

Cloud Computing

Definitions

- 1 [Introduction](#)
- 2 [Definitions](#)
- 3 [Characteristics](#)
- 4 [Service Models](#)
- 5 [Deployment Models](#)
- 6 [Virtualization and Elasticity](#)
- 7 [Typical Cloud Services](#)
 - [Data Storage in the Cloud](#)
 - [Communications: Publish/Subscribe](#)
 - [Batch Processing: Map/Reduce](#)
 - [Serverless Computing / Function as a Service](#)
- 8 [Edge Computing](#)

Define Cloud Computing

Distributed Computing

Distributed System: Definition ([Wikipedia](#))

*A system whose **components** are located on **different networked computers**, which then **communicate** and **coordinate their actions** by **passing messages** to one another.*

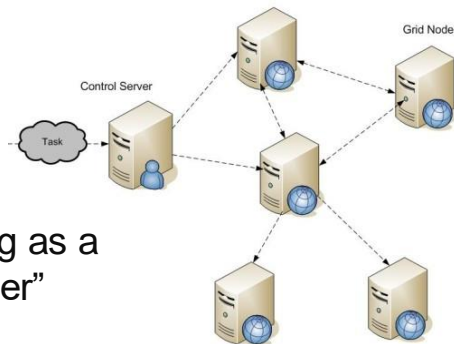
- Very broad! Different models are possible:
 - Centralized
 - Peer-to-peer/"mesh"
 - Hybrid

Grid Computing

7

Grid Computing: Definition ([Wikipedia](#))

*A **combination** of **computer resources** from multiple administrative domains applied to a **common task**.*



effectively functioning as a
“virtual supercomputer”

Utility Computing

Utility Computing: Definition (Wikipedia)

*The packaging of **computing resources** (computation, storage etc.) **as a metered service** similar to a traditional public utility.*

Users only pay for what they use

Cloud Computing?

Cloud Computing?

9

- Grid Computing + Utility Computing?
- Very hard to define – can mean so many different things to different parties!
- Many definitions

Cloud Computing: Definition (NIST)

*Cloud computing is a model for enabling ubiquitous, convenient, on-demand network access to a shared pool of configurable computing resources (e.g., networks, servers, storage, applications, and services) that can be rapidly provisioned and released with minimal management effort or service provider interaction. This cloud model is composed of 5 essential **characteristics**, 3 **service models**, and 4 **deployment models**.*

Cloud Characteristics

10

- 1 [Introduction](#)
- 2 [Definitions](#)
- 3 [Characteristics](#)
- 4 [Service Models](#)
- 5 [Deployment Models](#)
- 6 [Virtualization and Elasticity](#)
- 7 [Typical Cloud Services](#)
 - [Data Storage in the Cloud](#)
 - [Communications: Publish/Subscribe](#)
 - [Batch Processing: Map/Reduce](#)
 - [Serverless Computing / Function as a Service](#)
- 8 [Edge Computing](#)

Cloud Characteristics (1)

11

1. On-demand Self Service

- Ability to provision computing capabilities without intervention
 - Computation ("aka machine") time
 - Storage

Cloud Characteristics (1)

11

1. On-demand Self Service

- Ability to provision computing capabilities without intervention
 - Computation ("aka machine") time
 - Storage

2. Broad network access

- Capabilities available over the network
- Accessible by *thin* and *thick* clients (e.g., web, desktop/laptops, mobile devices, etc.)

Cloud Characteristics (2)

12

3. Resource pooling

- **Multi-tenancy:** the same cloud infrastructure can serve multiple customers, host multiple VMs, applications
- Computing resources are **pooled** (grouped as a single service)
 - Storage
 - Processing
 - Memory
 - Network
- *Physical and logical* resources are **dynamically** assigned and reassigned according to consumer **demand**
- **Location independence**
 - Precise location of the resources irrelevant
 - Only a general idea (e.g., Amazon EC2 US-east)

Cloud Characteristics (3)

13

4. Rapid elasticity

- *Elastic* provisioning – scaling up and down
- Can be done automatically
- To consumers: pool of resources might appear to be infinite

Cloud Characteristics (3)

