# Week 3 Virtualization and Containerization

Dr Zhi Zhang

---

## Overview

- Operating System Concepts
- Virtualization
- Containerization

---

## Operating System Concepts
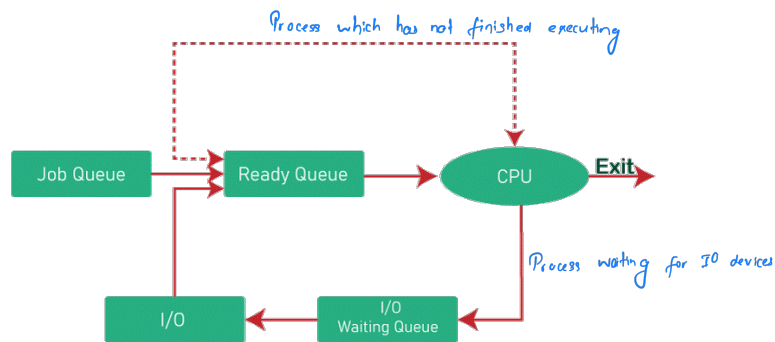
- Process scheduling
- Virtual memory

---

## What is a process?

- A process is a program that is executed
- It is a basic unit of execution in an OS.
- At a minimum, process execution requires <span style="color:red">CPU and memory without I/O devices.</span>

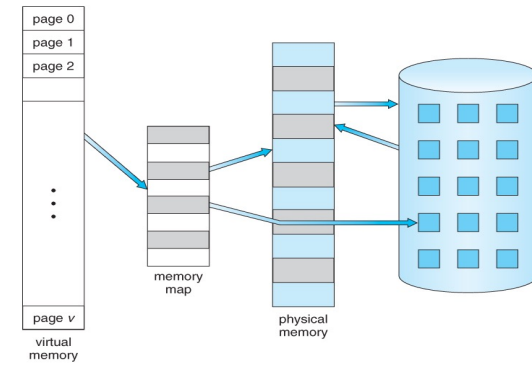## Process scheduling (time-sharing)

- Each process is scheduled and executed within their time share.
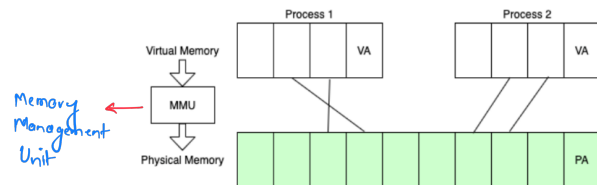- The number of processes can be much larger than that of the CPUs.

## Process scheduling

Process which has not finished executing



Job Queue → Ready Queue → CPU → **Exit**

Process waiting for IO devices

I/O ← I/O Waiting Queue

---

## Virtual memory



page 0
page 1
page 2
page v
virtual memory

memory map

physical memory

---

## Virtual memory

- Page table:
  - In a modern OS (e.g., Linux, Windows or MacOS), the OS uses a set of page tables to map virtual memory within a process to their corresponding physical memory in main memory.



Process 1        Process 2

Virtual Memory        VA        VA

Memory Management Unit ← MMU

Physical Memory        PA

  - VA refers to virtual address, pointing to a memory cell in a virtual page.
  - PA refers to physical address, pointing to a memory cell in a physical page.

---

## Virtual memory

- Page table:
  - In a modern OS (e.g., Linux, Windows or OS X), the OS uses a set of page tables to map virtual memory within a process to their corresponding physical memory in main memory.
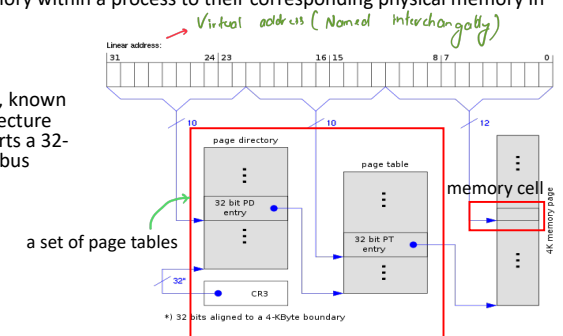
- 32-bit version of x86, known as IA-32 (Intel Architecture 32-bit), which supports a 32-bit memory address bus

Virtual address ( Named Interchangably )



