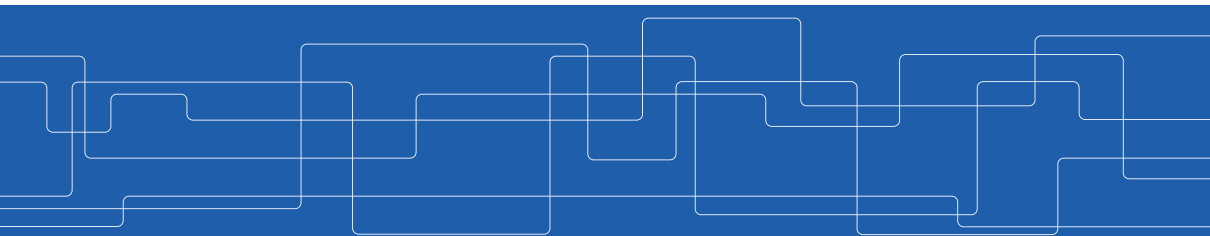




# An Introduction to Data Intensive Computing

Amir H. Payberah  
payberah@kth.se  
2023-08-29





# Course Information



## Course Objective

- ▶ Provide students with a solid foundation for **understanding** large scale distributed **systems** used for **storing and processing** massive data.
- ▶ Cover a wide variety of advanced topics in **data intensive computing platforms**, i.e., the frameworks to **store and process** big data.



## Intended Learning Outcomes (ILOs)

- ▶ **ILO1:** **Understand** the main concepts of **data-intensive computation platforms**.
- ▶ **ILO2:** **Apply** the grabbed knowledge to **store** and **process** massive data.
- ▶ **ILO3:** **Analyze** the **technical merits** of data-intensive computation platforms.



# The Course Assessment

- ▶ **Task1**: the **review** questions.
- ▶ **Task2**: the **lab** assignments.
- ▶ **Task3**: the **essay** and the **presentation**.
- ▶ **Task4**: the **project**.
- ▶ **Task5**: the final **exam**.
- ▶ All the assignments should be done in **groups** of **two/three** students.



## How Each ILO is Assessed?

	<b>Task1</b>	<b>Task2</b>	<b>Task3</b>	<b>Task4</b>	<b>Task5</b>
<b>ILO1</b>	x	x			x
<b>ILO2</b>		x		x	
<b>ILO3</b>			x		



## Task1: The Review Questions

- ▶ Five set of review questions, one set for each week.
- ▶ The review questions are graded P/F.



## Task2: The Lab Assignments

- ▶ Four lab assignments, each focuses on a specific topic.
- ▶ No deadline.





## Task3: The Essay and The Presentation

- ▶ One module for each group: writing an **essay** and **presenting** it to their **opponents** (another group).
- ▶ Grading of this task has the following parts:
  - *E*: **Essay** (5 points)
  - *P*: **Presentation** (2 points)
  - *Q*: **Reviewing essay** and **asking questions** (2 points)
  - *A*: **Answering questions** (1 point)
- ▶ Each part is graded **A-F**.
- ▶ The final grade: **A: 10, B: 9, C: 8, D: 7, E: 6, F: <5**.



## Task4: The Final Project

- ▶ One final project: **source code** and **oral presentation**.
- ▶ Proposed by students and confirmed by the teacher.
- ▶ It is graded **A-F**.



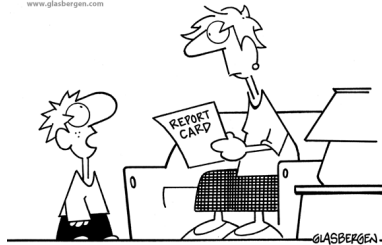
## Task5: The Final Exam

- ▶ The **final exam** covers **all the modules** presented during the course
- ▶ It is graded **A-F**.

# The Final Grade

- ▶ To pass the course: you must **pass Task 1** and get **at least E** in Task 3, Task 4, and Task 5.
- ▶ The **final grade** of the course is computed as  $0.3 \times \text{Task3} + 0.3 \times \text{Task4} + 0.4 \times \text{Task5}$ .