13

5. Measured service

- *Metering* of the different resources
 - CPU (e.g., \$/CPU time in ms)
 - Network bandwidth (e.g., \$/gb)
 - Processing (e.g., \$/X requests)
 - Storage (e.g., \$/gb)
- Monitoring, controlling, reporting
- Full transparency for cloud operator and consumer

Service Models

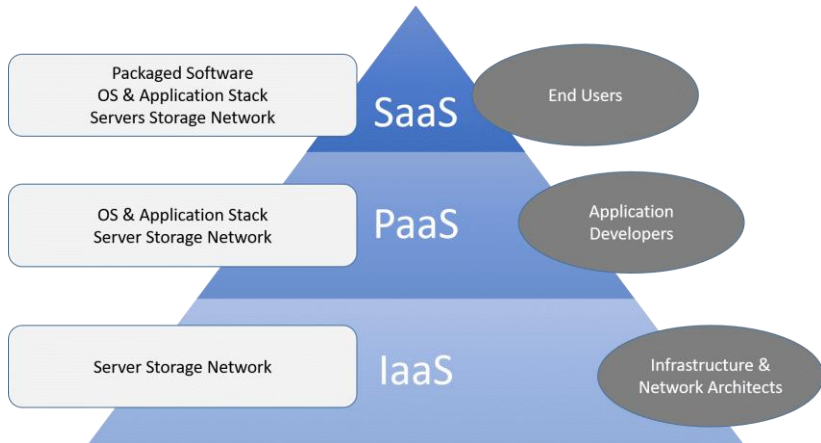
14

- 1 [Introduction](#)
- 2 [Definitions](#)
- 3 [Characteristics](#)
- 4 [Service Models](#)
- 5 [Deployment Models](#)
- 6 [Virtualization and Elasticity](#)
- 7 [Typical Cloud Services](#)
 - [Data Storage in the Cloud](#)
 - [Communications: Publish/Subscribe](#)
 - [Batch Processing: Map/Reduce](#)
 - [Serverless Computing / Function as a Service](#)
- 8 [Edge Computing](#)

Hierarchy of Service Models (Source)

15

Cloud Service Models



Infrastructure as a Service (IaaS) (1)

16

- Consumer can provision virtualized computing resources (aka VMs)
 - Processing, storage, network, GPU
- Can include OS and applications, or be bare metal
 - Example: Amazon EC2, Azure
- Consumer doesn't manage the hardware (physical or virtualized)
 - But has control over the OS, storage, applications, and limited network settings
 - e.g., firewall, port redirection, VLANs, etc

Infrastructure as a Service (IaaS) (2)

17

Virtualizing a machine implies that all of its components must be virtualized as well!

- Processor (CPU): virtualized from “real”, physical CPUs
 - Hardware acceleration is available on most recent processors
 - Support for multiple cores
 - “Standardized” metrics for modelling the performance of CPUs (e.g., Amazon vCPU)
 - ...
- Memory
- Storage
- Networking: SDNs
 - Network configuration is defined by software and not purely by the hardware (routers, switches)
 - Bandwidth, firewalls, subnets, etc.
- GPUs and other devices

The components of a VM are not all necessarily located on the same physical machine!

Data storage in the cloud

18

■ File storage:

- Common/typical storage abstractions (file systems, folders, files, etc.)
- Emulation of a “local” hard disk, but provided over the network
- Protocols: NFS, Google filesystem, etc.

Data storage in the cloud

18

■ File storage:

- Common/typical storage abstractions (file systems, folders, files, etc.)
- Emulation of a “local” hard disk, but provided over the network
- Protocols: NFS, Google filesystem, etc.

■ Object storage:

- Storage of objects and metadata (BLOB)
- ID for each “object”
- Typically accessed thru standard access protocols (e.g., HTTP)
- Version control systems (VCS) (e.g., Git, SVN) make use object storage
- Different storage systems are available based on customer needs (costs, frequency of data reads/writes, throughput, latency, etc.)
- Replication, versioning, encryption, availability in several “zones”, etc.
- e.g., Amazon S3, Google Cloud Storage

Platform as a Service (PaaS)

