

IT Infrastructure

Security and Resilience Solutions

Ralf Süß
Yannik Süß

Apress®

IT Infrastructure

Security and Resilience Solutions

Ralf Süß
Yannik Süß

Apress®

IT Infrastructure: Security and Resilience Solutions

Ralf Süß
Singapore, Singapur, Singapore

Yannik Süß
Unterhaching, Bayern, Germany

ISBN-13 (pbk): 979-8-8688-0076-4
<https://doi.org/10.1007/979-8-8688-0077-1>

ISBN-13 (electronic): 979-8-8688-0077-1

Copyright © 2024 by Ralf Süß and Yannik Süß

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

Trademarked names, logos, and images may appear in this book. Rather than use a trademark symbol with every occurrence of a trademarked name, logo, or image we use the names, logos, and images only in an editorial fashion and to the benefit of the trademark owner, with no intention of infringement of the trademark.

The use in this publication of trade names, trademarks, service marks, and similar terms, even if they are not identified as such, is not to be taken as an expression of opinion as to whether or not they are subject to proprietary rights.

While the advice and information in this book are believed to be true and accurate at the date of publication, neither the authors nor the editors nor the publisher can accept any legal responsibility for any errors or omissions that may be made. The publisher makes no warranty, express or implied, with respect to the material contained herein.

Managing Director, Apress Media LLC: Welmoed Spahr
Acquisitions Editor: Susan McDermott
Development Editor: Laura Berendson
Coordinating Editor: Shaul Elson

Cover designed by eStudioCalamar

Cover image from www.pixabay.com

Distributed to the book trade worldwide by Apress Media, LLC, 1 New York Plaza, New York, NY 10004, U.S.A. Phone 1-800-SPRINGER, fax (201) 348-4505, e-mail orders-ny@springer-sbm.com, or visit www.springeronline.com. Apress Media, LLC is a California LLC and the sole member (owner) is Springer Science + Business Media Finance Inc (SSBM Finance Inc). SSBM Finance Inc is a **Delaware** corporation.

For information on translations, please e-mail booktranslations@springernature.com; for reprint, paperback, or audio rights, please e-mail bookpermissions@springernature.com.

Apress titles may be purchased in bulk for academic, corporate, or promotional use. eBook versions and licenses are also available for most titles. For more information, reference our Print and eBook Bulk Sales web page at <http://www.apress.com/bulk-sales>.

Any source code or other supplementary material referenced by the author in this book is available to readers on GitHub (<https://github.com/Apress>). For more detailed information, please visit <https://www.apress.com/gp/services/source-code>.

Paper in this product is recyclable

*In gratitude for that initial laptop from father to son,
laying the foundation for our deeply cherished and
shared path in technology.*

Table of Contents

About the Authors	xiii
Introduction	xv
Chapter 1: The Architecture of IT Cloud Services	1
1.1 Cloud Services	4
1.2 IT Services Provided by CSP	5
Data as a Service (DaaS)	6
Communications as a Service (CaaS).....	7
Infrastructure as a Service (IaaS).....	7
Platform as a Service (PaaS).....	8
Software as a Service (SaaS).....	8
Business Process as a Service (BaaS)	9
X as a Service (XaaS).....	9
1.3 Deployment Models of Cloud Services	10
Deployment Models.....	11
1.4 Summary.....	13
Chapter 2: Data Center Facilities	15
2.1 Data Center Facility Design.....	17
2.2 Established Standards for Data Centers	18
Uptime Institute Tier Standard.....	18
EN 50600 Series.....	20

TABLE OF CONTENTS

- ASHRAE 20
- Other Examples 20
- 2.3 Space 21
 - Physical Space 21
 - Lighting..... 22
 - Noise 22
 - Weight 22
- 2.4 Facility Management..... 22
- 2.5 Infrastructure 23
 - Racks..... 23
- 2.6 Cooling 24
 - Temperature 24
 - Cooling..... 24
- 2.7 Power 26
 - Power Demands of DC..... 26
 - Uninterruptible Power Supply (UPS) 27
- 2.8 Security 30
 - Access Control..... 30
- 2.9 Summary..... 31
- Chapter 3: Compute and Virtualization..... 33**
 - 3.1 Hardware 34
 - 3.2 Software..... 35
 - 3.3 Types of Computer Systems..... 36
 - 3.4 Purpose of Computer Systems..... 38
 - General-Purpose Computer 38
 - Specialized Computer Systems 40
 - 3.5 Data Centers 41

3.6 Compute Resources	43
Compute Building Blocks.....	43
3.7 Computer Operating Systems	52
Key Functions of an Operating System.....	52
History of Operating Systems	53
Operating System Functions	56
Operating System Types	61
3.8 Compute Virtualization	70
Hardware Virtualization	70
Types of Hardware Virtualization	71
Autonomic Computing	72
Container	74
Kubernetes	76
3.9 Edge Computing.....	83
3.10 Compute Resiliency	87
Definition	87
Fault-Tolerant Computing.....	87
Resilient Computing	89
Federated Architecture	91
3.11 Provisioning and Administration of Compute Resources	92
Computer Workload	92
Types of Computer Workloads	92
Workload Deployment	95
Benefits and Challenges of Private and Public Clouds	96
Kubernetes Workload Management	99
3.12 Charging for Compute Resources	102
Terminology	102
Pricing Variables.....	103

TABLE OF CONTENTS

- Dashboards 104
- Charging Structure 104
- 3.13 Summary..... 105
- Chapter 4: Storage and Virtualization107**
- 4.1 Storage Resources..... 109
 - Primary Storage..... 109
 - Secondary Storage 111
- 4.2 External Storage Systems..... 116
 - Direct Attached Storage (DAS)..... 117
 - Network-Attached Storage (NAS) 118
 - Storage Area Network (SAN)..... 119
- 4.3 External Disk Configurations..... 121
 - JBOD..... 121
 - Disk Arrays 122
- 4.4 Storage Virtualization..... 132
 - Access Modes for Virtualized Storage 133
 - Types of Storage Virtualization 134
 - Symmetric and Asymmetric Virtualization..... 135
 - Virtualization Methods..... 136
 - Storage Virtualization Benefits 138
- 4.5 Storage Security and Resilience 139
 - Data Resiliency 140
 - Best Practices for Data Resiliency..... 146
 - Storage Security Implementation 149
- 4.6 Storage Provisioning and Administration..... 150
 - Classic SAN Provisioning 150
 - Storage Provisioning in Modern DCs 151

Storage Pools	152
Storage Allocation Tiering.....	155
Comparison of Public and Private Storage Infrastructure	155
4.7 Charging for Storage Resources	156
Cost-Saving Options for Cloud Storage	158
4.8 Summary.....	159
Chapter 5: Network.....	161
5.1 DC Network Components	162
Cables.....	162
Structured Cabling.....	164
Switches and Router	165
DC Switches	167
Router	168
DC Gateway	169
5.2 DC Network Topology	170
Centralized Topology.....	170
Zoned Network Topology	171
Top-of-Rack Topology	172
Mesh Network Topology	172
Multi-tier Network Topology	173
Software-Defined Networks	175
5.3 Network Resiliency	178
Network Fault Management.....	178
5.4 Network Provisioning and Administration	179
Network Provisioning	179
Network Administration.....	180

TABLE OF CONTENTS

- 5.5 Resilient Network for IT Data Center..... 184
 - Security Threats 185
 - Physical Network Security..... 185
 - Technical Network Security..... 185
 - Administrative Network Security..... 185
- 5.6 Resilient Network Architecture for IT Data Centers 186
- 5.7 Summary..... 188
- Chapter 6: Backup191**
 - 6.1 Evolution of Backup Systems..... 191
 - 6.2 Today’s Backup Systems..... 194
 - Backup in Private Data Centers or Private Clouds..... 194
 - Cloud-Based Backup Systems 194
 - 6.3 Types of Backup Methods 195
 - 6.4 Disaster Recovery Planning 197
 - Understanding Disaster Recovery 198
 - The Disaster Recovery Process 198
 - Key Considerations in Disaster Recovery Planning 199
 - The Role of Technology in Disaster Recovery 201
 - 6.5 Summary..... 201
- Chapter 7: Data Center Security and Resiliency203**
 - 7.1 Vulnerabilities of Computer Systems 204
 - Denial-of-Service Attack 204
 - Phishing..... 205
 - Spoofing Attack 205
 - Eavesdropping..... 206
 - Backdoor 207
 - Direct-Access Attacks 207

Privilege Escalation	208
Reverse Engineering.....	208
Multivector and Polymorphic Attacks	208
Social Engineering.....	208
Malware.....	209
7.2 Motivations and Impact of Attacks.....	210
Impact of Security Breaches	210
Attacker Motivation	210
7.3 Security by Design	212
Security Architecture	214
Security Infrastructure.....	214
Vulnerability Assessment and Management.....	215
Reducing Vulnerabilities	217
Hardware Protection.....	218
Access Control Lists	221
Security Tools	222
Security Training.....	223
Cyber Hygiene	223
Incident Response	224
Cybersecurity Planning.....	226
7.4 DC Resilience	227
DC Security.....	227
Critical Services.....	230
Achieving Data Center Resiliency	230
Improving Resilience	231
7.5 Summary.....	232

TABLE OF CONTENTS

- Chapter 8: IT Support Services235**
 - 8.1 IT Help Desk 236
 - Options to Contact a Help Desk 237
 - Trouble Ticketing Systems..... 245
 - 8.2 IT Service Desk 249
 - 8.3 Remote DC and Edge Computing Support 252
 - Remote DC Support 252
 - Edge Computing Support..... 252
 - 8.4 Summary..... 255
- Chapter 9: Summary257**
 - 9.1 Resilient IT Infrastructure 258
 - 9.2 IT Services Provided by Cloud Service Providers..... 258
 - 9.3 Data Center 259
 - 9.4 Compute..... 260
 - 9.5 Storage..... 262
 - 9.6 Network 264
 - 9.7 Backup 266
 - Backup System..... 266
 - 9.8 Resiliency..... 267
 - 9.9 IT Services 269
 - Help Desk 269
 - Service Desk..... 270
- References271**
- Index.....283**

About the Authors



Ralf Süß has dedicated over 40 years to the IT industry, deepening his expertise in all aspects of computing, data center design, and management. His experience is underscored by a significant role with Hewlett Packard’s Pacific Asia Technical Sales. Ralf’s professional journey has led him to collaborate with global giants in the realm of cloud services, including Amazon, Apple, and Microsoft. In tandem, he has catered to the needs of renowned network equipment providers, including Ericsson, Nokia, and Cisco. Melding theoretical acumen with practical experiences, Ralf has seamlessly adapted, ensuring his skills remain both relevant and innovative. This adaptability, combined with his rich experience, positions him as a knowledgeable figure who consistently offers insights and reflections from his vast tenure in IT infrastructure.



Yannik Süß has an extensive two-decade background in the Web and e-commerce. Driven by a deep passion for technology, he has managed critical web projects, spanning from site development to intricate e-commerce platforms. Proficient in data-driven decision making, he has implemented and overseen the development of robust data warehouses to optimize business intelligence and analytics. With a master’s degree in

ABOUT THE AUTHORS

Strategic IT Management and an ongoing doctorate, Yannik expertly bridges theoretical frameworks with practical challenges. He has also authored the second edition of *E-commerce for Small and Medium-sized Enterprises*, published by Springer, showcasing his comprehensive grasp of the digital commerce landscape.

Introduction

This book describes how IT systems have evolved from a relatively marginal role to perhaps the most essential pillar of modern infrastructure. Even in the 1950s, infrastructure was an element of roads, bridges, electricity, water supply, wastewater collection, and voice telecommunications. Today, we can hardly imagine living in a world without ubiquitous telecommunication and IT services. In this book, we explain how critical building blocks such as data centers, computer systems, storage systems, and IT security systems have evolved and how they provide the foundation to build a state-of-the-art resilient IT infrastructure.

How Did It All Begin?

Computing started in the fields of science and military applications. However, at the beginning of the computer era, the most famous historical event is probably the system that Alan Turing, an English mathematician, created during World War II. His machine was able to break the secret Nazi communication code. The Nazis had used the Enigma, a mechanical encoding machine, which created a changing encoding pattern every 24 hours. It was considered unbreakable. With Turing's "computer system," the British intelligence service was able to decode Nazi internal communication. This proved to be one of the contributing factors that eventually led to the end of World War II.

Until the 1950s, the word "computer" was mainly used to describe the actual people who did calculations for an organization, such as people calculating data in a research institute. The emergence of the

INTRODUCTION

word “computer” with a broader audience, as depicted in Figure 1, began to occur only in the 1960s, shaping our modern understanding of it. IT systems have come a long way during the last 60 years.

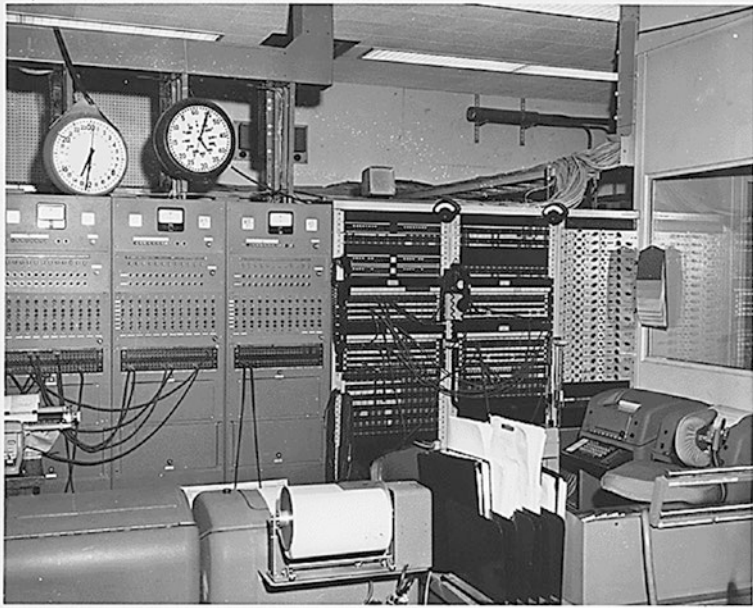


Figure 1. *Computers at NASA (NationalArchivesCatalog, 1952–1968)*

At the beginning of the 1960s, IT systems were just emerging and had a limited list of applications:

- **Military systems:** Military systems were important applications for computers at this time, mainly in the fields of surveillance, cryptography, missile defense systems, and logistic systems.
- **Scientific research:** The second leading application for computer systems, mainly used for modeling and simulations, was scientific research.

- **Government applications:** Governments used computers in several fields, such as tax collection and census analysis.
- **Business applications:** Business applications were limited to the accounting and payroll processing of larger organizations.
- **Education:** Computers were primarily used for research and teaching purposes for education.

Today's IT Infrastructures

Leaping forward 60 years to today, IT systems have become an integral part of modern-day society, and their importance cannot be overemphasized. There were two main phases. First, from the 1960s to the late 1990s, computer systems played an ever-increasing role in supporting a company's business. Computer systems were used for everything from hiring people, accounting purposes, sourcing materials, managing supply chains, providing marketing, driving sales, internal training, and so forth. In the second wave of computing, IT systems were not just supporting companies' business, they became an essential part of a new type of company. Google, Facebook, Amazon, Airbnb, Netflix, and Uber are companies that have built their core business with applications running on computer systems.

Today's world of social media is entirely built on computer systems and modern telecommunication infrastructure. At the core of the second wave of computerization is the evolution of the Internet, which enables people to connect to almost anybody in the world in real time from anywhere. In addition, the ubiquitous presence of the Internet made it possible to automate most business processes.

INTRODUCTION

In summary, IT systems have become an essential part of our daily lives. How we live, work, and communicate has drastically changed over the last 60 years. As technology continues to evolve, IT systems will continue to play an increasing role in our societies for the foreseeable future.

In the next chapter...

- IT evolution
 - CSP
 - Cloud services
-

CHAPTER 1

The Architecture of IT Cloud Services

From mainframes to cloud computing, cloud services have evolved over the years and were mainly driven by the advancements of three hardware pillars, namely: computer systems, storage systems, and data networks. We will look closer into the evolution of these three components in Chapters 3, 4, and 5. However, in retrospect, the eras in this technological evolution span in this way:

Mainframe computing (1960s to 1980s): Large, expensive, and complex machines were called mainframes, primarily used by large organizations for business applications like payroll and accounting. Mainframes are centralized computer systems where all processing and data storage occur on a single machine. During this period, IBM was the dominant computer company to the degree that the words “computer” and “IBM” were used synonymously.

Client-server computing (1980s to 1990s): Although the mainframe era was based on a fully centralized concept, the emergence of technical workstations for engineering and PCs for general

office purposes changed the landscape drastically. During this time, computing became more decentralized as data processing and storage were distributed across multiple machines. In addition, servers handle requests from numerous clients, such as a database server. An essential enabler for this era was the development of the Ethernet and the TCP/IP, which allowed for the first-time communication between computer systems manufactured by different vendors and running other operating systems.

Internet computing and utility computing (late 1990s to mid-2000s): The client-server era changed the architecture of applications. Until the 1990s, applications were running on a vertical stack from a vendor that included micro processes, operating systems, middleware, and application. Everything was proprietary, and moving an application from one vendor to another was a massive effort. During this era, the application architecture started to become horizontal, and many vendors used the same microprocessors from Intel and Motorola, standard operating systems such as Unix and Windows, and standard database software from vendors like Oracle and Informix. This made it possible to move an application from one vendor to another. In addition, the wide adoption of the Ethernet and TCP/IP enabled general connectivity between all computer systems that allowed such access.

The Internet emerged with even more enhancements and standardization of crucial protocols, including network protocols, and the development of browsers. The application became web based with early adoptions such as email, search engines, and e-commerce. These services were delivered over the Internet and were accessible from anywhere as long as there was a stable Internet connection.

These developments lead to the next step: utility computing. Amazon was the first starting to build large data centers and offered IT services like a utility under the AWS (Amazon Web Services) brand. Microsoft followed Amazon under the brand Azure. Later, Google entered the market with its Google Cloud Platform.

Cloud computing (late 2000s to present): Cloud computing evolved from utility computing. Cloud services are delivered on a pay-per-use basis, allowing organizations to scale up and down the IT resources they need. The data cloud center depicted in Figure 1-1 showcases a large facility comprising multiple server racks arranged in a highly organized manner.

Over time, cloud service providers (CSP) have built massive IT infrastructure resources that power a fair share of global IT applications. As a result, IT has become for modern societies what electricity was at the beginning of the industrial revolution.

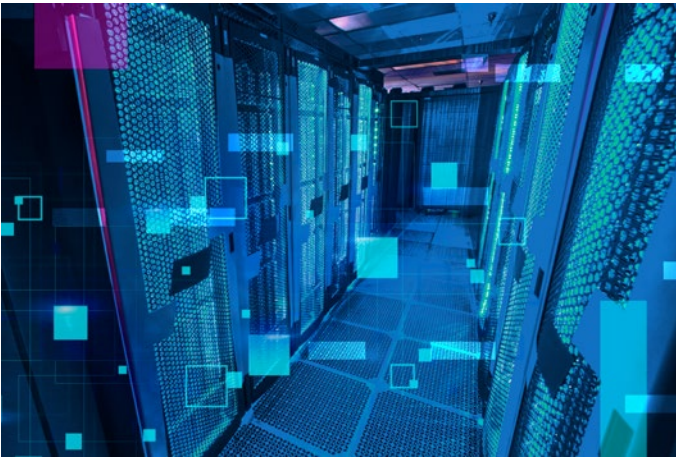


Figure 1-1. Cloud Data Center Example (rawpixel.com, 2023)

1.1 Cloud Services

The main characteristic of Cloud Services (CSPs [is this correct?]) is its distributed computing model. A vast number of servers are installed in racks and connected through a highly reliable and fast DC network infrastructure. This allows CSPs to provide services to organizations over the Internet. Although the cloud data center infrastructure varies depending on the benefits offered, several features appeared in most systems:

- **Virtualization:** Virtualization technology allows for the creation of virtual machines or containers that will enable multiple users to share the same physical hardware that will ensure efficient resource utilization by scaling their services to meet demand.
- **Data storage:** User data is stored in distributed storage systems that are designed to scale and provide high data security, such as file systems or object stores.

- **Load balancing:** In cloud environments, applications are distributed through a load balancer. This has two main advantages:
 - First, workloads can be distributed to resources that are not used.
 - Hardware problems become minor issues, as the load balancer can quickly move a workload from one defective system to another system.
- **Security:** Cloud services use various security measures to protect user data, including encryption, access controls, and monitoring. Processes to run and monitor a data center will generally be more advanced with larger CSP organizations than with smaller IT organizations.
 - **Front end:** Cloud services can be provisioned through several front ends, including a web browser, a mobile app, or a command-line interface.
 - **Back end:** The server infrastructure, storage systems, and the software that provides the service are called the back end of a CSP.

The CSPs have built their infrastructure based on an advanced cloud services architecture designed to provide users with scalable, reliable, and secure services over the Internet.

1.2 IT Services Provided by CSP

With CSPs' introduction of cloud data centers, a growing share of IT workloads was moved from private data centers, also called private clouds. At first glance, this looks like an either/or approach. An organization

operates its private cloud or uses a CSP. But organizations use private and public clouds at the same time. The CSPs offer several options to use their IT services, most commonly

1. Data as a Service
2. Communications as a Service
3. Infrastructure as a Service
4. Platform as a Service
5. Software as a Service
6. Business as a Service
7. X as a Service

Here is a description of each of them.

1.2.1. Data as a Service (DaaS)

Data as a Service is a model for distribution and information provision, which allows customers to access various data files via a network like the Internet. Data files include text, images, sounds, and videos. This model offers customers and client-oriented enterprises convenient and cost-effective solutions. Data as a Service enables decoupling data cost and utilization from platform or software cost. The prices DaaS providers offer can be volume or format-based.

Volume-based pricing contains a fixed price per megabyte of data in the whole repository. In format-based pricing, the charge is set based on the data format. Data as a Service can easily move data between platforms. This avoids the conflict and confusion arising from multiple copies of the same data or file. In addition, DaaS implements access control measures like strong encryption and passwords, avoiding “vendor lock-in,” ease of administration and collaboration, diverse platform compatibility, automatic updates, and global accessibility. These measures preserve data integrity.

1.2.2. Communications as a Service (CaaS)

Communications as a Service can contain Voice over IP (VoIP), collaboration, mobile and fixed device applications for videoconferencing, and instant messaging (IM). A Communications as a Service vendor assures guaranteed Quality of Service (QoS) and oversees all hardware and software management. Businesses operating on Communications as a Service can deploy communications and modes on a pay-per-use basis as and when they need it. The advantage of Communications as a Service is that it does not require significant capital investment and doesn't incur any ongoing overheads.

With CaaS, small and medium-sized businesses can be more flexible and expandable, which allows advantages like on-demand coverage and the addition of devices or modes. The network capacity and features are changed daily if necessary to avoid resource wastage and to ensure the functionalities are at par with the demand. This guarantees that the system is not outdated and would not require significant replacements or upgrades.

1.2.3. Infrastructure as a Service (IaaS)

Infrastructure as a Service is generally used by IT administrators, architects, and operators. IaaS provides virtual hardware resources on demand. It offers multiple on-demand features like individual email, domain name servers, messaging systems, and private networks. IaaS applications may incur OS license fees and require the installation of compatible software on the servers. Here, the customer has the flexibility to provision and de-provision the resources as per their business demand. An entire computing infrastructure can be built by an organization using an application. This type of service is beneficial for startup companies. It will help the company focus on its business process as IT infrastructure management and uptime are maintained by the IaaS cloud service provider.

1.2.4. Platform as a Service (PaaS)

Platform as a Service was mainly created for developers because it is required to run the software products, which need physical servers, web servers, and database software. The application is active on the software stack along with compilers, dependency files, etc. The disadvantage of Platform as a Service is that it can take very long and is rather complex to build the application's platform. PaaS must be updated and monitored regularly.

PaaS is directed at (because the definition of PaaS is Platform as a Service, it was repeating...) developers as it provides them with all the tools and dependency files required to develop new software (there was too much "developer"). The deployment platform, however, includes all controls to maintain the developer's customers for whom the product will be deployed. Platform as a Service provides an external platform to execute software applications without any administrative requirements of the lower-level components.

1.2.5. Software as a Service (SaaS)

Software as a Service is an application generally used by the end user. The advantage of this model is its lack of installation requirements. A network connection and a browser complete the Software as a Service requirements. The software is installed in the service provider's application server, and the customer gets the software service offshore (i.e., software delivered from the server), which is an outside organization. The service provider also achieves multi-instances, and although it has only one copy of the installation, multiple users can be handled with their dedicated storage space. Software as a Service is generally leased from a single vendor. It is an outsourced enterprise communication solution, which refers to the services hosted outside an organization.

1.2.6. Business Process as a Service (BaaS)

Any horizontal or vertical business process delivered based on the cloud service model is known as a Business Process as a Service. Software as a Service, Platform as a Service, and Infrastructure as a Service rely on this service. With the advent of cloud computing, companies prefer a more service-oriented approach. Instead of assuming that a packaged application is needed that includes business logic, data, and processes, selecting a process application not tied to a single application is possible.

A business is unable to forecast the future leverage of a business process. Hence, a Business Process as a Service must support multiple languages and deployment environments. In addition, a Business Process as a Service environment must be able to handle massive scaling. It must be able to progress from running a few to an increasing number of processes and customers. The service accomplishes that objective by optimizing the underlying cloud services to support this elasticity and scaling.

1.2.7. X as a Service (XaaS)

XaaS denotes the growing number of services distributed over the Internet instead of local or on-site provision. It is at the core of cloud computing. X as a Service uses hybrid cloud computing to deliver IT as a service. It refers to either a single or a blend of Software as a Service, Infrastructure as a Service, Platform as a Service, Communications as a Service, and Business Process as a Service. X as a Service is frequently used for the previously detached services on private or public clouds that are now integrating and becoming transparent. Table 1-1 provides a comprehensive overview of the service options.

Table 1-1. *Comprehensive Overview of the Service Options*

Function managed by	Onsite	DaaS	CaaS	IaaS	PaaS	SaaS
Application	Client	Client	Client	Client	Client	CSP
Data	Client	Client	Client	Client	Client	CSP
Runtime	Client	Client	Client	Client	CSP	CSP
Middleware	Client	Client	Client	Client	CSP	CSP
Operating system	Client	Client	Client	Client	CSP	CSP
Virtualization	Client	Client	Client	CSP	CSP	CSP
Server	Client	Client	Client	CSP	CSP	CSP
Storage	Client	CSP	Client	CSP	CSP	CSP
Networking	Client	Client	CSP	CSP	CSP	CSP

1.3 Deployment Models of Cloud Services

Cloud computing offers various deployment models for organizations: public, private, community, and hybrid. While these models save money, they can also pose security and management challenges. Businesses must therefore assess their needs before selecting a deployment model. For example, the public model allows public and organizational access, the private model ensures data security but is more expensive, and the community model involves shared resources. Finally, the hybrid model combines the benefits of both private and public models while maintaining data security.

1.3.1. Deployment Models

There are four main deployment models:

1. Public model
2. Private model
3. Community model
4. Hybrid model

Here is a comprehensive description of each deployment model. Considering these factors reduces the risk and assists in choosing the best option for each specific business.

1.3.2. Public Model

The public model is a cloud deployment model that allows both the public and organizations to access the IT infrastructure. It provides a shared environment where multiple users can access resources and services. The public model is typically cost-effective as the infrastructure and maintenance costs are spread across various users. However, since the infrastructure is shared among multiple entities, it may raise data security and privacy concerns.

1.3.3. Private Model

In this model, hosting is built and maintained specifically for each client, which ensures data security. The necessary infrastructure can be on-site or at a third-party location. Private models can be networks residing within an organization or hosted in another data center.

For instance, if it is hosted in another data-center leasing organization or by a network provider, it is termed a virtual private network. The private model is not cost-efficient, but the advantage of this model is

the level of security it offers. As data security has become a concern for many organizations, the private model ensures a secure-access VPN on the physical location within a client's firewall system. The significant advantage of the private model is the total cost of ownership for hardware and other components remaining with the organization.

1.3.4. Community Model

The community model involves multiple organizations sharing a common infrastructure, leading to cost reduction and resource optimization. It requires cooperation and standardized policies for data and application management. This collaborative approach saves money, is scalable, and fosters efficient resource utilization. However, successful implementation relies on coordination and adherence to shared governance frameworks and compliance standards. The community model offers organizations a cost-effective, streamlined resource-sharing data management approach.

1.3.5. Hybrid Model

This model enables different businesses to utilize secured applications and data hosting on a private model. However, companies continue to get cost benefits as the shared applications and data are kept on a public model, which has advantages over both private and public models. Migration of workloads between private and public clouds is assisted without inconveniencing the user. Several PaaS deployments expose application programming interfaces (API). This can be combined with internal or private cloud-hosted applications without compromising on security features. Hybrid models are more secure since customers can maintain highly sensitive data with their servers and less sensitive data with the network service provider's server.

Table 1-2 shows a comprehensive overview of all deployment models. CSP stands for Cloud Service Providers and DCH stands for Data Center Hub.

Table 1-2. *Deployment Options for Cloud Services*

Function properties	Public	Private	Community			Hybrid	
Application	Public	Private	Community			Noncritical	Critical
	CSP	CSP/DCH	CSP	CSP/DCH	CSP	Public	Private
Data	Public	Private	Community			Noncritical	Critical
	CSP	CSP/DCH	CSP	CSP/DCH	CSP	Public	Private
Management	Public	Private	Community			Noncritical	Critical
	CSP	CSP/DCH	CSP	CSP/DCH	CSP	Public	Private
Server and Storage	Public	Private	Public	Private	Noncritical	Critical	
	CSP	CSP/DCH	CSP	CSP/DCH	CSP	Public	Private
Networking	CSP	CSP	CSP			CSP	

1.4 Summary

In this chapter, the evolution of cloud computing was introduced, tracing its development from mainframe computing to the present day. The distinct eras in technological change were identified, including mainframe computing, client-server computing, Internet computing, utility computing, and finally, cloud computing. The role hardware advancements in computer systems, storage systems, and data networks had in driving the evolution of cloud services was highlighted.

Further, the characteristics of cloud services, which are based on a distributed computing model, were explained. This chapter also described the critical components of cloud data centers, including virtualization,

data storage, load balancing, security measures, front-end interfaces, and back-end infrastructure. It emphasized the scalability, reliability, and security that cloud service providers (CSPs) provide.

It also focused on the various IT services offered by CSPs, including Data as a Service (DaaS), Communications as a Service (CaaS), Infrastructure as a Service (IaaS), Platform as a Service (PaaS), Software as a Service (SaaS), Business Process as a Service (BaaS), and X as a Service (XaaS). Each service model has been explained and its benefits and target users highlighted. The chapter also discussed cloud service deployment models, including public, private, community, and hybrid models. This provides a concise overview of the progression of cloud computing, the key components comprising cloud data centers, the array of services offered by cloud service providers (CSPs), and the diverse deployment models accessible to organizations. In addition, it sets the foundation for understanding the subsequent chapters on cloud computing.

In the next chapter...

- What is a data center facility?
 - Crucial factors – design and components
-

CHAPTER 2

Data Center Facilities

A data center (DC) facility is the physical building used for equipment to run a data center. It used to be just a room in an office building. Modern large DCs, built by a DC facility provider or cloud service provider (CSP), as seen in Figure 2-1, cover full floors in a large building, or an entire building is solely built for data center activities.



Figure 2-1. Google Data Center, Iowa, USA (ChadDavis, 2017)

A DC facility design must consider the following factors:

- **Space:** A DC facility must provide sufficient floor space to hold all the elements of the IT infrastructure that an organization wants to deploy now and in the future. To

save costs and for security reasons, DC buildings are usually located outside central business district areas. These buildings do, however, need good access to general infrastructure in the form of roads, electricity, and telecommunications.

- **Facility management:** A building management system (BMS) allows the management of DC facilities from a centralized console. All critical data, such as temperature, humidity, power, and ventilation flow for each room, are measured and displayed on dashboards. The BMS also manages access control, security logging, and video surveillance of the DC rooms and the perimeter.
- **Infrastructure:** IT infrastructure refers to the IT components deployed into the DC, including the network equipment, server systems, and storage systems, all mounted into racks and the cabling connecting the systems. The elements of the infrastructure will be covered in detail in the following chapters.
- **Cooling:** The enormous amount of power delivered to a data center is primarily converted into heat, which must be removed from the IT infrastructure through cooling, as all IT equipment has a defined temperature range.
- **Power:** The DC requires sufficient power to operate the facility and the IT systems. The power consumption of a DC facility can be as much as 100 megawatts; hence, the local power provider requires an adequate power connection. It is vital to get power for the DC

with high reliability. Even though DCs usually have supplementary power systems in the form of an uninterruptible power supply (UPS), the stability of supplied power is an essential factor in securing highly reliable DC services.

- **Security:** A DC faces a range of threats to the security, which will be covered in Chapter 7. However, in the context of the DC facility, security mainly focuses on the building security itself and access control to the facility.

Before we look further into these factors, let's consider some aspects that must be taken into consideration when designing a data center facility, which have led to the creation of established standards that can be used as templates.

2.1 Data Center Facility Design

No DC design or construction standards are defined that must be followed. Instead, an organization's business and objectives determine a DC's design. For example, a private DC for a bank, a public CSP, and a university research lab will have different requirements and objectives to run their DC. The specific DC design must therefore follow the design objective of the specific organization.

However, organizations have built DCs for more than 60 years and have developed and evolved best practices to plan and run a DC covering the following building blocks:

- Conceptual design
- Determining layout and space
- Physical building requirements

- Building access control
- Building infrastructure, including plumbing, electricity, and fire alert systems

2.2 Established Standards for Data Centers

Data centers, integral to modern digital infrastructures, must adhere to various established standards. These standards provide templates for design, construction, and operation, ensuring data centers meet specific efficiency, reliability, and performance criteria. We will discuss some of the most applied standards in the industry.

Uptime Institute Tier Standard

The Uptime Institute Tier Standard is a commonly known standard to build DCs, and it is used to design, construct, and commission a DC. The Uptime Institute Tier Standard focuses on data center design, construction, and commissioning, defining four reliability levels, as seen in Table 2-1. In addition, IT professionals can obtain a certificate from the BICSI organization, qualifying them to design DCs and observe an organization's rules.

Table 2-1. Comparison of Data Center Tiers from the Uptime Institute

Requirements	Tier 1	Tier 2	Tier 3	Tier 4
Uptime guarantee	99.671%	99.741%	99.982%	99.995%
Downtime p.a.	<28.8 hours	<22 hours	<1.6 hours	<26.3 minutes
Component redundancy	None	Partial power and cooling redundancy (partial N+1 ¹)	Full N+1 ¹	Fault-tolerant (2N ² or 2N+1 ³)
Concurrently maintainable	No	No	Partially	Yes
Staffing	None	One shift	One+ shift	24/7/365
Typical customer	Small companies and startups	SMBs	Growing and large businesses	Government entities and enterprises
Usage	Budget-friendly data center tier	Good balance between cost and performance	High performance and affordability	High traffic and processing demands with fault tolerance
Requirements	Tier 1	Tier 2	Tier 3	Tier 4
Uptime guarantee	99.671%	99.741%	99.982%	99.995%

¹“N+1” denotes having one additional redundant component beyond what is necessary for normal operation (N). This redundancy ensures uninterrupted functionality by providing a backup if one component fails.

²“2N” signifies having double the necessary components for normal operation, ensuring a higher level of redundancy and fault tolerance.

³“2N+1” indicates having double the required components for normal operation, plus one additional spare component, enhancing redundancy and fault tolerance even further.

EN 50600 Series

These standards focus on data center design for IT cable and the network, emphasizing infrastructure redundancy and reliability, similar to the Uptime Institute's Tier Standard. They provide concepts and guidelines for efficient cable management, network design, and redundancy measures to ensure high performance and availability. Compliance with these standards promotes standardized design practices and enhances overall data center performance and customer satisfaction.

ASHRAE

ASHRAE is a standard that offers best practices for heating, ventilation, and air conditioning (HVAC) in buildings, although not specific to data centers. They provide guidelines for optimal HVAC performance, equipment selection, airflow management, and energy efficiency. Applying these standards ensures proper environmental conditions and energy conservation in data centers.

Other Examples

Other known standards include the following:

- ISO 9000 for operational quality
- ISO 14000 for environmental management
- ISO 28001 for information security regulatory standards, such as HIPAA
- Sarbanes-Oxley Act
- SAS 80 Type I or II
- Gramm-Leach-Bliley Act

- Payment Card Industry Data Security Standard for payment card security
- EN 50600-2-6 regarding management and operational information

Following the defined standards helps to ensure state-of-the-art data center design, construction, and operation of DCs. The design and construction process should be well documented and forms the foundation to build techniques for the security and resiliency of a DC.

2.3 Space

Physical Space

During the design phase, a DC is just an open space to host computer equipment; that is, a data center is a carefully prepared warehouse where IT systems are hosted and operated. Sizing a DC is often more complicated than it seems. However, it should have sufficient space for today's demand and allow for future expansion without being designed too big, thereby straining the organization's IT budget.

Several parameters change the physical demand for DC space over time:

- Changes in the business of an organization or adaptations to market changes.
- IT infrastructure evolves and becomes smaller for a given workload.
- IT technology such as punch cards, tape drives, or DVDs becomes obsolete.
- Organizations may decide to insource or outsource the function of the IT applications or move the workload to CSPs.

Lighting

No light is required in the DC during normal operations. Lights will only be turned on when equipment needs to be extended or removed or when maintenance work is performed.

Noise

The noise in large DCs can reach unhealthy levels for humans when hundreds of computer systems run their fans at high speed. It is therefore advisable that staff entering such a DC wear a noise protection headset.

Weight

The DC design must consider the weight of the equipment it will deploy. DC rooms are designed to carry a higher weight per square meter than most industrial buildings. Beyond weight, the construction planning needs to allow enough space to handle cooling airflows.

2.4 Facility Management

Key environmental data related to the DC must be measured and observed. Modern DC facilities use BMS to monitor all relevant data in a centralized system from various sensors measuring temperature, humidity, smoke, and flood indications. Any alarm will notify the building managers and the head of IT operations to allow quick reactions to the potential emergency.

2.5 Infrastructure

Racks

In today's DCs, the IT equipment is mounted in racks, which is an empty metal frame with standard spacing and mounting options (see Figure 2-2). It can hold standardized rack-mountable equipment, including computer systems, storage systems, networking equipment, cabling, and auxiliary power systems in the form of UPS devices, plus keyboards and monitors for the DC operation team.



Figure 2-2. Modern DC Racks (svstudioart, 2023)

2.6 Cooling

Temperature

DC temperatures are kept at low levels to keep the computer systems in their defined temperature range. These low temperatures are often uncomfortable for people performing maintenance work.

Cooling

To a large degree, the high levels of DC power consumption are converted into an undesirable by-product, namely, heat. The heat must be removed from the DC to keep the temperatures within the defined level. Ventilation and cooling are therefore critical components of a DC and must be designed by the planned IT capacity. Calculating the power consumption of a data center is not a simple exercise, as we'll see. Heat is directly related to power consumption, so the planning for the cooling and ventilation systems faces similar challenges. For example, if the cooling system is designed to be too small, the data center can't keep up with the heat dissipation of the computer systems or may run at its limits, so that DC expansion cannot be performed. On the other side, if the cooling system is over-dimensioned, the costs of the DC will not meet best-in-class standards.

Cooling DCs differs entirely from cooling offices and homes, where warm air is relatively evenly distributed. In general, the racks in DCs, and in particular frames holding computer systems, create heat in a very concentrated form. The airflow in a DC must be separated to maintain the room temperature at comfortable levels for humans and the airflow to cool the IT equipment. Thus, cool and hot air around the IT equipment must have different airflow paths. Modern DC design solves this problem by creating aisles of hot and cold air. Figure 2-3 is a good example on how

aligning the racks in rows with aisles in between to separate the cold and hot air provides adequate airflow. Cold air is supplied to the front row of the racks, while hot air is collected from the rear side of the racks.



Figure 2-3. *Data Center Layout Provides Adequate Airflow (Florian Hirzinger, 2009)*

A new cooling technique has evolved by using liquid cooling. This emerging technology is mainly used for systems that create very high heat levels, such as high-performance computer systems (HPC). With liquid cooling, the computer system is immersed in an electrically neutral liquid such as mineral oil. The liquid cooling process is much more effective than any air cooling because liquid transmits heat much better than air. When this technology matures and is more widely used, it can significantly improve the energy efficiency of DC and make DCs more sustainable.

2.7 Power

Power Demands of DC

The power supply is a crucial aspect of a data center plan, as it typically constitutes the largest item in the list of operational expenses. Large, modern data centers can consume over 100 megawatts, equivalent to the consumption of 80,000 households. Therefore, data center planners must consider the following parameters for the power supply:

- First and foremost, the local utility must provide adequate power to run the data center.
- As power consumption is the most significant item of a DC operational expense, the DC should be in an area where power is less expensive.
- Power must be “electrically clean.” This means it should be free of surges, spikes, noise, or voltage fluctuations, which are often indicators of the capacity constraints of a power grid.
- Reliable power is very important. Brownouts, fluctuations in the voltage levels, blackouts, total power disruptions, or other disruptions distract a smooth DC operation.

A business must be able to estimate its future power consumption during the planning phase. The power consumption planning for the DC facilities can then be easily calculated by adding the required elements for lighting, ventilation, and cooling. Estimating the power consumption of IT systems, however, is a difficult task. Power consumption of computer systems and, to a lesser extent, consumption of storage systems and

network equipment fluctuate with the workloads. The classic approach to calculating power consumption is based on average rack consumption, which often leads to inaccurate power consumption estimates. A modern approach is based on actual power measurements per server. The data is delivered by a power-handling device related to the computer system. Such intelligent power distribution units (PDUs) are deployed in each rack and create data to measure the power consumption of a specific workload.

Power outages happen even in countries with stable energy grids; therefore, a DC must have options to access power from a redundant power line or backup power. The best practice approach is to use several layers of secondary power, allowing it to react to different power interruptions. For example, racks can incorporate UPS devices for the IT infrastructure, covering short-term outages. Diesel- or gas-powered backup generators are required for the facility ensuring it remains operational even in more extended outages.

Uninterruptible Power Supply (UPS)

UPS plays a critical role in larger DCs by providing a continuous, reliable power supply to essential IT equipment, as seen in Figure 2-4. In addition, UPS systems are designed to protect the data center against power outages, voltage fluctuations, and other power-related disturbances that could cause damage to the IT equipment, resulting in data loss or downtime.

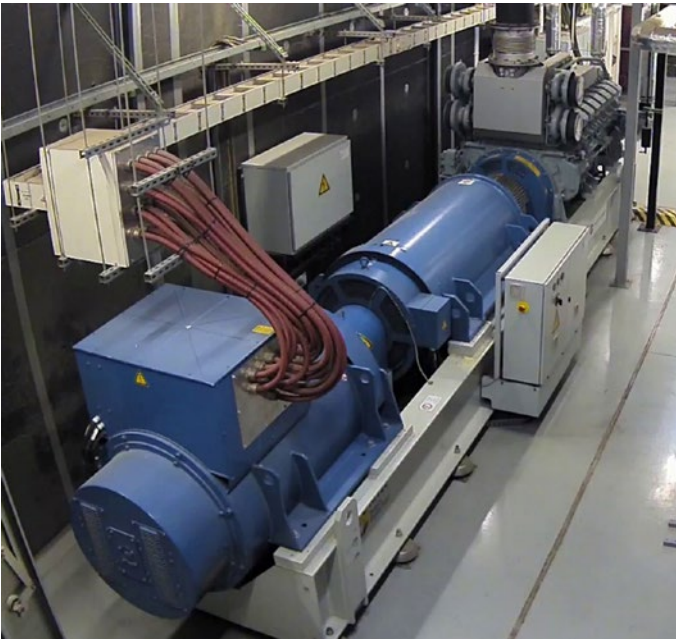


Figure 2-4. UPS for a Data Center (Bercik, 2020)

In a larger data center, UPS systems are typically deployed as a centralized system, often located in a dedicated room or area. The UPS system is connected to the main power supply of the data center. It provides a backup power source that automatically kicks in during a power outage. The role of UPS in a larger data center can be broken down into several vital functions:

Continuous power supply: The primary function of UPS is to provide constant and reliable power supply to critical IT equipment. This ensures that the IT equipment is protected against power outages, voltage fluctuations, and other power-related disturbances that could cause damage to the equipment or result in data loss or downtime.

Power conditioning: UPS systems also perform power conditioning, which involves regulating the voltage and frequency of the power supply to ensure that it meets the requirements of the IT equipment. This helps to protect against power surges and other power-related disturbances that could damage the equipment.

Battery backup: UPS systems typically include battery backup that can provide power for a certain amount of time in the event of a power outage. This allows the IT equipment to continue running for a period, allowing data center staff to address the issue and act fast to restore power to the equipment.

Monitoring and alerting: UPS systems are typically equipped with monitoring and alerting capabilities that allow data center staff to keep track of the system's performance and receive alerts in the event of an issue. This enables a team to quickly address the problems to ensure that the UPS system always functions correctly.

Overall, the role of UPS in a larger data center is critical to guarantee the continuous and reliable operation of critical IT equipment. By providing constant power supply, power conditioning, battery backup, and monitoring and alerting capabilities, UPS helps to protect against power-related disturbances so that the data center can operate smoothly and efficiently.

2.8 Security

Access Control

Although security for DC is usually associated with cybersecurity, the most fundamental security consideration is access control. This covers the protection of the facility itself as well as the hosted IT infrastructure. The physical security infrastructure and the associated processes must ensure that only authorized personnel have access to the facility and its systems. In addition, all human activities inside and around the DC must be documented. A set of devices, rules, and processes can achieve this:

- All entry points to the facility and the DC are electronically protected, and access is only permitted through individual badges.
- Separate access protection to the DC itself should only be accessible by personnel with the appropriate authorization.
- All entries and exits by staff, visitors, and vendors to the facility and the DC are logged into a security monitoring system.
- Strict enforcement that visitors must be escorted by authorized staff.
- Video surveillance at all facility entry points and in the DC itself.
- Dedicated on-site security personnel.

2.9 Summary

In this chapter, the design and best practices for data center facilities were discussed. That includes key factors to consider, such as space, power, cooling, security, facility management, and infrastructure. It also addressed the importance of designing a data center facility based on the specific requirements and objectives of the organization. Established standards for data center design were introduced, including the Uptime Institute Tier Standard, EN 50600 series, and ASHRAE guidelines.