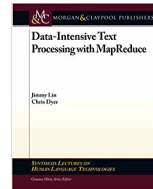
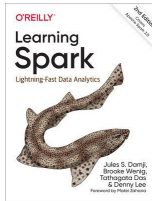
© Randy Glasbergen  
www.glasbergen.com



"Why is an A or B better than a C or D?  
Aren't all letters equal in the eyes of God?"

# The Course Material

- ▶ Mainly based on [research papers](#).
- ▶ We also cover the following books.





## The Course Web Page

`https://id2221kth.github.io`



## The Questions-Answers Page

<https://tinyurl.com/hk7hzpw5>



# The Course Overview





# Cloud Computing and Big Data

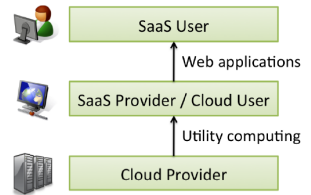
- ▶ The main trends:
  - Computers not getting any faster
  - Internet connections getting faster
  - More people connected to the Internet
- ▶ **Conclusion:** move the computation and storage of big data to the cloud!



# Cloud Computing

# Cloud Computing Definition

- ▶ **Cloud Computing** refers to both:
  1. The **applications** delivered as **services** over the Internet
  2. The **hardware and systems software** in the datacenters that provide those **services**
- ▶ The **services**: called **Software as a Service (SaaS)**
- ▶ The datacenter **hardware and software** is called **cloud**





▶ The **NIST** definition:

- Five **characteristics**
- Three **service models**
- Four **deployment models**

**NIST**

**National Institute of Standards and Technology**  
Technology Administration, U.S. Department of Commerce



# Cloud Characteristics

# Cloud Characteristics



On-demand  
self-service



Ubiquitous  
network  
access



Location  
transparent  
resource  
pooling



Rapid  
elasticity



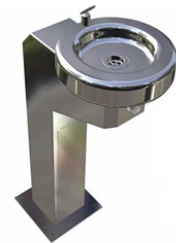
Measured  
service with  
pay per use

[<http://aka.ms/532>]



## Cloud Characteristics - On-demand Self-Service

- ▶ A consumer can **independently** provision **computing capabilities** without **human interaction** with the service provider.



On-demand self-service

## Cloud Characteristics - Ubiquitous Network Access

- ▶ Available over the **network**
- ▶ Accessed through mobile phones, laptops, ...



Ubiquitous  
network  
access



## Cloud Characteristics - Resource Pooling

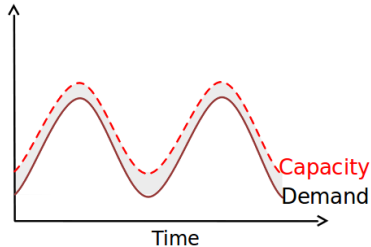
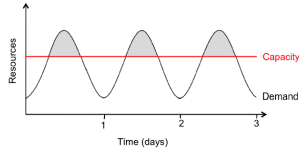
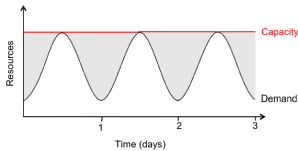
- ▶ Provider's computing resources are pooled to serve consumers
- ▶ Location transparent



Location  
transparent  
resource  
pooling

# Cloud Characteristics - Rapid Elasticity

- **Capabilities** can be rapidly and **elastically provisioned**, in some cases automatically.



Rapid elasticity

## Cloud Characteristics - Measured Service

- ▶ **Resource usage** can be monitored, controlled, and reported providing transparency for both the **provider** and **consumer**.



Measured  
service with  
pay per use

# Cloud Service Models

# Cloud Service Models



SaaS



PaaS



IaaS

[<http://aka.ms/532>]

- ▶ Assume, you just moved to a city and you are looking for a place to live.



- ▶ What is your choice?
  - Build a **new house**?
  - Buy an **empty house**?
  - Live in a **hotel**?



- ▶ Let's build a **new house!**
- ▶ You can **fully control** everything you like your new house to have.
- ▶ But that is a **hard work.**





- ▶ What if you buy an **empty house**?
- ▶ You can **customize** some part of your house.
- ▶ But never change the original architecture.