19

- Enables the deployment, management and execution of consumer or acquired **applications** onto cloud infrastructure (IaaS)
 - For customers: alleviates the need for managing the applications on their own infrastructure
 - Can be written into a variety of languages
 - Using a variety of libraries, services, tools supported by the provider
 - e.g., Web apps (Heroku, Google App Engine), APIs, microservices
- No control over underlying cloud infrastructure (IaaS)!
- Control over deployed applications
- Might have limited control over configuration settings of the hosting environment (e.g., config files)

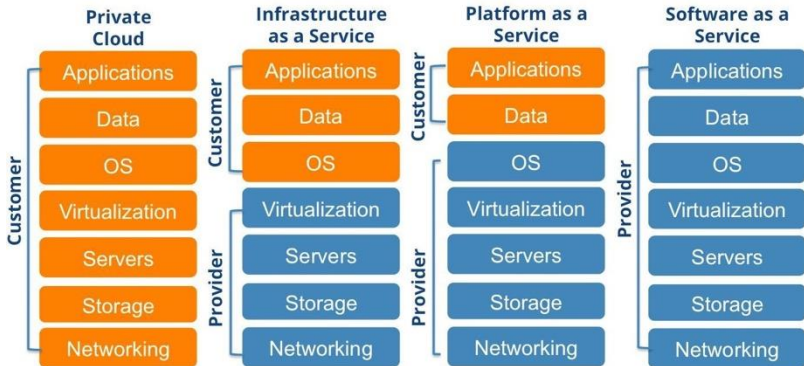
Software as a Service (SaaS)

20

- Use the provider's specific applications
 - Over the cloud provider's infrastructure (hardware + software/PaaS)
- Accessible from various clients
 - Thin & thick clients, mobile, web (e.g., web-based email)
- Consumer does not manage the underlying cloud infrastructure (network, servers, OS, storage, applications)
- Exception: limited user-specific application configuration settings (e.g., GMail settings)

Provisioning in Service Models

21



Exercise

To which service model (SAAS, PAAS, IAAS) does each of the following correspond?

- Editing a document online on Google Docs
- Testing a new Linux Kernel on an Amazon VM
- Accessing a MySQL database service
- Deploying a Python application
- Provisioning a virtual machine

Exercise

To which service model (SAAS, PAAS, IAAS) does each of the following correspond?

- Editing a document online on Google Docs **SaaS**
- Testing a new Linux Kernel on an Amazon VM **IaaS**
- Accessing a MySQL database service **PaaS**
- Deploying a Python application **PaaS**
- Provisioning a virtual machine **IaaS**

Deployment Models

22

- 1 [Introduction](#)
- 2 [Definitions](#)
- 3 [Characteristics](#)
- 4 [Service Models](#)
- 5 [Deployment Models](#)
- 6 [Virtualization and Elasticity](#)
- 7 [Typical Cloud Services](#)
 - [Data Storage in the Cloud](#)
 - [Communications: Publish/Subscribe](#)
 - [Batch Processing: Map/Reduce](#)
 - [Serverless Computing / Function as a Service](#)
- 8 [Edge Computing](#)

Public cloud

23

- Open use by general public
- Owned by business, academic, government organization, or a combination
- Exists on the premises the of cloud provider
- Example: Amazon, Google, MS Azure

Private cloud

24

- Exclusive use of a single organization with multiple “internal” consumers
 - e.g. different business units within a given organization
- Owned, managed, operated by organization, or a third-party, or a combination
- May exist on or off premises
- Example: a large company (e.g., Google Internal Cloud)

Community cloud

25

- Exclusive use of a specific community of consumers from **organizations with shared concerns**
 - Mission, security requirements, policy, compliance considerations
- Owned, managed, operated by one or more organizations in the community, a third party, or a combination of them
- May exist on or off premises
- Examples: Amazon Government Cloud, clouds that comply with BC data policies (e.g., UBC Workspace)

Hybrid cloud

26

- A composition of two or more distinct cloud infrastructure

Can include on-premise computing/storage/network infrastructure

Exercise

Choose the most accurate deployment model for each of the following (public cloud, private cloud, community cloud, hybrid cloud):

1. Due to strict privacy laws, UBC offers a storage service operated by the university and hosted on the university premises.
2. *Intel-ligent* operates a cloud service for its own internal needs – whenever the demand goes above the available capacity, the extra load is sent to the MS Azure cloud.
3. Your Android phone automatically stores the photos that you have taken onto the Google Photos service.
4. The healthcare network provides a cloud service that allows for storing patient information from multiple establishments (hospitals, clinics, pharmacies, etc.) scattered across the province. The information is stored on a layer of servers deployed/spread over the various establishments. The service respects various compliance policies pertaining to storing critical patient information.

