# Fundamentals of Data Engineering

## Trainer: Nilesh Ghule

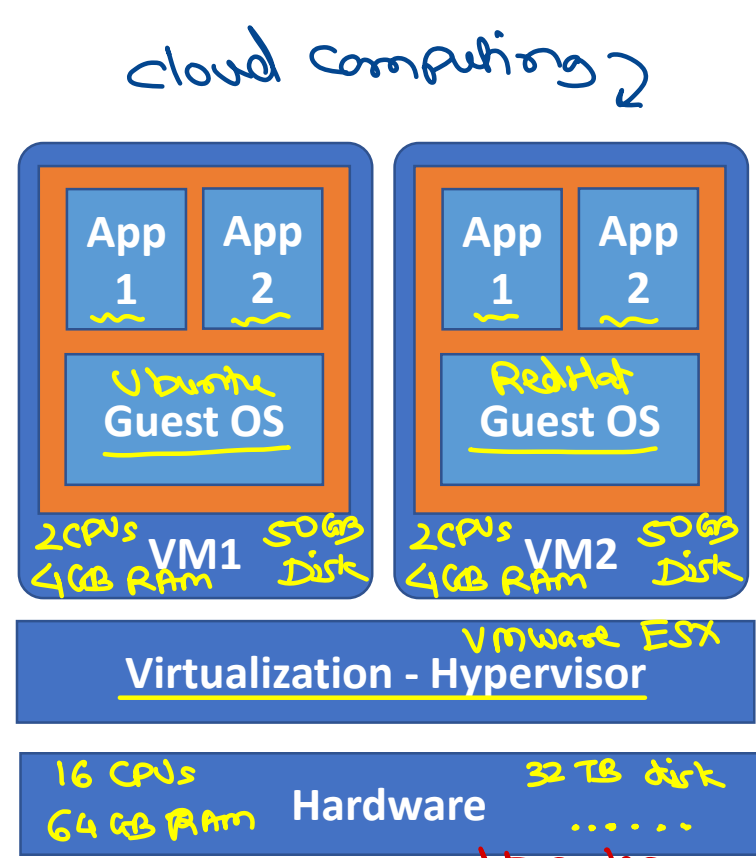# Virtualization vs Containerization

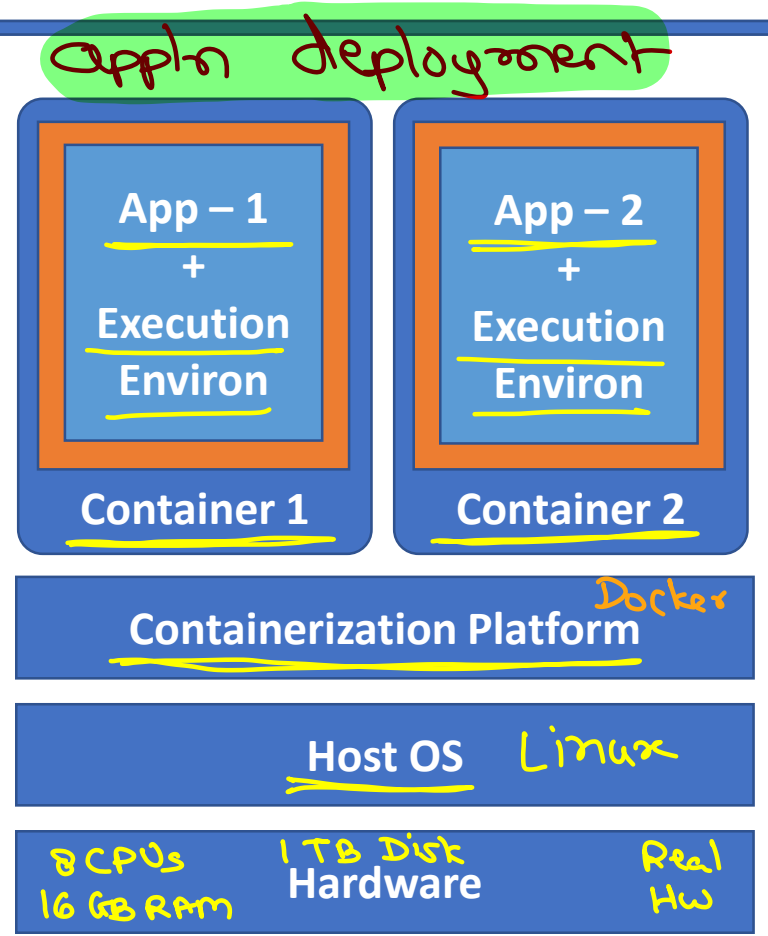*appln deployment*

**Type-II Virtualization**

| App 1 | App 2 |
| --- | --- |
| *Ubuntu* Guest OS | |

VM1 — *2CPUs 4GB RAM* — *50GB Disk*

| App 1 | App 2 |
| --- | --- |
| *Win XP* Guest OS | |

VM2 — *2CPUs 4GB RAM* — *50GB Disk*

Virtualization - Hypervisor — *VirtualBox*

Host OS *(windows)*

Hardware — *8 CPUs 16 GB RAM* — *1 TB Disk* — *Real HW*

**Type-II Virtualization**

VMWare, VirtualBox, KVM, …

---

*cloud computing*

| App 1 | App 2 |
| --- | --- |
| *Ubuntu* Guest OS | |

VM1 — *2CPUs 4GB RAM* — *50GB Disk*

| App 1 | App 2 |
| --- | --- |
| *RedHat* Guest OS | |

VM2 — *2CPUs 4GB RAM* — *50GB Disk*

Virtualization - Hypervisor — *VMWare ESX*

Hardware — *16 CPUs 64 GB RAM* — *32 TB disk* …….

