

Table of content

Table of content	1
C1: Introduction to IT Infrastructure	4
IT infrastructures become more complicated due to new types of applications:	4
Agile adaptations require infrastructure:	4
Infrastructure Architecture	4
Definition	4
IT Architecture	5
Architecture is crucial to control the infrastructure when it is:	5
Strategy vs technology focus	5
IT building blocks & NFA	5
Non-functional attributes (NFRs)	7
Cloud Computing	7
Cloud Characteristics	8
Cloud Deployment Models (4)	8
Cloud Service Models (3)	8
Edge computing	8
C2: Non-Functional Attributes	10
Availability	10
Performance	12
Security	16
C3: Datacenter	18
3.1 Functions of a Datacenter	18
3.2 Datacenter Building Blocks	18
a) Sub Equipment Room (SER)	18
b) Main Equipment Room (MER)	18
c) Organization Owned Datacenter (OOD)	18
d) Multi-Tenant Datacenter (MTD)	18
Factors Determining the Locations of Datacenters	18
Physical structure: FWWDP	18
Floors	18
Vents	19
Walls, Windows, Doors	19
Water and Gas Pipes	19
Layout of Datacenter – Rooms and Components	19
🔌 Power Density and Server Rack	20
⚡ Uninterruptible Power Supply (UPS) Systems	20
🔌 Power Distribution Units (PDUs)	20
❄️ Cooling Systems in Datacenters	21
🧑‍🔧 Environmental Guidelines for Datacenter Operation	21
🔥 Fire Prevention, Detection, and Suppression	21
📦 Equipment Racks	22
⚡ Datacenter Energy Efficiency	22
3.3 Datacenter APS	22
🏢 Datacenter Availability- 4 Tiers	22

Redundant Datacenters	22
🧹 Floor Management (by Floor Manager)	22
🚀 Datacenter Performance	22
🔒 Datacenter Security	23
C4: Networking	24
Network Topologies	24
① Physical Layer	24
② Data Link Layer	27
③ Network Layer	28
④ Transport Layer	28
⑤ Session Layer	29
⑦ Application Layer	29
4.2 Networking Virtualization	31
4.3 Networking: Availability, Performance & Security	32
Networking Availability	32
Networking Performance	34
Networking Security	35
C5: Storage	38
5.1 Storage building blocks	38
① Disks – command sets	38
① Disks - type (2)	38
② Tape (2)	39
③ Controller Implementation	39
④ Storage Architectures (4)	40
⑤ Software Defined Storage (SDS)	40
5.2 Storage: Availability, Performance & Security	41
① Improve Storage Availability (4)	41
② Storage Performance (3)	41
③ Improve Storage Security(2)	42
C6: Compute	43
6.1 Compute building blocks	43
6.2 Compute Virtualization	44
6.3 Types of Compute	45
6.4 Compute: Availability, Performance & Security	46
Chapter 7: Operating systems	49
7.1 Popular Operating Systems	49
7.2 Operating systems building blocks	49
7.3 Operating systems: Availability, Performance & Security	50
Chapter 8: End User Devices	52
8.1 End user devices building blocks	52
8.2 Desktop virtualization	52
8.3 End user devices: Availability, Performance & Security	53
C9: Infrastructure Management	55
9.1 Infrastructure Deployment Options	55
Infrastructure deployment models (4)	55
On-Premises	55

BMIS2113 Information Technology Infrastructure

Public Cloud	55
Private Cloud	55
Hybrid Cloud	55
9.2 Infrastructure Automation	55
Configuration Management Tools	55
Orchestration Tools	56
9.3 Infrastructure Documentation	56
Infrastructure documentation tools / techniques	56

C1: Introduction to IT Infrastructure

IT infrastructures become more complicated due to new types of applications:

E-Commerce with Personalized Shopping

- Mobile Computing: Users shop via mobile apps with personalized recommendations.
- Cloud Computing: E-commerce platform scales resources dynamically to handle millions of transactions.
- Big Data: Purchase histories, browsing patterns, and social media activities are analyzed.
- Artificial Intelligence: AI engines provide product recommendations and chatbot support.
- IoT: Smart inventory management in warehouses with RFID sensors and automated stock updates.

University Smart Campus System

- Mobile Computing: Campus apps provide schedules, exam results, navigation, and e-wallet services.
- Cloud Computing: Centralized academic management system on the cloud integrates admissions, finance, and library services.
- Big Data: Analyzes student performance trends, library usage, and resource allocation.
- Artificial Intelligence: AI-powered chatbots answer student queries about courses, timetables, or services.
- IoT: Smart ID cards, sensors in libraries/labs, and connected energy-efficient classrooms.

Smart City Traffic Management System

- Mobile Computing: Citizens access traffic updates, parking availability, and alternate routes via a mobile app.
- Cloud Computing: Centralized cloud platform integrates data from different city sensors.
- Big Data: Processes terabytes of data from CCTV, GPS, and sensors to analyze traffic patterns.
- Artificial Intelligence: AI optimizes traffic lights in real-time and predicts congestion hotspots.
- IoT: Roadside sensors, smart cameras, and connected vehicles feed live data.

Agile adaptations require infrastructure:

- **Solid:** Loosely coupled (min dependency) , highly cohesive (max relationship)
- **Scalable:** Ability to handle increased load without performance degradation
- **Modular:** Smaller, independent & reusable components / modules

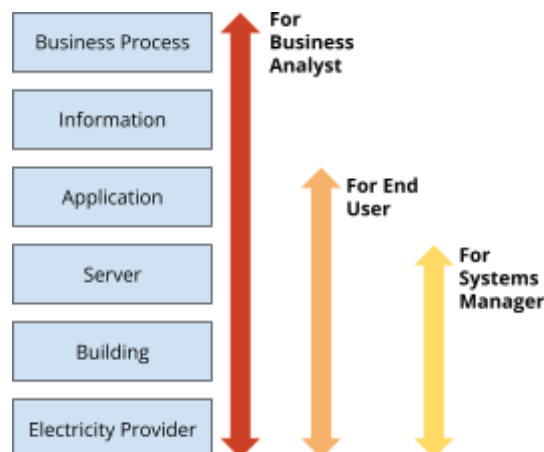
Infrastructure Architecture

Definition

IBM:

- IT infrastructure refers to the combined components needed for the operation and management of enterprise IT services and IT environments.
- Objective: For the operation & management of enterprise IT services and IT environments.

...



Infrastructure comprises 组成 depends on: who you ask and what their point of view

For most people, infrastructure is invisible and taken for granted 是看不见的, 是理所当然的

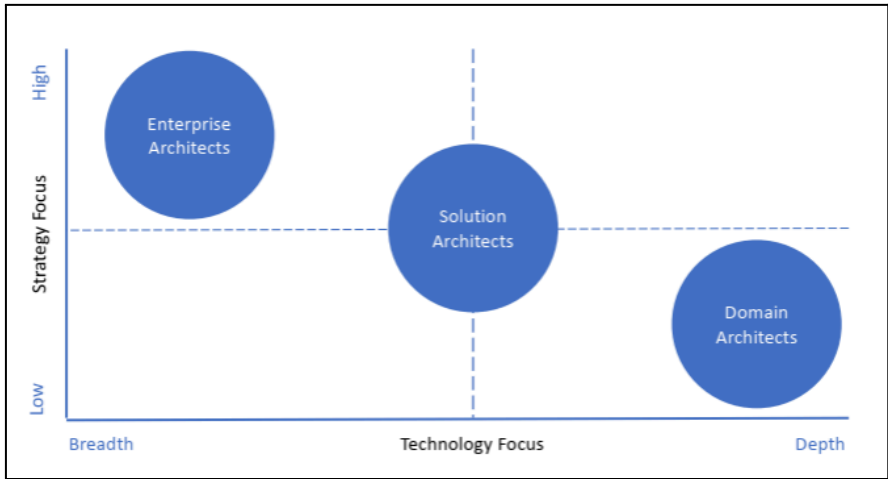
IT Architecture

Architecture is the philosophy that underlies a system and defines its purpose, intent, and structure.

Architecture is crucial to control the infrastructure when it is:

- Designed
- In use
- Changed

Strategy vs technology focus



	Enterprise architects	Solution architects	Domain architects
Strategy focus	High	Moderate	Low
Technology focus	Breadth (wide)	Moderate	Depth
Expert on	Align IT landscape with the business activities	IT solutions, technical & authority of a project	Business / technology topic
Work for	Align the business needs with current & future IT	Project architectural decisions	Infrastructure /software vendors
Assists	CIO & business units	Project manager	Solution architects

IT building blocks & NFA

Functional Management	Process/ information	(Non-Functional Attributes) <ul style="list-style-type: none">• Availability• Performance• Security
Application Management	Application	
Platform Management	Application platform	
Infrastructure Management	Infrastructure <ul style="list-style-type: none">• End User Devices• Operating Systems• Compute• Storage• Networking• Datacentres	

Functional Management	Process/ information <ul style="list-style-type: none"> • Business processes implemented to fulfil company's mission & vision • Business processes create & use information
Application Management	Application <ol style="list-style-type: none"> 1) <u>Usage</u> <ul style="list-style-type: none"> • Single-user application • Multi-user application 2) <u>Source</u> <ul style="list-style-type: none"> • Commercial off-the-shelf (COTS) • Custom software. 3) <u>Architecture</u> <ul style="list-style-type: none"> • Standalone applications • Multi-tier applications (front-end, API & backend (database)) 4) <u>Timeliness</u> <ul style="list-style-type: none"> • Real-time systems: timeliness is critical. • Interactive applications: respond to user actions • Batch-based systems: Regular processing
Platform Management	Application platform <ol style="list-style-type: none"> 1) <u>Application servers</u> <ul style="list-style-type: none"> • A server that hosts applications or software 2) <u>Container platforms</u> <ul style="list-style-type: none"> • A software solution that enables the management of containerized applications 3) <u>Connectivity</u> <ul style="list-style-type: none"> • Application server to database • Application server to container platform • Databases in container 4) <u>Databases</u> <ul style="list-style-type: none"> • Provides ways to store & retrieve data
Infrastructure Management	Infrastructure <ul style="list-style-type: none"> • End User Devices <ul style="list-style-type: none"> ◦ devices used by end users to work with applications ◦ PCs, Laptops, Thin clients, Mobile devices, Printers • Operating Systems <ul style="list-style-type: none"> ◦ collections of programs that manage a computer's internal workings: ◦ Memory, Processors, Devices, File system • Compute <ul style="list-style-type: none"> ◦ physical and virtual computers in the datacenter. Also known as servers • Storage <ul style="list-style-type: none"> ◦ Storage are systems that store data ◦ Hard disks, Tapes 磁带, Direct Attached Storage (DAS), Network Attached Storage (NAS), Storage Area Networks (SANs) • Networking <ul style="list-style-type: none"> ◦ Networking connects all components <ul style="list-style-type: none"> ■ Routers, Switches, Firewalls, WAN, LAN, Internet access, VPNs ◦ Includes infrastructure services <ul style="list-style-type: none"> ■ DNS, DHCP, Time services • Datacentres

	<ul style="list-style-type: none"> locations that host most IT infrastructure hardware <ul style="list-style-type: none"> Uninterruptible power supplies (UPSs) 不间断电源 Heating, Ventilation, and Air Conditioning (HVAC) 暖通空调 Computer racks Physical security measures
--	--

Infrastructure management are processes

- Information Technology Infrastructure Library (ITIL)
- Control Objectives for Information and Related Technology (COBIT)
- DevOps

Tools are used for:

- Monitoring
- Backup
- Logging

Infrastructure building blocks are not per definition hierarchically related!

- For instance, servers need both networking and storage
- Both are equally important

Non-functional attributes (NFRs)

- The name "Non-functional attributes" suggests they have no function 表明它们没有功能
- They are very important for the successful implementation and use of an IT infrastructure
- The term non-functional requirements or NFRs is frequently used and widely known
- The acceptance of a system is largely dependent on the implemented non-functional requirements

Conflicting NFRs

- Many of the non-functional attributes are delivered by the infrastructure
- Non-functional requirements are often conflicting:
 - Security versus user friendliness
 - Performance versus cost
- The infrastructure architect should present stakeholders with these conflicting requirements and their consequences 应该展示这些相互冲突的需求及其后果

Cloud Computing	model that provide on-demand network access to a shared pool of configurable computing resources , users can easily access and use shared IT resources from anywhere, anytime, without managing the infrastructure themselves
Objective	Enables ubiquitous, convenient, on-demand network access to a shared pool of computing resources for rapid provision & release with minimal management effort / service provider interaction. Outsourcing外包模式. To cut cost while focusing on core business
Datacenters	<ul style="list-style-type: none"> On premises On cloud On hybrid mode: premises + cloud
Popular public cloud providers	<ul style="list-style-type: none"> Amazon Web Services (AWS) Microsoft Azure Google Cloud Platform (GCP)

Cloud is new infrastructure

Most organizations will be using on-premises infrastructure本地基础设施 for many years to come

In many cases, there will be a hybrid situation: on-premises + clouds

The cloud is just a number of datacenters that are still filled with hardware: compute, networking and storage.

Cloud computing is a model for enabling ubiquitous无处不在, convenient, **on-demand按需 network access to a shared pool** of configurable computing resources (e.g., networks, servers, storage, applications, and services) that can be **rapidly provisioned快速配置** and **released with minimal management effort** or service provider interaction.

Cloud computing is an outsourcing business model外包业务模式. It enables organizations to cut costs while at the same time focusing on their primary business.