Exercise

Choose the most accurate deployment model for each of the following (public cloud, private cloud, community cloud, hybrid cloud):

1. Due to strict privacy laws, UBC offers a storage service operated by the university and hosted on the university premises. **Private cloud or community cloud**
2. *Intel-ligent* operates a cloud service for its own internal needs – whenever the demands goes above the available capacity, the extra load is sent to the MS Azure cloud. **Hybrid cloud**
3. Your Android phone automatically stores the photos that you have taken onto the Google Photos service. **Public cloud**
4. The healthcare network provides a cloud service that allows for storing patient information from multiple establishments (hospitals, clinics, pharmacies, etc.) scattered across the province. The information is stored on a layer of servers deployed/spread over the various establishments. The service respects various compliance policies pertaining to storing critical patient information. **Community cloud**

Virtualization

27

- 1 [Introduction](#)
- 2 [Definitions](#)
- 3 [Characteristics](#)
- 4 [Service Models](#)
- 5 [Deployment Models](#)
- 6 [Virtualization and Elasticity](#)
- 7 [Typical Cloud Services](#)
 - [Data Storage in the Cloud](#)
 - [Communications: Publish/Subscribe](#)
 - [Batch Processing: Map/Reduce](#)
 - [Serverless Computing / Function as a Service](#)
- 8 [Edge Computing](#)

What is virtualization?

28

Decoupling the **physical resources (physical hardware)** into **virtual resources**

Why virtualize?

- Cloud provider might have heterogeneous hardware
- Offering a consistent configuration to the different customers of the cloud
 - CPU performance
 - Amount of memory
 - Storage
 - Network bandwidth
- Offering additional isolation (reliability)
- **Virtualization of resources happen at different levels based on the service model!**

Virtualization: IaaS (1)

29

Hardware-level virtualization (lowest level)

- 1 Physical machine \Rightarrow n *virtual* machines
- Hypervisor: VMWare, VirtualBox, MS HyperV, Xen, etc.
- Run over an OS or “bare-metal”
- Nowadays, virtualization is hardware-assisted: can run at near-native speeds

Virtualization (IaaS) (2)

30

Hardware-level virtualization (lowest level)

Virtualized Hardware

- CPU (modern CPUs support virtualization extensions - e.g., Intel VT, AMD-V)
- Memory: portions of the RAM of the host machine are reserved
- Storage: virtual hard drives and other I/O peripherals, data center storage
- Network: virtual network adapters, virtualized networks/subnets
- GPU: for specific (AI) applications

Virtualization: PaaS / SaaS

31

Virtualization of the **combined** resources of a pool of machines (VMs)

- Build over IaaS virtualization layer
- Processing power (CPU)
- Pools of memory
- Distributed data storage
- Virtualized networking and addressing

Elasticity

32

- Allocating a pool of resources from the provider in an “elastic” manner, according to current needs
- Resource allocation can be made directly according to user requirements (e.g., Amazon EC2 dashboard, or thru the command-line)
- Can be triggered by the needs of higher-level apps & services deployed over higher-level layers (e.g., PaaS)

Elasticity: IaaS

33

Two approaches to scalability:

- Vertical: more powerful hardware (limited)
- **Horizontal: partitioning / sharding**

Elasticity: IaaS (Infrastructure as a Service)

- Allocating new VM instances
- Deallocating instances which aren't needed anymore
- Allocating storage, RAM, network (or a specific network configuration), etc. (can be properties of the VMs)

Elasticity: PaaS

34

Elasticity: PaaS (Platform as a Service):

- Automatic provisioning of VM/physical resources (IaaS layer) to execute the PaaS application
- The elasticity of the application itself might or might not be done automatically

Elasticity: PaaS

34

Elasticity: PaaS (Platform as a Service):

- Automatic provisioning of VM/physical resources (IaaS layer) to execute the PaaS application
- The elasticity of the application itself might or might not be done automatically

Example: for a request-based application, the PaaS “execution layer” could provision / allocate enough resources as necessary from the IaaS layer to satisfy the current volume of requests