*data centers*

**Type-I Virtualization**

VMWare ESX, XEN, Hyper-V, …

---

**Containerization**

| App – 1 + Execution Environ |
| --- |
| Container 1 |

| App – 2 + Execution Environ |
| --- |
| Container 2 |

Containerization Platform — *Docker*

Host OS — *Linux*

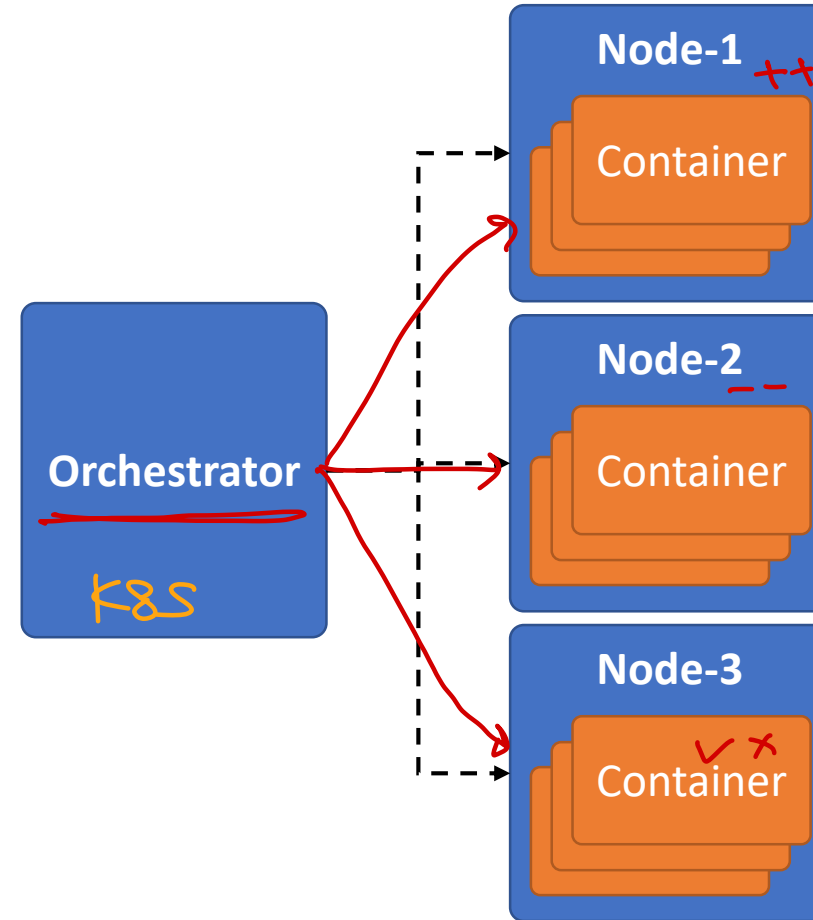Hardware — *8 CPUs 16 GB RAM* — *1 TB Disk* — *Real HW*

**Containerization**

Docker, Podman, rkt, …

# Orchestration

- **Container Orchestration auto increase or decrease containers to handle change in workloads/demands. It also handles container failure (re-start).**
- **Ex: Docker swarm, Kubernetes, …**