Cloud Characteristics

1. On demand self-service	<ul style="list-style-type: none"> Min systems management effort is needed for deployment End users can configure, deploy, start & stop systems on demand
2. Rapid Elasticity	<ul style="list-style-type: none"> Able to quickly scale-up & scale-down resources
3. Resource Pooling	<ul style="list-style-type: none"> Provides resources from a shared pool (using virtualization technologies)
4. Measured service	<ul style="list-style-type: none"> The actual resource usage is measured and billed There are no capital expenses, only operational expenses
5. Broad network access	<ul style="list-style-type: none"> Capabilities are available over the network

Cloud Deployment Models (4)

Public cloud: delivered by a cloud service provider, accessible through internet

Private cloud: operated solely for a single organization, managed internally

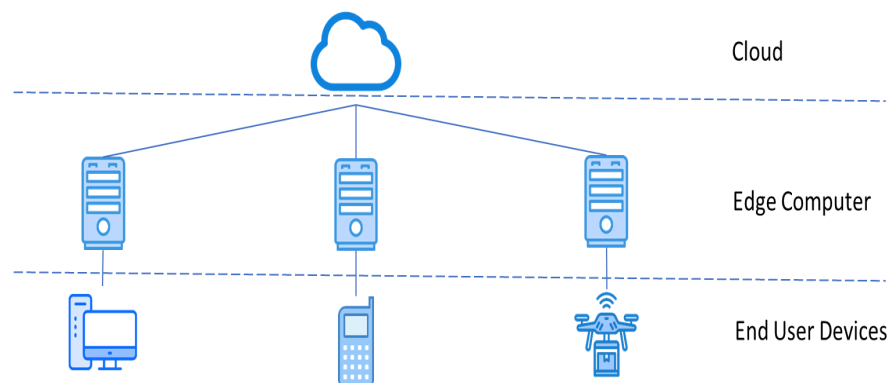
Community cloud: operated by one or more parties in communities with shared concerns,

Hybrid cloud: combine Public + (community / private) cloud; Public: Run generic services (email servers)

Community/ Private: Host specialized services (specific apps)

Cloud Service Models (3)

Software-as-a-Service (SaaS)	Consume	Application Building block	<ul style="list-style-type: none"> Delivers full applications Little / no configuration needed E.g.: Microsoft Office365, LinkedIn
Platform-as-a-Service (PaaS)	Build	Application Platform Building block	<ul style="list-style-type: none"> Delivers a scalable, high available, open programming platform Used by developers to build applications E.g.: Microsoft Azure Cloud Service, Google App Engine
Infrastructure-as-a-Service (IaaS)	Host	Infrastructure Building block(Compute[Operating, Virtualization, Servers], Storage, Networking, Datacentres)	<ul style="list-style-type: none"> Delivers virtual machines, networking, storage Needs to install and maintain the OSs & the layers above that. E.g.: Amazon Elastic Cloud (EC2 and S3) and Microsoft Azure IaaS

**Edge computing**

Brings computing power and data storage closer to where it is needed

	such as at the edge of the network
Objective	<u>Min cloud / on-premises datacenter access</u>
Components	Routers, gateways, switches & sensors
Pros	<ul style="list-style-type: none"> • Low latency <ul style="list-style-type: none"> ◦ Data is processed close to the source (e.g., sensors, devices), reducing the delay to send data/ get response with distant cloud servers • Low bandwidth needs <ul style="list-style-type: none"> ◦ Only necessary/ summarised data sent to the cloud, minimize network usage • Real-time processing <ul style="list-style-type: none"> ◦ data is analyzed locally at the edge, responses and decisions can be made instantly <p># Most processing happens locally at the edge node (closer to the device), so only small or non-urgent data goes to the cloud. Cloud used only for heavy or historical tasks</p>
Application	<u>IoT applications</u> (a large number of devices generate data to be processed in real time)

C2: Non-Functional Attributes

Non-functional attributes	<ul style="list-style-type: none"> To describe the qualitative behaviour of a system For the successful implementation, use & acceptance of an IT infrastructure
A.k.a.	Non-functional requirements or NFRs
Components	<ul style="list-style-type: none"> <u>Availability</u> <u>Performance</u> <u>Security</u>
Conflicts	<ul style="list-style-type: none"> Security vs. user friendliness Performance vs. cost

Availability

Definition	Uptime/ Ability to provide products or services without interruption / downtime.		
Measured by	Expressed as a percentage of uptime (in one year / one month basis)		
Characteristics	Cannot be guaranteed upfront		
Carrier grade availability	<ul style="list-style-type: none">99.999% uptime (for one component)For A full IT system: 99.8% or 99.9% (per month)For the IT infrastructure: 99.99% or higher		
Calculating availability	<ul style="list-style-type: none">a) Mean Time Between Failures (MTBF): Uptime b) Mean time to Repair (MTTR): DowntimeAvailability<ul style="list-style-type: none">Single component: One defect leads to downtime Availability = (MTBF / (MTBF+MTTR)) x 100%Serial components: One defect leads to downtime Total Availability = Availability x AvailabilityParallel components: One defect leads to no downtime But beware of SPOFs Total Availability = 1 - (1 - availability) x (1 - availability)		
	To min MTTR / downtime	<ul style="list-style-type: none">1) Having a service contract with the supplier2) Having spare parts on-site3) Automated redundancy & failover	
	Steps complete repairs to	<ul style="list-style-type: none">1) Notification of the fault (time before seeing an alarm message)2) Processing the alarm3) Finding the root cause of the error4) Looking up repair information5) Getting spare components from storage6) Having technician come to the datacenter with the spare component7) Physically repairing the fault8) Restarting & testing the component	

Sources of unavailability	<div>1) <u>Human errors</u> (accident)</div> <div>2) <u>Software bugs</u> (software complexity)</div> <div>3) Planned <u>maintenance</u> (upgrade, migration, changes)</div> <div>4) Physical <u>defects</u> (mechanical part likely to break first)</div> <div>5) <u>Bathtub curve</u> 类似一开始不会坏, 用着也不会坏, 但是长时间后坏(new component likely to fail)</div> <div>6) <u>Environmental issues</u> (failing facilities, disaster)</div> <div>7) <u>Complexity</u> of the infrastructure (complex system)</div>																																													
To achieve high availability	<div><div>1. Redundancy</div><div><div>○ The duplication of critical components in a single system</div><div>○ E.g. A single component having 2 power supplies</div></div></div> <div><div>2. Failover</div><div><div>○ A (semi) automatic switch-over to a standby system from the same location</div><div>○ E.g. Windows Server failover clustering</div></div></div> <div><div>3. Fallback (3)</div><div><div>○ A manual switchover to an identical standby computer system in a different location</div><div>○ Used for disaster recovery</div></div></div> <div>Solutions (3)</div> <table><tr><td></td><td>Hot site</td><td>Cold site</td><td>Warm site</td></tr><tr><td>Hardware , Power & cooling</td><td>Fully equipped & configured</td><td>Hardware only</td><td>Ready</td></tr><tr><td>Applications</td><td>Installed on servers</td><td>Needs to be installed</td><td>Not be installed</td></tr><tr><td rowspan="2">Data kept</td><td>Up-to-date</td><td>Current data</td><td>Need to restore</td></tr><tr><td>For full mirror</td><td colspan="2">Restored from backup</td></tr><tr><td>Pros</td><td>Accurately mirrors</td><td>Low cost</td><td>Hot + cold site</td></tr><tr><td>Cons</td><td>Requires constant maintenance</td><td>Least preference method</td><td>Testing needed. Slow.</td></tr></table> <div><div>4. Availability in the cloud</div><div><div>○ Regions & availability zones</div><table><tr><td></td><td>Regions</td><td>Availability zone</td></tr><tr><td>Definition</td><td>A <u>geographically defined area</u> that contains multiple Availability Zones.</td><td>A physically <u>isolated location</u> within a Region</td></tr><tr><td>Scope</td><td><u>Broader</u></td><td><u>Localized within region</u></td></tr><tr><td>To maintain availability</td><td><u>Fallback</u></td><td><u>Local failover</u> (using <u>update & fault domains</u>)</td></tr></table></div><div><div>○ Fault & update domains</div><table><tr><td></td><td>Fault domains</td><td>Update domains</td></tr><tr><td>Definition</td><td>A group of virtual machines that share a common power</td><td>A group of virtual machines that can be rebooted during</td></tr></table></div></div>		Hot site	Cold site	Warm site	Hardware , Power & cooling	Fully equipped & configured	Hardware only	Ready	Applications	Installed on servers	Needs to be installed	Not be installed	Data kept	Up-to-date	Current data	Need to restore	For full mirror	Restored from backup		Pros	Accurately mirrors	Low cost	Hot + cold site	Cons	Requires constant maintenance	Least preference method	Testing needed. Slow.		Regions	Availability zone	Definition	A <u>geographically defined area</u> that contains multiple Availability Zones.	A physically <u>isolated location</u> within a Region	Scope	<u>Broader</u>	<u>Localized within region</u>	To maintain availability	<u>Fallback</u>	<u>Local failover</u> (using <u>update & fault domains</u>)		Fault domains	Update domains	Definition	A group of virtual machines that share a common power	A group of virtual machines that can be rebooted during
	Hot site	Cold site	Warm site																																											
Hardware , Power & cooling	Fully equipped & configured	Hardware only	Ready																																											
Applications	Installed on servers	Needs to be installed	Not be installed																																											
Data kept	Up-to-date	Current data	Need to restore																																											
	For full mirror	Restored from backup																																												
Pros	Accurately mirrors	Low cost	Hot + cold site																																											
Cons	Requires constant maintenance	Least preference method	Testing needed. Slow.																																											
	Regions	Availability zone																																												
Definition	A <u>geographically defined area</u> that contains multiple Availability Zones.	A physically <u>isolated location</u> within a Region																																												
Scope	<u>Broader</u>	<u>Localized within region</u>																																												
To maintain availability	<u>Fallback</u>	<u>Local failover</u> (using <u>update & fault domains</u>)																																												
	Fault domains	Update domains																																												
Definition	A group of virtual machines that share a common power	A group of virtual machines that can be rebooted during																																												

	source & network switch	planned maintenance
Protection against	Hardware and infrastructure failures	Downtime during planned maintenance and updates
E.g.	When a fault domain fails (e.g., a rack outage), VMs within that domain are affected, but VMs in other fault domains remain online.	When an update domain is rebooted, it's given a recovery time (e.g., 30 minutes) before the next update domain is affected

5. Business continuity - IT disaster

- An irreparable problem in a datacenter due to unusable
- Types
 - i. Natural disaster - Floods, hurricanes, tornadoes, earthquakes
 - ii. Man-made disaster - Hazardous material spills, infrastructure failure, bio-terrorism
- Impact: Infrastructure become unavailable
- Solutions

1) Business continuity management (BCM)	<ul style="list-style-type: none">● IT● Managing business processes● Availability of people and workplaces in disaster situations
2) Disaster recovery plan (DRP)	<ul style="list-style-type: none">● A set of measures to take if disaster● To accommodate the IT infrastructure in an alternative location

RTO & RPO are critical components of BCM & DRP

RTO	<ul style="list-style-type: none">● Recovery Time Objective (RTO)● The maximum duration to be restored after a disaster● To avoid unacceptable consequences● A <u>shorter RTO</u> implies a faster <u>recovery</u> and often requires <u>more resources</u>.
RPO	<ul style="list-style-type: none">● Recovery Point Objective (RPO)● The point (in time) to which data must be recovered considering some "acceptable loss" in a disaster situation● A <u>lower RPO</u> means <u>less data loss</u> and usually requires <u>more frequent backup</u>

Performance

Definition	Perceived performance refers to how <u>quickly</u> a system appears to perform its task
Indicator	Inform the user about how long a task will take, using <u>progress bars</u> , <u>splash screens</u>

Calculating performance (2)	A. Performance during infrastructure design phase	
	Nature	Complexity, extremely difficult, unreliable
	Considerations	<ul style="list-style-type: none"> • When the system works as expected - normal • When the system is in a special state (e.g.: failing parts, maintenance state, performing backup, and running batch job)
	<p>To evaluation performance:</p> <p>Benchmark 跑分对比</p> <ul style="list-style-type: none"> • Uses a test program to assess the relative performance of an infrastructure component • Scope <ul style="list-style-type: none"> ◦ Performance of various subsystems ◦ Across different system architectures • E.g. <ul style="list-style-type: none"> ◦ Compare the raw speed of parts of an infrastructure (processor) <p>👍 Vendor experience</p> <ul style="list-style-type: none"> • Vendors have experience running their products in various configurations • Vendors can provide: Tools, Figures & Best practices <p>Prototype (aka PoC)</p> <ul style="list-style-type: none"> • To measure the performance of a system at an early stage, for part with highest risk • How? <ul style="list-style-type: none"> ◦ Hiring equipment from suppliers ◦ Using datacenter capacity at a vendor's premise ◦ Using cloud computing resources <p>User profiling 预计负载 在还没正在开发</p> <ul style="list-style-type: none"> • Predict the load a new software system to the infrastructure before the software is actually built to get expected usage of the system • How? <ul style="list-style-type: none"> ◦ Define a number of typical user groups ◦ Create a list of tasks personas will perform on the new system ◦ Decompose tasks to infrastructure actions ◦ Estimate the load per infrastructure action ◦ Calculate the total load <p>How?</p> <p>Scalable cloud environment</p> <ul style="list-style-type: none"> • In cloud environments, it offers rapidly elasticity • Cloud environments have extensive logging and monitoring capabilities <p>B. Performance of a running system(2)</p> <p>1. Manage bottleneck</p> <ol style="list-style-type: none"> A component causing the system to reach limit, that negatively influence performance Every system has at least 1 bottleneck that limits its performance If the bottleneck has no negative impact to performance of the system under the highest expected load, it is OK BASED ON <ul style="list-style-type: none"> The performance of all its components The interoperability互通性 of various components <p>2. Performance test (3)</p> <ol style="list-style-type: none"> Load test <ol style="list-style-type: none"> Shows how a system performs under the expected load Stress test <ol style="list-style-type: none"> Shows how a system reacts when it is under extreme load Endurance test 耐力测试 	

