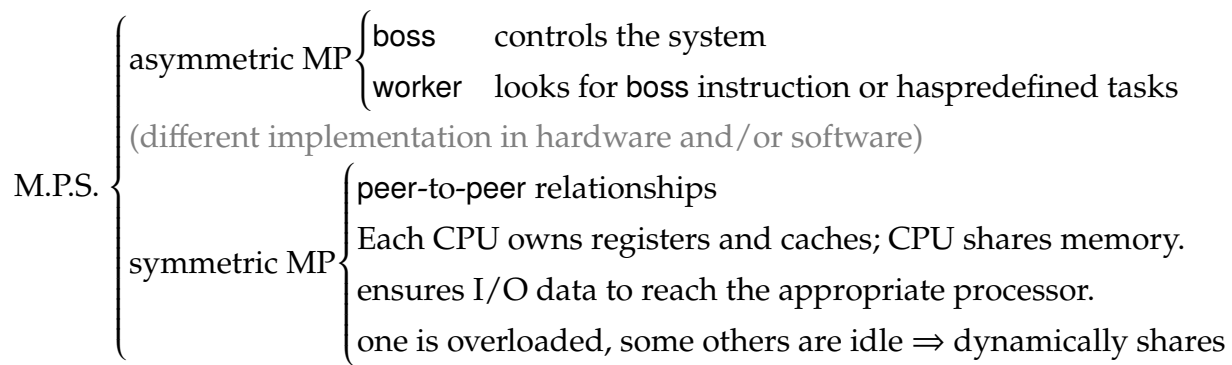
### 1.3.1 Single-Processor Systems

$$\text{S.P.S.} \begin{cases} \text{1 GP CPU \& GP registers} & \text{GP instruction set} & \text{runs user program} \\ \text{SP registers (from processors)} & \text{limited SP instruction set} & \begin{cases} \text{is managed by OS} \\ \text{only autonomously} \end{cases} \end{cases}$$

### 1.3.2 Multiprocessor Systems



Figure: Multiprocessor System Concepts

$$\text{M.P.S.} \begin{cases} \text{asymmetric MP} \begin{cases} \text{boss} & \text{controls the system} \\ \text{worker} & \text{looks for boss instruction or haspredefined tasks} \end{cases} \\ \text{(different implementation in hardware and/or software)} \\ \text{symmetric MP} \begin{cases} \text{peer-to-peer relationships} \\ \text{Each CPU owns registers and caches; CPU shares memory.} \\ \text{ensures I/O data to reach the appropriate processor.} \\ \text{one is overloaded, some others are idle} \Rightarrow \text{dynamically shares} \end{cases} \end{cases}$$

Types of Multiprocessor System

**Multicore**: includes multiple computing cores (owns registers and local cache) on a single chips; on-chip communication is faster and uses significantly less power than between-chip communication.

### 1.3.3 Clustered Systems

A **clustred sytem** are composed of two or more individual systems, or nodes, joined together.