Node-1 ++

Container

Node-2 --
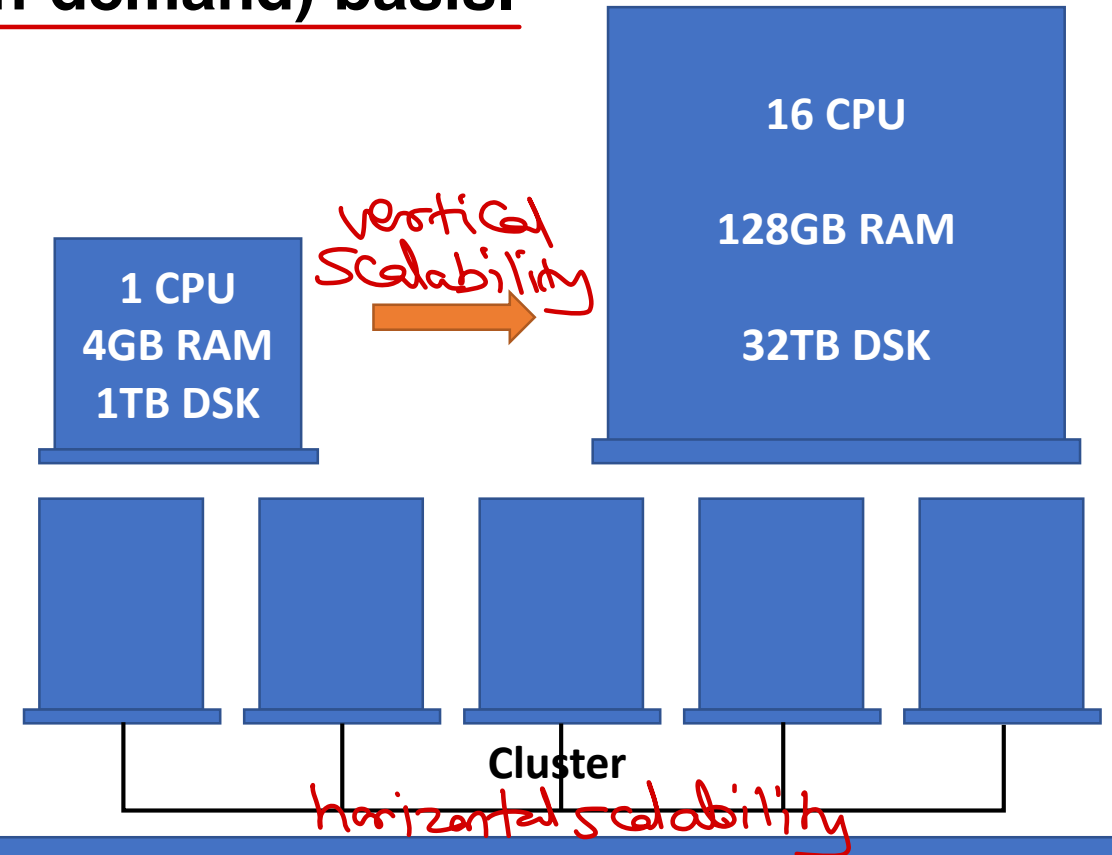
Container

**Orchestrator**

K8S

Node-3

Container ✓ ✗

DevOps = Docker + K8S
+ GIT + Jenkins + …
CI/CD pipeline
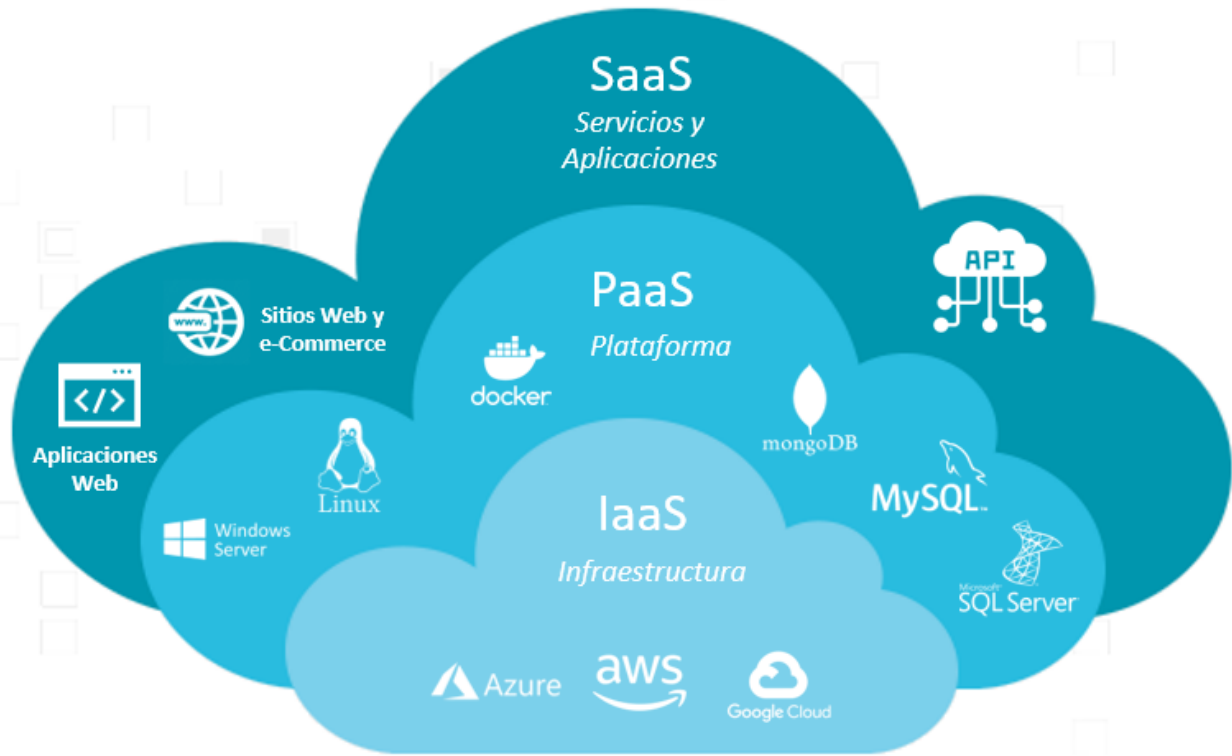
# Scalability and Elasticity

- Scalability is "ability of system / application to perform well under an increased or expanding workload".

- The resource usage is increased or decreased as per workload.

- Vertical scalability / Up scaling:
  - Increasing single system (hardware) resources in order to handle higher loads.
  - Need to handle SPOF (single point of failure) by adding backup system.

- Horizontal scalability / Out scaling:
  - Adding new systems/notes into the cluster in order to handle higher loads.
  - More economical solution with higher complexity.