- ▶ How about living in a **hotel**?
- ▶ Living in a hotel will be a good idea if the only thing you care is about enjoying your life.
- ▶ There is **nothing you can** do with the house except living in it.





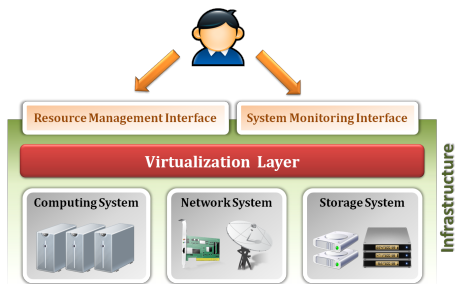
Let's translate it to Cloud Computing



## Service Models

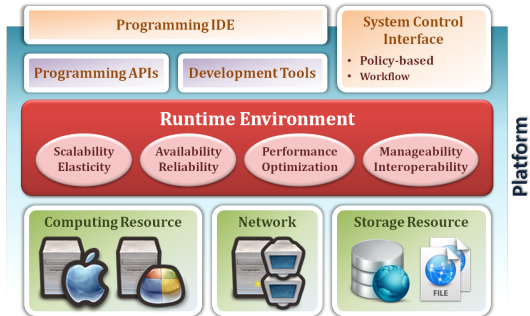
- ▶ Infrastructure as a Service (**IaaS**): similar to **building a new house**.
- ▶ Platform as a Service (**PaaS**): similar to **buying an empty house**.
- ▶ Software as a Service (**SaaS**): similar to **living in a hotel**.

- ▶ Vendor provides **resources**, e.g., processing, storage, network, ...
- ▶ Consumer is provided customized **virtual machines**.
- ▶ **Example: Amazon Web Services (EC2 instances and S3 storage)**



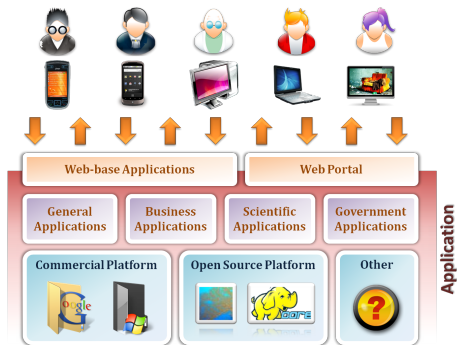
# PaaS

- ▶ Vendor provides hardware and **development environment**.
- ▶ **Example: Google app engine**

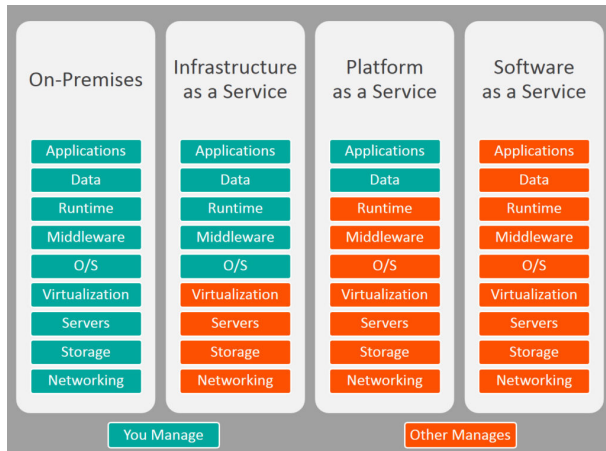


# SaaS

- ▶ Vendor provides **applications** accessed over the network.
- ▶ Example: Gmail, Github



# IaaS - PaaS - SaaS



[<https://goo.gl/xMko1z>]





# Deployment Models

# Deployment Models




**VS**






 Publically Shared  
Virtualised Resources

 Privately Shared  
Virtualised Resources

Supports multiple  
customers 

Cluster of dedicated  
customers 

 Supports connectivity  
over the internet

 Connectivity over  
internet, fibre and private network 

Suited for less  
confidential information 

 Suited for secured  
confidential information  
& core systems 

[<https://goo.gl/fWmcGK>]

# Public Cloud Infrastructure Vendors

- ▶ Amazon Web Services (AWS)
- ▶ Microsoft Azure
- ▶ Google Cloud Platform
- ▶ IBM Bluemix
- ▶ ...





