

## Lecture 3

# Big Data Processing Frameworks (Hadoop and Spark)

Dr. Lydia Wahid

## Agenda

### Basic Concepts

### Apache Hadoop

### Hadoop ecosystem

### Hadoop YARN

### Apache Spark

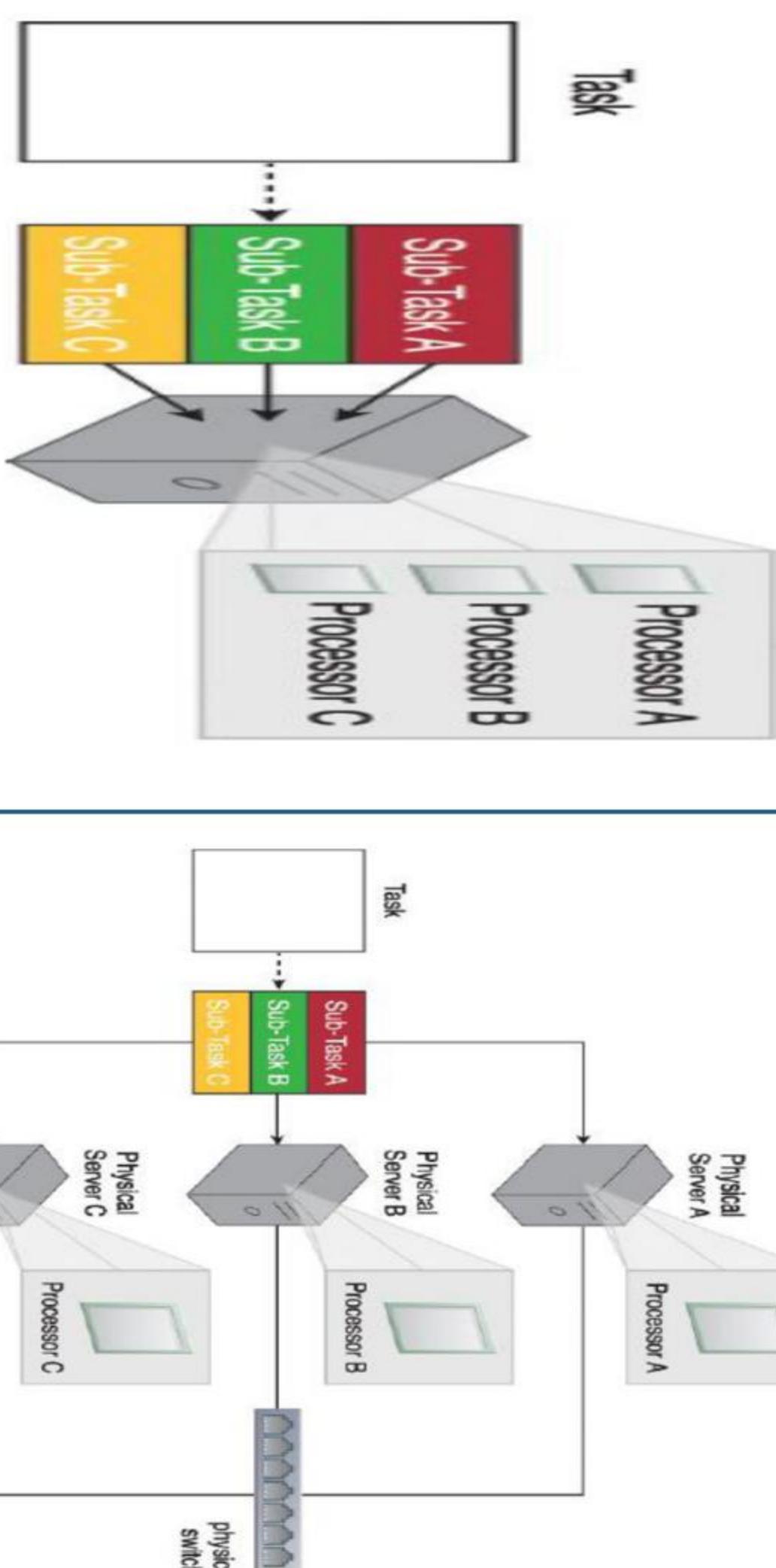
### Hadoop VS. Spark

## Basic Concepts

### Basic Concepts

#### Parallel Data Processing

- Parallel data processing involves the **simultaneous execution** of multiple sub-tasks that collectively comprise a larger task.
- The goal is to reduce the execution time by dividing a single larger task into multiple smaller tasks that run concurrently.