- Elastic: Cloud systems are designed to increase/decrease load as per workload.
  *resource*

- Cloud payments are usually pay-per-use (on-demand) basis.

1 CPU
4GB RAM
1TB DSK

*vertical Scalability*

16 CPU

128GB RAM

32TB DSK

Cluster

*horizontal scalability*

# Cloud Service Models

AWS (Amazon), GCP (Google), Azure (Microsoft), Smart Cloud (IBM), ...



- **IaaS: Infrastructure as-a Service**
  - **AWS EC2, S3, VPC**
    - vm     storage    network
- **PaaS: Platform as-a Service**
  - **Beanstalk, SageMaker,**
- **SaaS: Software as-a Service**
  - **Gmail, Drive, Facebook, LinkedIn, Netflix**
- **DaaS: Database as-a Service**
  - **RDS, Aurora, Atlas, DynamoDb**
- **FaaS: Function as-a Service**
  - **Lambda, Google functions**

Data centers → Huge infrastructure with lots of real computers and network, ...
→ Computer/machine
→ Network
→ Storage

# Big Data & Analytics Spectrum

PG-DBDA course

- **Data storage**
  - RDBMS & NoSQL databases
  - Data warehouse
  - S3, DFS, …
- **Data Analysis & visualizations**
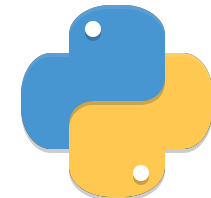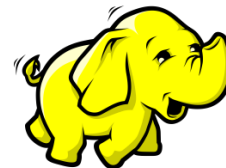  - Data Visualizations
  - Business reports
- **Artificial Intelligence, Data Science & Data mining**
  - Mathematics, Statistics & Computer algorithms
  - Machine learning & Deep learning
  - R Programming, Python
- **(Big) Data Engineering**
  - Hadoop, Hive, Spark, Kafka, BigTable, …
  - Java, Scala, Python, SQL
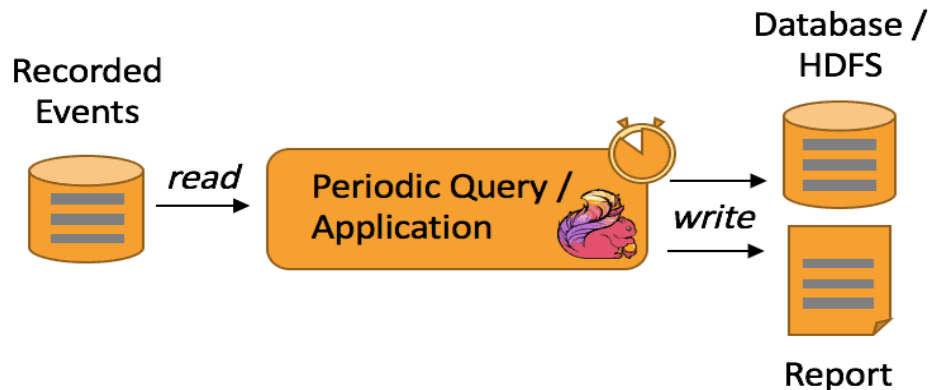- **Infrastructure**
  - Linux, Cloud Computing
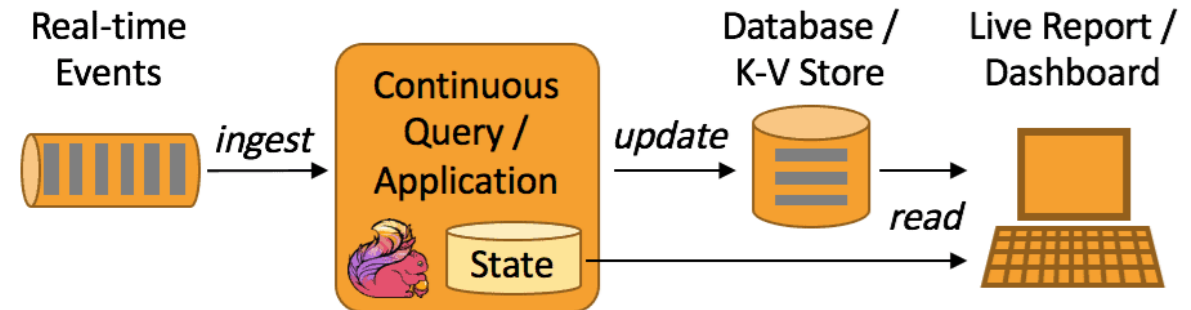
# Batch processing vs Stream processing