	<ul style="list-style-type: none"> i. Shows how a system behaves when it is used at the expected load for a long period of time d. Breakpoint - How e. Ramp up / increase the load f. Ideal using Cloud environment, due to: Rapidly elasticity, Reduce the cost, Simulating a very large number of users g. Software <ul style="list-style-type: none"> i. One or more servers to act as injectors <ul style="list-style-type: none"> 1. Each emulating a number of users 2. Each running a sequence of interactions ii. A test conductor <ul style="list-style-type: none"> 1. Coordinating tasks 2. Gathering metrics from each of the injectors 3. Collecting performance data for reporting purposes h. Where <ul style="list-style-type: none"> i. Production-like environment (for reliability) ii. Temporary (hired) test environment (for min cost)
Performance patterns (to improve performance)	<p>1. Increase Upper Layer Performance</p> <ul style="list-style-type: none"> • Definition/Working: Improve performance at application & database level using tuning, task prioritization, working on memory rather than disk , make good use of queues and schedulers. • Pros: Low-cost, flexible, faster response. • Cons: Limited impact if hardware is weak. • Suitable for: General systems where quick software optimization is needed. <p>2. Caching</p> <p><u>a. Disk Caching</u></p> <ul style="list-style-type: none"> • Definition: Stores recently read disk data in cache. • Pros: Faster access to common data. • Cons: Cache miss still requires slow disk access. • Suitable for: File systems, operating systems. <p><u>b. Web Proxies</u></p> <ul style="list-style-type: none"> • Definition: Proxy server storing frequently accessed web content. • Pros: <u>Faster browsing, more bandwidth to use (due to no download needed)</u> . • Cons: May serve outdated content if not refreshed. • Suitable for: Organizations with many users accessing similar web content. <p><u>c. Front-End Server (Web Server)</u></p> <ul style="list-style-type: none"> • Definition: Serves static data and acts as reverse proxy(auto-cache most requested data). • Pros: Reduces backend load, caches static content. • Cons: Limited for dynamic data. • Suitable for: Websites with many static resources (images, CSS, JS). <p><u>d. In-Memory Databases</u></p> <ul style="list-style-type: none"> • Definition: Entire DB stored in RAM. • Pros: Extremely fast performance. • Cons: Risk of data loss if power fails; costly RAM. • Suitable for: Real-time financial trading, gaming, analytics. <p><u>e. Edge Servers</u></p> <ul style="list-style-type: none"> • Definition: Cache data close proximity to end-users (edge locations). • Pros: Low latency, reduced backbone traffic. • Cons: Expensive to deploy globally. • Suitable for: Cloud providers, CDNs (Netflix, Cloudflare). <p>3. Operational Data Stores (ODS)</p> <ul style="list-style-type: none"> • Definition: Database <u>integrating real-time operational data from multiple sources</u>.

- Pros: Real-time insights, reduces load on main DB.
- Cons: Extra system to manage, requires integration.
- Suitable for: Businesses needing consolidated, real-time operational reporting.

4. Scalability

a. Vertical Scaling (Scale Up)

- Definition: Add resources (CPU, RAM) to a single server/components.
- Pros: Simple, no software change.
- Cons: Hardware upgrade limit, costly.
- Suitable for: Small/medium systems with moderate growth.

b. Horizontal Scaling (Scale Out)

- Definition: Add more servers/components to infrastructure.
- Pros: High scalability, fault-tolerant.
- Cons: Complexity, requires distributed system design.
- Suitable for: Cloud computing, big data, SaaS platforms.

5. Load Balancing

- Definition: Distributes incoming requests across multiple servers.
- Pros: Improves availability, avoids single-server overload.
- Cons: Extra cost, requires application to support multi-server.
- Suitable for: Web apps, e-commerce, cloud services.
- **Load balance**
 - Spreads the load to available machines:
 - Checks the current load on each server in the farm
 - Sends incoming requests to the least busy server
- **Advanced load balancers**
 - Spread the load based on Server side's:
 - Number of connections
 - Response time

6. High-Performance Computing

a. High Performance Cluster

- Definition: Combines many servers into one powerful computing unit.
- Pros: Massive computing power at relatively low cost (off-the-shelf servers).
- Cons: Needs high-speed interconnects, complex management.
- Suitable for: Scientific research, simulations, weather forecasting.

b. Grid Computing

- Definition: Geographically distributed computing cluster.
- Pros: Utilizes underused resources across locations.
- Cons: Bandwidth limitation, security risks.
- Suitable for: Research collaborations across universities/companies.

7. Design for Use

- Definition: Optimize system based on workload purpose.
- Tips:
 - Match design to purpose (batch, online, real-time).
 - Spread load across systems.
 - Move rarely used data off the main system.
- Pros: Tailored, efficient performance.
- Cons: Requires deep understanding of workload.
- Suitable for: Enterprises with mixed workloads (e.g., ERP + analytics).

8. Capacity Management

- Definition: Monitor and plan resources for long-term high performance.
- Pros: Prevents future bottlenecks, aligns with business growth.

- Cons: Requires continuous monitoring & forecasting.
- Suitable for: Any business-critical IT system.

Security

Definition	<ul style="list-style-type: none"> • Security is the combination of: <ul style="list-style-type: none"> - <u>Availability</u> : Authorised user has reliable access when needed - <u>Confidentiality</u> : Sensitive info is accessible by authorised user only - <u>Integrity</u> : Maintain accuracy & trustworthiness of data 	
Characteristics	Focused on the <u>recognition</u> and <u>resistance</u> of attacks	
Core infrastructure security	<ul style="list-style-type: none"> • Irreversibility of <u>hash keys</u> • Practical unbreakable <u>encryption</u> • Unbreachable <u>virtualization</u> 	
Motive /reason for crime against IT infrastructure	<ul style="list-style-type: none"> • <u>Personal</u> exposure and prestige (visibility / perceiveability by others) • Creating <u>damage</u> • Financial <u>gain</u> • <u>Terrorism</u> • <u>Warfare</u> 	
Cloud security	<ul style="list-style-type: none"> • The <u>public cloud</u> applies a shared responsibility model: <ul style="list-style-type: none"> - The cloud provider takes care of security <u>of</u> the cloud, with many specialists - The customer takes care of security <u>in</u> the cloud 	
Prevention	<ul style="list-style-type: none"> • <u>Design</u> for <u>minimum</u> risk (using source code analysis, standalone system) • Incorporate <u>safety devices</u> (using firewall, hardened screened routers) • Implement <u>training</u> & <u>procedures</u> (to mitigate risk, ensure proper use) 	
Implementation	Zero trust	Assumes no implicit trust regardless of location & requires continuous verification of users and devices before granting access
	Segregation of duties 职责分离	Divides critical tasks among multiple individuals to prevent fraud, errors & unauthorized access
	Least privilege	Granting users only the minimum necessary access rights to perform their job functions
	Privileged Access Management (PAM)	Manage & monitor access to sensitive resources and privileged accounts.
	Layered security	Multiple security controls to protect against different types of threats
	Identity and Access Management (IAM)	Manage user identities & access rights across an organization
	Authentication	Verify the identity of a user or device before granting access
	Password	A secret string of characters used to authenticate a user's identity
	Role Based Access Control (RBAC)	Control access to resources based on user roles
	Cryptography	Secure communication & data by transforming it into an unreadable format, making it incomprehensible to unauthorized

BMIS2113 Information Technology Infrastructure

	Encryption	Encode data using algorithms & keys, converting readable information (plaintext) into an unreadable format (ciphertext)
	Computer Emergency Response Team (CERT)	A team that responds to computer security incidents

C3: Datacenter

3.1 Functions of a Datacenter

Definition: Most IT infrastructure hardware, except for end-user devices, are hosted in datacenters.

Functions A datacenter provides:

- Power supply
- Cooling
- Fire prevention and detection
- Equipment racks

History Early datacenters (computer rooms): Designed and built for large mainframe systems.

Today's Datacenters

- Equipped with standardized 19" racks: House CNS equipment
- Contains shipping containers packed with thousands of servers each (for large scale)

3.2 Datacenter Building Blocks

a) Sub Equipment Room (SER)

- smaller (office building)
- Size: Patch closet, small room, or closet that houses networking/electrical equipment
- Key Function: Acts as a distribution point for network connections

b) Main Equipment Room (MER)

- larger (centralized room)
- Size: A small datacenter located in the organization's subsidiaries or buildings
- Key Function: Houses core network infrastructure

c) Organization Owned Datacenter (OOD)

- larger (purpose-built facility)
- Size: A datacenter that contains all central IT equipment for the organization
- Key Function: Houses the organization's entire IT infrastructure

d) Multi-Tenant Datacenter (MTD)

- largest (facility shared by multiple org)
- Size: Used by service providers to offer services to multiple organizations
- Key Function: Hosts IT equipment for multiple tenants

Factors Determining the Locations of Datacenters

- Environment of the datacenter
- Visibility of the datacenter
- Utilities available to the datacenter
- Datacenters located in foreign countries

Physical structure: FWWDP

Floors

1. Datacenter Floor Load
 - Can carry 1500 to 2000 kg/m²
 - A fully filled 19" computer rack weighs approximately 700 kg
 - Rack footprint (60x100 cm) results in a load of 1166 kg/m²
2. Office Floor Load
 - Can carry approximately 500 kg/m²
3. Raised Floor
 - Constructed with metal framework and removable tiles (60x60 cm)
 - Disadvantages:
 - Expensive

- Reduces total available height in the datacenter
- Limits maximum floor load
- Difficult to install doors and equipment loading slopes
- Fire can easily spread through
- **Overhead cable trays** are a safer alternative

Vents

- provide cool airflow to racks on the floor

Walls, Windows, Doors

- Walls
 - Must extend from floor to ceiling
 - Require adequate fire rating to act as physical firewalls
- Windows
 - Not desirable in datacenters
 - If present, must be:
 - Translucent
 - Shatterproof
 - Impossible to open
- Doors
 - Must be large enough for equipment
 - Should resist forced entry

Water and Gas Pipes

- Leakage from ceiling-mounted water pipes can damage equipment
- Datacenter operators must know the location of shutoff valves

Layout of Datacenter – Rooms and Components

1. Computer Room
 - Location where actual IT infrastructure components are installed
2. UPS or Generator
 - Diesel generator provides backup electrical power
 - Fuel should be stored outside or in an isolated, secure room nearby
3. Input Power Transformers
 - Transformers from the power utility company
4. UPS (Uninterruptible Power Supply)
 - Provides continuous power during short outages
5. UPS Batteries
 - Battery set for short-term power supply
6. Cooling
 - Systems for maintaining optimal temperature
7. Fire Extinction
 - Fire suppression systems
8. Operator Room
 - Includes a large window facing the computer room for monitoring activity
9. Storage Room
 - Stores spare hardware and equipment
10. Entrance
 - Access point to other rooms; no windows
11. Meeting Room

- Used for staff and visitor meetings
- Includes a secured (shatterproof) window for sunlight

Power Density and Server Rack

Power Density (kilowatts per m²):

- Power drawn:
 - 1 rack of servers: measured in kilowatts (kW)
 - Large facilities: measured in megawatts (MW)
- Normal-density datacenter: 2–6 kW/m²
- High-density datacenter: 10–20 kW/m² (racks filled with 40–80 servers)

Server Rack Considerations:

- Cannot be fully equipped due to:
 - Stability and transport
 - Power and cooling limitations
 - Accessibility and maintenance needs
 - Future expansion
 - Cable management

Uninterruptible Power Supply (UPS) Systems

Objective

- Prevent power issues that lead to downtime and equipment damage

Characteristics

- Operates independently of utility power
- Delivers high-quality electrical power

Installation Components

- Filters
- Diesel power generator
- Battery set or flywheel system

Types of UPS Systems

- Battery Powered UPSs
 - Also called Standby UPS or Off-line systems
 - Suitable for small setups (few workstations or servers)
- Line Interactive UPS Systems
 - Use a transformer between utility power and IT equipment
 - Filter many power issues
- Double Conversion UPS Systems
 - Convert AC to DC and back to high-quality AC
 - Ensures stable and clean power supply

Power Distribution Units (PDUs)

Definition

- A device with multiple outlets that distributes power to datacenter equipment

Types

- Floor PDUs:
 - Used in large datacenters
 - Distribute power to multiple racks or equipment
- Rack PDUs:
 - Mounted in standard 19-inch racks
 - Offer various power distribution options

High Availability Strategy

- Equip infrastructure with two power supplies for redundancy
- Use at least two power strips per rack for reliable power delivery

Cooling Systems in Datacenters

Objective

- To dissipate heat generated by IT equipment

Types of Cooling Systems

- CRAC (Computer Room Air Conditioners)
 - Use refrigerants
 - Connected to external condensing units
- CRAH (Computer Room Air Handlers)
 - Use chilled water
 - Connected to outside chillers
 - Chillers produce chilled water via refrigeration

Efficiency Metrics

- EER (Energy Efficiency Ratio)
 - Measures efficiency at maximum load
 - Ratio: BTU/hour output cooling ÷ Watts input
- SEER (Seasonal Energy Efficiency Ratio)
 - Similar to EER but based on seasonal usage (typically summer)
- COP (Coefficient of Performance)
 - Ratio: Cooling load in kW ÷ Electric energy input in kW
 - Typical values: 3 to 10

Operating Temperature

- Normal range: 18–27 °C
- Server shutdown threshold: 40 °C at air inlet
- Raising temperature by 1 °C reduces cooling cost by 5%

Environmental Guidelines for Datacenter Operation

Operating Temperature

- Maintain air temperature between 18–27 °C
- Servers shut down at 40 °C inlet temperature
- Increasing temperature by 1 °C saves 5% in cooling costs

Airflow

- Optimize airflow to eliminate hot spots
- Avoid excessive cooling of the entire datacenter

Liquid Cooling

- Suitable for large datacenters
- Components immersed in non-conductive, non-corrosive fluid
- Allows closer placement of system boards, increasing CPU density

Humidity and Dust Control

- Maintain humidity between 40% and 60%
- Minimize dust by:
 - Restricting visitor access
 - Wearing dust-free clothing and protective shoe sleeves

Fire Prevention, Detection, and Suppression

- Fires may be caused by short circuits or defective equipment
- Airflow and raised floors can accelerate fire spread
- Smoke can damage equipment

Four Levels of Fire Safety:

1. Fire Protection – Prevent fire occurrence
2. Passive Fire Protection – Limit fire exposure once started
3. Fire Detection – Identify smoke and fire
4. Fire Suppression – Extinguish fire after detection

Equipment Racks

- A 19-inch rack is a standardized metal enclosure for IT components
- Rack height is measured in rack units (U)
 - 1U = 44.5 mm
 - Typical rack height = 42U

Datacenter Energy Efficiency

- Energy cost often exceeds server cost
- Power Usage Effectiveness (PUE) measures total datacenter power usage:
 - $PUE = \text{Total datacenter power} \div \text{IT equipment power}$
 - Typical PUE values range from 1.1 to 2.0

Example:

- A datacenter with PUE = 1.5 means:
 - 1 watt used by IT equipment
 - 0.5 watt used by the rest of the datacenter (cooling, lighting, etc.)