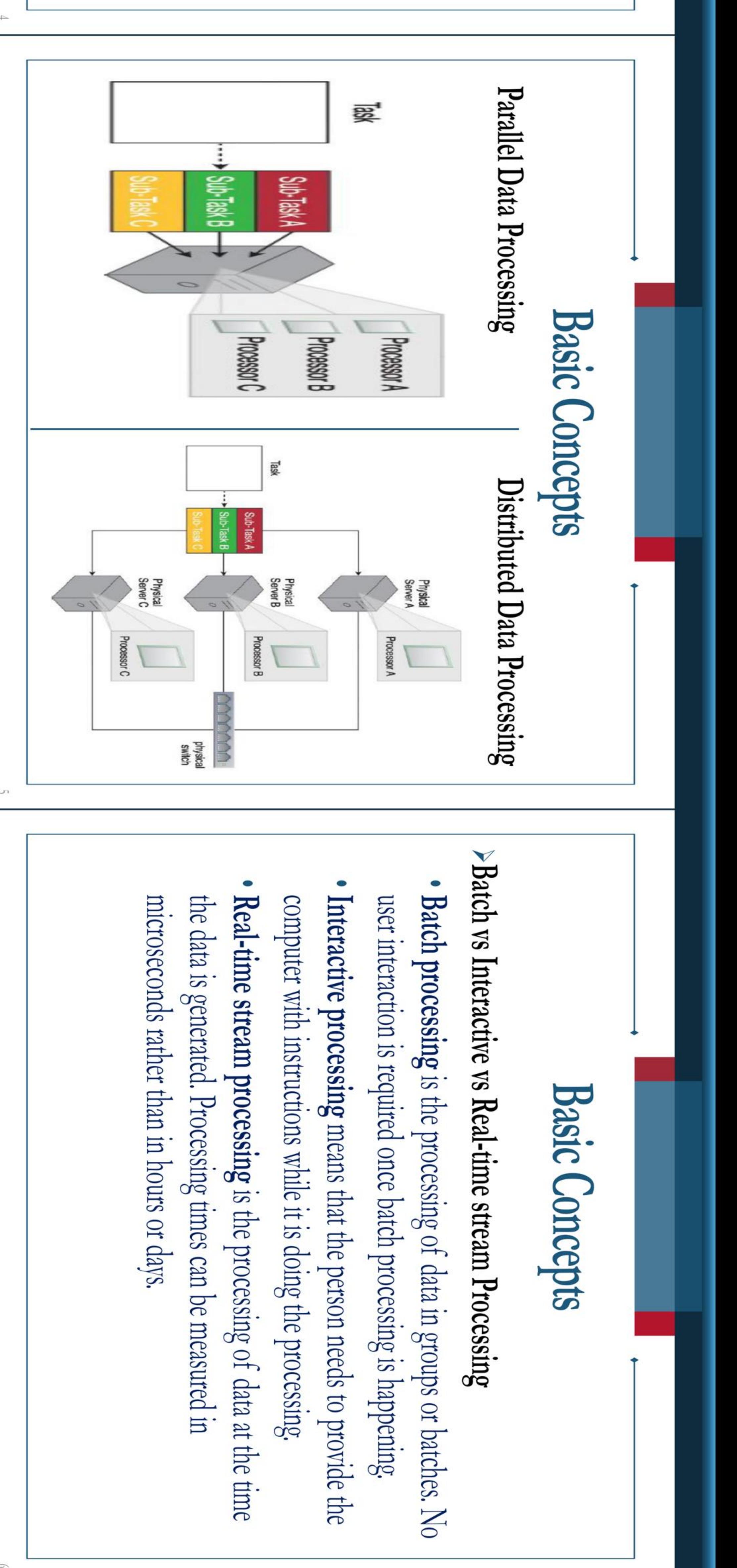
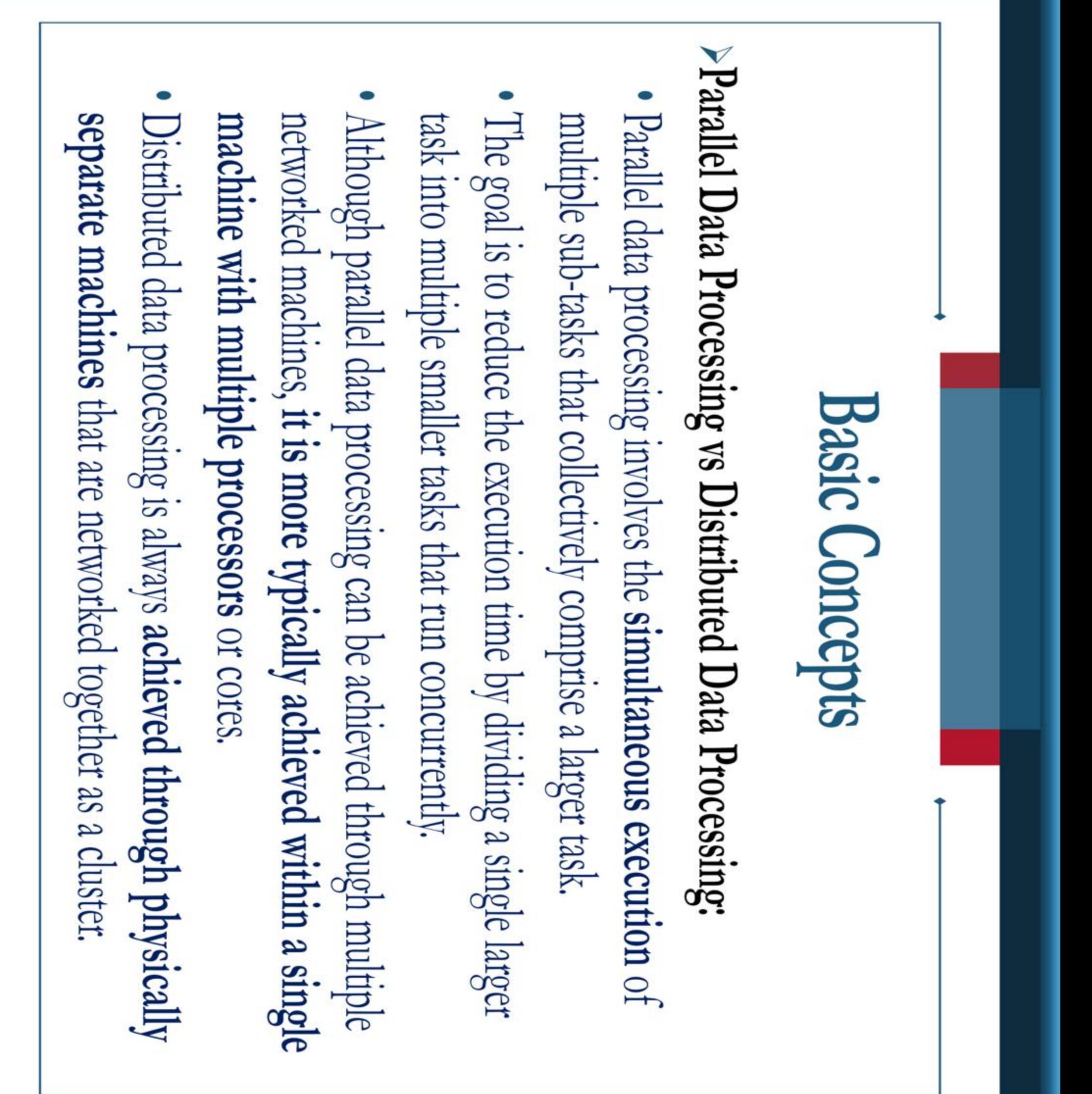
#### Distributed Data Processing