- **Processing finite set of data (data at rest).**

- **Incremental data load is managed by programmer.**

- **Cluster planned as per data size. High throughput.**

- **Job run once per batch.**

- **Processing live stream of data (data in motion).**

- **Data processing is managed by the framework.**

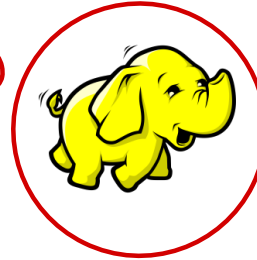- **Less throughput.**

- **Job is running forever.**

**Batch Processing**

Recorded Events → read → Periodic Query / Application → write → Database / HDFS

→ Report

**Stream Processing**

Real-time Events → ingest → Continuous Query / Application (State) → update → Database / K-V Store → read → Live Report / Dashboard
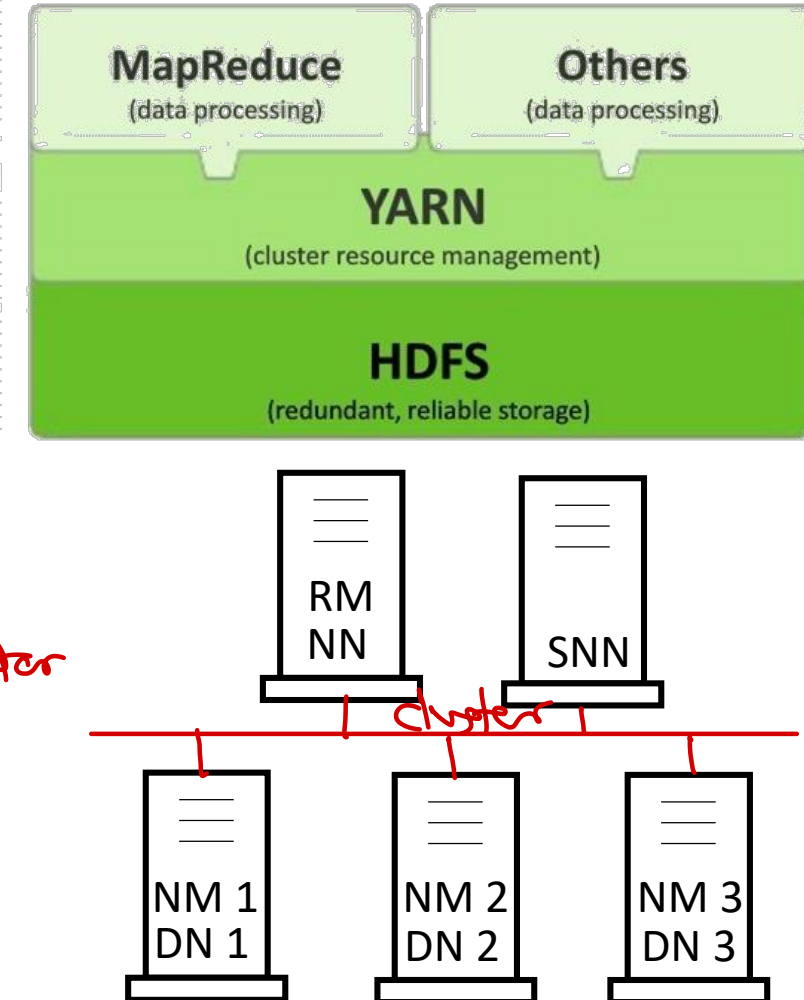
# Apache Hadoop

*Google* → Distributed Storage = Google File System - 2003
→ Distributed Computing = Map Reduce - 2004

- **Hadoop is developed by <u>Doug cutting.</u>**
  - **<u>Web crawler – Nutch</u>**
  - **<u>Distributed computing and storage</u> needed to process huge data produced by the crawler.**
  - **Joined <u>Yahoo.</u> Developed and <u>open sourced</u> under Apache license.** → 2006

- **Hadoop** — Hadoop Distributed File System (like GFS)
  - **Distributed storage: HDFS**
  - **Distributed computing : Map-reduce**
  - **Cluster manager: YARN** - Yet Another Resource Negotiator

- **Hadoop is like a Kernel/Platform on which many different applications are built (eco-systems).**
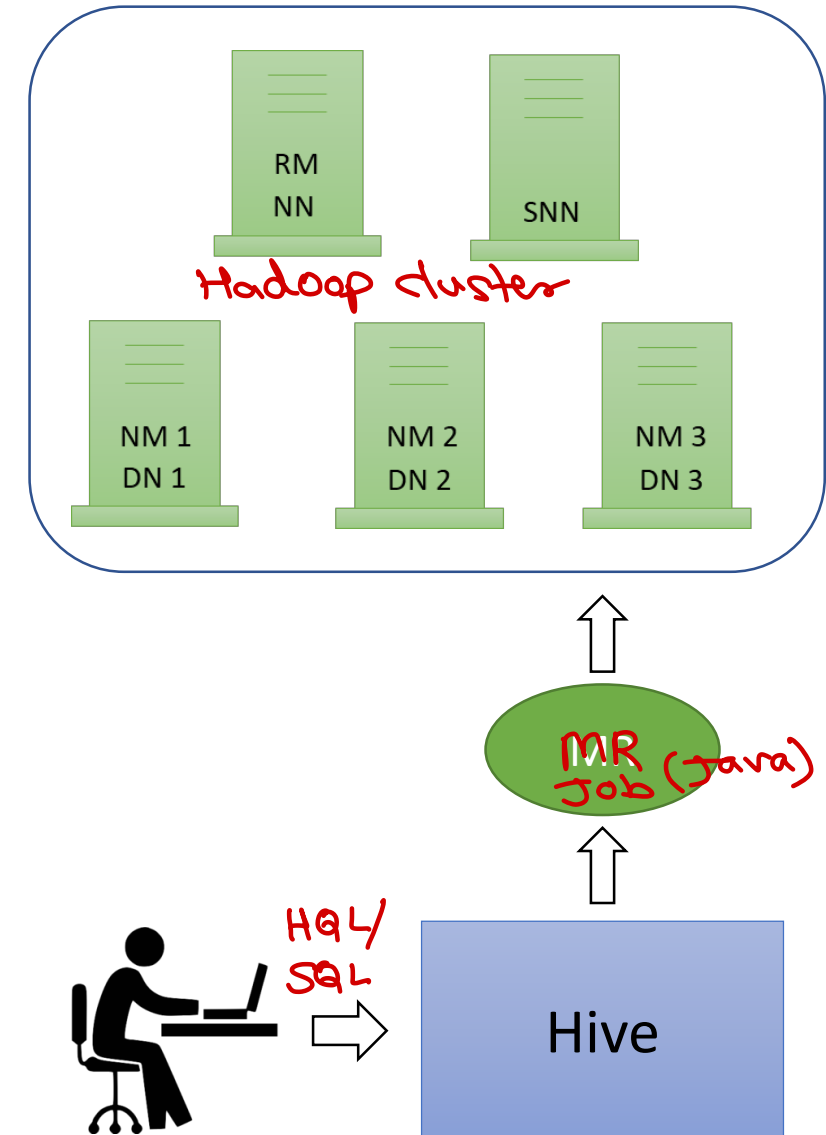
**HADOOP 2.0**

