# Scalable Machine Learning and Cloud Computing

Xianyuan Liu

xianyuan.liu@sheffield.ac.uk

# In Today's Lecture:

HPC = High-performance computing cluster
AWS = Amazon Web Services

Today we will,

- introduce local, HPC and Cloud Computing.
- look at **whether** to distribute a given problem
- have a quick run through of **services** and **tools** from the field of Cloud Computing.
- have several **walk-throughs** (AWS **EC2** and **lambda**).

# In the Lab on Friday:

On Friday you will,

- **create** and **connect** to an **instance** on AWS
- set up an **Elastic Map Reduce cluster**
- use it to run **Apache Spark** to solve a matrix factorisation problem
- set up a **Lambda function**

# Local vs HPC vs Cloud Computing

- Local Computing
    - Physically and directly managed by users, normally in the same geographic location as the user.
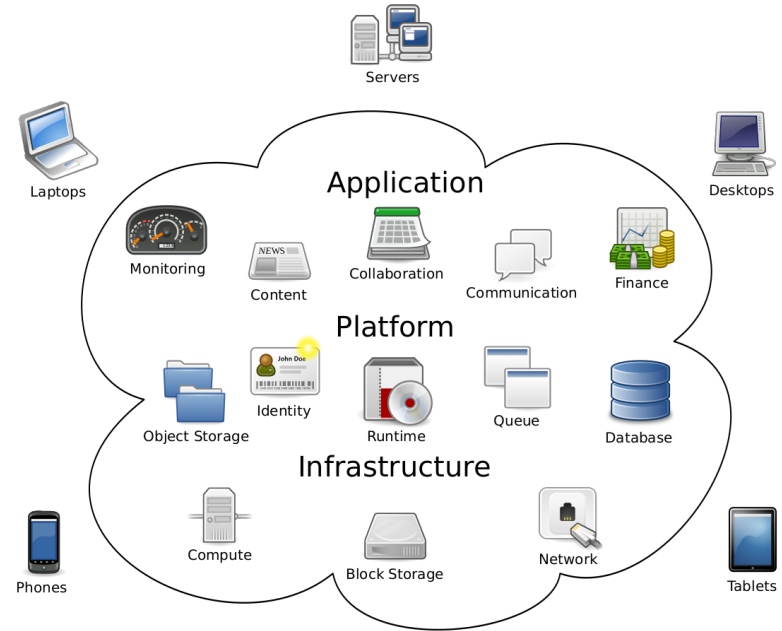
# Local vs HPC vs Cloud Computing

- Local Computing
    - Physically and directly managed by users, normally in the same geographic location as the user.

- HPC
    - Clusters of computers combined to form a powerful computing environment (Stanage - Manchester, Bessemer - Leeds)

- Cloud Computing



*The old Iceberg system server racks.*

Image from: https://docs.hpc.shef.ac.uk/en/latest/hpc/what-is-hpc.html#gsc.tab=0

# What is Cloud Computing

- Cloud computing is the **on-demand** availability of computer system resources, without direct active management by the user.

- Cloud computing relies on the **sharing of resources** to achieve coherence and typically uses a **pay-as-you-go** model.

- Large clouds often have functions **distributed** over multiple locations.

# HPC vs Cloud?

Until now we've used the university's own HPC system.

HPC:
- Often provided for **free** by institution
- highly **interconnected** nodes - vital for many problems
- Most HPCs use SLURM (includes tools like job array, which let you run embarrassingly parallel arrays of jobs).
- Provide **support**

# HPC vs Cloud?

HPC:

- Often provided for **free** by institution
- highly **interconnected** nodes - vital for many problems
- Most HPCs use SLURM (includes tools like job array, which let you run embarrassingly parallel arrays of jobs).
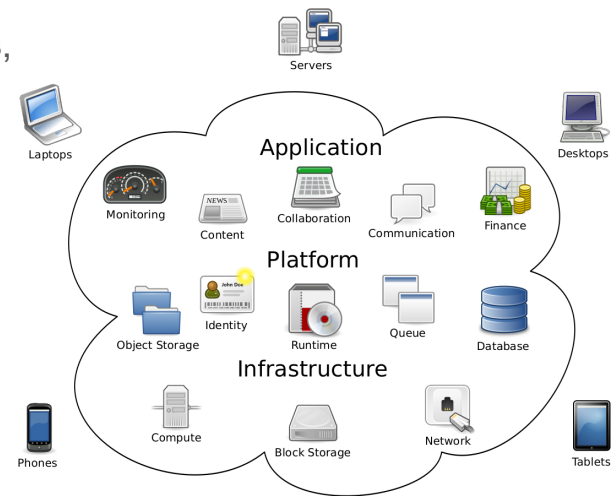- Provide **support**

Cloud:

- **No queuing**
- **Quick** to set up
- Often **only option** outside of universities
- Well suited for **embarrassingly parallel** problems
- Almost all **platforms** and features **supported**
- Appropriate if **hosting a front-end server**
- Can be very expensive (+ accidental spend!)