Linear address:
31    24 23    16 15    8 7    0

page directory

32 bit PD entry

page table

32 bit PT entry

memory cell

CR3

a set of page tables

*) 32 bits aligned to a 4-KByte boundary

4K memory page

## What is virtualization?

- Virtualization is the ability to run multiple operating systems on a single physical system and share the underlying hardware resources [1]
- Allows one computer software (called Virtual Machine Monitors or Hypervisor) to provide the appearance of many computers (called virtual machines).
- Goals:
  - Provide flexibility for users
  - Amortize hardware costs
  - Isolate completely separate users

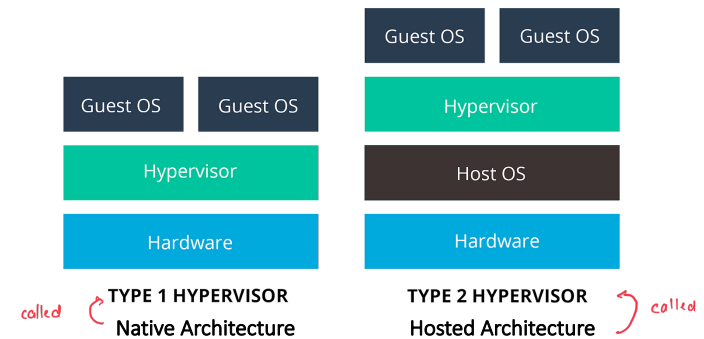[1] VMWare white paper, Virtualization Overview

## Characteristics for Virtualizable Computer Architectures

- First, the VMM provides an environment for VM. The environment is essentially identical with the original machine.
- Second, VM running in this environment shows minor decreases in speed.
- Last, the VMM is in complete control of all system resources, such as CPUs, memory, disk, etc.

## VMM Types

- **Hosted Architecture**
  - Install as an application on existing x86 "host" OS, *e.g.* Windows, Linux, OS X
  - Leverage host I/O stack and resource management
  - Examples: Virtualbox, VMware
- **Bare-Metal/Native Architecture**
  - VMM or Hypervisor is installed directly on hardware
  - Acknowledged as preferred architecture for mainstream public clouds
  - Examples: KVM, Xen, Hyper-V

## VMM Types

| Guest OS | Guest OS |
|----------|----------|
| Hypervisor | |
| Hardware | |

called

**TYPE 1 HYPERVISOR**
Native Architecture

| Guest OS | Guest OS |
|----------|----------|
| Hypervisor | |
| Host OS | |
| Hardware | |

**TYPE 2 HYPERVISOR**
Hosted Architecture

called

## Virtualization: history

- 1960's: first track of virtualization
  - Time and resource sharing on expensive computers
  - IBM VM/370
- Late 1970s and early 1980s: became unpopular
  - Cheap hardware and multiprocessing OS
- **Late 1990s:** became popular again
  - Wide variety of OS and hardware configurations (compatibility issue)
  - VMWare
- **Since 2000:** hot and important
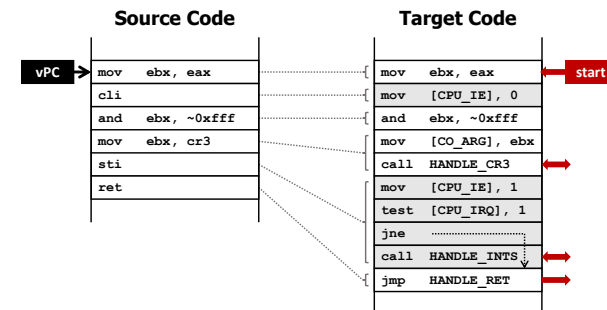  - Cloud computing
  - Containers

## Virtualization on x86 architecture

- Solutions: → *Generations of virtualization solutions*
  - Full virtualization implemented by Dynamic binary translation & Two-stage translation
  - Para-virtualization (Xen)
  - Full virtualization implemented by Hardware extension

## Dynamic Binary Translation

- The entire program (e.g., ARM format) is translated into a binary of another architecture (e.g., x86 format).
- The translation occurs on-the-fly during program execution.
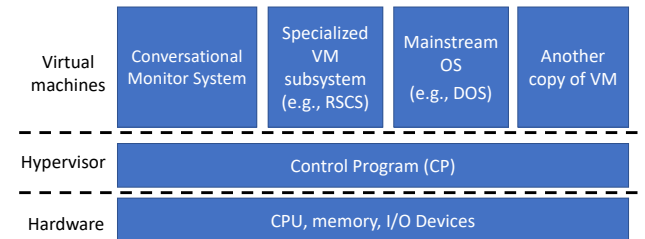
## Dynamic Binary Translation