Also, various aspects of data center facilities, including physical space, lighting, temperature, noise, weight, racks, access control, facility management, and power demands were covered. This chapter explained the significance of power supply and the need for reliability and redundancy. Additionally, it discussed the cooling requirements for data centers and the use of liquid cooling technology. Finally, the role of UPS (uninterruptible power supply) systems in providing continuous and reliable power supply to critical IT equipment in data centers was highlighted.

In the next chapter...

- What makes a computer?
 - What is computer virtualization?
-

CHAPTER 3

Compute and Virtualization

A system that can be programmed to carry out arithmetic or logical operations is automatically called a computer. Over the last 80 years, computers have evolved from machines that were very expensive, slow, only able to calculate simple functions, cumbersome, and difficult to manage to a ubiquitous building block of almost every modern electronic device, machine, and system.

Most industrial and consumer products above a certain complexity use computers as control systems. It ranges from simple special-purpose devices like microwave ovens, washing machines, and refrigerators to modern robots, autonomous vehicles, and general-purpose devices like personal computers, smartphones, and compute servers. Computers are the key building block of the Internet, connecting billions of people and devices today.

3.1 Hardware

The hardware of a computer system comprises all its physical elements:

CPU	Central processing unit, the actual compute engine, today deployed as microprocessors
Memory	Semiconductors to store programs and data for the CPU
Storage	Devices to store programs and data permanently
Power Supply	Devices to provide electricity to all computer elements
Connectivity	Enables the computer to connect with LAN, Wi-Fi, and Bluetooth networks, as well as to other external devices via USB or other protocols

Most of these components are included in the computer's motherboard. This piece of hardware is represented in Figure 3-1.

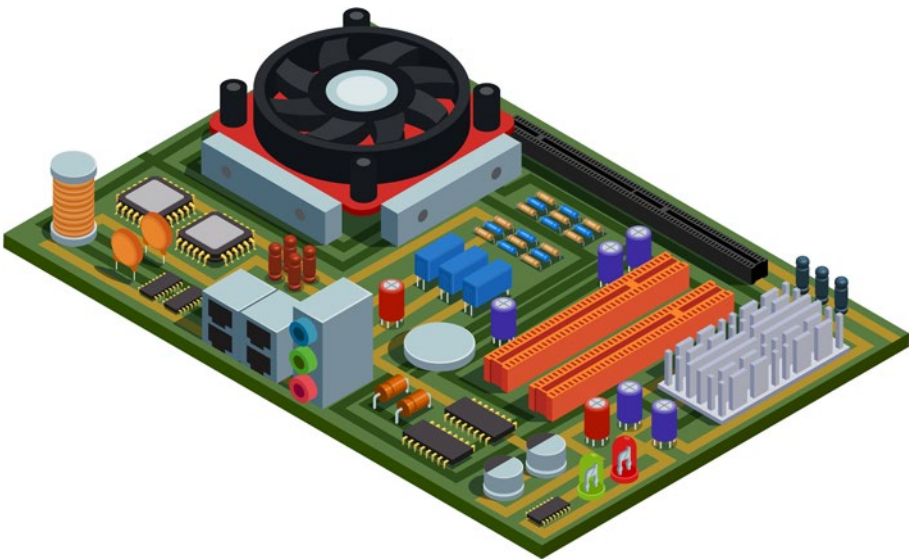


Figure 3-1. *Motherboard of a Computer System (Macrovector, 2023)*

3.2 Software

“Computer software” is a collective term that encompasses the entirety of digital instructions that can be loaded and executed by a computer system. This includes everything from simple scripts to complex programs. At its core, software is what allows hardware to perform tasks and functions.

One of the most fundamental pieces of software within any computer system is the operating system (OS). The OS acts as an intermediary between the computer hardware and the computer user. Its primary function is to control and coordinate the use of hardware resources, ensuring smooth operation. When a user interacts with a computer – whether it’s typing on a keyboard, clicking with a mouse, or running an application – it’s the operating system that interprets these actions and translates them into instructions that the hardware can understand. In essence, the OS serves as a bridge, allowing software applications to function on diverse hardware platforms without the need for developers to understand the intricacies of each specific hardware type.

Software applications, often simply referred to as “apps,” are designed to perform specific tasks or solve particular problems. The scope of these tasks varies widely, from basic functions like drafting a document in a word processor or calculating a spreadsheet to more advanced tasks such as simulating real-world environments for research, piloting drones, or even facilitating augmented reality experiences. Some applications, like those used for steering an autonomous car or predicting weekly weather patterns, employ advanced algorithms and require significant computational power.

The creation of these software applications hinges on the detailed description of tasks in the form of algorithms. An algorithm, in the realm of computer science, is a finite set of well-defined instructions for completing a task or solving a problem. Once an algorithm is conceptualized, it’s up to programmers to translate it into a computer-readable format. This translation is done using programming languages. There’s a multitude of

languages available, each with its unique strengths, characteristics, and applications. Popular languages like C, Java, and Python offer a blend of flexibility, efficiency, and user-friendliness, making them favored choices among developers worldwide.

As technology continues to advance, the relationship between software and hardware will only deepen, enabling even more complex and transformative applications that can reshape how we interact with the digital world.

3.3 Types of Computer Systems

1. **Personal systems:** These are devices primarily designed for individual use. They cater to personal computing needs, entertainment, and communication.
 - **PC (personal computer):** This usually refers to desktop computers used at homes or offices.
 - **Laptop:** Portable computers that combine the monitor, keyboard, mouse (or touchpad), and the computer itself into a single unit.
 - **Tablet:** A portable device that typically uses a touch interface. It's larger than a mobile phone but offers similar functionality.
 - **Mobile phone:** Primarily used for communication, modern smartphones can also run a variety of applications, making them mini-computers in their own right.

- **Smartwatch:** Wearable devices that offer features like notifications, health tracking, and even some lightweight apps.
2. **Embedded systems:** These are computers integrated into other devices, primarily to control certain functionalities of that device.
- **Cars:** Modern vehicles contain multiple embedded systems for functions like entertainment, navigation, and safety.
 - **Microwaves and refrigerators:** These home appliances often have simple embedded systems for control functions.
 - **Robots:** May contain more advanced embedded systems to control movement, recognize patterns, or perform specific tasks.
 - **Drones:** Use embedded systems for navigation, control, and, sometimes, real-time communication.
3. **Server:** These are powerful machines designed to provide services, data, and resources to other computers, usually over a network.
- **Standalone:** Typical tower servers that are not mounted in racks.
 - **Rackmount:** Designed to be mounted in standard racks to optimize space.
 - **Blades:** A server architecture that houses multiple server modules (“blades”) in a single chassis. They share resources like power and cooling.

4. **Data centers:** Massive facilities used to house large numbers of servers and related equipment. They provide the backbone for many online services and cloud computing.
 - **Private DC (data center):** Owned and operated by single businesses and caters exclusively to their needs.
 - **Public DC (like clouds):** Provides services to multiple clients or businesses. Cloud providers like Amazon AWS, Google Cloud, and Microsoft Azure fall under this category.
 - **Leased DC:** Companies might lease space or infrastructure within a data center rather than owning and maintaining their own.

3.4 Purpose of Computer Systems

Computer systems can be one of the following.

General-Purpose Computer

Most computer systems deployed today are general-purpose computers using microprocessors from key semiconductor manufacturers such as Intel, AMD, or TCM.

In the early era of computing (between 1960 and 1990), general-purpose computers were not fast enough for certain types of applications. In 1965, Gordon Moore, an engineer from Intel, recognized that the number of transistors on a CPU board had doubled every two years, and he predicted that this trend would continue.

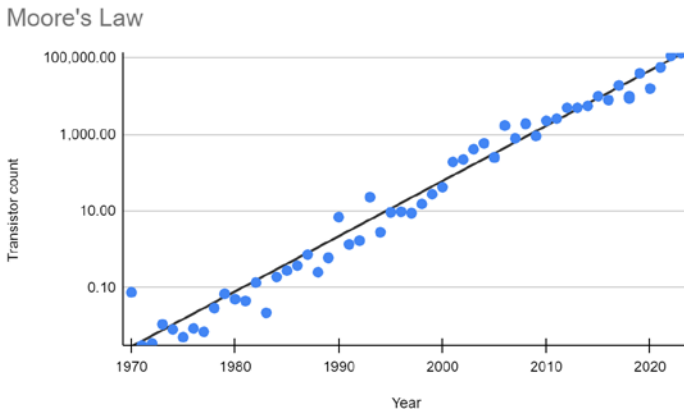


Figure 3-2. *Moore's Law (Source: Wikipedia, 2023)*

This prediction is known as “Moore’s law,” and today, 58 years later, it is still valid (see Figure 3-2). The actual compute performance increased even at a higher speed than many other parameters, such as the CPU frequency that was raised at the same time and thereby leads to a doubling of CPU performance every 18 months, or thousand times in 15 years, a million times in 30 years, and a trillion times in 60 years. The computer industry’s rapid advancement in its relatively brief history is unparalleled in most other technological domains. For years, Moore’s law has consistently held true, forecasting the pace at which this evolution has unfolded. However, as transistors edge closer to atomic scales, they face increasing physical and economic challenges, casting doubts on the continued applicability of Moore’s law in its traditional sense. Nonetheless, groundbreaking innovations and emergent computing paradigms, such as quantum computing, may reshape the future trajectory of computational advancements.

Through this evolution, modern general-purpose microprocessors have become so fast that most computer applications run fast enough on general-purpose CPUs, and there are only a few application areas where specialized computers are still required.

Specialized Computer Systems

In the early phases of computer evolution, computers were used to run simple applications mainly to automate companies' finance and administration processes, such as payrolls. The performance for floating-point computation, however, needed to be improved to run complex applications like weather forecasts or flight simulators. In the 1970s, specialized computing for high-end floating-point calculations was developed by vendors like FPS, Convex, and Cray. These systems were called supercomputers. Real-time computing was another domain where special computers were built. To control machines and certain processes, the computer system must be able to gather data in real time and react to it quickly. In the 1980s, specialized systems were created for image processing and artificial intelligence.

Over time, the performance evolution of general-purpose computers following Moore's law made most of those specialized computer systems obsolete. Even applications at the edge of computing, like high-performance computers, are nowadays deployed with large arrays of standard general-purpose microprocessors. In some domains, however, additional compute elements are used if standard computer performance is insufficient, such as the graphic processor in high-end gaming consoles.

The latest field of specialized computers that promises to leapfrog current general-purpose computers is quantum computers that will outperform standard computers by factors of thousands for certain applications such as encryption when they become commercially available in the coming years.

3.5 Data Centers

As we've seen in Chapter 2, a data center is a dedicated space within a building, an entire building, or several buildings. It houses an IT infrastructure with the following components:

- Computers
- Storage systems
- Network systems

Data centers emerged in the 1960s when companies started automating standard processes on computers. All data centers were private, usually a room or a floor of a company's building, and were part of a company's assets. Today, most companies still maintain their own data centers. But with the rise of the Internet, companies like Amazon and Microsoft began to build very large DCs and started offering computer resources to enterprises.



Figure 3-3. Cloud Data Center (Benzoix, 2023)

Such large publicly available DCs are called “clouds.” Cloud services were initially mainly used for development and testing as it was easy to add compute resource quickly, shortcutting the slow internal investment processes that most organizations must follow. As the demand for such services grew, cloud DCs became large and began to develop the cost-of-scale benefits that modern cloud systems can offer. Today, many cloud service providers (CSP) offer cost-effective and easy-to-provision compute resources to any company, resulting in a continuously rising share of compute capacity in the cloud DCs, while the capacity in private DCs is stagnating and even shrinking. Figure 3-3 represents what a cloud data center looks like.

For much of the early history of computing, computer systems were primarily employed by companies for fundamental business support processes, such as accounting and payroll management. However, with the evolution of the Internet, these systems now underpin nearly every facet of modern businesses. This encompasses everything from sales and marketing websites to internal communication platforms, as well as comprehensive design and manufacturing processes. A new wave of companies has emerged, wherein their primary business model revolves around applications hosted in the cloud. Prominent examples of such companies include Salesforce.com, LinkedIn, Uber, and Airbnb, among others.

As a result, IT operations have become indispensable for ensuring business continuity. Typically, this involves integrating redundant or backup components and infrastructure, encompassing power supplies, data communication connections, cooling systems, and security measures. Large-scale data centers can rival small cities in terms of electricity consumption. Presently, international standards like EN 50600 and ISO 22248 exist to classify data centers in the realm of information technology.

Data center facilities and infrastructures are the following:

- *Class 1* single path solution
- *Class 2* single path with redundancy solution
- *Class 3* multiple paths that provide a concurrent repair/operation solution
- *Class 4* multiple paths that provide a fault-tolerant solution (except during maintenance)

3.6 Compute Resources

Compute Building Blocks

Data centers (DC) can be built in many different forms. Some DCs are built by using a unified compute platform, such as blades or rackmount servers. Other DCs have deployed different types of computers, including

1. Rackmount server in different performance tiers (low, medium, high)
2. Blade server
3. High-performance computer

Rackmount Server

A rackmount server, or rack server, is any server that is specifically designed to fit and be mounted within a server rack. Most rack servers are general-purpose computers deployed in data center environments. Depending on the design, the rack server can be secured into the rack using mounting screws or rails.

The size or “height” of the rack server, measured in rack units (Us), often corresponds to its capacity for components rather than its computational power. For example, a 1 U server might be optimized for space but can still be very performant. In contrast, larger servers, like those occupying 2 Us or more, offer room for additional components, such as extra CPUs or memory. To optimize space, servers are stacked vertically within a rack.

Benefits of a Rack Server

- **Powerful**

Rack servers have built-in components to operate as a stand-alone computer system. Rack servers are offered in a wide range from an entry-level system to larger, more powerful high-end systems.

- **Convenience**

The server can be easily mounted into a rack, which is convenient and saves a lot of space, especially when compared to a traditional tower-style server.

- **Cooling**

Each rack server has its own cooling resources. They are usually equipped with internal fans. Placing them in a rack also increases airflow.

In many DCs, the rackmount servers are the main building block of compute resources. These are standard servers usually equipped with microprocessors from Intel or AMD. Larger racks with a height of about 2 m have a capacity of 42 Us and can hold up to 42 rackmount servers with a height of 1 U each (see a representation of this in Figure 3-4).



Figure 3-4. *DC Racks with Rackmount Server (Macrovector, 2023)*

Using a unified rackmount server system for a large cloud DC offers an enormous economy of scales. CSPs can buy large quantities from branded vendors such as HPE and Dell at a very favorable price point. Some CSPs use nonbranded servers, and some buy compute boards from component vendors such as TCM and assemble their own servers.

As each server is an independent entity, the networking of the rack must connect each server with sufficient bandwidth and redundant connectivity to avoid a single point of failure.

This strategy allows a simple deployment of the compute resources, and the compute density is very high. A modern DC room of 2,000 square meters can accommodate 800 racks with 20,000–25,000 servers, plus storage and network equipment. A large DC can have 20 or more such rooms, with a total capacity of more than 500,000 servers.

Blade Server

A blade server is a modular compute system that allows multiple units to be housed in a small area. The blade servers are thin and consist only of CPUs, memory, and integrated network controllers. Storage is usually separated; some blade vendors offer storage blades with the same physical properties as the blade server to fit into the chassis. See an example of a blade enclosure in Figure 3-5.



Figure 3-5. Blade Enclosure (ColossusCloud, 2016)

One large chassis will be mounted into a server rack for a larger configuration. The blade servers are added to the chassis. The chassis can then provide power to all blade servers and connect them within the blade chassis.

The shared resources in a blade chassis allow each blade server to operate efficiently, requiring fewer internal components. Blade servers are generally used for applications with high-performance requirements.

A key advantage of blade servers is their rapid serviceability by allowing components to be swapped out without taking the blade chassis offline. Blade architecture features a much higher processor density than a rackmount server, but it requires the facility to support a much higher thermal and electrical load per square foot.

Benefits of a Blade Server

- **Power consumption**

The total power consumption is lower than a rackmount server as the chassis for the blade server supplies the power to multiple servers.

- **Hot swappable**

Blade servers are replaceable without a shutdown of the blade chassis when the blade chassis is operating in hot-swappable mode. This helps to reduce downtime in case of system failures.

- **Less need for cables**

Rackmount servers require individual cables for each server. In contrast, blade servers can have only one cable, usually a fiber cable, connected to the chassis, thus significantly reducing the number of cables in a rack.

- **High processing power density**

Blade servers can provide extremely high processing power while taking up minimal space.

High-Performance Computer or Supercomputer

In the past, high-end computer systems were called supercomputers. Now, the term high-performance computer (HPC) is more common. HPCs are designed to compute complex algorithms for floating-point calculations. Hence, the performance of HPC is mainly measured and benchmarked in floating-point operations per second (FLOPS). For general-purpose computers, the number of instructions per second is the common performance index measured in million instructions per second (MIPS). To benchmark the computer system's floating-point performance, a standardized benchmark is used, called LINPACK.

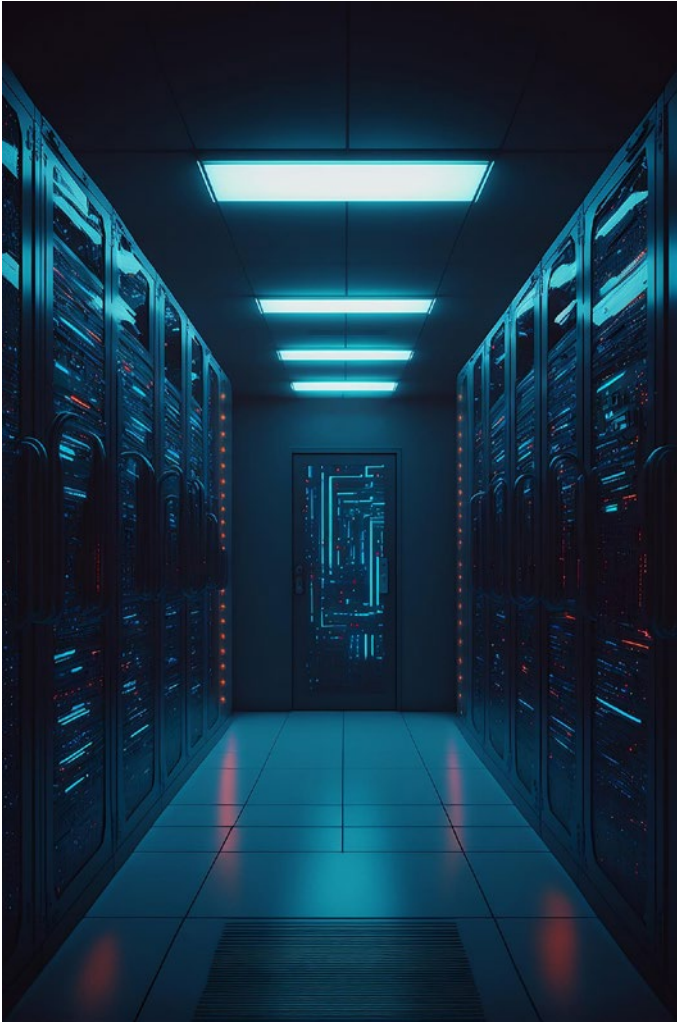


Figure 3-6. *High-Performance Computer (HPC) (Svstudioart, 2023)*

Figure 3-6 shows a representation of a high-performance computer (HPC). Since 2018, these computing systems have surpassed a significant milestone: 100 petaflops. To put this in perspective, 1 petaflop equates to 1 quadrillion (1,000,000,000,000,000) floating-point operations per second (FLOPS). In comparison, an average desktop computer typically

performs at a rate of hundreds of gigaflops, where 1 gigaflop is a billion (1,000,000,000) flops. This means the world's leading HPC systems are approximately a million times more powerful than a standard desktop. The majority of HPC systems operate on Linux-based platforms.

Table 3-1 shows the top five HPC systems in the world as of November 2021. The performance in the second column is measured in the number of LINPACK-based floating-point calculations in petaflops.

Table 3-1. *Top Five HPC Systems Based on LINPACK Benchmark*
(Source: top500.org, 2023)

Rank	Rmax (PFlop/s)	Name	Characteristics	Site – country	Year
1	1194	Frontier	HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot 11, HPE	DOE/SC/Oak Ridge National Laboratory United States	2021
2	442.01	Supercomputer Fugaku	Supercomputer Fugaku, A64FX 48C 2.2GHz, Tofu interconnect D, Fujitsu	RIKEN Center for Computational Science Japan	2020
3	309.10	LUMI	HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot 11, HPE	EuroHPC/CSC Finland	2023

(continued)

Table 3-1. (continued)

Rank	Rmax (PFlop/s)	Name	Characteristics	Site – country	Year
4	238.70	Leonardo	BullSequana XH2000, Xeon Platinum 8358 32C 2.6GHz, NVIDIA A100 SXM4 64GB, Quad-rail NVIDIA HDR100 Infiniband, Atos	EuroHPC/ CINECA Italy	2023
5	148.60	Summit	IBM Power System AC922, IBM POWER9 22C 3.07GHz, NVIDIA Volta GV100, Dual- rail Mellanox EDR Infiniband, IBM	DOE/SC/Oak Ridge National Laboratory United States	2018

HPC systems are mainly used in the field of computational science. These applications are computationally intensive tasks in various fields, including

- Quantum mechanics
- Weather forecast
- Climate research
- Oil and gas exploration
- Molecular modeling, which includes computing the structures and properties of various chemical compounds

- Analyzing biological macromolecules, polymers, and crystals
- Physical simulations, spanning areas like astrophysics and the aerodynamics of airplanes and spacecraft
- Cryptanalysis

3.7 Computer Operating Systems

Key Functions of an Operating System

An operating system (OS) serves as a fundamental component of computer systems. Its primary role is to manage and facilitate interactions between the computer's hardware – including the central processing unit (CPU), cache memory, main memory, and I/O devices like disks, screens, keyboards, and network controllers. This management ensures that applications can seamlessly access and utilize these hardware components.

Today, operating systems aren't just confined to traditional computers. Many household appliances, such as refrigerators, washing machines, and dishwashers, now incorporate microprocessors. When you select a function like “quick wash” on a washing machine, you're essentially initiating an application. This “quick-wash” application communicates with the washing machine's microprocessor via its embedded operating system. While the OS in household devices tends to be straightforward, at the opposite spectrum lie the complex operating systems and load-balancing software in cloud data centers. These sophisticated systems orchestrate and route the workloads of myriad applications across hundreds of thousands of computer systems.

The building blocks of an operating system include

- Process management
- Memory management
- File management
- I/O management
- Security

History of Operating Systems

The history of operating systems can be divided into the following periods:

- **1950s and 1960s**

The area of batch processing systems: During this period, computer systems were controlled by applications and operating systems loaded from punch cards and paper stripes. Creating and running an application was a slow and tedious process, as it required one punch card per line of code. All programs were run in batch mode, which means the computer systems were processing one application at a time and each “compute job” was queued into a batch system. During this time, computing was dominated by IBM, and the most used OS was OS/360 from IBM.

- **1970s**

In the late 1960s, the foundations of time-sharing systems were laid, leading to the emergence of time-shared computer systems by the 1970s. This advancement ushered in an era where, for the first time, programmers could interactively create and

modify their code on screens, significantly boosting their productivity. Initially, due to the limited performance of these computer systems, applications largely operated in batch mode. However, by the latter part of the 1970s, technological enhancements enabled computer systems to support direct on-screen editing and execution of software. These systems could accommodate multiple programmers and users simultaneously. Large companies predominantly sourced their IT systems from giants like IBM and Digital Equipment, with MVS (from IBM) and VMS (from Digital Equipment) emerging as the most prevalent operating systems of that period.

- **1980s**

Time-sharing systems, which enabled programmers to execute applications directly on their screens, paved the way for the development of multitasking systems. These advanced systems could simultaneously run multiple applications. The 1980s marked the rise of workstations, with Apollo and Sun Microsystems leading the charge. These companies introduced specialized personal computers for engineers, aptly termed “workstations.” Equipped with new operating systems and enhanced networking capabilities – Apollo’s “Domain” and Sun’s “Solaris” – these workstations transformed collaborative engineering, allowing professionals to operate applications collaboratively as though they were on a single system.

- **1990s**

The workstation for engineers evolved into PCs for all office workers and later even for consumers.

In addition to computer specialists, some scientists acquired a foundational understanding of computing and were able to work with the operating systems (OS) used during this period.

A major milestone in democratizing computer systems was the creation of graphical user interfaces (GUIs). Early developments by Xerox did not become widely used. Only when Apple and Microsoft adopted GUIs into their operating systems (MS Windows and Apple macOS), computer systems in the form of PCs started to become widely used. They evolved to the systems that we use today.

- **2000s**

Until the late 1990s, the concept of a computer system was based on a centralized approach. With the emergence of workstations and PCs, the notion shifted to a more decentralized approach, but it wasn't more than client/server computing. PCs and workstations were used to perform tasks locally. Still, most organizations would run all their key business applications on their central computer system, which soon became the bottleneck for the application performance. The approach of networked systems started to address this issue. In the 2000s, network operating systems evolved, allowing multiple computers to connect and share resources.

- **2010s**

Starting in the 1990s, mobile phones became more popular. While the initial purpose of mobile phones was to make phone calls while traveling, they gradually became devices that could do more. The first data application was SMS. From there, more and more applications were developed, and mobile phones increasingly became the ultimate device that allowed us to access the Internet from anywhere at any time. The evolution was enabled by the development of increasingly sophisticated mobile phone operating systems, such as iOS and Android.

The most popular operating systems in use today are Windows, macOS, and Linux for computer systems, while iOS and Android are the dominant OS for mobile phones.

Operating System Functions

- **Booting the computer**

Every computer journey begins with the booting process. The motherboard's firmware, known as BIOS (Basic Input/Output System) or UEFI (Unified Extensible Firmware Interface) in modern systems, oversees this. Initially, it performs a POST (Power-On Self-Test) to verify the system's hardware integrity. Once cleared, the BIOS/UEFI locates the boot loader on the storage medium (like an SSD or HDD), which then loads the operating system into memory, paving the way for user interaction.

- **Processor management**

In the heart of every computer lies the CPU (central processing unit), orchestrating every command and application. Modern systems often multitask, meaning they run multiple applications simultaneously. The OS efficiently manages this using “process scheduling.” It utilizes algorithms to determine the execution sequence, ensuring that high-priority tasks receive more immediate attention and resources than lower-priority ones. The OS also ensures equitable distribution of CPU time, maintaining system responsiveness and stability.

- **Memory management**

Every byte of data, every application, and even the OS itself reside within a computer’s memory at some point. The OS supervises this using memory management techniques. This includes “paging” and “segmentation,” strategies that allow the OS to utilize memory efficiently. The OS also manages “virtual memory,” a technique that uses a section of the hard drive to emulate RAM, ensuring that applications run smoothly even if physical RAM is fully utilized.

- **I/O system and driver management**

A computer’s true prowess lies in its ability to communicate with a myriad of devices, from keyboards to high-end graphics cards. Every device speaks a different “language,” and drivers serve as translators. The OS manages these drivers, storing them, updating

them when necessary, and ensuring they function harmoniously. Features like “plug and play” simplify this, allowing users to connect a device and use it almost instantly, with the OS autonomously handling driver installation and configuration. See Figure 3-7 for a glimpse into the architecture of a typical operating system.

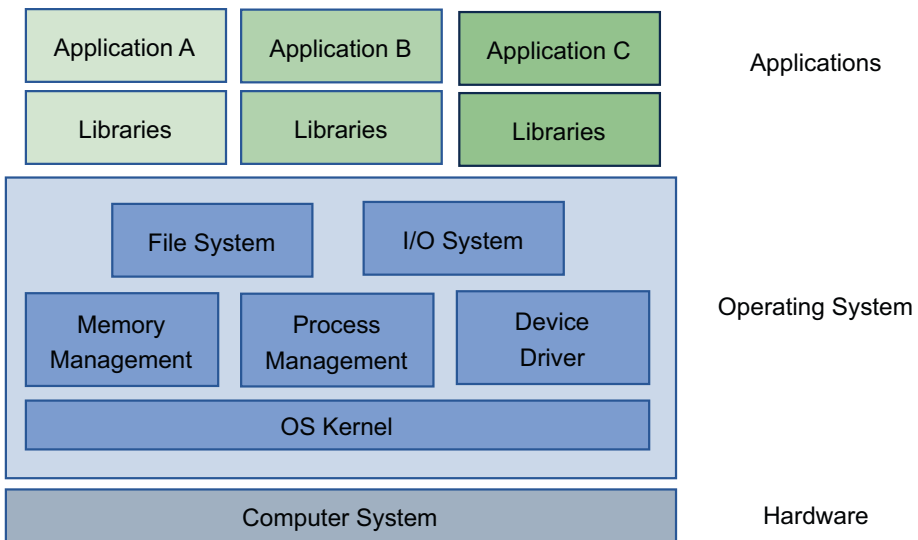


Figure 3-7. *Typical Operating System Architecture*

- **File system**

Imagine a library without a cataloging system – chaos! Computers store vast amounts of data, and the OS ensures this data is methodically organized and swiftly accessible. Whether it’s the FAT32, NTFS, ext4, or APFS file system, the OS tracks where every piece of data is stored, manages permissions, and ensures data integrity.

The file system on drives used to be only the classic spinning magnetic disks, also known as hard disk drives (HDDs). Nowadays, however, HDDs are often replaced by semiconductor-based flash memory, which is called a solid-state disk or SSD.

The OS file system organizes the available space in drives and directories. Directories contain objects including directories, files, and file links. These features enable users to build hierarchical data storage systems to manage files of applications and user data. The file system manages the resources of the devices as well as the access rights to an object that the system administrator defines.

- **Job control and accounting**

The system operator will use the functions of the OS to perform job control and accounting. The job control system allows the operator to define the sequence and priorities of processes that run on the computer systems. OS job accounting allows IT-related services to charge internally and externally. As an IT organization is usually operated as a cost center in an organization, a department will be charged according to the resources that were used by its members from the computer systems, such as computer time, disk-space usage, and pages printed. For CSPs, job accounting is their key tool to charge their clients for the services provided.

- **Resource monitoring**

The OS monitors all its resources and collects relevant data about the performance and resource allocation of applications. If certain processes abort frequently,

the system will alert the system operator to analyze the problem. The collected data provides the foundation for the system operator to analyze the usage pattern of applications and the system itself to identify problem causes of application errors. It also identifies overall issues with the operating system or the computer system and its peripherals.

- **Security**

Several functions in the operating system support the system operator and security specialist to keep the computer system secure. This includes functions such as access-control lists to files and applications, and user and password management.

- **OS command interface**

Each IT user needs to interact with the computer system. This function of the OS is called the command interface. Before graphical user interfaces were developed, the command interface was a text-based interface like the that in Microsoft DOS or Unix shell. With the development of graphical user interfaces (GUIs), systems like Apple Macintosh or Microsoft Windows offered graphical user interfaces. Instead of entering the command “del file” in DOS, or the “rm file” in Unix to delete a file, one could just click a file in Windows and move it to the waste basket. The development of GUI was the key milestone to make computer systems usable by non-IT experts.

Operating System Types

There are several types of operating systems, each designed for specific types of devices and applications. Here are some of the main types of operating systems.

Embedded Operating Systems

To control the functions of built-in computer systems, embedded operating systems are used, which are usually stripped-down versions of general-purpose operating systems. An embedded OS is mainly used in smaller, lightweight hardware platforms with limited functionalities and designed to serve the hardware of an embedded system. Examples of embedded systems are presented in Table 3-2.

Table 3-2. Embedded Systems

Segment	Embedded system		
Consumer goods	Washing machine	Dishwasher	Microwave
Automotive	ABS (anti-lock brake)	Airbags	Engine management
Medical devices	Pacemaker	Hearing aids	ECG monitors
Industrial application	Industrial sensor networks	Industrial control systems (ICS)	Programmable logic controller (PLC)

As embedded OS are often used to build machines, they usually provide real-time capabilities, required in military, manufacturing, and robotics applications. Embedded operating systems are specialized OS that are tailored for hardware platforms to perform specific tasks in embedded systems.

Real-Time Operating Systems

Most embedded OS have real-time features. But real-time features are required in many applications beyond embedded systems such as military equipment, airplanes, robotics, and complex machines. Real-time operating systems (RTOS) provide predictable and deterministic behavior for the computer systems that they control.

Key features of an RTOS include the following:

- **Scheduling of tasks**

The system operator assigns priorities to each task. The RTOS scheduler ensures that high-priority tasks are executed first and that lower-priority tasks will wait until the higher-priority tasks have completed their functions.

- **Task synchronization**

Synchronization and communication between tasks are very important features in real-time applications. RTOS provides synchronization functions to control the sequence of schemes and protocols for intertask communications.

- **Interrupt handling**

Another important feature to control the behavior of real-time applications is handling of interruptions. If a robot detects an obstacle in its way, it needs to stop the current movement immediately. This can be managed through the interrupt function of RTOS. The visual control system of the robot will send an interruption to the central control systems, which then will stop the motion control system immediately.

- **Timers and clocks**

In real-time applications, timers and an accurate clock are important. As tasks are often time driven, accurate and reliable timers and clocks are essential. RTOS provides functions to the application to measure time intervals and trigger real-time events.

- **Resource management**

To ensure that the real-time applications function according to their tasks, the computer must operate in a very reliable form. This requires the RTOS to manage key system resources efficiently, including memory, CPU time, and I/O devices.

- **Tools for profiling and debugging**

Testing of real-time applications is much more complex than applications with linear features. Therefore, real-time application developers need RTOS tools to profile the application. Debugging is another complex task for real-time applications as problems sometimes occur if a certain sequence of events has happened. Hence, the debugging functions of RTOS are important for testing during the development of applications and to figure out the cause if real-time applications create unexpected results.

RTOS are designed to ensure that real-time tasks are executed in a deterministic and timely manner controlled by the priority of tasks and that system resources are managed efficiently to optimize system performance.

Server Operating Systems

OS for server systems have some distinct characteristics that define them as server OS systems. These include the following:

- **Scalability**

Servers used in data centers of larger organizations or CSPs need to cope with the workload of a large number of concurrent users and processes. Therefore, a server OS system must be able to scale. This involves managing the key resources of a computer system in the form of CPUs, memory, and disks per the requirements of users and applications without compromising availability or performance. Scalability is one of the main factors to ensure an overall high performance of the systems.

- **Reliability**

High reliability is one of the key features of a server OS system, as it must be able to cope with events such as failures from hardware components, crashing applications, and users that make mistakes. This can be achieved by building many fault-tolerant features into the OS, such as disabling faulty elements like a disk drive or a memory board, running applications in an environment that is encapsulated from the OS, and many other internal OS recovery mechanisms. These features ensure that servers can operate continuously without any unplanned downtime.

- **Security**

Server OS systems must provide sufficient features to protect against many forms of threats from unauthorized access, data breaches, and other security threats. As security requires a holistic approach to the overall system, we have dedicated a full chapter to this topic (see Chapter 7).

- **Manageability**

Server OS systems feature several tools to monitor and configure the OS. It will include some basic maintenance tools to check the functions of the underlying hardware of the computer system. An important function of the OS is allowing fast upgrades to newer versions, the ability to return to a previous version, and the ability to include patches, which is a method to fix an existing problem without going through the process of a software upgrade.

Server operating systems are designed to meet the specific requirements of servers in DCs.

The main characteristics of scalability, reliability, security, and manageability ensure that servers operate continuously, using the resources of the underlying computer hardware in the most efficient way.

Distributed Operating Systems

Distributed operating systems (DOS) are designed for managing a network of computers. The DOS combine the networked computers into a single system. There is no central controlling element in a DOS; each computer system runs autonomously. But it provides means to communicate

and synchronize between the independent systems to share data and coordinate tasks. It is used in applications that are complex and require many resources, including the following:

- High-performance computing, mainly used for simulation and scientific modeling applications
- Large-scale data processing for data analytics, scientific computing, and AI applications
- Distributed databases
- Internet of Things
- Telecommunications, to manage large-scale network resources of the telecommunication network

Some of the key characteristics of distributed operating systems include the following:

- **Resource sharing**

All resources, including CPU, memory, and storage, are shared among the connected computer systems. The DOS provides commands and interfaces that allow the management and allocation of resources to run complex applications.

- **Concurrency**

The DOS allows applications to be distributed to many computer systems. The applications run in concurrent processes across multiple computers. The DOS provides a mechanism to manage access to shared resources and allows the synchronization of processes across the computer system used for the execution of the applications.

- **Transparency**

The system operator of large complex applications must be transparent about the resources that are in use to perform the applications and includes parameters such as location, access, and concurrency.

- **Fault tolerance**

The DOS must provide a certain degree of fault tolerance functions for the system. A complex application should continue to run – even if failures occur, such as faults from hardware elements, network connections, and errors in the application code. The DOS must be able to detect such problems and provide means to recover from those issues. Examples of fault tolerance in this context are re-assigning a subtask to another computer system if a hardware failure occurs or restarting an application if it has aborted.

- **Security**

All operating systems that manage DC resources require sufficient features to protect against many forms of threats from unauthorized access, data breaches, and other security threats. As security requires a holistic approach to the overall system, we have dedicated a full chapter to this topic (see [Chapter 7](#)).

Distributed operating systems are designed to manage compute resources for complex application that can also be performed on larger distributed systems. As we have considered, key features of a DOS include resource sharing, concurrency, transparency, fault tolerance, and security.

Cloud Operating Systems

Cloud operating systems are designed to manage cloud infrastructures. Cloud data centers can be very large, with 100,000s of computer systems, disk space measured in petabytes, and large complex internal data networks. The cloud operating system should manage large-scale resources and provide resources and concepts that are useful for the users of such cloud data center resources. The concepts of virtual machines and containers evolved with the progress of cloud operating systems.

Large cloud service providers such as Amazon (AWS), Microsoft (Azure), and Google (GCP) have developed their own cloud operating systems.

The features of a cloud operating system are very similar to those of a server and distributed operating system. But there is the notion of scale, which makes a difference between a server OS, a DOS, and a cloud OS that manifests in the following features:

- **Scalability**

The scalability of a cloud OS is on a much higher level than the other OS for server systems. Cloud operating systems are designed to be highly scalable, with the ability to scale resources up or down, depending on demand.

- **Elasticity**

Elasticity is a feature that is only available in a cloud OS.

This feature provides an automatic adjustment of resources based on workloads. System resources are only allocated when needed and will be released when they are not required any longer.

- **Resiliency**

Resiliency is, of course, a vital feature of a cloud OS. But in comparison to server and distributed operating systems, resiliency can be built on a higher level. For a server OS, it is important that a computer system will continue to run should there be an error in the hardware, network, or application.

In a DOS, a single computer system is not that critical anymore, as tasks can be redistributed from one computer system to another. For CSPs with data centers at many different locations, the outage of an entire DC can be tolerated, as workloads can be reallocated to another DC. Therefore, the functions of the various OS are similar, but the implementation can be very different as fault tolerance can be built on different levels in the system.

- **Security**

Security for CSP has many similar requirements as server OS and DOS, but the scale of CSP requires many additional security features. We cover this in a separate chapter (see Chapter 7).

- **Management**

Beyond the standard management tools for server OS and DOS, cloud OS provides management tools to manage and monitor IT resources across multiple data centers and locations.

In summary, cloud operating systems are designed to manage and monitor resources of the large-scale data centers of CSPs. They provide elasticity, scalability, resiliency, and security to manage the resources of a CSP.

Operating Systems for Client Systems

Users connect for their private and business needs to access applications through their mobile phones and personal computers. These end devices run their own OS. For tablets and mobile phones, Android from Google, iOS from Apple, and Windows Phone from Microsoft are the most used operating systems. For desktops and laptop computers, Windows from Microsoft, macOS from Apple, and Linux from Open Source are the most popular OS suppliers.

3.8 Compute Virtualization

Hardware Virtualization

Through virtualization technics, applications can run on virtual machines that act like real computer systems with an operating system. This concept of a virtual machine allows a hypervisor to distribute all applications to a pool of compute-server resources (see Figure 3-8). Failure recovery becomes simple through this technology, as the hypervisor can re-assign an application to a different compute resource if a failure occurs, or if the physical machine encounters a hardware problem. It allows system operators to run applications on a different operating system. For example, a computer running on a Microsoft Windows operating system may host a virtual machine that looks like a computer with the Linux operating system.

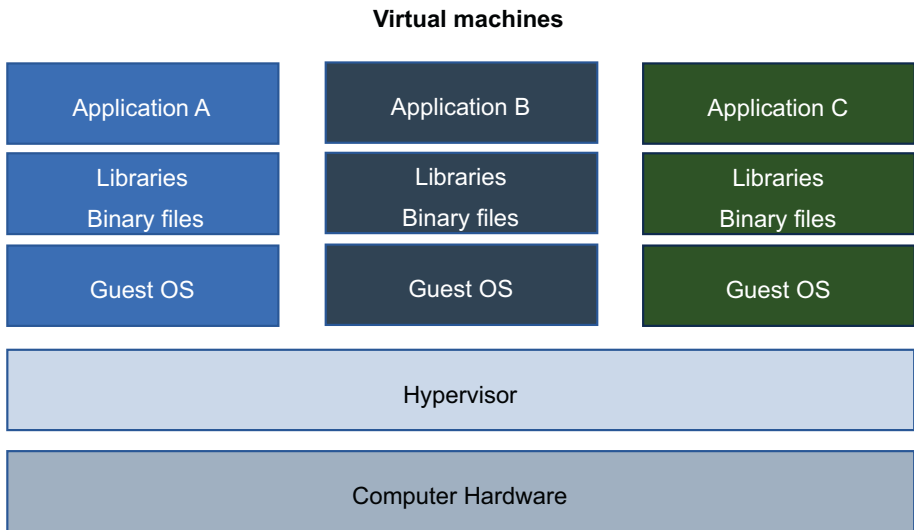


Figure 3-8. *Virtual Machines*

Virtualization is different from hardware emulation, where a piece of hardware imitates another type of hardware. In virtualization, a hypervisor, which is software, imitates parts or even the entire computer of another system. In emulation, the application from another system runs on the emulator, while the hypervisor runs the application from another application on the CPU itself, emulating the differences in software.

Types of Hardware Virtualization

- **Full virtualization**

Full virtualization provides a complete environment for the guest application, which includes hardware, operating systems, libraries, and external software components. This allows the guest application to run unmodified.

- **Paravirtualization**

In a paravirtualization environment, guest applications are executed in their own isolated domains, but the guest hardware environment is not simulated. Hence, guest programs need to be adapted to the host environment.

- **Hardware-assisted virtualization**

Hardware emulation of a guest machine can be ineffective if the architecture of the guest and host machines is very different. Therefore, some CPUs have the feature to support virtualization in hardware and other hardware components to improve the performance of a guest environment.

Autonomic Computing

Hardware virtualization is an important stepping stone toward autonomic computing. Up to the late 1980s, when software standards such as the Unix operating system, Ethernet and TCP/IP, and SQL emerged, applications could only run on the computer system on which they were developed. Each hardware vendor had a fully proprietary application stack from the computer microprocessor to its operating system, to the database and all software libraries.

With the emergence of standards like Unix, the TCP/IP communication protocols, and SQL, the language to program databases, applications became portable from one computer vendor to another.

It was a tedious process, and a programmer usually spent several weeks to perform software migration from one computer vendor to another.

Only with the emergence of the Internet and its standards like HTTP and HTML was it possible to access applications from any device, even though the applications were still bound to a certain hardware and software stack.

The concept of virtual machines eventually broke the dependency of the application on the underlying hardware and operating system. It was an important step in the overall trend in enterprise IT toward autonomic computing, a scenario in which the IT environment will be able to manage itself based on the requirements of all applications running on the system.

Virtualization benefits:

- Scalability of the compute resources
- Improved utilization of the compute resources
- Efficient, centralized administration

Multiple operating systems can run in parallel on a single CPU through virtualization. This parallelism reduces overhead costs and is more efficient than multitasking, where several applications use the same operating system. Another important part of virtualization is the easier way to manage updates. Software upgrades for the operating system and application can be more easily managed without disrupting the user. Overall, virtualization significantly improves the efficiency and availability of resources and applications. In the world of proprietary vertical stacks, a computer is assigned to one or a set of applications and was dimensioned for peak capacity. Therefore, over longer periods, the average utilization in most DCs was around 40%. With virtualization, all computer resources can be dynamically applied to all applications, which drives average utilization to 80% or higher.

Container

The concept of container is a relatively new idea to simplify the virtualization of applications further. It was basically copied from the logistic industry, where the introduction of containers revolutionized the way goods were transported and was a very important enabler in developing today's global market for goods.

Following the idea of a physical shipping container, a standardized large box in which all types of goods can be packed, a container for IT systems wraps an application into an isolated box (see Figure 3-9).

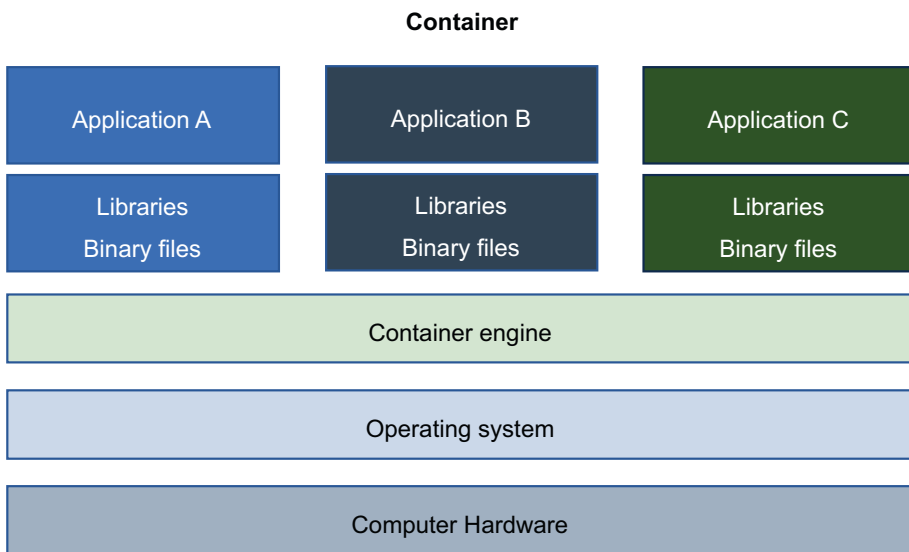


Figure 3-9. *Container*

If an application is “packed” into a container, everything is included that the application requires to run without any dependencies on other applications. When the container is moved to another location, like a different computer system, the application can run at the new location as all dependencies are inside the container.

Fundamentally containers virtualize on a higher level than a hypervisor. A hypervisor virtualizes the underlying computer hardware, while a container virtualizes an operating system. Containers are run on a physical server and its host operating system, mainly a Linux or Windows Server variant. Containers share as many components as possible, including binaries, libraries, and the host kernel operating system. Containers are therefore much smaller than virtual machines, with a size of a few megabytes. Containers have the advantage of a startup time of seconds, while virtual machines often need a few minutes to start.