| MapReduce (data processing) | Others (data processing) |
|---|---|

YARN (cluster resource management)

HDFS (redundant, reliable storage)

RM NN    SNN

Cluster

NM 1 DN 1    NM 2 DN 2    NM 3 DN 3

# Apache Hive

- **Developed by Facebook (2007)**
- **Client software that convert Hive QL queries to MapReduce.**
- **Hive QL is similar to SQL with many extended features.**
- **Hive manages structured data.**
- **Hive is data warehouse (OLAP) built for Hadoop.**
  - **Data storage = HDFS**
  - **Metadata = RDBMS**
  - **Data processing = Map-reduce or Spark or Tez.**

*High speed execution*

RM NN

SNN

Hadoop cluster

NM 1 DN 1

NM 2 DN 2

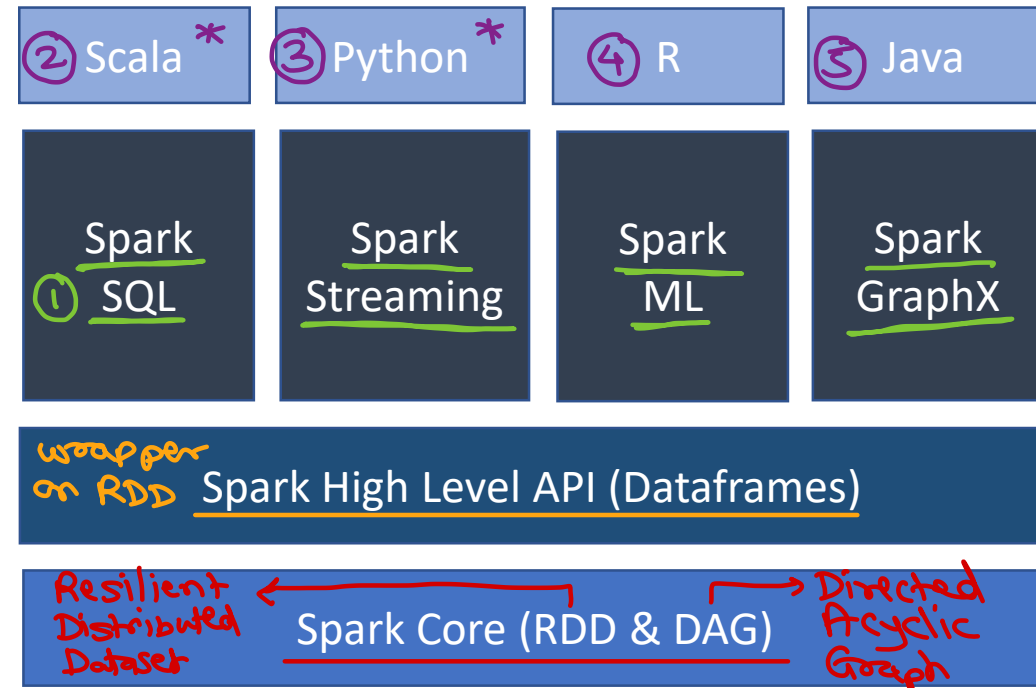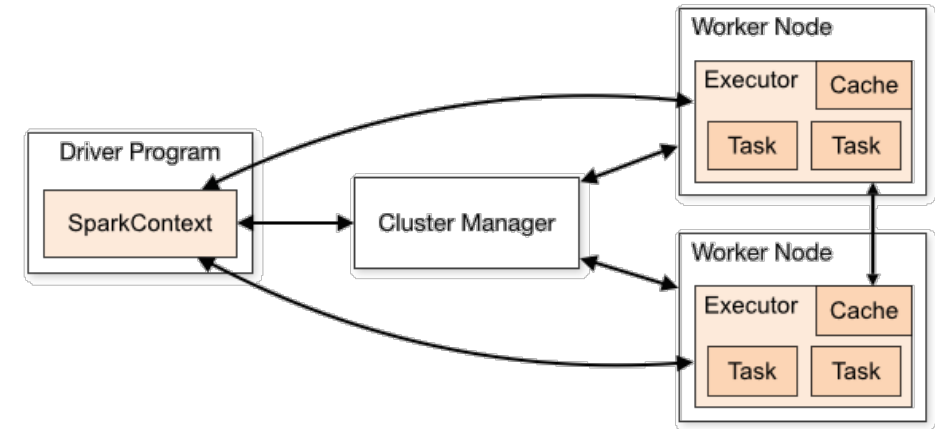NM 3 DN 3