- The unit of measure for control & billing purposes can then be different (higher-level) compared to the billing metrics for the IaaS layer
- For instance, the customer can be billed by the number of requests or for the execution time allotted for handling the requests (as opposed to billing for the “raw” CPU usage, memory, etc. of the VM)

35

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡

Data Storage in the Cloud

36

- 1 [Introduction](#)
- 2 [Definitions](#)
- 3 [Characteristics](#)
- 4 [Service Models](#)
- 5 [Deployment Models](#)
- 6 [Virtualization and Elasticity](#)
- 7 [Typical Cloud Services](#)
 - [Data Storage in the Cloud](#)
 - [Communications: Publish/Subscribe](#)
 - [Batch Processing: Map/Reduce](#)
 - [Serverless Computing / Function as a Service](#)
- 8 [Edge Computing](#)

Data Storage in the Cloud

37

How can data be stored across different nodes?

Data Storage in the Cloud

37

How can data be stored across different nodes?

- Distributed File Systems
 - Google FS, Hadoop
 - Provides file-system like abstractions in a distributed manner
- Block Storage
 - Amazon S3 (storage of objects, can be files)
- Databases:
 - SQL
 - NoSQL (e.g., Key-value Stores, MongoDB, etc.)

Data Storage in the Cloud: Properties

38

- Scalability
- High availability
- Low latency
- Durability
- Fault tolerant
- Predictable costs

38

- Predictable costs

Partition tolerance

- Cloud storage systems often opt for **eventual consistency**

Communications: Publish/Subscribe

39

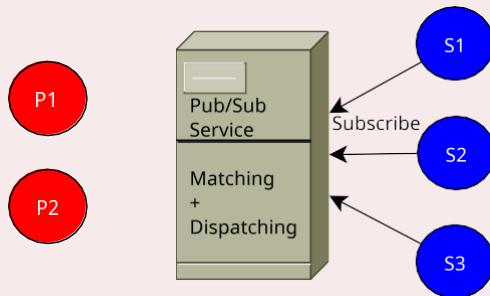
- 1 [Introduction](#)
- 2 [Definitions](#)
- 3 [Characteristics](#)
- 4 [Service Models](#)
- 5 [Deployment Models](#)
- 6 [Virtualization and Elasticity](#)
- 7 [Typical Cloud Services](#)
 - [Data Storage in the Cloud](#)
 - [Communications: Publish/Subscribe](#)
 - [Batch Processing: Map/Reduce](#)
 - [Serverless Computing / Function as a Service](#)
- 8 [Edge Computing](#)

Publish/Subscribe Paradigm

40

- Provides an elegant way to decouple content producers (publishers) from content consumers (subscribers)
- Publications are matched against subscriptions
- Many *flavours* of publish/subscribe

Publish/Subscribe - Example

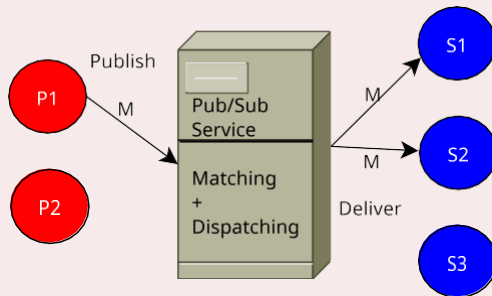


Publish/Subscribe Paradigm

40

- Provides an elegant way to decouple content producers (publishers) from content consumers (subscribers)
- Publications are matched against subscriptions
- Many *flavours* of publish/subscribe

Publish/Subscribe - Example



Topic-Based Publish/Subscribe

41

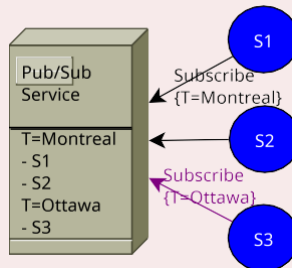
- Very common flavour of pub/sub
- Subscription language: a key (topic name)
- Publications tagged with a topic T , sent to all subscribers of T

