

Data Integration and Large Scale Analysis

07- Cloud Computing Fundamentals

Dr. Lucas Iacono - 2025

Part B

Large-Scale Data Management & Analysis

- **Cloud Computing**
 - Fundamentals [Nov 21]
 - Resource Management and Scheduling [Nov 28]
 - Distributed Data Storage[Dec 05]

Part B

Large-Scale Data Management & Analysis

- **Large-Scale Data Analysis**
 - Distributed, Data-Parallel Computation [Dec 12]
 - Distributed Stream Processing [Dec 19]
 - Distributed Machine Learning Systems [Jan 16]

Agenda

- Motivation and Terminology
- Cloud Computing Service Models
- Cloud, Fog, and Edge Computing

Motivation and Terminology

Motivation and Terminology

- How new it is?

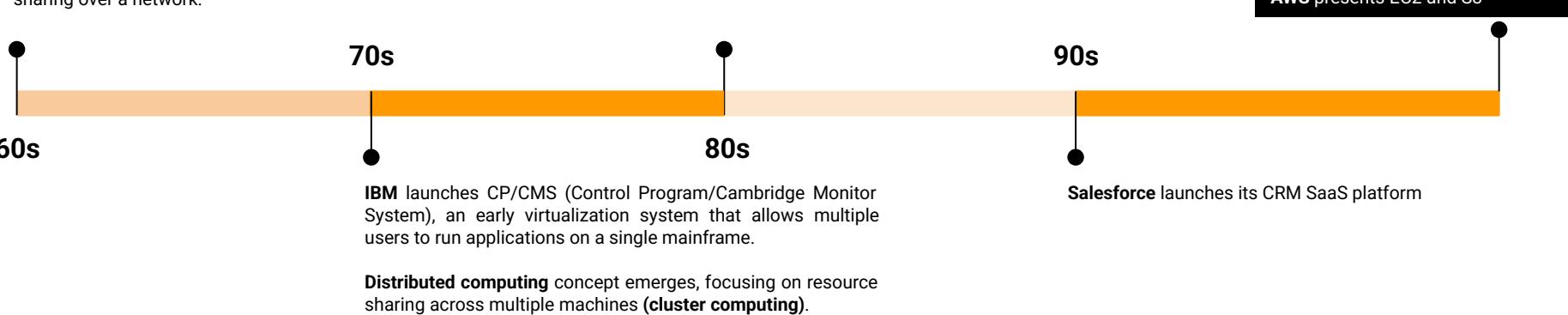
Motivation and Terminology

John McCarthy conceptualize the idea of computing as a **public utility** at a lecture at MIT, envisioning "computing on demand."

ARPANET is developed, enabling computational resources sharing over a network.

CompuServe starts offering small-scale cloud-like storage

Tim Berners-Lee invents the World Wide Web



Motivation and Terminology

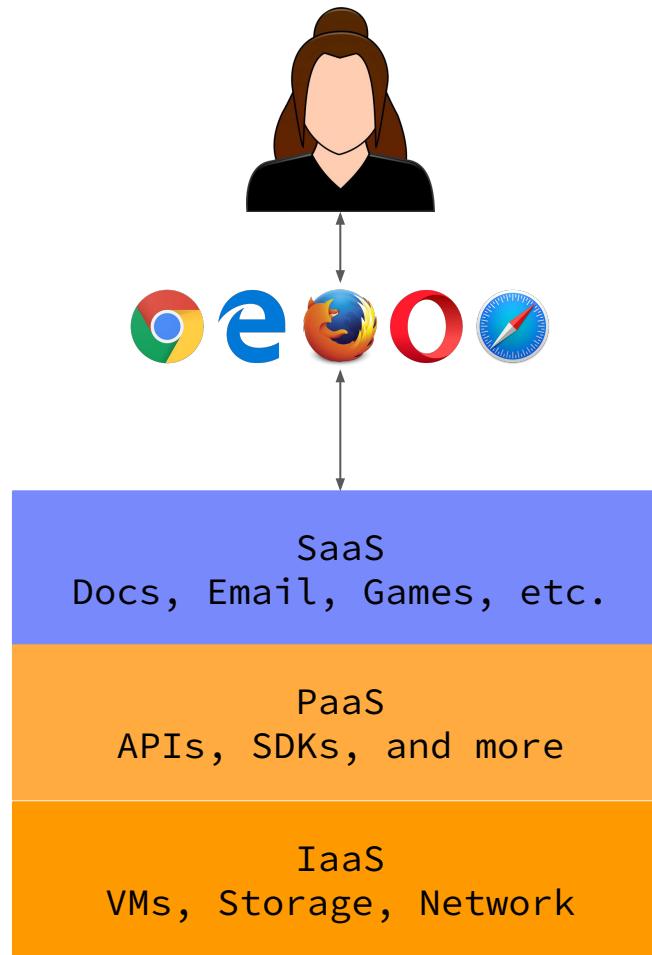
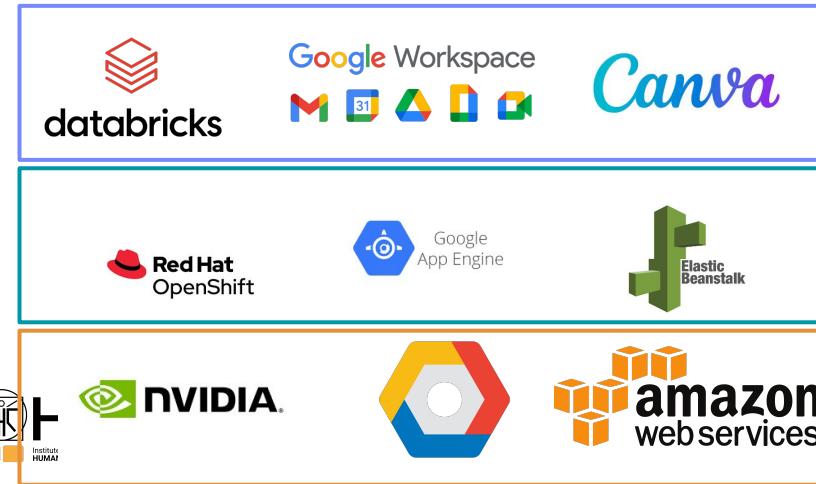
● Definition

- “A Cloud is a type of **parallel** and **distributed system** consisting of a collection of interconnected and **virtualized computers** that are dynamically provisioned and **presented as one or more unified computing resource(s)** based on service-level **agreements** established through negotiation between the **service provider and consumers**” [*].