MR Job (Java)

HQL/SQL

Hive

# Apache Spark

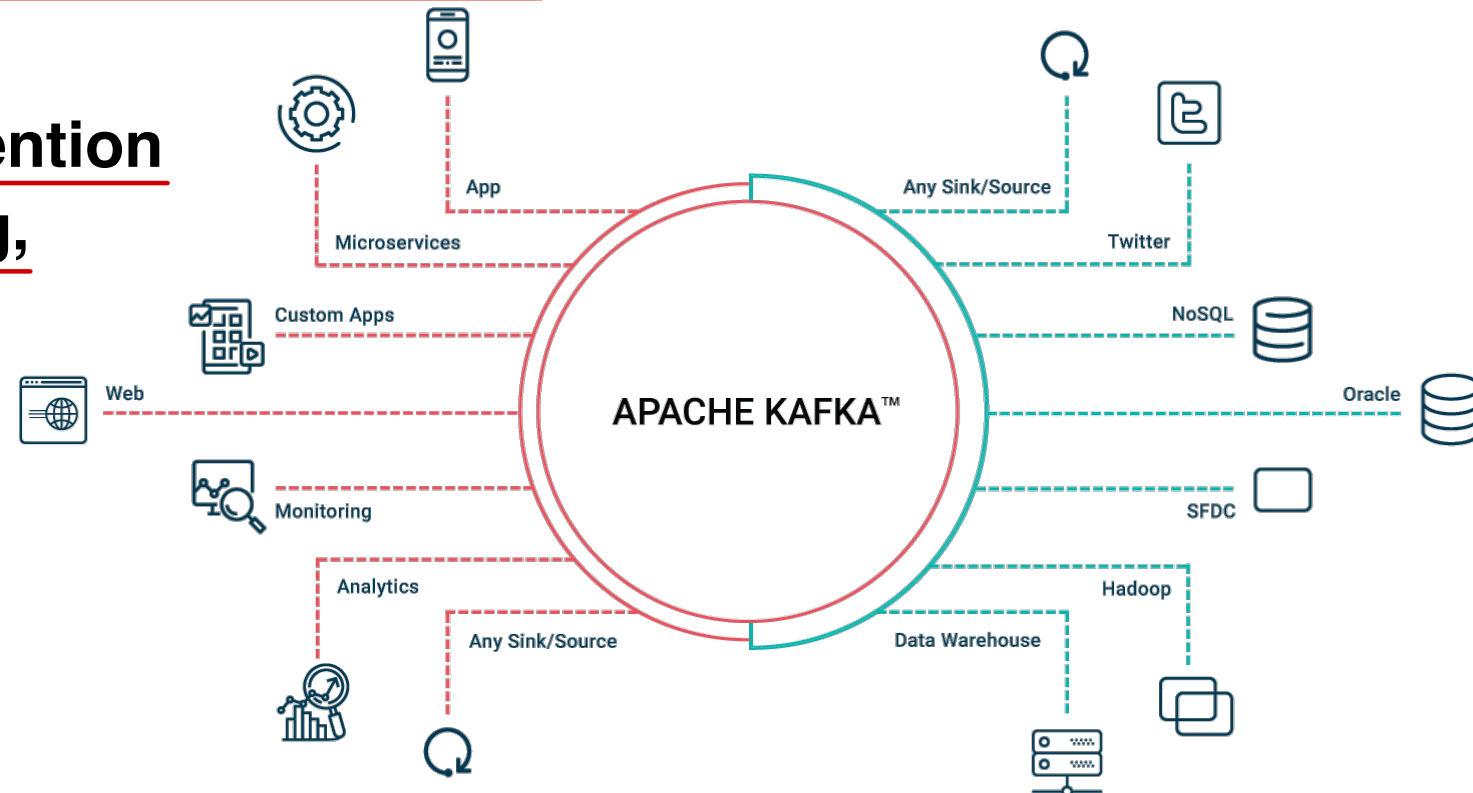Can work with any Storage: HDFS, S3, Local, ...