- 
- ```
graph LR; Task[Task] --> SubTaskA[Sub-Task A]; SubTaskA --> PhysicalServerA[Physical Server A]; PhysicalServerA --> ProcessorA[Processor A]; SubTaskB[Sub-Task B] --> PhysicalServerB[Physical Server B]; PhysicalServerB --> ProcessorB[Processor B]; SubTaskC[Sub-Task C] --> PhysicalServerC[Physical Server C]; PhysicalServerC --> ProcessorC[Processor C];
```
- Although parallel data processing can be achieved through multiple networked machines, it is more typically achieved **within a single machine** with multiple processors or cores.
  - Distributed data processing is always achieved **through physically separate machines** that are networked together as a cluster.

### Basic Concepts

#### Batch vs Interactive vs Real-time stream Processing

- **Batch processing** is the processing of data in groups or batches. No user interaction is required once batch processing is happening.
- **Interactive processing** means that the person needs to provide the computer with instructions while it is doing the processing.
- **Real-time stream processing** is the processing of data at the time the data is generated. Processing times can be measured in microseconds rather than in hours or days.



# Apache Hadoop

# Apache Hadoop - Architecture

- Hadoop is an Apache open-source framework written in java that allows distributed processing of large datasets across clusters of computers.

- Apache Hadoop**

➤ A cluster is a configuration of nodes that interact to perform a specific task.



Apache Hadoop

➤ Hadoop is designed to scale up from single server to thousands of machines, each offering local computation and storage.

- # Hadoop Architecture:

  - It consists of two main components:
    1. Storage (Hadoop Distributed File System)
      - MapReduce consume data from HDFS. HDFS creates multiple replicas of data blocks and distributes them on the nodes in a cluster. This distribution enables reliable and extremely rapid computations.
    2. Processing/Computation (MapReduce)
      - Capable of processing enormous data in parallel on large clusters of computation nodes. It performs two main tasks: *map* and *reduce*

# Apache Hadoop - Architecture

- Hadoop Distributed File System (HDFS) is managed with the *master-slave* architecture.

- **NameNode:** This is the *master* of the HDFS system. It maintains the file system tree and the metadata for all the files and directories present in the system.
  - **DataNode:** These are *slaves* that are deployed on each machine and provide actual storage. They are responsible for serving read-and-write data requests from the clients. The internal mechanism of HDFS divides the file into one or more blocks; these blocks are stored in a set of data nodes.

Note: HDFS is not a database, but it is a distributed file system that can store huge volume of data sets across a cluster of computers to be processed.

# Apache Hadoop - Architecture

- MapReduce is managed with *master-slave* architecture included with the following components:

- **JobTracker:** This is the *master* node of the MapReduce system, which manages the jobs and resources in the cluster. The JobTracker tries to schedule the maps to specific nodes in the cluster, ideally the nodes that have the data.
  - **TaskTracker:** These are the *slaves* that are deployed on each machine. They are responsible for running the map and reduce tasks as instructed by the JobTracker.