3.3 Datacenter APS

Datacenter Availability- 4 Tiers

Availability tiers describe the reliability of datacenter facilities, not the IT infrastructure components. They are based on uptime, redundancy, and fault tolerance.

Tier	Infrastructure Description	Expected Uptime
Tier 1	Basic infrastructure	99.671%
Tier 2	Adds some redundancy and backup paths for power and cooling	99.741%
Tier 3	Multiple redundant paths for power and cooling; allows maintenance without downtime	99.982%
Tier 4	Fully fault-tolerant; redundant systems for power, cooling, and all IT equipment	99.995%

Redundant Datacenters

- Used to increase availability beyond Tier 4 (99.995%)
- Multiple datacenters are required for ultra-high availability
- Redundant datacenters should be located at least 5 km apart

Floor Management (by Floor Manager)

- Minimize personnel walking around server racks
- Keep the datacenter floor tidy
- Change backup tapes
- Provide power connections to racks
- Maintain and test fire extinguishing and UPS systems

Datacenter Performance

- Datacenter does not directly affect IT infrastructure performance
- Performance contributions are limited to:
 - Bandwidth of internet connectivity
 - Scalability of the location

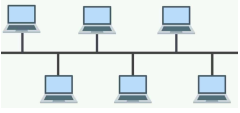
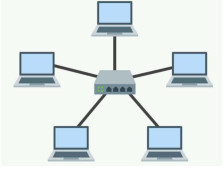
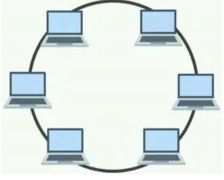
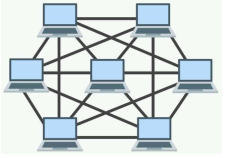
Datacenter Security

Physical Security Measures:

- Restrict physical access to selected, qualified staff
- Use entry registration systems and maintain access logs
- Secure doors with conventional or electronic locks (with authentication)
- Entry points may include:
 - Regular doors
 - Mantraps
 - Revolving doors
 - Weighing scales to ensure only one person enters at a time

C4: Networking

Network Topologies

	Bus	Star	Ring	Mesh
				
Connection by	Single cable / bus	Central hub / switch	A circular fashion	By connected to multiple devices
Data transmission	Along the bus, to intended recipient	Through the central point	In 1 direction around the ring	Routed through different paths
Pros	Simple & inexpensive to implement	Easy to manage and troubleshoot	Simple & efficient	Fault-tolerant
Cons	Prone to performance issues	Central hub as a SPOF	Disruptive to add / remove devices	Expensive to implement

Networking Building Blocks

- OSI Reference Model (OSI-RM) defines the different stages that data must go through to travel from one host to another over a network

1 Physical Layer

<ul style="list-style-type: none"> • Deal with the physical components that carry the data • Data transmission media and connectors • Signal characteristics • Network topology 		
Cables	Twisted pair cables	<p><u>Unshielded Twisted Pair (UTP)</u></p> <ul style="list-style-type: none"> • Consists of twisted wire pairs without any additional shielding • Used in Ethernet networks and for telephone systems • Pros: <ul style="list-style-type: none"> ◦ Cost-effective ◦ Easy to install ◦ More flexible • Cons: <ul style="list-style-type: none"> ◦ Susceptible to interference ◦ Lower data rate <p><u>Shielded Twisted Pair (STP)</u></p> <ul style="list-style-type: none"> • Suitable for most home and office networks where interference isn't a major concern • Used in factories, datacenters, or areas with high electrical noise • Pros: <ul style="list-style-type: none"> ◦ Reduce interference ◦ Higher data rates • Cons: <ul style="list-style-type: none"> ◦ More expensive ◦ More complex installation ◦ Less flexible

	Fiber Optic Cables	Multimode <ul style="list-style-type: none">Core diameter: 50 or 62.5 microns (larger)Light source: LEDs (multiple paths)Used in LAN, datacenter, short distanceCost: Less expensive Single-mode <ul style="list-style-type: none">Core diameter: 9 microns (smaller)Light source: Laser (single path)Used in telecommunication, long distanceCost: More expensive				
	Coaxial cables	For TV & Internet connection, outside Malaysia				
Patch Panel	Function	Provides a centralized location to organize and manage network cables				
	Pros	<ul style="list-style-type: none">Organise cablesEase changesProtect connectionsProvide clear reference for labelling & documentation				
Cabling	Vertical Cabling	<table><tr><th>Vertical Cabling</th><th>Horizontal Cabling</th></tr><tr><td><ul style="list-style-type: none">connects different floors or buildings within a network infrastructureConnection: Patch panels to data centerLong distance</td><td><ul style="list-style-type: none">Connects devices within a floorEndpoints in the walls to patch panelsUp to 90 meters</td></tr></table> <p>Pros</p> <ul style="list-style-type: none">Enhance network reliability: fault isolation, increase uptime, disaster recoveryImprove network performance: Load balancing, reduce latencySimplify network management: Centralized management, ease troubleshooting <p>Connection</p> <ul style="list-style-type: none">The main distribution cabling connects the patch panels on the floors to the datacenterConnects floors and buildings <p>Cable Types</p> <ul style="list-style-type: none">Fiber Optic (primary)Twisted pairCoaxial	Vertical Cabling	Horizontal Cabling	<ul style="list-style-type: none">connects different floors or buildings within a network infrastructureConnection: Patch panels to data centerLong distance	<ul style="list-style-type: none">Connects devices within a floorEndpoints in the walls to patch panelsUp to 90 meters
	Vertical Cabling	Horizontal Cabling				
<ul style="list-style-type: none">connects different floors or buildings within a network infrastructureConnection: Patch panels to data centerLong distance	<ul style="list-style-type: none">Connects devices within a floorEndpoints in the walls to patch panelsUp to 90 meters					
Horizontal Cabling	<p>Pros</p> <ul style="list-style-type: none">Same with vertical cabling <p>Connection</p> <ul style="list-style-type: none">Endpoints in the walls are connected to the patch panelsWithin a specific floor / building <p>Cable Types</p> <ul style="list-style-type: none">Twisted pair (primary)Fiber optic					

Leased lines	Definition	<ul style="list-style-type: none">• A dedicated data connections between two locations• Provided by a telecom provider									
	Types	<p><u>T or E lease line</u></p> <ul style="list-style-type: none">• Definition: Not standard industry terms for leased lines• Pros:<ul style="list-style-type: none">◦ Low data rate: T (1.544 Mbps), E (2.048 Mbps)• Used in<ul style="list-style-type: none">◦ T (USA, Canada and Japan)◦ E (most other countries) <p><u>Synchronous Optical Network (SONET)</u></p> <ul style="list-style-type: none">• Definition: Transmit large amounts of data using fiber optic• Pros:<ul style="list-style-type: none">◦ High data rate◦ Low latency• Used in:<ul style="list-style-type: none">◦ Connecting corporate offices & datacenters◦ Enabling high-speed data replication & backup◦ Supporting real-time applications with low latency requirements <p><u>Synchronous Digital Hierarchy (SDH)</u></p> <ul style="list-style-type: none">• Definition: The international standard equivalent to SONET• Pros:<ul style="list-style-type: none">◦ High reliability◦ Scalability• Used in:<ul style="list-style-type: none">◦ Telecommunications to transmit multiple digital data stream over a fiber optic <p>Dark Fiber</p> <table><tr><th>Dark Fiber</th><th>Regular Fiber (Lit Fiber)</th></tr><tr><td>Not actively used by the provider Installed but unused fiber optic cables</td><td>Actively used by the provider to deliver services</td></tr><tr><td>Fully controlled by the customer(lessee)</td><td>Controlled and managed by the provider</td></tr><tr><td>Lessee installs and manages own equipment</td><td>Provider supplies and manages equipment</td></tr><tr><td>Private, point-to-point network</td><td>Shared infrastructure managed by ISP</td></tr></table> <ul style="list-style-type: none">• Definition: Unused fiber optic cables that have been installed but are not actively transmitting data• Pros:<ul style="list-style-type: none">◦ High bandwidth, Low latency◦ Security, Scalability, Reliability• Used in:<ul style="list-style-type: none">◦ To accommodate future growth◦ Lessee then installs their own equipment to "light up" the fiber and create their own private network	Dark Fiber	Regular Fiber (Lit Fiber)	Not actively used by the provider Installed but unused fiber optic cables	Actively used by the provider to deliver services	Fully controlled by the customer(lessee)	Controlled and managed by the provider	Lessee installs and manages own equipment	Provider supplies and manages equipment	Private, point-to-point network
Dark Fiber	Regular Fiber (Lit Fiber)										
Not actively used by the provider Installed but unused fiber optic cables	Actively used by the provider to deliver services										
Fully controlled by the customer(lessee)	Controlled and managed by the provider										
Lessee installs and manages own equipment	Provider supplies and manages equipment										
Private, point-to-point network	Shared infrastructure managed by ISP										
Internet Access	DIA	<p>Characteristics</p> <ul style="list-style-type: none">• High speed• Reliable connectivity• Guaranteed bandwidth <p>Connection</p> <ul style="list-style-type: none">• Dedicated cable connection to single business <p>Ideal for Business with critical applications</p>									

	Broadband	Characteristics <ul style="list-style-type: none"> • High-speed over cable, Digital Subscriber Line (DSL) or fiber optic networks Connection <ul style="list-style-type: none"> • Shared among multiple users Ideal for Internet browsing, Email
	WiFi	Characteristics <ul style="list-style-type: none"> • Includes Wi-Fi, cellular data and satellite Internet Connection <ul style="list-style-type: none"> • Shared connection without physical cable Ideal for Wirelessly

2 Data Link Layer

- Ensure reliable & error-free data transfer between devices
- Takes packets received & breaks them into frames
- Detect errors in the frames
- Manage how devices access the physical medium to prevent collisions

Network

	PAN	LAN	MAN	WAN
Distance	Within a few meters	Home, office building or school	A city or a large campus	Large geographical area
Used in	Connecting personal devices within a few meters	Setting up a network in a home office, or campus for fast local access	Linking multiple LANs across a city or large institution	Connecting networks across cities, countries or globally
Protocols	Bluetooth (L2CAP), Zigbee, Infrared	Ethernet (CSMA/CD), Wi-Fi (CSMA/CA)	Metro Ethernet, DQDB	PPP, HDLC, Frame Relay

Common Protocols Used

	Ethernet	WiFi	PPP
Technology	Wired	Wireless	Direct connection between two devices
Used for	Connecting devices in a LAN	Connecting to the Internet wirelessly	VPNs, dial-up & Internet connection

Implementation

- Ethernet
- Switching: Split a single network segment into multiple segments, for each device
- Public wireless networks:

	Features	Used in
1G	Limited data speeds and security features	(Analog) Support voice calls
2G	Improved call quality and security	(Digital) Improved call quality and security
3G	Increased data speeds, allowing for mobile internet access	Support email, web browsing

4G	Further speed improvements	Streaming videos and online gaming
5G	Faster speeds, lower latency (delay) and greater capacity	Internet of Things (IoT), virtual reality and self-driving cars

3 Network Layer

- Define the route the data is sent to the recipient
- Implementations:
 - IPv4, IPv6, Routing and addressing

IP Protocol

A unique numerical label assigned to each device connected to a network		
Server can have multiple IP addresses, both public and private		
<ul style="list-style-type: none"> • Load balancing • Hosting multiple websites • Providing different services on the same server 		
	Public IP Address	Private IP Address
Visible to	Entire network	Within a local server within a private network
Used for	<ul style="list-style-type: none"> • Web server • Email server • Public-facing applications 	<ul style="list-style-type: none"> • Internal communication

Addressing

	Static IP Address	Dynamic IP Address
Nature	Constant	Change periodically
Ideal for	Public servers	Internal servers

Routing

<ul style="list-style-type: none"> • Routers compile routing tables to make IP packet forwarding decisions • Routing in the context of a server, covers: <ul style="list-style-type: none"> ◦ Network Routing <ul style="list-style-type: none"> ■ Direct data packets are from the server to the other network devices ◦ Application-Level Routing <ul style="list-style-type: none"> ■ Handle incoming requests within the server itself: load balancing, Content Delivery Networks (CDNs), application-specific routing 	
---	--

4 Transport Layer

- Maintain flow control
- Provide error checking
- Recovery of data between network devices

Protocols Used

	Transmission Control Protocol (TCP)	User Datagram Protocol (UDP)
Function	Provides reliable delivery of a stream of data between applications	Reduced latency over reliability by sending data without checking if the data arrived

Used in	FTP Web browsing Email	Live stream Online games VoIP
Used for	Large DNS queries	Small DNS queries

DNS - Domain Name System

- Port 53 is specifically for DNS. Other services should not use this port.
- Firewalls should allow traffic on port 53 for DNS to function correctly

5 Session Layer

- Provides mechanisms for opening, closing and managing a session between end-user application processes

Virtual Private Network (VPN)

- Uses a public network to interconnect private sites in a secure way, a.k.a. VPN tunnel
- Uses "virtual" connections based on IPsec / SSL
- Most network providers also offer private VPNs based on MPLS
- VPNs use strong encryption and strong user authentication. Using the Internet for transmitting sensitive data is considered safe
- VPN tunnels are often used for remote access to the LAN by users outside of the organization's premises
- Most common VPN communication protocol standards:

	Point-to-Point Tunneling Protocol (PPTP)	Layer 2 Tunneling Protocol (L2TP)	IPsec
Uses	For individual client to server connections	For individual client to server connections	For network-to-network connectivity. IPsec is built into IPv6 standard and is implemented as an add-on to IPv4
Security	Least	No encryption and relies on IPsec for encryption and authentication	Provides encryption and authentication for network traffic
Used in	Legacy system or when speed is prioritized over security	Paired with IPsec for a more secure VPN connection	VPNs and other secure network applications

6 Presentation Layer

- Takes the data provided by the application layer and converts it into a standard format that the other layers can understand

	Secure Socket Layer (SSL)	Transport Layer Security (TLS)
Nature	SSL is considered insecure and should not be used	TLS is securing WWW traffic carried by HTTP to form HTTPS

7 Application Layer

- Interacts with the OS or application

Roles of application layer in servers

- Protocol implementation
- Data formatting
- Error handling

- Security
- Session management

Protocols used

Protocol	Definition	Function
Domain Name System (DNS)	DNS is a distributed database that links IP addresses with domain names DNS was not designed with security in mind Updates to DNS records are done in non-encrypted clear text Authorization is based on IP addresses only	Translates human-readable domain names (like google.com) into machine-readable addresses
DNS Security Extensions (DNSSEC)	Provides origin authentication of DNS data for data integrity	Verifies the authenticity & integrity of DNS data
IP Address Management (IPAM)	IPAM systems are appliances that can be used to plan, track and manage IP addresses in a network IPAM systems integrate DNS, DHCP and IP address administration in one high available redundant set of appliances	Plans, tracks and manages the IP address space within a network
Network Time Protocol (NTP)	NTP ensures all infrastructure components use the same time in their real-time clocks Particularly important for: Log file analysis, clustering software, Kerberos authentication	Synchronization clocks of computers in a network to a common time source
Post Office Protocol (POP)	Used by email clients to retrieve email messages from a mail server	Retrieves emails from a mail server
Simple Mail Transfer Protocol (SMTP)	Used to send email messages from a mail client to a mail server or between mail servers	Sends emails from a mail client to a mail server or between mail server
Multipurpose Internet Mail Extensions (MIME)	Enables SMTP to support file attachments in email messages	Transmits non-text data (like images, audio, video) in email
File Transfer Protocol (FTP)	A protocol for transferring files between computers	Transfers files between computers over a network
Hypertext Transfer Protocol (HTTP)	Defines how messages are formatted and transmitted, and how web servers and browsers should respond to various commands	Transfers hypertext (HTML, CSS, JavaScript, images, etc)
Hypertext Transfer Protocol Secure (HTTPS)	Used when browsing the web with a web browser	Protected data during transfer

* Example refers to images

4.2 Networking Virtualization

Approaches	Used for / as	How
Virtual LAN (VLAN)	Logical grouping (for network segmentation)	Logically divides a single physical network into multiple broadcast domains
Virtual Extensible LAN (VXLAN)	Network virtualization technology	Uses encapsulation to create virtual networks that can span across physical networks, allowing for greater scalability and flexibility compared to VLANs
Virtual Routing and Forwarding (VFR)	Network routing technology (for segmentation)	Hosts multiple independent routing tables in a single physical router
Virtual Network Interface Controllers (VNIC)	Software-based representation of a network interface	Enables virtual machines to connect to the network and communicate with other virtual machines or physical devices
Virtual Switch (VS)	Software-based equivalents of physical network switches	Manages network traffic within a virtualized environment, providing functionalities like VLAN tagging, traffic shaping and connection to physical networks
Software Defined Networking (SDN)	Software-based controllers (Abstract control plan from data plane)	Manages network traffic and resources, enabling dynamic and programmable network configurations
Network Function Virtualization (NFV)	Network architecture (virtualize network function)	Replaces traditional dedicated network hardware appliances with virtualized software instances running on commodity servers

* Example refers to images

4.3 Networking: Availability, Performance & Security

Networking Availability

Layered network topology

Definition	Divides a network into multiple layers, each with specific functions and responsibilities
A.k.a.	Tiered or hierarchical topology
Components	<p>Core Layer</p> <ul style="list-style-type: none"> • Definition <ul style="list-style-type: none"> ◦ The center of the network ◦ The backbone of the network ◦ It typically uses high-capacity routers & switches • Nature <ul style="list-style-type: none"> ◦ High-speed network backbone between network segments • Function <ul style="list-style-type: none"> ◦ Offer rapid & reliable transfer for large volume of data <p>Distribution / Aggregation Layer</p> <ul style="list-style-type: none"> • Definition <ul style="list-style-type: none"> ◦ It combines the access layer data and sends its combined data to one or two ports on the core switches • Nature <ul style="list-style-type: none"> ◦ Sits between the core (in datacenter) & access layers (in patch closet) • Function <ul style="list-style-type: none"> ◦ Performs routing functions, implements security policies, access and core layers <p>Access Layer</p> <ul style="list-style-type: none"> • Definition <ul style="list-style-type: none"> ◦ For servers, located at server racks / in blade enclosures ◦ For workstations, placed in patch closets • Nature <ul style="list-style-type: none"> ◦ Connect workstations & servers to the distribution layer • Function <ul style="list-style-type: none"> ◦ Connect end-users and devices to the network
Benefits	<ul style="list-style-type: none"> • Improve availability & performance (with multiple paths to any piece of equipment) • Provides scalability • Provides deterministic routing: advance determination of the routes between given pairs of nodes • Avoids unmanaged ad-hoc data streams

* Example refers to images

Spines and Leaf topology

Definition	A modern datacenter network architecture designed to address the limitations of traditional three-tier hierarchical designs
Natures	<ul style="list-style-type: none"> • The spine switches are not interconnected • Each leaf switch is connected to all spine switches • Each server is connected to two leaf switches • The connections between spine and leaf switches typically have 10 times the bandwidth of the connectivity between the leaf switches and the servers
Key components	Spine Switches

	<ul style="list-style-type: none"> • Function <ul style="list-style-type: none"> ◦ Form the core of the network • Nature <ul style="list-style-type: none"> ◦ Interconnected in a full-mesh topology, providing redundancy & load balancing • Working (Retail Industry) <ul style="list-style-type: none"> ◦ The spine layer serves as the core of the network, connecting all leaf switches <p>Leaf Switches</p> <ul style="list-style-type: none"> • Function <ul style="list-style-type: none"> ◦ Connect to servers and other end devices • Nature <ul style="list-style-type: none"> ◦ Each leaf switch connects to all spine switches, creating multiple paths for traffic • Working (Retail Industry) <ul style="list-style-type: none"> ◦ Leaf switches act as access points, connecting servers, point-of-sale (POS) systems, inventory management systems, and other devices to the network
Working	A simple physical network is used that can be programmed to act as a complex virtual network. Such a network can be organized in a spine and leaf topology
Benefits	<ul style="list-style-type: none"> • Highly scalable: There are no interconnects between the spine switches • Simple to scale: Just add spine or leaf servers • Physical servers can be connected using relatively few switches • Predictable latency: Each server is always exactly four hops away from every other server

* Example refers to images

Network Teaming

A.k.a.	Link aggregation, port trunking, network bonding
Definition	A method of combining multiple network interface cards (NICs) into a single logical interface
Objective	<ul style="list-style-type: none"> • Provides a virtual network connections using multiple physical cables • To achieve high availability and increased bandwidth • To improve network performance and reliability
Working	<p>Bonds physical NICs together to form a logical network team:</p> <ul style="list-style-type: none"> • Sends traffic to the team's destination to all NICs in the team • Allows a single NIC, cable or switch to be unavailable without interrupting traffic

* Example refers to images

Spanning Tree Protocol (STP)

Definition	An Ethernet level protocol that runs on switches
Working	<ul style="list-style-type: none"> • Guarantees only one path is active between two network endpoints at any given time • Redundant paths are automatically activated when the active path experiences problems • Ensures no loops are created when redundant paths are available in the network

Pros	<ul style="list-style-type: none"> • Shortest Path Bridging (SPB) allows all paths to be active simultaneously, enables much larger topologies, supports faster convergence times • Improves efficiency by allowing traffic to be load balanced across all paths. While STP can take 30 to 60 seconds to respond to a topology change, SPB can respond to changes in less than a second
Cons	It is not using half of the network links in a network, since it blocks redundant paths

* Example refers to images

Multihoming

Definition	Connecting a network to two different Internet Service Providers (ISPs)
Methods	<ul style="list-style-type: none"> • Single router with dual links to a single ISP • Single router with dual links to two separate ISPs • Dual routers each with its own link to a single ISP • Dual routers each with its own link to a separate ISP
Pros	<ul style="list-style-type: none"> • Improve availability
Cons	It is not always guaranteed that multiple network paths actually run on a different set of cables. Cables are used by multiple carrier providers.

Networking Performance

Factors affect the speed of a connection

- Throughput and bandwidth

Definition	Amount of data that is transferred through the network during a specific time interval
Constraint	Throughput is limited by the available bandwidth
Working	<p>When an application requires more throughput than a network connection can deliver:</p> <ul style="list-style-type: none"> • Queues in the network components temporarily buffer data • Buffered data is sent as soon as the network connection is free again • When more data arrives than the queries can store in the buffer, packet loss occurs

- Latency

Definition	The time from the start of packet transmission to the start of packet reception
Depend on	<ul style="list-style-type: none"> • The physical distance a packet has to travel • The number of switches and routers the packet has to pass
Rules	<ul style="list-style-type: none"> • 6 ms latency per 100 km • WANs: Each switch in the path adds 10 ms to the one-way delay • LANs: Add 1 ms for each switch
Types	<p>One way latency</p> <ul style="list-style-type: none"> • The time from the source sending a packet to the destination receiving it <p>Round-trip latency</p> <ul style="list-style-type: none"> • The one-way latency from source to destination plus the one-way latency from the destination back to the source <p>A "ping" can be used to measure round-trip latency</p>

- Quality of Service (QoS)

Definition	Ability to provide different data flow priority to different applications, users or types of data
Nature	Allows better service to certain important data flows compared to less important data flows
Used for	Real-time applications like video and audio streams and VoIP telephony
Implementations methods	<p>Congestion Management</p> <ul style="list-style-type: none"> • To prioritize traffic • Defines action when the amount of data to be sent exceeds the bandwidth of the network • Packets can either be dropped or queued <p>Queue Management</p> <ul style="list-style-type: none"> • To make the wait time more manageable and transparent • Defines criteria for dropping packets that are of lower priority before dropping higher priority packets • When queue are full, packets will be dropped <p>Link efficiency</p> <ul style="list-style-type: none"> • To ensure efficient use of available resources • Ensures the link is used in an optimized way • By fragmenting large packets with low QoS, allowing packets with a high QoS to be sent between the fragments of low QoS packets <p>Traffic shaping</p> <ul style="list-style-type: none"> • To prioritize certain types of traffic over others • Limit the full bandwidth of streams with a low QoS • Have QoS streams have a reserved amount of bandwidth

- WAN Link Compression

Definition	Data compression reduces data size before it is transmitted
Pros	<p>WAN acceleration appliances:</p> <ul style="list-style-type: none"> • Provide compression • Perform some caching of regularly used data at remote data

* Example refers to images

Networking Security

Network Encryption

Definition	Encryption is often a feature in the datacenter
Types	<ul style="list-style-type: none"> • Encrypting data in transit: Encrypting data on the network • Encrypting data at rest: Encrypting data in the storage • End-to-end encryption: Encrypting data between 2 end-points, with network traffic encryption

Firewalls

Definition	<ul style="list-style-type: none"> • Firewalls separate 2 / more LAN / WAN segments for security reasons • Firewalls block all unpermitted network traffic between network segments • Permitted traffic must be enabled by configuring the firewall to allow it
------------	--

Implementation	<ul style="list-style-type: none"> • In hardware appliances • As an application on physical servers • In virtual machines
Host based firewall	<ul style="list-style-type: none"> • Protect a server or end user computer against network based attacks • Part of the operating system
Traffic control methods	<p>Packet Filtering</p> <ul style="list-style-type: none"> • Data packets are analyzed using preconfigured filters • This function is always available on routers & most OS <p>Proxy</p> <ul style="list-style-type: none"> • A proxy terminates the session on the application level on behalf of the server (proxy) or the client (reverse proxy) and creates a new sessions to the client or server <p>Stateful Inspection</p> <ul style="list-style-type: none"> • Inspects the placement of each individual packet within a packet stream • Maintains records of all connections passing through the firewall and determines whether a packet is the start of a new connection, part of an existing connection, or is an invalid packet

Intrusion Detection System (IDS) or Intrusion Prevention System (IPS)

Definition	Detects & prevents activities that compromises system security (a hacking attempt)
Working	<ul style="list-style-type: none"> • Monitors for suspicious activity • Alerts the systems manager when these activities are detected • Stop attacks by changing firewall rules on the fly
Types	<p>Network-based IDS (NIDS)</p> <ul style="list-style-type: none"> • Placed at a strategic point in the network • Monitors traffic to and from all devices on that network • The NIDS is not part of the network flows, but just "looks at it", to avoid detection of the NIDS by hackers <p>Host-based IDS (HTDS)</p> <ul style="list-style-type: none"> • Runs on individual servers or network devices • It monitors the network traffic of that device • It also monitor user behavior and the alteration of critical (system) files

Example (Retail)	It can detect unusual login attempts, large data transfers, or traffic originating from known malicious sources.
-------------------------	--

* Example

De-Militarized Zone (DMZ)

Definition	A DMZ is a network that serves as a buffer between a secure protected internal network and the inaccurate internet
------------	--

Example (Retail)	Public-facing servers (web servers, email servers, and payment gateways) are placed in the DMZ, while sensitive data (customer databases) is kept on the internal network behind a firewall. This prevents direct access to the internal network if a server in the DMZ is compromised
-------------------------	--

* Example

Remote Authentication Dial In User Service (RADIUS)

Definition	RADIUS is a networking protocol that provides centralized user and authorization management for network devices
------------	---