- **Spark is Distributed computing framework, that can process huge amount of data.**
- **Spark can be used as eco-system of Hadoop or can be used as independent distributed computing framework.**
- **Developed by UCB AMPlabs division.**
- **Further developed/maintained by DataBricks.**

  → Algorithms,
  → Machines,
  → People

- **Popular Spark vendors**
  - **DataBricks, AWS EMR, Cloudera, MapR**
- **Spark Toolkit**

### Diagram (right side)

Driver Program
SparkContext ↔ Cluster Manager

Worker Node
Executor | Cache
Task | Task

Worker Node
Executor | Cache
Task | Task

| ② Scala * | ③ Python * | ④ R | ⑤ Java |
|---|---|---|---|
| Spark ① SQL | Spark Streaming | Spark ML | Spark GraphX |

wrapper on RDD — Spark High Level API (Dataframes)

Resilient Distributed Dataset — Spark Core (RDD & DAG) — Directed Acyclic Graph
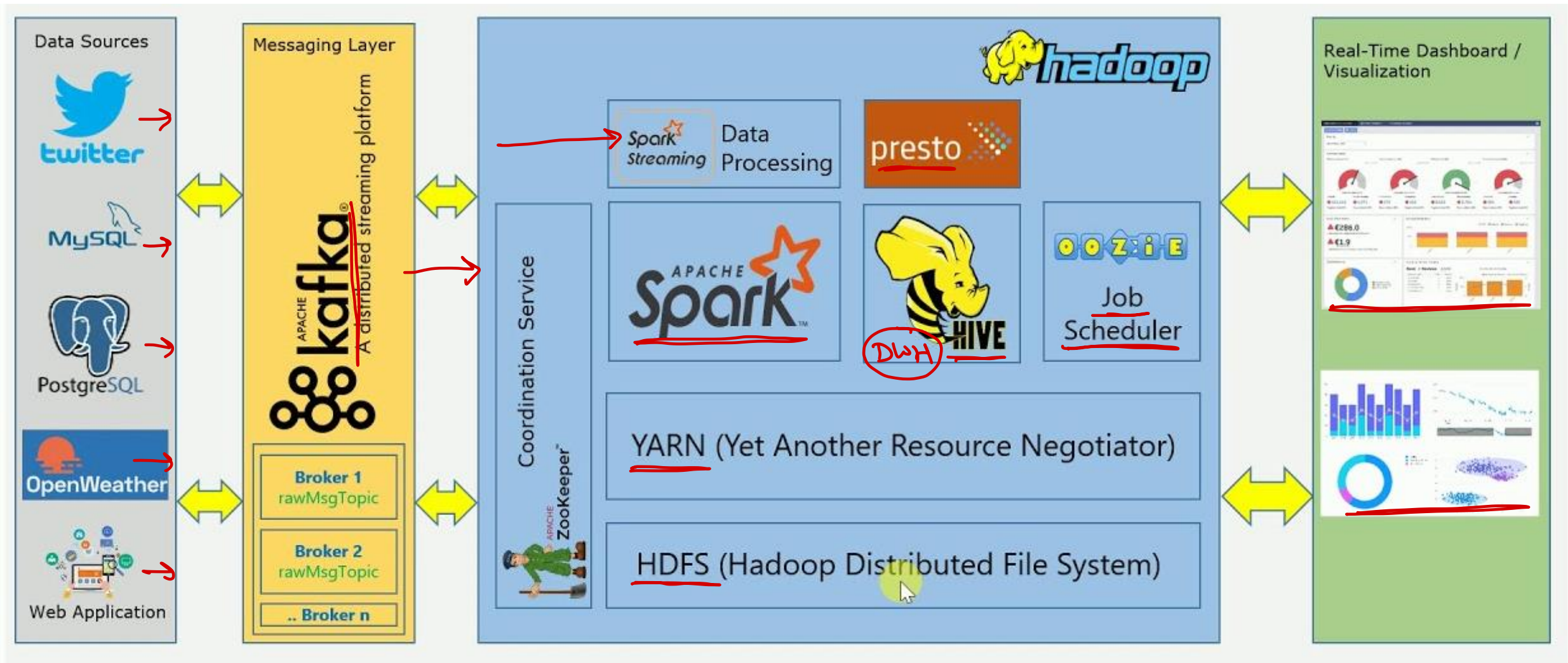
# Apache Kafka

- **Kafka is a distributed messaging system.**

- **Developed at LinkedIn and open sourced in 2011.**

- **Used by LinkedIn, Twitter, Uber, airbnb, …**

- **Advantages**
  - **Scalable, Durable, Finite retention**
  - **Low latency, Strong ordering,**
  - **Exact once delivery**

- **Applications**
  - **Stream processing**
  - **Notifications.**

# Real time dashboard reference architecture

# Big Data domains & opportunities

- **Domains: Health-care, Retails, Trading/Share market, Finance, Security, Fraud, Search engines, Log Analysis, Telecom, Traffic Control, Manufacturing and lot more.**

- **Big Data is all about :- Think, Collect, Manage, Analyze, Summarize, Visualize, Discover Knowledge and Take Decisions.**

- **Job profiles:**
  - **Business Analyst/Intelligence**
  - **Database engineer / DWH**
  - **Big Data engineer**
  - **Data operations**
  - **Big Data Architect**

- **The sexiest job in the 21st century require a mixture of multidisciplinary abilities and suitable candidates must be prepared to learn & develop constantly.                                             -Ronald Van Loon**

https://www.youtube.com/live/BxwpqnQ6BgQ?si=55cmOUDfiIGDAsLY

*SunBeam*
*Big Data* | *Important.*
*Webinar* |



# Thank you!

## Nilesh Ghule <nilesh@sunbeaminfo.com>