# Apache Hadoop - Architecture

- How MapReduce tasks are assigned to specific nodes in the cluster:

1. Client applications submit jobs to the JobTracker.
  2. The JobTracker talks to the NameNode to determine the location of the data (DataNode)
  3. The JobTracker locates TaskTracker nodes with available slots at or near the data (DataNode)
  4. The JobTracker submits the work to the chosen TaskTracker nodes.
  5. The TaskTracker nodes are monitored. If they do not submit signals often enough, they are considered to have failed and the work is scheduled on a different TaskTracker.
  6. A TaskTracker will notify the JobTracker when a task fails. The JobTracker decides what to do then: it may resubmit the job elsewhere, it may mark that specific record as something to avoid, and it may even blacklist the TaskTracker as unreliable.
  7. When the work is completed, the JobTracker updates its status.
  8. Client applications can poll the JobTracker for information.

## Number of Maps

- The number of maps is usually driven by the number of blocks of the input files.
- 10TB of input data and a blocksize of 128MB, will end up with 82,000 maps, unless **Configuration.set(MRJobConfig.NUM\_MAPS, int)** is used to set it even higher.
- The right level of parallelism for maps seems to be around **10-100 maps per-node**, although it has been set up to 300 maps for very cpu-light map tasks.
- Block size can also be changed to adjust the number of blocks.

## Apache Hadoop - HDFS Features

- **Replication** - Due to some unfavorable conditions, the node containing the data may be lost. To overcome such problems, HDFS always maintains the copy of data on more than one machine.
- **Fault tolerance** - In HDFS, the fault tolerance signifies the robustness of the system in the event of failure. Due to Replication, HDFS is highly fault-tolerant that if any machine fails, the other machine containing the copy of that data automatically become active.
- **Portability** - HDFS is designed in such a way that it can be easily portable from platform to another.

## Number of Reducers

- In Hadoop, the default number of reducers is **one**.
- In this phase the `reduce(WritableComparable, Iterable<Writable>, Context)` method is called for each `<key, (list of values)>` pair in the grouped inputs.
- The number of reducers for the job is set by the user via **Job.setNumReduceTasks(int)**
- Increasing the number of reducers increases the framework overhead, but increases load balancing and lowers the cost of failures.
- It is legal to set the number of reduce-tasks to zero if no reduction is desired.

## Apache Hadoop - Hadoop Operation Modes

- **Hadoop can run in 3 different modes:**
  1. **Standalone Mode** (Single machine Single process):  
By default, Hadoop is configured to run in a non-distributed mode. It runs as a single Java process. Instead of HDFS, this mode utilizes the local file system. This mode is useful for debugging.
  2. **Pseudo-Distributed Mode** (Single machine Multiple processes):  
Hadoop can also run on a single machine in a Pseudo Distributed mode. In this mode, each Hadoop process runs as separate java process. Here HDFS is utilized for input and output.
  3. **Fully Distributed Mode** (Multiple machines Multiple processes):  
This is the production mode of Hadoop. In this mode, Hadoop is distributed across multiple machines. Therefore, separate java processes are present. This mode offers fully distributed computing capability, reliability, fault tolerance and scalability.

## Apache Hadoop - HDFS Features

- **High Scalability** - HDFS is highly scalable as it can have hundreds of nodes in a single cluster.
- **Handling Huge datasets** – Due to the high scalability, HDFS can manage applications having huge datasets.
- **Distributed data storage** - This is one of the most important features of HDFS that makes Hadoop very powerful. Here, data is divided into multiple blocks and stored into nodes.
- **Hardware at data** – A requested task is done efficiently as the computation takes place near the data. This reduces the network traffic and increases the throughput, especially where huge datasets are involved.

## Hadoop ecosystem



# Apache Hadoop - Hadoop ecosystem

## Apache Hive

- It is a software developed by Facebook that facilitates reading, writing, and managing large datasets residing in distributed storage using SQL.
- It allows users to fire queries in SQL-like languages, such as **HiveQL**.

## Apache Pig

- It is platform for creating programs that run on Apache Hadoop. The language for this platform is called **Pig Latin**.
- Apache Pig has been developed by Yahoo. Currently, Yahoo and Twitter are the primary users of Pig.

19

# Apache Hadoop - Hadoop ecosystem

## Apache Hbase:

- It is the Hadoop database, a distributed, scalable, big data store. This allows random, real-time read/write access to Big Data. Apache HBase is an open-source, distributed, non-relational database modeled after Google's Bigtable.
- The following are some companies using HBase: Yahoo, Twitter, and stumble upon (This is a personalized recommender system, realtime data storage, and data analytics platform).