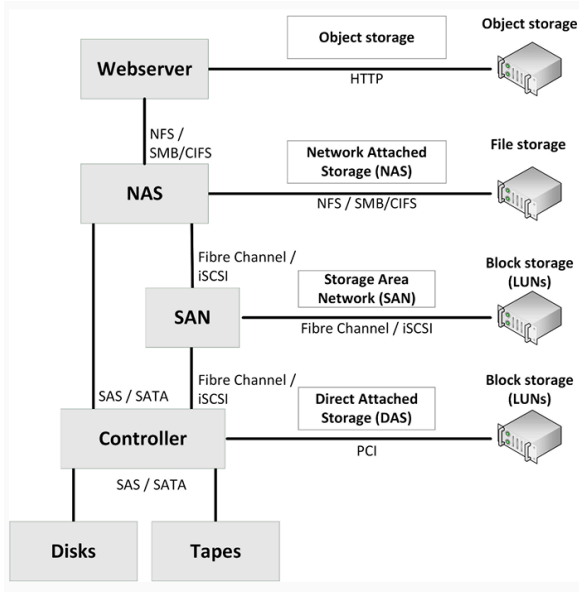
Example (Retail)	Many retail stores offer free Wi-Fi to customers to enhance their shopping experience. RADIUS can be used to authenticate customers attempting to
---------------------	---

* Example

C5: Storage

Storage Model

Server: external + internal



5.1 Storage building blocks

1Disks – command sets

Advanced Technology Attachment (ATA)	<ul style="list-style-type: none">• Aka IDE• simple hardware and communication protocol (disks to PCs)• (Common interface standard for connecting storage devices like HDDs and SSDs to computers.)
Small Computer System Interface (SCSI)	<ul style="list-style-type: none">• set of standards for physically connecting and transferring data between computers (mostly servers) and peripheral devices, like disks and tapes• (High-performance command set used mainly in servers and enterprise systems.)
command sets (ATA or SCSI) stayed mostly the same, but the physical connection type changed from parallel interface(many wires at once) to serial interface(one wire sends data in order) to improve speed, reliability, and scalability	

1Disks - type (2)

Mechanical Hard Disk (HDD)	<p>SATA (Serial ATA): Standard interface for consumer PCs; cost-effective and widely used.</p> <p>SAS (Serial Attached SCSI): Enterprise-grade interface offering higher speed and reliability.</p> <p>NL-SAS (Nearline SAS): Combines large capacity of SATA with reliability of SAS; used for backup or archival storage.</p>
Solid State Drive (SSD)	<p>Uses flash memory with no moving parts.</p> <p>Provides faster data access, lower latency, and higher durability compared to HDDs.</p>

2Tape (2)

Cheapest option for long-term storage and archiving.
Life expectancy up to 30 years (e.g., DLT, SDLT, LTO).
Disadvantages: fragile, slow sequential access, prone to damage.

Tape Cartridge Formats	SDLT (Super Digital Linear Tape): Up to 300 GB. LTO (Linear Tape Open): Most common; 80%+ market share. LTO-9: 18 TB uncompressed, ~200 MB/s throughput.
Library Types	Tape Library: Automates tape handling using drives, slots, barcode/RFID, and robotic loaders. Virtual Tape Library (VTL): Disk-based backup emulating tape systems; faster, supports parallel processing; typically uses SATA/NL-SAS arrays.

3Controller Implementation

Connect disks/tapes to servers via PCI or network (NAS/SAN).
Provide virtualization, high availability, deduplication, cloning, and thin provisioning

RAID (Redundant Array of Independent Disks)	<ul style="list-style-type: none">Improves performance and data availability using multiple disks.Implemented in hardware (controller) or software (OS). <p>Common RAID Levels:</p> <ul style="list-style-type: none">RAID 0 (Striping): Speed only; no redundancy.RAID 1 (Mirroring): Two copies; high reliability, 50% efficiency.RAID 10: Combines striping + mirroring; high performance & availability.RAID 5: Striping + single parity; balanced speed and protection.RAID 6: Dual parity; tolerates two disk failures.				
	RAID Level	Description	Benefits	Drawbacks	Use Case
	RAID 0 (Striping)	Data split evenly across multiple disks (no redundancy)	Increases performance	No redundancy; if one disk fails, all data is lost	Temporary or non-critical data where speed matters
	RAID 1 (Mirroring)	uses two disks that contain the same data Duplicates data exactly	High availability; most reliable	Uses 50% of disk space for redundancy; expensive	Critical data needing high reliability
	RAID 10 (Striping + Mirroring)	Combination of RAID 0 and RAID 1	High performance and availability	Uses 50% disk space for redundancy; costly	Databases and critical systems needing speed + redundancy
	RAID 5 (Striping + Distributed Parity)	Data and parity blocks striped across disks	Good balance of performance, availability, and storage efficiency	Single disk failure tolerated; slower write speed due to parity calculations	General purpose with moderate reliability needs
	RAID 6 (Striping + Double Parity)	Like RAID 5 but with two parity blocks	Can tolerate two disk failures	More overhead than RAID 5; slower writes	Systems requiring extra fault tolerance during rebuilds

Compression	<ul style="list-style-type: none"> Reduces storage needs (2–2.5× typical), depending on data type.
Data Deduplication	<p>Removes duplicate data segments to save space.</p> <p>Inline deduplication: immediate but slower.</p> <p>Post-process deduplication: done later, reduces performance impact.</p>
Cloning & Snapshots	<p>Clone: full copy of data (like RAID 1).</p> <p>Snapshot: point-in-time version; read-only original, writes go to a new volume.</p> <p>Used for backups, testing, or quick restore.</p>
Thin Provisioning	<p>Allocates more logical space than physically available.</p> <p>Dynamically assigns real capacity as needed; ideal for home directories, email storage.</p>

4 Storage Architectures (4)

Direct Attached Storage (DAS)	<ul style="list-style-type: none"> Local disks connected via SATA/SAS controller through PCI bus. OS manages LUNs and file systems. Simple and low-cost but limited scalability and sharing
Storage Area Network (SAN)	<ul style="list-style-type: none"> Dedicated high-speed network connecting servers and storage (block-level). Uses Fibre Channel, FCoE, or iSCSI. Centralized storage with large disk arrays (petabyte scale). <p>Key Components:</p> <ul style="list-style-type: none"> Fabric: SAN switch network. HBA (Host Bus Adapter): Server interface for Fibre Channel or iSCSI. <p>Protocols:</p> <p>Fibre Channel (FC): High reliability, 1–128 Gbps, zero data loss.</p> <p>Topologies: Point-to-Point, Arbitrated Loop, Switched Fabric (most common).</p> <p>FCoE (Fibre Channel over Ethernet): Runs FC over Ethernet ≥10 Gbps using DCB/CEE extensions.</p> <p>iSCSI: SCSI over TCP/IP; cost-effective and easy to deploy; lower performance but closing gap with 10/40 Gbps Ethernet.</p>
Network Attached Storage (NAS)	<ul style="list-style-type: none"> Provides shared file-level access over TCP/IP (NFS, SMB/CIFS). Acts as a file server, may use SAN as backend storage. Supports file-level snapshots and permissions (LDAP/AD integration). Clustered NAS: Distributes data across multiple nodes for scalability. Cloud examples: AWS File Gateway, Azure Storage Account, GCP Filestore.
Object Storage	<p>Stores data as objects (file + metadata + unique ID) accessed via REST API (HTTP).</p> <ul style="list-style-type: none"> Immutable — modified files are replaced with new versions. <p>Ideal for static content (archives, media, backups).</p> <p>Used by AWS S3, Azure Blob, GCP Object Storage.</p> <p>Supports massive scalability and geographic redundancy.</p>

5 Software Defined Storage (SDS)

- Separates control plane (management) from data plane (physical storage).
- Pools all physical storage into a virtualized shared resource.
- Works across SAN, NAS, and Object systems using commodity servers.
- All cloud providers' storage systems are SDS-based.

Functions:

- Provides data services: deduplication, compression, caching, snapshotting, cloning, replication, and tiering.
- Enables policy-based provisioning (e.g., automatic setup for databases with snapshots and tiering).

- Managed via API, CLI, or GUI, ensuring desired performance, availability, and security.

5.2 Storage: Availability, Performance & Security

1- Improve Storage Availability (4)

Refers to how often storage systems remain accessible for data operations (store, retrieve, manage).

1	<u>RAID (Redundant Array of Independent Disks)</u> <ul style="list-style-type: none"> • Combines multiple disks to provide redundancy and fault tolerance.
2	<u>Redundancy and Data Replication (2)</u> Synchronous replication: <ul style="list-style-type: none"> • Data is written to both primary and secondary storage before confirming write completion. • Ensures data consistency; risk of latency or downtime if connection fails. Asynchronous replication: <ul style="list-style-type: none"> • Data is first written to primary storage and later copied to secondary storage. • Reduces latency but may risk small data loss if failure occurs.
3	<u>Backup and Recovery</u> <ul style="list-style-type: none"> • Protects data from deletion, corruption, or disasters. <p>Follows the 3-2-1 rule: 3 copies, 2 media types, 1 offsite.</p> <p>Backup types:</p> <ul style="list-style-type: none"> i) Full Backup – Complete data copy. ii) Incremental Backup – Changes since last backup. iii) Differential Backup – Changes since last full backup. iv) Incremental Forever Backup – One full backup + continuous incrementals. v) Continuous Data Protection (CDP) – Captures every data change in real time (zero RPO).
4	<u>Archiving</u> Long-term storage for compliance and regulations. Data is read-only, encrypted, and stored on durable media (e.g., WORM tapes, optical disks). Use open formats (e.g., XML) for future readability; migrate to new media every 10 years.

2- Storage Performance (3)

Refers to how often storage systems remain accessible for data operations (store, retrieve, manage).

To improve storage availability:

1	<u>Factors affecting disk performance:</u> <ul style="list-style-type: none"> i) Disk Rotation Speed (RPM): Faster rotation → lower delay (e.g., 5400 RPM = 5.6ms; 15000 RPM = 2ms). ii) Seek Time: Time for disk head to locate track (3-9ms typical). iii) Interface Protocol: Determines data transfer rate (e.g., SATA 6Gb/s, SAS 12Gb/s, NVMe 7GB/s).
2	<u>Metrics to measure performance:</u> IOPS (Input/Output Operations Per Second): Number of reads/writes per second. <ul style="list-style-type: none"> • HDD: <500 IOPS • SSD: ~400,000 (read), 150,000 (write) • NVMe: up to 1,000,000 (read)
3	<u>To improve storage performance:</u> <ul style="list-style-type: none"> i) Cache: <ul style="list-style-type: none"> • Read cache stores frequently accessed data.

	<ul style="list-style-type: none"> • Write-back cache speeds up write operations. <p>ii) Storage Tiering:</p> <ul style="list-style-type: none"> • Assigns data to storage based on importance and usage (e.g., SSD for Tier 1, SAS for Tier 2, tape for Tier 4). • Automated tiering moves data between tiers based on access patterns. <p>iii) Load Optimization:</p> <ul style="list-style-type: none"> • Balances workloads to prevent bottlenecks. • Example: Oracle uses a mix of RAID 1 & 5 for database optimization.
--	--

3- Improve Storage Security(2)

1	<p><u>Data at Rest (on disk or tape)</u></p> <p>Data Encryption: Prevents unauthorized access without a decryption key.</p> <p>Self-Encrypting Drives (SEDs): Built-in encryption requiring password at startup.</p> <p>Cryptographic Disk Erasure (CDE): Deletes encryption keys to render data unreadable.</p>
2	<p><u>Data in Transit</u></p> <p>SAN Zoning: Divides Fibre Channel SANs into logical groups; only devices in the same zone can communicate.</p> <p>SAN LUN Masking: Restricts Logical Unit Numbers (LUNs) to specific hosts; implemented at HBA level.</p> <p>Combining zoning and masking enhances access control.</p>

C6: Compute

6.1 Compute building blocks

1	<p><u>Compute Housing</u></p> <p>Tower (Pedestal): Standalone; placed on floor.</p> <p>Rack Servers: Standardized frames for multiple systems.</p> <p>Blade Servers: Shared enclosure components (power, fans, network, SAN); reduced wiring, cost, and failure points. Blade enclosure: houses 8–16 blades; shared redundant power, Ethernet, SAN, and management modules.</p> <ul style="list-style-type: none"> • Reduced wiring then fewer failure points: it only need to run one cable to the enclosure/chassis机箱,外壳, making it easier to manage and avoid safety concerns from tangled wires电线缠绕. • Lower initial deployment costs: it uses the enclosure's shared components like power supplies and fans, also shared redundant components (Power supplies, Backplane for interconnection, Network switches (redundant Ethernet connections), SAN switches (redundant Fibre Channel connections)).
2	<p><u>Processor</u></p> <p><u>1</u> Central Processing Unit (CPU)</p> <ul style="list-style-type: none"> • Executes program instructions (arithmetic, logic, I/O). • Works using an instruction set (machine code). • Operates via clock cycles (GHz = billions of ticks/sec). • Word size: data handled per operation (modern = 64-bit). <p>Processor Families</p> <p>Intel x86: Standard CPU architecture.</p> <p>AMD x86: Competitor; Ryzen 9950X (16 cores, 5.7 GHz).</p> <p>x86-64: 64-bit architecture used by Intel & AMD.</p> <p>ARM: RISC-based; used in smartphones, tablets, IoT; licensed to many manufacturers (e.g., Apple M1–M3, AWS Graviton, Microsoft Ampere, Qualcomm Snapdragon X).</p> <p><u>2</u> Graphics Processing Unit (GPU)</p> <ul style="list-style-type: none"> • Thousands of cores for parallel processing. • Accelerates AI/ML workloads and graphics rendering. <p>Example: NVIDIA Tesla GP100 (3840 cores, 150B transistors).</p> <p>#GPU usage increases?</p> <p>growing demand for the artificial intelligence, machine learning, big data analytics application</p> <p>GPUs are highly parallel processors that much faster than CPUs for workloads of these application</p>
3	<p><u>Memory</u></p> <p>Evolution</p> <p>Early: vacuum tubes → relays → magnetic core memory → RAM chips.</p> <p>Core memory was replaced by transistor-based RAM in the 1970s.</p> <p><u>1</u> RAM (Random Access Memory)</p> <p>Volatile; data lost without power.</p> <p>SRAM(Static RAM): 6 transistors/bit, implement for cache</p> <ul style="list-style-type: none"> • fast, costly. <p>DRAM(Dynamic RAM): 1 transistor + 1 capacitor/bit, implement for main memory</p> <p>DRAM loses its data after a short time due to the leakage of the capacitors, refresh regularly to keep data available, cheaper.</p> <p><u>2</u> BIOS (Basic Input/Output System)</p>