Buyya, R., Yeo, C. S., Venugopal, S., Broberg, J., & Brandic, I. (2009). Cloud computing and emerging IT platforms: Vision, hype, and reality for delivering computing as the 5th utility. Future Generation computer systems, 25(6), 599-616.

Motivation and Terminology



Motivation and Terminology

- **Transforming IT Industry/Landscape**

- Since ~2010 increasing move from on-prem to cloud resources
- System **software licenses** become increasingly **irrelevant**
- Few cloud providers dominate IaaS/PaaS/SaaS markets
- **2023 revenue:**
 - Microsoft (\$ 111.6B)
 - Amazon AWS (\$ 88B)
 - Oracle Cloud (\$ 35.3B)
 - IBM Cloud (\$ 20.8B)
 - Google Cloud (8.41B)
 - Alibaba Cloud (\$ 3.789M)

Motivation and Terminology

- **Argument #1: Pay as you go:**
 - **No upfront cost** for infrastructure
 - **Variable utilization** → over-provisioning
 - **Pay per use** or acquired resources
- **Argument #2: Economies of Scale**
 - **Lower cost** for purchasing and managing IT infrastructure at scale → lower cost (applies to both HW resources and IT infrastructure/system experts)
 - Focus on **scale-out on commodity HW** over scale-up → lower cost

Motivation and Terminology

- **Argument #3: Elasticity**
 - System can scale up according to demand
 - **Virtually unlimited resources** allows to reduce time as necessary ((Task time = Time x Resources))
- **Argument #4: Availability**
 - Resources are available 24x7

Characteristics and Deployment Models

- **On-demand self service.** Users can access and manage computing resources themselves without requiring human assistance.
- **Broad network access.** Services are accessible from anywhere with an internet connection, using various devices.
- **Resource pooling.** Computing resources are shared among multiple users through virtualization, securely and efficiently

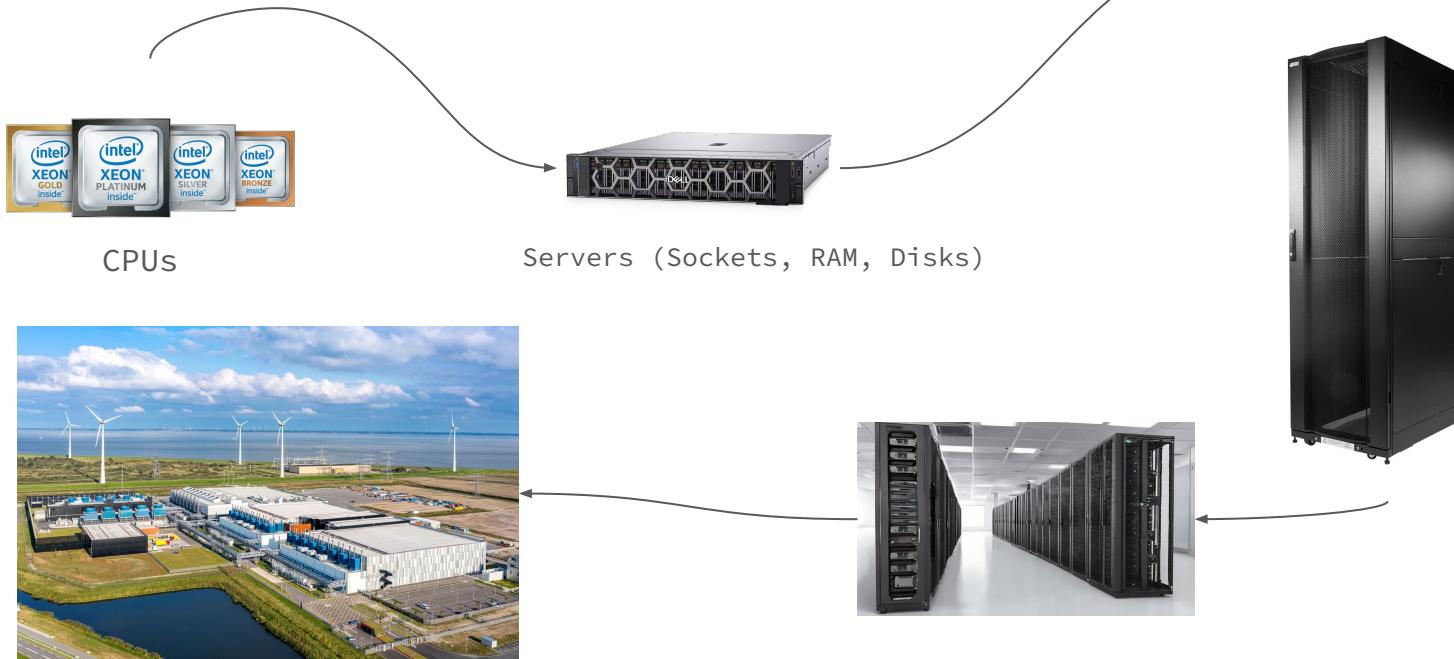
Characteristics and Deployment Models

- **Rapid elasticity.** Resources can be scaled up or down quickly based on demand, in real time.
- **Measured service.** Resource usage is monitored and recorded to optimize consumption and bill only for what you use.