## Apache Impala:

- With Impala, you can query data, whether stored in HDFS or Apache HBase – including SELECT, JOIN, and aggregate functions – in real time.
- Impala directly access the data through a specialized distributed query engine. The result is order-of-magnitude faster performance than Hive

## Apache Sqoop:

- Apache Sqoop is a mutual data tool for importing data from the relational databases to Hadoop HDFS and exporting data from HDFS to relational databases.
- It works together with most modern relational databases, such as Microsoft SQL Server, MySQL, and Oracle.

20

# Apache Hadoop - Hadoop ecosystem

## Apache Solr:

- It is an open-source enterprise search platform.
- Solr is highly reliable, scalable and fault tolerant, providing distributed indexing, replication and load-balanced querying, automated failover and recovery, centralized configuration and more.
- This allows building web application with powerful search capabilities.
- Solr powers the search and navigation features of many of the world's largest internet sites

21

## Apache Zookeeper:

- It is a service that enables highly reliable distributed coordination.
- It is a centralized service for maintaining configuration information, naming, and providing distributed synchronization.

## .....and others

22

# Apache Hadoop - Hadoop ecosystem

## Apache Mahout:

- It is a popular data mining library. It includes the most popular data mining scalable machine learning algorithms. Also, it is a scalable machine-learning library.
- The following are some companies that are using Mahout: Amazon, Twitter, Yahoo, and LucidWorks Big Data (This is an analytics firm, which uses Mahout for clustering, duplicate document detection, phase extraction, and classification).

23

## Apache Mahout:

- It is a popular data mining library. It includes the most popular data mining scalable machine learning algorithms. Also, it is a scalable machine-learning library.
- The following are some companies that are using Mahout: Amazon, Twitter, Yahoo, and LucidWorks Big Data (This is an analytics firm, which uses Mahout for clustering, duplicate document detection, phase extraction, and classification).

24

## Hadoop YARN

- The JobTracker in v1.0 is the single master that allocates resources for applications, performs scheduling for demand and also monitors the jobs of processing in the system. (master for resources + processing)

### Hadoop YARN

- YARN stands for “**Y**et **A**nother **R**esource **N**egotiator”. It was introduced in Hadoop version 2.0.

- The fundamental idea of YARN is to split up the functionalities of resource management and job scheduling/monitoring.

25

26

27

28

## Hadoop YARN - Architecture

- The main components of YARN architecture include:

1. **Resource Manager:** It is the master of YARN and is responsible for **resource assignment and management** among all the applications.
  - Whenever it receives a processing request, it forwards it to the corresponding node manager and allocates resources for the completion of the request accordingly.
- The ResourceManager has two main components: *Scheduler* and *ApplicationsManager*.

- The main components of YARN architecture include:
  1. **Resource Manager:**
    - ApplicationsManager is responsible for accepting job-submissions, negotiating the first container for executing the application specific ApplicationMaster and provides the service for restarting the ApplicationMaster container on failure.

2. **Container:** In the Container, there are the physical resources like a disk, CPU cores, RAM.
3. **Application Master:** An application is a single job submitted to a framework. The application master is responsible for **negotiating resources** with the resource manager, **executing, tracking the status**, and **monitoring progress** of a single application.
4. **Node manager:** It **sends each node's health status** to the Resource Manager, stating if the node process has finished working with the resource. Node manager is also responsible for **monitoring resource usage** by individual Container and reporting it to the Resource manager.