Courtesy Scott Define VMware Inc

## Dynamic Binary Translation

- For sensitive/privileged instructions, trap-and-emulate is needed.
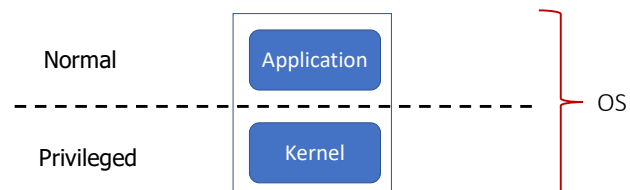  - e.g., IBM VM/370: the first virtualized system

---

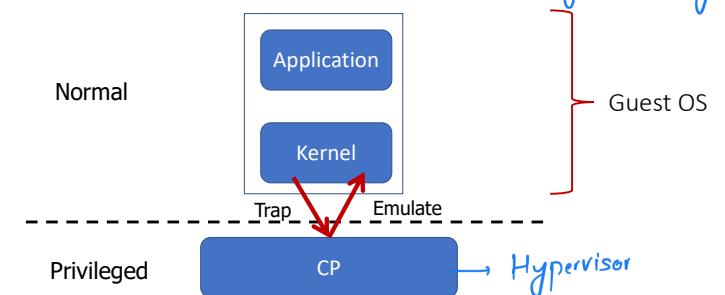## IBM VM/370: the first virtualized system

| Virtual machines | Conversational Monitor System | Specialized VM subsystem (e.g., RSCS) | Mainstream OS (e.g., DOS) | Another copy of VM |
|---|---|---|---|---|
| Hypervisor | Control Program (CP) | | | |
| Hardware | CPU, memory, I/O Devices | | | |

---

## IBM VM/370

- Virtualization technology: trap-and-emulate

Normal

Application

Privileged

Kernel

OS

This is how normal OS would look. Forget about Guest OS here.

---

## IBM VM/370

When guest OS wants to run some privileged instructions, those instructions will create a trap that goes into hypervisor and hypervisor will emulate the intended functionality of the guest os

- Virtualization technology: trap-and-emulate

Normal

Application

Kernel
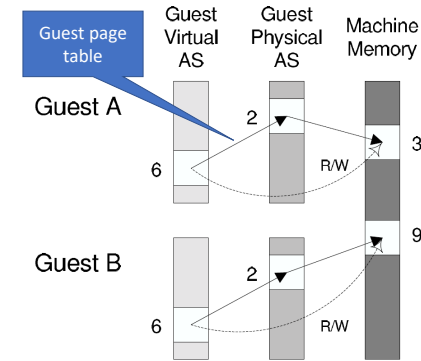
Guest OS

Trap   Emulate

Privileged

CP → Hypervisor

## Two-stage translation

- Guest page table:
  - A virtual machine has its own guest page table (GPT), managed by the guest OS. The GPT translates guest virtual addresses to guest physical addresses.
- Shadow page table:
  - The hypervisor creates a shadow page table for each guest page table. The shadow page table translates guest virtual addresses to machine addresses.
- Two-stage translation:
  - When a virtual machine running on the hypervisor performs a memory access, the virtual address goes through two stages of translation.
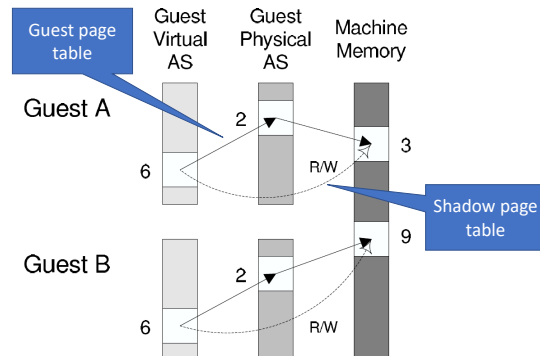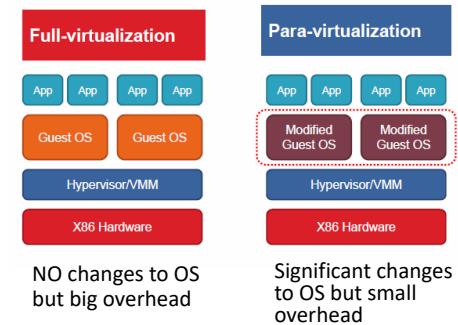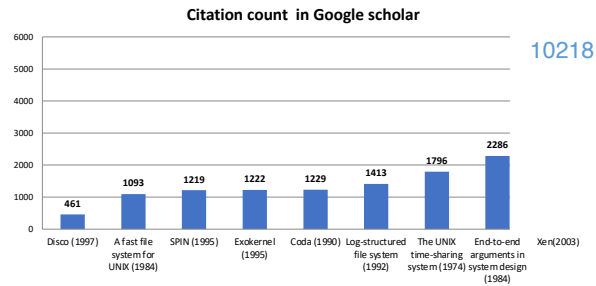