## Main Services

- ▶ Computing
- ▶ Storage
- ▶ Database
- ▶ Big data analytics
- ▶ ...

# Computing Services

- ▶ Virtual machines
- ▶ Container services
- ▶ Serverless compute



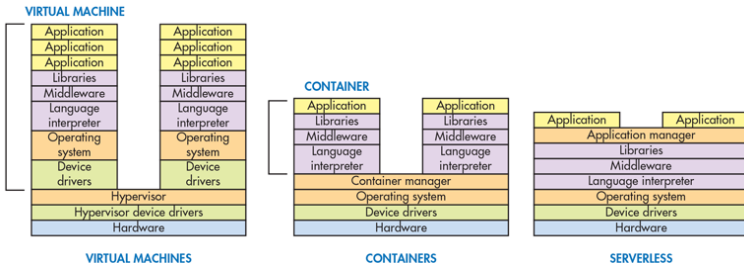
VM



Container



Serverless



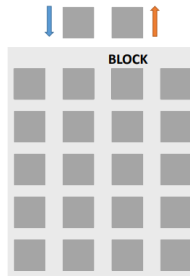
# Storage Services

- ▶ File storage
- ▶ Block storage
- ▶ Object storage

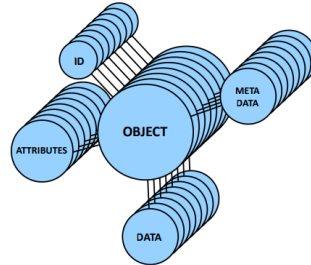
**File Storage**



**Block Storage**

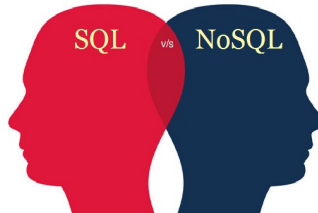


**Object Storage**



# Database Services

- ▶ Relational Database Management Services (RDBMS)
- ▶ NoSQL databases
- ▶ In-Memory data services









# Big Data

# What is Big Data?



[<https://www.sue-anderson.com.au/index.php/2017/08/18/cursing-curious-work>]



# Big Data

Big data is the data characterized by 4 key attributes: volume, variety, velocity and value.

**Buzzwords**

**ORACLE®**



# Big Data in Simple Words



**DevOps Borat**

@DEVOPS\_BORAT

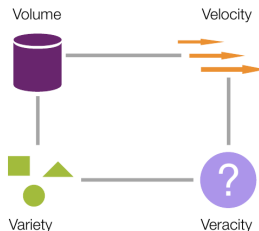
Small Data is when is fit in RAM.  
Big Data is when is crash because  
is not fit in RAM.

2/6/13, 8:22 AM



# The Four Dimensions of Big Data

- ▶ **Volume:** data size
- ▶ **Velocity:** data generation rate
- ▶ **Variety:** data heterogeneity
- ▶ This 4th **V** is for **V**acillation:  
Veracity/Variability/Value



# How Much Data?

## THE INTERNET IN **2023** EVERY MINUTE



Created by: eDiscovery Today & LTMG



# How To Store and Process Big Data?





# Problem

- ▶ Traditional platforms **fail** to show the expected performance.
- ▶ Need **new systems** to **store and process** large-scale data

## Scale Up vs. Scale Out (1/2)

- ▶ Scale **up** or scale **vertically**: adding **resources** to a **single** node in a system.
- ▶ Scale **out** or scale **horizontally**: adding **more nodes** to a system.

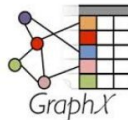


## Scale Up vs. Scale Out (2/2)

- ▶ Scale **up**: more **expensive** than scaling out.
- ▶ Scale **out**: more challenging for **fault tolerance** and **software development**.