Often 'rent' (virtual) hardware while you need it. Large clusters can be created briefly.
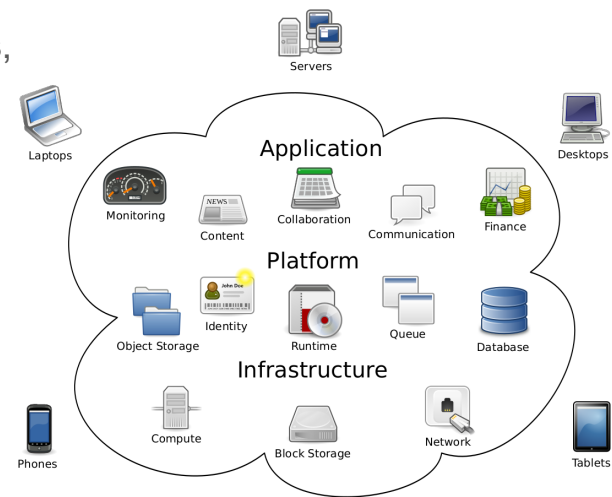
# Why Cloud Computing for Machine Learning

- Powerful Hardware
  - **On-demand** access to state-of-the-art hardware resources, instead of expensive investments in physical hardware.
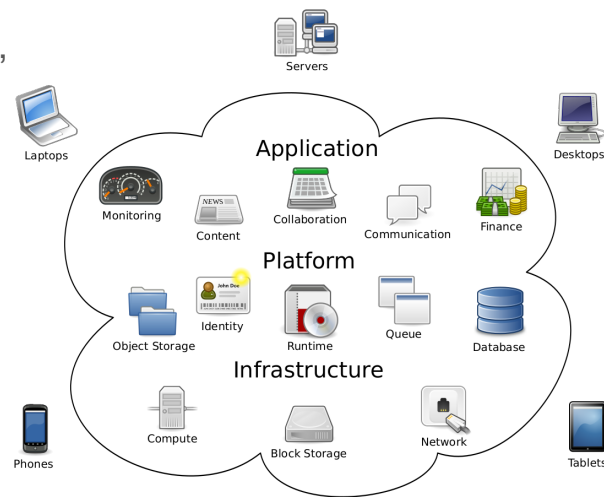
# Why Cloud Computing for Machine Learning

- Powerful Hardware
  - **On-demand** access to state-of-the-art hardware resources, instead of expensive investments in physical hardware.
- Scalability
  - **Dynamically** adjust resources (e.g., hardware, memory, storage) as your needs change.



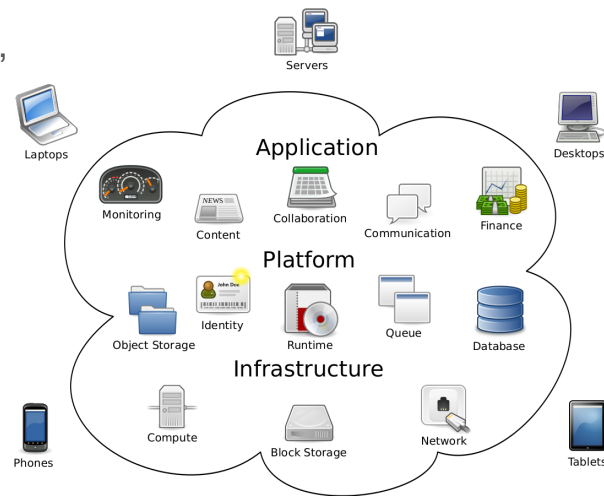Content and image from: https://en.wikipedia.org/wiki/Cloud_computing

# Why Cloud Computing for Machine Learning

- Powerful Hardware
  - **On-demand** access to state-of-the-art hardware resources, instead of expensive investments in physical hardware.
- Scalability
  - **Dynamically** adjust resources  (e.g., hardware, memory, storage) as your needs change.
- Pay-as-you-go Model
  - Pay **only** for the compute time and storage used.
  - Achieve significant savings on energy consumption, facility management, and ongoing hardware maintenance.

# Why Cloud Computing for Machine Learning

- Powerful Hardware
  - **On-demand** access to state-of-the-art hardware resources, instead of expensive investments in physical hardware.
- Scalability
  - **Dynamically** adjust resources  (e.g., hardware, memory, storage) as your needs change.
- Pay-as-you-go Model
  - Pay **only** for the compute time and storage used.
  - Achieve significant savings on energy consumption, facility management, and ongoing hardware maintenance.
- Flexibility
  - Access from **anywhere** with an internet connection.
  - Share resources with **multiple** users.