The light weight of a container allows you to run up to three times as many applications on a single server with containers compared to a server running virtual machines. Containers provide a portable, consistent operating environment for development, testing, and deployment with the benefit of keeping the environment consistent.

Benefits of Containers

Advantages of containers are as follows:

- Containers provide reliable container image build and deployment. Rollbacks can be performed quickly due to image immutability.
- The application-centric management raises the level of abstraction from running an OS on virtual hardware to running an application on an OS.
- Loosely coupled, distributed, elastic, liberated micro-services.
- Applications are broken into smaller, independent pieces and can be deployed and managed dynamically in comparison to hypervisor running a monolithic stack on one big single-purpose machine.

- Containers are much smaller than virtual machines.
- Containers need fewer resources.
- The startup time is in seconds compared to minutes for virtual machines.
- Multiple instances of a single instance, such as a SQL Server, are much faster.
- Feature an isolated environment from the host server.
- Provide a consistent environment for development, testing, and production: it can run in the same way on a PC as in a public cloud.
- Highly portable across operating systems and different CSPs.
- Better observability as OS-level information and metrics are available, as well as information about application health.
- Resource isolation: predictable application performance.

For some applications, however, like a full web server, a virtual machine may be the better choice. A virtual machine provides greater flexibility to choose and upgrade an operating system.

Kubernetes

Kubernetes, also called K8s, is an open source container system. Kubernetes allows you to manage, scale, and automate containerized application deployments. A Kubernetes deployment is called a cluster and consists of at least one main plane, the control plane, and one or more worker machines, called nodes. Both the control planes and node instances can be physical devices, virtual machines, or instances in the cloud.

Control Plane

The control plane, also called a master node or head node, manages the worker nodes and the Pods in the cluster. The control plane runs across multiple computer systems in most DCs. In production environments, a cluster runs on various nodes and hereby provides a high availability environment for the application. The control plane is operated through a command-line interface (CLI), or a program by an application programming interface (API). Master nodes are the control units of a cluster and therefore should not run user applications.

The control plane's components make decisions about the cluster. It also detects and responds to cluster events. Kubernetes relies on several administrative services running on the control plane managing:

- Cluster component communication
- Workload scheduling
- Cluster state persistence

The control plane tracks the state of all cluster components and manages the interaction between them. It is designed to scale horizontally.

Key Value Store

It is a consistent, distributed, and highly available area to store cluster data. The value store is a stateful, persistent storage that stores all Kubernetes cluster data, including its status and all configuration data. It acts as a “source of truth” for the cluster. The data is stored in the control plane or can be configured externally.

Nodes

A virtual or physical machine that has the resources to run containerized applications is called a node. It can be referenced as a worker node or compute node. Kubernetes cluster nodes usually have many nodes. The worker nodes host the Pods that are the components of the application workload. A cluster can be scaled by adding and removing nodes.

Pods

Pods are the smallest deployable computing units created and managed in Kubernetes. A Pod is a group of one or more containers with shared storage and network resources that includes specifications on how to run the containers. The control node schedules and orchestrates Pods to run on nodes.

Figure 3-10 shows the architecture of Kubernetes.

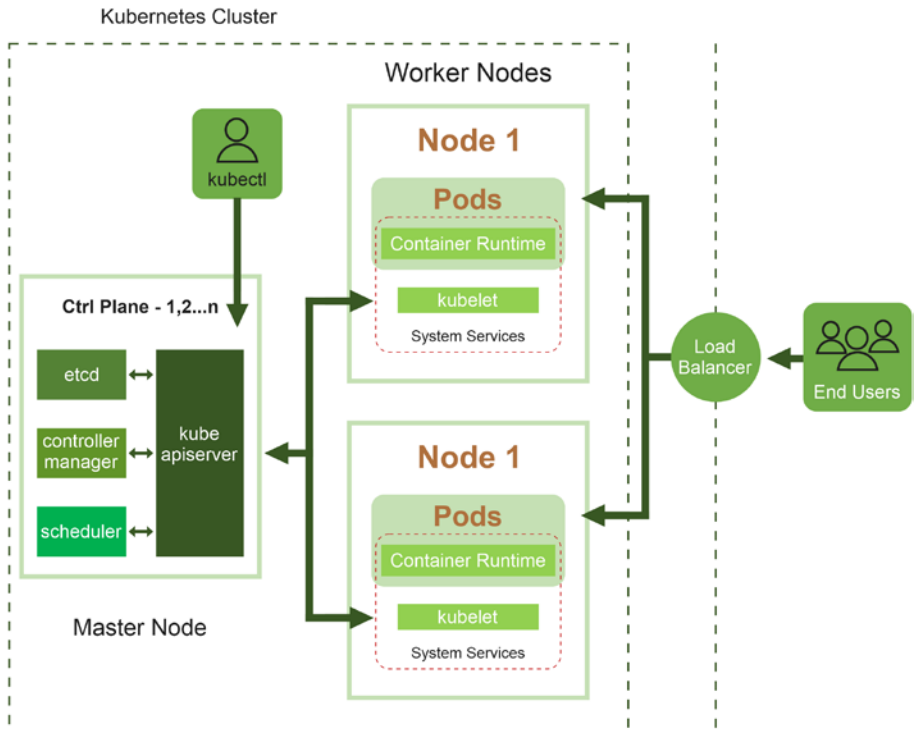


Figure 3-10. *Kubernetes Architecture*

API Server

The Kubernetes control plane can be accessed through the Kubernetes API.

Kubernetes Scheduler

The Kubernetes scheduler performs the following tasks:

- Schedules Pods to worker nodes.
- Watches the API server for newly created Pods with no assigned node and selects a healthy node for them to run on. If there are no suitable nodes, the Pods are put in a pending state until a healthy node appears.

- Watches the API server for new work tasks.
- Makes scheduling decisions based on individual and collective resource requirements and constrains related hardware, software, and policies, specifications for affinity, data locality, inter-workload interferences, and deadlines.

Kubernetes Controller Manager

The controller manager is the controller of controllers. It watches the desired and current state of its managed objects through the API server and can perform corrective actions to bring the current state of an object to the desired state. While each controller is a separate process, for simplification, all processes are compiled into a single binary running in a single process.

Types of Kubernetes Controllers

- **Node controller**
Responsible for noticing and responding when nodes go down
- **Job controller**
Watches for job objects that represent one-off tasks and then creates Pods to run those tasks to completion
- **Endpoints' controller**
Populates the Endpoints object (it therefore joins Services and Pods)
- **Service account and token controllers**
Create default accounts and API access tokens for new namespaces

Cloud Controller Manager

If a cluster is running in a cloud environment, the cloud controller manager integrates the cluster into a cloud by adapting to the underlying cloud technologies of this cloud. The cloud controller manager therefore only runs controllers that are specific to the cloud service and links the cluster to the cloud service provider's API. It separates components that interact with the cloud platform from those components interacting only with the cluster.

Controllers with Cloud Provider Dependencies

- **Node controller**

This controller is used for checking the cloud provider to determine if a node has been deleted in the cloud after it stops responding.

- **Route controller**

The route controller sets up routes in the underlying cloud infrastructure.

- **Service controller**

The service controller is used for creating, updating, and deleting cloud provider load balancers.

Kubernetes Node Components

Node components run on every node, maintaining running Pods and providing the Kubernetes runtime environment.

- **Kubelet**

A Kubelet is an agent that runs on each node in the cluster and acts as a conduit between the API server and the node. It ensures that containers are running healthy in a Pod. It instantiates and executes Pods, watches the API server for work tasks, and gets instructions from the master.

- **Kube-proxy**

A kube-proxy is a networking component that manages IP translation and routing. It is a network proxy that runs on each node in a cluster and maintains the network rules on nodes. These network rules allow network communication from inside or outside of a cluster to Pods and manage to assign a unique IP address to a Pod, ensuring that all containers in a Pod share a single IP. It manages Kubernetes networking services and load-balancing across all Pods in a service, deals with individual host sub-netting, and ensures that the services are available to external parties.

- **Container runtime**

The container runtime refers to the software responsible for running containers in Pods. To run the containers, each worker node has a container runtime engine that pulls images from a container image registry and starts and stops containers. Kubernetes supports several container runtimes: Docker, containerd, CRI-O, and other implementations of the Kubernetes Container Runtime Interface.

- **Container registry**

The container images that Kubernetes relies on are stored in a container registry. This can be a registry you configure, or a third-party registry like Docker Hub, Amazon Elastic Container Registry (ECR), Azure Container Registry (ACR), and Google Container Registry (GCR).

- **Persistent storage**

Managing the containers that run an application, Kubernetes can also manage application data attached to a cluster. Kubernetes allows users to request storage resources without having to know the details of the underlying storage infrastructure.

- **Underlying infrastructure**

One of the key advantages of Kubernetes is that it works on many kinds of compute infrastructures. This can be bare metal servers, virtual machines, public cloud providers, private clouds, and hybrid cloud environments.

3.9 Edge Computing

Over the last 60 years, there has been a constant change in the approach of centralized vs. decentralized computing. It started with the centralized approach of mainframes in the 1960s. All compute resources were part of the main computer, and IT users had to use simple terminals that had no local computing power. In the 1980s, the paradigm shifted with the invention of technical workstations. Workstations brought engineers

a computer system with good graphics to their desks to run their applications for mechanical and electrical design. Later, workstations were used for software development too. PCs entered in the 1990s, and they were largely deployed due to their low price. PCs gave all users full computing power on their desk, shifting the balance toward decentralized computing.

In the 2000s, cloud computing started to gradually reverse this trend again, and very large DCs with 100,000s as before, should it be changed to “hundreds of thousands”? of computers were built. This shifted a large share of the compute resources to the central cloud. While cloud computing provided a very cost-effective way to run workload, it was recognized that cloud computing had difficulties to deal with real-time applications as the data transport to a central system causes delays and is therefore often too slow to react timely to critical events. Video streaming services that require fast and uninterrupted transmission of large amounts of data are another application where the centralized cloud approach does not work well.

This brought the concept of edge computing into the picture. As the name suggests, edge computing moves some compute resources from the center to the edge. This shifted back the emphasis to a more distributed computing infrastructure where computer systems and storage are moved closer to the end users and sources of data. Early in the 2000s, content distribution networks (CDN) were built to improve the download speed of popular websites and video content.

With the emergence of IoT, billions of devices will be connected to the Internet or private networks over time. This will create huge amounts of data from devices such as video cameras, machinery sensors, and autonomous vehicles, which could be considered as an edge computing system by itself. The trend to build smarter infrastructures in the future will require large deployments of sensors, actuators, video cameras, and edge computing.

A good example of an IoT edge system is an intelligent lamppost, which was deployed in some cities such as Milan, Hong Kong, and Singapore. An intelligent lamppost is a multifunctional device that provides the following:

- Intelligent lighting for the surroundings
- Video camera for traffic control and general security
- Data collection of environmental parameters such as noise, pollution, and weather sensing
- Loudspeaker for traffic warning and public announcements
- Wi-Fi connection
- Charging for electric vehicles

All these functions require that a computing device is part of the intelligent lamppost and hereby makes it an edge computing device of the smart infrastructure (see Figure 3-11).

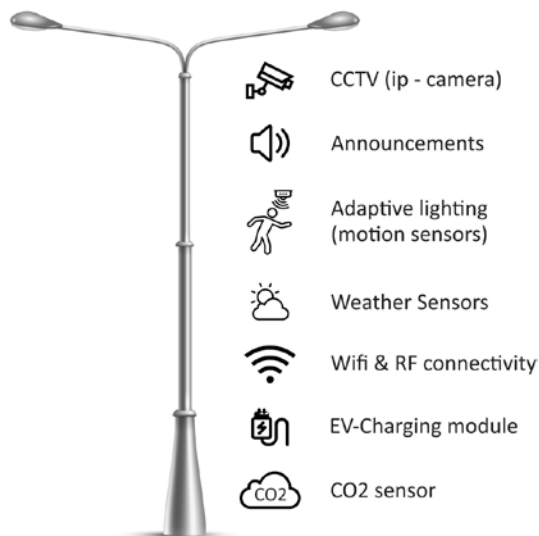


Figure 3-11. Smart Lamppost (*macrovector; yoyonpujiono; veryicon.com, 2023*)

A newer trend in edge computing is edge nodes used for game streaming services, also referred to as gamelets. Gamelets are required to meet the features of interactive gaming with its demand of high data transfer in real time.

Edge computing can drive the efficiency of an IT infrastructure as analytical resources can be deployed at the edge. Data collection and data analytics by AI tools can run on the edge of the system and will therefore increase the operational efficiency of the overall infrastructure.

The following are a few good examples of edge computing devices:

- **Intelligent lamppost**

Delivers intelligent lighting, traffic control, and video surveillance.

- **Content distribution network**

A popular video is downloaded to an edge device from where hundreds of users can stream it. Without the edge device, all the data must be delivered from a central system located in another country or even another continent, creating a hundred times higher load to the central system and to its data network.

- **IVR systems**

Voice recognition systems have a much better quality when long data transfers with delays can be avoided. IVR running in edge devices therefore provides better voice quality services.

- **Autonomous cars**

Autonomous cars require very fast reaction times to external events. Only computer systems in a car are fast enough to perform such reactions and communicate with other centralized systems such as traffic control systems and traffic routing services.

3.10 Compute Resiliency

Definition

A system is classified as resilient if it remains intact and dependable when facing changes in the form of attacks, threats, disasters, failures, and updates.

Fault-Tolerant Computing

The first fault-tolerant computer system was built in the 1960s by NASA and other organizations for their space program. Fault-tolerant systems became generally available when Tandem created the first generation of commercial fault-tolerant computer systems in the 1980s. The design of such systems ensures the system continues its operation in case of single or multiple failures of its components. Providing an unplanned downtime of close to 0% enables a fault-tolerant computer to power applications where any downtime could be life-threatening, dangerous, or could cause large financial damages. Examples of such applications are computer systems used in hospitals, air control systems, nuclear power plant control systems, systems for chemical processes, and exchanges for stocks or currencies.

When a fault-tolerant computing system faces a problem, the operating quality will decrease only proportional to the severity of the failure, while standard computer systems may stop working entirely even if a small failure occurs. The gradual reduction of functionality during a fault situation is called “graceful degradation.”

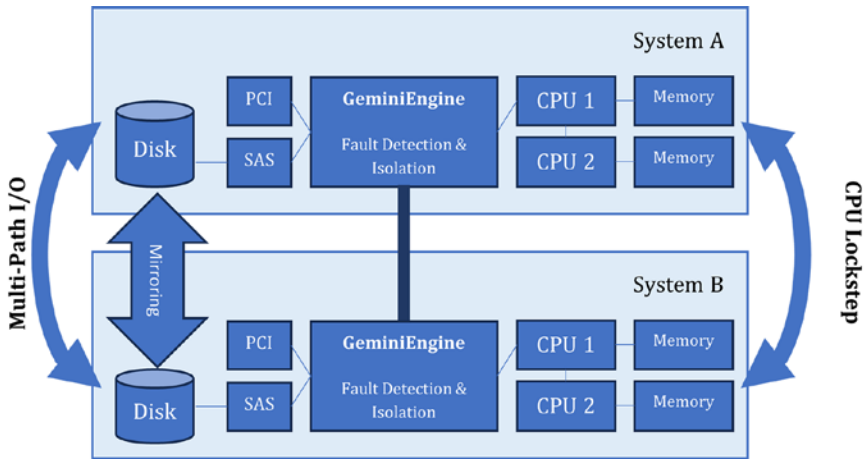


Figure 3-12. Architecture of a Fault-Tolerant Computer from NEC

When parts of the system fall, a fault-tolerant system continues its intended operation. It may reduce its performance but will not fail completely. Figure 3-12 shows a typical system architecture of a fault-tolerant system and reveals that a fault-tolerant computer is created from two independent computer systems that share the same system design with separate external connection and power supplies.

The two systems can operate in different modes, namely:

- **Cold standby**

In this configuration, system A runs the application, while system B is in standby mode. System B monitors system A and takes over the applications if system A does not function properly.

- **Hot standby**

In this configuration, both systems run applications. System A may run the critical applications, and system B may run noncritical applications. If system A fails, system B will suspend its applications and take over the critical applications from system A.

- **Parallel computing**

In this configuration, both systems run the same application with the same data so that, should there be a failure, at least one system continues to function. For super-critical applications, more than two systems operate in parallel, which is also a method to verify the results of the computation process.

Fault-tolerant computing is a proven solution and assumes that applications run on one or a few computers. The system is still vulnerable to events that disrupt the function of a DC such as longer power disruption, sabotage acts, or natural events such as earthquakes and flooding.

As modern DCs may have installed thousands or even ten thousands of computer systems, other strategies can be applied to ensure that applications run with a 100% uptime.

Resilient Computing

Resilient computing goes beyond the concept of fault-tolerant computing (see Figure 3-13). While the objective of fault-tolerant design is defined by providing a very high uptime level for a given set of applications, resilient computing provides a system that is designed for a very high level of reliability and tolerance for changes that can be introduced by application updates. The deployment of large amounts of servers in a DC makes it

much easier to deal with hardware and software failures. The control plane of a hypervisor distributes an application to a virtual machine and checks the execution status.

If the application fails due to a software or hardware problem, the control plane reassigns the application to another virtual machine. In modern DCs that have 10,000 or more computers, failure of a couple of systems does not matter much as there are still plenty of resources available. Therefore, such DCs have very low-level support agreements with their hardware vendor, or may not even have one at all. System repairs can be performed weekly or biweekly by sweeping through each DC room and simply replacing all the failed hardware components. The control plane is the critical component that needs to be replicated with a fault-tolerant design in such a highly distributed architecture.

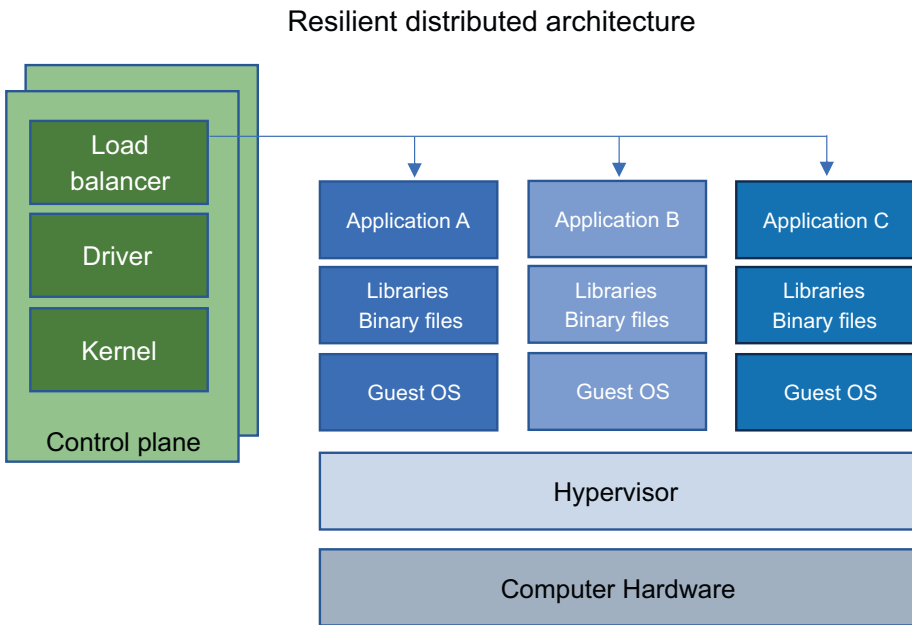


Figure 3-13. Resilient Distributed Architecture

For DCs using containers, the same approach can be applied. The control plane distributes applications in containers and monitors them. In case of failure, the container is moved to another computer system.

Federated Architecture

As described in the previous paragraph, a hypervisor or container-based control plane can easily re-distribute applications to another computer system, making the system fault tolerant to hardware and software failures on single machines. But the DC itself could become unavailable for several reasons, such as power outages, natural disasters, or criminal acts. A higher level of resilience can be achieved if multiple DCs can be connected through a federated architecture at different locations, preferably different cities, or countries. This elevates the concept of distributed hypervisors to a higher level. In federated architecture, applications can be distributed to different systems within a DC and to computer systems in another DC, which provides a much higher level of security than a single DC (See Figure 3-14).

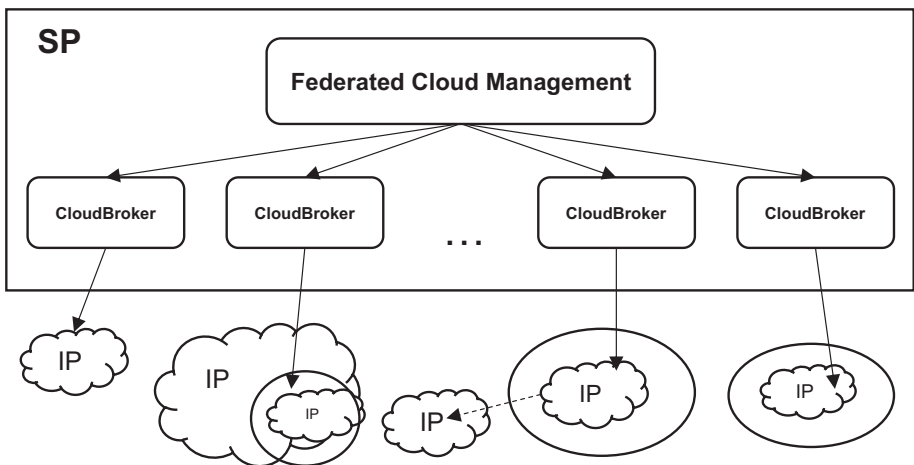


Figure 3-14. Federated Computer Architecture

3.11 Provisioning and Administration of Compute Resources

Computer Workload

A computer workload is simply any program or application running on a computer system. It has the same meaning as application, software, and program.

In today's context of large private DCs and compute clouds, the term "workload" is used to define the amount of work that an application imposes on the underlying computing resources. A light workload refers to applications that only require small amounts of the compute infrastructure, while a heavy workload will require large amounts of computing resources to perform the application.

Types of Computer Workloads

Workloads come in many different shapes and forms. The following workload terms are common:

- **A static workload**

A static workload is defined by applications that are always running, such as the operating system of the computer system itself, an email system, most websites, enterprise resource planning (ERP), customer relationship management (CRM), and most other applications supporting the business of a company.

- **A dynamic workload**

Dynamic workloads are applications that run when there is a demand for it. Almost all applications related to the development and testing of products and

services are in this category. Dynamic workload can be recurrent, like running a certain application such as programs creating monthly bills or payrolls.

- **Batch and transactional workload**

Batch processing is an old computing term emanated in the era of mainframes where compute resources were expensive and where the compute capacity was divided between batch workloads and transactional workloads. Transactional processes required user interactions such as programming or order data entry, while other applications run as background batch processes for tasks such as daily accounting or a monthly payroll.

- **Database workload**

This refers to enterprise applications such as ERP, CRM, and Order Management Systems using data that is stored in databases. Database workloads require specific compute and storage resources and need to be tuned and optimized to maximize the search performance for all applications, depending on the database.

- **HPC workload**

High-performance computing (HPC) workloads are mainly used in the field of computational science. These applications perform computationally intensive tasks in various fields, such as weather forecasting, climate research, and physical simulations. In certain application areas such as daily weather forecast, the application must be able to run the computation

within a certain time frame. To be meaningful, a daily weather forecast should not take more than a few hours compute time.

- **Analytical workload**

Analytical workloads are related to applications processing large amounts of data. These types of applications have evolved from simple data analytics processing of a relatively small amount of data to sizable data warehouse-based data analytics starting in the 1990s and, today, machine learning applications where a very large amount of data is processed.

- **Real-time workload**

Real-time workloads are different from other workloads. Computer systems controlling critical processes or powering simulators require the feature to react in “real time,” which means that the reaction time to events must be fast. The systems must therefore have a low latency between an event at the controlled object and the time the information reaches the computer system. To stay within the requirements of the processes, the time to compute the event and process the algorithm must be short. Such “real-time computer systems” are therefore dedicated to the real-time application and will not perform any other workloads.

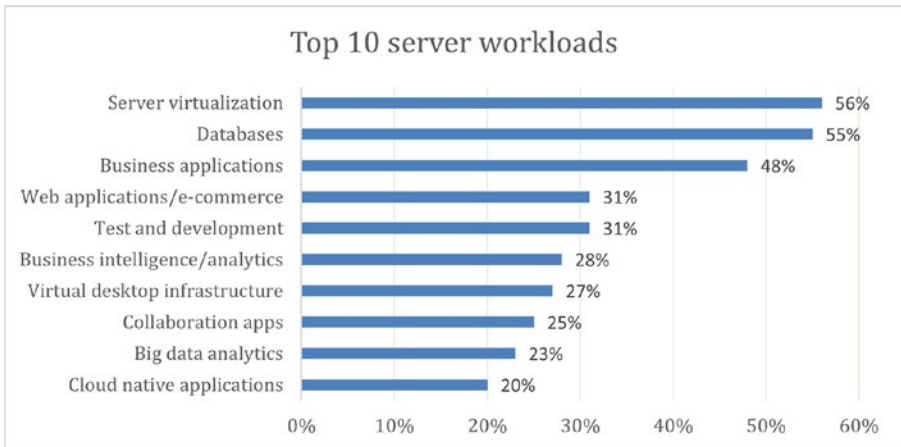


Figure 3-15. Top Ten Server Workloads (Source: techtarget, 2023)

The graphic in Figure 3-15 shows the average workloads distribution in DCs.

Workload Deployment

In the past, all workloads were running in the DC of the organization. With the emergence of cloud DCs, organizations have the choice to shift some or all of their applications to a cloud service provider. The DC of an organization is also referred to as the private DC, private cloud, or on-premises, while DCs of CSPs are called cloud DC, public cloud, or off-premises. With the option of CSPs for their workload, organizations need to define rules for how they want to run their workloads.

In Chapter 1, Section 1.3, we had described the various deployment scenarios that are available for organizations to choose from.

Benefits and Challenges of Private and Public Clouds

CSP

Benefits

- **Scalability**

CSPs have built large computer resources. Using a CSP resource offers organizations a very elastic way to use computer systems. Compute resources can be increased and scaled down on very short notice.

- **Performance**

CSPs have many resources often distributed over several DCs. If organizations need higher performance computing power, a CSP can easily provide such resources. This, however, may not be economical for an organization running a private cloud.

- **Cost management**

CSPs may charge based on a monthly fixed fee or on a consumption-based model. This means all expenses related to an organization's computing are variable costs without any capital expenses otherwise involved in building and maintaining a private cloud.

- **Data residency**

Businesses are subjected to varied data protection and data residency regulations in the country where they operate. For example, if an organization has a private cloud only at their HQ and operates in a different country, it may be challenging to comply with the local

data protection regulations of the country. Keeping the data with a local CSP is an easy way to overcome this problem.

Risks

- **Operational transparency**

For workloads running with a CSP, there is a general lack of transparency about the underlying, often multi-tenant infrastructure of the CSP, let alone any control. This makes it very difficult for organizations to validate or audit regulatory demands for workloads running at a CSP.

- **Service disruption**

If there are outages caused by service disruptions of the CSP, there are few options to act for the organization using the CSP. The management of the incident is completely in the hands of the CSP, and depending on the contracted service-level agreement (SLA) with the CSP, there might not be enough escalation power to ensure fast remediation.

- **Changes with the CSP**

The CSP chosen by an organization will evolve over time. It may merge or be acquired and may change its business direction. This can lead to changes and disruptions for the organization's workload running at the CSP. Hence, the organization needs a workload failover plan to handle cloud disruptions caused by the CSP.

Private Clouds

Benefits

- **Control over cloud resources**

The DC of an organization is in full control, and therefore, it has access to all relevant data for the operation, including log files. Activities such as troubleshooting, correcting, and audit can be performed within its own premises. As the organization has full control over the DC, proactive steps can be taken to protect the DC environment, and SLAs for the DC operations can be established.

- **Visibility and compliance**

With a private cloud, an organization has complete control and visibility of the data center infrastructure including all servers, storage, network, and other components such as the operating systems and other elements of the software stack. Private clouds are the preferred deployment method for business-critical applications, demanding workloads or workloads with stringent regulatory compliance or high-security requirements.

Challenges

- **Costs**

Building a traditional data center can be a significant undertaking that requires substantial investment into a DC infrastructure. There are high capital expenses for construction and outfitting a DC, as well as ongoing operational expenses to keep the system running. This includes maintenance and regular updates of all DC components - regardless of the workload utilization.

- **Protection**

Operating a private cloud includes the complete responsibility to keep the DC secure and resilient. For smaller companies, it may be difficult to afford the hardware and software to sufficiently secure the private cloud and develop the expertise to run a private cloud with their own staff.

Kubernetes Workload Management

All workloads are running inside a set of Pods. Pods represent a set of running contents on a compute node and have defined life cycles (see Figure 3-16).

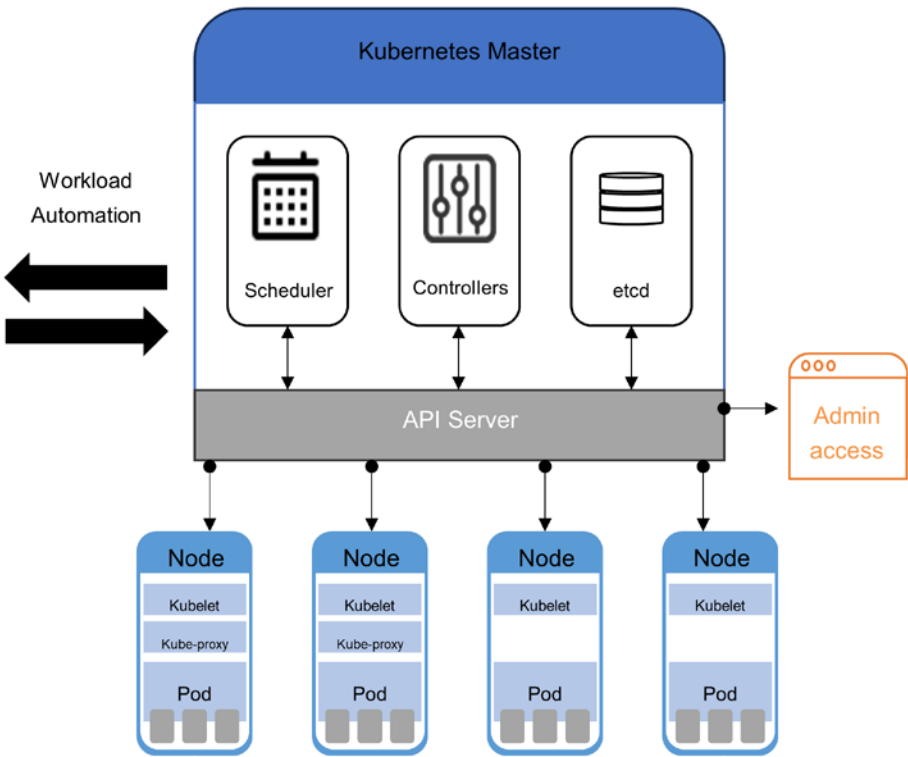


Figure 3-16. Automated Workload Distribution with Kubernetes (Icons from veryicon.com, 2023)

Pods are constantly surveilled. If the system detects a critical fault on a node in one of the containers, Kubernetes considers that all Pods running on the node have failed. New Pods need to be created to restart the workload. To avoid manual intervention, the Kubernetes Workload Resource Manager can be used to automatically manage the workloads and Pods.

Kubernetes provides several built-in workload resources:

- **“Deployment” and “ReplicaSet”**

“Deployment” is a workload resource to manage a stateless application workload on a cluster, where any Pod in the “Deployment” is interchangeable and, if needed, replaceable. What about “ReplicaSet”?

- **“StatefulSet”**

“StatefulSet” is used as a workload resource to run one or more related Pods tracking a status. If a workload records data persistently, the “StatefulSet” function matches each Pod with a “Persistent Volume.” This allows a Pod to replicate data to other Pods in the same “StatefulSet” to improve overall resilience.

- **“DaemonSet”**

“DaemonSet” creates Pods with node-local facilities. Every time a node is added to the cluster that matches the specification in a “DaemonSet,” the control plane schedules a Pod for that “DaemonSet” onto the new node.

- **“Job” and “CronJob”**

A “Job” represents workload as one-off tasks that run to its completion and then stop, whereas “CronJob” provides the same functionality to recurring workloads with a planned schedule.

Other third-party workload resources are available from software companies with different behaviors in the wider Kubernetes ecosystem. Using a custom resource definition allows you to use workload resources with specific behavior that’s not part of Kubernetes’ core.

3.12 Charging for Compute Resources

If an organization moves workloads to a CSP, they will be charged based on usage base pricing. While this seems to be a straightforward approach, the charging for compute services is not as transparent as most organizations would expect.

Terminology

The lack of transparency starts with the terminology that is used: A virtual machine (VM) is referred to by Amazon Web Services (AWS) as an “instance”; in Google Cloud, a VM is called a “VM instance,” while Microsoft Azure calls it a “VM.” In AWS and Google, a group of VMs is called an “auto-scaling group,” while Azure refers to it as “scale-sets.” AWS offers “on-demand instances,” which are unscheduled execution of a workload. Azure calls this “pay as you go,” and Google refers to it as “sustained use.” For noncritical workloads, a discounted pricing scheme is offered. It is called “spot instances” in AWS, “preemptible instances” in Google, and “low-priority VMs” in Azure.

The different terminology is associated with pricing models, which usually makes comparison difficult. The pricing often includes elements of charges for the reservation of resources. This is called “reserved instances” in AWS, “reserved VM instances” in Azure, and “committed use” in Google.

See examples of Microsoft Azure’s cloud services charges in Table 3-3.

Table 3-3. *Example of Microsoft Azure Cloud Services Charging Structure (Source: cloudmore, 2023)*

Microsoft product	Price list name	Price list currencies
Office365	License-based services	18 currencies
Azure pay-as-you-go	Azure Plan consumption pricing	USD
Reserved Instances	Microsoft Azure reserved instances	USD
Microsoft Software Subscriptions	Software subscriptions	18 currencies
Marketplace services	Software subscriptions	Global

Pricing Variables

Most CSPs use parameters such as operating systems, compute type, network speed, memory size, and disk space, and they are all different factors to determine the sizing and pricing of an instance to run a workload. Most CSPs use the following categories for VMs:

- General purpose
- Compute optimized
- Memory optimized
- Disk optimized

On the next level, families are defined; AWS uses the following types: t2, t4, m2, and m4.

Then the processor family must be defined in small to very large categories.

All together it is a complicated charging structure that needs to be applied to each workload and application set.

Dashboards

Most CSPs provide consoles and dashboards to monitor the resources in use. For most organizations, it is difficult to interpret the data and relate it to the workload they are running with the CSP, and therefore, they are not able to assess the costs that are incurred from their current workloads. Many organizations have built their own dashboards by using application programming interfaces (API) from the CSP to get a better understanding of how much the CSP charges for their workloads.

Alternatively, an organization can use commercially available software to manage its interaction with CSP. Tools like ParkMyCloud provide a platform to monitor usage patterns and costs of CSP services and are often a worthy investment to save costs from reallocation workloads to CSPs truly.

Charging Structure

CSP pricing may be based on different time scales, such as charging per second, per minute, or per hour. The clock for charging runs when the workload is deployed. Organizations often forget to turn static workloads off. This is an issue for organizations using CSP resources in different regions, or on different accounts.

It is a common mistake that organizations overprovision or oversize resources with a CSP; this typically happens when resources are deployed for maximum capacity, like an order system for a company with peak sales before Christmas.

Another issue that makes budget planning for CSP resources difficult is the fact that CSPs offer to change their pricing. AWS's pricing model has already changed 60 times since it was introduced.

While many CSPs promote their services as an easy-to-use resource pool like water or electricity from a Utility, all the described items suggest that using a CSP is not so simple after all. This means organizations wanting to move parts of their workloads to CSPs need sufficient expertise in cloud computing to select the right mode of operations for them and monitor the associated costs.

3.13 Summary

Computers are programmable systems that have evolved over 80 years from basic devices to essential components of various modern gadgets and industrial equipment. Their physical elements include the CPU, memory, storage, power, and connectivity modules. Software, the set of digital instructions run by computers, is broad-ranging, from basic apps to intricate programs. The operating system acts as a bridge between the user and hardware. Based on usage, computer systems can be personal (like PCs and smartphones), embedded in devices (like cars and robots) and servers or be part of large data centers. While many computers today are general purpose, following the predictions of “Moore’s law,” specialized systems were more common in the past. Data centers, facilities housing IT infrastructures, started in the 1960s and have transformed with the rise of the Internet into vast public clouds offered by companies like Amazon and Microsoft. As the digital era advances, the interplay between software and hardware intensifies, shaping our interaction with the technological world.

In the next chapter...

- Storage resources, pools, and pricing
 - How to save money – useful strategies
-

CHAPTER 4

Storage and Virtualization

All computer systems need a storage facility where operating systems, applications, and data can be stored and retrieved. In a broader definition, computer storage includes the following:

- Main memory of a computer system.
- Internal and external storage in the form of disks and tapes.
- External storage facilities and solutions in the form of remote storage.
- Storage and backup management software utilities.
- Networking technologies for storage, such as NAS, DAS and SAN.
- It also includes data center storage policies and procedures that govern the entire data storage and retrieval process. Moreover, data center storage may also incorporate data center storage security and access control procedures and methodologies.

The following paragraphs will focus on the hardware and software for the storage systems, while the topic of backup will be covered in a separate chapter (Chapter 6).

The features of the storage systems are defined by the following parameters:

- **Latency**

The time it takes a storage medium to find a particular location for the storage operation is called latency.

With current storage technology, the latency of primary storage is measured in nanoseconds, semiconductor-based storage like memory sticks in the microsystem, disk drive latency in milliseconds, and tapes in seconds or minutes. Latency for the operation of reading and writing may vary for each medium.

- **Throughput**

Throughput measures the data rate at which information can be stored or retrieved from a storage medium and is usually measured in megabytes per second. For spinning devices such as tapes, magnetic, and optical drives, there is a big difference between the throughput for sequential access in comparison to random access.

- **Granularity**

Data transfers to a storage medium are transmitted in blocks. The size of the blocks, usually measured in kilobytes, determines the speed of the data access. Granularity needs to be optimized for certain applications to ensure a high throughput.

- **Reliability**

Data storage reliability is determined by the accuracy of the data that is stored on a storage medium.

4.1 Storage Resources

Computer systems use a hierarchical structure for data to distinguish between categories: primary, secondary, and tertiary storage. The storage capacity of registers and cache is small. In today's computers, a register is measured in bytes and cache in kilobytes, while main memory reaches gigabytes.

Primary Storage

The main memory, also called the internal memory or just memory, is the direct accessible storage resources of the central processing unit (CPU). In the early days of computers, magnetic cores were used, and memory was referred to as "core memory." Since computer memory was made from integrated circuits (IC) in the 1980s, the term "random-access memory" (RAM) is used. The main memory contains all the instructions and data for a CPU to run an application. To improve computer performance, the primary storage includes a small faster layer of memory, called cache.

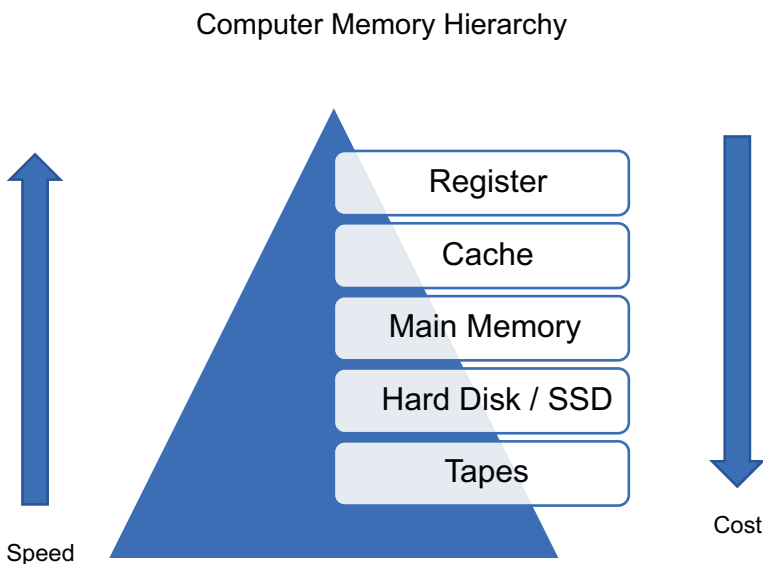


Figure 4-1. *Computer Memory Hierarchy*

To start from the top of the storage hierarchy (see Figure 4-1), the CPU has registers with a size of 32 or 64 bits, where the data for computation is loaded and stored. Actual computer operations like adding two variables can only be performed with registers. On the next level, there is the processor's cache memory, used as an intermediate medium between the very fast registers and the much slower main memory. As cache memory is much more expensive than main memory, only a small part of the memory capacity is cache memory. There are multiple cache layers in modern computers. Table 4-1 shows a typical hierarchy for primary storage.

Table 4-1. *Typical Hierarchy for Primary Storage*

Type	Typical size per CPU core	Speed
L1-Cache	32–64 KB	Super-fast
L2-Cache	256–512 KB	Very fast
L3-Cache	8–32 MB	Fast
RAM	1–8 GB	Slow

Only the most actively used information of applications, instructions, and data is copied from the main memory into the cache memory. Cache memory constantly gets overwritten when a particular section of program and data is no longer in use, and this is called the “least recently used” (LRU) cache management algorithm.

The main memory is connected to the CPU by a fast memory bus. The data that is stored in the RAM is volatile and is lost when a computer system is switched off, or when the operating system is shut down. However, computer systems use integrated circuits that contain permanent information like the small memory unit, which is used to start a computer system, called the BIOS. For many appliances, the operating system and application code are static and are therefore stored in a permanent read-only memory (ROM). While the term “ROM” is still used

for memory areas with static content, most of them can be re-written; hence, such ROM is not “read-only.” When a firmware upgrade is performed, it usually means to overwrite the content that is stored in a ROM.

Secondary Storage

Secondary storage is only connected via the input/output channel (IO channel) of a computer system and is therefore not directly accessible by the CPU. The operating system, applications, and data are stored in secondary storage. All data stored in secondary storage is static and remains intact after the computer system or storage system is switched off.

Modern computer systems use secondary storage hard disk drives (HDDs), solid-state drives (SSDs), or both. Compared to RAM, HDDs and SSDs have a much bigger storage capacity, measured in GB for RAM and TB for disk. HDDs are much cheaper than SSDs, but the access time is much longer, measured in milliseconds for HDDs and nanoseconds for RAM in SSDs. Computer systems using laser-based CDs or DVDs media for secondary storage have an even slower access time.



Figure 4-2. HDD (Unsplash, 2018)

HDD

HDDs are rotating magnetic disks as shown in Figure 4-2. Rotating magnetic disks were already developed in the 1980s, replacing punch card and paper tapes. Early disk systems had a diameter of 40 cm and a capacity of 5 MB. Modern HDDs have a diameter of 10 cm and a capacity of up to 10 TB. This is an increase in storage capacity of 140 million times per cm^2 , which shows that the evolution in disk technology happened at a comparable speed as computer technology following Moore's law over the last 40 years.

Over time, many methods were developed to increase the speed of HDDs. Positioning the read-write head of an HDD is a relatively slow process compared to the actual read and write function. Seek time and rotational latency of the disk can be reduced when data is transferred to and from disks in large contiguous blocks. This makes random access of data on HDDs much slower than sequential read or write access. HDDs

are usually formatted to comply with the file system that is being used by the operating system. The file system stores the actual data of a file as well as information about the file, which is called metadata, and includes information such as file name, file size, file owner, and access rights.



Figure 4-3. SSD

SSD

SSDs are integrated circuits that can store data permanently (see Figure 4-3). Unlike HDDs, SSDs have no moving parts, and therefore, they are not really drives. SSDs are much faster than HDDs; they are more durable and lighter and can be used in moving objects like a mobile phone, but they are more expensive. SSDs are the latest evolution in storage technology.

SSDs provide a data access speed of ten times faster compared to HDDs, but the cost per MB is almost ten times higher. Therefore, HDDs are better for storing larger amounts of data and archiving. With continuously falling prices of semiconductors, SSDs will soon overtake HDDs as the primary storage devices in computers.

Tertiary Storage

Data stored on media, which is not immediately available, is called tertiary storage. It used to be magnetic tapes that were shown as the characteristic image of a computer system in video footage of a computer system in the 1960s to the late 1980s. Modern tape libraries, operated by a robot, are still in operation for very large data storage, as the cost per MB is still the cheapest. But with the emergence of CSPs with very large storage capacities and constantly falling prices for disks, tertiary storage plays a diminishing role in storage systems.

Offline Storage

Offline storage refers to media detached from computer systems. Media such as tapes or tape libraries can be physically transported to another location. Tape libraries, tapes, or DVDs allow you to store data at a highly secured location outside of a DC. This is required for several applications:

- **Disaster recovery**

When a DC is damaged by disasters such as fire, flooding, an earthquake, or criminal acts, the DC environment of an organization can be restored from the stored media that is kept at a safe location.

- **Confidential data**

To prevent staff or hackers from accessing data, very confidential data may not reside on a computer system. Offline media stored at a secured location are much harder to access for unauthorized users than data that is stored on a computer system.

- **Data integrity**

To validate critical information, a copy can be recorded (just a suggestion) on a tape and stored at a secured location. If doubt arises about the integrity of data in files, the offline storage can be used to validate the data.

Table 4-2 compares different storage systems.

Table 4-2. *Comparison of Storage Systems*

Media	HDD	SDD and flash	Optical disk	Tapes
Technology	Magnetic disk	Semiconductor	Laser beam	Magnetic tape
Random access	Yes	Yes	Yes, but limited	No, only sequential
Controller	Integrated	Integrated	External	External
Latency	15 ms	Microseconds	150 ms	Very slow
Failure cause	Head crash	Circuit failure	Scratches, dust, age	Magnets or cuts
Price per MB	Moderate	Decreasing but typically higher than HDD	Moderate	Lowest
Application	Main storage in desktops and servers	SDD – main storage in a computer Flash – portable storage and devices	Long-term archive, distribution	Long-term large-scale backup

4.2 External Storage Systems

Computer systems have a limited capacity to house storage devices in the form of HDDs or SSDs in their chassis. If more secondary storage capacity is needed, external disk systems are required. External storage ranges from simply connecting a PC to an external USB disk to high-end storage system providing petabytes of data that is implemented with sophisticated storage technologies.