Cloud Computing Service Models

(computing as a utility)

Anatomy of a Data Center



Fault Tolerance

- **Yearly Data Center Failures**

- ~0.5 **overheating** (power down most machines in <5 mins, ~1-2 days)
- ~1 **PDU failure** (~500-1000 machines suddenly disappear, ~6 hrs)
- ~1 **rack-move** (plenty of warning, ~500-1000 machines powered down, ~6 hrs)
- ~1 **network rewiring** (rolling ~5% of machines down over 2-day span)
- ~20 **rack failures** (40-80 machines instantly disappear, 1-6 hrs)
- ~5 **racks go wonky** (40-80 machines see 50% packet loss)
- ~8 **network maintenances** (~30-minute random connectivity losses)
- ~12 **router reloads** (takes out DNS and external vIPs for a couple minutes)
- ~3 **router failures** (traffic rerouting, ~ 1 hour to stabilize.)
- ~dozens of minor 30-second DNS issues
- ~1000 **individual machine failures** (2-4% failure rate, at least twice)
- ~thousands of hard drive failures (**1-5%** of all disks will **die**)

Fault Tolerance

- **Other Common Issues**

- **Configuration issues**, partial SW updates, SW bugs.
- **Transient errors**: no space left on device, memory corruption, stragglers.

- **Error Rates at Scale**

- Cost-effective commodity hardware.
- **Error rate increases** with increasing **scale**.
- **Fault Tolerance** for distributed/cloud storage and data analysis.

Failures are inevitable. large-scale systems are designed with **fault tolerance** to ensure they can **detect**, **recover** from, and **mitigate** the impact of errors, ensuring **reliability** and **minimal downtime**.

Fault Tolerance

Cost-effective Fault Tolerance (BASE)

- **B**Asically available: system mostly operational even during failures, providing partial functionality if necessary.
- **S**oft state: state of the system may change over time, even without new input, due to replication or recovery processes.
- **E**ventual consistency: data may not be instantly consistent across all nodes but will eventually synchronize to the correct state

Fault Tolerance

Cost-effective Fault Tolerance (BASE)

- **Data corruption prevention**
 - ECC (error correction codes)
 - CRC (cyclic redundancy check)
- **Resilient storage**
 - Replication (multiple data copies stored across nodes)
 - Checkpointing (save application state)
 - Lineage (data origin + dependencies)
 - Resilient compute: task re-execution / speculative execution

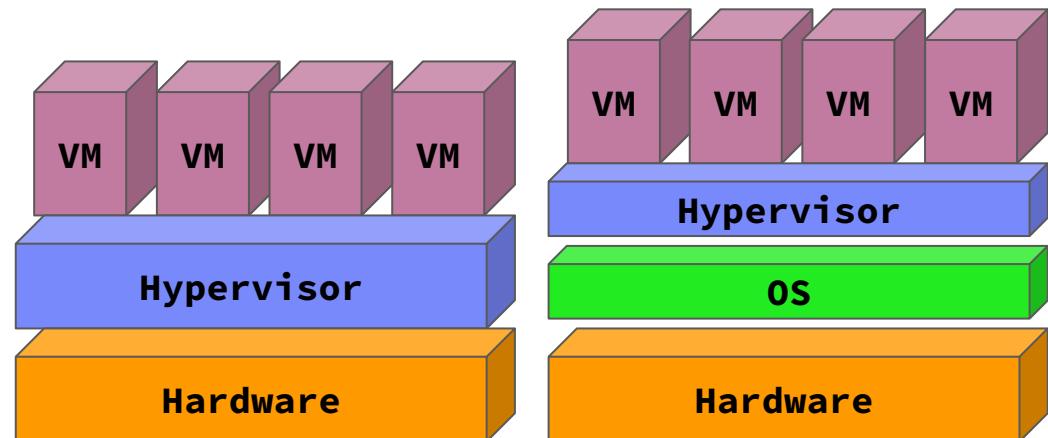
Virtualization I

- **Native Virtualization**

- Simulates most of the HW interface
- Unmodified guest OS to run as if it were on real HW
- Examples: VMWare, Virtualbox, AMI (HVM)

- **Advantages**

- Guess OS as it is!
- High compatibility
- High security



Virtualization II

- **Para-virtualization**

- No HW interface simulation, but special API (hypervcalls) replacing hardware instructions.
- Requires modified guest OS to use hyper calls, trapped by hypervisor
- Examples: Xen, KVM, Hyper-V, AMI (PV)

- **Advantages**

- Faster and more efficient than native virtualization (skips the overhead of hardware simulation).
- Ideal for environments where performance is critical, and modifying the OS is acceptable.

Virtualization III

- **OS-level Virtualization**
 - Allows multiple isolated environments (virtual servers or containers) to run on the OS.
 - **Guest OS** appears isolated but **same as host OS**
 - Examples: **Solaris/Linux containers, Docker**
- **Application-level Virtualization**
 - App executed within a virtualized runtime environment that abstracts away dependencies on the host system.
 - Examples: Java VM (JVM), Ethereum VM (EVM), Python virtualenv

Virtualization: summary

Topic	Native Virtualization	Para Virtualization	OS-level Virtualization	Application-level Virtualization
Definition	Simulates hardware fully	Uses hypercalls, no HW emulation	Containers sharing host OS	Virtualizes specific applications
Guest OS	Unmodified	Modified	Same as host OS	NA
Performance	Moderate (HW overhead)	High	Very high	High
Examples	VMware, VBox	Xen, KVM	Docker, Linux Containers	JVM, Python virtualenv
Use Case	Mixed OS environments	High-performance VMs	Cloud-native apps	App portability across platforms

Containerization

- **Docker Containers**

- Shipping container analogy
- Arbitrary, self-contained goods, standardized units
- Containers reduced loading times → efficient international trade
- Self-contained package of necessary SW and data (read-only image)
- Lightweight virtualization w/ shared OS and resource isolation via control groups



Containerization

Cluster Schedulers

- Container orchestration: scheduling, deployment, and management
- Resource negotiation with clients
- Typical resource bundles (CPU, memory, device)
- Examples: Kubernetes, Mesos, (YARN), Amazon ECS, Microsoft ACS, Docker Swarm