# Local vs HPC vs Cloud Computing

- Local Computing
  - Physically and directly managed by users, normally in the same geographic location as the user.

- HPC
  - Clusters of computers combined to form a powerful computing environment (Stanage - Manchester, Bessemer - Leeds)

- Cloud Computing

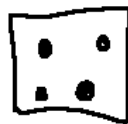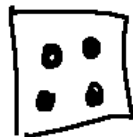  Distributed system - Parallel computing



*The old Iceberg system server racks.*

Image from: https://docs.hpc.shef.ac.uk/en/latest/hpc/what-is-hpc.html#gsc.tab=0

# Distributed system - Parallel computing

Different levels of parallelism:

- **Single process** (e.g. python GIL code)

- Well parallelised multiprocess code (might be good to deploy on a single instance with lots of cores on AWS or HPC)

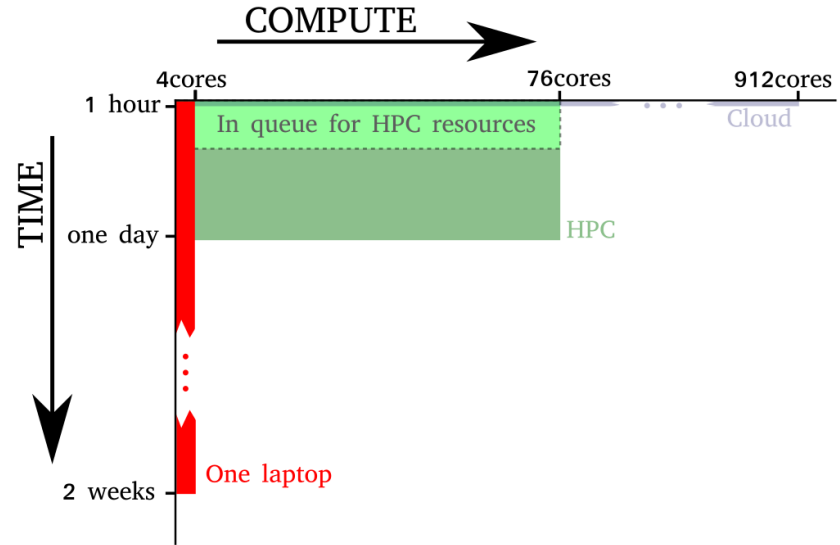- Distributed over **multiple nodes** (either on AWS or HPC).

# Distributed system - Parallel computing

Different levels of parallelism:

- **Single process** (e.g. python GIL code)

- Well parallelised multiprocess code (might be good to deploy on a single instance with lots of cores on AWS or HPC)

- Distributed over **multiple nodes** (either on AWS or HPC).

# Local vs Distributed

- It's not always worth parallelising your application!

# Local vs Distributed

- It's not always worth parallelising your application!
- **Human time** is worth more than **computer time**

# Local vs Distributed

- It's not always worth parallelising your application!
- **Human time** is worth more than **computer time**, so if you can do something else and be patient it might be best not to bother with rewriting it to be parallel.
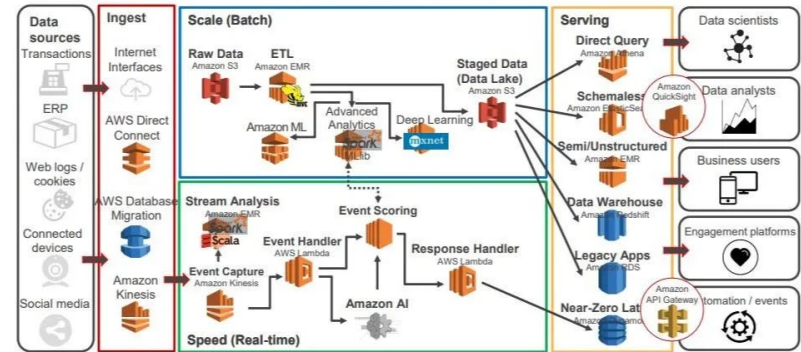
# Why stay Local…

-   May need to **redesign** algorithm

# Why stay Local…

- May need to **redesign** algorithm
- Managing the **pipeline**.

# Why stay Local…

- May need to **redesign** algorithm
- Managing the **pipeline**.
- Typically such environments are less **reliable** and less **predictable**.
  - Far more difficult to **debug**.