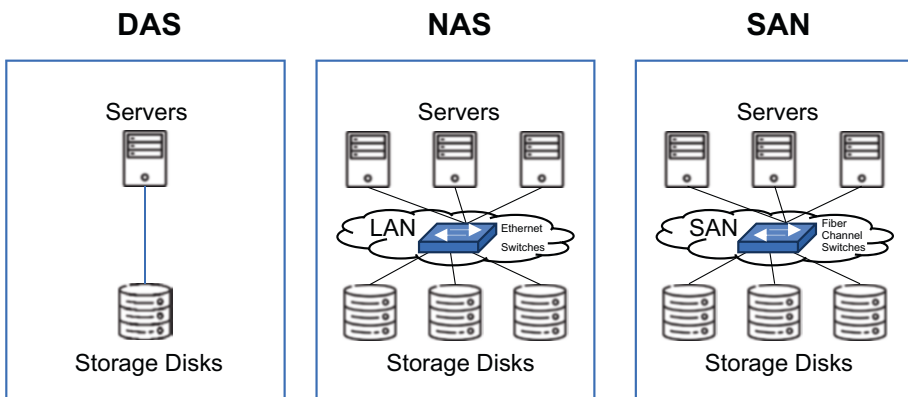


Figure 4-4. Comparison of Storage Types (Icons from Flaticon, 2023)

Types of external computer storage are

- Direct attached storage (DAS)
- Network-attached storage (NAS)
- Storage area networks (SAN)

These three architectures distinguish how an external storage device is connected to a computer system (see Figure 4-4). Direct attached storage directly connects to a computer system via its I/O system. Today, network-attached systems connecting storage network systems through

the networking infrastructure are mainly based on the TCP/IP. Storage-attached networks use a storage network system with its protocols to connect external storage to the computer systems.

Direct Attached Storage (DAS)

Direct attached storage (DAS) refers to digital storage directly attached to a computer system (see Figure 4-5). A DAS consists of storage units such as hard drives, solid-state drives, or optical disk drives that are housed in an enclosure outside of the computer chassis.

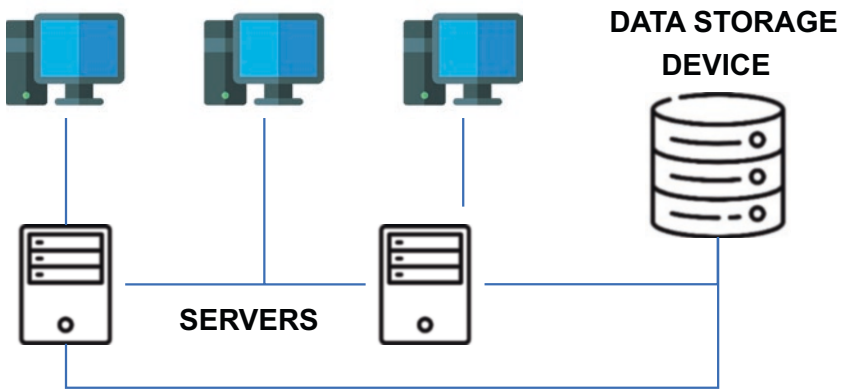


Figure 4-5. *Direct Attached Storage (Icons from Flaticon, 2023)*

A DAS system connects directly to the I/O system of a computer using a network. It must therefore be in proximity to the computer system. Most computer systems are connected to a DAS storage system by protocols such as FC, SCSI, SATA, or USB. Higher-level DAS systems can be connected to more than one computer system.

Network-Attached Storage (NAS)

Network-attached storage (NAS) consists of network storage facilities that contain one or more storage drives (see Figure 4-6). Data access is provided on a file level, unlike DAS systems where access is performed on block level. Most NAS systems are manufactured as an autonomous appliance, which includes a computer system designed for storage, storage devices, and a network controller. NAS systems can be accessed by one or multiple computers, usually through LAN-based TCP/IP connections and file-sharing protocols such as SMB, AFP, or NFS. As NAS systems are based on purpose-built computer systems for storage, they usually provide advantages compared to general-purpose computers dedicated for storage purposes. In comparison, NAS systems provide faster access to data and easier administration and are simpler to configure. The disk drives used in NAS systems are usually the same as those used in general-purpose computers, but their firmware is often modified to provide better features for a NAS, including a better vibration tolerance and power distribution to support RAID configurations.

Network Attached Storage (NAS)

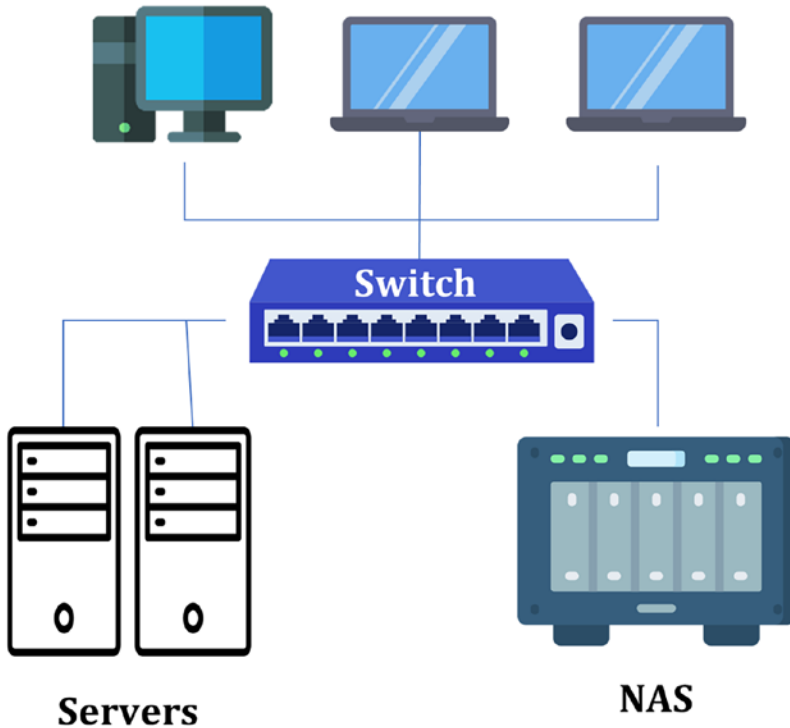


Figure 4-6. Network-Attached Storage (Icons from Flaticon, 2023)

Storage Area Network (SAN)

A storage area network (SAN), also called a storage network, provides access to a group of storage devices that include disk arrays and, sometimes, tape libraries. A SAN is a network version of a DAS. Access is provided to block-level storage that appears as direct attached storage to a computer's operating system. Unlike NAS, the storage devices are not

connected through the LAN, but through a dedicated storage network. Most SAN systems use the Fibre Channel Protocol (FCP), where the physical connection is a fiber cable.

SAN systems were historically built as a centralized data storage model where the storage system is in the center while the computer systems are peripherals (see Figure 4-7). This is in opposition to the classic model, where the computer system is in the center and storage is considered as peripheral. This model comes from the concept that data is most important for applications, while computing is a function. In configurations where disk systems were connected to a single computer system, the computer becomes the single point of failure. When SAN-based configurations are used, multiple computer systems can be connected. If the data within the SAN system is protected by methods such as duplication or RAID, the system becomes tolerant to single or even multiple failures.

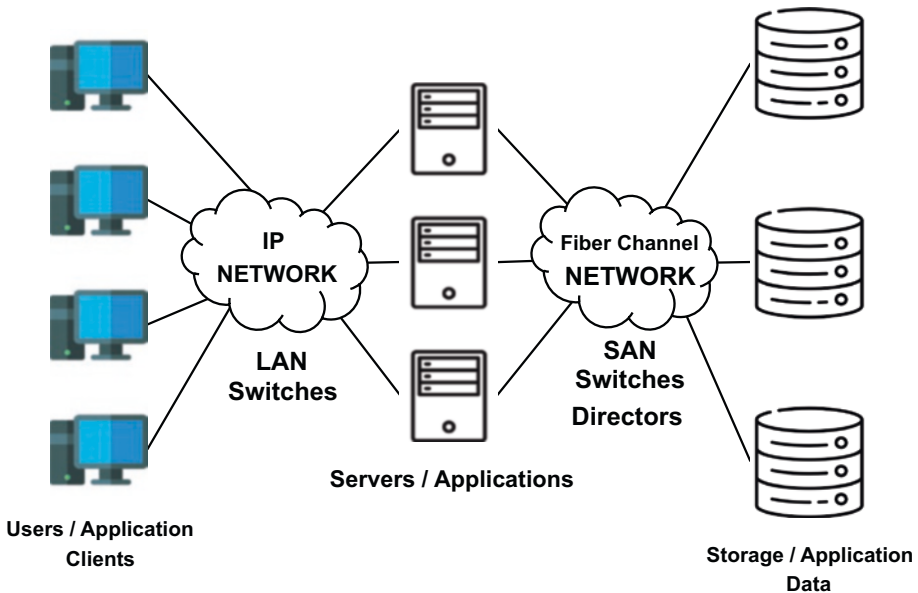


Figure 4-7. Storage Area Network (Icons from Flaticon, 2023)

SAN systems use specialized protocols such as Fibre Channel, iSCSI, and InfiniBand. SANs therefore need to use more expensive equipment, including a network switch, such as a fiber switch for fiber cables, specialized computer systems, and storage devices supporting the SAN protocol. This makes SANs more expensive than NAS, positioning them as high-end storage systems.

4.3 External Disk Configurations

DAS, NAS, and SAN systems can be added to computer systems to provide storage capacity. The storage systems can be deployed with different types of devices, such as HDDs and SSDs. The devices can be built in many different configurations, from JBOD to variations of RAID levels.

JBOD

If a disk enclosure only contains standard HDD or SDD devices without additional components, it is called JBOD, short for “Just a Bunch of Disks,” which explains what it is (see Figure 4-8). Some vendors offer software products as operating system add-ons that emulate the logic of disk array controllers. This allows configuration of JBOD devices in RAID configurations, but this method does not reach the performance of disk arrays.

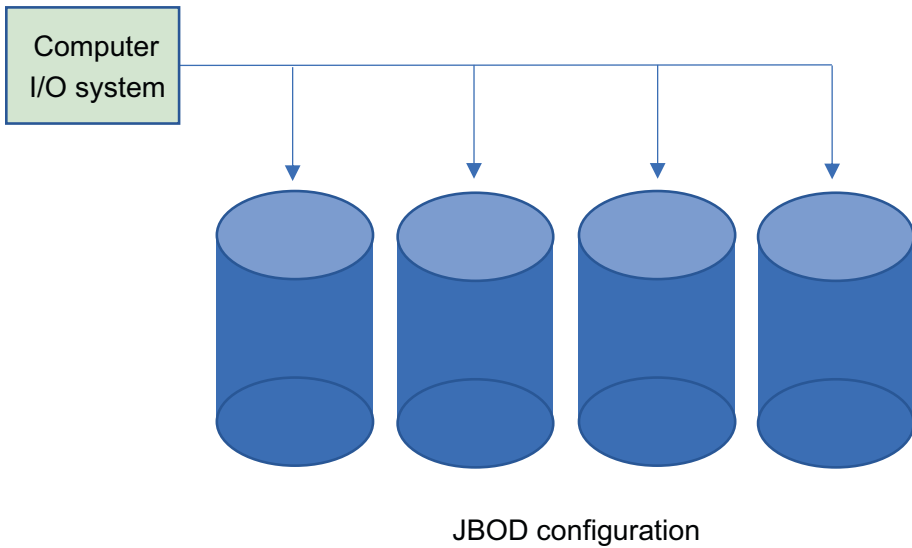


Figure 4-8. JBOD

Disk Arrays

Disk arrays contain the following components:

- Disk array controller
- Cache
- HDD or SSD drives
- Power supply

Disk arrays use additional redundant components to increase secondary storage's availability, resiliency, and maintainability (see Figure 4-9). Higher-end disk arrays include so much redundancy that the system does not have a single point of failure and is, therefore, very reliable. Disk arrays with hot-swappable features allow to replace components without stopping the system, which further improves the uptime.



Figure 4-9. *Disk Array (Macrovector, 2023)*

The disk array controller provides functionality that allows configuring the array in different ways to optimize it in respect of performance and data security. Inexpensive standard drives can be used in arrays with functionalities that only high-end storage systems offer otherwise.

RAID

A common expression for a disk array is “Redundant Array of Independent Disks,” in short, RAID. The disk array controller’s firmware allows operating the disks in different modes of operations, which are called RAID levels and hereby feature several configurations to optimize redundancy, capacity, and speed. The Storage Networking Industry Association standardized the various RAID levels. It is important to note that disk arrays offer several methods to protect data from hardware errors but do not protect against data loss caused by events affecting a DC, like fire or flooding.

RAID Levels

The next few paragraphs describe the various RAID levels that were defined. The most used RAID levels are 0, 1, 5, and 10.

RAID 0

In RAID 0, a set of disks of a disk array is striped, which means the data is evenly distributed over the disk set. It does not protect against any data failure. This configuration is also called “stripe set.” RAID 0 configuration provides the highest speed of all RAID levels and is used for applications where data access speed is most important, but data protection is not. The speed increases almost proportionally with the number of devices. The other advantage of RAID 0 is to provide a large logical file system volume. The size is as big as the sum of the disk capacity of all striped disks, provided they have the same size, which is the usual case in disk array configurations.

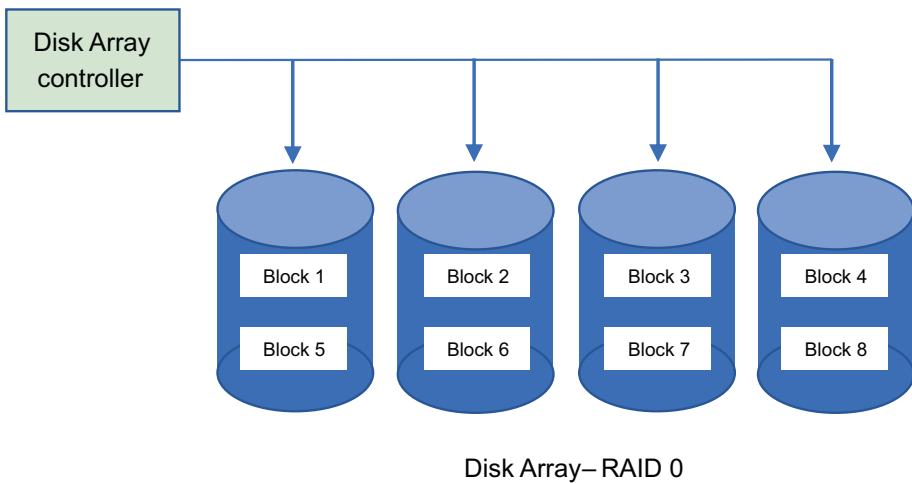


Figure 4-10. RAID 0

Figure 4-10 shows a disk array with four disks in a RAID 0 configuration. If each disk has a capacity of 2 TB, the configuration provides a single file system with a capacity of 8 TB. In RAID 0, data is accessed in blocks. If a block has 8 KB in a four-disk configuration, a stripe contains four blocks with a size of 32 KB. The first stripe containing Blocks 1-4 is distributed over the four disks, and the data access speed can be up to four times the single disk performance. This configuration does not feature redundancy. But this is acceptable in applications such as scientific computing or computer gaming where access speed is most important.

RAID 1

In RAID 1 configuration for disk arrays, data from one disk is mirrored onto a second disk (see Figure 4-11). It does not offer striping or spanning over multiple disks. More than two disks can be configured in a RAID 1 configuration. Still, additional disks will only increase the reliability but not increase the available disk size, because all disks in this configuration will be mirrors of the first disk.

RAID 1 configurations are used if reliability and read performance are most important while write performance and disk capacity can be compromised. Most RAID 1 configurations use two disks for mirroring, but for some very critical applications, more than two disks can be used. The array will function if one drive is operational.

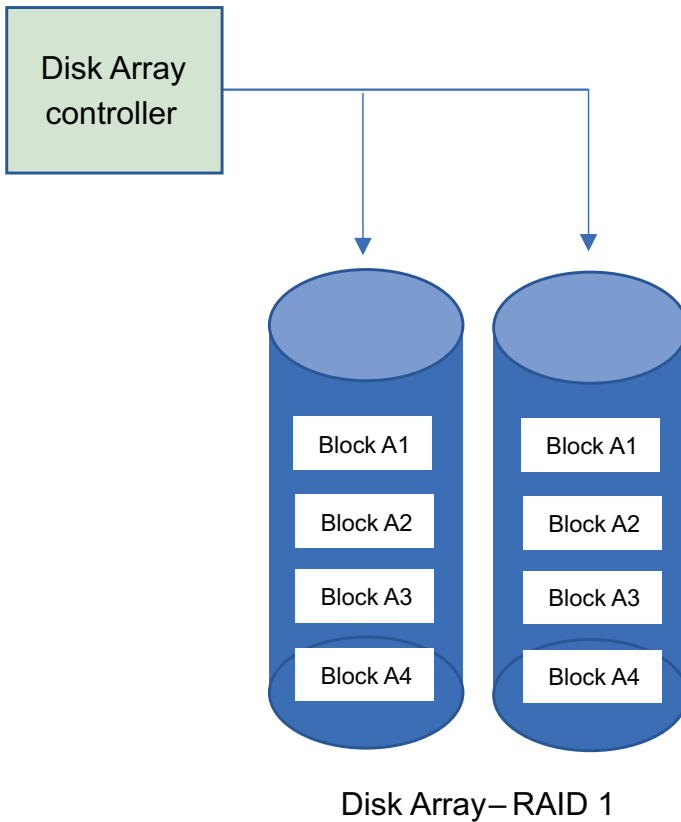


Figure 4-11. RAID 1

Read requests can be distributed over the available number of disks in the array, increasing the read performance similar to RAID 0. But the write speed is limited by the write speed of one disk. If different types of disks are used in RAID 1, which is uncommon, the disk size is determined by the smallest disk, and the write speed by the slowest unit.

RAID 2

RAID 2 uses disk striping like RAID 0, but the data is striped on bit level instead of on block level. A hamming code is used for error correction. It allows many disks to operate in an array with very high data transfer

rates. Some Data Vault disk arrays use 32 disks in parallel. But the synchronization of the disks, which have their own internal error correction logic, makes the overall error correction very complex, resulting in small advantages over other RAID configurations. RAID 2 is only used in a few disk configurations.

RAID 3

RAID 3 uses byte-level striping with an additional parity disk. RAID 3 configurations cannot service multiple requests because each block of data is spread across all disks of the array and resides at the same physical location on each disk. RAID 3 is the most suitable for applications requiring high data rates in long sequential reads and writes like video editing. On the other hand, applications using short data transfers may randomly experience a very low performance level. RAID 3 is therefore not widely used.

RAID 4

Similar to RAID 3, RAID 4 uses data striping with an additional parity disk. The difference is that RAID 4 uses block-level striping instead of byte-level striping. It provides good performance of random reads, but random writes performance [looks like there might be something missing here.] Also, it can quickly be extended online without the need to recompute the parity of existing data. RAID 4 is not widely used.

RAID 5

Disk arrays in RAID 5 configurations use block-level striping with distributed parity. The configuration tolerates the failures of one disk. If a drive fails, subsequent reads are calculated from the distributed parity blocks.

RAID 5 can be implemented in different layouts, namely:

- The data blocks can be written from left to right on the disk array, or from right to left.
- The parity block can be written at the beginning or at the end of the stripe.
- The location of the first block of a stripe with respect to the parity of the previous stripe.

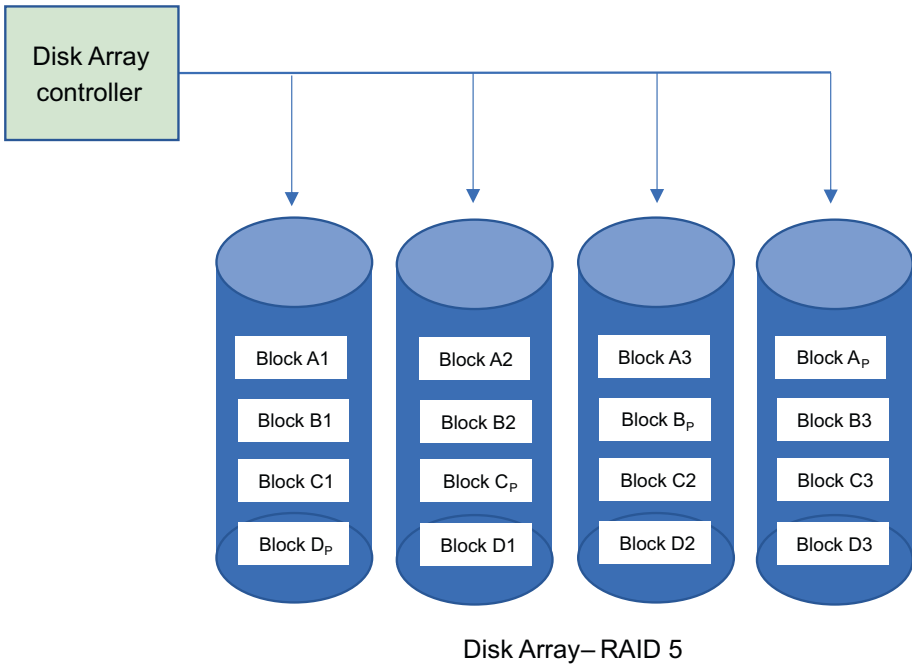


Figure 4-12. RAID 5

Figure 4-12 shows a RAID 5 configuration with four disks. A stripe contains four data blocks (A1:A3) and one parity block (AP) in left-to-right layout. The second stripe (B1:B3, BP) is written in the same way, but the parity block moves to another disk.

The distribution of the parity blocks over all disks eliminates the constraints of RAID 4, which is caused by its dedicated parity disk. As all RAID members participate in the serving of write requests, the write performance is better than arrays in a RAID 4 configuration. RAID 5 is a popular configuration, offering a good compromise between transfer rates, available array capacity, and data reliability.

RAID 6

RAID 6 provides a higher level of data security than RAID 5 by using two parity blocks instead of one. It is configured by block-level striping with two parity blocks distributed to all disks. This can be implemented with the layout choices that were described for RAID 5. Configurations using RAID 6 provide higher data reliability. Failures of up to two disks can be mitigated, but it reduces array storage capacity and transfer rates for write operations.

Nested RAID Levels

The combination of RAID levels, also called hybrid RAID, combines two or three RAID levels. The first number of a combined RAID describes the lowest level. This means RAID 01 is on the lowest level striped (RAID 0) and on the next level mirrored (RAID 1).

Nested RAID levels include RAID 01, RAID 10, RAID 100, RAID 50, and RAID 60, all resulting in a combination of disk striping with other RAID levels. The most used nested configurations are RAID 01 and RAID 10.

RAID 01

RAID 01, also called RAID 0+1, is a configuration that uses a mirror of striped disks. RAID 01 configurations have the same capacity as RAID 1. At least four disks are required for RAID 01 configuration, but it can be extended to more disks.

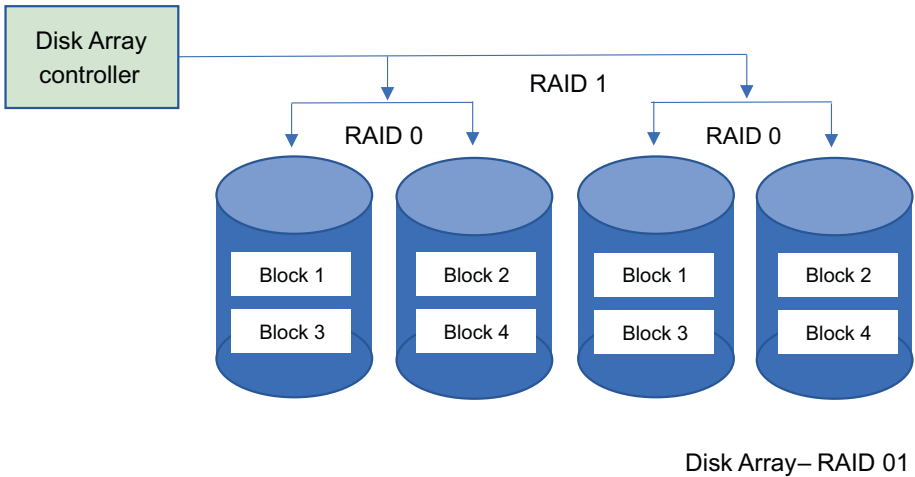


Figure 4-13. Disk Array – RAID 01

Figure 4-13 shows how data is stored in a RAID 01 configuration. With four available disks, a stripe includes two blocks. The first stripe contains blocks B1 and B2 that are stored on the first two disks. At the same time, the blocks are mirrored on disks 3 and 4.

RAID 10

RAID 10, also called RAID 1+0, is like RAID 01, but the layer is in the opposite order. RAID 10 is a stripe of mirrored disks that provides better throughput and latency than all other RAID levels, except for RAID 0. RAID 10 is therefore a widely used configuration well suited for I/O-intensive applications that require high data transmission rates, such as web server, email back-office, and database server.

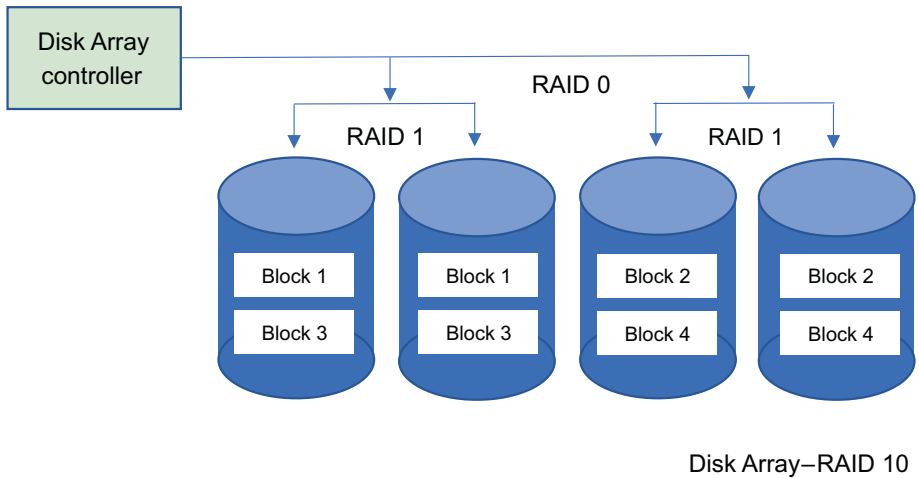


Figure 4-14. Disk Array – RAID 10

Figure 4-14 shows how data is stored in a RAID 10 configuration. With four available disks, a stripe includes two blocks. The first stripe contains blocks B1 and B2; B1 is stored on disk 1 and mirrored on disk 2, while B2 is stored on disk 3 and mirrored on disk 4.

RAID 50

RAID 50, also called RAID 5+0, provides a striped layer on top of the underlying RAID 5 distributed parity configurations. RAID 50 can be configured with a minimum of six disks. It can tolerate up to four disk failures, but only one per RAID 5 set. RAID 50 provides a higher performance than a RAID 5 configuration for write operations. It is suitable for applications that require a higher level of fault tolerance, capacity, and random-access performance, but as more disks are involved, the recovery time in case of faults is longer.

Figure 4-15 shows a RAID 50 configuration with eight disks. A stripe contains six data blocks (A1:A4, BP) and two parity blocks (AP) in left-to-right layout. The parity block is stored in both RAID 5 sublevels. The second stripe (B1:B4, BP) is written in the same way, but the parity block moves to another disk. This pattern repeats for stripes C and D.

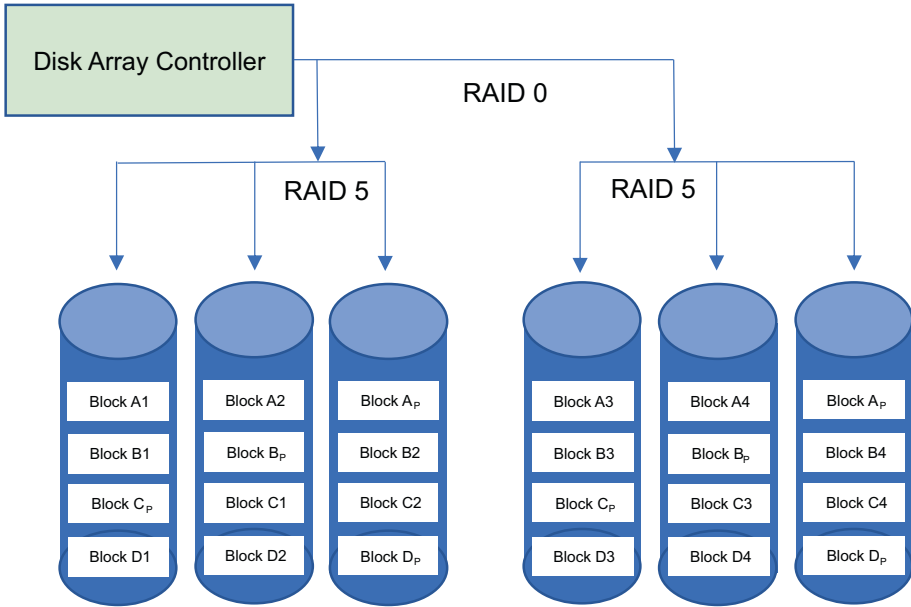


Figure 4-15. Disk Array – RAID 50

4.4 Storage Virtualization

Storage virtualization provides a logical view of the physical storage resources of a computer system. It creates a single resource pool of storage that includes all storage media that are accessible in the form of HDDs, SSDs, optical drives, or tapes.

Storage virtualization is based on software that can identify available storage capacity from a range of physical storage devices and aggregate the storage capacity into a single storage reservoir, which can be managed

from a central console. The aggregated storage can then be separated into multiple virtual volumes, which looks like physical external storage to a computer's operating system. RAID systems provide a low-level form of disk virtualization when any form of striping is configured. Any RAID level that uses striping combines physical disk space of multiple disks into a unified logical volume.

The functionality of storage virtualization goes far beyond RAID striping. Any form of disk systems, and even tape libraries can be combined into virtual disks. The virtual storage software intercepts data transfer from computer systems or virtual machines to a storage system. It then replaces it with an actual data transfer to a physical disk or a disk array. The physical storage resources are invisible to the operating system and the applications, and they can only interact with the logical storage volume provided by the virtualization layer.

Access Modes for Virtualized Storage

Storage virtualization can operate in file-access mode and block-access mode.

- File-based storage virtualization, which is used for NAS, supports file-sharing protocols such as Server Message Block (SMB), Common Internet File System (CIFS) (both used for Windows operating system environments), or Network File System (NFS) for Unix/Linux-based operating systems. Storage virtualization breaks the direct link between a computer system and NAS systems by pooling all NAS units into a single storage volume. This makes it easier to perform file migrations without interrupting the applications. It also simplifies the overall administration of NAS systems.

- Block-based disk access is used for data transfer to storage in SANs, which are usually connected by protocol stacks such as FC and iSCSI. Block-based disk access is more frequently used in virtual storage systems than in file-based storage systems. External block storage, in comparison with file-based storage systems, is like internal computer storage systems, reducing overheads for read and write operations. Block-based operation enables virtualization management software to collect the capacity of the available blocks of storage space across all virtualized arrays. This virtualized storage capacity can be accessed by conventional computer systems, nowadays called bare-metal, virtual machines, or containers.

Types of Storage Virtualization

Storage virtualization can be configured for different aggregates (see Figure 4-16). The lowest aggregation level is block virtualization, where blocks from multiple storage units are unified into a logical unit, but the transaction can be performed on a block. File virtualization, which is the next level, is where files from multiple NAS systems are unified into unified namespace. The next level is disk and tape virtualization. A virtual device in the form of a disk or a tape can be created from multiple disks.

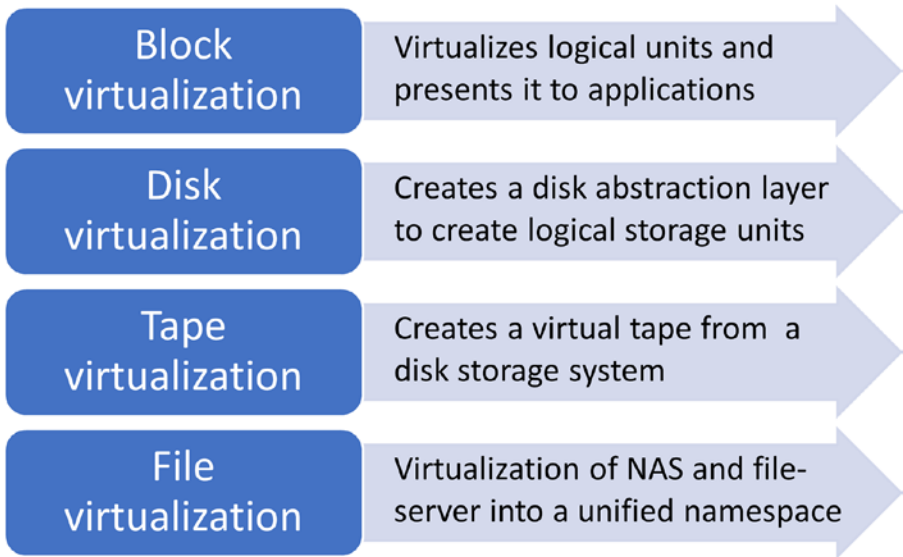


Figure 4-16. *Storage Virtualization Types*

Symmetric and Asymmetric Virtualization

Storage virtualization can operate in symmetric or asymmetric configurations.

- Symmetric virtualization, also called in-band virtualization, manages to read or write data as well as the related control information, such as I/O instructions or metadata, in the same layer.
- Asymmetric virtualization, also called out-of-band virtualization, splits the control path from the data paths.

Virtualization Methods

There are three main storage virtualization methods available: host-based storage, array-based storage, and network-based storage.

- **Host-based storage virtualization**

As the headline suggests, storage virtualization began in the 1980s when “host” was the word commonly used for computer systems, which was predominantly an IBM mainframe system then. Today, the term “host” applies to classic central computers such as mainframes and other types of computer systems, including large, virtualized server systems in cloud data centers.

“Host-based virtualization” refers to software that runs directly on a computer system. The software creates a virtualized file system for the “host” that includes internal and external storage resources. Virtualization and management are performed by virtualization software. Some operating systems such as Windows Servers have integrated storage virtualization functionality.

- **Array-based storage**

When the storage virtualization software runs on a storage array itself, the virtualization method is called “array-based storage.” In this configuration, the storage array controller can pool storage resources of other arrays to provide a logical file system from the combined storage arrays. The actual physical infrastructure of the storage system in the form of

HDDs, SSDs, and tapes is not visible to the computer systems accessing the virtualized array-based storage.

- **Network-based storage**

Today, the most used form of storage virtualization is “network-based storage.” In this virtualization method, a network device, such as a smart switch or a purpose-built server, connects all storage devices of a SAN and presents the combined storage resources as a single, virtual pool to the computer system.

Figure 4-17 shows the architecture of storage virtualization.

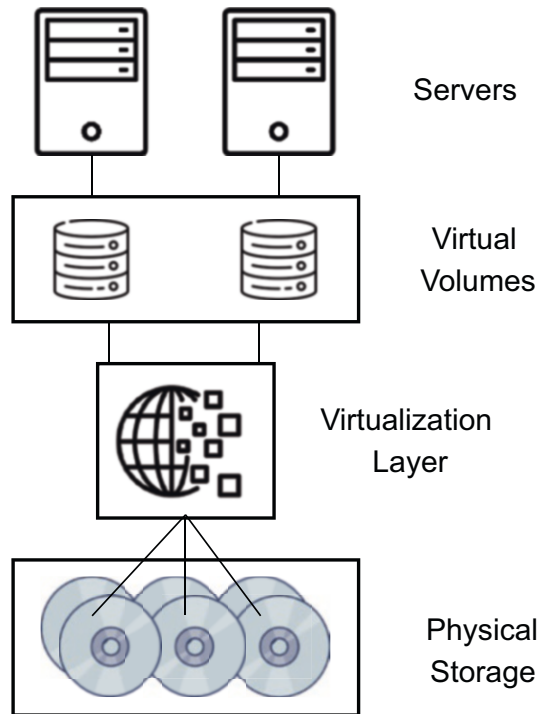


Figure 4-17. Storage Virtualization Architecture (Icons from Flaticon, 2023)

Storage Virtualization Benefits

The early versions of storage virtualization were all host-based storage virtualizations. These implementations were challenging as all computer systems were based on proprietary technology. A storage system for an IBM mainframe only worked with a range of such mainframe types and a few compatible third-party vendors. Computer and storage systems from companies like Digital or HP had different technologies and protocols for everything, including CPUs, I/O interfaces, communication protocols, and operating systems. Storage virtualization software had to be installed and maintained on all servers, which was only possible for a few computer systems from the same vendor.

As the IT industry matured, standards became more common and computer systems as well as storage systems could be mixed and matched. This was driven through the evolution of standardized protocols like SCSI, FC, Ethernet, TCP/IP, Unix, and NFS, to name a few.

Over time, virtualization software evolved, and standards such as Storage Management Initiative Specification (SMI-S) were developed. This allowed the implementation of storage virtualization with a wider range of storage systems.

Storage virtualization provides the following benefits:

- **Simplified management**

Storage virtualization allows the monitoring and maintaining of multiple virtualized storage arrays from a centralized console. This reduces the time to manage physical storage and makes it possible to create a storage pool from multiple storage vendors in an efficient way.

- **Better storage utilization**

Through storage virtualization, storage capacity can be allocated and used more efficiently by pooling storage capacity across multiple systems. Without the resource pooling features of a storage virtualization system, storage systems will be directly allocated to computer systems. The effect is that some storage systems were used at their limit, while others had a low utilization rate.

- **Storage life cycle extension**

Virtualization can extend the life cycle of storage. Through resource pooling, older storage systems can be used for archiving or handling less critical data.

- **Uniform storage features**

Many advanced storage features like tiering, caching, or replication may not be available on all storage systems, or the actual feature may differ substantially. With storage virtualization, such advanced features can be implemented at the virtualization layer, hereby unifying the features of the storage system.

4.5 Storage Security and Resilience

ISO defines storage security as the “application of physical, technical, and administrative controls to protect storage systems and infrastructure as well as the data stored within them” (see Figure 4-18).



Figure 4-18. *Data Security and Resiliency (Freepik, 2023)*

Storage security is focused on protecting data and its underlying storage infrastructure against any form of unauthorized access, copying, modification, or destruction of data while assuring that the data is available to authorized users. This means storage security must first prevent and discourage any form of unauthorized access.

Data Resiliency

Suppose the main objective of data security has failed and unauthorized personnel is able to access protected data. In that case, the system must be able to detect the breach and prevent further damages. A resilient system can detect breaches immediately and is able to restore the data quickly. After the data is restored, measurements are implemented to prevent a repetition of such data breaches.

Storage security comprises all manual or automated technologies and processes to ensure the integrity and security of data in storage systems. Data stored in computer systems, portable storage devices, or external

storage in the form of NAS, SAN, or DAS are protected by the storage security system using physical hardware protection and security software (see an example in Figure 4-19). For every organization that handles sensitive data, protection is essential to avoid data theft and to ensure continuous operations of its business activities.

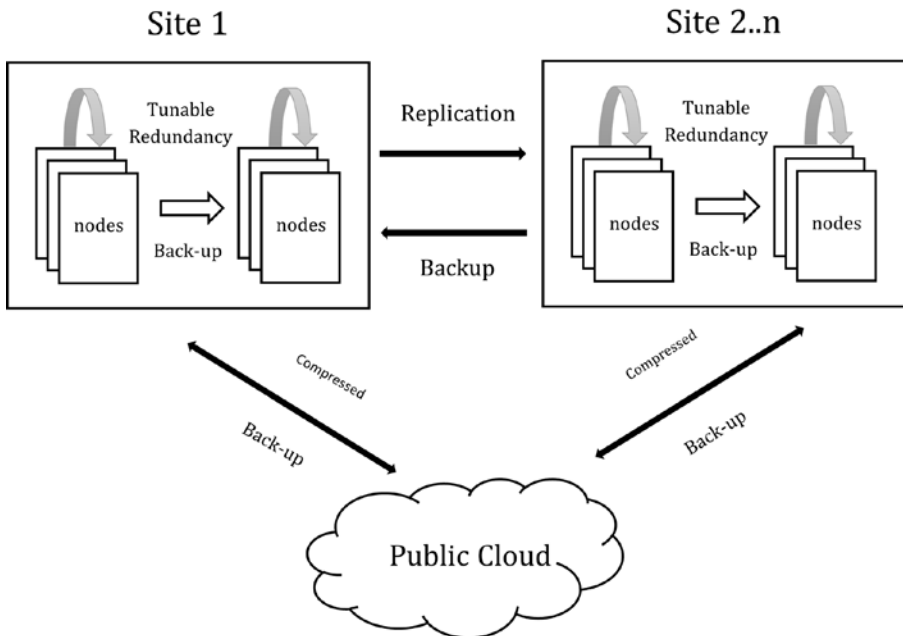


Figure 4-19. Example of a Resilient Storage Architecture

Key Methods to Ensure Storage Data Security

- Infrastructure security
- Preventing physical access to storage systems by unauthorized personnel
- Data encryption
- Access control mechanisms at all levels of data storage systems

- Protection against malicious software in the form of viruses, worms, and other threats
- Security as an important design principle of a storage architecture

Data protection and data security are key functions to ensure storage protection.

Data Protection

The objective of data protection is to ensure data availability in the event of hardware failures in computer and storage systems, power outages, or natural disasters.

Data Security

The key goal of data security is ensuring that any unauthorized user will not access private and confidential information. Private data must be protected against any form of unauthorized reading, copying, altering, or deleting. See Figure 4-20 to learn the three principles of data protection.



Data Security

Are you who you say you are?



Access Control

Prove you are who you say you are



Data Availability

Ensuring data is ready to support business operations

Figure 4-20. *Data Protection Principles (Icons from Flaticon, 2023)*

Effective data protection systems and procedures can only be built if an organization understands potential sources of threats.

Threats to Data Security

External

- Espionage by foreign countries
- Terrorist attacks
- Organized crime
- Hackers and cybercriminals
- Industrial espionage by competitors
- Power outages
- Flooding and other natural disasters

Internal

- Weak security procedures
- Lack of awareness of cybercriminal threats
- Lack of staff training for data security
- Malicious insiders
- Angry employees

Any internal and external attempt to breach data security will start with exploring the vulnerabilities of an organization's data security system and procedures.

Areas of Vulnerabilities of Storage System

Cloud Storage

A rising number of organizations use CSPs to perform their workloads and therefore have data stored in the cloud. Some CSPs execute best-in-class systems and procedures to protect data, and those CSPs may provide better data security than many private DCs. But it may be difficult for an organization to assess the data security quality of a CSP. Cloud storage may pose a vulnerability to an organization if its staff lacks the training or knowledge necessary to efficiently utilize the tools and procedures provided by a cloud service provider (CSP), thereby creating security risks.

Physical Security

There are several aspects to consider ensuring that data cannot be accessed physically:

- Access control of DC buildings and rooms
- Insiders that have access to the facilities for other purposes, such as security guards or cleaning crews
- Internal staff that do not have authorization to access DC facilities

Physical security vulnerabilities are often overlooked, and data protection may often be too focused on preventing cybercrime attacks.

Data Encryption

If an attacker could breach the access control systems, the next level of defense is data encryption, which is important to protect against physical access breaches. While most high-end storage systems automatically perform data encryption, most other storage systems do not provide this functionality. Hence, organizations must install encryption software to ensure that the stored data is encrypted. This includes all desktop and

laptop computers used by an organization's workforce and requires that all employees follow compulsory IT security procedures for their personal computer devices when they transfer data of their organization.

Deletion of Data

If files or even a disk volume is deleted, the file system only deletes the entries in the file system tables and not the actual data blocks on the disk. Software is available from several vendors that can recover such data. People often use this type of software when they experience a disk failure without having created a recently updated backup. Therefore, the IT organization must ensure that data is fully erased when disks or tapes are replaced or disposed of.

Data Storage Security Principles

The principles of data security are availability, integrity, and confidentiality (see Figure 4-21).

Availability

To ensure that data is always available. Data must therefore be protected against internal and external attacks. In case of a computer or storage system failure, or other disruption like a power failure, the data must be restored quickly.

Integrity

To ensure that data cannot be changed or be tampered with.

Confidentiality

To ensure that data is kept confidential by ensuring that an unauthorized person cannot access it through physical access or through network access from an internal or external network.

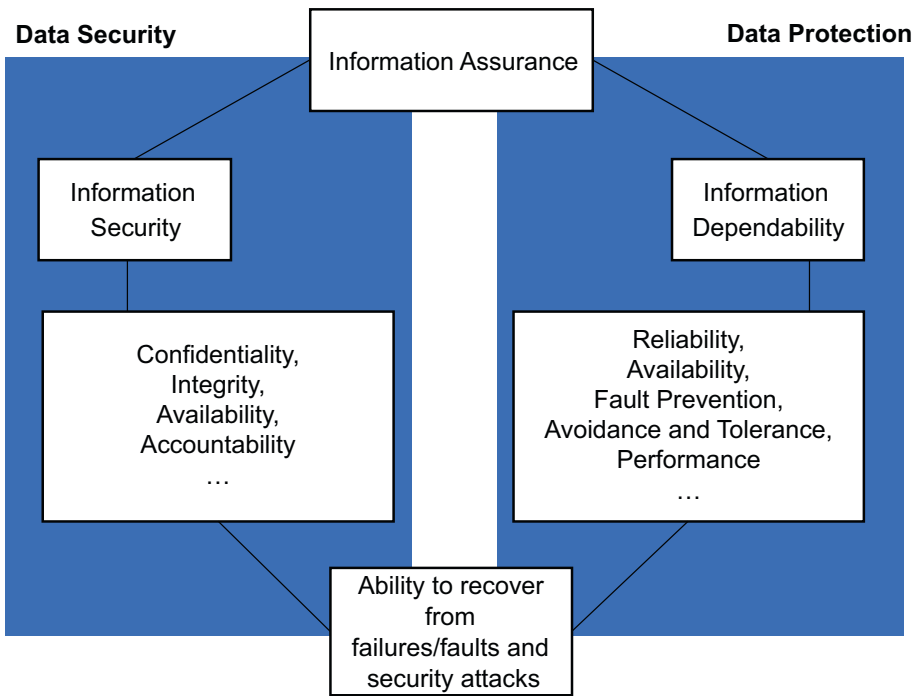


Figure 4-21. Processes for Data Resiliency

Best Practices for Data Resiliency

To provide best-in-class data security and resiliency, the following best practices should be applied:

- **Strong access security**

Data stored in private DCs or clouds must be protected against any form of physical access by unauthorized personnel. Data can be accessed through internal or external networks besides physical access. Security systems, such as firewalls, anti-malware protection, security gateways, and intrusion detection systems, are therefore important elements of any form of data security. More details are provided in the chapter about DC security (see Chapter 7).