	<p>Firmware initializing hardware before OS loads. Stored in flash memory; regularly updated (BIOS flashing).</p>
4	<p><u>Interfaces</u> <u>External Interfaces:</u></p> <p>① USB: Replaced older interfaces (since 1996). USB 3.1 = 10 Gbps; USB-C = 40 Gbps (USB 4) & up to 20V power.</p> <p>② Thunderbolt (Light Peak): Up to 80 Gbps; 100W power; uses USB-C connector.</p> <p><u>Internal Interfaces:</u></p> <p>③ PCI: (cheaper and wide adoption) slower, uses shared parallel bus, limited bandwidth.</p> <p>④ PCIe: (routed by a hub that allows multiple devices to communicate simultaneously) faster, uses point-to-point serial links, and scalable bandwidth</p>

6.2 Compute Virtualization

LPAR (Logical Partition)	<ul style="list-style-type: none"> • Hardware-based virtualization • full and can different OS per LPAR • divides single physical sys into multiple isolated env. • Common on IBM mainframes and Power Systems. 	Enterprise systems, mainframes, midrange
VM (Virtual Machine)	<ul style="list-style-type: none"> • Software-based virtualization • full and can different OS per VM (<i>Hypervisor</i> create vm) • creates and manages multiple virtual machines (VMs) on one physical host. 	Data centers, servers
Container	<ul style="list-style-type: none"> • OS-level • Shared host OS (Container Engine - Docker,) • Pay for container runtime or host instance 	Cloud-native apps
Serverless	<ul style="list-style-type: none"> • Function-level • by third-party cloud providers (AWS,) • Pay per function execution 	Event-driven cloud tasks

1	<u>Compute Virtualization vs. Public Cloud</u>		
	Aspect	Compute Virtualization (Virtual Machine)	Public Cloud (include serverless)
	Definition	Uses hypervisor 虚拟机管理程序 to create multiple virtual machines (VMs) on a single physical server.	Provides virtualized compute resources over the Internet, managed by cloud providers.
	Ownership	Managed by organization on its own servers. Full control over infrastructure.	Managed by third-party providers (AWS, Azure, GCP).
	Scalability	Limited by on-premise hardware capacity.	Highly scalable on demand.
	Cost Model	High initial cost, lower long-term cost.	Pay-as-you-go model, no upfront hardware.

	<table border="1"><tr><td>Use Case (Takoyaki Example)</td><td>A local takoyaki shop uses virtualization to run POS and inventory VMs on one physical server.</td><td>A takoyaki franchise uses cloud VMs to host its online ordering system that scales automatically during lunch rush hours.</td></tr></table>	Use Case (Takoyaki Example)	A local takoyaki shop uses virtualization to run POS and inventory VMs on one physical server.	A takoyaki franchise uses cloud VMs to host its online ordering system that scales automatically during lunch rush hours.
Use Case (Takoyaki Example)	A local takoyaki shop uses virtualization to run POS and inventory VMs on one physical server.	A takoyaki franchise uses cloud VMs to host its online ordering system that scales automatically during lunch rush hours.		
2	<p><u>Compute Virtualization Technology</u></p> <p>1 Emulation</p> <ul style="list-style-type: none">• Simulates different hardware architecture entirely through software.• Allows running software for one system on another (e.g., PlayStation emulator on PC). <p>Example (Takoyaki): A takoyaki shop runs an old POS app made for PowerPC on a modern x86 server using emulation.</p> <p>2 Logical Partitions (LPARs)</p> <ul style="list-style-type: none">• Hardware-based virtualization that divides a single physical system into multiple isolated env.• Common on IBM mainframes and Power Systems.• Each LPAR can run its own OS. <p>Example (Takoyaki): A large takoyaki chain uses an IBM Power server — one LPAR for sales, another for payroll, another for supply management — all isolated but sharing hardware.</p> <p>3 Hypervisor</p> <ul style="list-style-type: none">• Software-based virtualization layer that creates and manages multiple virtual machines (VMs) on one physical host.• Each VM has its own virtual CPU, memory, storage, and OS. <p>Types:</p> <ul style="list-style-type: none">• Type 1 (Bare Metal): Runs directly on hardware (e.g., VMware ESXi, Microsoft Hyper-V).• Type 2 (Hosted): Runs on top of existing OS (e.g., Oracle VirtualBox). <p>Example (Takoyaki): The main takoyaki shop uses a hypervisor to host separate VMs — one for accounting, one for recipe management, and one for online orders.</p>			

6.3 Types of Compute

1	<p><u>Compute: Cloud – Containers Technology</u></p> <p>i) Technology: Container</p> <ul style="list-style-type: none"> Containers package applications with their dependencies to ensure consistency across environments. <p>ii) Implementation</p> <ul style="list-style-type: none"> Chroot or Jail: Isolate processes and file systems for lightweight containment. Namespaces: Provide process isolation by separating system resources (PID, network, users). Cgroups: Control and limit resource usage (CPU, memory, I/O) for containers. <p>iii) Container Orchestration</p> <ul style="list-style-type: none"> Tools like Kubernetes automate deployment, scaling, and management of containerized applications.
2	<p><u>Compute: Cloud – Serverless Computing</u></p> <ul style="list-style-type: none"> Allows running code without managing servers. Automatically scales by providers based on demand <ul style="list-style-type: none"> charges only for actual execution time. <p>#Serverless computing means running code without managing servers yourself. The cloud provider (like AWS Lambda, Azure Functions, or Google Cloud Functions) handles server provisioning, scaling, and maintenance automatically. Just upload your function/code → provider runs it when triggered → you pay only for runtime.</p>

3	<p><u>Mainframe</u></p> <p>#x86 = general-purpose; Mainframe = high reliability, enterprise-scale workloads.</p> <p>i) Components</p> <ul style="list-style-type: none"> • PU (Processing Unit): Executes instructions and handles computations. • Memory: Stores active data and instructions. • I/O Channels: Manage data transfer between mainframe and peripherals. • Control Units: Direct and coordinate input/output operations. <p>ii) Mainframe Virtualization: Logical Partitions (LPARs)</p> <ul style="list-style-type: none"> • Divides mainframe hardware into multiple isolated logical systems, each running its own OS.
4	<p><u>Midrange Systems</u></p> <p>i) Midrange Virtualization: Logical Partitions (LPARs)</p> <p>Similar to mainframes but designed for medium-scale workloads in enterprises.</p>
5	<p><u>X86 Servers</u></p> <p>Industry-standard servers based on x86 architecture, commonly used for virtualization, cloud, and enterprise workloads.</p>
6	<p><u>Supercomputers</u></p> <p>#Supercomputer = classical parallel processing; Quantum = quantum bit-based, experimental.</p> <p>Extremely powerful systems designed for intensive scientific and engineering computations requiring parallel processing.</p>
7	<p><u>Quantum Computers</u></p> <p>Use quantum bits (qubits) and quantum mechanics principles to perform complex calculations exponentially faster than classical computers.</p>

6.4 Compute: Availability, Performance & Security

1	<p><u>Improve Compute Availability</u></p> <ul style="list-style-type: none"> • Hot swappable components <ul style="list-style-type: none"> ○ Components (e.g., memory, CPUs, interface cards, power supplies) can be installed, replaced, or upgraded while the server is running. ○ Why it improves availability: It prevents downtime because the system stays operational during • Parity memory 奇偶校验内存 <ul style="list-style-type: none"> ○ Uses parity bits as a simple error detection code to detect memory data corruption. ○ Why it improves availability: Early detection of memory faults prevents crashes and data corruption that could cause service interruptions. ○ it can only detect errors, not correct them. • Error-Correcting Code (ECC) memory <ul style="list-style-type: none"> ○ Detects and corrects single-bit memory errors automatically using Hamming Code or TMR. ○ Why it improves availability: Continuous error correction keeps servers stable and prevents system failure due to memory faults, which are more frequent in 24/7 environments. • Virtualization availability <ul style="list-style-type: none"> ○ Provides failover clustering where VMs automatically restart on another host if a hardware or OS failure occurs. ○ • Why it improves availability: Ensures minimal service downtime and automatic recovery during failures. ○ • Clustering requires spare or underloaded hosts to take over workloads, maintaining service continuity.
---	---

2 Performance

Factors affect compute performance

- Architecture of the server
 - Refers to CPU design, memory hierarchy, and bus interconnections.
 - • How it affects performance: Efficient architecture improves data throughput and reduces bottlenecks between components.
- Speed of the memory and CPU
 - Determines how fast data is processed and accessed.
 - • How it affects performance: Faster CPUs execute instructions more quickly; faster memory reduces data wait time.
- Bus speed
 - Refers to the data transfer rate between CPU, memory, and I/O devices.
 - • How it affects performance: A faster bus minimizes latency when moving data between components.

To improve compute performance

Moore's Law

- States that transistor count doubles every ~2 years.
- Why it affects performance: More transistors enable more processing units and higher efficiency, increasing CPU power — although physical limits are being reached.
- Increasing clock speed (more instructions per second)
 - CPU executes instructions at each clock pulse.
 - • Why it improves performance: Higher clock rates mean more instructions per second, improving compute speed — limited by heat and power constraints.
- Cache memory
 - Small, high-speed memory located on the CPU (L1/L2/L3).
 - • Why it improves performance: Stores frequently accessed data, reducing delays from slower main memory access.
- Pipeline 多个指令阶段可以重叠执行
 - Allows overlapping execution of multiple instruction stages.
 - • Why it improves performance: Increases instruction throughput by utilizing CPU resources continuously.
- Prefetching 加载即将到来的指令 and branch prediction 预测猜测下一个执行路径
 - Prefetching loads upcoming instructions; branch prediction guesses the next execution path.
 - • Why it improves performance: Reduces idle CPU cycles and cache misses, enabling faster instruction delivery.
- Superscalar CPUs
 - Executes multiple instructions per clock cycle via parallel execution units.
 - • Why it improves performance: Maximizes CPU utilization and throughput.
- Multi-cores CPUs
 - Integrates multiple cores into a single chip.
 - • Why it improves performance: Enables true parallel processing; improves multitasking and workload distribution while reducing heat and energy per operation.
- Hyper-threading
 - Allows one core to execute multiple threads simultaneously.
 - • Why it improves performance: Keeps execution pipelines active, improving CPU efficiency for multi-threaded workloads.
- Virtualization performance
 - Multiple VMs share one physical machine.
 - • Why it affects performance: Efficient consolidation reduces idle time but can cause I/O bottlenecks if overloaded.
 - • High CPU/memory capacity and fast storage improve performance.
 - • Overhead from hypervisor processing is typically <10%.
 - • Databases are I/O-intensive and may require dedicated physical servers or Raw Device Mapping to maintain speed.

3 Security

Physical security

i) Physical server

- USB port
- BIOS setting
- Server housing

Disable external USB ports and **secure BIOS with passwords.**

- Why it improves security: Prevents unauthorized access, data theft, or malware introduction via external devices.
- Detecting case openings can trigger alerts for tampering attempts.

ii) Data in use

- Trusted Execution Environment (TEE)
- Virtualization security

i) To minimize attacks

- Firewalls
- IDS
- Patching
- Minimize complexity of hypervisor

ii) To improve virtualization security

- **De-militarized Zone (DMZ)**

- The DMZ **isolates external-facing services** (e.g., web servers) **from the internal private network.**
- How it affects virtualization: It **prevents direct access from the internet to virtual machines** in the internal network, reducing the attack surface and protecting sensitive data or management interfaces.

- **Systems management console**

- Provides **centralized control and monitoring of all virtual machines and hosts.**
- How it affects virtualization: **Enables administrators to manage** user permissions, track configurations, detect unauthorized changes, and apply patches or updates consistently, ensuring the virtualized environment stays secure and compliant.

Chapter 7: Operating systems

7.1 Popular Operating Systems

IBM z/OS
 IBM i (OS/400)
 UNIX
 Linux
 Berkeley Software Distribution (BSD)
 Windows
 MacOS
 OS for mobile
 Special purpose OSs

7.2 Operating systems building blocks

1	<p><u>OSs building blocks</u></p> <ul style="list-style-type: none"> An operating system allows multiple users, processes, and applications to share hardware and hides hardware complexity. <p>Kernel</p> <ul style="list-style-type: none"> The heart of the OS — starts/stops programs, manages files, and controls hardware access to prevent conflicts. <p>Drivers</p> <ul style="list-style-type: none"> Connect hardware devices (printers, NICs, keyboards, video) to the kernel. <p>Utilities</p> <ul style="list-style-type: none"> Built-in tools like user interfaces, editors, loggers, and update managers. <p>Applications</p> <ul style="list-style-type: none"> Software that communicates with the OS through system calls and APIs.
2	<p><u>OSs functions</u></p> <p>Process Scheduling</p> <ul style="list-style-type: none"> Simulates parallelism by rapidly switching between processes (preemptive multitasking). Ensures fair CPU allocation through complex scheduling algorithms. <p>File Systems</p> <ul style="list-style-type: none"> Organize data in directories and files, managing disk communication and permissions. Support multiple file system types: FAT, NTFS, UFS, VxFS, Ext (Linux). Journaling file systems track changes for fast recovery. File systems must be mounted before use and may use drive letters (Windows) or mount points (UNIX/Linux). File sharing via NFS (UNIX) or SMB/CIFS (Windows). <p>APIs and System Calls</p> <ul style="list-style-type: none"> System calls provide a hardware-independent interface for applications. APIs define how software can use these system calls (e.g., POSIX for UNIX/Linux, Windows API). <p>Device Drivers</p> <ul style="list-style-type: none"> Software that allows the OS to communicate with hardware components through defined APIs. <p>Memory Management</p>

- Allocates/deallocates memory for applications.
- Includes cache, paging, swapping, and DMA (Direct Memory Access) for efficient transfers.
- Prevents memory shortage by moving pages to/from disk as needed.