# Why stay Local…

- May need to **redesign** algorithm
- Managing the **pipeline**.
- Typically such environments are less **reliable** and less **predictable**.
    - Far more difficult to **debug**.
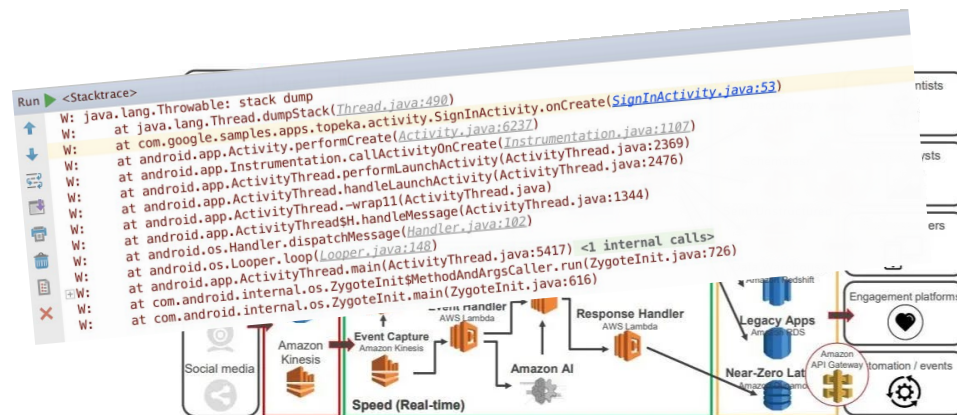- **Shifting data** to/from the servers might be a bottleneck.

# Why stay Local…

- May need to **redesign** algorithm
- Managing the **pipeline**.
- Typically such environments are less **reliable** and less **predictable**.
  - Far more difficult to **debug**.
- **Shifting data** to/from the servers might be a bottleneck.
- **Security** and **privacy** (where are you sending the data? Is it private?)

# Privacy and the GDPR…

The **EU GDPR** (General Data Protection Regulation) restricts the transfer of **personal data** to countries outside of the EU.

# Privacy and the GDPR…

Post-brexit: UK GDPR is now in law, and is basically the same as the EU GDPR for now.

The **EU GDPR** (General Data Protection Regulation) restricts the transfer of **personal data** to countries outside of the EU.

# Privacy and the GDPR…

Post-brexit: UK GDPR is now in law, and is basically the same as the EU GDPR for now.

The **EU GDPR** (General Data Protection Regulation) restricts the transfer of **personal data** to countries outside of the EU.

One can process the data if it's done in a country covered by an EU/UK "**adequacy decision**"

The UK and EU consider each other to have this.

Otherwise one needs a **standard contractual clause** (these have begun to diverge between the EU and UK).

# Privacy and the GDPR…

Post-brexit: UK GDPR is now in law, and is basically the same as the EU GDPR for now.
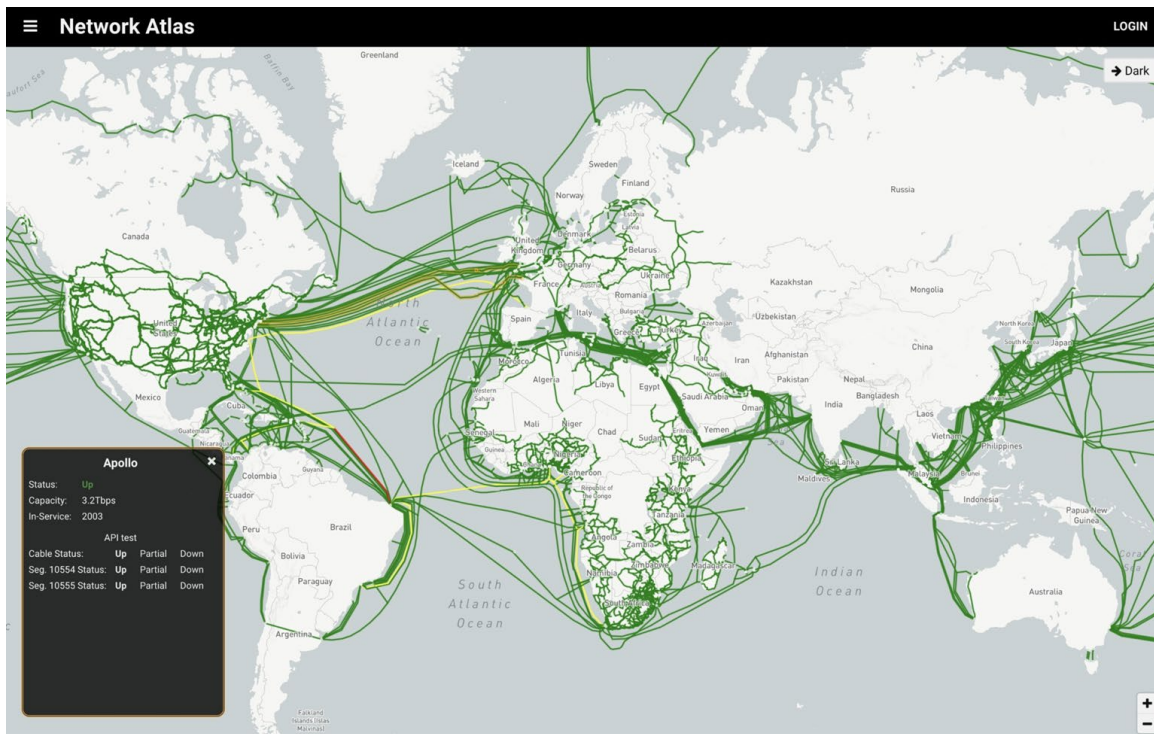
The **EU GDPR** (General Data Protection Regulation) restricts the transfer of **personal data** to countries outside of the EU.