APACHE  
**HBASE**



**Storm**



**S4** distributed stream  
computing platform

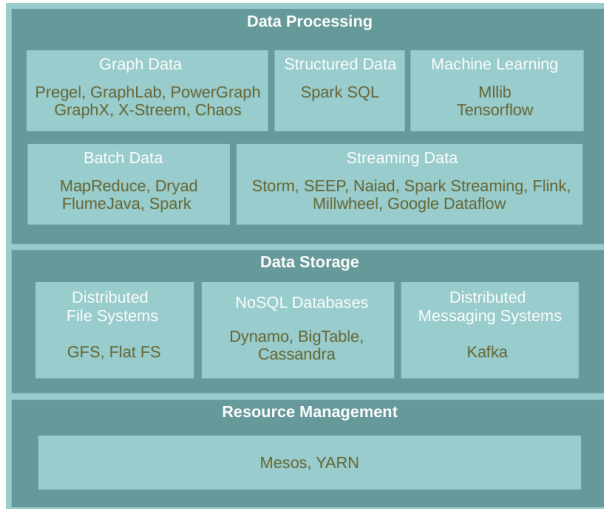


Google Cloud Platform



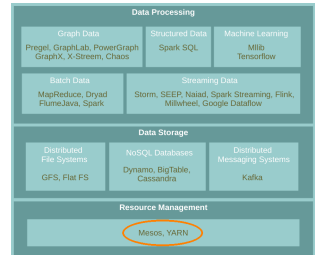


# Big Data Stack



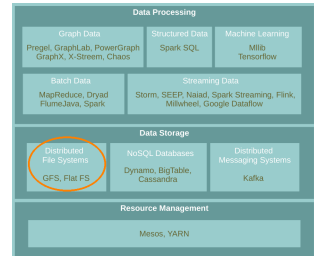
# Resource Management

- ▶ Manage resources of a cluster
- ▶ Share them among the platforms
- ▶ Mesos, YARN, Borg, ...



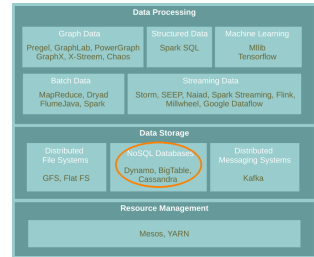
# Data Storage - Distributed File Systems

- ▶ Store and retrieve **files** on/from distributed disks
- ▶ **GFS, HDFS, FlatFS, ...**



# Data Storage - NoSQL Databases

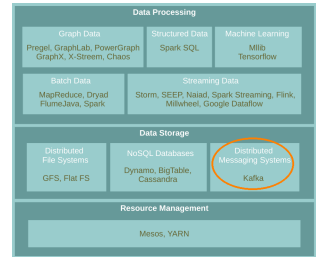
- ▶ BASE instead of ACID
- ▶ BigTable, Dynamo, Cassandra, ...





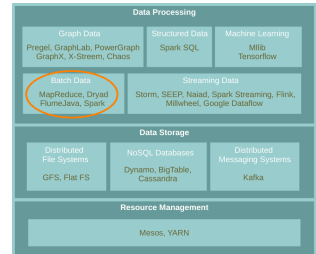
# Data Storage - Messaging Systems

- ▶ Store streaming data
- ▶ Kafka, Flume, ActiveMQ, ...



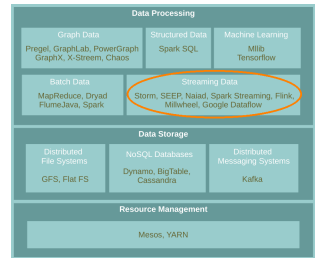
# Data Processing - Batch Data

- ▶ Process data-at-rest
- ▶ Data-parallel processing model
- ▶ MapReduce, FlumeJava, Spark, ...



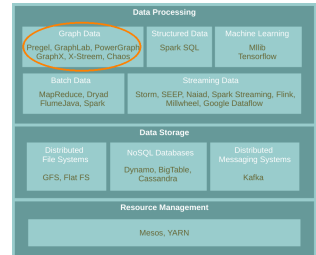
# Data Processing - Streaming Data

- ▶ Process data-in-motion
- ▶ Storm, Flink, Spark Streaming, ...