Example - Weather Reports

Publishers



Subscribers

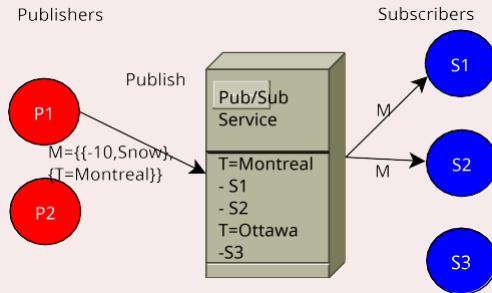


Topic-Based Publish/Subscribe

41

- Very common flavour of pub/sub
- Subscription language: a key (topic name)
- Publications tagged with a topic T , sent to all subscribers of T

Example - Weather Reports



Applications of Topic-Based Pub/Sub

42

Traffic alert systems



Weather alert systems



Mobile notif. frameworks



Social networks



IoT



Multiplayer Games



Desirable properties:

- Scalability & Elasticity
- Low Latency
- Reduced & Predictable Costs

Batch Processing: Map/Reduce

43

- 1 [Introduction](#)
- 2 [Definitions](#)
- 3 [Characteristics](#)
- 4 [Service Models](#)
- 5 [Deployment Models](#)
- 6 [Virtualization and Elasticity](#)
- 7 [Typical Cloud Services](#)
 - [Data Storage in the Cloud](#)
 - [Communications: Publish/Subscribe](#)
 - [Batch Processing: Map/Reduce](#)
 - [Serverless Computing / Function as a Service](#)
- 8 [Edge Computing](#)

44

- **grep**
- **wc (word count)**

grep with MapReduce

45



- Partitioning the files to be searched onto several nodes (map)
- Executing "grep" on each instance
- Partitioning the intermediate results to send them to a "reducer"
- Concatenation of all intermediate results

wc with MapReduce

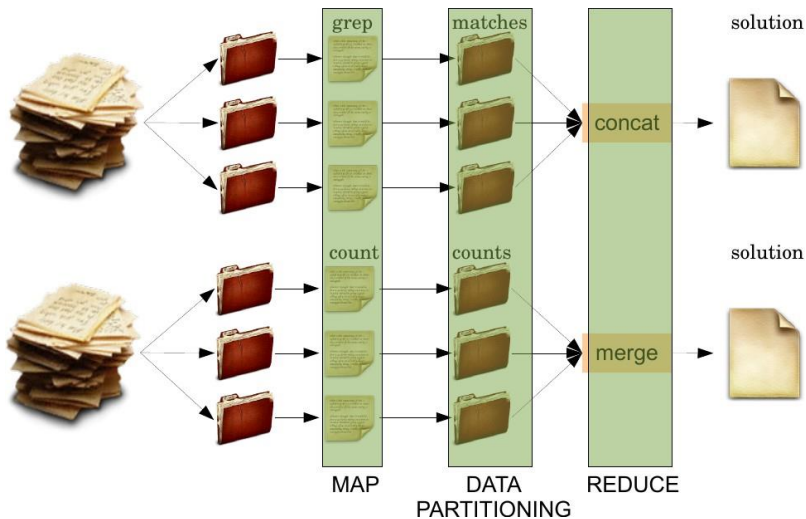
46



- Partitioning the files on which the wordcount should be performed onto several nodes (map)
- Executing "count" on each instance to compute the number of occurrences of each word
- Partitioning the intermediate "counts" to send them to a "reducer" (e.g., by hashing the words)
- Merging of all results (adding the partial counts for each word)

grep and wc with MapReduce

47



Serverless Computing / Function as a Service

48

- 1 [Introduction](#)
- 2 [Definitions](#)
- 3 [Characteristics](#)
- 4 [Service Models](#)
- 5 [Deployment Models](#)
- 6 [Virtualization and Elasticity](#)
- 7 [Typical Cloud Services](#)
 - [Data Storage in the Cloud](#)
 - [Communications: Publish/Subscribe](#)
 - [Batch Processing: Map/Reduce](#)
 - [Serverless Computing / Function as a Service](#)
- 8 [Edge Computing](#)

Function as a Service (FaaS)

49

- Application made of a set of functions
- Executed upon certain events being triggered
 - Web request
 - File upload
 - Change to DB
 - Timer
- Executed within containers (thin VMs)
 - Full isolation
 - FaaS functions are **stateless!**
 - Changes in state must be persisted to durable storage
- Example: Amazon Lambda, Google Cloud Functions, MS Azure Functions
- Some “functions” can be executed “at the edge”: e.g., lambda@edge