---

## Two-stage translation

---

## Two-stage translation

---

## Para-virtualization

- Full vs. Para virtualization



NO changes to OS but big overhead

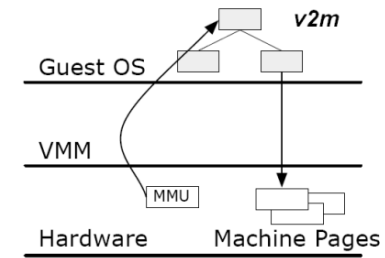Significant changes to OS but small overhead

## A typical example: Xen

- SOSP (ACM Symposium on Operating Systems Principles ) in 2003.
- Very high impact (data collected in 2013 and 2023)



**Citation count in Google scholar**

10218

| | |
|---|---|

Disco (1997): 461, A fast file system for UNIX (1984): 1093, SPIN (1995): 1219, Exokernel (1995): 1222, Coda (1990): 1229, Log-structured file system (1992): 1413, The UNIX time-sharing system (1974): 1796, End-to-end arguments in system design (1984): 2286, Xen(2003)
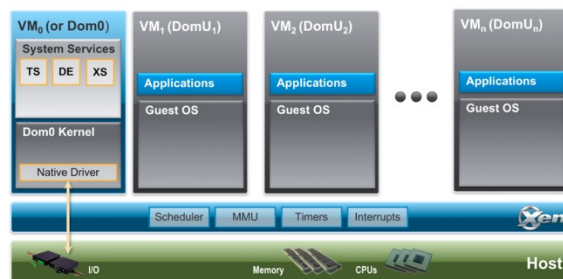
## Overview of Xen

- Modify guest OS to be aware of virtualization
  - No emulation of privileged operations from guest OS
  - Better performance

**(a) MMU Para-virtualiztion**



*Single Translation stage but compatibility issue*

## Xen architecture



- TS is Toolstack, providing a user interface to the Xen hypervisor administrator.
- DE is Device Emulation (DE), emulating devices when the devices are not available.
- XS is XenStore/XenBus, managing information that are shared between some domains.

https://wiki.xenproject.org/wiki/Xen_Project_Software_Overview
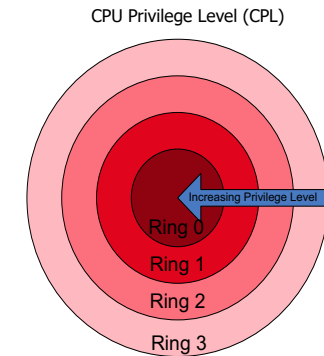
## A short conclusion

- Full virtualization
  - Unmodified guest OS
  - Performance issue
- Para virtualization:
  - Better performance
  - Modified guest OS
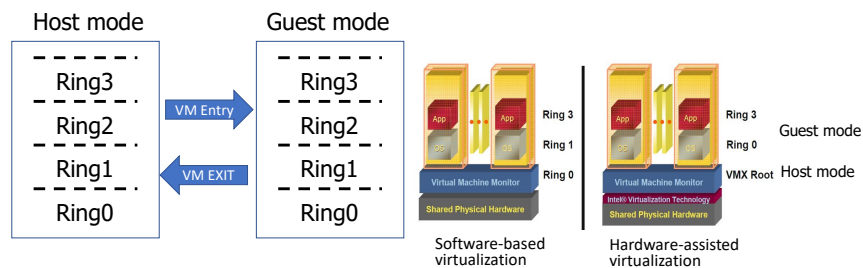
## Support from hardware extensions

- Full/Para virtualization: software-based virtualization

- Hardware support: Intel and AMD assist virtualization

## x86 Protection Rings (CPL)

CPU Privilege Level (CPL)