## Hadoop YARN

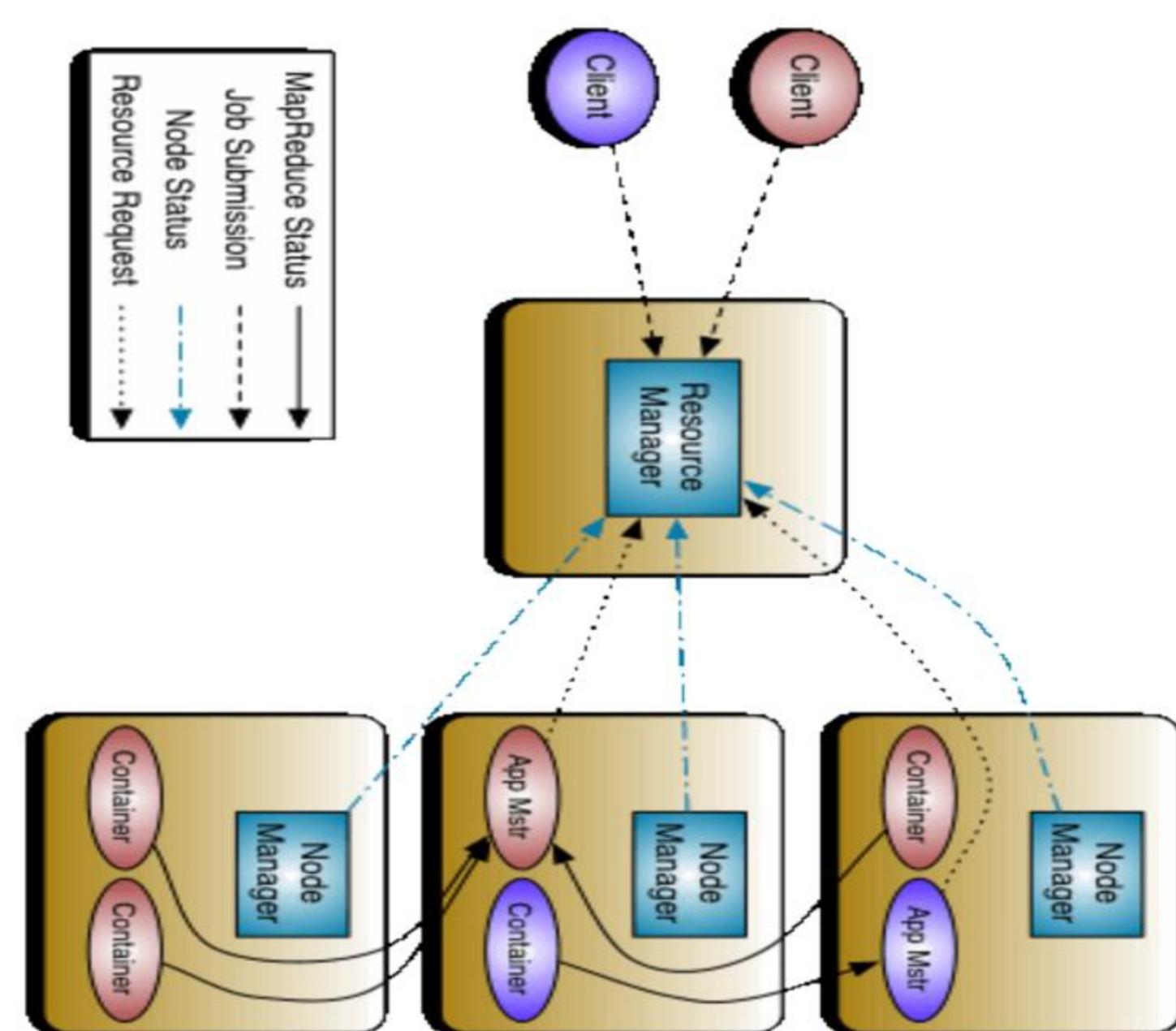
- Thus, YARN is now responsible for Resource Management and Job scheduling.

- Through its various components, YARN can dynamically allocate various resources and schedule the application processing.

- The idea is to have a global ResourceManager (RM) and per-application ApplicationMaster (AM). An application is either a single job or a DAG of jobs.

29

## Hadoop YARN - Architecture



31

## Hadoop YARN - Features

- **Multi-tenancy:** YARN has allowed access to multiple data processing engines such batch, interactive, and real-time stream processing.
- **Scalability:** The scheduler in Resource manager of YARN architecture allows Hadoop to extend and manage thousands of nodes and clusters.
- **Cluster utilization:** YARN allocates all cluster resources in an efficient and dynamic manner, which leads to better utilization.
- **Compatibility:** YARN supports the existing map-reduce applications without disruptions thus making it compatible with Hadoop 1.0 as well.

## Apache Spark

- Apache Spark is a data processing framework that can quickly perform processing tasks on very large data sets, and can also distribute data processing tasks across multiple computers.

- Apache Spark is a prime example of how YARN enables customers to build a real-time analytics platform.

- Spark also takes some of the programming burdens of these tasks off the shoulders of developers with an easy-to-use API that abstracts away much of the work of distributed computing and big data processing.

32

## Hadoop Version 2.0 vs Version 1.0

| Criteria                  | Version 2.0                                                                               | Version 1.0                                                                |
|---------------------------|-------------------------------------------------------------------------------------------|----------------------------------------------------------------------------|
| Components                | <ul style="list-style-type: none"><li>• HDFS</li><li>• MapReduce</li><li>• YARN</li></ul> | <ul style="list-style-type: none"><li>• HDFS</li><li>• MapReduce</li></ul> |
| Managing cluster resource | Excellent due to central resource management                                              | Average due to fixed Map and Reduce slots                                  |

33

## Apache Spark - Features

- **In-memory computing** – The main feature of Spark is its in-memory cluster computing that increases the processing speed of an application.

- **Speed** – Spark helps to run an application in Hadoop cluster faster when running on disk and even faster when run in memory. This is possible by reducing number of read/write operations to disk. It stores the intermediate processing data in memory.

- **Supports multiple languages** – Spark provides built-in APIs in Java, R, Python and Scala.
- **Advanced Analytics** – Spark not only supports Map and reduce, but also supports SQL queries, Streaming data, Machine learning (ML), and Graph algorithms.

34

## Apache Spark - Data structure

➤ **Resilient Distributed Datasets (RDD)** is a fundamental data structure of Spark. It is a distributed collection of objects.

➤ RDDs are the building blocks of any Spark application. RDDs Stands for:

- **Resilient:** Fault tolerant and is capable of rebuilding data on failure
  - **Distributed:** Distributed data among the multiple nodes in a cluster
  - **Dataset:** Collection of partitioned data with values
- Each dataset is divided into logical partitions, which may be computed on different nodes of the cluster.