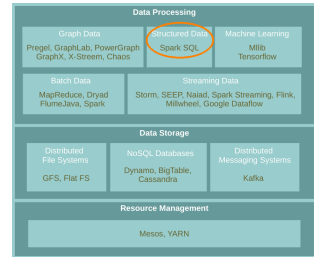
# Data Processing - Linked Data (Graph)

- ▶ Graph-parallel processing model
- ▶ Vertex-centric and Edge-centric programming model
- ▶ Pregel, GraphLab, GraphX, ...



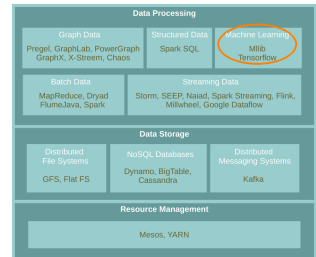
# Data Processing - Structured Data

- ▶ Take advantage of **schemas** in data to process
- ▶ **Hive, Spark SQL, ...**



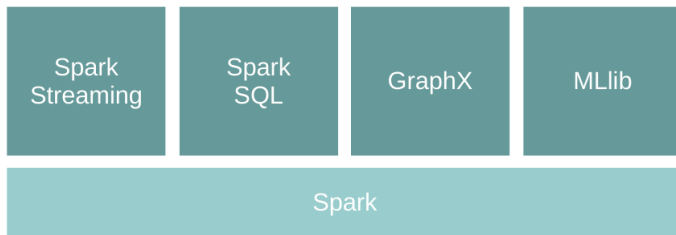
# Data Processing - Machine Learning

- ▶ Data analysis, e.g., supervised and unsupervised learning
- ▶ Mahout, TensorFlow, MLlib, ...





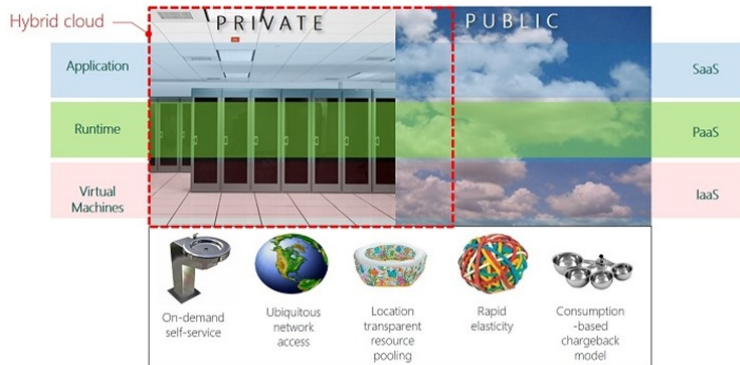
# Spark Processing Engine



# Summary

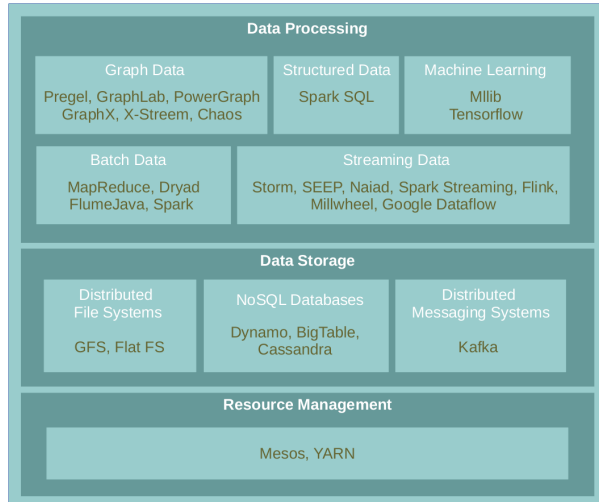


# Summary



[<http://aka.ms/532>]

# Summary





## References

- ▶ D. Sikeridis et al., A Comparative Taxonomy and Survey of Public Cloud Infrastructure Vendors, arXiv preprint arXiv:1710.01476, 2017.
- ▶ A. Fox et al., Above the clouds: A berkeley view of cloud computing, UCB/EECS 28.13 (2009): 2009.
- ▶ P. Mell et al., The NIST definition of cloud computing, 2011.

Questions?