- **Strong endpoint security**

Data security includes all devices that the staff of the organization uses to access its data. Hence, all security procedures must apply to such devices, including desktops, laptops, tablets, and smartphones. This is often a big challenge for organizations that allow staff to “bring your own devices,” which is called BYOD.

- **Access controls**

The implementation of a strong access control mechanism is important to secure stored data. Role-based access control must be strictly implemented; for confidential and finance-related data, multifactor authentication is strongly recommended. The IT organization’s administration must ensure that all users define strong passwords to access their accounts and that frequent password changes are mandated.

- **Data loss prevention**

Similar to an antivirus system for a computing system, data loss prevention (DLP) solutions can detect malicious attacks on data and stop them.

- **Encryption**

All data needs to be encrypted during any data transfer and when the data is stored on the devices. The system administrator must deploy secure key management to track their encryption keys.

- **Redundancy**

Data redundancy means that data is protected against system failures, such as disk crashes. RAID-based storage systems offer many options to implement redundancy. Redundancy can be achieved in different forms and on different levels. RAID systems provide redundancy on an array level, and other methods copy data on JBOD storage to another JBOD system. On the highest level, data can be duplicated to a DC at a different location, preferably in another region or country.

- **Backup and recovery**

In some incidents, such as flooding, earthquakes, successful malware, or ransomware attacks, the data may no longer be restorable from the storage systems. In such cases, data can only be restored from backups. An appropriate backup system must therefore be part of any data security practice. The IT administration must frequently test if the backup system works correctly and if data can be restored. An incident can turn into a disaster if the restoration from backup systems fails due to errors during the backup process and data is permanently lost.

The IT administration must ensure that backup processes and systems are adequate for any type of event. Backup systems must be treated with the same level of security severity as the storage process itself. Backup and recovery are essential elements to provide data resiliency beyond data security.

- **Security policies**

Data storage protection starts with proper design and implementation of data security systems. Architecture, procedures, and policies must be well documented. Data protection will be based on levels of security such as highly confidential, confidential, restricted, private, and public. The organization needs to have security models, procedures, and tools in place to protect data accordingly. Data security policies should include details on the security measures that need to be deployed on each storage device in use.

Storage Security Implementation

The first line of defense to protect data is the security of the DC network and the DC facilities itself. To provide a second layer of data protection, the best practices to apply include the following procedures:

- Strong protection of the DC network and protection against cyberattacks.
- Strong access security to the physical infrastructure, preferably by using biometrical access controls. Full video surveillance should cover all parts of the DCs and their entry points.
- Implementation of strong endpoint security by controlling the end devices of the network and securing the data through private data encryption methods.
- Strict implementation of role-based access controls for all users and functions of the IT systems and application.

- All data that is stored in HDD and SSD devices should be encrypted.
- Running the DC infrastructure in a federated network infrastructure is preferred.
- On the device level, all disk systems should be configured in RAID configuration with mirroring. Raid 10 implementation is one of the recommended RAIDs to provide performance and data security.
- A strong data backup system should be available. Details are described in Chapter 6.

4.6 Storage Provisioning and Administration

Storage provisioning is the process of assigning storage to workloads and to optimize the performance of the available storage resources.

Storage provisioning can be performed based on the underlying storage systems in the form of DAS, SAN, and NAS systems. Provisioning of storage is performed on block levels for DAS and SAN systems, while for NAS systems, storage is provisioned on a file level by a file-sharing system such as NFS.

Classic SAN Provisioning

When SAN systems were introduced, the storage controller offered a rich set of features, but the actual SAN administration was tedious. Most SAN systems were provisioned by creating a storage configuration template that defined the attributes and conditions to allocate the storage.

Storage Provisioning in Modern DCs

The introduction of cloud data centers, hypervisors, and container-based computing has changed the methods to provide storage as the SAN approach was a storage-centric approach. Virtualization suggests an application-centric storage provision and is based on a virtualized storage environment. This means virtual storage is assigned to computer systems, hypervisors, or containers during provisioning. The storage allocation is based on application-defined capacity, availability, and performance requirements and is called “thin provisioning.” Alternatively, storage provisioning can be performed traditionally, called “thick provisioning.” See a representation of both storage provisioning types in Figure 4-22.

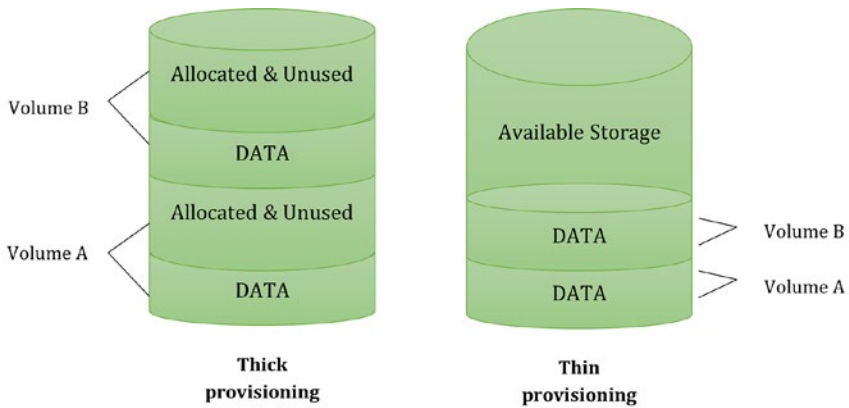


Figure 4-22. *Thick and Thin Provisioning*

Thick Provisioning

Disks are provisioned into storage pools that can host both block and file data. Connectivity is offered for both block and file protocols. A storage pool is an aggregated group of physical disks where different RAID levels

can be applied to provide redundancy. RAID protection is directly applied to discrete groups of drives within the storage pool.

Traditionally, storage volumes are partitioned into logical disks, which fundamentally create a representation layer on top of physical disks by presenting them as a logical volume to operating systems and applications. In this context, provisioning of disk space means the allocation of logical disk space from a logical volume. Configuring, changing, or moving disk space is easier and quicker to implement in comparison to thick provisioning.

An example of thick provision is the formatting of the disk using the Windows operating system. An external disk will usually be formatted into several logical volumes, which will be associated to applications. Once the disk allocation is performed and data is loaded, it is difficult to change the size of the volumes.

Thin Provisioning

In thin provisioning, disk space is strategically pre-allocated to a server or a virtual machine. This means that the logical space provided by partitioning equals the amount of actual physical space set aside on the disk.

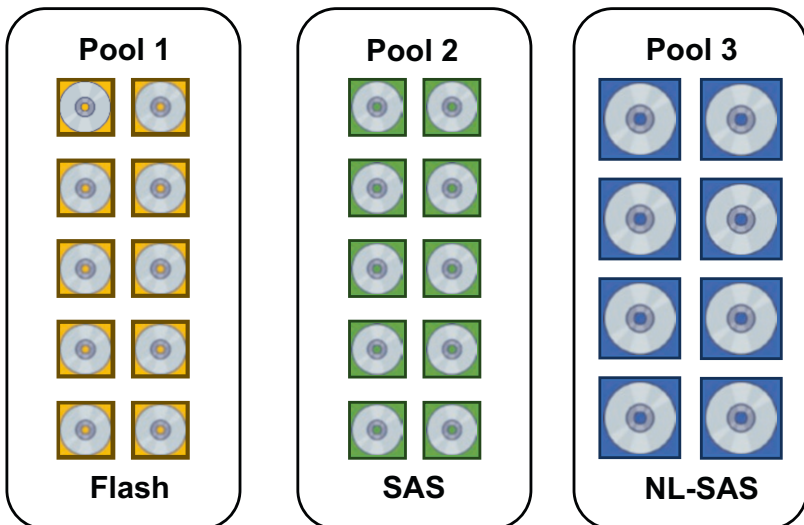
Storage Pools

Storage pools can consist of storage with different performance characteristics. Take a look at [Table 4-3](#).

Table 4-3. *Storage Pools*

Very fast	Flash-based storage
Fast	Disk system connected through an I/O channel, such as SCSI, called SAS
Less fast	Nearline SAS, larger disks with slower access speed

These three disks can be sorted into pools for flash, SAS, and NL-SAS, as shown in Figure 4-23. This is called a “Homogeneous Storage Pool.”

**Figure 4-23.** *Homogeneous Storage Pool (Icons from Flaticon, 2023)*

Some high-end storage systems provide an automatic storage tiering feature for virtual pools using a technology called “FAST-VP.” Through this technology, heterogeneous storage pools can provide efficient balancing of data between tiers without requiring user intervention. This is called a “Heterogeneous Storage Pool” (see Figure 4-24).

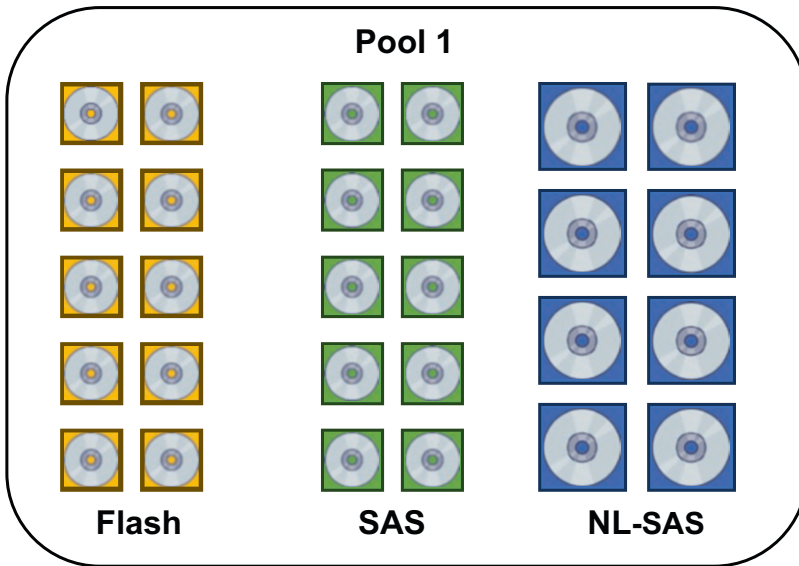


Figure 4-24. *Heterogeneous Storage Pool (Icons from Flaticon, 2023)*

It is recommended that a storage pool always has at least 10% free capacity to maintain proper functioning. The number of storage pools should be limited to reduce complexity and to increase flexibility.

Guidelines for good provisioning of storage pools are as follows:

- Configuration of all-flash pools
- Resource separation for multitenancy
- Separation of workloads with different I/O profiles
- Separate pools for workloads using fast cache
- Dedicated storage resources for specific performance goals

Storage Allocation Tiering

Table 4-4 shows storage allocation tiering.

Table 4-4. *Storage Allocation Tiering*

Policy tier	Initial tier	Description
Start high, then auto	Highest available	Initial allocation to the highest available tier and then automatically classified according to the I/O activity of the workload
Auto	Optimal tier	Allocation to the optimal available tier based on the I/O activity of the workload
Highest	Highest available	The workload is allocated to the storage pool with the highest performance
Lowest	Lowest available	The workload is allocated to the storage pool with the lowest performance

Comparison of Public and Private Storage Infrastructure

Table 4-5 provides a comparison between the options of an organization to use storage for their workloads in a private cloud, public cloud, and hybrid cloud deployment.

Table 4-5. *Comparison Between Private, Public, and Hybrid Cloud Storage*

Feature	Private cloud	Public cloud	Hybrid cloud
Costs	Higher due to costs for room, equipment, and staff	Lower, linear incremental costs	Medium, depending on the mix between private and public
Scalable	Limited, expansion requires capital investment	Very high	High for the workloads running in the public cloud
Performance	Higher	Low to medium	Medium, depending on the workload mix
Reliability	Higher as all equipment is controlled by the organization	Medium on average, depending on the expertise and resources of the CSP	High, important data can be copied to a secured domain at the CSP
Security	High, as equipment and personnel are controlled by the organization	Medium on average, depending on the expertise and resources of the CSP	Higher than CSP as critical data is kept in the premises of the organization

4.7 Charging for Storage Resources

Most CSPs offer their services for storage capacity in flexible form. The main charging feature is obviously the storage space usually measured and charged monthly. There are, however, several other features determining the prices such as service types, location of the storage, security levels, and extra features such as data retention and deduplication. As a result, the costs of cloud storage can be untransparent.

The way storage services are offered and bundled varies widely. For example, Microsoft Azure offers durable and available storage products, like blob, queue, file, and disk storage. Amazon Web Services provides storage services, such as Simple Storage Service, Elastic Block Storage, and Amazon Glacier, while Google Cloud Storage offers premium-priced high-performance regional, multiregional storage and less expensive but slower options. This makes the comparison of storage pricing a rather complicated undertaking if an organization wants to use different types of storage services for their workloads (see Figure 4-25).

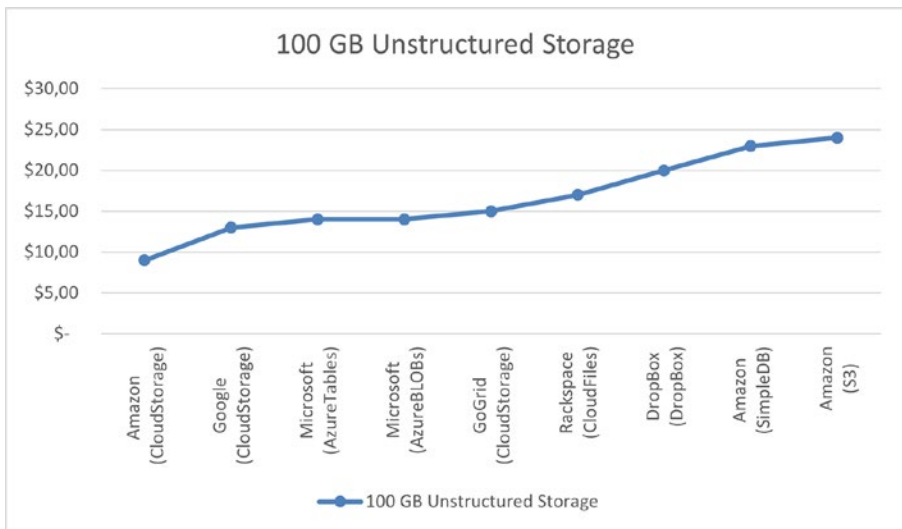


Figure 4-25. Comparison of Cloud Storage Charges for 100 GB

The typical pricing for cloud storage is currently in the region of 0.02 USD per GB for a capacity of 50 TB. This pricing only covers the storage capacity. If data is accessed or moved, additional fees will be charged. Data movement costs vary and depend on the service, tier, and resiliency options that an organization chooses for its workloads.

For example, Azure Hot Blob storage users with geographically redundant storage cost incur 1 USD per 100,000 for read or write operations, 0.04 USD per 100,000 for other operations, and 0.02 USD per GB of data transfers for geographical replications. Google Cloud Storage charges 0.12 USD per GB for up to 1 TB of data read; the tier until 10 TB is priced at 0.11 USD per GB, while the tier above 10 TB is priced at 0.08 USD per GB.

Some CSPs charge for storage space that is deleted when the minimum time frame for storage allocation is not reached. For example, Google Cloud Storage charges 0.008 USD per GB for data at rest and 0.05 USD per GB for data retrieval.

Cost-Saving Options for Cloud Storage

Most of the CSPs offer a range of storage services such as online, nearline, offline, archive, and other storage types and provide options of services for various workloads. If an organization accesses data infrequently, storage archiving and retrieving can be a more cost-effective method for some workloads. Some CSPs require a minimum billable service period of 30 days for many services. For data from certain workloads, it may be more cost-effective to use a different service level, like offline or nearline, if data is only used for a few days.

Thin-provisioning storage facilities are cheaper than thick-provision services. CSPs usually require defining the maximum capacity for a storage instance. Therefore, it is important to check that the storage billing reflects the actual used storage capacity. Some higher-end storage services, such as the Azure Premium Storage service, may bill for the entire allocation from the start.

To keep storage costs under control, it is important to establish well-defined policies outlining the life cycle of cloud resources. Monitoring tools, such as those offered by the cloud providers, are important for IT administrators to oversee public cloud usage and expenses. Excess usage must be prevented by quickly removing older, unused storage instances.

Data movement is another source of costs that can be avoided or reduced. Public cloud storage should not be used as the primary storage location as the movement of large amounts of cloud storage can be expensive. Therefore, it is prudent to perform data deduplication of local data. For some workloads an application redesign may require reducing the interaction of storage data and hereby reduce the costs.?

4.8 Summary

In this chapter, storage provisioning and administration were explored, focusing on optimizing storage resources for performance, particularly in the context of modern data centers. The shift from traditional, storage-centric provisioning to more application-centric approaches, such as thin provisioning, was discussed the concept of thick provisioning, where physical space is rigidly allocated, was contrasted with thin provisioning's flexibility.

Additionally, storage pools, which can consist of disks with varying performance characteristics, including flash-based storage, SAS disks, and NL-SAS disks were examined. These pools can be either homogeneous or heterogeneous, offering different levels of redundancy and performance.

There were also insights into the key considerations when charging for storage resources, highlighting the complexities associated with cloud storage pricing. Various factors, including storage capacity, service types, location, security levels, and additional features, impact the cost of cloud storage. The importance of understanding the pricing models of different cloud service providers to make informed decisions about storage was emphasized.

Furthermore, cost-saving strategies, such as choosing the right storage types for specific workloads, monitoring usage to prevent excess costs, and minimizing data movement expenses were discussed. These strategies can help organizations effectively manage their storage resources in a cost-efficient manner.

Overall, this chapter equips organizations with valuable knowledge and strategies for optimizing storage provisioning, administration, and cost management in today's dynamic IT landscape.

In the next chapter...

- Data center networking
 - Software-defined networking
 - Essential elements of a robust and resilient network architecture
-

CHAPTER 5

Network

Data center networking provides the resources to facilitate the storage and processing of applications by connecting servers and storage systems with each other and to external parties. Networking uses hardware in the form of switches and routers, and software, such as load balancers and analytic tools, to perform networking operations.

Modern data center (DC) networking architectures are leveraging full-stack networking and security virtualization platforms such as virtual machines, containers, and bare-metal configurations. This represents a significant shift from the classic networking model in DCs. The DC has evolved from containing a single system to many physical compute servers, to virtualized DCs, and eventually to the distributed DC infrastructures of today. The network components of the DC have evolved with similar progress as server and storage technology, and they are the third building block in the immense evolution of DCs.

DC networking platforms (see Figure 5-1) control the physical networking elements in the form of a software-defined network (SDN). It runs all network services required to run traditional enterprise workloads and enables the automatization of provisioning. The system allows capacity planning, security policy planning, and network troubleshooting. At the end of the workload life cycle, the networking platform executes the associated de-provisioning policies, hereby ensuring a high level of manageability, security, connectivity, and compliance of the overall DC operations.



Figure 5-1. DC Network (Victor217, 2023)

5.1 DC Network Components

DC network equipment consists of hardware in the form of cables, modems, switches, routers and software that is used to configure, monitor, and manage the DC network and its connectivity to the outside world. Its role is to transport data between storage and computer systems in the DC.

Cables

Cables are an essential part of every DC; they are the “pipes” through which the data flows.

DCs use mainly two types of cables.

Cables for Power Distribution

Most DCs use standard alternating current (AC) power supplies, but in some areas like the telco industry, direct current (DC) power supplies are used. The power cabling will be deployed according to the type of power that is used. Due to the high-power consumption of modern data centers, power cabling must be planned very well to provide sufficient power to the DC and all its electrical components.

Cables for Data Connection

Most of the data transporting cables in DCs are copper based. The classic Ethernet cable with its striking blue or yellow color is familiar to everyone who has connected a server or a PC to a LAN port (see Figure 5-2). Due to the constantly rising amount of data that is transported in DCs, fiber cables with a much higher bandwidth are gradually replacing the classic copper-based cables (see Figure 5-3). DCs are connecting to the Internet with fiber links, and the upper layer of data networks is connected through fiber connections.



Figure 5-2. Ethernet Cable (jannoon028, 2023)

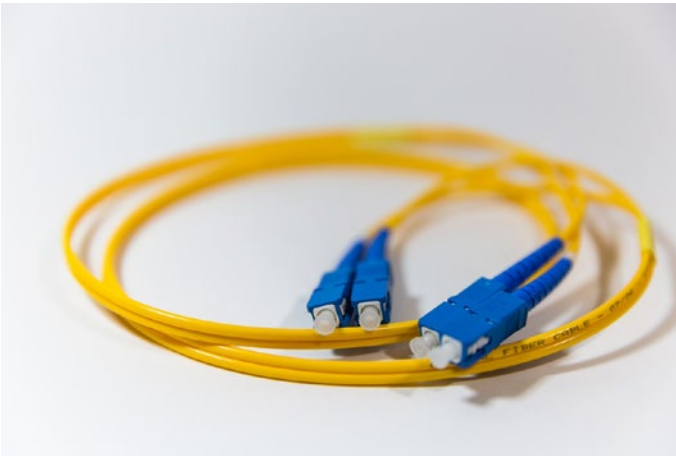


Figure 5-3. *Fiber-Optic Cable (PawinG, 2017)*

Structured Cabling

Structured cabling is used to avoid messy cabling that was often found in larger DCs in the past. It uses a main distribution area (MDA) as a central point of cable connections (see Figure 5-4). ANSI and the Telecommunications Industry Association (TIA) have defined the TIA-942 standard to define MDA as the central point of distribution for the data center structured cabling system. Using MDAs provides orderly structures for the cabling and allows easier changes to the DC and extensions.

The principles of structured cabling design and installation define the type of cables that should be used including LAN cable (Cat5e and Cat6), fiber cables, and modular connectors. The cables can be laid in several topologies depending on the size and features of a DC. It includes a central patch panel.

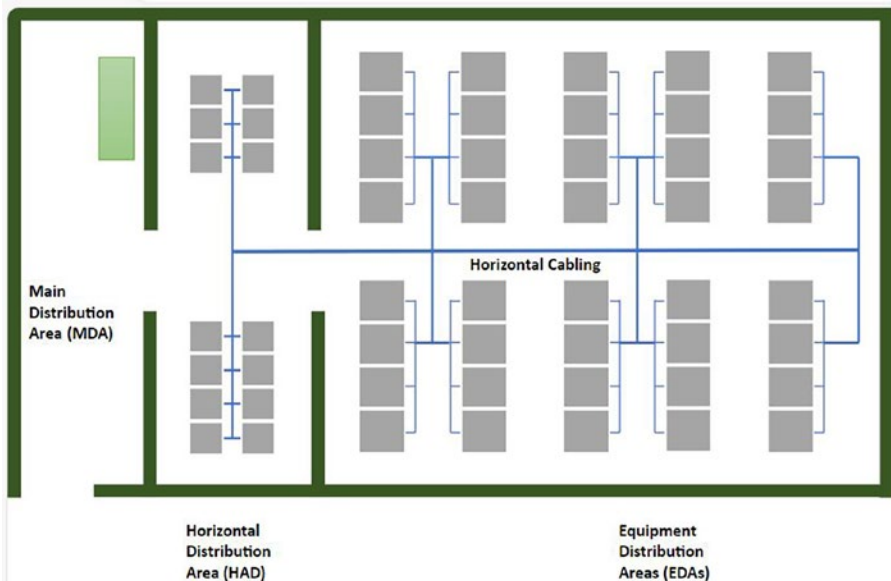


Figure 5-4. DC with MDA Cabling Structure

The voice communication system can be integrated into the MDA as well. With IP telephony, the network can run voice and data communication on the same cables, eliminating the need for separate phone wiring. For all copper-based cables, the maximum length is 90 meters from the MDA to the rack switches, plus an allowance of up to 10 meters for the patch cords on both ends of the cable.

Switches and Router

Even though switches and routers look quite similar as they use the same type of cabling and the same transmission protocols, their functions are different. Switches connect end devices such as server and storage systems, PCs, and printers. The main purpose of a switch is to receive a data packet, determine its destination, and forward the packet to the destination address. The role of a router is to find the fastest and best route for a packet to reach its destination.

Each connected network or computing device has a unique MAC (Media Access Control) address. When a device or a computer sends an IP packet to another device, the switch creates an IP packet that includes the MAC of the source and destination devices and encapsulates it. Therefore, a switch uses the MAC address to transport data packages, while a router uses IP addresses of devices to transport packages.

The functions of switches and routers can be classified with the well-defined seven-layer communication architecture from OSI: a switch works on the data link layer, while a router works on the network layer (see Figure 5-5).

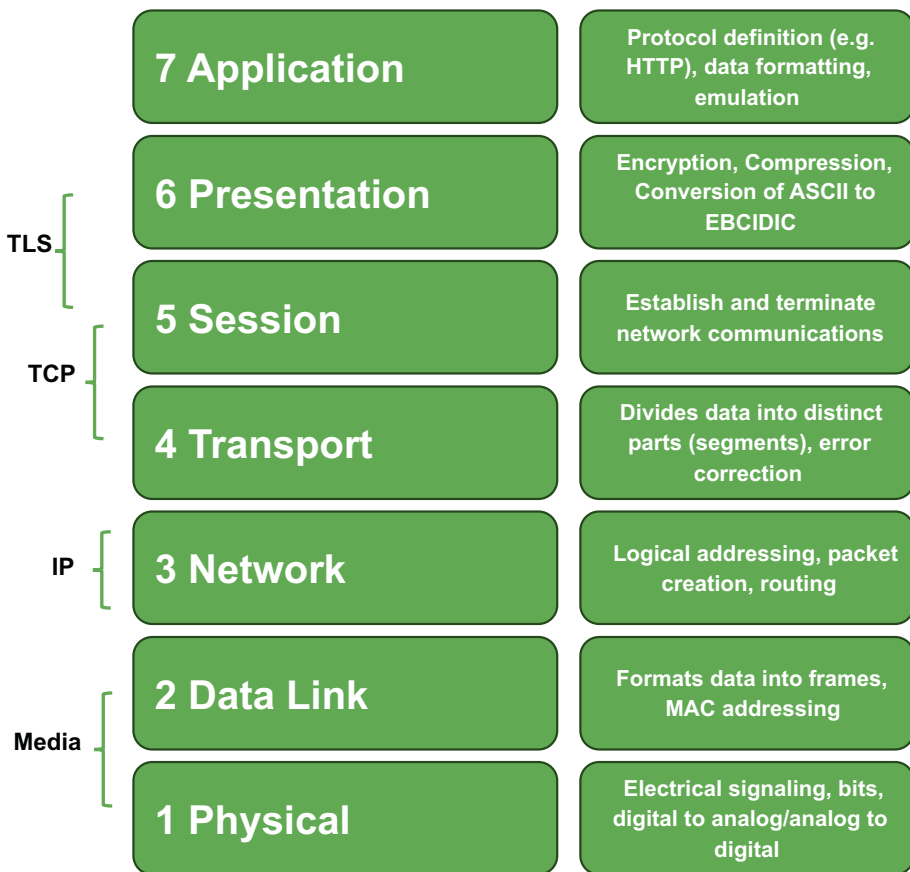


Figure 5-5. OSI Seven-Layer Communication Model

DC Switches

A large switch for a DC is high-performance network equipment used by large enterprises and cloud providers who have configured a virtualized DC environment. It can be deployed throughout the data center in several topologies, including single-tier and multi-tier configurations. The switches can effectively transport data within a rack row (east-west) and between rows (north-south). DC switches are deployed in top-of-rack and end-of-row configurations. The DC switches provide high bandwidth connections between all connected components. Standard TCP/IP LAN Ethernet protocol and SAN protocols are used for data transmission.

Functions

All switches use standards-based protocols such as TCP/IP and maintain tables for ports and MAC addresses. Different switches can be distinguished by their features such as redundancy, support of a virtualized infrastructure, and support for protocols such as EVPN, VXLAN, or OpenFlow.

Capacities

DC switches are available in smaller and bigger configurations. A larger switch will be on level 2 of a multi-tier architecture or in the center of a centralized topology. Smaller switches such as rack switches will be used in the lowest tier of a multi-tier architecture or in distributed topology like a mesh network.

Connection Speed

Switches are connected via fiber connections. Smaller DCs typically use fiber links with transmission rates of 10, 25, or 40 Gbps (gigabits per second). Larger DCs require faster connection speed and will deploy switches with a bandwidth of up to 4,000 Gbps.

Ports

A key switch parameter is the number of ports it can connect.

Network Management

An important consideration is the choice of management software for the DC network. In the past, most DCs were deployed from vendors like Cisco or Juniper Networks. The management software was provided along with equipment and created a proprietary network system. While computer and storage systems moved away from proprietary stacks to open standards in the 1990s, networking is only slowly changing to open standards.

Router

A router connects switches and their corresponding networks to build a large network (see Figure 5-6). In multi-tier network topologies, switches are deployed in the lower tiers, while routers connect the lower tiers to higher tiers. Switches connected to a router can be located at single or multiple locations.

Routers are intelligent devices routing the data packets from source to destination over a network. They transmit data packets from one network to another or within a network. The router provides a local IP address for each connected device and finds the best and fastest path to transport it.



Figure 5-6. Enterprise-Class Router (Unsplash – Albert Stoynov, 2023)

DC Gateway

A gateway (GW) transmits data from one discrete network to another. Gateways can be implemented in hardware or in software. The main difference between gateways and routers/switches is that a GW can transport data across different protocols. Gateways can also operate on all seven layers of the OSI communication model.

In the past, GWs were used to transmit data between different networks such as TCP/IP, SNA from IBM, or DECnet from Digital Equipment. With the proliferation of the Internet, TCP/IP has become the unified network protocol. Today, most GWs are used to separate networks such as corporate networks and the Internet.

A DC GW connects a DC to the Internet and provides access to virtual private networks (VPN) customers. They come in the form of simple routers connecting to the Internet and IP VPN devices.

The Internet connection is provided by an Internet service provider (ISP), a company that has a network connection to the backbones of the Internet. For most organizations, it is much easier to use the service of an

ISP to connect to the immense networks of the Internet. But very large cloud service providers such as Google and Amazon have built themselves very large private networks connecting to most countries in the world.

5.2 DC Network Topology

Depending on the requirements and size of a DC, several network topologies are available.

Centralized Topology

A centralized network topology is often used for smaller data centers (see Figure 5-7). The storage area network (SAN) and the local area network (LAN) are separated. The notion of centralized comes from the core switches in the main distribution area. All servers are cabled to the core switches. This provides very efficient utilization of the LAN and SAN switches. Managing and adding network, server, or storage components are easy to perform up to a certain size. Beyond this size, the centralized topology does not scale up well. Expansions become difficult as more cables with extended lengths are required, causing congestion in the cable pathways and cabinets and becoming increasingly costly.

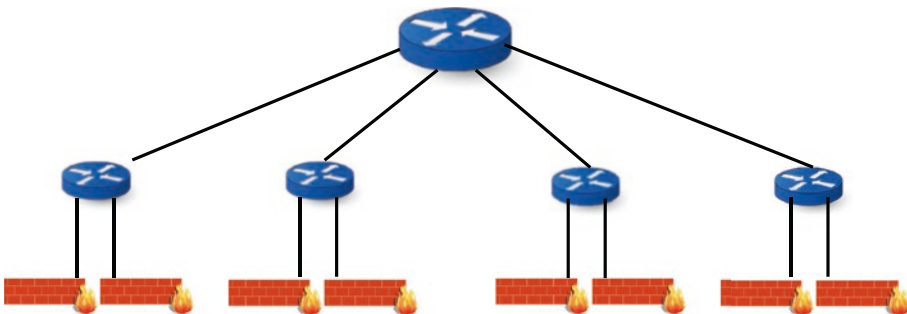


Figure 5-7. Centralized Network Topology

In some larger DCs, the LAN is configured with different topologies, while the SAN is still configured in the centralized topology. The reasons are that the cost of SAN switch ports is much higher than that of LAN ports and port utilization is an important cost factor.

Zoned Network Topology

A zoned topology consists of distributed network switches (see Figure 5-8). If the DC is organized into rows of racks, the network switches can be distributed in the middle of row (MoR) or at the end of row (EoR). The Data Center Standard organization recommends the topology for larger DCs and provides an easily repeatable and highly scalable architecture. Networks with zoned topology have the highest switch and port utilization level while causing lower cabling costs.

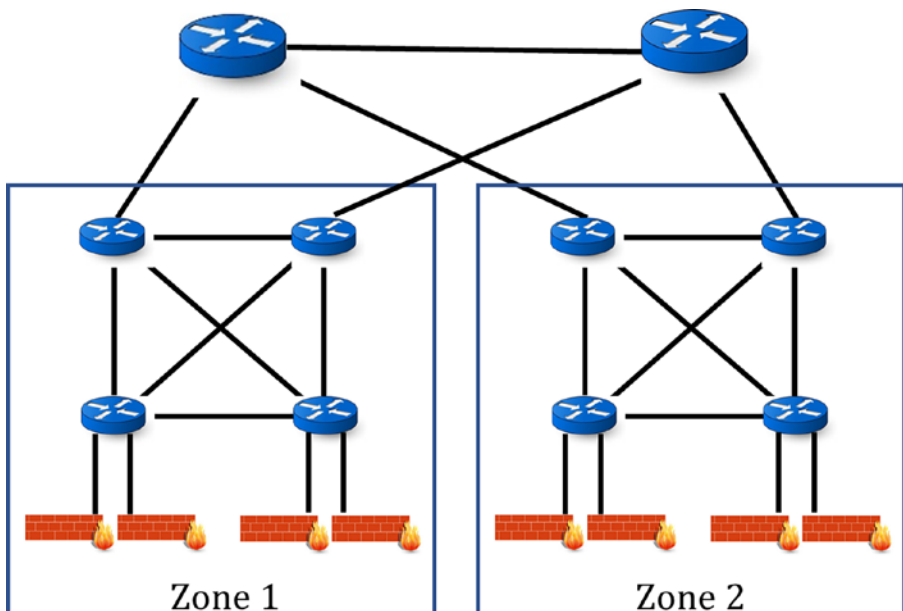


Figure 5-8. Zoned Network Topology

EoR switching often has performance advantages. It provides port-to-port switching with low latency between two racks that are needed to exchange large volumes of data. It may have disadvantages for cabling longer distances in comparison to a MoR configuration.

Top-of-Rack Topology

Top-of-rack (ToR) topology is a suggested choice for server configurations with very high rack density (see Figure 5-9). In this configuration, two or more switches are placed at the top of the rack. Each rack server is connected to both ToR switches. All ToR switches have at least two links to the higher-level networking switches. ToR configurations provide simpler cable management than other configurations and utilize cable more efficiently.

The disadvantages of ToR configurations are the higher costs of switches and the lower utilization of ports. The management of ToR configurations can be challenging in very large DC deployments.

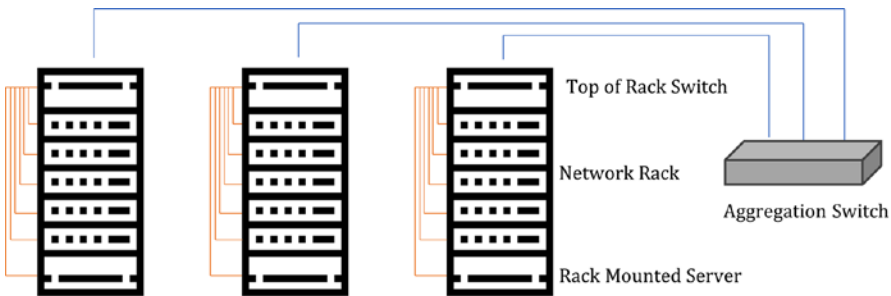


Figure 5-9. *Top-of-Rack Network Topology (Icons from Flaticon, 2023)*

Mesh Network Topology

In a mesh network topology, each element is connected to each other (see Figure 5-10). The architecture, also called network fabric, provides predictable capacity and data transport with low latency. These features

make it well suited for universal cloud services. Mesh networks offer a high level of redundancy and are often less expensive than traditional centralized switching equipment.

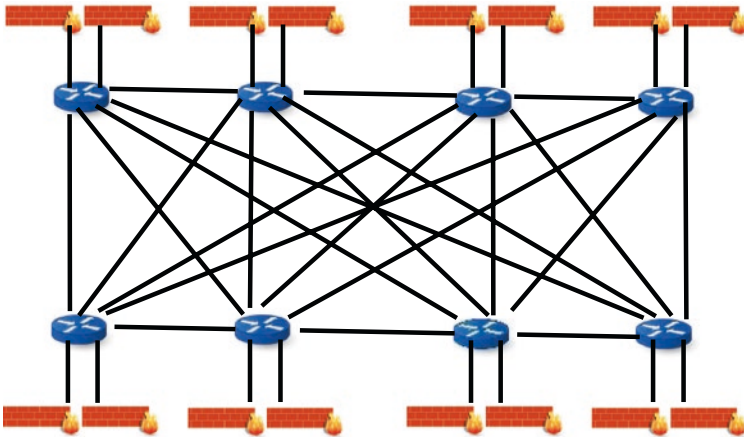


Figure 5-10. Mesh Network Topology

Multi-tier Network Topology

Most modern large DCs use a multi-tier network topology to deploy their networks for their campus-style DCs. Certain applications are often clustered in rows of racks, such as database servers, web services, and high-performance computer systems. Cross-connection switches connect cables, subsystems, and equipment to the racks. A cross-connection provides excellent cable management and design flexibility to support future growth. This topology offers operational advantages, as all connections are managed from one location. The major disadvantage is higher implementation costs due to increased cabling requirements.

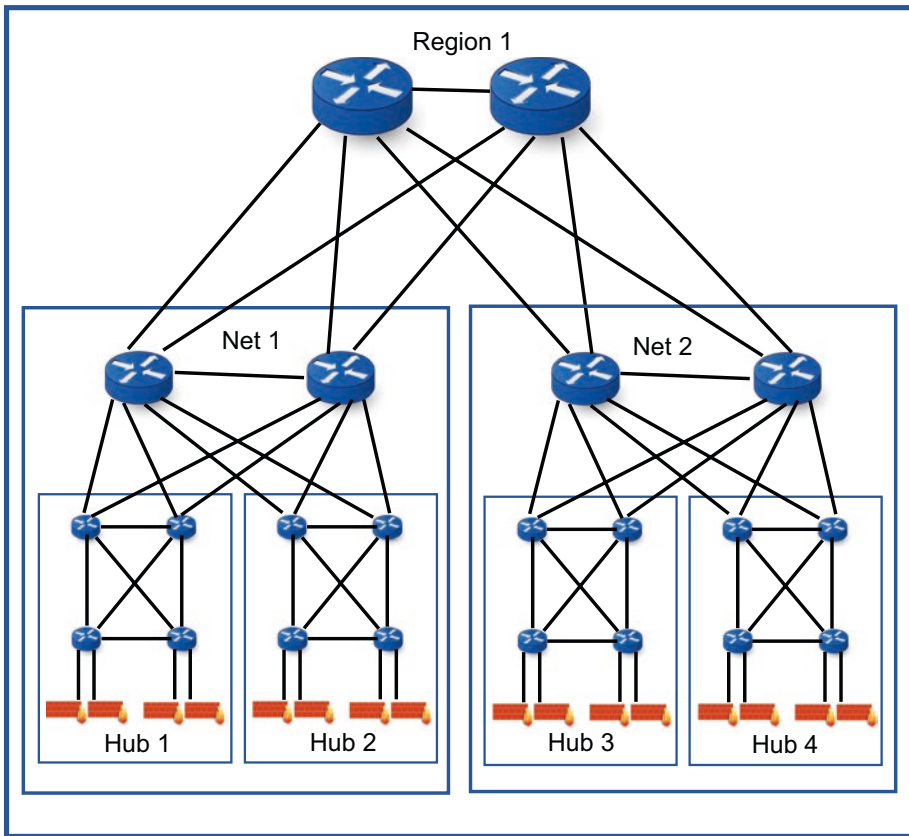


Figure 5-11. Tiered Network Topology

Multi-tier configurations allow multilevel network architectures. Figure 5-11 shows three layers: hub, net, and region. To better understand these three layers, please check Table 5-1.

Table 5-1.

LAYER	DEFINITION
HUB	A single room within a DC
NET	A DC building with several DC rooms
REGION	Several DCs across different locations

Software-Defined Networks

The trend toward open standards, the intent to reduce proprietary solutions, and tendencies toward cloud architectures led to the development of software-defined networking (SDN) technologies. SDN technology provides a network architecture that enables dynamic and programmatically efficient network configuration. With SDN, the static architecture of DC networks is transformed into a dynamic, configurable network design. In SDN networks, the transport of data packages is separated from the routing process of data packages. Data transport is handled in the SDN data plane, also called the data layer, while routing and management are performed in the control plane, also called the control layer. The control plane is the brain in an SDN architecture and contains one or several SDN controllers.

The SDN concept was first implemented with the OpenFlow protocol that was developed at Stanford University in the USA. Several technology companies, including Nicira, Google, NEC, and HP, developed OpenFlow-based solutions. In 2011, the Open Networking Foundation (ONF) was founded to promote SDN, OpenFlow, and drive its standardization. Many organizations have changed their network architecture to SDN by replacing their own, feature-rich, intelligent switches with switches using the simpler SDN architecture. In SDN networks, switches (hardware) are mainly used to transport data packages, while software performs the data routing and the network configuration. SDN network can be deployed with inexpensive, commoditized white label switches.

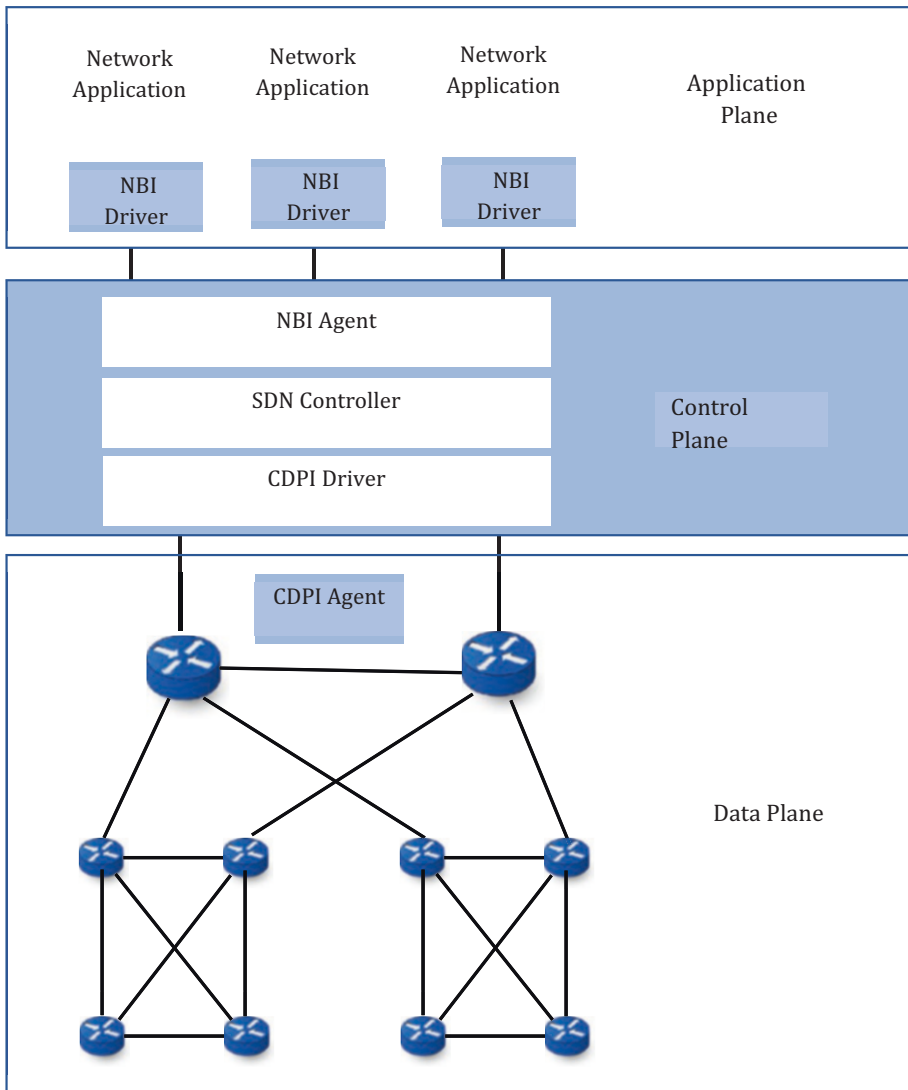


Figure 5-12. SDN Architecture

SDN Architecture

As shown in Figure 5-12, an SDN architecture contains the following components:

SDN Application

An SDN application is software that provides voice or data-related services to an organization. Examples are network virtualization, network monitoring, data flow balancing, virtual private networks, and intrusion detection. SDN applications communicate directly with the northbound interface (NBI) of the SDN control plane through its NBI driver. SDN applications can aggregate services and may offer interfaces to higher-level SDN applications.

SDN Control Plane

NBIs are interfaces between SDN applications (NBI driver) and SDN controllers (NBI agent) and hereby provide an abstraction layer of the underlying network data plane. The application interfaces are standardized by ONF (Open Networking Foundation), ensuring that network applications can run on other SDN platforms. This ensures interoperable and vendor-neutral communication services. The southbound interface of the SDN control plane is called the “control-data-plane interface” (CDPI). CDPI provides the SDN control plane access to the SDN data plane and its physical network resources. The main function of the data plane is to transport data. The standardization of CDPI by ONF ensures that the control plane is interoperable with network elements from any vendor with CDPI interfaces that comply with the ONF standard definitions.

SDN Data Plane

All physical network resources are in the data plane. The network elements need to run CDPI agent software to allow the SDN control plane to perform its functions. CDPI agents provide a logical representation of the physical network resources. The data network elements provide the key function of transporting the data, mainly forwarding data packages.

The SDN data plane provides a logical network device interface to the SDN control plane, including control over the network elements. The SDN interface is provided through CDPI agents that run on the network components. The logical representation of the network element may provide all or a subset of the physical resources.

An SDN data plane may be contained in a single (physical) network element – an integrated physical combination of communications resources, managed as a unit.

5.3 Network Resiliency

Network Fault Management

Network fault management is divided into three categories of accountability in case of an outage. The telecommunications industry defines the categories as follows:

- Product-attributable service outages
- Customer or service-provider-attributable outages
- External attributable outages

Product-attributable service outages are failures in software or hardware, outages attributed to customers or service providers, external outages like natural disasters, or some malicious acts. They are mainly offset by a system design, hardware, software, components, or other parts of the system. The design of the system makes scheduled outages a necessity. Customer or service provider-attributable outages are mainly offset by errors in the procedure, the environment of the office, etc. The outages caused by natural disasters or third parties are known as external attributable disasters. Examples of natural disasters are floods, earthquakes, and severe storms. The third-party-caused outages may or may not be related to the customer or suppliers.

5.4 Network Provisioning and Administration

Network Provisioning

Provisioning is the process of putting resources into use and making them available for a specific purpose. Network provision can benefit organizations to deploy networks with greater efficiency and security by reducing time for the setup and configuration of DC networks.