One can process the data if it's done in a country covered by an EU/UK "**adequacy decision**"

The UK and EU consider each other to have this.

Otherwise one needs a **standard contractual clause** (these have begun to diverge between the EU and UK).



≡ **Network Atlas**     LOGIN

Countries in the eyes of the GDPR
- ■ Adequate Protection
- ■ EU Member States
- ■ Norway, Liechtenstein and Iceland

Created with mapchart.net

Falkland Islands Islas Malvinas

# Privacy and the GDPR…

Post-brexit: UK GDPR is now in law, and is basically the same as the EU GDPR for now.
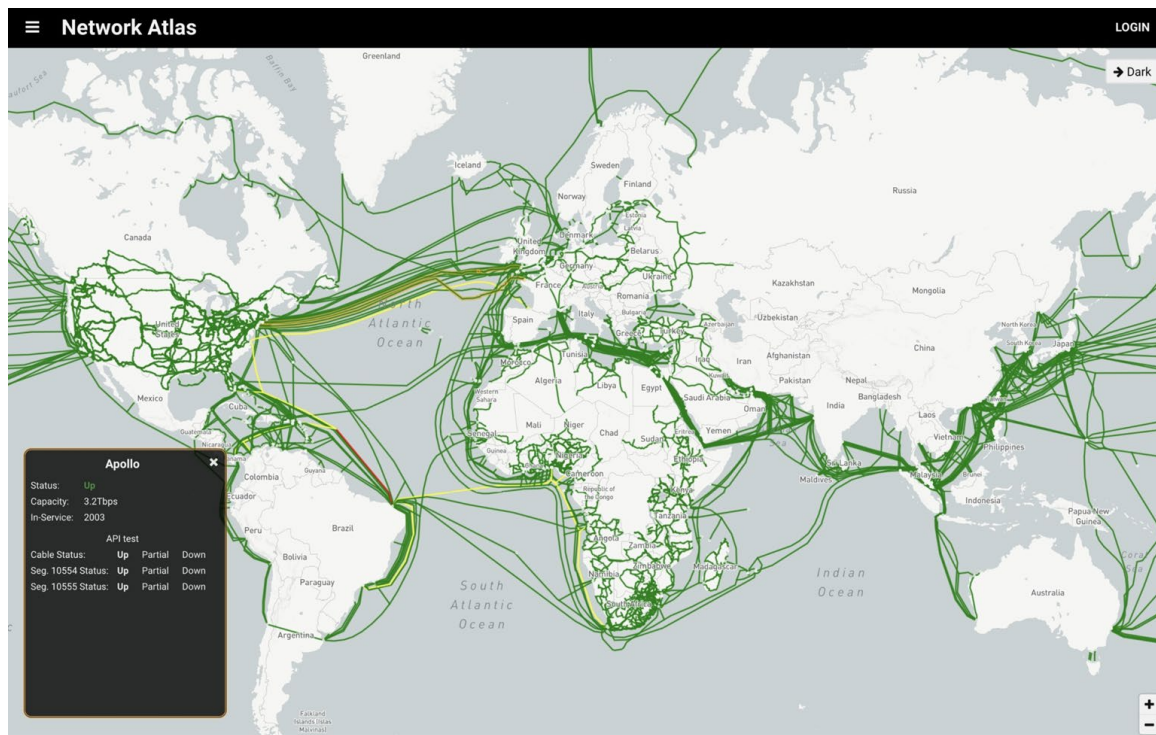
The **EU GDPR** (General Data Protection Regulation) restricts the transfer of **personal data** to countries outside of the EU.

One can process the data if it's done in a country covered by an EU/UK "**adequacy decision**"

The UK and EU consider each other to have this.

Otherwise one needs a **standard contractual clause** (these have begun to diverge between the EU and UK).



If working with UK/EU data: Best to pick AWS regions inside the relevant jurisdiction.