## Apache Spark - RDD

➤ RDD is a read-only, partitioned collection of records. RDD is a fault-tolerant collection of elements that can be operated on in parallel.

➤ There are two ways to create RDDs:

- **Parallelizing** an existing collection in your driver program.
- **Referencing a dataset** in an external storage system, such as a shared file system, HDFS, HBase, or other data sources.

## Apache Spark - RDD

➤ Spark makes use of the concept of RDD to achieve faster and efficient MapReduce operations.

➤ With RDDs, you can perform two types of operations:

- **Transformations:** They are the operations that are applied to create a new RDD. They are *lazy*, their result RDD is not immediately computed.
- **Actions:** They are applied on an RDD to instruct Apache Spark to apply computation and pass the result back to the driver. They are *eager*, their result is immediately computed.

## Apache Spark - RDD

➤ Transformations examples:

- **map(func):** Return a new distributed dataset formed by passing each element of the source through a function func.
- **filter(func):** Return a new dataset formed by selecting those elements of the source on which func returns true.
- **intersection(otherDataset):** Return a new RDD that contains the intersection of elements in the source dataset and the argument.
- **reduceByKey(func, [numPartitions]):** When called on a dataset of (K, V) pairs, returns a dataset of (K, V) pairs where the values for each key are aggregated using the given reduce function func.

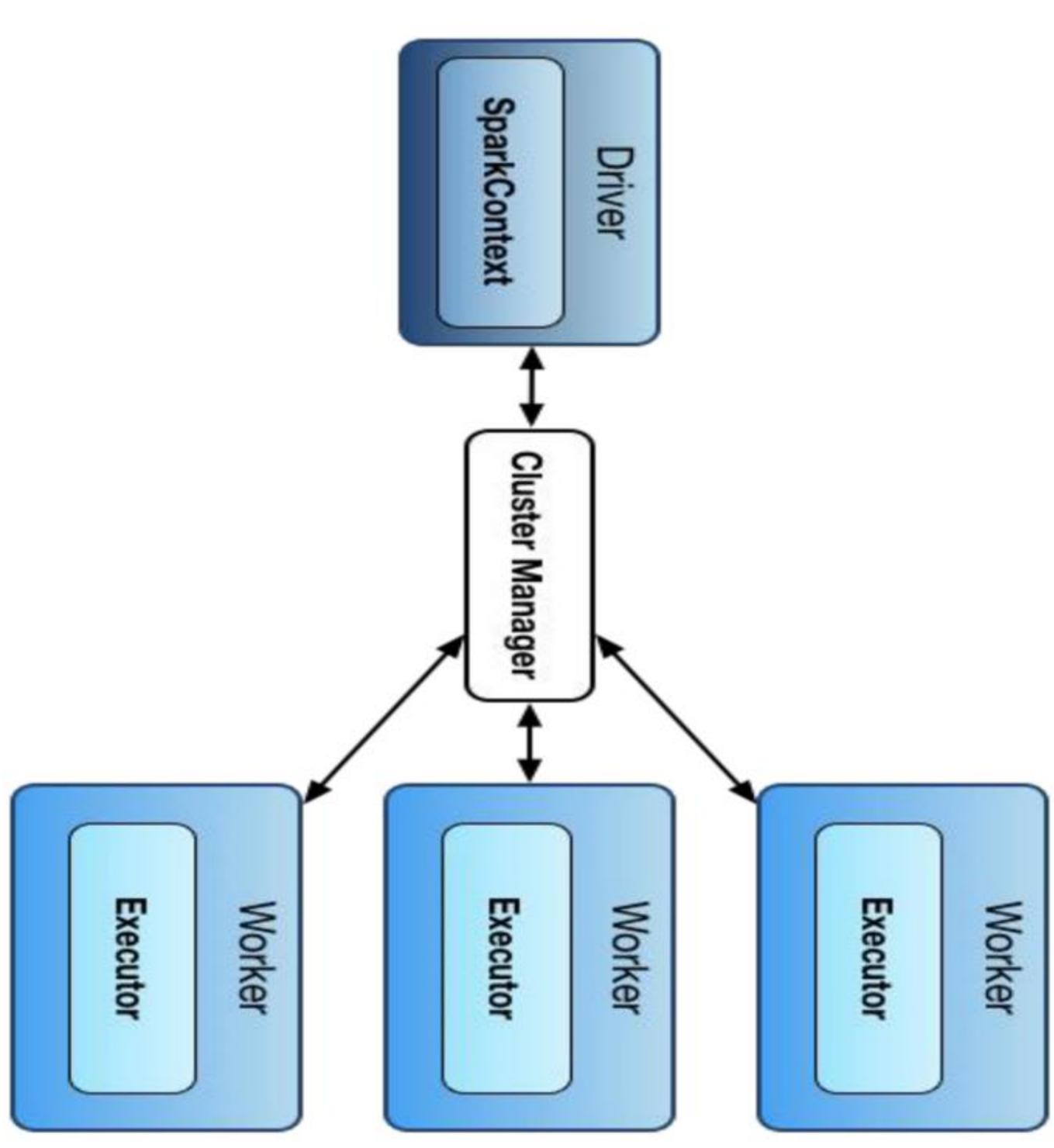
## Apache Spark - Architecture

➤ At a fundamental level, an Apache Spark follows a master/slave architecture with two main components: a *driver* and *workers*.