- x86 has four protection levels (rings) /CPU privilege levels (CPLs).
  - ring 0 – "Kernel mode" (most privileged)
  - ring 3 – "User mode"
  - ring 1 & 2 – Other
- Linux only uses ring 0 and ring 3.
  - "Kernel vs. user mode"
- For previous full virtualization, guest OS applications run in ring 3, guest OS kernel runs in ring 1, and the hypervisor runs in ring 0.
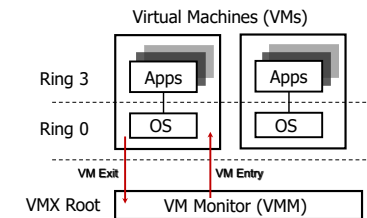


Increasing Privilege Level

Ring 0
Ring 1
Ring 2
Ring 3

## CPU virtualization: Intel VT-x & AMD SVM



Host mode          Guest mode

Ring3              Ring3
Ring2   VM Entry   Ring2
Ring1   VM EXIT    Ring1
Ring0              Ring0

Software-based virtualization

Hardware-assisted virtualization

## An example: Intel VT-x

- Two VT-x operating modes
  - Less-privileged mode (guest or VMX non-root) for guest OSes
  - More-privileged mode (host or VMX root) for VMM
- Two transitions
  - VM entry to non-root mode
  - VM exit to root mode



Virtual Machines (VMs)

Ring 3        Apps        Apps

Ring 0        OS          OS

VM Exit       VM Entry
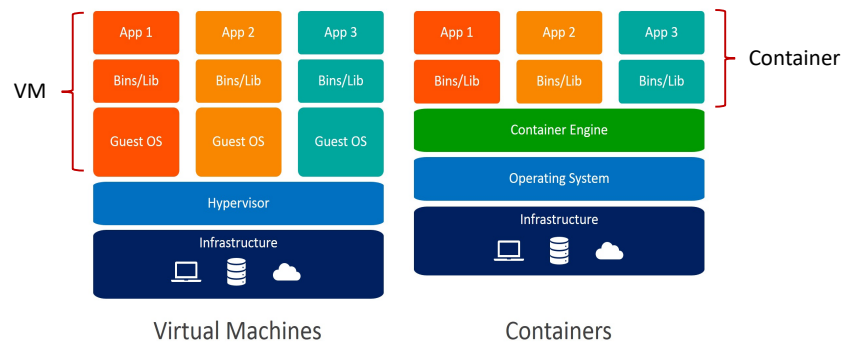
VMX Root      VM Monitor (VMM)

## CPU virtualization

- New VM mode bit
  - Create two orthogonal modes
- If VM mode bit is clear → host mode
- If VM mode bit is set → guest mode

- Benefits:
  - Eliminating the need of binary translation: significant performance improvement.
  - Requiring no changes to privileged operations in guest OSes: compatible with existing modern OSes.

## Containers: concepts

- A container is a sandboxed process running on a host OS that is isolated from all other processes running on that host OS [2]
- A container image is a stand-alone and executable software package (e.g., dependencies, binaries, etc) that contains everything needed to run an application.
- A container is a runnable instance of the image.

[2] https://docs.docker.com/get-started/

## Containers



| App 1 | App 2 | App 3 |
|-------|-------|-------|
| Bins/Lib | Bins/Lib | Bins/Lib |
| Guest OS | Guest OS | Guest OS |

Hypervisor

Infrastructure

Virtual Machines

| App 1 | App 2 | App 3 |
|-------|-------|-------|
| Bins/Lib | Bins/Lib | Bins/Lib |

Container Engine

Operating System

Infrastructure

Containers

https://www.weave.works/blog/a-practical-guide-to-choosing-between-docker-containers-and-vms

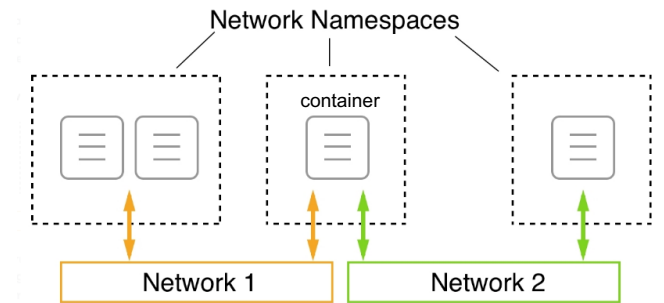## Containerization vs. Virtualization

- Granularity
  - Containers are an abstraction of the process layer and VMs are a simulation of the hardware layer.
- Overhead
  - Required resources: containers are created to run one application and VMs support a whole OS.
  - Efficiency: containers are launched to run an application. VMs need to boot up an entire OS.
- Security/Isolation
  - Containers are isolated from each other at the process level. VMs are isolated at the OS level.