# Why stay Local…

- May need to **redesign** algorithm
- Managing the **pipeline**.
- Typically such environments are less **reliable** and less **predictable**.
    - Far more difficult to **debug**.
- **Shifting data** to/from the servers might be a bottleneck.
- **Security** and **privacy** (where are you sending the data? Is it private?)
- **Cheaper**.

# Service Models

There are different **service models** available: the trade off is typically complexity vs flexibility:

- **Infrastructure** as a Service (IaaS) - high level of control
    - Provides virtual servers, storage, networking resources, and other essential elements.
    - E.g. **AWS EC2**, Microsoft Azure Virtual Machines, Google Compute Engine.

# Service Models

There are different **service models** available: the trade off is typically complexity vs flexibility:

- **Infrastructure** as a Service (IaaS) - high level of control
    - Provides virtual servers, storage, networking resources, and other essential elements.
    - E.g. **AWS EC2**, Microsoft Azure Virtual Machines, Google Compute Engine.

- **Platform** as a Service (PaaS) [always on] - medium
- **Function** as a Service (FaaS) [event driven] - low
    - Offers a development and deployment platform on top of the infrastructure.
    - You provide software/code, but the service starts up virtual machines, allocates storage, runs databases, etc.
    - E.g. website with database / **AWS Lambda**, Azure Functions, Google Cloud Functions.

# Service Models

There are different **service models** available: the trade off is typically complexity vs flexibility:

- **Infrastructure** as a Service (IaaS) - high level of control
    - Provides virtual servers, storage, networking resources, and other essential elements.
    - E.g. **AWS EC2**, Microsoft Azure Virtual Machines, Google Compute Engine.

- **Platform** as a Service (PaaS) [always on] - medium
- **Function** as a Service (FaaS) [event driven] - low
    - Offers a development and deployment platform on top of the infrastructure.
    - You provide software/code, but the service starts up virtual machines, allocates storage, runs databases, etc.
    - E.g. website with database / **AWS Lambda**, Azure Functions, Google Cloud Functions.

- **Software** as a Service (SaaS) - low
    - Provides ready-to-use software applications over the internet
    - E.g. Google Docs

# Service Models

There are different **service models** available: the trade off is typically complexity vs flexibility:

- **Infrastructure** as a Service (IaaS) - high level of control
    - Provides virtual servers, storage, networking resources, and other essential elements.
    - E.g. **AWS EC2**, Microsoft Azure Virtual Machines, Google Compute Engine.

- **Platform** as a Service (PaaS) [always on] - medium
- **Function** as a Service (FaaS) [event driven] - low
    - Offers a development and deployment platform on top of the infrastructure.
    - You provide software/code, but the service starts up virtual machines, allocates storage, runs databases, etc.
    - E.g. website with database / **AWS Lambda**, Azure Functions, Google Cloud Functions.

Infrastructure as Code (IaC) - cloud configuration specified in code.

- **Software** as a Service (SaaS) - low
    - Provides ready-to-use software applications over the internet
    - E.g. Google Docs

# Cloud Computing Platforms

We will focus on **AWS** (others similar).



**Amazon Maintains Cloud Lead as Microsoft Edges Closer**

Worldwide market share of leading cloud infrastructure service providers in Q4 2023*

| Provider | Market Share |
|----------|-------------|
| aws | 31% |
| Azure | 24% |
| Google Cloud | 11% |
| Alibaba Cloud | 4% |
| salesforce | 3% |
| IBM Cloud | 2% |
| ORACLE | 2% |
| Tencent Cloud | 2% |

Cloud infrastructure service revenues in Q4 2023
**$73.7B**

* Includes platform as a service (PaaS) and infrastructure as a service (IaaS) as well as hosted private cloud services
Source: Synergy Research Group

statista

# Cloud Computing Platforms

We will focus on **AWS** (others similar).

Huge range of services.

Amazon Elastic Compute Cloud (**EC2**) allows users to rent virtual machines.

**Amazon Maintains Cloud Lead as Microsoft Edges Closer**

Worldwide market share of leading cloud infrastructure service providers in Q4 2023*

## All services

### Services by category

**Compute**
- EC2
- Lightsail
- Lambda
- Batch
- Elastic Beanstalk
- Serverless Application Repository
- AWS Outposts
- EC2 Image Builder
- AWS App Runner
- AWS SimSpace Weaver

**Containers**
- Elastic Container Registry
- Elastic Container Service
- Elastic Kubernetes Service
- Red Hat OpenShift Service on AWS