$$\textbf{Clustering} \begin{cases} \text{provides} \begin{cases} \text{high-availability service (one fails} \Rightarrow \text{continues)} \Leftarrow \text{Redundancy} \\ \text{high-performance computing environment} \Leftarrow \text{Parallelization} \end{cases} \\ \text{type} \begin{cases} \text{assymmetric} \begin{cases} \text{hot-standby} & \text{does nothing, monitors active server} \\ & \text{fails} \Rightarrow \text{becomes another active server} \\ \text{other} & \text{runs applications} \end{cases} \\ \text{symmetric: } \geq 2 \text{ machines run and monitor each other.} \end{cases} \end{cases}$$
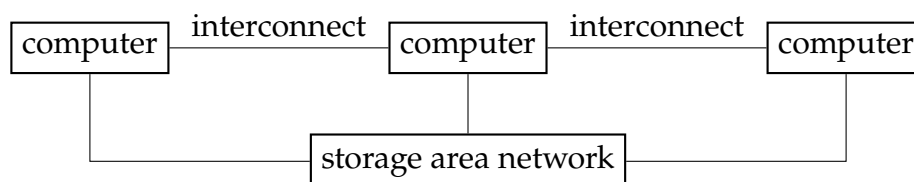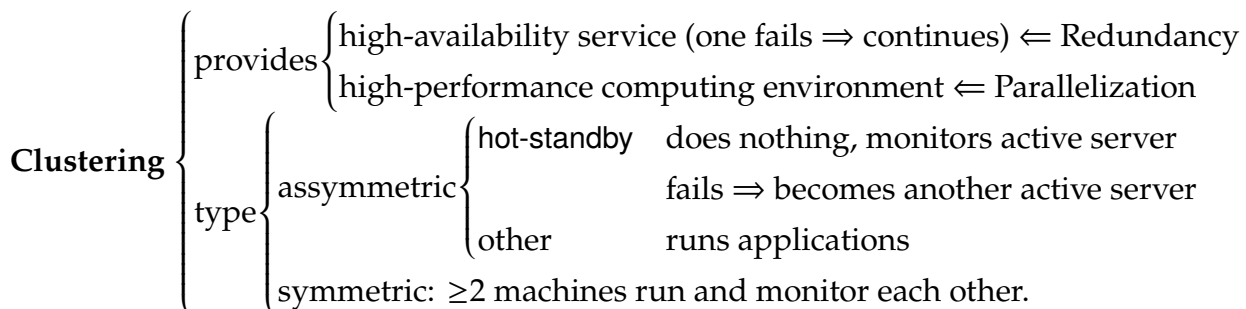


Figure: General structure of a clustered system

**Parallelization**: divides a program into separate components to run on individual computers in the cluster $\implies$ much greater computational power (significantly greater than multiple single-processor systems or even symmetric multiprocessor systems).
**Parallel clusters**: multiple hosts to access data on shared storage $\Rightarrow$ access control and locks

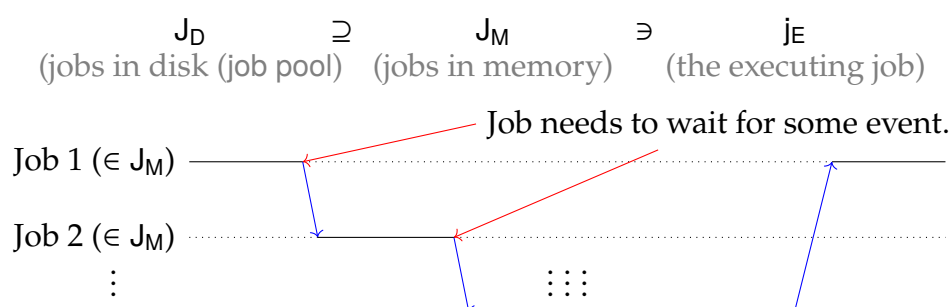## 1.4 Operating-System Structure
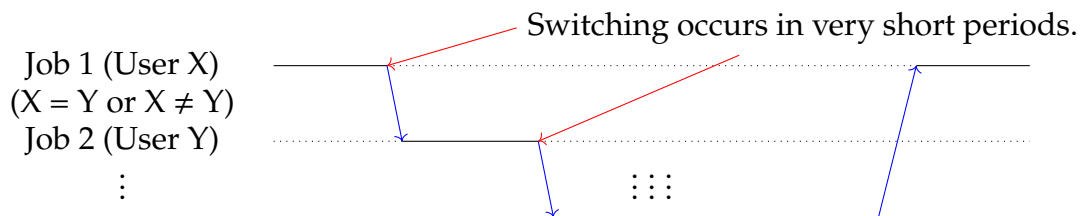


Figure: Multiprogramming Idea Interpretation

Figure: Multitasking Idea Interpretation

(Multitasking requires an **interactive** system, with short device responsive time.)

Job scheduling: OS needs to choose ready jobs from disk because memory is too small.
CPU scheduling: OS needs to choose a job in memory when multiple jobs are ready to run.
Swapping: processes are swapped in and out of main memory.