Due to the number and variety of network elements and the changing IT-related demands of organizations, deploying and managing larger DC networks can be a complex task. Automated provisioning will reduce the administrative work to create and deploy policies, add network addresses, configure, and deploy voice and data network services, and reduce errors in the provisioning process. Using automated tools for networking provision reduces the administrative staff for managing the DC network and records the process itself. The recorded data can be used to identify network problems and to perform network system audits.

DC network provisioning needs to include the following:

- **User provisioning**

An important measurement to improve cybersecurity is identity management that monitors access rights and authorization privileges. The provision process ensures that every IT user is classified in types like IT security staff, IT staff, employees, vendors, and contractors. For each type of user, access rights are defined to objects and group of objects. This is also called a role-based access control system. For example, a vendor may have reading rights to the procurement system of an

organization but not to data from the financial and HR department.

- **Network provisioning**

Network connectivity must be provided to all computer, storage, and network elements of the DC, plus all its internal and external users. This requires network provisioning to all the elements in the form of a physical and logical network connection.

- **Service provisioning**

In general, service provision ensures that DC services are deployed. An example of service provisioning: an organization wants to move an application to a cloud service provider (CSP). Before this process can be performed, the CSP must have provisioned several resources, including a server system, sufficient storage space, and software to load the application and data and send the output. This will be performed in a service provisioning process at the CSP, nowadays a highly automated process that can be initiated by a CSP client.

Network Administration

The role of a network administrator is to deploy, manage, monitor, and secure the network of an organization.

Monitoring

Network monitoring is essential to ensure the smooth operation of the network. The monitoring includes surveilling all network-connected devices, ensuring the overall health, and recognizing unusual traffic patterns. Unusual situations can be detected through monitoring, such as network connection issues, excessive bandwidth consumption, and abnormal activities. The network system administrator will act upon such events and take preventive and remedial actions to secure network quality and integrity.

Network management includes a list of tasks that a network administrator needs to perform:

- Network planning
- Network implementation
- Network configurations
- Data traffic management
- Software and firmware updates for all network elements
- Implementing security procedures and protocols
- Detecting vulnerabilities
- Change management of DC network capacity
- Network security

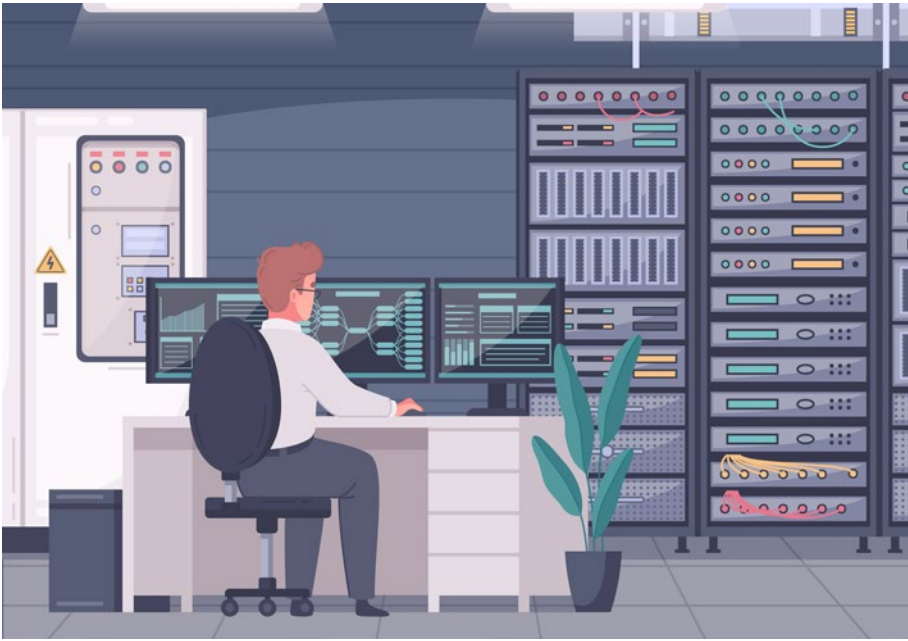


Figure 5-13. *Network Management Console (Macrovector, 2023)*

Network administration needs to ensure that the data network of an organization runs efficiently and is secure (see Figure 5-13).

Key tasks of a network administrator:

- Maintains a resilient, high-quality network
- Plans network capacity
- Ensures seamless network access
- Uses network management tools to administer the network
- Documents and tracks changes
- Risks assessments
- Plans risk mitigation strategies

- Prevents activities that could compromise the network integrity
- Prevents network security breaches
- Networks administration key areas

Network Administration Areas

Security Management

The network administrator needs to ensure that connected devices are authorized and that users are authentic and authorized for their specific activities. Processes must be put in place to enforce disciplined access for all authorized users. All required security software must be implemented and kept up to date, including tools for firewall management, threat management, intrusion detection, and virus scanners. Intrusion attempts must be detected, any suspicious activity must be investigated, and malicious activities must be prevented.

Account Management

The role of the system administrator is to collect network usage information to track and optimize network usage and use the data for billing purposes.

Fault Management

Network administrators monitor all network activities. Data transmission faults caused by defects in network elements or cables are identified, and remediation processes are initiated. Most network elements use the standards of the Simple Network Management Protocol (SNMP) and send SNMP alarms if network errors occur that trigger a notification to the network monitoring system.

Configuration Management

A configuration management system supports the network administrator to keep an up-to-date inventory of all network components and their configuration, including switches, firewalls, hubs, and routers. Configuration management is essential to streamline, track, and manage all configuration changes.

Performance Management

Another important part of the network management system is performance management. The system collects various metrics and analytical data to assess network performance, including transmission speeds, response times, packet loss, and utilization of links and ports.

5.5 Resilient Network for IT Data Center

As many organizations use IT services from a cloud service provider, they rely heavily on the Internet to function. Using CSPs offers several benefits, including access to efficient, scalable, and cost-effective IT resources. However, this comes at the price of security risks for IT services. Compared to the usage of a well-secured and managed private network, the Internet is prone to data breaches, hacking, and unauthorized access.

To ensure the safe and secure use of IT data and public clouds, businesses must implement robust security measures, such as encryption, authentication, and access controls. Organizations need to select trustworthy cloud providers with a proven track record of security and compliance.

Security Threats

Network architecture is faced with a list of threats caused by changing environments and attackers trying to find and exploit vulnerabilities. These vulnerabilities can exist in a broad number of areas, including devices, data, applications, users, and locations. Therefore, many network security management tools and applications are in use today that address individual threats, exploits, and regulatory noncompliance.

When just a few minutes of downtime can cause widespread disruption and massive damage to an organization's bottom line and reputation, it is essential that these protection measures are in place.

Physical Network Security

Physical security controls are designed to prevent unauthorized personnel from gaining physical access to network components such as routers, cabling cupboards, and so on. Controlled access, such as locks, biometric authentication, and other devices, is essential to any organization operating a data center.

Technical Network Security

Technical security controls protect data stored on the Web or in transit across, into, or out of the network. Protection is twofold; it needs to protect data and systems from unauthorized personnel, and it also needs to protect against malicious activities from employees.

Administrative Network Security

Administrative security controls consist of security policies and processes that control user behavior, including how users are authenticated, their level of access, and how IT staff members implement changes to the infrastructure.

5.6 Resilient Network Architecture for IT Data Centers

A well-designed IT data center network architecture must be developed and constantly updated to provide a highly efficient and secure infrastructure. The architecture should be designed to accommodate constant changes that are driven by the replacement of hardware components such as server, storage, and network systems; changes in the version of operating systems, virtualization software, and network management applications; and, finally, the constantly changing demand from organizations toward the resources of a data center or a cloud of data centers.

Resiliency is a feature that can be achieved by several contributing factors:

- Reliable hardware components for server, storage, and network components.
- Redundancy on each level. An example: a disk uses RAID format that provides redundancy on a device level, a file is copied to another RAID system in the same data center, to another RAID in another data center, and to another RAID system into a data center in another country, and so on.
- Virtualization software and load balancing monitor all processes. When a process encounters a problem, the current process is stopped, and the process will be moved to another system in the same data center or in another location.
- Tight security on all levels from physical access control of the data center to the components of a data center to all forms of cybercrime prevention software.

- Active IT management of all the hardware components, systems, and processes through state-of-the-art IT operation software and a well-trained data center operation team.

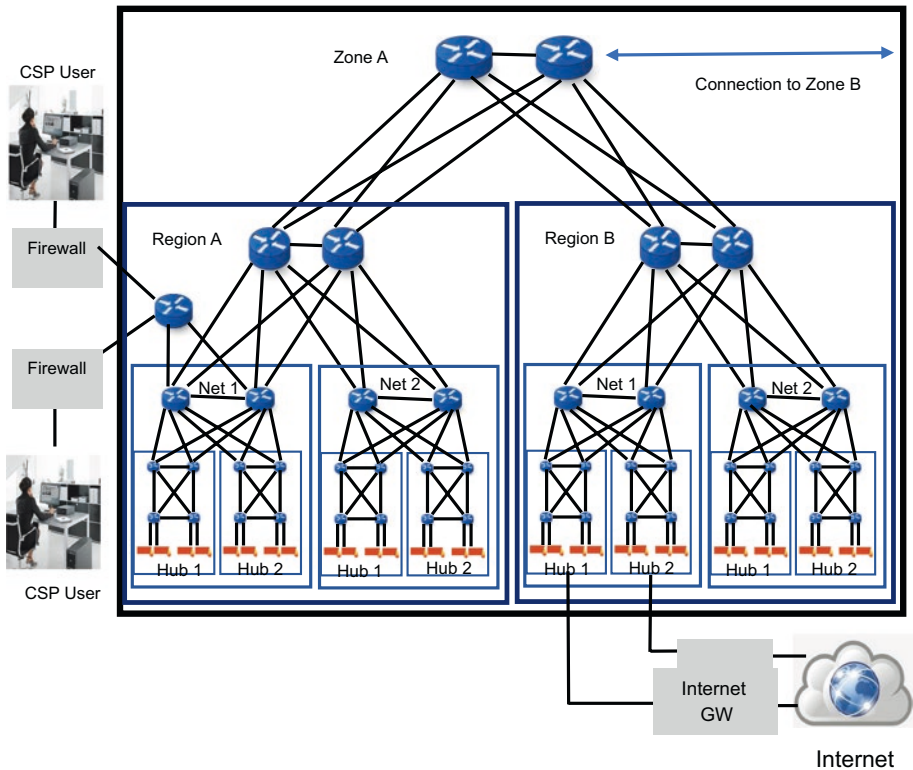


Figure 5-14. Tier 4 Data Center Network Architecture

Figure 5-14 represents the network topology of a four-level tier architecture that offers the highest level of resiliency. Table 5-2 shows the domain and territory correspondent to each tier.

Table 5-2. *Tier, Domain, and Territory*

Tier	Domain	Territory
1	Hub	Data center
2	Net	City or subregion
3	Region	Large region or country
4	Zone	Continent

5.7 Summary

This chapter provided a comprehensive exploration of data center networking and resilient network architectures. It commenced by introducing the fundamental concept of data center networking, detailing its hardware components, including cables, switches, and routers.

A significant shift toward software-defined networking (SDN) was discussed, emphasizing its impact on modern network management. The chapter delved into network resiliency, categorizing network fault management, and outlined the importance of network provisioning and administration.

Network administration tasks, such as monitoring, security management, and performance management, were explained, leading to a discussion on the critical elements of a resilient network architecture. The tiered architecture was presented as a blueprint for achieving network resiliency.

In conclusion, this chapter equips readers with a profound understanding of data center networking, its ever-evolving landscape, and the essential elements of a robust and resilient network architecture. It empowers individuals to navigate the intricacies of this vital technological domain, enabling the creation of network infrastructures that uphold the digital age.

In the next chapter...

- Backups
 - Disaster recovery planning
-

CHAPTER 6

Backup

Data backup is the process of copying data to another medium with the objective of retrieving the data in the case of an event that caused the loss of data on the primary medium. The process of restoring the data from the secondary medium is called recovery. Data backup requires the copying and archiving of computer data to make it accessible in case of data corruption or deletion.

Data backup is an important element in making the IT environment resilient and essential to any sensible disaster recovery plan. Today, data backup is often performed by using cloud storage. Therefore, data can be archived on a local system's hard drive, an external storage system, or by using cloud storage.

6.1 Evolution of Backup Systems

The history of computer backup systems started in the early days of computing when backup systems were not sophisticated and were using methods such as creating duplicate copies of data on magnetic tapes or other storage devices (see Figure 6-1).



Figure 6-1. *Classic Tape Drive for IBM Mainframe Computer (NationalArchivesCatalog, 1969)*

As data center applications and the associated data became more important to organizations, the budget was increased to purchase more advanced backup systems. In the 1980s, tape backup systems became popular, and backup software was widely available. This allowed an automatic approach to creating backups, drastically improving recovery times.

Networked storage systems started to emerge in the 1990s. These disk-based backup systems slowly started to replace the classic tape drives. The use of disks instead of tapes provided a more efficient backup process and, even more important, a much faster recovery process.

Virtual tape libraries (VTLs) started to emerge in the early 2000s. VTLs simulated the process of tape drives at a much higher speed and allowed existing backup software applications to continue with their existing tape backup processes (see Figure 6-2).



Figure 6-2. *Storage Tape Library (Cory Doctorow, 2008)*

Storage experts realized that in larger organizations, the same files were stored many times. A hundred employees could have stored company brochures or PowerPoint presentations. This recognition led to the development of intelligent storage software that can recognize such duplication and applies mechanisms to reduce the multiplications of the same data. This technique is called deduplication and helps organizations to remove a significant amount of redundant data.

Cloud-based backup services started to emerge in the later part of the 2000s as a new viable backup option. Cloud-based backup solutions allow backup data in a remote data center. Before this option was available, there was a threat that all the data of an organization could be destroyed in certain events such as flooding or fire, during which disk drives and the backup tapes could be destroyed. Cloud backup services provide more efficient backup management, automatically creating and storing backups remotely.

Backup systems are continuing to evolve. Today many organizations apply a hybrid approach to backup using cloud backup services and on-premises backup solutions.

6.2 Today's Backup Systems

Most organizations today use a combination of on-premises backup solutions and cloud-based backup services.

Backup in Private Data Centers or Private Clouds

The function of the backup system is to protect an organization against data loss caused by a variety of factors, including human error, natural disasters, hardware failure, software corruption, and cybercrime. Most IT organizations use a combination of tape drives or disk-based storage, offsite storage facilities, and backup management software. The backup process is typically automated to ensure that data is regularly backed up and stored in a secure location. In case of a data loss event, the private DC backup system must ensure that critical data can be restored quickly and easily. At times, the backup system must be able to recover data from a specific point in time. Another important feature is the ability to restore individual files, folders, and applications.

Cloud-Based Backup Systems

Unlike private data center backup solutions, cloud backup services are designed to provide backup solutions for cloud-based applications and data. While the fundamental functions of a backup system to back up and restore data efficiently are the same, there are fundamental differences. Cloud backup systems are more complex due to the distributed nature of cloud computing. Compared to a private data center, a different set

of backup tools is used for cloud backup solutions. Cloud-based backup solutions are more distributed in nature and require backup data to be stored in multiple locations, often across multiple data centers. This ensures a much higher level of data reliability, and data losses can therefore be entirely avoided.

Cloud-based backup solutions not only protect against data loss but also ensure the availability and reliability of cloud-based applications and services. Cloud-based backup systems can quickly recover from a wide range of data loss events.

In summary, backup is a critical function for a data center of any size and location. A well-designed backup solution is essential for ensuring the availability and reliability of critical data and applications. The specific approach to backup depends on many factors such as the size of the data center, if it's a private or a public cloud, and to which degree an organization must protect its data. Not all the data in an organization is critical, but data that is critical must be recoverable in any circumstances.

Backup technologies have become more sophisticated thanks to the latest features such as snapshotting, replication, and instant recovery. The rise of big data is another important event that has led to backup technologies that are specifically designed for handling large amounts of data.

6.3 Types of Backup Methods

Data backup is a critical component of any IT infrastructure, and choosing the right backup method is essential to ensure data integrity and efficient recovery in case of data loss. There are several types of backup methods, each with its advantages and disadvantages. Understanding these methods is crucial for designing a backup strategy that aligns with an organization's specific needs. Table 6-1 explores the most common types of backup methods.

Table 6-1. *Types of Backup Methods*

Type of backup	Description	Advantages	Disadvantages
Full backup	Copies all data in a selected location	Complete snapshot, easy restoration with the latest full backup	Consumes large storage capacity and is a slow process; this method may not be practical to backup large datasets frequently
Incremental backup	Captures changes since the last backup	Efficient storage and backup duration, keeps backups up to date	Restoration can be complex, requires correct order of full backup and subsequent incrementals
Differential backup	Captures changes since the last full backup. No reliance on previous differentials	Simplifies restoration compared to incrementals. Requires the latest full and recent differential for data restoration	Over time, can become larger and consume more storage than incrementals
Synthetic full backup	Combines full backup with incrementals to create a new full backup without copying all data	Time and space-efficient, provides full backup benefits with fewer resources	Implementation may need specialized backup software or appliances

(continued)

Table 6-1. *(continued)*

Type of backup	Description	Advantages	Disadvantages
Mirror backup	Creates an exact copy of source data on another storage medium or location	Real-time redundancy, quick recovery from hardware failures or corruption	Storage-intensive, requires twice the storage space, may not protect against accidental deletions or changes
Snapshot backup	Captures data state at a specific time without a full copy, recording differences from the previous snapshot	Efficient storage, useful for multiple recovery points	May not be suitable for long-term archiving, relies on underlying storage tech, may not be application-aware

Choosing the right backup method depends on factors like data volume, recovery time objectives, and available storage resources. Many organizations use a combination of these backup methods to create a comprehensive data protection strategy that addresses their unique needs.

6.4 Disaster Recovery Planning

Disaster recovery planning is a critical aspect of IT management and organizational resilience. It involves creating a strategy and set of procedures to ensure that critical systems and data can be rapidly recovered in the event of a disaster or unforeseen incident. In today's technology-driven world, where data is often the lifeblood of businesses, having a well-thought-out disaster recovery plan is not a luxury but a necessity.

Understanding Disaster Recovery

Disasters come in many forms, from natural events like hurricanes, earthquakes, and floods to human-made incidents such as cyberattacks, data breaches, or hardware failures. Regardless of the cause, the impact of a disaster on an organization can be severe. Without proper planning and preparation, the consequences can include data loss, extended downtime, financial losses, and damage to reputation.

The Disaster Recovery Process

A robust disaster recovery plan typically follows a well-defined process, which can be summarized in several key steps:

1. **Risk assessment:** The first step in disaster recovery planning is to conduct a thorough risk assessment. This involves identifying potential risks and vulnerabilities that could disrupt business operations. Risks can vary widely depending on your geographic location, industry, and the nature of your IT infrastructure.
2. **Business impact analysis:** After identifying risks, it's essential to assess the potential impact of a disaster on your organization. What systems, applications, and data are critical to your business operations? Understanding these dependencies is crucial for prioritizing recovery efforts.
3. **Strategy development:** Once you've assessed risks and their impact, you can develop a disaster recovery strategy. This strategy should outline the methods and technologies you'll use to mitigate

risks and recover from disasters. It might include offsite data backups, redundant hardware, failover systems, and cloud-based solutions.

4. **Plan documentation:** Documenting your disaster recovery plan is essential. This document should detail the roles and responsibilities of team members, the step-by-step recovery procedures, contact information, and any external resources you may need during a recovery effort.
5. **Testing and training:** A disaster recovery plan is only effective if it's regularly tested and the team knows how to execute it. Conducting tabletop exercises and drills helps identify weaknesses and areas for improvement.
6. **Execution:** In the event of a disaster, your team should be ready to execute the plan immediately. This may involve relocating to a backup site, restoring data from backups, and implementing failover systems.
7. **Monitoring and review:** Disaster recovery planning is an ongoing process. It's essential to continually monitor and review your plan, making updates as your organization evolves or new risks emerge.

Key Considerations in Disaster Recovery Planning

While the specific details of a disaster recovery plan will vary from one organization to another, several key considerations apply universally:

1. **Recovery time objectives (RTOs) and recovery point objectives (RPOs):** RTOs and RPOs define how quickly you need to recover systems and data after a disaster. These metrics help guide your planning efforts and determine the appropriate technologies and strategies.
2. **Data backup and storage:** Implementing robust data backup and storage solutions is a cornerstone of disaster recovery. Regularly backing up critical data to offsite or cloud locations ensures that you have a copy to restore in case of data loss.
3. **Redundancy and failover:** Redundancy involves duplicating critical systems or components to eliminate single points of failure. Failover systems automatically take over when primary systems fail, minimizing downtime.
4. **Communication plans:** Effective communication is crucial during a disaster. Establish clear communication plans, including how to notify employees, customers, and stakeholders about the situation and recovery efforts.
5. **Resource identification:** Identify external resources you may need during a disaster, such as alternate workspaces, hardware suppliers, or data recovery specialists.
6. **Security considerations:** Ensure that your disaster recovery plan includes cybersecurity measures to protect data during recovery efforts. This is especially important in the case of cyberattacks or data breaches.

7. **Compliance and regulations:** Depending on your industry, you may need to comply with specific regulations related to data protection and disaster recovery. Ensure your plan aligns with these requirements.

The Role of Technology in Disaster Recovery

Technology plays a pivotal role in disaster recovery planning and execution. Advancements in cloud computing, virtualization, and data replication have transformed the way organizations approach disaster recovery. Cloud-based disaster recovery solutions, for example, offer flexibility, scalability, and cost-effectiveness, making them an attractive option for businesses of all sizes.

6.5 Summary

In this chapter, the critical aspects of data backup and disaster recovery planning were explored. Beginning with an introduction to data backup, its role in safeguarding data and enabling recovery in the face of data loss was emphasized. The evolution of backup systems traces from early magnetic tape methods to the emergence of cloud-based solutions in the 2000s.

Today's backup systems involve a combination of on-premises and cloud-based approaches, catering to various organizational needs. The types of backup methods, including full backup, incremental backup, differential backup, synthetic full backup, mirror backup, and snapshot backup were examined, outlining their advantages and disadvantages.

This chapter concluded by emphasizing the significance of disaster recovery planning. The importance of understanding disasters, conducting risk assessments, defining recovery time and point objectives, and developing comprehensive strategies was discussed. Key considerations

such as data backup, redundancy, communication plans, resource identification, security, and compliance were highlighted. Technology's pivotal role in disaster recovery, including cloud computing and data replication, is also acknowledged.

This chapter equips organizations with a thorough understanding of data backup and disaster recovery essentials, emphasizing their role in modern IT environments for data protection and business continuity.

In the next chapter...

- Data center security
 - Vulnerabilities, attack motivations, prevention strategies, security tools, etc.
-

CHAPTER 7

Data Center Security and Resiliency

Data center security or IT security describes the methods to protect data centers including their computer, storage, and network systems from theft, damage, or disruptions to their hardware, software, or data. If computer systems are connected to the Internet, the term “cybersecurity” is often used.

The field of cybersecurity is becoming increasingly important for all companies and public agencies as so many business and public processes are powered by applications running in a private DC, also called a private cloud, or running on computers powered by public clouds through the Internet. Cybersecurity is an arms race between hackers and IT security staff. Hackers use increasingly sophisticated technologies to break into systems, while IT security staff try to protect the data centers by constantly managing and enhancing hardware and software protection resources (see Figure 7-1).



Figure 7-1. Computer Security (Freepik, 2023)

7.1 Vulnerabilities of Computer Systems

Fundamentally, vulnerabilities are caused by design, implementation, operation, or internal control weaknesses. The vulnerability can be exploited after somebody has researched, reverse-engineered, hunted, or exploited the attacked system often by using customized scripts or automated tools. Attacks on computer systems can be classified into one of the following categories.

Denial-of-Service Attack

Through denial-of-service attacks (DoS), computer systems become unavailable in general or to a specific group of users. Individual users can be attacked by deliberately entering a wrong password multiple times until the victim's account is locked by the authorization components of the computer system. DoS attacks are designed to overwhelm a computer system with service requests. A good firewall can detect and stop such an

attack when it comes from a single source by identifying its IP address. Distributed denial-of-service (DDoS) attacks that are coming from many points are more difficult to prevent. Such attacks can originate from so-called “zombie computers” or through other means, where other computer systems are misused to send traffic to the victim.

Phishing

Through this technique, the attacker tries to obtain information from its target by pretending to be a legitimate source like an official company or a bank. For example, the attacker sends an email or an instant message and asks for “confirmation” of information such as a home address, username, credit card information, or password. Alternatively, the attacker sends an invoice about goods or a service, and the target is asked to confirm this purchase. The user is asked to click the offered URL to the company’s website, which is a fake website of a legitimate company. For example, if the attacked person has an account at a bank with the URL 124bank.com, the attackers may use 124-bank.com instead. This link then goes to the phisher’s website from where private information is collected and subsequently used for malicious acts. The phisher’s website will look very similar to the website of a legitimate company, and the user may not recognize that data is entered on a system that is not legitimate.

Spoofing Attack

Through spoofing, the computer system of an attacker can disguise its identity and pretend to be a known, trusted source. There are several ways spoofing can apply:

- **Email**

Forging of an email address.

- **Phone call**

Forging a phone number.

- **Website**

Attacking a computer system forges a legitimate website.

- **IP address**

IP address of the attacking computer system is forged.

- **Domain Name System (DNS)**

Changing the link between a URL and the actual IP address.

- **Media Access Control (MAC)**

The MAC address of the attacking computer system is forged.

- **Biometric spoofing**

The attacking system provides a fake biometric sample posing as another user.

Spoofing is performed to gain access to information or resources that one is otherwise unauthorized to obtain.

Eavesdropping

Eavesdropping is a technique by which communications between computer systems are tapped and listened to. There are several techniques available. It is common that national service organizations, like the FBI, have software that can “listen” to the Internet communication of computers and can decode any information that is transmitted. Even if

computer systems are operated as closed systems without contact with the outside world, they can be eavesdropped upon by monitoring the electromagnetic transmissions submitted by the hardware.

Backdoor

A backdoor is a method to bypass a computer system's normal authentication and security controls for legitimate or nonlegitimate purposes. Backdoors to system may exist for different reasons:

- Designed by the programmer to check or test a certain part of an application;
- Created for legitimate access through lawful interception;
- Flaws in the design of an application;
- Flaws in the security architecture of the operating system.

Backdoors can be difficult to detect. The creation and usage of a backdoor requires someone with access to the application source code or intimate knowledge of the computer's operating system.

Direct-Access Attacks

If the access control system to a DC has weaknesses, an unauthorized user may be able to gain physical access to a computer system. With physical access, external devices can be added to copy confidential data. A person with physical access can compromise the system security by making operating system modifications, installing software such as worms or keyloggers, or using covert listening devices or wireless microphones. If the security of the operating system is not strong enough to prevent such attacks, the attacker may be able to modify the operating system by rebooting the computer

system from an external device that the person brings along, such as a CD-ROM drive, a USB disk, or a USB memory stick. Only trusted operating systems with encrypted data and security keys can prevent such attacks.

Privilege Escalation

A privilege escalation is defined as a scenario where an attacker with limited access rights to a system can elevate their privileges or access level without authorization.

Reverse Engineering

During reverse engineering, a man-made object is decomposed to reveal its software code and architecture or to extract general knowledge from the object. Through the knowledge of the internal design or by reverse-engineering the source code, an attacker can create modified versions of the software with attached malware.

Multivector and Polymorphic Attacks

Multivector attacks are initiated simultaneously from many attacking computing systems and can be combined with multiple types of attacks, often referred to as polymorphic attacks. Most cybersecurity systems may not be able to control such attacks.

Social Engineering

Social engineering targets the people using or administering computer systems. The goal is to convince the attacked person to disclose secrets such as passwords and card numbers or grant physical access to the targeted computer system. Social engineering involves exploiting people's trust and requires a high level of cognitive skills by the attacker.

Malware

Malware is malicious software that is installed on a computer to leak personal or system administrator information to the attacker to gain control of the system. Malware is often installed on a computer system through viruses that a user gets from visiting malicious websites.

Figure 7-2 shows computer security threats.

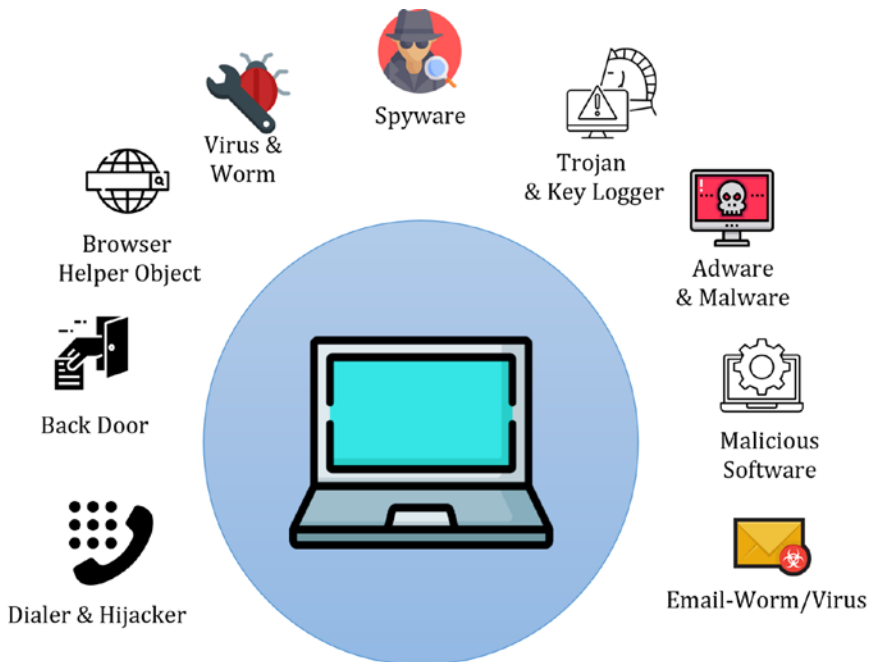


Figure 7-2. Computer Security Threats (Flaticon, 2023)

7.2 Motivations and Impact of Attacks

Impact of Security Breaches

It is often hard to assess the actual financial impact of an incident. The total damage will often go far beyond the actual breach. Even if only a small number of customers were affected by a malicious attack with limited financial damage, the attacked company must perform thorough security investigation of all client data after the attack was detected. This is mandated in many countries nowadays. The attacked company must notify all its clients about the security breach, which often comes with serious impact on the company's reputation and may even impact its market valuation.

Damages from computer attacks are constantly rising. It is estimated that in the United States alone, the financial implications of cybercrimes have risen from USD 1.4 billion in 2018 to USD 4.2 billion in 2020. This is threefold in only four years. Reasonable estimates of the financial cost of security breaches can motivate organizations to sufficiently invest in information security. Cybersecurity analysts have repeatedly reported that the amount a firm spends to protect its applications and data will generally be only a small fraction of the expected losses from cybercrime.

Attacker Motivation

As with physical security, the motivations for breaches of computer security vary between attackers. In Figure 7-3, you can see some of the sources of cybersecurity threats. There are several motivations for attackers for their actions:

- Thrill-seeking vandals;
- Criminals looking for financial gains;

- State-supported attackers for political espionage and cyber warfare;
- Industrial espionage to steal confidential information and trade secrets from companies;
- Terrorist organizations attack critical infrastructure such as the electricity grid;
- Hackers help organizations to find security vulnerabilities;
- Whistleblowers and political organizations expose information that these organizations don't want to disclose.



Figure 7-3. Sources of Cybersecurity Threats (Icons from Flaticon, 2023)

Anticipating the motivation of a potential attacker is an important part of the security assessment. Another important aspect is the “value” of the attacked objects for the attacker. A personal home computer has a relatively low “attack value,” while computers of companies and institutions are on a higher level, followed by companies dealing with financial transactions, such as banks. On the highest level are computer systems with access to critical infrastructure like an energy grid, an airlight control system, or classified military computer systems.

7.3 Security by Design

Applications and the underlying infrastructure, such as libraries and operating systems, need to be designed from the ground up to be secure. This is called “security by design.” This contrasts the classic approach of creating a system and adding security features after the application is completed. Important principles to implement “security by design”:

- **Least privilege**

This principle defines that each part of the system has only the privilege that it requires to operate its function.

- **Automated theorem**

This headline summarizes attempts to automate the validation of software codes. These methods can be successfully applied in the range of simple to moderate complex applications. But when a certain complexity threshold is exceeded, the system will have at least one assumption that cannot be proved within the system. For such software, validation cannot be performed, but code reviews and unit testing can be done to improve security.

- **In-depth defense**

This methodology creates multiple layers of security control for an application with the intent to provide redundancy if a system vulnerability is exploited by an attacker or if a security control fails. Controls are implemented in three domains:

1. **Technical**

Technical methods to protect the system include encryption, access control, and file-integrity checks.

2. **Physical**

Preventing physical access to the IT systems, such as fences, locks, or CCTV cameras.

3. **Administrative**

Execution of security protocols, such as password guidelines, hiring practices, and security audits.

- **Default secure settings**

This fundamental concept is to make systems “fail-secure,” which means that a secure system requires the legitimate administrator to make a deliberate, conscious, knowledgeable, and free decision to make it insecure.

- **Audit trails**

A system that permanently records and audits all activities on a system is called an audit trail tracking system. When a security breach occurs, the system records all the information about, the mechanism, and the extent of the breach. The audit trails are recorded remotely to prevent intruders from covering their activities.

Security Architecture

A security architecture is defined by Techopedia as “a unified security design that addresses the necessities and potential risks involved in a certain scenario or environment. It also specifies when and where to apply security controls. The design process is generally reproducible (...)”

The key attributes of security architecture are

- The relationship and dependencies of components;
- Controls based on risk assessment, good practices, finances, and legal matters;
- Control standards.

Practicing security architecture provides the right foundation to systematically address business, IT, and security concerns in an organization.

Security Infrastructure

All security systems are based on the three key processes: prevention, detection, and response.

Several processes can be implemented:

- Strong access controls and cryptography can protect the files and data of systems.
- Hardware- or software-based firewalls are commonly used as prevention to improve network security. Through packet filtering, firewalls can shield access to internal network services and block certain kinds of attacks.

- Attacks and intrusion of a network can be detected by an Intrusion Detection System (IDS). An IDS provides detailed information about the attacks and assists in post-attack forensics.
- The “Response” process consists of a range of means that can be executed when a system is attacked. This can range from simple updating or upgrading the security software in use, notification of authorities, counterattacks, to, as an extreme measure, the complete destruction and deletion of all data of a compromised system.

The strong increase in damages from cybercrimes has led to the situation today where companies incur more damages from cybercrime than from the classic crime of stealing assets. This suggests that relatively few organizations operate DCs with effective detection systems, and even fewer have organized response mechanisms in place. Most organizations just rely on simple firewall and detection systems, which is not sufficient. All security infrastructure elements must be properly implemented following the idea of the “theory of constraints,” which defines that a system is only as strong as its weakest link. In many companies, the weakest link is the administrative part due to lack of sufficient security protocols and processes, as well as the lack of recognition that all users of a company’s IT system must be security-aware and diligent in their operation with computer systems.

Vulnerability Assessment and Management

Assessing the vulnerability of a DC infrastructure is an important first step to improve cybersecurity. Figure 7-4 represents a security operation center.

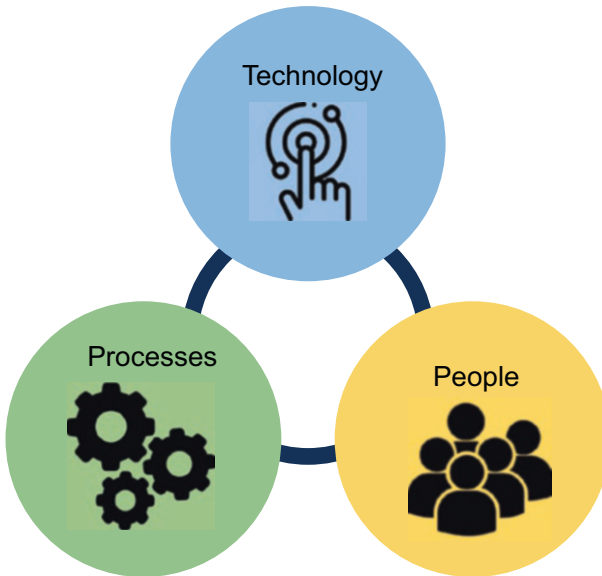


Figure 7-4. *Security Operation Center (Icons from Flaticon, 2023)*

Several software companies offer tools to assess a DC’s status by applying a vulnerability scanner. This scanner analyzes a computer system searching for vulnerabilities, such as open ports, insecure software configurations, and installed malware.

If vulnerabilities are detected, the appropriate mitigation will be identified, resulting in changes to hardware, software, or processes.

This assessment must be performed on a regular basis and with the latest version of the vulnerability scanner as the attackers constantly create new methods to perform their malicious activities. Many organizations use external organizations to audit their cybersecurity through independent, external experts, if companies do not have their own sufficient internal resources to conduct audits. In some sectors, external audits may be mandated by regulations or by contractual requirements.

Reducing Vulnerabilities

There are several means available to reduce DC vulnerabilities:

- **Making system access more difficult**

Through “two-factor authorization,” a system is protected by requesting something the user knows, like a password or a PIN, and something the person has, like a dongle or a cellphone. Access is only granted when data from both sources is correctly entered.

- **Increasing security awareness**

Social engineering attacks can only be prevented by noncomputer means, which is often difficult to impose. Training is important to mitigate this risk, but even in highly disciplined environments, like military organizations, it remains challenging to prevent social engineering.

- **Vigilant security culture**

Creating a more secure culture by instilling more nontrusting behavior and building resistance against social engineering through information and training about actual cases of cybersecurity breaches.

- **Keeping the system up to date**

As cybersecurity is an ongoing “arms race,” it must be mandated that all systems, particularly the security systems, are kept up to date by installing patches and updates quickly after they are available.

Hardware Protection

Computer hardware can be a source of insecurity when malicious components are added, or when items are modified during the computer system's manufacturing process.

On the other hand, hardware-based computer security also offers additional security layers compared to software-only computer security.

Hardware devices for enhanced security include the following:

- **Trusted platform modules**

A microprocessor system with an integrated cryptographic chip that controls access to the system is called a Trusted Platform Module (TPM). The application or operating system must provide the correct key to unlock access to the microprocessor (see Figure 7-5).

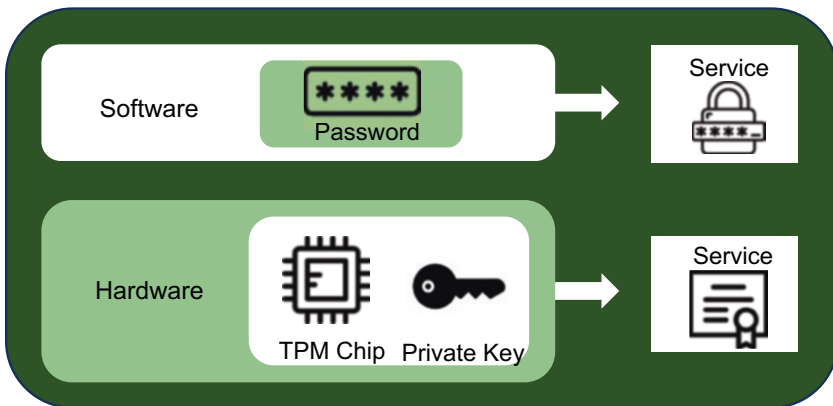


Figure 7-5. *Trusted Platform Module (Icons from Flaticon, 2023)*

- **Intrusion detection device**

A computer intrusion detection device detects when a computer chassis is opened, which can be implemented by a push-button switch.

The microprocessor firmware is designed to use the status from the intrusion detection device by alerting the administrator the moment the computer is rebooted.

- **Disabled USB ports**

Disabling USB ports of a computer system prevents attackers or staff from copying data from the computer system to a USB disk or memory stick, or copying malicious software from a USB device to the computer system.

- **Disconnecting unused devices**

Hereby all devices like built-in cameras, GPR systems, Bluetooth connectivity, or removal storage are disabled to close vulnerability windows of unused computer devices.

- **Drive locks**

Drives can be locked by using software tools that encrypt all the data stored on hard drives, making them inaccessible to thieves.

- **Mobile access control**

A new set of access control systems was developed with general proliferation of mobile phones. Mobile phones can connect directly to a computer system via protocols such as near-field communication (NFC) or Bluetooth. Security features of the phone, such as a thumb print and eye or voice recognition, can be used to unlock applications (see Figure 7-6). Other types of methods use the phone and its connection to a cellular network to enable applications with a QR code, or by sending a

PIN via SMS or message applications like WhatsApp. These mobile phone-based methods have become quite a common access protection tool for consumer applications that require higher levels of security such as mobile banking and mobile payments.

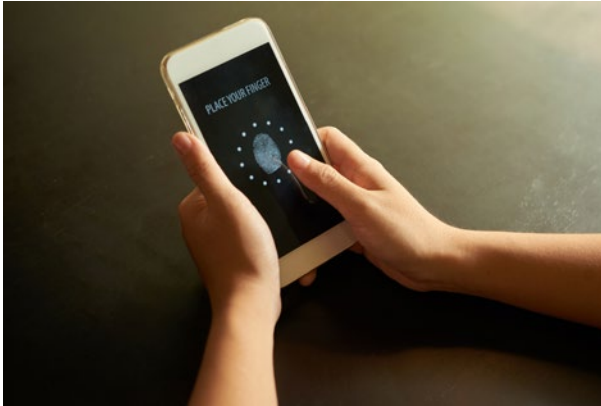


Figure 7-6. *Access Control by Mobile Phones (Freepik, 2023)*

- **USB dongles**

A simple way to use a dongle is with computer systems that are configured to require a specific dongle to unlock. Another method uses protected software components that require a key from a dongle to run. The dongle's key creates a "secure encrypted tunnel" between the software application and the key. Hacking or replicating the dongle is difficult. Dongles can also be used to protect the network access through cloud software or a virtual private network (VPN). Some dongles use fingerprint recognition to unlock the computer system.

Access Control Lists

An access control list (ACL) defines the permission associated with an object. A user can be defined individually or classified in groups such as end user, developer, or administrator. The ACL contains the access right of each user or user group for the specific file.

Another method using ACLs is role-based access control (RBAC) that restricts system access to authorized users (see Figure 7-7). It provides an access-control mechanism defined around roles and privileges. Key components of RBAC are role-permissions, user-role, and role-role relationships.

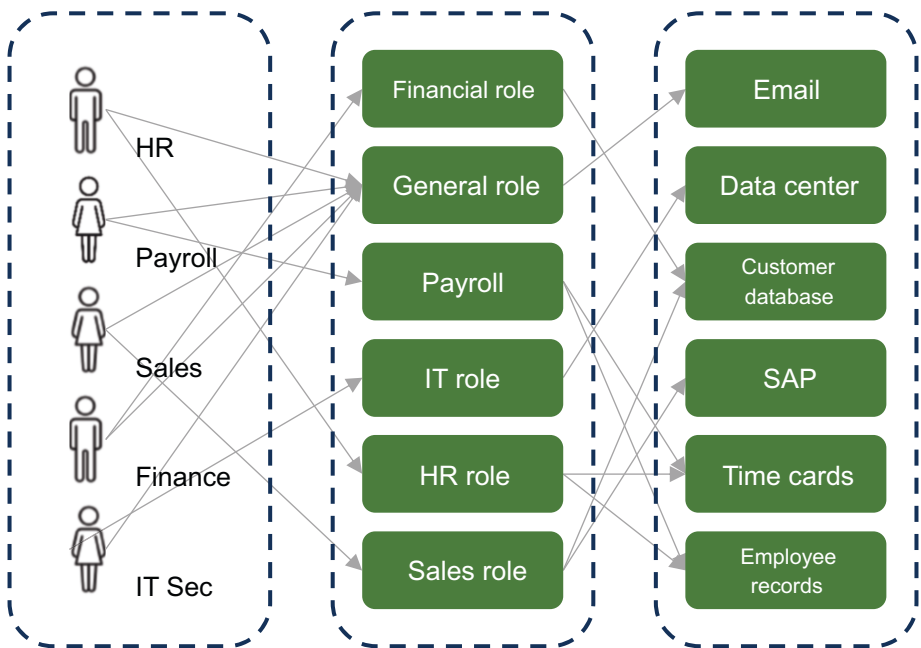


Figure 7-7. Example of an RBAC (Icons from Flaticon, 2023)

RBAC is commonly used in organizations with more than 500 IT users.

Security Tools

There is a set of software available to implement security on computer systems:

- Access control;
- Anti-keyloggers;
- Anti-malware;
- Anti-spyware;
- Anti-subversion software;
- Anti-tamper software;
- Anti-theft;
- Antivirus software;
- Cryptographic software;
- Computer-aided dispatch (CAD);
- Firewalls;
- Intrusion detection system (IDS);
- Intrusion prevention system (IPS);
- Log management software;
- Parental control;
- Records management;
- Sandbox.

Security Training

More than 80% of security incidents and breaches are caused by a set of common mistakes by human users, including the following:

- Easy-to-guess passwords;
- Passwords were written down at a place that is accessible;
- Sending emails containing sensitive data to unauthorized recipients;
- Attaching confidential files to emails and sending them to unauthorized recipients;
- Entering user IDs and passwords to fake websites;
- Not recognizing misleading URLs;
- Saving user IDs and passwords in browsers;
- Revealing user IDs and passwords to non-authorized personnel.

As IT systems are used by almost every person in an organization in some form, the approach to cybersecurity must change. IT security used to be the job of a few IT and security specialists in the IT organization. Cybersecurity is a task that needs to be owned by the entire organization. Hence, it is very important for every organization to mandate security awareness training at all levels.

Cyber Hygiene

Cyber hygiene follows the concept of using personal hygiene to prevent humans from being infected by bacteria or viruses. Cyber hygiene provides users of IT systems a layer of protection to reduce the risk of spreading a virus from one vulnerable computer system to the organization's system of network computers.

Just as personal hygiene, cyber hygiene imposes simple measures that are based on discipline and education. Guidelines are provided regarding the “dos” and “don’ts” of behaviors related to security, including password strength, password protection, regular password changes, using different passwords for different systems, vigilance when accessing the external website, and caution when unknown people ask for security-related data.

Incident Response

A security breach requires an organized approach to addressing and managing the aftermath of a computer security incident or compromise in the form of an incident report. The report aims to identify the attack’s structure and methods and detect the exploited weaknesses of the attacked system. Any form of intrusion or security breach must be determined and managed immediately to prevent the escalation to a more damaging event such as a data breach or system failure.