Edge Computing

50

- 1 [Introduction](#)
- 2 [Definitions](#)
- 3 [Characteristics](#)
- 4 [Service Models](#)
- 5 [Deployment Models](#)
- 6 [Virtualization and Elasticity](#)
- 7 [Typical Cloud Services](#)
 - [Data Storage in the Cloud](#)
 - [Communications: Publish/Subscribe](#)
 - [Batch Processing: Map/Reduce](#)
 - [Serverless Computing / Function as a Service](#)
- 8 [Edge Computing](#)

Edge Computing

51

Edge Computing, from a systems point of view, aims at processing the data as close as possible to where the data is **produced** and **consumed**.

The definition by itself is very vague, and can refer to various models that are described in the literature.

- For instance, for cloud providers, the “edge” can refer to smaller/micro cloud deployments that are more “localized”
- In other contexts, the “edge” can refer to processing the data onto the devices themselves that are becoming more and more powerful
 - For instance, Raspberry Pi devices run a full Linux OS

Edge Computing, why is it important?

52

The cloud is widespread and is typically “reachable” from anywhere – however, in some cases, it might not be the best option (in the IoT landscape). **Why?**

Edge Computing, why is it important?

52

The cloud is widespread and is typically “reachable” from anywhere – however, in some cases, it might not be the best option (in the IoT landscape). **Why?**

1 Limited connectivity

- For devices that are deployed in an isolated environment, or if the wireless communication is intermittent

Edge Computing, why is it important?

52

The cloud is widespread and is typically “reachable” from anywhere – however, in some cases, it might not be the best option (in the IoT landscape). **Why?**

1 Limited connectivity

- For devices that are deployed in an isolated environment, or if the wireless communication is intermittent

2 Lower latency

- Some critical apps must react very quickly to specific events – in that case, sending the data to the cloud (back and forth) can push the latency above a critical level
- e.g., smart vehicles

52

The cloud is widespread and is typically “reachable” from anywhere – however, in some cases, it might not be the best option (in the IoT landscape). **Why?**

- 1 Limited connectivity
 - For devices that are deployed in an isolated environment, or if the wireless communication is intermittent
- 2 Lower latency
 - Some critical apps must react very quickly to specific events – in that case, sending the data to the cloud (back and forth) can push the latency above a critical level
 - e.g., smart vehicles
- 3 Alleviating the dependence to a third party
 - Some IoT devices are too strongly coupled to a specific cloud-based service, and can stop working if the cloud service ceases to exist, or if the operator changes its terms of use.

Edge Computing, why is it important?

52

The cloud is widespread and is typically “reachable” from anywhere – however, in some cases, it might not be the best option (in the IoT landscape). **Why?**

1 Limited connectivity

- For devices that are deployed in an isolated environment, or if the wireless communication is intermittent

2 Lower latency

- Some critical apps must react very quickly to specific events – in that case, sending the data to the cloud (back and forth) can push the latency above a critical level
- e.g., smart vehicles

3 Alleviating the dependence to a third party

- Some IoT devices are too strongly coupled to a specific cloud-based service, and can stop working if the cloud service ceases to exist, or if the operator changes its terms of use.

4 Security requirements

- It might be preferable to process the data locally to meet the needs of security requirements and policies.

Edge Computing, why is it important?

52

The cloud is widespread and is typically “reachable” from anywhere – however, in some cases, it might not be the best option (in the IoT landscape). **Why?**

1 Limited connectivity

- For devices that are deployed in an isolated environment, or if the wireless communication is intermittent

2 Lower latency

- Some critical apps must react very quickly to specific events – in that case, sending the data to the cloud (back and forth) can push the latency above a critical level
- e.g., smart vehicles

3 Alleviating the dependence to a third party

- Some IoT devices are too strongly coupled to a specific cloud-based service, and can stop working if the cloud service ceases to exist, or if the operator changes its terms of use.

4 Security requirements

- It might be preferable to process the data locally to meet the needs of security requirements and policies.

5 Reducing the costs

- In a context in which a lot of data is produced and consumed, processing some of the data locally can reduce the volume of data that is sent and processed in the cloud.