## Backbones of containers

- namespaces
- cgroups
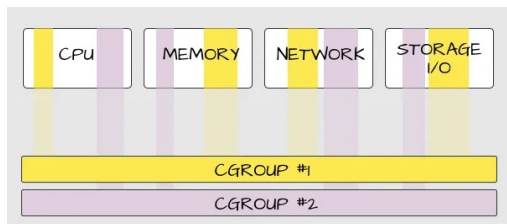
## namespaces: limit what a container can see

- Provide containers with their own view of system resources (e.g., network)
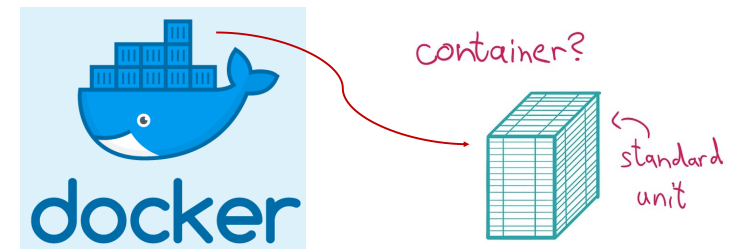
### Network Namespaces

## cgroups: limits how much a container can use

- Control Groups: manage resource usage and limits for processes within a container.
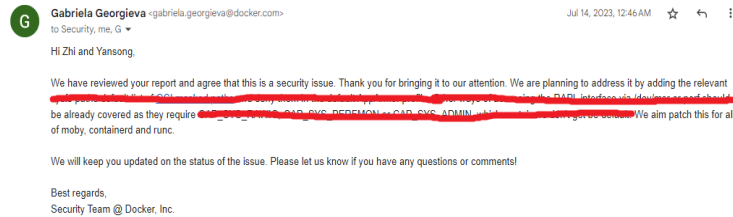  - Example: a container is limited to use only 1 Mbps of network bandwidth



## Docker container

- Docker provides an interface on top of the techniques
  - Popularized containers (a standard unit of software)

## Docker container



Gabriela Georgieva <gabriela.georgieva@docker.com>   Jul 14, 2023, 12:46 AM
to Security, me, G

Hi Zhi and Yansong,

We have reviewed your report and agree that this is a security issue. Thank you for bringing it to our attention. We are planning to address it by adding the relevant ~~... BPF interface via /dev/... and should~~ be already covered as they require ~~... CAP_SYS_PERFMON or CAP_SYS_ADMIN ...~~ We aim patch this for all of moby, containerd and runc.

We will keep you updated on the status of the issue. Please let us know if you have any questions or comments!

Best regards,
Security Team @ Docker, Inc.

## Docker cheatsheet

- docker run – Runs a command in a new container.
- docker start – Starts one or more stopped containers
- docker stop – Stops one or more running containers
- docker build – Builds an image form a Docker file
- docker pull – Pulls an image or a repository from a registry
- docker push – Pushes an image or a repository to a registry
- docker export – Exports a container's filesystem as a tar archive
- docker exec – Runs a command in a run-time container
- docker search – Searches the Docker Hub for images
- docker attach – Attaches to a running container
- docker commit – Creates a new image from a container's changes

## Practice Questions

- [7 marks] Q1: Describe what virtualization is and describe the characteristic attributes of the different types of virtualization (Language, Operating System and Hardware).

- [1 mark] Virtualization is a technology that allows multiple virtual instances of resources (e.g., operating systems) to run on a single physical machine.
- [2 marks] **Language virtualization**: provides a virtual runtime environment (e.g., JVM) that allows code written in a specific programming language (e.g., Java) to execute.
  - Cross-OS: code can be executed in a runtime environment across OSes without modification.
- [2 marks] **Operating system virtualization**: allows multiple user-space instances (e.g., docker containers) within a single operating-system kernel (e.g., Linux kernel).
  - Limited Isolation: instances share the host OS kernel, which means a vulnerability in the kernel could affect all instances.
- [2 marks] **Hardware virtualization**: creates multiple virtual machines sharing a single physical server.
  - VM Snapshots: allows saving the current state of a VM
  - Live VM Migration: enables moving runtime VMs between physical hosts seamlessly.