## 1.5   Operating-System Operations

The OS must ensure that incorrect or malicious behaviors in a program cannot cause other programs and the OS to execute incorrectly.
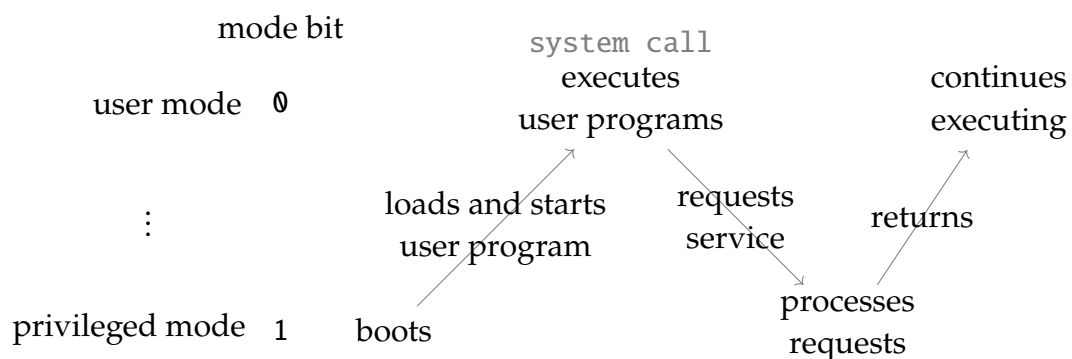
### 1.5.1   Dual-Mode and Multimode Operation



Figure: Dual-Mode Interpretation

Some instructions that might cause harm are called privileged instructions. The hardware only supports privileged instructions and some other operations, like I/O controls, interrupt management, etc., in privileged mode. Some modern architectures support more than two modes.

The lack of hardware-supported dual mode causes serious shortcomings. For example, a user program can overwrite the OS, which might make the system crash or behave oddly.
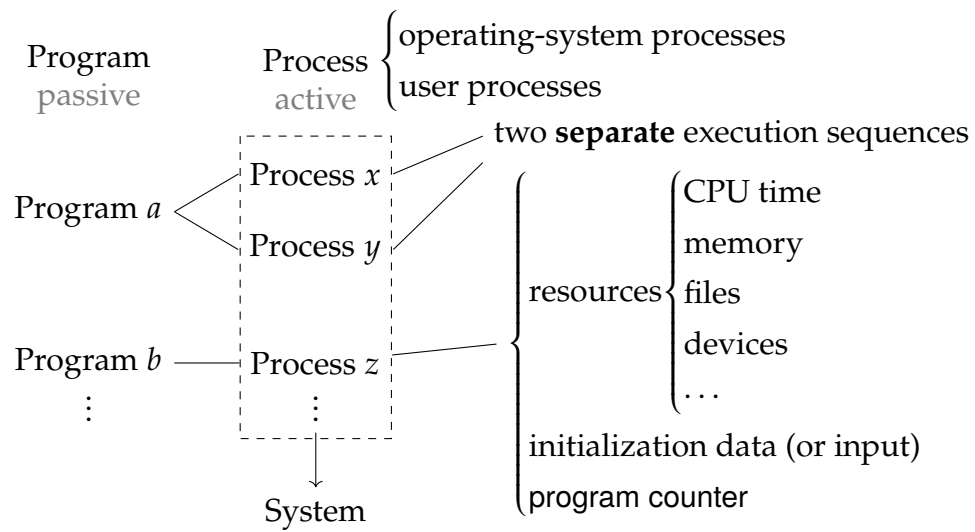
If a privilege violation is occurred, these errors are handled by the OS. The OS must terminate the program abnormally, probably with an appropriate error message.

### 1.5.2   Timer

Timer is to detect infinite loops, failed returns, user programs running too long, etc. It can interrupt the computer after a specified period. It is usually implemented by a fixed-rate clock and a counter set by the OS. The OS might terminate the program or give more time. Instructions modifying the timer are privileged.

## 1.6  Process Management

A process is a program **in execution**.



## 1.7  Memory Management