**Storage**

**Management & Governance**
- AWS Organizations
- CloudWatch
- AWS Auto Scaling
- CloudFormation
- AWS Config
- OpsWorks
- Service Catalog
- Systems Manager
- Trusted Advisor
- Control Tower
- AWS Well-Architected Tool
- AWS Chatbot
- Launch Wizard
- AWS Compute Optimizer
- Resource Groups & Tag Editor
- Amazon Grafana
- Amazon Prometheus
- AWS Resilience Hub

**Security, Identity, & Compliance**
- Resource Access Manager
- Cognito
- Secrets Manager
- GuardDuty
- Amazon Inspector
- Amazon Macie
- IAM Identity Center
- Certificate Manager
- Key Management Service
- CloudHSM
- Directory Service
- WAF & Shield
- AWS Firewall Manager
- AWS Artifact
- Detective
- AWS Signer
- AWS Private Certificate Authority
- Security Hub

# AWS: Examples of Instance Types

On-demand (pay per hour) or Reserved

# AWS: Examples of Instance Types

On-demand (pay per hour) or Reserved

**Cheapest**: t4g.nano [burstable, smallest instance] for 3 years $41 (ex. VAT).
- Also need to pay a few pence for storage, data transfer, etc.

# AWS: Examples of Instance Types

On-demand (pay per hour) or Reserved

**Cheapest**: t4g.nano [burstable, smallest instance] for 3 years $41 (ex. VAT).
- Also need to pay a few pence for storage, data transfer, etc.

**Compute** optimised

| Instance name ▽ | On-Demand hourly rate ▽ | vCPU ▽ | Memory ▽ | Storage ▽ | Network performance ▼ |
|---|---|---|---|---|---|
| c7a.metal-48xl | $9.85344 | 192 | 384 GiB | EBS Only | 50000 Megabit |

**Memory** optimised

| Instance name ▽ | On-Demand hourly rate ▽ | vCPU ▽ | Memory ▽ | Storage ▽ | Network performance ▼ |
|---|---|---|---|---|---|
| u-24tb1.112xlarge | $261.73 | 448 | 24576 GiB | EBS Only | 100 Gigabit |

# AWS: Examples of Instance Types

**Other decisions:** Storage type (e.g. HDD vs SSD), what OS will you install, etc. We will explore this in more detail in the lab.

On-demand (pay per hour) or Reserved

**Cheapest**: t4g.nano [burstable, smallest instance] for 3 years $41 (ex. VAT).
- Also need to pay a few pence for storage, data transfer, etc.

**Compute** optimised

| Instance name ▽ | On-Demand hourly rate ▽ | vCPU ▽ | Memory ▽ | Storage ▽ | Network performance ▽ |
|---|---|---|---|---|---|
| c7a.metal-48xl | $9.85344 | 192 | 384 GiB | EBS Only | 50000 Megabit |

**Memory** optimised

| Instance name ▽ | On-Demand hourly rate ▽ | vCPU ▽ | Memory ▽ | Storage ▽ | Network performance ▽ |
|---|---|---|---|---|---|
| u-24tb1.112xlarge | $261.73 | 448 | 24576 GiB | EBS Only | 100 Gigabit |

There are also instances that have one or more **GPUs**.

# Some tools widely used…

- MapReduce

- Apache Spark

- Terraform

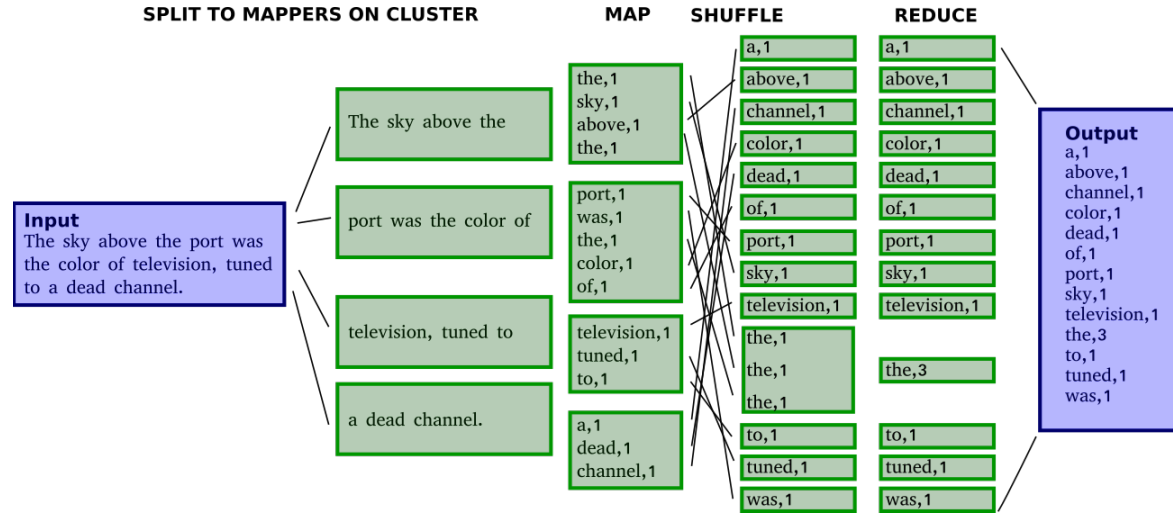- Docker

# MapReduce

- A programming model for processing **large datasets** in parallel across clusters.