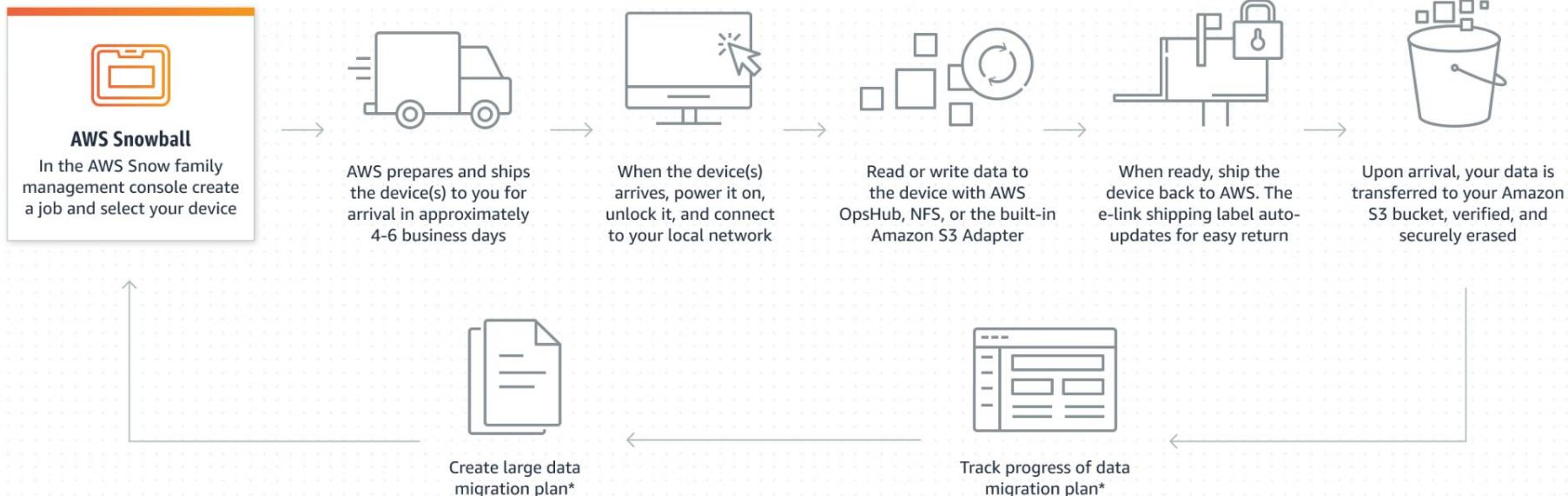




Cloud: How far can we go?

Snowmobile Service:

- Data transfer on-premise → cloud via 100PB trucks w/ 1Gb link



How far can we go?

Microsoft Underwater Datacenter



How far can we go?

Data centers in the space: why and how?

UBIC-DIMW 2019 - The Ninth IEEE International Conference on Mobile Ubiquitous Computing, Systems, Services and Technologies

Space Data Centers – Future Development and Application Perspectives

Tassa Daniel
Côte d'Azur Research School
Rue Fabrice Baudis
Lagord, France
Email: tassa.daniel@enscences.fr

Olawuyomi P. Balogun
Africa Space Innovation Centre
Cape Peninsula Univ of Technology
Cape Town, South Africa
Email: obalogun@capetech.ac.za

Ayodele A. Porilua
Africa Space Innovation Centre
Cape Peninsula University of Technology
Cape Town, South Africa

Abstract: The exponential growth of digital data presents escalating challenges for terrestrial data centers, including energy consumption, ecological impact, and cybersecurity risks. This paper explores the potential of space as an alternative, leveraging microgravity, extreme temperatures, and solar energy to achieve high performance, efficiency, and environmental sustainability. This study explores the technological challenges and opportunities of space-based Data Center Concepts. It also highlights the environmental benefits of space-based data centers, such as reduced energy consumption, high water and land usage for hosting web services. In response to the challenges of terrestrial data centers, this paper advocates for low-cost, eco-friendly web development and content delivery. The research underscores space data centers as a forward-looking solution to the challenges of traditional data centers.

Keywords: Data centers; satellite; space; Internet; computing networks; photons.

1. INTRODUCTION

Data centers are vital for hosting digital content across the Internet. They store and process data in social media and data storage. Terrestrial data centers consist of multiple servers requiring substantial power and cooling, often housed in large buildings. These centers also require extensive cooling systems, such as chillers and indirect environmental cooling systems [1, 2], while power demands increase with data center scale [3, 4]. These constraints lead developers in search of alternative locations to build data centers, resulting in latency and power issues. The challenges of building data centers on Earth have led to the development of space-based data centers, which have been proposed in [5, 6]. The discussion in [5] and [6] has highlighted some limitations in the construction of space-based data centers, such as the challenges in ensuring that data centers deploy their own power generation and cooling systems. However, these approaches do not eliminate the competition for natural resources between terrestrial data center operators and other entities seeking to utilize space resources for other applications.

Terrestrial data centers process mass-based data [1]. To meet the growing demand for data and latency in downloading such data hinder timely decision-making. To address these challenges, this paper proposes a space-based infrastructure with lower environmental impact and reduced latency. The proposed space-based data centers will reduce dependence on Earth's land and water resources and are positioned to process satellite data in orbit, enabling rapid processing and delivery of data to various mobile applications. They also host caches to support low-latency delivery of data to mobile devices. This paper also discusses the challenges of establishing space-based data centers over terrestrial centers for satellite-based communications.