Shells, CLIs, and GUIs

- User interfaces to interact with the OS.
 - CLI: text-based (bash, sh, cmd.exe).
 - GUI: graphical (Windows, X Windows).

OS Configuration

- Stores settings in databases or text files (Windows Registry, Linux /etc, AIX ODM).
- Editable via configuration tools that simplify text-based settings.

7.3 Operating systems: Availability, Performance & Security

Availability

Failover Clustering – Improves system uptime by automatically transferring workloads from a failed node to a standby node.

- Failover cluster is a group of independent servers running identical operating systems, known as nodes, connected through a network and managed by cluster software.
- Every active application has a standby counterpart on a passive node, which remains idle until a failover occurs.
- When a failure happens, the standby application automatically becomes active and continues serving clients [ensures minimal downtime and uninterrupted service].
- The cluster manages each running application within a node as a package of application components, called a resource pool or application package [organizes applications for automatic monitoring and failover control].

i) Shared storage

- All nodes can access the same data; ensures data continuity during failover.

ii) Shared node

- Provides redundancy; another node takes over automatically during failure.

iii) Voting and quorum disks

- Prevents “split-brain” errors by deciding which node stays active during network disconnections.

iv) Cluster-aware applications

- Run on multiple nodes simultaneously; improve both scalability and failover recovery time.

Performance

Factors affecting OS performance:

- Hardware performance: Faster CPU, RAM, and storage directly enhance OS responsiveness.
- Application load: Heavy workloads slow down the OS; optimization ensures smoother multitasking.
- OS configuration: Inefficient settings or unnecessary services waste system resources.

To improve performance:

- Increase memory: Reduces paging/swapping to disk, enabling faster data access and smoother multitasking.
- Decrease kernel size: Frees up RAM, shortens boot time, reduces crash risk, and lowers attack surface.

Why it improves performance:

Optimized memory and kernel management maximize hardware utilization, making the system faster and more stable.

Security

To improve OS security:

- **Patching:** Fixes vulnerabilities, bugs, and design flaws缺陷 to close potential attack vectors.
- **Hardening:** Disables unnecessary services, users, and protocols to reduce attack surface.
 - step by step process of configuring an operating system to protect it against security threats
 - used to instantiate new operating systems. Ensure security is optimal and is consistent in all deployment
- Virus Scanning: Detects and removes malicious software that can harm data or performance.
- Host-Based Firewalls: Filters incoming/outgoing traffic to block unauthorized access.
 - Most operating systems, including Windows, Linux, and UNIX, provide a built-in host-based firewall
- Limiting User Accounts: Minimizes risk from privileged misuse; enforces principle of least privilege.
- Hashed Passwords: Protects credentials by preventing recovery of original passwords.

Why it improves security:

Each measure reduces the likelihood of unauthorized access, malware infection, and privilege abuse, ensuring system integrity and reliability.

Chapter 8: End User Devices

8.1 End user devices building blocks

End user devices are tools humans use to interact with applications.

Examples: desktops, laptops, virtual desktops, mobile devices, printers.

[They connect users to IT infrastructure and deliver input/output functions.]

Desktop & laptop

i) Desktop

- High performance, can store large data and run complex software.
- Issues: complex management, high maintenance cost, and local security risks.
- [Best for stationary, high-power tasks.]

ii) Laptop

- Portable and as powerful as desktops.
- Common risks: theft, damage, and malware from external networks.
- Usually connected to docking stations for power and peripherals.
- [Ideal for mobility but less secure.]

Mobile devices

- Examples: smartphones, tablets, smartwatches, cameras, cars.
- Connect via wireless networks (UMTS/LTE/Wi-Fi).
- Features: small form factor, low bandwidth, variable reliability.
- [Enable mobile access but require adaptive app design and secure connectivity.]

Bring Your Own Device (BYOD)

- Users bring personal devices to access organizational data, instead of using only company-provided devices.
- user paid for the device (they brought their own device), it will not be acceptable to:
 - Have systems managers erase the device (including all family photos or purchased music) in case of an incident
 - Have personal data visible to the systems managers
- Conflict: user freedom vs. IT security control.
- **Solution:**
 - **virtualization(personal + work vm)** + Mobile Device Management (MDM) → separates personal and work environments.
 - **Mobile Device Management (MDM)**
 - to monitor organization data on personal devices.
 - It can **enforce password strength, device encryption, and remote wipe** in case of device loss or theft.
 - **Network Access Control and segmentation**
 - require devices to **meet security compliance standards (updated OS, antivirus, encryption)** before connecting to internal systems.
 - **Provide DMZ or VLAN for personal devices to connect,** instead of the core banking network.

Printers

1. Laser Printers: use toner + drum + laser → high quality, fast.
2. Inkjet Printers: spray ink droplets → cheaper, energy-efficient, high-quality color prints.
3. Multi-Functional Printers (MFPs): combine printer, scanner, copier, fax; include OS, storage, and network features → require patching and data protection.
4. Specialized Printers:
 - a. Dot Matrix: uses pins; reliable, low cost, noisy.
 - b. Line Printer: prints full lines; durable for industrial use.
 - c. Thermal Printer: uses heat-sensitive paper; quiet, compact, fast, but prints fade.

8.2 Desktop virtualization

Technology	Runs On	User Rss	Cost Model	Key Benefit
Application Virtualization	Local or Server	Shared OS	License-based	Compatibility & easy deployment
SBC	Server	Shared VM	Per-user/server license	Centralized management
VDI	Server(VM per user)	Dedicated VM	Pay-per-use (cloud)	Full OS isolation
Thin/Zero Client	Minimal hardware	Server resources	Low device cost	Easy to manage, low maintenance
PXE Boot	Network OS image	Network-based	Shared infra	Diskless, fast deployment

Application virtualization

- Runs apps in isolated virtual environments.
- OS resources are virtualized, not the app itself.
- Requests and file operations are redirected to virtualized locations.
- **[Allows running old/incompatible apps side-by-side and easier migration.]**
- Examples: Microsoft App-V, VMware ThinApp.

Server based computing (SBC)

- Applications/desktops **run on remote servers.**
- **Only display updates, keyboard, and mouse data are exchanged.**
- [Low bandwidth, centralized management.]
- Products: Windows RDS, Citrix XenApp.
- Advantages: Easier maintenance and consistent configuration.
- Disadvantages: **Depends on network latency, limited local performance.**

Virtual Desktop Infrastructure (VDI)

- Each user runs a full desktop OS in their own VM.
- Managed by a hypervisor that allocates resources.
- **[Provides isolation but uses more CPU/storage than SBC.]**
- Available as cloud services (Azure Virtual Desktop, Amazon WorkSpaces, Google Cloud VDI).
- Issue: ***"Logon storm" when many users start VMs simultaneously.***

Thin clients (2)

- Lightweight devices connecting to SBC/VDI servers.
- **No local storage or configuration; easy to replace.**
- [Reduce cost, maintenance, and security risk.]

i) Zero client

- No local OS; relies entirely on the server for boot and operation.

ii) Preboot eXecution Environment (PXE) boot

- Boots OS over network instead of local disk; used for diskless workstations.
- Requires constant network connection.
- [Useful for managed office environments but unsuitable for mobile devices.]

8.3 End user devices: Availability, Performance & Security

Availability

To improve end user device's availability:

- Use **high-quality components** to reduce failure rates.
- Perform **regular maintenance** (software updates, cleaning, battery checks).
- Add **redundancy** (spare devices, backup data copies).

Performance

To improve end user devices' performance:

- **Add more RAM** – improves multitasking and system speed.
- **Use SSD instead of HDD** – faster data access and startup.
- **Ensure sufficient network bandwidth** for all users.
- **Optimize mobile apps** for low-bandwidth or offline operation.
- Use **Server-Based Computing (SBC)** to reduce the effect of slow connections.

Security

Challenges:

- Devices are **spread across offices, homes, and client sites**.
- Not locked in a datacenter, thus prone to **theft, malware, and misuse**.

To improve security:

i) General Protection

- Use **laptop cable locks** to prevent theft.
- Install **malware protection** (antivirus, firewall).
- Enable **full-disk encryption** to protect data.
- Erase disks completely before disposal or repair.

ii) Mobile Device Management (MDM)

- **Monitor, maintain, and secure** mobile devices remotely.
- **Remote wipe** feature to delete data on lost or stolen devices.
- **Locate** stolen devices via tracking software.

iii) End User Authorization & Awareness

- Restrict user privileges (no admin rights).
- Set **BIOS passwords** and disable USB/DVD booting.
- Conduct **security awareness training** (social engineering, strong passwords, handling sensitive data).



iv) Network Access Control (NAC)

- Control network access based on:
 - **Device identity** (known/trusted device)
 - **User/group roles**
 - **Compliance status** (latest antivirus, OS patches)
- Non-compliant devices are placed in an **isolated network** until updated.

C9: Infrastructure Management

9.1 Infrastructure Deployment Options

Infrastructure deployment models (4)

On-Premises You own and manage all hardware and software in your own data center. Requires sufficient space, UPS, cooling, fire safety, redundant networks, and strong floor load capacity. <ul style="list-style-type: none"> • Large enterprises with strict data control needs • Industries with regulatory requirements (e.g., finance, defense) 	 Full control and high security  1. High Costs – Expensive setup and maintenance (hardware, licenses, IT staff, energy, cooling). Vendor lock-in risk. 2. Poor Scalability – Scaling requires physical installation and manual setup. 3. Limited Redundancy & Disaster Recovery – Single points of failure; complex and costly DR solutions. 4. Security Burden – Needs dedicated IT security, tools, and staff.
Public Cloud You rent computing resources (e.g., servers, storage) from providers like AWS, Azure, or GCP. Management level depends on the model (IaaS, PaaS, SaaS). <ul style="list-style-type: none"> • Startups or new businesses (greenfield projects) • Workloads that need rapid scaling or global reach 	<ul style="list-style-type: none"> • Cost-effective and scalable • Fast innovation and deployment • Pay-as-you-go model • Less control and data security • Vendor dependency • Requires reliable internet connection
Private Cloud Infrastructure dedicated to one organization, either hosted internally or by a provider. Functions as a software-defined data center (SDDC) with automation and orchestration. <ul style="list-style-type: none"> • Enterprises requiring strict security and compliance • Organizations needing isolated environments with virtualization 	<ul style="list-style-type: none"> • High control, privacy, and customization • Automated management (costing, reporting, scaling) • Comparable to IaaS model • High initial investment • Limited scalability • No true pay-per-use model
Hybrid Cloud Combines on-premises or private cloud with public cloud for flexibility and phased migration. Sensitive data stays private, while non-critical workloads use public resources. <ul style="list-style-type: none"> • Enterprises transitioning to cloud • Workloads needing both secure storage and dynamic scaling 	<ul style="list-style-type: none"> • Flexibility and cost optimization • Enables gradual migration • Best of both environments • Complex integration and management • Requires dual-environment expertise • Mixed cost model (cloud + on-premises)

9.2 Infrastructure Automation

Configuration Management Tools <ul style="list-style-type: none"> • Automate the configuration of servers, network devices, and infrastructure components, ensuring consistency. • Continuously check systems against blueprints; if deviations occur (e.g., missing software, changed settings), tools automatically correct them.
Infrastructure as Code (IaC)

Treats infrastructure like software code — configurations are defined in code files to automate deployment and management.

Terraform (by HashiCorp) is an open-source IaC tool using HashiCorp Configuration Language (HCL) or JSON to manage infrastructure across cloud and on-premises environments.

Version Control

Tracks code changes over time using repositories that create new versions automatically when updates are pushed.

Tools: Git (distributed CLI tool), GitHub (web-hosted Git), GitLab (self-hosted Git platform).

Orchestration Tools

Coordinate and automate complex workflows across multiple infrastructure components.

Act like a conductor ensuring all parts work in harmony (e.g., for deployments or scaling).

Key Functions:

- Workflow Automation – Automates provisioning, configuration, deployment, and scaling.
- Dependency Management – Handles retries, rollbacks, and alternative paths when tasks fail.
- Resource Coordination – Manages resources efficiently across cloud and on-premises platforms.

Cloud Management Platforms (CMPs)

Provide centralized automation for provisioning, monitoring, and managing cloud resources.

9.3 Infrastructure Documentation

Infrastructure documentation tools / techniques

- Documentation ensures reliability, security, and effective management.
- Preserves knowledge when staff leave, aids troubleshooting, and supports disaster recovery by showing how infrastructure is configured and functions.

Configuration Management Database (CMDB)	<p>Inventory of hardware, software, and networking components, detailing make, model, location, and function.</p> <p>Supports ITIL processes; should be updated regularly (preferably automated).</p> <p>Enables correlation of components to identify changes or root causes of failures.</p>
Diagrams	<p>Topology Diagrams – Show relationships and connections between components. Tools: Microsoft Visio, Diagrams.net.</p> <p>ArchiMate – Open standard for modeling enterprise architecture; describes IT systems, business processes, and information flows.</p>
IaCs tools	<p>IaC can serve as living documentation, showing how infrastructure is built and why. Documentation updates can be made alongside code changes.</p> <p>Advantage: Always up-to-date with modifications.</p> <p>Disadvantage: Requires reading code; doesn't provide instant visual overview like diagrams.</p>
Documenting procedures	<ol style="list-style-type: none"> 1. Procedures for routine tasks 2. Infrastructure naming convention 3. IP addressing plan 4. DNS naming convention 5. Fallback procedure 6. Disaster recovery plan