- First, **breaks down** a large task into smaller, independent subtasks (**map** phase)

- **Shuffles** the results to be in the right locations and then aggregates them (**reduce** phase).



Example of a distributed word-count, using MapReduce

# Apache Spark

- Apache SPARK is a **distributed cluster** computing framework.

- More flexible and faster than MapReduce
    - Functionalities for batch data processing (similar to MapReduce), real-time data processing, and machine learning.
    - In-memory processing rather than disk-based processing.

- Particularly useful for **iterative algorithms** (e.g. ML gradient descent).

    We'll look at this in the lab in more detail.

Spark's RDDs function as a working set for distributed programs that offers a (deliberately) restricted form of distributed shared memory.

Inside Apache Spark the workflow is managed as a directed acyclic graph (DAG). Nodes represent RDDs while edges represent the operations on the RDDs.

- from wikipedia's article.

# Terraform



- Defines the cloud infrastructure (virtual servers, storage, networking) as code files and can automate the creation and management across different cloud providers.

- Essential to manage in a scalable and automated way, especially for complex deployments.

- **Infrastructure** as code (IaC)

https://docs.aws.amazon.com/cli/latest/userguide/getting-started-install.html
https://app.terraform.io/app/getting-started/example

# Docker



- A platform allows to put all the **software, libraries and configuration** files that you need to run an application into a **container**.

- It sort of replaces **virtual machines**, in a 'lighter' way, and makes it easier to **move data** between them.

- Simplifies development workflows, and ensures consistent application behavior across different environments.

- Mainly for linux, but typically not on HPC.

# Other service… Lambda



- AWS Lambda allows you to put code on the cloud, **without having to deal with running a server**. The code is **triggered** by a variety of **events**, and can perform basic tasks in response. Typically **stateless**.

- **Function** as a Service (FaaS)

- Examples of where Lambda might be useful:
    - IoT deployments                              Parsing incoming data from sensors
    - Automated backups      An hourly script run and checking the backups!
    - ML deployment                              Refresh air pollution map every 10 minutes
    - Media conversion      Managing the deployment of infrastructure, Netflix
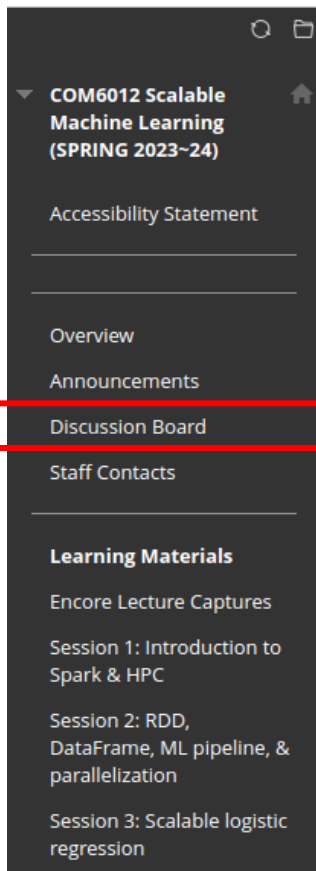    - Websites                              Low traffic sites with scripts

# Activities

-   AWS EC2 - Infrastructure as a Service (IaaS)

-   AWS Lambda - Function as a Service (FaaS)

# Take Home Messages

- Not always worth distributing your code (and lots of issues, including privacy).

- Service Models (Infrastructure-, Platform-/Function-, Software- as-a-Service).
    - EC2 is an example of Infrastructure-as-a-Service.
    - Lambda is an example of a Function-as-a-Service.

- Technology is changing quickly.
    - AWS EC2 is only about 15 years old.
    - Docker and Lambda are only about 10 years old.
    - Need to keep yourself updated if you enter this field.

# Please note:



Assignment Brief:
https://docs.google.com/document/d/1QSBkfnLLgf5qM0KWkeRayeAbrITRKY4JZ8aYZA8npro/edit

FAQ
https://docs.google.com/document/d/1QSBkfnLLgf5qM0KWkeRayeAbrITRKY4JZ8aYZA8npro/edit#bookmark=id.3ukytk47kosf

You are required to complete all assignment questions using **batch mode**.

# Acknowledgement

The slides for this lecture are adapted from <u>Dr Michael Smith</u>'s guest lecture slides in 2023.