The concept of SDG's is gaining global attention, with the UN's Sustainable Development Goals (SDGs) [7] and the Advanced Space Cluster for European Net Zero Emission (ASCEND) [8] both aiming to develop space-based data centers powered by solar energy to reduce the environmental impact of data centers. This paper explores the feasibility, technical challenges, and environmental implications of deploying data centers in space. It also highlights the challenges of space-based data center development tasks, where the reliance on cloud-based tools and software makes it difficult to estimate the environmental costs. By shifting such responsibilities to SDCs, the potential exists to achieve more sustainable computing. This paper also highlights the importance of research and application of SDG's in the area of web development. The research presented in this paper shows that the web application is that web development utilizes a significant proportion of existing computing resources. If it is proposed to move web development to SDGs, the potential exists for the benefit of reducing the environmental toll due to the use of cloud-based tools and software for web development. The research presents a network architecture for space-based data centers to support web development. This is done to achieve the goal of executing web development aboard SDCs. In addition, the research recognizes the SDG



nature electronics

Perspective

<https://doi.org/10.1038/s41928-019-0516-z>

The development of carbon-neutral data centres in space

Received 1 November 2023

Accepted 8 August 2023

Published online 21 August 2023

Check for updates

Abstract                                                                         <img alt="ORCID icon" data-bbox="8380 210

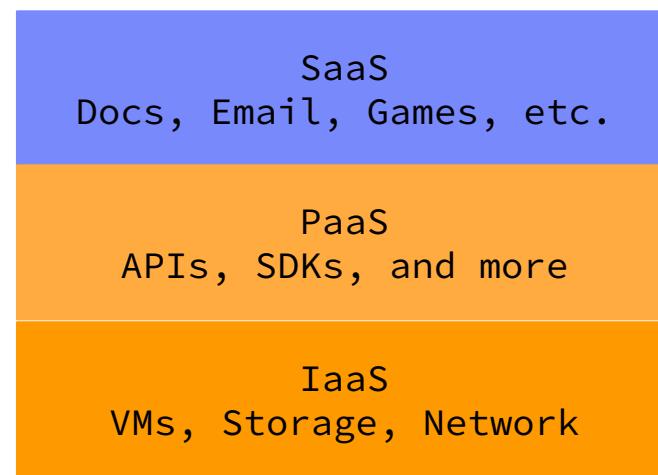
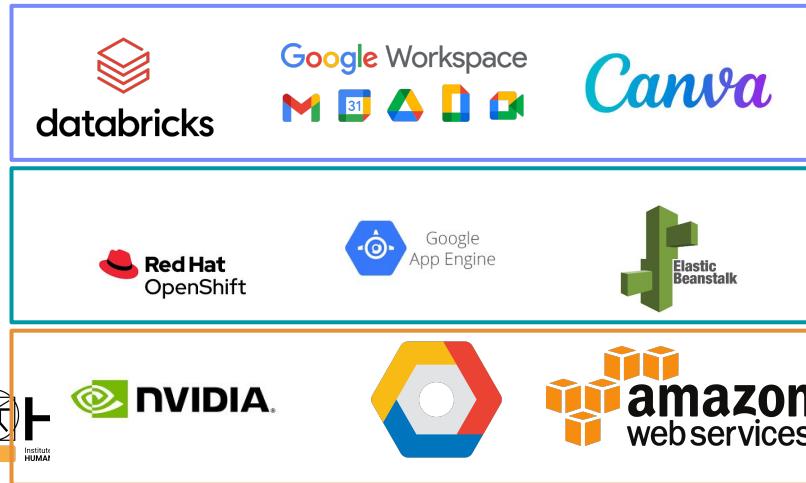
Cloud vs other HPC Technologies

Computing Techniques	Features
Cloud Computing	Cost efficient, almost unlimited storage, backup and recovery, easy deployment
Grid Computing	Efficient use of idle resources, parallelism can be achieved, handles complexity
Cluster Computing	Reduced cost, processing power, improved network technology, scalability, availability

Manvi, S. S., & Shyam, G. K. (2014). Resource management for Infrastructure as a Service (IaaS) in cloud computing: A survey. *Journal of network and computer applications*, 41, 424-440.



Cloud Computing Service Models



IaaS



IaaS
VMs, Storage, Network

Combination of hosting, hardware provisioning, basic services, and other services needed to run a cloud.

Uses of IaaS:

- Provides access to shared resources on need basis, without revealing details like location and hardware to clients.
- Provides details like server images on demand, storage.
- Offers full control of server infrastructure, not limited specifically to applications, instances and containers.

IaaS

Major issues associated with IaaS:

- Resource management.
- Internet access.
- Virtualization.
- Data management.
- APIs.
- Interoperability.



IaaS
VMs, Storage, Network

PaaS

Provision of a computing platform and the provision and deployment of the associated set of software applications (called a solution stack).

Uses of PaaS:

- Provides tools, frameworks, and runtime environments to build and deploy applications efficiently.
- Removes the need to worry about hardware, operating systems, or server management.
- While IaaS provides control over server infrastructure, PaaS limits control to what's essential for applications, offering an environment tuned specifically for application development and management.

PaaS
APIs, SDKs, and more

PaaS

Major issues associated with PaaS:

- Provider reliability
- Resource management
- Internet dependency
- Performance
- Scalability costs



PaaS
APIs, SDKs, and more

SaaS

Software as a Service is a software distribution model in which applications are hosted by a vendor or service provider and made available to customers over a network.

Uses of SaaS:

- On-demand software access.
- Scalability
- Cost-effective solution
- Cross device availability
- Collaboration and real-time updates

SaaS
Docs, Email, Games, etc.

SaaS

Major issues associated with SaaS:

- Internet dependency
- Performance
- Long-term cost
- Data Security and Privacy

SaaS
Docs, Email, Games, etc.

Serverless Computing

A cloud computing paradigm where developers can build and run applications without managing server infrastructure. Operational responsibilities like fault tolerance, scaling, and resource allocation are handled by cloud providers.

Uses of SaaS:

- **Function-as-a-Service (FaaS):** Runs stateless, event-driven functions in response to triggers.
- **Serverless Databases:** Auto-scale capacity as needed and hibernate during periods of inactivity, reducing costs while maintaining availability.

Kounev, S., Herbst, N., Abad, C. L., Iosup, A., Foster, I., Shenoy, P., ... & Chien, A. A. (2023). Serverless computing: What it is, and what it is not?. Communications of the ACM, 66(9), 80-92.