Computer security incident response has the objective to contain the incident, limit damage, and assist recovery to business as usual. Incident response planning gives organizations a framework to establish best practices to stop intrusions early before it escalates to actual damages. Plans for incident responses contain a list of written instructions that outline the organization’s processes and escalation that will be performed during a cyberattack. Without defined processes, responsibilities, and escalation procedures, an organization may not be able to detect an intrusion or compromise quickly enough. Stakeholders may know their roles immediately, slowing the organization’s response with the danger that an incident escalates into a severe security breach with all its potential damages.

Key components of a computer security incident response plan:

- **Preparation**

It is important to prepare all stakeholders for a security breach by training procedures for handling computer security incidents or compromises.

- **Detection and analysis**

Constantly monitoring all system activities with the ability to identify and investigate suspicious activities to confirm a security incident and subsequently prioritize the response, based on the impact of the incident.

- **Containment and restoration**

An incident must trigger fast responses. An affected system must be isolated quickly to prevent escalations and limit the attack's impact. After the infected computer system is isolated, the genesis of the incident must be analyzed, and the malware needs to be identified and removed. Security measures may also be applied to the actors involved in the attack. When the incident is resolved, the system can be restored to its healthy state.

- **Post-incident activity**

A thorough postmortem analysis of the incident requires a deep root cause analysis of the incident and the underlying weakness in the cybersecurity that was exposed. Identified weaknesses will lead to improvements in security methods and protocols.

Cybersecurity Planning

Cybersecurity planning involves several components:

- **Strategic planning**

Clear targets need to be set to develop a better awareness program. Assembling a team of skilled professionals is helpful to create a plan and keep it up to date.

- **Operative planning**

A good security culture can be established based on internal communication, management buy-in, security awareness, and a training program.

- **Implementation**

Four stages should be used to implement the information security culture. They are:

- Commitment of the management,
- Communication with organizational members,
- Courses for all members of the organization,
- Commitment of the employees.

- **Post-evaluation**

To assess the success of the planning and implementation and to identify unresolved areas of concern.

7.4 DC Resilience

DC Security

According to the Cambridge dictionary, resilient is defined as “able to improve quickly after being hurt or being ill.” In the context of a DC, resiliency can be defined as the ability of an IT organization to recover quickly from any form of failure, natural event, or malicious attack and continue the operations. The foundation for DC resilience is DC security. All the key building blocks that were described before must be applied:

- **Infrastructure**
 - Reliable hardware components for the server, storage, and networking;
 - Redundancy of resources;
 - Uninterruptable power supply;
 - No single point of failure;
 - Adequate backup facilities;
 - Secured locations for backup media;
 - Distributed and, if possible, federated resources;
 - Physical access protection of the DC.
- **Design and tools**
 - Security by design on all levels of hardware, operating systems, middleware, and applications;
 - Firewalls;
 - Virus detection and anti-malware software;
 - Intrusion detection;
 - Access control management.

- **Processes and policies**

- Training of all IT users about IT security and security awareness;
- Strong policies and enforcement of access control, including password management;
- Access control management;
- Disciplined approach for data backup and recovery.

If an organization is following the well-defined rules of DC security, a high level of security will be achieved, but it will never be 100%. There could be an unexpected event such as a major power outage, an earthquake, a physical breach of the DC access, or the antivirus software in use is unable to detect a new software virus. The concepts for DC resiliency must take such events into consideration, develop plans and implement them to recover from a disaster. This is often called a “disaster recovery plan.”

DC resiliency can be achieved using redundancy for its critical resources, including facilities, systems, and components. If an element fails or is disrupted, the redundant element takes over seamlessly and continues to provide the same service as before the event happened. The overall resiliency of an organization is built by processes including business continuity, emergency responses, and incident responses (see Figure 7-8). The goal of resiliency is to avoid downtime. Ideally, users of a resilient system do not recognize when a disruption has occurred.

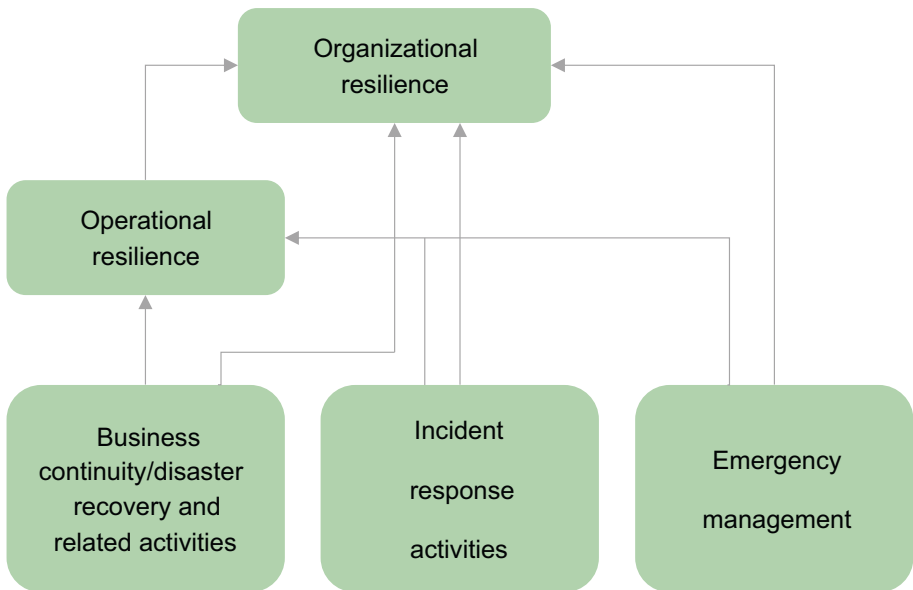


Figure 7-8. *Processes for DC Resiliency*

In the earlier days of computing, a DC had one computer called a “mainframe.” In the 1980s, more computers were deployed in DCs – often for different purposes. With the emergence and rise of workstations and later PCs, computing became more decentralized. Modern DCs, especially those from CSPs, have thousands or ten thousands of computer systems. This DC evolution from one to many computer systems has changed the way DC resiliency is built. If a DC has only one computer system, resiliency can only be achieved if this computer is fully redundant and does not have a single point of failure. If several computers are used in a DC, the emphasis on resiliency shifts from the single computer to the compute cluster. A single computer failure can be tolerated if the other computer systems can take over and if other threats such as power interruption or flooding can be prevented. In today’s DCs of CSPs, the focus has shifted further. A failure of one or several computer systems is not a critical event anymore if a higher-level load-balancing software is active that can

redistribute workloads to other active computer systems. As workloads can be easily distributed within a DC, a shift to another DC can be done as well. This offers an even higher level of resiliency: the usage of software virtualization and load balancer allows building a federated architecture comprising of several DCs. This structure provides the highest level of resilient compute services, so that even a natural disaster such as an earthquake will not cause a disruption.

Critical Services

Building a high level of redundancy is costly. Not all applications have the same importance for an organization. Therefore, it is prudent to classify the workloads in classes of importance such as

- Very critical,
- Critical,
- Standard,
- Low priority.

The cost of application downtime for an organization provides a good indicator to classify the criticality of workloads. The other crucial factor to consider is the impact of workload delays. A monthly inventory report can run a few hours later if a computer system is not available. The transaction of an equity broker must be concluded in a few milliseconds; hence, such computer systems used by exchanges are in the group of the most critical systems.

Achieving Data Center Resiliency

To achieve a high level of resilience without exhausting IT budgets, an organization must first assess its existing workloads and classify them. Afterward, it must determine what level of resiliency each workload

requires by considering business and technical factors. The cost of resiliency can be high because more resilience requires more resources and is therefore more expensive. Table 7-1 shows how the risk can be reduced by building up resiliency from a noncritical DC to a federated DC infrastructure.

Table 7-1. *Building Up Resiliency*

Features	Risk level
DC with no redundancy	High
DC with redundant components	Medium
DC with backup DC at the same location	Low
DC with multiple DC at different locations	Very low

Improving Resilience

To ensure the continuous operations of a DC, it is prudent to monitor its operating conditions. These include environmental conditions by monitoring temperature and humidity within the DC.

The DC should have additional monitoring equipment and software to continuously observe important data about the DC, such as server operations, application processing, data backup, and power levels. Monitoring the operational data from the DC can often indicate early if something is not working as it should. Early recognition of such situations, and corrective actions, can often prevent disruptive problems before they escalate into outages.

Redundancy is required not only for the computer and storage systems but also for networking and security. Using only one Internet service provider (ISP) creates a single point of failure; hence, it is suggested to have an Internet connection from two different ISPs.

When the monitoring systems detect anomalies or when a failure is detected, the alarm system must be activated, and the DC operation team must be alerted.

The DC operation team must frequently get training to be prepared for critical situations. Simulation of security issues such as hardware failures, power interruption, or cybersecurity attacks can identify any vulnerabilities that could cause a real incident. Maintaining the resiliency by constantly updating and enhancing security tools of a DC is an ongoing process and must be part of the daily data center operations and not only be part of an occasional exercise, like a fire drill.

7.5 Summary

This chapter delved into the crucial aspects of data center security and resiliency. It began by introducing the concept of data center security and its significance in safeguarding computer systems, storage, and network infrastructure against theft, damage, or disruptions. The term “cybersecurity” was highlighted, particularly when computer systems are connected to the Internet, emphasizing its growing importance in today’s digital landscape.

The chapter proceeded to explore various vulnerabilities that computer systems face, primarily attributed to design, implementation, operation, or internal control weaknesses. It categorized attacks into distinct types, including denial-of-service (DoS) attacks, phishing, spoofing, eavesdropping, backdoors, direct-access attacks, privilege escalation, reverse engineering, multivector and polymorphic attacks, social engineering, and malware.

The motivations behind these attacks were discussed, ranging from thrill-seeking vandals to state-supported attackers, industrial espionage, and hackers aiding organizations to find security vulnerabilities. The financial impact of security breaches was underscored, with an increasing trend in cybercrime damages over the years.

The concept of “security by design” was introduced, emphasizing the need for secure application and infrastructure design from the ground up. Key principles such as least privilege, automated theorem, in-depth defense, default secure settings, and audit trails are explained. Security architecture, including the relationship and dependencies of components and control standards, is also discussed, emphasizing its role in addressing business, IT, and security concerns.

Then, the security infrastructure was explored, including prevention, detection, and response mechanisms, with a focus on strong access controls, firewalls, intrusion detection systems, and response strategies. It highlighted the rising damages from cybercrimes and the need for effective detection and response mechanisms.

Vulnerability assessment and management were discussed as crucial steps in improving cybersecurity. The use of vulnerability scanners and external audits was highlighted as essential practices to identify and mitigate vulnerabilities effectively.

Following that, strategies to reduce vulnerabilities were addressed, including two-factor authentication, security awareness, cultivating a vigilant security culture, keeping systems up to date, and hardware protection. Various hardware-based security measures such as trusted platform modules, intrusion detection devices, disabled USB ports, and mobile access control were explained.

Access control lists (ACLs) and role-based access control (RBAC) were introduced as methods to manage permissions and restrict system access effectively.

The chapter further explored security tools available for implementing security on computer systems, such as access control, anti-keyloggers, anti-malware, firewalls, intrusion detection and prevention systems, and more.

Security training is emphasized as a critical component, with over 80% of security incidents resulting from common human errors. Cyber hygiene was introduced, focusing on simple measures and guidelines to reduce the risk of security breaches.

Incident response was discussed, highlighting the need for an organized approach to addressing and managing security incidents. The key components of a computer security incident response plan, including preparation, detection and analysis, containment and restoration, and post-incident activities, were explained.

This chapter concluded by delving into cybersecurity planning, covering strategic planning, operative planning, implementation, and post-evaluation. It stressed the importance of assessing existing workloads, classifying them based on criticality, and determining the appropriate level of resiliency required for each workload.

Data center resilience was introduced as the ability to recover quickly from failures, natural events, or malicious attacks while maintaining operations. The foundations of data center resilience, including infrastructure, design and tools, processes and policies, and redundancy, were discussed.

Also, the significance of improving resilience through continuous monitoring, redundancy in networking and security, alarm systems, regular training, and the constant enhancement of security tools was underscored.

In summary, this section provided a comprehensive exploration of data center security and resiliency, offering insights into vulnerabilities, attack motivations, prevention strategies, security tools, and the crucial role of preparedness in the modern digital landscape. It emphasized the importance of building a strong security culture and resilient infrastructure to protect against evolving cybersecurity threats.

In the next chapter...

- IT support services – from foundational aspects to modern evolutions
-

CHAPTER 8

IT Support Services

In general, IT services are provided by each organization that operates a data center (see Figure 8-1). These services can be provided by the internal IT organization of a company or by external companies such as cloud service providers (CSPs). CSPs offer IT services over the Internet, such as IaaS, SaaS, and PaaS that were described earlier in Chapter 1.



Figure 8-1. Help Desk Call Center (Freepik, 2023)

8.1 IT Help Desk

A Help Desk is a function that ensures that users of an IT system get support when they have issues they can't fix themselves. The services of an IT Help Desk therefore refer to the support provided by an IT team to end users of technology within an organization. Help Desk services are usually offered to IT users as a single point of contact to report IT issues, request IT services, and receive technical support. The function of the Help Desk is to resolve technical problems and provide guidance to operate IT functions efficiently.

IT Help Desk services include functions such as the following:

- **Login to systems and applications**

Users require to retrieve their password or login ID to access their system or applications.

- **Hardware support**

Users require support for hardware-related issues, including hardware diagnostics, replacing faulty components, and coordinating repairs.

- **Training and guidance**

The Help Desk team offers training and guidance on how to effectively use software applications, hardware components, and other IT systems.

- **Technical support**

Users require help to use their systems without being able to describe whether the encountered issues are caused by hardware, software, network, or any other IT-related problems.

- **Troubleshooting**

Users require a Help Desk support technician to investigate the cause of problems in their attempts to use systems and applications.

IT Help Desk services are a critical function in every IT organization as they ensure that end users use technology effectively and efficiently. An effective Help Desk with a short response time to user requests will reduce downtime of IT applications, increase user productivity, and improve end-user satisfaction. The person who responds to a service request from a client is typically called a customer service representative (CSR). A CSR is responsible for addressing the client's questions, concerns, or issues with the product or service that the company provides.

Options to Contact a Help Desk

Call the Help Desk

Most IT Help Desk organizations provide a phone number to call them, often in the form of a local toll-free number. For most larger Help Desk organizations, it has become common practice to use an IVR (Interactive Voice Response) system to respond first to a user Help Desk request.

IVR Systems

Users can interact with a call center's automated system through voice or touch-tone keypad inputs (see Figure 8-2).

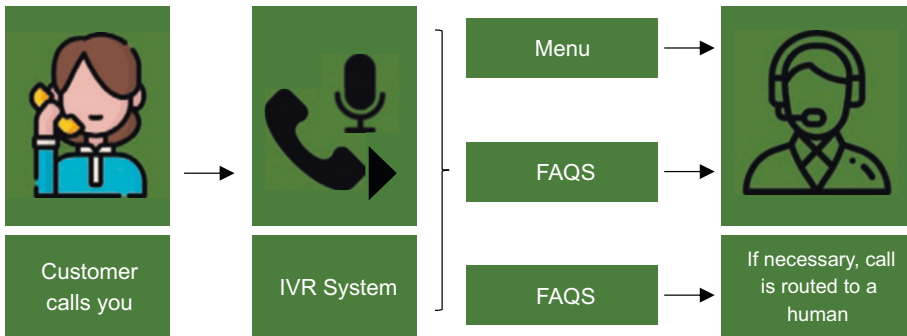


Figure 8-2. Example of an IVR Call Flow (Icons from Flaticon, 2023)

Several Functions Offered by IVR Systems

- Route calls to the appropriate department
- Self-service options to customers
- Gathering information from callers before connecting them to a live CSR

Benefits of Using an IVR System for a Call Center

- **Increased efficiency**

IVR systems can automate routine tasks such as routing a call to the most suitable CSR, scheduling appointments, and notifying clients about broader outages of systems. These IVR functions reduce the call center’s overall workload, allowing the IT Help Desk to handle more complex issues.

- **Saving costs**

IVR systems automate routine tasks and can directly address topics, such as frequently asked questions and notifications about outages, hereby leading to reduce staffing costs for the call center. Additionally, the ability to handle more calls with fewer CSRs reduces the IT costs of an organization.

- **Improved customer experience**

IVR systems allow callers to use its voice or the phone's keyboard navigating the IVR system to find the right department quickly and easily, or to access the self-service options they need. This results in shorter waiting times, faster resolution of issues, and a more positive overall customer experience.

- **Data collection**

IVR systems gather data about customer behaviors and preferences with a positive impact on the overall customer experience and provide data for business decisions.

Drawbacks of the IVR System

- **Complexity**

The call flows of IVR systems can be difficult to set up and maintain. This is challenging for IT processes with complex call routing requirements or for Help Desk support provided in multiple languages.

- **Limited personalization**

IVR systems often feel impersonal to customers, as they are interacting with a machine rather than a live person. This can lead to frustration or dissatisfaction with the service.

- **Language barriers**

IVR systems may face difficulties in understanding voice commands due to the accent of the speaker. The same English sentence spoken by a person from England, France, Korea, Brazil, and Senegal will sound quite different. Many IVR systems cannot cope with the variety of accents, resulting in misunderstanding of customer requests and in frustration with IVR systems.

- **Technical issues**

IVR systems are basically special-purpose computer systems reliant on technology, which can occasionally malfunction or experience downtime. This can lead to disruptions in service and negative customer experiences.

Overall, an IVR system can be a valuable tool for a call center as it provides cost savings and increases the capacity of the Help Desk Center. But it is important to consider the benefits and drawbacks carefully before implementing one. The design of the call flows should provide more efficiency to the call handling process but needs to avoid that users are discouraged from using the Help Desk function.

Help Desk

Here are the pros and cons of calling a Help Desk.

Pros

- Calling a Help Desk allows users to immediately interact with a live CSR.
- Direct communication allows the CSR to ask questions if the problem is not clearly described or if the CSR needs more information to qualify the cause of the issue.
- In general, a call to a Help Desk is appropriate for more complex issues that require a back-and-forth dialogue.

Cons

- It can be a frustrating experience for a Help Desk user to wait for a long time until the next CSR is available.
- The Help Desk user may need to leave a voicemail and wait for a callback. This is a frustrating experience if the voicemail system does not provide information in which time frame a callback will happen, or if it is vague like “We will return your call during the next business day.”

Email

Most Help Desk organizations offer an email address to which users can write their issues. This option of communication is often useful when a problem does not require an urgent solution but should be rectified at some stage. Email is a good option for users who are frustrated with IVR systems or do not have the patience to wait for an unpredictable period when calling a Help Desk Call Center.

Here are some of the pros and cons of Help Desk email support.

Pros

- Emails allow customers to describe the issue or question in a clear and concise manner and in their own time.
- Emails provide a reference to the customer and the CSR in further communications, and it is always clear to both the CSR and client what the subject is.
- Emails have the convenience of time, so both parties can send and reply when they have time to do so.

Cons

- The time to resolve a problem can be long if many questions are required to qualify the issues and one or both sides take a long time to respond to the email.
- Emails are prone to miscommunication due to a lack of tone and body language.

Online Chat

Online chat provides a communication channel for a user to interact with a CSR in real time through an online chat platform. It enables customers to ask questions, report issues, and seek assistance from the CSR. The Help Desk online chat can be accessed through a PC via a browser or by phone via a mobile application. The chat provides customers with immediate assistance without the need to wait, which is often the case when waiting at a call line or for the response to an email. Help Desk Chat can be either synchronous, where the conversation takes place in real time, or asynchronous, where the customer can send a message and

receive delayed replies from the CSR. Online chats are often a convenient and efficient way to provide customer support and ensure customer satisfaction.

Here are the pros and cons of Help Desk communication by chat.

Pros

- A chat provides immediate interaction between a Help Desk CSR and a user.
- Interaction between the user and the Help Desk agent is in real time.
- Chat communication is often a fast path to solve lesser complex problems.

Cons

- Can be prone to miscommunication if the customer and the representative are not clear in their messages.
- Chats are less suitable for more complex issues that require a longer dialogue to qualify and solve.

Chatbots

IVR systems were developed with the goal of automating some functions of a Help Desk Call Center. Chatbots fundamentally have a similar objective, which is to automate communication in a chat line. They are increasingly used for many applications in vocations, such as e-commerce, healthcare, education, banking, tourism, and entertainment.

A chatbot uses artificial intelligence (AI) and natural language processing (NLP) to interact with customers and provide support. Chatbot application for a Help Desk includes the following:

- Answer common questions;
- Route support inquiries to the appropriate CSR;
- Provide self-service options.

Benefits of Chatbots

- Provide 24/7 support without the need for human intervention.
- Provide answers to common questions.
- Can interact with users at any time.
- Guide customers through self-service options.
- Reduce the overall Help Desk workload on support by handling routine inquiries.
- Efficiently route support inquiries to the appropriate CSR.
- Improve response times to customer issues and ensure that customers receive timely support.
- Provide a more personalized support experience for customers.
- NLP-based chatbots can understand natural languages and therefore often provide a human-like interaction.

But chatbots are not created to replace human CSRs. Even though chatbots can handle routine inquiries, they may not be able to handle complex issues and lack the empathy of a human being.

Summary – Options to Communicate with a Help Desk

The choice of communication method for a Help Desk depends on factors such as the complexity of a problem, the urgency to solve it, and the reliability of the interaction in the chosen medium. Email is best suited for detailed issues that do not require an immediate response, while a call or chat may be more appropriate for urgent or less complex issues that require immediate attention.

Trouble Ticketing Systems

Most IT Help Desk organizations use trouble ticketing systems to manage Help Desk requests from their client in a structured form. A trouble ticketing system is a tool in the form of a software application that allows CSRs to track and manage customer issues or requests. It provides an organized way for Help Desk personnel to receive, document, and resolve issues or requests from customers.

With a trouble ticketing system, a Help Desk organization can create a defined work process to handle requests from clients (see Figure 8-3). This ensures that calls are processed in a consistent form and can be handed over from one CSR to another (if this is required due to the domain expertise of the CSR) or to hand over critical trouble tickets to the next shift in the Help Desk organization that provides support around the clock.



Figure 8-3. *Trouble Ticketing Systems (rawpixel.com, 2023)*

When a customer contacts the Help Desk, the CSR creates a ticket that contains detailed information about the issue, including its severity and urgency. The trouble ticket will include important information associated with the customer request, such as screenshots, error messages, and data about the IT systems used, which include details about the hardware, operating system, and its versions. All relevant information about the customer's problem is then stored in a database, and the CSR and other members of the Help Desk organization can easily track the ticket's status, monitor progress, and collaborate with colleagues to resolve the issue.

Trouble Ticketing Dashboard

The trouble ticketing system allows CSRs to process customer requests efficiently and consistently. As all information is available in a database, the Help Desk team, including its management, can monitor the overall performance of the Help Desk through the trouble ticketing dashboard.

The dashboard provides a centralized view for managing and monitoring customer support inquiries. It provides a quick and easy way for the Help Desk team to view the status of open tickets, track progress, and prioritize work based on urgency and severity.

The dashboard displays key metrics including

- The number of open tickets;
- The average response time;
- The number of tickets resolved within a certain period.

The data displayed in the dashboard provides a snapshot of the overall performance of the Help Desk team and allows managers to identify areas for improvement (see Figure 8-4).

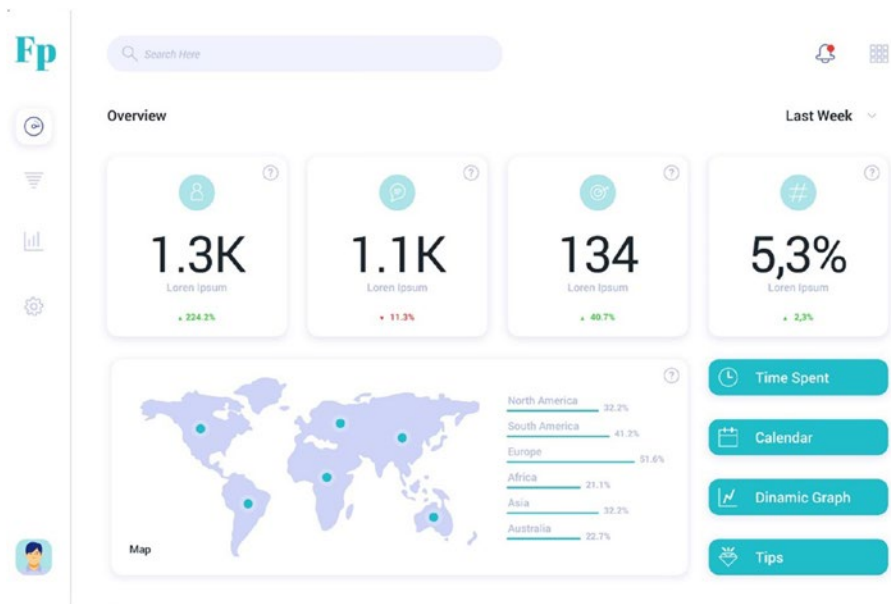


Figure 8-4. Help Desk Dashboard Representation (pikisuperstar, 2023)

The dashboard is an important tool to prioritize the overall workload of the Help Desk team. From the dashboard, the team can assess the open tickets, organized by priority, severity, or date. This allows them to quickly assess the most pressing issues and prioritize their work accordingly. Trouble tickets include details such as the customer's name, the date the ticket was created, and the status of the pending issues.

An important feature of a trouble ticketing system is the ability to assign tickets to a specific CSR or domain specialist and track the status of escalations. Beyond the management of user problems, the system can collect feedback from customers and ask for ratings from clients.

Overall, the dashboard of a trouble ticketing system is an essential tool for managing customer support inquiries. It provides a real-time view of key metrics and allows support agents to quickly assess and prioritize their work. This leads to faster response times, improved customer satisfaction, and a more efficient support operation.

A trouble ticketing system is an essential tool for any Help Desk team. It can help them to streamline customer support operations, improve communication and collaboration, and enhance customer satisfaction.

Every organization with an online presence needs to offer some form of Help Desk. If the Help Desk function is part of the business, the service will not be charged to the client. A bank offering Internet banking services must have a Help Desk function for struggling clients, and therefore, the bank cannot charge for the service. To keep the costs for a Help Desk function at bay, organizations have introduced Interactive Voice Response (IVR) systems. The IVR system replies to clients' calls and obtains basic information, such as name, account number, language, and type of question, and then passes it on to the Help Desk agent when the call is connected to a Help Desk agent. The other means for organizations to reduce their Help Desk workload is to offer online information about the services and provide information about frequently asked questions related to a topic.

IT Help Desk services provided by an external service provider will always be chargeable. The charges will depend on the complexity of the application, the response time, the availability of the Help Desk, and the communication channel.

Typical relationship between costs and service level is represented in Table 8-1.

Table 8-1. *Typical Relationship Between Costs and Service Level*

Cost of service	Low	Medium	High
Application complexity	Low	Medium	High
Access to service by/via	Email	Chatline	Phone
Service coverage	Mon–Fri 9 a.m.–5 p.m.	Mon–Sat 7 a.m.–7 p.m.	Mon–Sun 24 hours
Response time	Next day	1–4 hours	<5 minutes

IT Help Desk clients can choose from a variety of service-level agreements (SLAs). Most clients choose an SLA that provides online telephone support during their business hours, which is a five- or six-day week with a coverage of eight to ten business hours. Some customers with critical IT systems, or organizations operating 24 hours daily, require a Help Desk coverage of 24/7 from the IT Help Desk.

8.2 IT Service Desk

The Information Technology Infrastructure Library (ITIL) defines a Service Desk as follows: “The single point of contact between the service provider and the users. A typical Service Desk manages incidents and service requests and handles communication with the users.” This sounds a little

unclear and is similar to what was described in Section 8.1 about the Help Desk. When ITIL developed the concept of “managing IT as a service,” it evolved the concept of a Help Desk from an end-user-centric support model to a more generic support function that includes all services delivered by a DC, also called IT Service Management (ITSM).

Service Desk services include multiple ITSM activities:

- Service request management;
- Incident management;
- Knowledge management;
- Self-service;
- Reporting;
- Change management.

The classic Help Desk function dealt mainly with break-fix of a function or service that is requested by an IT user. In ITSM, a Help Desk function is called “incident management.” The preceding list shows that a Service Desk has several more functions than just incident management. Some other functions are represented in Figure 8-5.

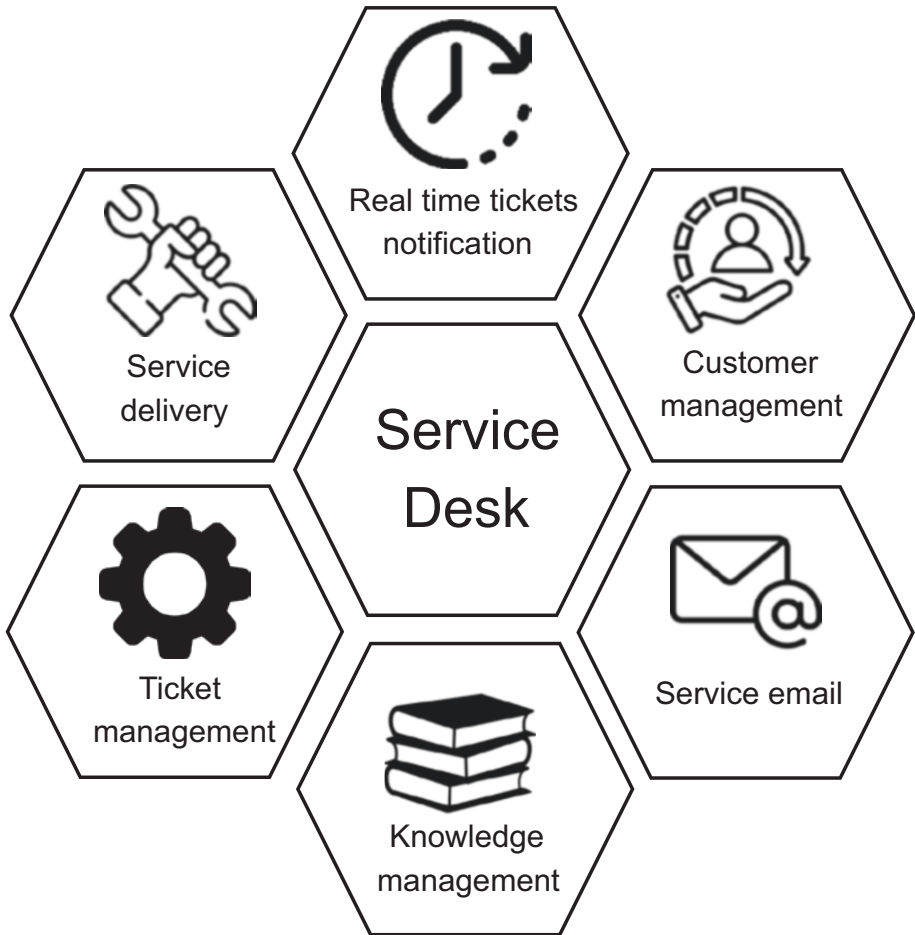


Figure 8-5. ITSM Service Desk Functions (Icons from Flaticon, 2023)

In addition to incident report, an IT Service Desk assists customers with service request management. It creates and manages IT knowledge for an organization and offers self-service for customers who have sufficient IT knowledge to resolve their issues quickly and independently.

A critical function of a Service Desk is the monitoring of the overall reliability of the IT systems of an organization and the measuring of the frequency of incidents. Most organizations define key performance

indicators (KPI) to measure the overall reliability and user-friendliness of the supported IT infrastructure and applications to deliver a robust, service-focused, and customer-centric IT environment for the organization.

Depending on the role of the IT Help Desk to an organization, the Help Desk support of a Service Desk support can be chosen.

8.3 Remote DC and Edge Computing Support

Remote DC Support

Some organizations run their IT systems in hybrid mode where some workloads are designated to a CSP or a private IT service provider, while other workloads run in their own private DC. This split may be mandated by data privacy laws of a country or by regulations in a specific industry such as banking, where banks are obliged to keep their data private. This requires that at least the bank's data storage systems are located at an office or DC on their premises.

Even though the local DC or local storage systems are on the premises of the clients, the organization may still want to get support for this equipment as they may not have the right staff to run a DC in a fully secured environment, or if they want to focus their staff on their core business.

Edge Computing Support

As described in Section 3.9, edge computing is strongly evolving to overcome limitations of the centralized cloud infrastructure (see Figure 8-6). It will play an even more important role in the emerging smart infrastructure of the public and private sectors.

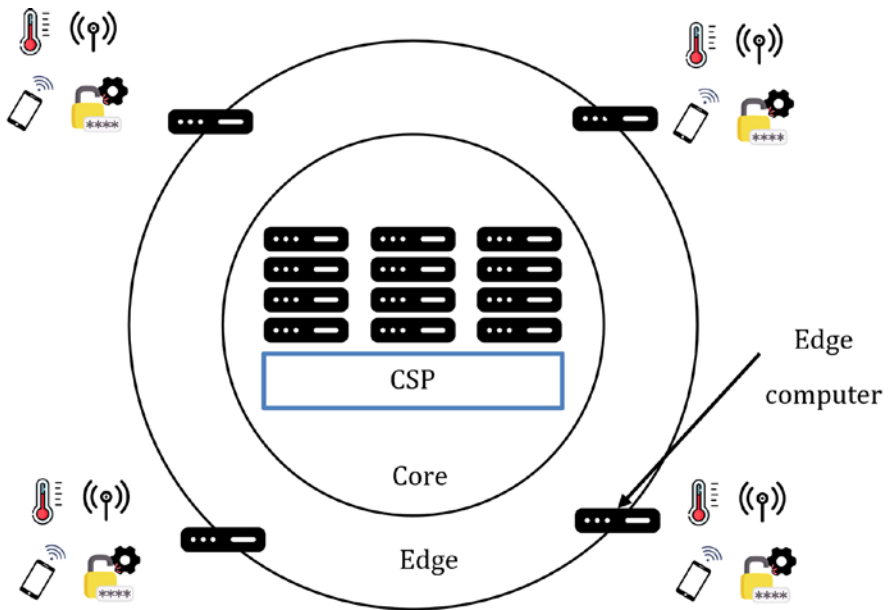


Figure 8-6. *Edge Computing (Icons from Flaticon, 2023)*

Services offered for remote DC and Edge computers are extensions of the Service Desk services, which include the following:

- **Service request management**

Authorized users can request the deployment of new services. In the service management process, resources in the form of HW components, operating systems, and applications are deployed.

- **Incident management**

This includes the monitoring of remote compute resources. HW and SW component failures can be detected and solved through a trouble ticketing process. This may include activating repair services for HW components.

Users can report incidents of applications that do not perform as specified at the remote DC or Edge computer systems.

- **Knowledge management**

Inventory management is performed in all HW and SW components deployed in the remote DC or Edge systems, including revision levels of the operating system, firmware, applications, and installed patches. Attempted cybersecurity breaches are logged, and preventive measurements are put into action.

- **Self-service**

Key functions of the remote DC or Edge systems are documented and can be accessed by authorized users including key functions like resetting and rebooting.

- **Reporting**

The status of the remote DC or Edge system is constantly monitored, and key parameters are stored, including load patterns of CPUs, memory, storage systems, power consumptions, and temperatures. Reports are created for the ITSM administrators on a regular basis to assess if resources are adequately deployed, and if trends like declining storage capacity lead to issues.

- **Change management**

All changes to the remote DC or Edge systems are deployed through a change management process. The process starts with a change request for an HW or SW component. The system administrator will ensure that all data is backed up and a secure point of return is created in case a change leads to service disruptions.

The changes are implemented and thoroughly tested afterward. If the tests are successful, the change management process is concluded, and all changes are properly documented.

8.4 Summary

This chapter explored the intricacies of IT support services within modern organizations. Beginning with the IT Help Desk, its pivotal role as the single point of contact for users needing IT support was underscored. This encompassed tasks ranging from simple logins to intricate troubleshooting, with a special emphasis on the growing use of Interactive Voice Response (IVR) systems and their benefits and limitations. Additionally, the emergence and role of chatbots in aiding IT support were examined. Transitioning to the IT Service Desk, the content expanded on the concept by integrating IT Service Management (ITSM). While incident management remained at its core, a Service Desk took on more functions, including knowledge management and self-service, all aimed at a comprehensive IT support approach. Lastly, the chapter touched upon the hybrid mode in some organizations, necessitating remote DC and edge computing support. This support ensured the smooth operation of IT systems in various environments, highlighting the unique challenges and solutions presented by remote and edge computing. Through this chapter, readers gained a holistic understanding of the multifaceted world of IT support services, from foundational aspects to modern evolutions.

In the next chapter...

- A full summary – to help you consolidate all the information
-

CHAPTER 9

Summary

Over the last 60 years, the IT industry has evolved from humble beginnings to one of the core building blocks of modern society. We can hardly imagine living in a world without ubiquitous telecommunication and IT services anymore. In this book, we have described how the key building blocks, such as data centers, computer systems, storage systems, and IT security systems, have evolved and how these elements can be used to build a state-of-the-art resilient IT infrastructure.

The process happened in two phases: from the 1960s to the late 1990s, computer systems played an ever-increasing role in supporting a company's business. In the second phase, which started in the late 1990s, IT systems were not only supporting the business of companies but became an essential part of a new type of company. Google, Facebook, Amazon, Netflix, and Uber are companies that have built their core business with applications running on computer systems.

Today's world of social media is built entirely on computer systems and a modern telecommunication infrastructure. In the core of the second wave of computerization is the evolution of the Internet, which enabled people to connect from anywhere to almost anybody elsewhere in the world in real time. The ubiquitous presence of the Internet made it possible to automate most business processes.

The wonderful world of the modern Internet would not be available without the underlying resilient IT infrastructure.

9.1 Resilient IT Infrastructure

A resilient IT infrastructure consists of the following elements:

- DC building and facilities;
- Network systems;
- Computer systems;
- Storage systems;
- Backup systems;
- DC resiliency;
- IT support services.

9.2 IT Services Provided by Cloud Service Providers

With the introduction of cloud data centers by cloud service providers (CSPs), a growing share of IT workloads was moved from private data centers, which are also called private clouds. CSPs offer several options to use their IT services, which include services such as Data as a Service, Communications as a Service, Infrastructure as a Service, and Platform as a Service.

Cloud data centers and private data centers are used by most organizations at the same time. A common approach is to use a private data center to run the core applications, while public cloud data centers are used for less critical applications, such as websites and email. Organizations need to decide which application to run and where to run them and can choose from the following deployment options: public model, private model, community model, and hybrid model.

9.3 Data Center

A DC facility is a physical building used for equipment to run a data center. DCs used to be a single room in an office building. Today, a modern DC of CSPs covers full floors in a large building, or an entire building is just a large DC.

The DC site must provide sufficient floor space to hold all the IT infrastructure elements that an organization wants to deploy. To save costs and for security reasons, DC buildings are usually located outside of CBD areas but need to have good access to the public infrastructure in the form of road access, electricity, and fiber connection.

The power consumption of a DC facility can reach up to 100 megawatts; hence, an adequate power connection is required from a local power service provider. It is important to get high-reliability power for the DC. Even though DCs usually have supplementary power systems in the form of a UPS, the stability of the supplied power is an important factor in securing highly reliable DC services.

The enormous amount of power delivered to a data center is to a large degree converted into heat, which must be removed from the IT infrastructure through cooling, as all IT equipment has a defined temperature range.

DC security in the context of the facility focuses mainly on the building security itself and the access control to the facility.

Building management systems (BMS) allow to manage DC facilities from a centralized console. All critical data from the facility is measured and displayed in dashboards such as temperature, humidity, power, and ventilation flow for each of the rooms. Access control, security logging, and video surveillance of the DC rooms and the perimeter are managed by the BMS as well.

IT infrastructure refers to the IT components that are deployed into the DC, which include the network equipment, server, and storage systems, all mounted into racks, and the cabling connecting the systems.

9.4 Compute

A system that can be programmed to carry out arithmetic or logical operations automatically is called a computer. A computer system that is used by multiple users is called a server. Over the last 80 years, computers have evolved to the ubiquitous building blocks of almost every modern electronic device, machine, and system. Computers are the key building block of the Internet, connecting today billions of people and devices.

Computer systems used to be designed for certain types of applications. Through the exponential growth of computer performance over the last 60 years, computers have become so fast and cheap that specialized computer systems are only used in a few niche areas such as high-performance computing. Most modern DCs are deployed with large amounts of rack-mounted servers powered by microprocessors from Intel or AMD.

The operating system (OS) is an important component of each computer system. An OS manages and interacts with the hardware components of a computer system including the CPU, memory, and I/O devices such as disks, screens, keyboards, and network controllers. Operating systems are running on all types of computer systems. The OS features are dependent on the type of computer system they are running (machines, PCs, servers, compute cluster) and the type of applications. All operating systems manage key functions of a computer system including compute processes, memory management, file management, I/O management, and security.

A software layer was developed to distribute applications, nowadays called workloads, which separate workloads from the actual physical computer system. This allows us to distribute workloads over a large range of servers of the same or different types. The two main concepts that are widely in use are virtualization and containerization.

Through virtualization, workloads can run on a virtual machine (VM). To the workload, it seems it's running on the operating system of a real computer while it is running on a software layer. This concept of a virtual

machine allows a hypervisor to distribute workloads to a pool of compute servers, and hereby workloads can be moved around quickly within a DC or to a DC at another location. This makes the system very robust as neither a failure of a server, a rack, or even an entire DC stops a service, as the hypervisor can reroute the workload to a system that is operational.

A newer virtualization concept is based on container, reducing the images' size and simplifying applications' virtualization. Following the idea of a physical shipping container, a standardized large box in which all types of goods can be packed, a container for IT systems wraps an application into an isolated box. A container-based workload is much lighter than the code for a VM, as it does not include the operating systems and libraries.

Security and resiliency of a computer system are an important consideration in running a DC. Compute security describes how computer systems protect from theft, damage, or disruptions to their hardware, software, or data. If computer systems are connected to the Internet, the term "cybersecurity" is often used.

Resilient computing goes beyond the concept of secured computing by providing a system that is designed for a very high level of reliability and tolerance for changes that application updates can introduce. The concept of modern DCs with large amounts of servers makes it much easier to deal with hardware and software failures. The control plane of a hypervisor distributes a workload to a virtual machine and checks the execution status.

A hypervisor or container-based control plane can easily re-distribute applications to another computer system, thereby making the system fault tolerant to hardware and software failures of machines. But service interruptions can occur when the DC itself could become unavailable. A higher level of resilience can be achieved if multiple DCs, at different locations, preferably different cities, or countries, can be connected through a federated architecture. This elevates the concepts of distributed hypervisors to a higher level.

The introduction of VMs and containers allows workloads to move around easily. With the emergence of cloud DCs, organizations have the choice to shift some or all their applications to a cloud service provider. This provides organizations with choices about how they want to run their workloads by distributing them between their private cloud and CSPs. It also provides options to move workloads from their private cloud in case of system problems, high customer demand, and resources for a special project.

9.5 Storage

Storage systems are only connected via the I/O channel of a computer system and therefore not directly accessible by the CPU. Therefore, storage systems are called secondary storage, while primary storage refers to semiconductor storage directly connected to a CPU. Operating systems, applications, and data are stored in storage systems. All data stored in secondary storage is static and remains intact after the computer system or storage system is switched off.

Modern computer systems use HDDs or SSDs, or both. Compared to RAM, HDDs are much cheaper than SSDs on a per-gigabyte basis, but the access time is much slower, measured in milliseconds for HDDs and nanoseconds for RAM and SSDs. Types of external computer storage are direct attached storage, network-attached storage, and storage area networks.

Like the trend in computing, storage is not deployed by a few expensive storage systems, but by a large number of inexpensive disks that are held in racks in so-called JBOD configurations or bundled into disk arrays, also called RAID. Disk arrays use additional redundant components to increase availability, resiliency, and maintainability. Higher-end disk arrays include so much redundancy that the system does not have a single point of failure and is therefore very reliable. Disk arrays with hot-swappable features allow the replacement of components without stopping the system. This further improves the uptime.

Disk arrays can be configured in different modes of operations, which are called RAID levels, featuring several configurations to optimize redundancy, capacity, and speed. The various RAID levels were standardized by the Storage Networking Industry Association. The most used RAID levels are 0 for block-level disk striping, 1 for block-level disk mirroring, 5 for block-level disk striping with distributed parity, and 10 for block-level disk striping of mirrored disks.

Storage virtualization provides a logical view of the physical storage resources to a computer system. It creates a single resource pool of storage that includes all storage media that are accessible in the form of HDDs, SSDs, optical drives, or tapes. Benefits from storage virtualization are simplified management, better storage utilization, storage life cycle extension, and uniform storage features.

Storage security comprises of all manual or automated technologies and processes to ensure the integrity and security of data in storage systems. Data stored in computer systems, portable storage devices, and external storage in the form of NAS, SAN, or DAS are protected by the storage security system using physical hardware protection and security software. For every organization that handles sensitive data, protection is essential to avoid data theft and ensure continuous operations of its business activities.

Most threats to data security are like those of DC security and cybersecurity, which are described in Chapter 4 Section 4.5, such as access control and disaster prevention. To bolster data security and ensure data resiliency, it is imperative for organizations to enforce mandatory encryption across all data, irrespective of its location – be it on servers, storage systems, personal computers, or mobile devices. Deletion of data is part of these measurements. If files or even a disk volume is deleted, the file system only deletes the entries in the file system tables, but the actual data blocks on a disk remain unchanged. There is software available that can recover such data. Therefore, the IT organization must ensure that data is fully erased when disks or tapes are replaced or disposed of.

Storage provisioning is the process of assigning storage to workloads and optimizing the performance of the available storage resources.

Storage provisioning can be performed based on the underlying storage systems. The introduction of cloud data centers, hypervisors, and container-based computing has changed the methods to provision storage. Virtualization uses an application-centric storage provision and is based on a virtualized storage environment. Virtual storage is assigned to computer systems, hypervisors, or containers during the provisioning process. The storage allocation is based on application-defined capacity, availability, and performance requirements and is called “thin provisioning.” Alternatively, storage provisioning can be performed traditionally, which is called “thick provisioning.”

9.6 Network

Data center networking provides the resources to connect server and storage system with each other and to external parties. Networking uses hardware, in the form of switches, routers, gateways, and software to perform the networking operations.

Modern DC networking architectures provide connectivity for virtualization platforms such as virtual machines, containers, and bare-metal configurations. DCs have evolved from a single computer in a room to today’s highly distributed DC infrastructures. The network components of the DC have evolved with a similar progress as server and storage technology; they are the third building block in the immense evolution of DCs.

The key HW components of a DC network are cables, switches, routers, and gateways. Structured cabling is a method to avoid the messy cabling that was often found in larger DCs. It uses a main distribution area as a central point of cable connections.