• **Driver (Master Node):** It is the entry point that runs the main () function of the application and is responsible for the translation of spark user code into actual spark jobs executed on the cluster.

- **Workers (Slave Nodes):** They contain the *executors* that run tasks assigned to them. These tasks are executed on the partitioned RDDs in the worker node. Executors perform all the data processing.
- **Cluster Manager:** This is an external service responsible for acquiring resources on the spark cluster and allocating them to a spark job.

➤ There is also the **SparkContext** which is a gateway to all the Spark functionalities. It is similar to your database connection.



## Apache Spark - Architecture

37

38

39

40

41

42

## Apache Spark - RDD

### Actions examples:

- `reduce(func)`: Aggregate the elements of the dataset using a function func.
- `count()`: Return the number of elements in the dataset.
- `first()`: Return the first element of the dataset.
- `take(n)`: Return an array with the first n elements of the dataset.
- `takeOrdered([n, [ordering]])`: Return the first n elements of the RDD using either their natural order or a custom comparator.

43

## Apache Spark - Workflow

- **STEP 1:** The client submits spark user application code. When an application code is submitted, the driver implicitly converts user code into a logically *directed acyclic graph* called **DAG**. At this stage, it also performs optimizations.

- **STEP 2:** After that, it converts the logical graph called DAG into physical execution plan with many stages. After converting into a physical execution plan, it creates physical execution units called tasks under each stage. Then the tasks are bundled and sent to the cluster.

44

## Apache Spark - Workflow

- **STEP 3:** Now the driver talks to the cluster manager and negotiates the resources. Cluster manager launches executors in worker nodes on behalf of the driver. At this point, the driver will send the tasks to the executors based on data placement. The driver will have a complete view of executors that are executing the task.

- **STEP 4:** During the course of execution of tasks, driver will monitor the set of executors that runs. Driver node also schedules future tasks based on data placement.

45

## Apache Spark - ecosystem

➤ The following illustration depicts the different components of Spark:



## Apache Spark - ecosystem

➤ Apache Spark Core:

- Spark Core is the underlying general execution engine for spark platform that all other functionality is built upon. It provides In-Memory computing and referencing datasets in external storage systems.

- Spark SQL:
  - Spark SQL enables users to run SQL queries. Alongside standard SQL support, Spark SQL provides a standard interface for reading from and writing to other datastores including JSON, HDFS, Apache Hive and others.

## Apache Spark - ecosystem

➤ Spark Streaming:

- Spark Streaming extended the Apache Spark concept of batch processing into streaming by breaking the stream down into a continuous series of microbatches, which could then be manipulated using the Apache Spark API.

### ➤ MLlib (Machine Learning Library)

- Machine learning library delivers both efficiencies as well as the high-quality algorithms. It is capable of in-memory data processing that improves the performance of iterative algorithm drastically.

### ➤ GraphX

- Spark GraphX is the graph computation engine built on top of Apache Spark that enables to process graph data at scale.

46

## Hadoop VS Spark

|                                | Hadoop                                                                                                                     | Spark                                                                                                                                                       |
|--------------------------------|----------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Processing Speed & Performance | Hadoop reads and writes from a disk, thus slowing down the processing speed.                                               | Spark reduces the number of read/write cycles to disk and stores intermediate data in memory, hence <b>faster</b> -processing speed.                        |
| Usage                          | Hadoop is designed to handle <b>batch</b> processing efficiently.                                                          | Spark is designed to handle <b>real-time</b> data efficiently.                                                                                              |
| Fault Tolerance                | Fault-tolerance achieved by <b>replicating blocks of data</b> . If a node goes down, the data can be found on another node | Fault-tolerance achieved by <b>storing chain of transformations</b> . If data is lost, the chain of transformations can be recomputed on the original data. |

50

## Hadoop VS Spark

|           | Hadoop                                                                                                                                                                                                                                                                                                                              | Spark                                                                                                                                                                                                                                                                               |
|-----------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Use Cases | <ul style="list-style-type: none"><li>• Processing big data sets in environments where <b>data size exceeds available memory</b></li><li>• Building data analysis infrastructure with a <b>limited budget</b></li><li>• Completing jobs that are <b>not time-sensitive</b></li><li>• Historical and archive data analysis</li></ul> | <ul style="list-style-type: none"><li>• Dealing with chains of parallel operations by using <b>iterative algorithms</b></li><li>• Analyzing <b>stream data</b> in real time</li><li>• <b>Graph-parallel processing</b> to model data</li><li>• All <b>ML</b> applications</li></ul> |

51

Thank You

52