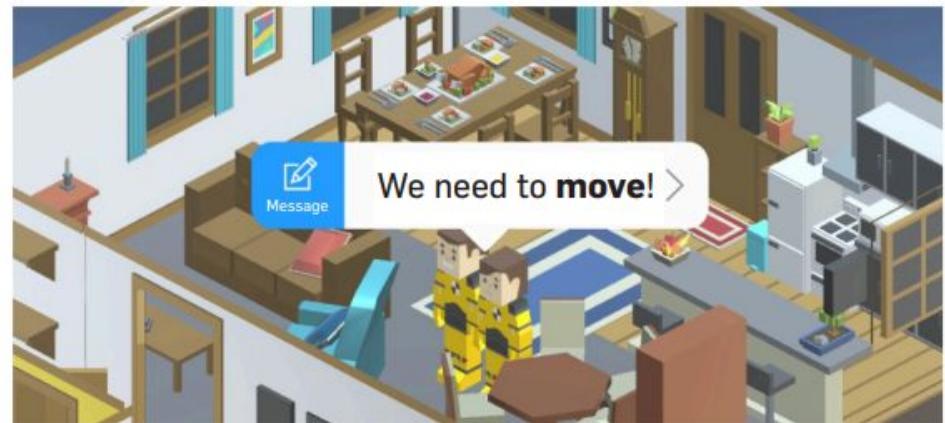


Other “layers”: Serverless

Serverless

IaaS/PaaS/SaaS

Self-hosting



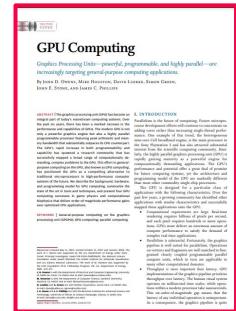
GPU, Fog, and Edge Computing

GPU Computing

GPU processor used as a **highly parallel programmable processor** to handle computationally intensive tasks.

Uses:

- **Machine Learning and AI:**
 - **Training and inference** for deep learning models.
 - Accelerating **neural networks** computations.
- **Scientific Applications:**
 - Protein folding simulations.
 - Computational biophysics and molecular dynamics.
 - Fluid dynamics and heat transfer.

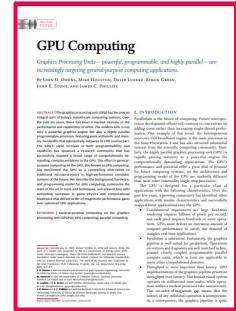


Owens, J. D., Houston, M., Luebke, D., Green, S., Stone, J. E., & Phillips, J. C. (2008). GPU computing. *Proceedings of the IEEE*, 96(5), 879-899.

GPU Computing

Uses:

- **Real-Time Graphics:**
 - Game physics and rendering.
 - Simulation of physical environments in games.
- **Data Processing:**
 - Sorting and searching large datasets.
- **General Numerical Computations:**
 - Differential equation solvers.



Owens, J. D., Houston, M., Luebke, D., Green, S., Stone, J. E., & Phillips, J. C. (2008). GPU computing. *Proceedings of the IEEE*, 96(5), 879-899.



Providers (AWS, NVIDIA, Microsoft, Google, IBM, Oracle, Huawei, Tencent)

AWS:

- EC2 with GPU Instances (NVIDIA Tesla V100, T4, K80)
- Deep Learning AMIs (Amazon Machine Images).
- AWS Lambda for GPU-powered serverless workflows.

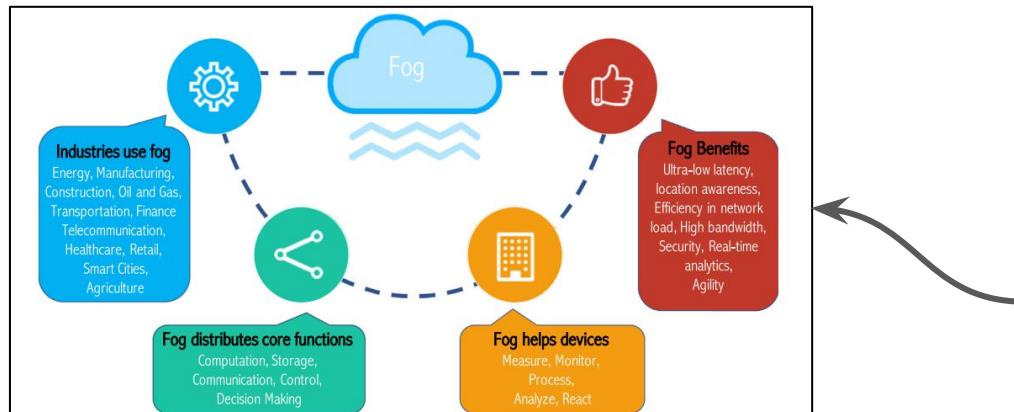
NVIDIA Cloud

- NVIDIA GPU Cloud (Focused on NVIDIA GPUs, including A100 and V100)
- Optimized software stacks for AI, HPC, and data analytics.
- Pre-built containers for machine learning frameworks.



Fog Computing

- Enables computing, storage, networking, and data management on network nodes **within the close vicinity** of IoT devices.
- Computation, storage, networking, decision making, and data management not only occur in the cloud → also **along the IoT-to-Cloud path** (preferably close to the IoT devices).



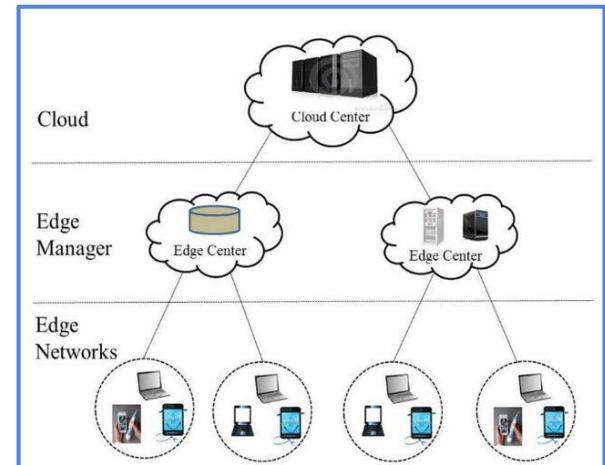
Yousefpour, A., Fung, C., Nguyen, T., Kadiyala, K., Jalali, F., Niakanlahiji, A., ... & Jue, J. P. (2019). All one needs to know about fog computing and related edge computing paradigms: A complete survey. *Journal of Systems Architecture*, 98, 289-330.