Switches and router may look quite similar, but they have different functions. Switches connect end devices such as server and storage systems, PCs, and printers. The main purpose of a switch is to receive

a data packet, determine its destination, and forward the packet to the destination address. The role of a router is to find the fastest and best route for a packet to reach its destination.

A gateway transmits data from one discrete network to another. Gateways can be implemented in hardware or software. The main difference between gateways and routers/switches is that a gateway can transport data across different protocols. Gateways can operate on different levels of the communication stack defined by the OSI standard model. With the unification of networks under the standardized TCP/IP gateways are less commonly used today and are primarily found at the boundaries between private corporate networks and the public Internet.

There are several network topologies available depending on the size and features of a DC. Most modern large DCs use a multi-tier network topology to deploy their networks for their large DCs. Cross-connection switches connect cables, subsystems, and equipment to the racks. A cross-connection provides excellent cable management and design flexibility to adapt to changing requirements. This topology provides operational advantages, as all connections are managed from one location.

The trend toward open standards and cloud architectures led to the development of software-defined networking (SDN) technology. In SDN networks, the transport of data packages is separated from the routing process of data packages. Data transport is handled in the SDN data plane, while routing and management are performed in the control plane. The Open Networking Foundation drives the standardization of OpenFlow, the commonly used SDN protocol. Many organizations have changed their network architecture to SDN by replacing their proprietary switches with switches using the simpler SDN architecture. SDN network can be deployed with inexpensive, commoditized white label switches.

Network fault management is an important network component to secure a high availability of the network. Network failures can be caused by mistake in the network design, malfunctions of cables or network equipment, or external factors such as natural disasters.

9.7 Backup

Backup System

Data backup is the process where data is copied to another medium with the objective of retrieving the data in case of an event that caused the loss of data on the primary medium. The process of restoring the data from the secondary medium is called recovery. Data backup requires the copying and archiving of computer data to make it accessible in case of data corruption or deletion.

Data backup is an important element to make an IT environment resilient, making it an essential part of any sensible disaster recovery plan. Today, data backup is often performed by using cloud storage. Data can therefore be archived on a local system's hard drive, an external storage system, or by using cloud storage.

Remote backup can be performed through the resources of public cloud storage or by a private cloud if an organization has backup resources in multiple locations, preferably in several regions or countries. A remote backup system comprising access point network bases that build on a global scale security network using existing communication networks functions as access points for accessing user machines existing in appropriate countries. It is installed in one or more locations in each country. A main remote backup system installed in one country can connect with other access points and can store and retrieve data from the remote access point.

The remote backup system has a configuration in which the backup system is installed in one country and the network thereof uses a security network, excluding the Internet, so that unauthorized entry is almost impossible. In addition, it has the effect of building a backup system that can prevent natural disasters. In this way, the base to be backed up is set up in a country that has few disasters due to earthquakes, good power supply, network satellites, and its own security network. Remote backup can be performed very safely and reliably among backups.

9.8 Resiliency

In the interconnected age where computer systems span the globe, ensuring the safety and integrity of data centers is of paramount importance. This protective endeavor, often encapsulated by the term “cybersecurity,” revolves around shielding data centers – their computing, storage, and network systems – from potential theft, damage, or disruptions. Vulnerabilities can emerge due to shortcomings in design, implementation, operation, or internal protocols. Once these vulnerabilities come to light, they might be preyed upon by malicious entities who, fortified with research, reverse engineering, or automated tools, initiate their malevolent endeavors. Notably, while disruptions might arise from commonplace threats like viruses and malware, the digital landscape also witnesses sophisticated incursions like ransomware and DDoS attacks.

At the heart of a resilient cybersecurity strategy lie three foundational pillars: prevention, where protective measures and barriers deter unauthorized access; detection, emphasizing real-time monitoring to identify threats; and response, a predefined counteraction plan for any security anomalies. Delving deeper, one realizes that hardware components, too, play dual roles. Malicious hardware modifications can introduce security breaches, but, contrastingly, hardware-based protective solutions like secure enclaves and trusted platform modules offer an added protective layer. In the software realm, the mantra of “security by design” reigns supreme. Instead of retrofitting security as an afterthought, this approach stresses its integration right from the developmental onset.

Ensuring that only authorized personnel can access pivotal information is crucial, and this is where access control lists step in. By dictating varied permissions for different objects, they streamline access based on user classifications – be it end users, developers, or administrators. The human element, though, remains a wildcard. In today’s digitized organizations, IT access is widespread, rendering

everyone a potential security loophole. This underscores the need for universal awareness, fostering a vigilant security ethos and mandating periodic security training across all ranks. However, even the most fortified bastions might experience a breach, and during such instances, the response's swiftness and organization are pivotal. Comprehensive incident reports should chronicle the breach, followed by a meticulous probe to discern the attack's nature and the vulnerabilities it exploited. Prompt and decisive action can halt an incident from snowballing into catastrophic escalations, like extensive data leaks or complete system outages.

Segueing to the notion of data center (DC) resiliency, it embodies a DC's prowess to spring back from adversities – whether stemming from natural calamities or cyberattacks – and maintain seamless operations. This resilience is anchored firmly in robust DC security. To realize this resilience without overburdening IT coffers, it's imperative for organizations to introspect their workloads, calibrating the resiliency levels based on both business and technical paradigms. In conclusion, as we plunge deeper into the digital epoch, the intricate dance between security and resiliency will gain prominence, highlighting the need for a unified, holistic approach to safeguarding data centers.

Beyond the technical components, a significant factor in the resilience and security of data centers is their physical security and location. Strategically positioning data centers in regions less susceptible to natural calamities (such as floods, earthquakes, or hurricanes) aids in disaster mitigation. Furthermore, on-site security measures, including biometric access controls, surveillance systems, and security personnel, are essential in warding off physical intrusions or thefts. Integration with modern cloud infrastructures also brings forth hybrid solutions, adding another layer of resiliency. By leveraging cloud-based backups and disaster recovery solutions, organizations can ensure quicker recovery times and data availability even if the primary data center faces disruptions.

Lastly, a continual evolution in the cyber threat landscape mandates that cybersecurity strategies remain adaptive. Organizations should invest in regular threat intelligence and vulnerability assessments, ensuring that they are always prepared for emerging threats. Cyber insurance is another evolving domain, providing a safety net against potential financial implications of significant breaches. In essence, the intricate dance between security and resiliency in the realm of data centers is an ongoing journey, not a destination. As technology evolves and cyber threats morph, the strategies and tools to combat them must likewise adapt, underscoring the imperative of a proactive and ever-evolving approach to data center security and resilience.

9.9 IT Services

Help Desk

A Help Desk serves as the primary point of contact for users experiencing issues with IT systems. This support function enables users to report incidents they can't rectify independently. Initially, users reached the Help Desk via toll-free numbers, speaking directly with customer support representatives (CSRs). As communication technologies evolved, email support emerged for less time-sensitive queries or those requiring documentation. Social media ushered in live support chats, increasing efficiency as CSRs could multitask across several chats concurrently. The cutting edge in Help Desk communications now includes AI-driven chatbots and mobile apps; the former assists with common issues, while the latter offers convenient reporting and tracking tools for users on the go.

Service Desk

Building upon the foundation of the Help Desk, the Service Desk offers an expansive support spectrum, embracing all services an IT organization delivers. Beyond troubleshooting, the Service Desk facilitates service requests and knowledge management and even provides self-service portals for IT-savvy users to swiftly address their issues. These portals often feature a “service catalog,” detailing the array of available services and setting clear expectations. One pivotal aspect of the Service Desk is its commitment to proactive IT system maintenance. Through practices like system updates, patch management, and preventive measures, it mitigates potential system vulnerabilities. To ensure optimal performance, Service Desks adhere to specific metrics, KPIs, and service-level agreements (SLAs). Such benchmarks help gauge their efficacy and efficiency.

Furthermore, with the integration of structured approaches like ITIL, Service Desks achieve a harmonized IT service management, ensuring seamless operations even in hybrid environments where workloads are distributed between in-house IT and external service providers. Given the ever-evolving IT landscape, Service Desks must anticipate the needs of the future. The rapid growth of edge computing and IoT demands adaptability. As these technologies become integral to infrastructures, Service Desks will have to evolve their support mechanisms to cater to this new breed of devices and systems.

Lastly, while technology is a major component of the Help and Service Desks, the human element remains indispensable. The continuous training of personnel, coupled with their ability to empathize and communicate effectively, ensures that the technical support ecosystem remains both efficient and user-friendly.

References

1. Computers at NASA (NationalArchivesCatalog, 1952–1968)

<https://catalog.archives.gov/id/278196>

[Online; accessed Oct. 4, 2023]

1-1. Cloud Data Center Example (rawpixel.com, 2023)

www.freepik.com/free-photo/cloud-storage-banner-background_16016448.htm#query=data%20center&position=0&from_view=search&track=ais

[Online; accessed Oct. 4, 2023]

2-1. Google Data Center, Iowa, USA (ChadDavis, 2017)

<https://flickr.com/photos/146321178@N05/49062863796>

[Online; accessed Oct. 4, 2023]

2-2. Modern DC Racks (svstudioart, 2023)

www.freepik.com/free-ai-image/cloud-rack-equipment-room-with-big-data-cyber-internet-content-blue-light-server-interior-modern-communication-storage-hardware-system-hub_40583169.htm#page=2&query=server%20rack&position=13&from_view=search&track=ais

[Online; accessed Oct. 4, 2023]

2-3. Data Center Layout Provides Adequate Airflow (Florian Hirzinger, 2009)

https://commons.wikimedia.org/wiki/File:CERN_Server_03.jpg

[Online; accessed Dec. 15, 2023]

2-4. UPS for a Data Center (Bercik, 2020)

<https://commons.wikimedia.org/wiki/File:DRUPS.jpg>

[Online; accessed Oct. 4, 2023]

REFERENCES

3-1. Motherboard of a Computer System (Macrovector, 2023)

www.freepik.com/free-vector/system-plate-pc-isometric-illustration-with-semiconductor-elements-slots-microchips-capacitors-diodes-transistors_7496526.htm#query=motherboard&position=17&from_view=search&track=sph

[Online; accessed Sep. 21, 2023]

3-2. Moore's Law (Source: Wikipedia, 2023)

https://en.wikipedia.org/wiki/Transistor_count

[Online; accessed Dec. 9, 2023]

3-3. Cloud Data Center (Benzoix, 2023)

www.freepik.com/free-photo/server-racks-computer-network-security-server-room-data-center-d-render-dark-blue-generative-ai_38095230.htm#query=server&position=0&from_view=keyword&track=sph

[Online; accessed Sep. 21, 2023]

3-4. DC Racks with Rackmount Server (Macrovector, 2023)

www.freepik.com/free-vector/network-servers-isolated_4658640.htm#query=server%20racks&position=0&from_view=search&track=ais

[Online; accessed Sep. 21, 2023]

3-5. Blade Enclosure (ColossusCloud, 2016)

<https://pixabay.com/photos/server-cloud-development-business-1235959>

[Online; accessed Sep. 22, 2023]

3-6. High-Performance Computer (HPC) (Svstudioart, 2023)

www.freepik.com/free-photo/data-server-racks-hub-room-with-big-data-computer-center-blue-interior-hosting-storage-hardware_37159301.htm#query=super%20computer&position=4&from_view=search&track=ais

[Online; accessed Sep. 22, 2023]

Table 3-1. Top Five HPC Systems Based on LINPACK Benchmark

(Source: top500.org, 2023)

www.top500.org/lists/top500/2023/06/

[Online; accessed June 2023]

3-11. Smart Lamppost (Macrovector; yoyonpujiono; veryicon.com, 2023)

Background Lamp: www.freepik.com/free-vector/metallic-street-lamp-isolated_10603682.htm#query=street%20lamp&position=4&from_view=keyword&track=ais

Motion Sensor Icon: www.flaticon.com/free-icon/motion-sensor_4883053

Other icons from veryicon.com

[Online; accessed Sep. 22, 2023]

3-15. Top Ten Server Workloads (Source: techtarget, 2023)

www.techtarget.com/searchdatacenter/definition/workload

3-16. Automated Workload Distribution with Kubernetes (Icons from veryicon.com, 2023)

Calendar Icon: www.veryicon.com/icons/miscellaneous/data-product-icon-library/calendar-412.html

Controller Icon: www.veryicon.com/icons/miscellaneous/manufacturing/controller-9.html

[Online; accessed Sep. 23, 2023]

Table 3-3. Example of Microsoft Azure Cloud Services Charging Structure (Source: cloudmore, 2023)

<https://cloudmore.com/content-hub/cloudmore-supports-your-azure-revenue-growth-using-azure-plan-and-microsofts-new-commerce-experience>

[Online; accessed 2023]

REFERENCES

4-2. HDD (Unsplash, 2018)

https://unsplash.com/pt-br/fotografias/wYD_wfifJV5

[Online; accessed Sep. 25, 2023]

4-4. Comparison of Storage Types (Icons from Flaticon, 2023)

Server Icon: www.flaticon.com/free-icon/server_969438?term=server&page=1&position=23&origin=search&related_id=969438

Storage Disks Icon: www.flaticon.com/free-icon/server_689401?term=database&page=1&position=8&origin=tag&related_id=689401

[Online; accessed Sep. 25, 2023]

4-5. Direct Attached Storage (Icons from Flaticon, 2023)

Server Icon: www.flaticon.com/free-icon/server_689401?term=database&page=1&position=8&origin=tag&related_id=689401

Storage Disk Icon: www.flaticon.com/free-icon/server_689401?term=database&page=1&position=8&origin=tag&related_id=689401

Computer Icon: www.flaticon.com/free-icon/computer_1865273?term=computer&page=1&position=3&origin=search&related_id=1865273

[Online; accessed Sep. 25, 2023]

4-6. Network-Attached Storage (Icons from Flaticon, 2023)

Computer Icon: www.flaticon.com/free-icon/computer_1865273?term=computer&page=1&position=3&origin=search&related_id=1865273

Laptop Icon: www.flaticon.com/free-icon/laptop-screen_3616437?term=laptop&page=1&position=25&origin=search&related_id=3616437

Server Icon: www.flaticon.com/free-icon/server_969438?term=server&page=1&position=23&origin=search&related_id=969438

Switch Icon: www.flaticon.com/free-icon/network-switch_9694716?term=switch&page=1&position=9&origin=search&related_id=9694716

NAS Icon: www.flaticon.com/free-icon/nas_4943962?term=nas&page=1&position=1&origin=search&related_id=4943962

[Online; accessed Sep. 25, 2023]

4-7. Storage Area Network (Icons from Flaticon, 2023)

Computer Icon: www.flaticon.com/free-icon/computer_1865273?term=computer&page=1&position=3&origin=search&related_id=1865273

Server Icon: www.flaticon.com/free-icon/server_969438?term=server&page=1&position=23&origin=search&related_id=969438

Storage Icon: www.flaticon.com/free-icon/server_689401?term=database&page=1&position=8&origin=tag&related_id=689401

[Online; accessed Sep. 25, 2023]

4-9. Disk Array (Macrovector, 2023)

www.freepik.com/free-vector/datacenter-equipment-isometric-set_5970811.htm#query=computer%20Array&position=36&from_view=search&track=ais

[Online; accessed Sep. 25, 2023]

4-17. Storage Virtualization Architecture (Icons from Flaticon, 2023)

Server Icon: www.flaticon.com/free-icon/server_969438?term=server&page=1&position=23&origin=search&related_id=969438

Storage Icon: www.flaticon.com/free-icon/server_689401?term=database&page=1&position=8&origin=tag&related_id=689401

Virtualization Icon: www.flaticon.com/free-icon/digital_7035582?term=virtualization&page=1&position=29&origin=search&related_id=7035582

CD Icon: www.flaticon.com/free-icon/dvd_3440089?term=cd&page=1&position=50&origin=search&related_id=3440089

[Online; accessed Sep. 25, 2023]

4-18. Data Security and Resiliency (Freepik, 2023)

www.freepik.com/free-photo/standard-quality-control-collage-concept_30589259.htm#query=data%20security&position=4&from_view=search&track=ais

[Online; accessed Sep. 25, 2023]

REFERENCES

4-20. Data Protection Principles (Icons from Flaticon, 2023)

Data Security Icon: www.flaticon.com/free-icon/secure-data_1035086?term=data+security&page=1&position=2&origin=search&related_id=1035086

Fingerprint Icon: www.flaticon.com/free-icon/fingerprint-scan_6692271?term=finger+print&page=1&position=1&origin=search&related_id=6692271

Internet Icon: www.flaticon.com/free-icon/click_2807258?k=1695709362602&sign-up=google
[Online; accessed Sep. 26, 2023]

4-23. Homogeneous Storage Pool (Icons from Flaticon, 2023)

CD Icon: www.flaticon.com/free-icon/dvd_3440089?term=cd&page=1&position=50&origin=search&related_id=3440089
[Online; accessed Sep. 25, 2023]

4-24. Heterogeneous Storage Pool (Icons from Flaticon, 2023)

CD Icon: www.flaticon.com/free-icon/dvd_3440089?term=cd&page=1&position=50&origin=search&related_id=3440089
[Online; accessed Sep. 25, 2023]

5-1. DC Network (Victor217, 2023)

www.freepik.com/free-photo/network-switch-with-cables_902013.htm#query=server%20cables&position=0&from_view=search&track=ais
[Online; accessed Sep. 27, 2023]

5-2. Ethernet Cable (jannoon028, 2023)

www.freepik.com/free-photo/carrying-noise-category-white-plastic_1103234.htm#query=internet%20cable&position=37&from_view=search&track=ais
[Online; accessed Sep. 27, 2023]

5-3. Fiber-Optic Cable (PawinG, 2017)

<https://pixabay.com/photos/networking-fiber-optics-2633600>
[Online; accessed Sep. 27, 2023]

- 5-6. Enterprise-Class Router** (Unsplash – Albert Stoynov, 2023)
<https://unsplash.com/pt-br/fotografias/um-close-up-de-uma-rede-com-fios-conectados-a-ela-dyUp7WPu5q4>
[Online; accessed Sep. 27, 2023]
- 5-9. Top-of-Rack Network Topology** (Icons from Flaticon, 2023)
Server Icon: www.flaticon.com/free-icon/server_1792493?related_id=1792493
[Online; accessed Sep. 28, 2023]
- 5-13. Network Management Console** (Macrovector, 2023)
www.freepik.com/free-vector/datacenter-cartoon-composition-with-indoor-view-data-analysts-working-place-with-desktop-computer-servers-illustration_21744328.htm#query=computer%20server%20room&position=1&from_view=search&track=ais
[Online; accessed Sep. 28, 2023]
- 6-1. Classic Tape Drive for IBM Mainframe Computer**
(NationalArchivesCatalog, 1969)
<https://catalog.archives.gov/id/532417>
[Online; accessed Sep. 28, 2023]
- 6-2. Storage Tape Library** (Cory Doctorow, 2008)
www.flickr.com/photos/doctorow/2711081060
[Online; accessed Sep. 28, 2023]
- 7-1. Computer Security** (Freepik, 2023)
www.freepik.com/free-photo/standard-quality-control-collage-concept_30589259.htm#query=computer%20security&position=1&from_view=search&track=ais
[Online; accessed Sep. 28, 2023]
- 7-2. Computer Security Threats** (Flaticon, 2023)
Dial Icon: www.flaticon.com/free-icon/dial_3849780?term=dialing&page=1&position=6&origin=search&related_id=3849780

REFERENCES

Back Door Icon: www.flaticon.com/free-icon/backdoor_2008915?term=backdoor&page=1&position=10&origin=search&related_id=2008915

Browser Icon: www.flaticon.com/free-icon/web-search-engine_3003388?term=browser+search&page=1&position=3&origin=search&related_id=3003388

Bug Icon: www.flaticon.com/free-icon/bug_576509?term=computer+bug&page=1&position=6&origin=search&related_id=576509

Spyware Icon: www.flaticon.com/free-icon/spyware_1059644?term=spy&page=1&position=18&origin=search&related_id=1059644

Trojan Horse Icon: www.flaticon.com/free-icon/trojan-horse_6916900?term=computer+backdoor&page=1&position=3&origin=search&related_id=6916900

Malware Icon: www.flaticon.com/free-icon/malware_556261?term=malware&page=1&position=2&origin=search&related_id=556261

Malicious Software Icon: www.flaticon.com/free-icon/content-management-system_8278793?term=malicious+software&page=1&position=2&origin=search&related_id=8278793

Email Virus Icon: www.flaticon.com/free-icon/email_11891962?term=malicious+software&page=1&position=1&origin=search&related_id=11891962

Laptop Icon: www.flaticon.com/free-icon/laptop_689355?term=laptop&page=1&position=14&origin=search&related_id=689355

[Online; accessed Sep. 28, 2023]

7-3. Sources of Cybersecurity Threats (Icons from Flaticon, 2023)

Criminal Group Icon: www.flaticon.com/free-icon/group_10304099?term=criminal+group&page=1&position=3&origin=search&related_id=10304099

Detective Icon: www.flaticon.com/free-icon/detective_1320570?term=spy&page=1&position=10&origin=search&related_id=1320570

Nationalist Icon: www.flaticon.com/free-icon/soldier_6496383?term=salute&page=1&position=4&origin=search&related_id=6496383

Terrorist Icon: www.flaticon.com/free-icon/terrorist_2099819?term=terrorist&page=1&position=29&origin=search&related_id=2099819

Hacktivists Icon: www.flaticon.com/free-icon/movement_11895488?term=hacktivists&page=1&position=8&origin=search&related_id=11895488

Malware Icon: www.flaticon.com/free-icon/malware_556261?term=malware&page=1&position=2&origin=search&related_id=556261

Hacker Icon: www.flaticon.com/free-icon/hacker_924915?term=hacker&page=1&position=1&origin=search&related_id=924915

[Online; accessed Sep. 29, 2023]

7-4. Security Operation Center (Icons from Flaticon, 2023)

Tech Icon: www.flaticon.com/free-icon/touch_1835955?term=tech&page=1&position=9&origin=search&related_id=1835955

People Icon: www.flaticon.com/free-icon/users-group_32441?term=people&page=1&position=39&origin=search&related_id=32441

Gears Icon: www.flaticon.com/free-icon/settings-gears_60473?term=gears&page=1&position=1&origin=search&related_id=60473

[Online; accessed Sep. 29, 2023]

7-5. Trusted Platform Module (Icons from Flaticon, 2023)

Padlock Icon: www.flaticon.com/free-icon/padlock_3257787?term=password&page=1&position=3&origin=search&related_id=3257787

Key Icon: www.flaticon.com/free-icon/key_9932807?term=key&page=2&position=32&origin=search&related_id=9932807

TPM Chip Icon: www.flaticon.com/free-icon/cpu_597874?term=chip&page=1&position=4&origin=search&related_id=597874

Password Icon: www.flaticon.com/free-icon/password_482624?term=password&page=1&position=22&origin=search&related_id=482624

Diploma Icon: www.flaticon.com/free-icon/diploma_1916131?term=diploma&page=1&position=2&origin=search&related_id=1916131

[Online; accessed Sep. 29, 2023]

REFERENCES

7-6. Access Control by Mobile Phones (Freepik, 2023)

www.freepik.com/free-photo/cropped-hands-placing-finger-identification-spot-touchscreen_5698878.htm#query=smartphone%20fingerprint&position=1&from_view=search&track=ais

[Online; accessed Sep. 29, 2023]

7-7. Example of an RBAC (Icons from Flaticon, 2023)

Man and Woman Icon: www.flaticon.com/free-icon/toilet_828863?term=man+and+woman&page=1&position=7&origin=search&related_id=828863

[Online; accessed Sep. 29, 2023]

8-1. Help Desk Call Center (Freepik, 2023)

www.freepik.com/free-photo/people-working-call-center_22896178.htm#query=call%20center&position=0&from_view=search&track=ais

[Online; accessed Oct. 2, 2023]

8-2. Example of an IVR Call Flow (Icons from Flaticon, 2023)

Call Center Person Icon: www.flaticon.com/free-icon/support_1067566?term=people+calling&page=1&position=2&origin=search&related_id=1067566

Customer Calling Icon: www.flaticon.com/free-icon/woman_1841251?term=people+calling&page=1&position=89&origin=search&related_id=1841251

Microphone Icon: www.flaticon.com/free-icon/microphone_709682?term=recording&page=1&position=3&origin=search&related_id=709682

Phone Call Icon: www.flaticon.com/free-icon/phone-call_597177?term=phone&page=1&position=4&origin=search&related_id=597177

[Online; accessed Oct. 2, 2023]

8-3. Trouble Ticketing Systems (rawpixel.com, 2023)

www.freepik.com/free-vector/illustration-avatar-social-network-concept_2945077.htm#query=computer%20forum&position=6&from_view=search&track=ais

[Online; accessed Oct. 2, 2023]

8-4. Help Desk Dashboard Representation (pikisuperstar, 2023)

www.freepik.com/free-vector/professional-dashboard-user-panel_5633137.htm#query=ticket%20dashboard&position=3&from_view=search&track=ais

[Online; accessed Oct. 2, 2023]

8-5. ITSM Service Desk Functions (Icons from Flaticon, 2023)

Clock Icon: www.flaticon.com/free-icon/wall-clock_833602?term=clock&page=1&position=3&origin=search&related_id=833602

Wrench Icon: www.flaticon.com/free-icon/wrench_2593066?term=wrench&page=1&position=9&origin=search&related_id=2593066

Customer Service Icon: www.flaticon.com/free-icon/customer-service_9759919?term=customer+management&page=1&position=1&origin=search&related_id=9759919

Settings Icon: www.flaticon.com/free-icon/settings_3524659?term=settings&page=1&position=4&origin=search&related_id=3524659

Books Icon: www.flaticon.com/free-icon/books-stack-of-three_29302?term=book&page=1&position=12&origin=search&related_id=29302

Email Icon: www.flaticon.com/free-icon/email_482138?term=email&page=1&position=11&origin=search&related_id=482138

[Online; accessed Oct. 2, 2023]

8-6. Edge Computing (Icons from Flaticon, 2023)

Thermometer Icon: www.flaticon.com/free-icon/thermometer_2100100?term=temperature&page=1&position=9&origin=search&related_id=2100100

REFERENCES

Server Icon: www.flaticon.com/free-icon/server_4144964?term=server&page=1&position=23&origin=search&related_id=4144964
Smartphone Icon: www.flaticon.com/free-icon/smartphone-call_15874?term=phone&page=1&position=7&origin=search&related_id=15874
Waves Icon: www.flaticon.com/free-icon/phone_11041856?term=phone+waves&page=1&position=19&origin=search&related_id=11041856
Antena Icon: www.flaticon.com/free-icon/antena_6861894?term=antena&page=1&position=2&origin=search&related_id=6861894
Setting Icon: www.flaticon.com/free-icon/setting_3019014?term=settings&page=1&position=8&origin=search&related_id=3019014
[Online; accessed Oct. 4, 2023]

Index

A

- Access control list (ACL), 60, 221, 233
- Access controls, 5, 147, 149, 184
- Account management, 183
- Administrative network
 - security, 185
- Amazon Web Services (AWS),
 - 3, 68, 102, 157
- Analytical workloads, 94
- Application-centric
 - approaches, 159
- Application programming
 - interfaces (API), 12, 77, 104
- Array-based storage, 136, 137
- Artificial intelligence (AI), 244
- ASHRAE, 20
- Attacker motivation, 210–212
- Attacks, motivations/impact
 - attackers, 210, 212
 - security breaches, 210
- Audit trails tracking system, 213
- Automated theorem, 212
- Autonomic computing, 72–73

B

- Backdoor, 207, 232
- Backup systems, 148, 266

- cloud-based applications,
 - 194, 195
 - disaster recovery
 - planning, 197–201
 - evolution, 191–194
 - private data centers/clouds, 194
 - recovery, 191
 - types, 195–197
- Big data, 195
- Basic Input/Output System (BIOS), 110
- Blade server, 46–47
- Block-based disk access, 134
- Bring your own devices (BYOD), 147
- Building management system (BMS), 16, 259
- Business Process as a service (BaaS), 9, 14

C

- Cables, 47
 - data connection, 163
 - DC switches, 167, 168
 - power distribution, 163
 - structured cabling, 164, 165
 - switches and routers, 165, 166
- Cache memory, 52, 109, 110

INDEX

- Centralized network topology,
 - 170, 171
- Central processing unit (CPU), 34,
 - 52, 57, 109
- Change management, 250, 254
- Charging structure, 103–105
- Chatbots, 243–245
- Classic SAN provisioning, 150
- Client-server computing, 1, 13
- Client systems, 70
- Cloud-based backup
 - systems, 194–195
- Cloud computing, 1, 3, 38, 84, 105
- Cloud data center, 4, 5, 14, 42,
 - 52, 68, 258
- Cloud operating systems, 68–71
- Cloud service provider (CSP), 3, 15,
 - 144, 180, 184, 235, 258
 - benefits, 96, 97
 - risks, 97
- Cloud services, 15, 102
 - back-end, 5
 - by IT CSP, 5
 - data storage, 4
 - front-end, 5
 - load balancing, 5
 - security, 5
 - virtualization, 4
 - See also* Data center (DC)
- Cloud storage, 144, 157–159,
 - 191, 266
- Cold standby, 88
- Common Internet File System (CIFS), 133
- Communication plans, 200, 202
- Communications as a service (CaaS), 6–9, 258
- Community model, 10–12, 258
- Compliance, 12, 20, 98, 161, 201
- Compute resources
 - pricing variables, 103
 - terminology, 102, 103
- Computer memory hierarchy, 109
- Computer security, 204
 - incident response, 224
 - threats, 209
- Computer software, 35, 36
- Computer system, 257, 260–262
 - hardware, 34–36
 - purpose
 - general-purpose, 38, 39
 - specialized, 40, 41
 - software, 35, 36
 - types, 36–38
 - vulnerabilities, 204–209
- Computer workload, 92
 - types, 92–95
- Compute security, 261
- Confidential data, 114, 207
- Configuration management, 184
- Container
 - benefits, 75, 76
 - hypervisor, 75
 - virtualization, 74
- Containerization, 261
- Container runtime, 82
- Content distribution networks (CDN), 84, 86

Control-data-plane interface
 (CDPI), 177, 178
 Control layer, 175
 Control plane, 76, 77, 90, 91, 101,
 261, 265
 Cooling, 24, 25, 37, 42, 44
 Core memory, 109
 Cost-saving strategies, 160
 Critical services, 230
 Customer service representative
 (CSR), 237
 Customer support representatives
 (CSRs), 269
 Cyberattacks, 198, 200, 224
 Cybercrimes, 210, 215
 Cyber hygiene, 223, 224, 233
 Cyber insurance, 269
 Cybersecurity, 200, 203, 210, 216,
 223, 232, 261, 267, 269
 planning, 226
 threats, 211

D

DaemonSet, 101
 Dashboards, 104
 Data as a Service (DaaS), 6, 258
 Data backup, 266
 Data backup/storage, 200
 Data breaches, 65, 198, 200
 Data center (DC), 15, 38, 259
 blade server, 46, 47
 cloud, 42
 components, 41
 cooling, 16
 facilities and infrastructures, 43
 facility design, 17, 18
 facility management, 16, 22
 infrastructure, 23, 24
 IT infrastructure, 16
 IT operations, 42
 power, 16
 demands, 26, 27
 UPS, 27–29
 rackmount server, 43–46
 security, 17, 30
 space, 15
 lighting, 22
 noise, 22
 physical space, 21
 weight, 22
 standards established, 18–21
 temperatures, 24
 Data connection, 163
 Data deletion, 145, 191,
 263, 266
 Data encryption, 144, 149
 Data integrity, 6, 58, 115, 195
 Data layer, 175
 Data loss prevention (DLP), 147
 Data movement, 157, 159, 160
 Data protection, 96, 97, 124, 142,
 143, 149, 202
 Data redundancy, 148
 Data residency, 96
 Data security, 263
 goals, 142
 methods, 141

INDEX

- Data security (*cont.*)
 - principles, data protection,
 - 142, 143
 - resiliency, 140
 - threats, 143
 - vulnerabilities, 144, 145
- DC networking, 162
 - administration, 180–184
 - components, 162–170
 - load balancers and analytic tools, 161
 - network resiliency, 178
 - provisioning, 179, 180
 - resilient network, 184, 185
 - resilient network architecture,
 - 186, 187
 - server and storage technology, 161
 - tier, domain and territory, 188
 - topology, 170–178
 - workload life cycle, 161
- DC switches
 - capacities, 167
 - connection speed, 167
 - functions, 167
 - network management, 168
 - ports, 168
 - single-tier and multi-tier configurations, 167
- Deduplication, 159, 193
- Default secure settings, 213, 233
- Denial-of-service attacks (DoS),
 - 204, 205
- Deployment models, 10
 - community model, 12
 - hybrid models, 12, 13
 - private, 11
 - public model, 11
- Direct-access attacks, 207, 208
- Direct attached storage (DAS),
 - 116–119, 228, 262
- Disabling USB ports, 219
- Disaster mitigation, 268
- Disaster recovery, 114, 196, 200
- Disaster recovery plan, 191,
 - 197–199, 228
- Disaster recovery planning, 191
 - considerations, 199, 201
 - IT management and organizational resilience, 197
 - natural events, 198
 - steps
 - business impact analysis, 198
 - execution, 199
 - monitoring and review, 199
 - plan documentation, 199
 - risk assessment, 198
 - strategy development, 198
 - testing and training, 199
 - technology, 201
- Disconnecting unused devices, 219
- Disk arrays, 123, 262, 263
 - components, 122
 - RAID, 123–132
 - redundancy, 122

Disk crashes, 148
 Disk technology, 112
 Distributed denial-of-service
 (DDoS), 205, 267
 Distributed operating
 systems (DOS), 65–69
 Drive locks, 219
 Dynamic workloads, 92, 93

E

Eavesdropping, 206–207
 Edge computing, 83–86
 Edge computing support, 252–255
 Elasticity, 68
 Email, 241, 242
 Embedded operating
 systems, 52, 61
 Embedded systems, 37
 EN 50600 series, 20
 EN 50600-2-6, 21
 Encryption, 5, 6, 40, 144, 147, 184
 Enterprise-class router, 169
 Ethernet cable, 163
 External attributable
 disasters, 178
 External disk configurations
 device types, 121
 disk arrays, 122–132
 JBOD, 121
 External storage systems
 architectures, 116
 DAS, 117
 NAS, 118

SAN, 119, 121
 types, 116
 External threats, 143

F

Failover systems, 200
 FAST-VP, 153
 Fault management, 178, 183, 188
 Fault-tolerant computer
 system, 87–89
 Fiber-optic cable, 163, 164
 Fibre channel protocol (FCP), 119
 File-based storage
 virtualization, 133
 Floating-point operations per
 second (FLOPS), 48, 49

G

Gateways, 169, 265
 Gramm-Leach-Bliley Act, 20
 Graphical user interfaces
 (GUI), 55, 60

H

Hard disk drives (HDDs),
 59, 111–113
 Hardware, 34–36
 Hardware failures, 67, 142, 194,
 198, 232
 Hardware protection, 141, 218–220,
 233, 263

INDEX

Hardware virtualization
 compute-server resources, 70
 types, 71, 72
 virtual machines, 71

Heating, ventilation, and air
 conditioning (HVAC), 20

Help Desk, 269
 benefits, 241
 call center, 235
 calling, 237
 chatbots, 243, 244
 communication
 method, 245
 drawbacks, 241
 email, 241, 242
 functions, 236

IVR systems
 benefits, 238, 239
 call flow, 238
 drawbacks, 239, 240
 functions, 238
 online chat, 242, 243
 trouble ticketing
 systems, 245–249

Heterogeneous storage pool,
 153, 154

High-performance computing
 (HPC), 48–53, 66, 93, 260

Homogeneous storage pool, 153

Host-based storage virtualization,
 136, 138

Hot standby, 89

Hybrid models, 12, 13, 258

Hybrid RAID, 129

Hypervisor, 70, 71, 75, 91, 151,
 261, 264

Hypervisor/container-based
 control plane, 91, 261

I

In-band virtualization, *see*
 Symmetric virtualization

Incident management, 250,
 253, 255

Incident response, 224, 225, 234

In-depth defense, 213, 233

Information Technology
 Infrastructure Library
 (ITIL), 249, 250, 270

Infrastructure as service (IaaS),
 6, 7, 258

Instant messaging (IM), 7

Integrated circuits (IC), 109

Integrated Service Management
 (ITSM), 250

Interactive voice response (IVR),
 237, 248, 255

Internal/just memory, 109

Internal threats, 143

Internet computing, 2

Internet service provider (ISP),
 169, 231

Intrusion detection system (IDS),
 215, 219

Inventory management, 254

ISO 9000, 20

ISO 14000, 20

ISO 28001, 20
 IT security systems, 257

J

JBOD configurations, 262
 Just a Bunch of Disks (JBOD), 121

K

Key performance indicators
 (KPI), 251
 Knowledge management,
 254, 255
 Kubelet, 82
 Kube-proxy, 82
 Kubernetes (K8s), 76
 API server, 79, 80
 cloud controller manager, 81
 cloud provider
 dependencies, 81
 controller manager, 80
 control plane, 77
 key value store, 77
 node components, 82–84
 nodes, 78
 pods, 78, 79
 types, 80
 workload management, 99–101

L

Least privilege, 212, 233
 Least recently used (LRU), 110

LINPACK Benchmark, 50, 51
 Linux operating system, 70
 Load balancer, 5, 81, 230
 Local area network (LAN), 170

M

Magnetic tapes, 191, 192
 Main distribution area (MDA),
 164, 165
 Mainframe computing, 1, 229
 Main memory, 109, 110
 Malware, 209
 Media access control (MAC),
 166, 206
 Mesh network topology, 172, 173
 Metadata, 113
 Million instructions per second
 (MIPS), 48
 Mobile access control, 219,
 220, 233
 Moore's Law, 39, 40, 105, 112
 Multi-tier network
 topology, 173–174, 265
 Multivector attacks, 208

N

Natural language processing
 (NLP), 244
 Nested RAID levels
 RAID 01, 129, 130
 RAID 10, 130, 131
 RAID 50, 131, 132

INDEX

Network administration
 account management, 183
 configuration management, 184
 fault management, 183
 monitoring, 181–183
 performance management, 184
 security management, 183
Network-attached storage (NAS),
 116–118, 262
Network-based storage, 136, 137
Networked storage systems, 192
Network fabric architecture, 172
Network fault management,
 178, 265
Network File System (NFS), 133
Networking, 264, 265
Network management, 168, 182
Network monitoring, 181–183
Network provisioning, 179, 180
Network topologies, 265
Nodes, 78
Northbound interface (NBI), 177

O

Offline storage, 114, 115
Online chat, 242, 243
OpenFlow protocol, 175
Open Networking Foundation
 (ONF), 175, 265
Operating system (OS), 35, 260
 building blocks, 53, 54
 functions
 booting, 56

 command interface, 60
 driver management, 57, 58
 file system, 58, 59
 I/O system, 57, 58
 job control and
 accounting, 59
 memory management, 57
 processor management, 57
 resource monitoring, 59
 security, 60
 history, 53–56
 key functions, 52, 53
 types, 61–71

OSI seven-layer communication
 model, 166

OSI standard model, 265

Out-of-band virtualization, *see*
 Asymmetric virtualization

P

Parallel computing, 89
Payment Card Industry Data
 Security Standard, 21
Performance management,
 184, 188
Personal systems, 36
Phishing, 205
Physical network security, 185
Physical security, 144
Platform as a service (PaaS), 8, 258
Pods, 78, 79
Polymorphic attacks, 208
Power consumption, 47

Power demands of DC, 26, 27

Power distribution, 163

Primary storage, 109–111

Private clouds, 194, 203, 258

benefits, 98, 99

challenges, 99, 100

Private data centers, 5, 154, 194, 258

Private model, 10–12, 258

Privilege escalation, 208, 232

Product-attributable service

outages, 178

Public model, 258

Q

Quality of service (QoS), 7

R

Rackmount server, 43

convenience, 44

cooling, 44

DC Racks, 45

powerful, 44

Racks, dc, 23, 24

Rack server, *see* Rackmount server

RAID 0, 124, 125

RAID 1, 125, 126, 129, 130

RAID 2, 126

RAID 3, 127

RAID 4, 127

RAID 5, 127–129

RAID 6, 129

RAID 10, 130, 131

RAID 50, 131, 132

Random-access memory
(RAM), 109

Read-only memory (ROM), 110

Real-time operating systems
(RTOS), 62–64

Real-time workloads, 94, 95

Recovery, 148, 191, 266

Recovery point objectives
(RPOs), 200

Recovery time objectives
(RTOs), 200

Redundancy, 200, 231

Redundant array of independent
disks (RAID), 262

definition, 123

levels, 124–132

Regulations, 201

Remote backup, 266

Remote DC support, 252

Reporting, 254

Resilience, 186, 267–269

best practices

access controls, 147

access security, 146

backup and recovery, 148

DLP, 147

encryption, 147

endpoint security, 147

redundancy, 148

security policies, 149

building up, 231

critical services, 230

data breaches, 140

INDEX

Resilience (*cont.*)

- data center, 230
- data protection, 142
- data security, 142–145
- definition, 227
- factors, 186, 187
- improvement, 231, 232
- physical hardware protection
 - and security software, 141
- processes, 146, 229–232
- security, 227–230
- storage architecture, 141
- storage data security, 141
- See also* Resilient network

Resiliency, 69

- definition, 87
- fault-tolerant computer
 - system, 87–89
- federated architecture, 91

Resilient computing, 89–91, 261

Resilient cybersecurity

- strategy, 267

Resilient IT infrastructure, 258

Resilient network

- architecture, 186, 187
- benefits, 184
- security threats, 185

Resource identification, 200

Resources

- computer workload, 92

Reverse engineering, 208

Role-based access control (RBAC), 179, 221, 233

Routers, 165, 166, 168, 265

S

Sarbanes-Oxley Act, 20

SAS 80 Type I/II, 20

Scalability, 68

Secondary storage, 111–113

Security, 30, 69

Security architecture, 214

Security breaches, 210

Security by design, 267

ACL, 221

architecture, 214

cyber hygiene, 223, 224

cybersecurity planning, 226

definition, 212

hardware protection, 218–220

incident response, 224, 225

infrastructure, 214, 215

principles, 212, 213

tools, 222

training, 223

vulnerabilities reduction, 217

vulnerability assessment and management, 215, 216

Security considerations, 200

Security infrastructure, 214, 215

Security management, 183

Security Operation Center, 216

Security policies, 149

Security threats, 185

Security tools, 222

Security training, 223, 233

Self-service, 254

Server, 37, 260

- Server Message Block (SMB), 133
- Server operating systems, 64–66
- Service Desk, 249–252, 270
- Service interruptions, 261
- Service-level agreements (SLAs), 249, 270
- Service options, 9
- Service provisioning, 180
- Service request management, 253
- Simple Network Management Protocol (SNMP), 183
- Social engineering, 208
- Social media, 257
- Software, 35, 36
- Software as a service (SaaS), 8
- Software-defined networking (SDN), 161, 265
 - architecture, 176
 - application, 177
 - control plane, 177
 - data plane, 177, 178
 - cloud architectures, 175
 - data transport, 175
- Software virtualization, 230
- Solid-state drives (SSDs), 111, 113
- Special-purpose devices, 33
- Spoofing attack, 205, 206
- Static workload, 92
- Storage allocation tiering, 155
- Storage area network (SAN), 119, 121, 170
- Storage Management Initiative Specification (SMI-S), 138
- Storage network, *see* Storage area network (SAN)
- Storage Networking Industry Association, 123, 263
- Storage pools, 152–154
- Storage provisioning, 264
 - allocation tiering, 155
 - classic SAN systems, 150
 - defined, 150
 - modern DCs, 151, 152
 - private and public infrastructure, 155, 156
 - storage pools, 152–154
- Storage resources, 156–159
 - comparison, 115
 - offline storage, 114, 115
 - primary storage, 109–111
 - secondary storage, 111–113
 - tertiary storage, 114
- Storage security, 145, 263
 - definition, 139
 - implementation, 149, 150
 - principles
 - availability, 145
 - confidentiality, 145
 - integrity, 145
 - provisioning and administration, 150–155
 - resources, charging for, 156–159
 - See also* Resilience
- Storage systems, 109, 121, 262–264
 - definition, 107
 - external, 116–121
 - granularity, 108

INDEX

Storage systems (*cont.*)

- latency, 108
- reliability, 108
- security and resilience, 139–150
- throughput, 108
- virtualization, 132–139
 - See also* External disk configurations; Storage resources

Storage Tape Library, 193

Storage virtualization, 263

- access modes, 133, 134
- architecture, 137, 138
- benefits, 138, 139
- functionality, 133
- methods, 136, 137
- RAID systems, 133
- symmetric/asymmetric configurations, 135
- types, 134

Stripe set, 124

Structured cabling, 164, 165, 264

Supercomputer, 48–53

Switches, 165, 166, 264

Symmetric virtualization, 135

T

Technical network security, 185

Techopedia, 214

Telecommunications Industry

- Association (TIA), 164

Temperatures, 24

Tertiary storage, 114

Theory of constraints, 215

The Uptime Institute Tier

- Standard, 18–20

Thick provisioning, 151, 152, 158, 159, 264

Thin provisioning, 151, 152, 158, 159, 264

Tier 4 data center network

- architecture, 187

Top-of-rack (ToR) topology, 172

Trouble ticketing systems, 246

- costs and service level, 249
- customer contacts, 246
- customer issues/requests, 245
- dashboard, 246–249

Trusted Platform

- Module (TPM), 218

U

Uninterruptible power

- supply (UPS), 17

- battery backups, 29
- continuous power supply, 28
- monitoring and alerting, 29
- power conditioning, 29

USB dongles, 220

User provisioning, 179

Utility computing, 2

V

Virtualization, 4, 260, 264

- hardware, 70–72

Virtual machine (VM), 260
Virtual private networks (VPN),
 169, 220
Virtual storage, 264
Virtual tape libraries (VTLs), 192
Voice communication system, 165
Voice over IP (VoIP), 7
Vulnerability
 assessment and management,
 215, 216
 cloud storage, 144
 computer systems
 backdoor, 207
 direct-access attacks,
 207, 208
 DoS, 204, 205
 eavesdropping, 206, 207
 malware, 209
 multivector and
 polymorphic attacks, 208
 phishing, 205

 privilege escalation, 208
 reverse engineering, 208
 social engineering, 208
 spoofing attack, 205, 206
data deletion, 145
data encryption, 144
physical security, 144
reduction, 217

W

Workload deployment, 95
Workloads, 158, 230, 260

X, Y

X as a Service (XaaS), 9, 10

Z

Zoned network topology, 171, 172