Edge Computing

- Enhances → management, storage, and processing for data generated by connected devices.
- Edge computing operates at the edge of the network, positioned near IoT devices (one hop away.)
- **The edge is not directly on the IoT devices!**

Cao, K., Liu, Y., Meng, G., & Sun, Q. (2020). An overview on edge computing research. *IEEE access*, 8, 85714-85728.



Ren, Y., Zhu, F., Qi, J., Wang, J., & Sangaiah, A. K. (2019). Identity management and access control based on blockchain under edge computing for the industrial internet of things. *Applied Sciences*, 9(10), 2058.

A short view about a public cloud...

Connecting...

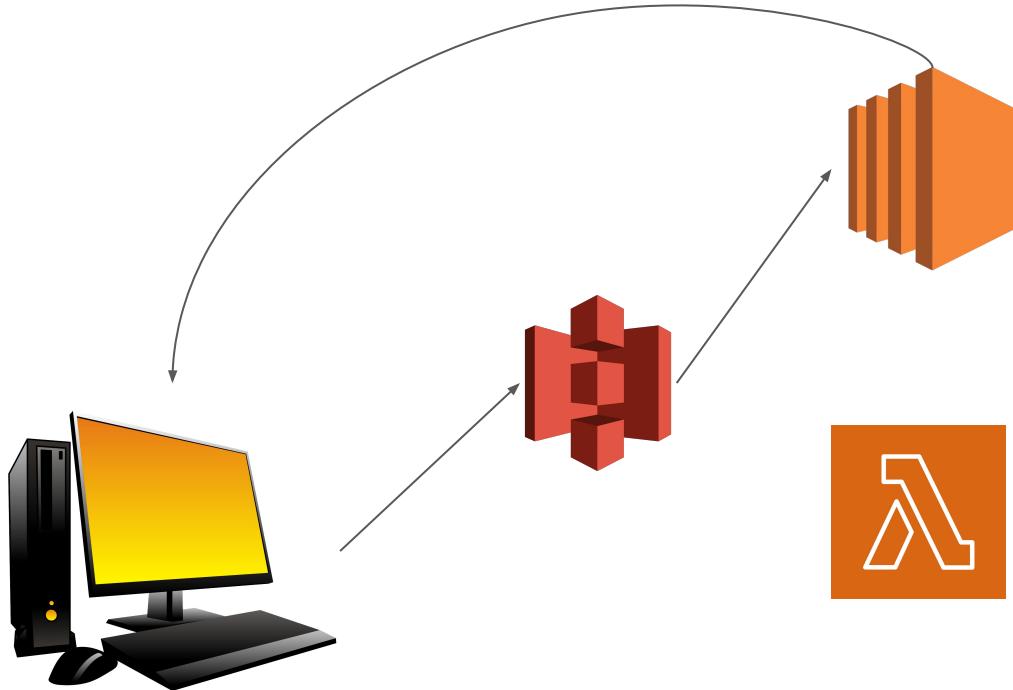
<https://us-west-2.console.aws.amazon.com/>

```
ssh -i <my-key> ubuntu@<my-instance-ip-public-address>
```

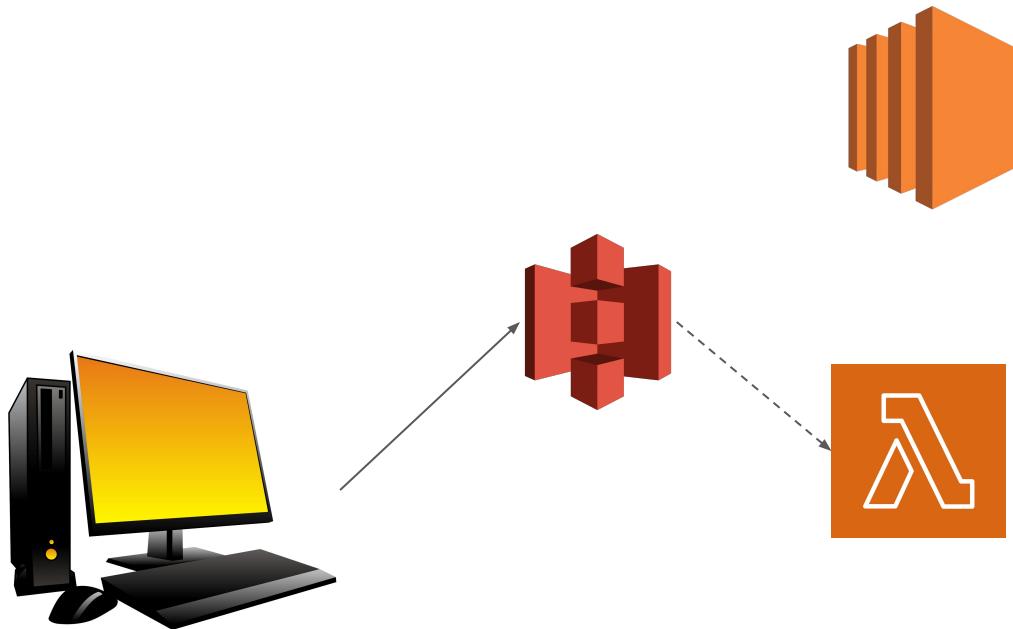
Cloud pipeline



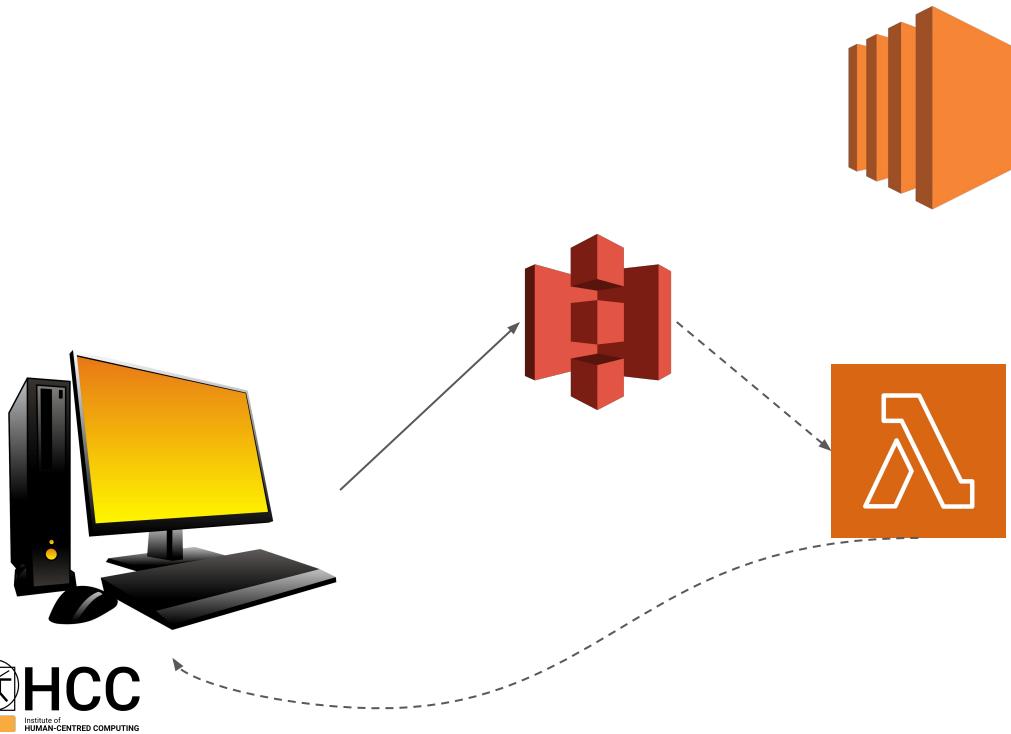
Cloud pipeline



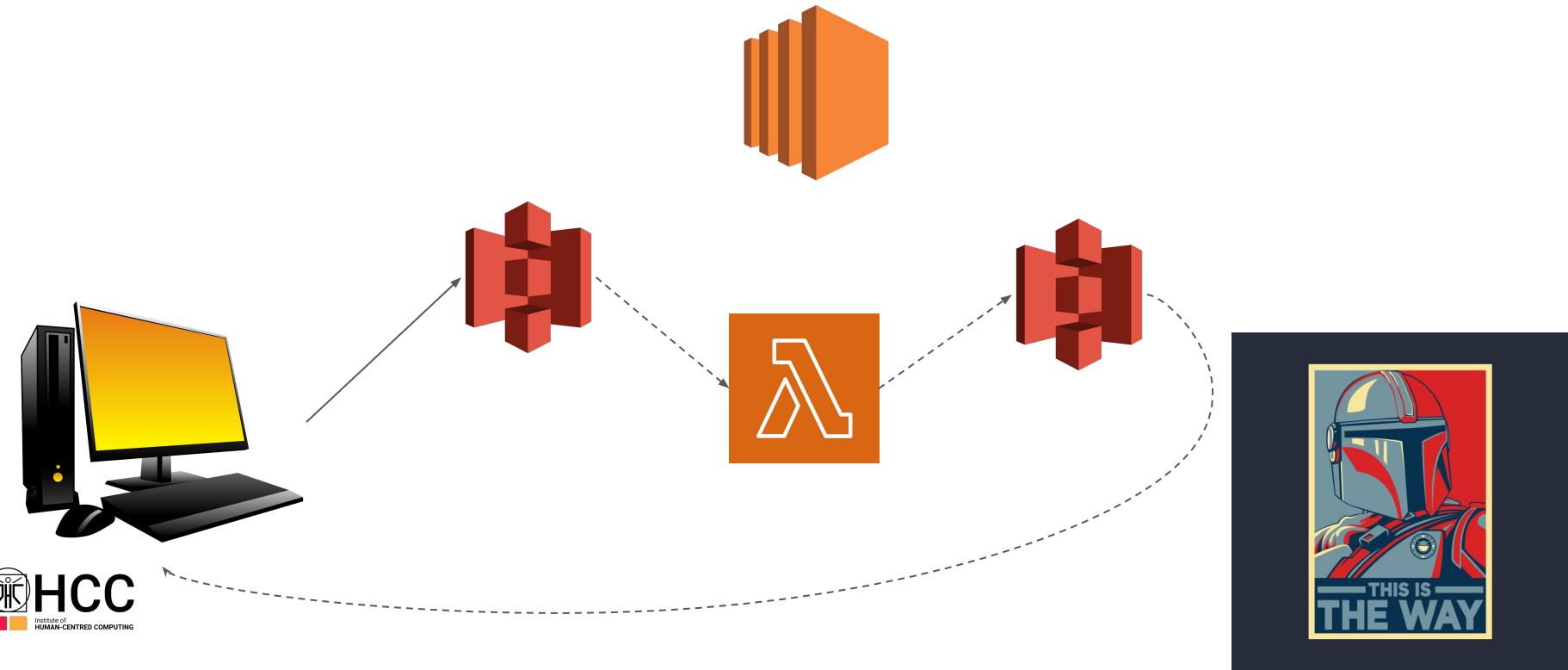
Cloud pipeline



Cloud pipeline



Cloud pipeline



Summary and Q&A

Summary and Q&A

- **Summary and Q&A**
 - Cloud Computing Motivation and Terminology
 - Cloud Computing Service Models
 - Cloud, Fog, and Edge Computing
- **Next Lectures**
 - 08 Cloud Resource Management and Scheduling [Nov 28]