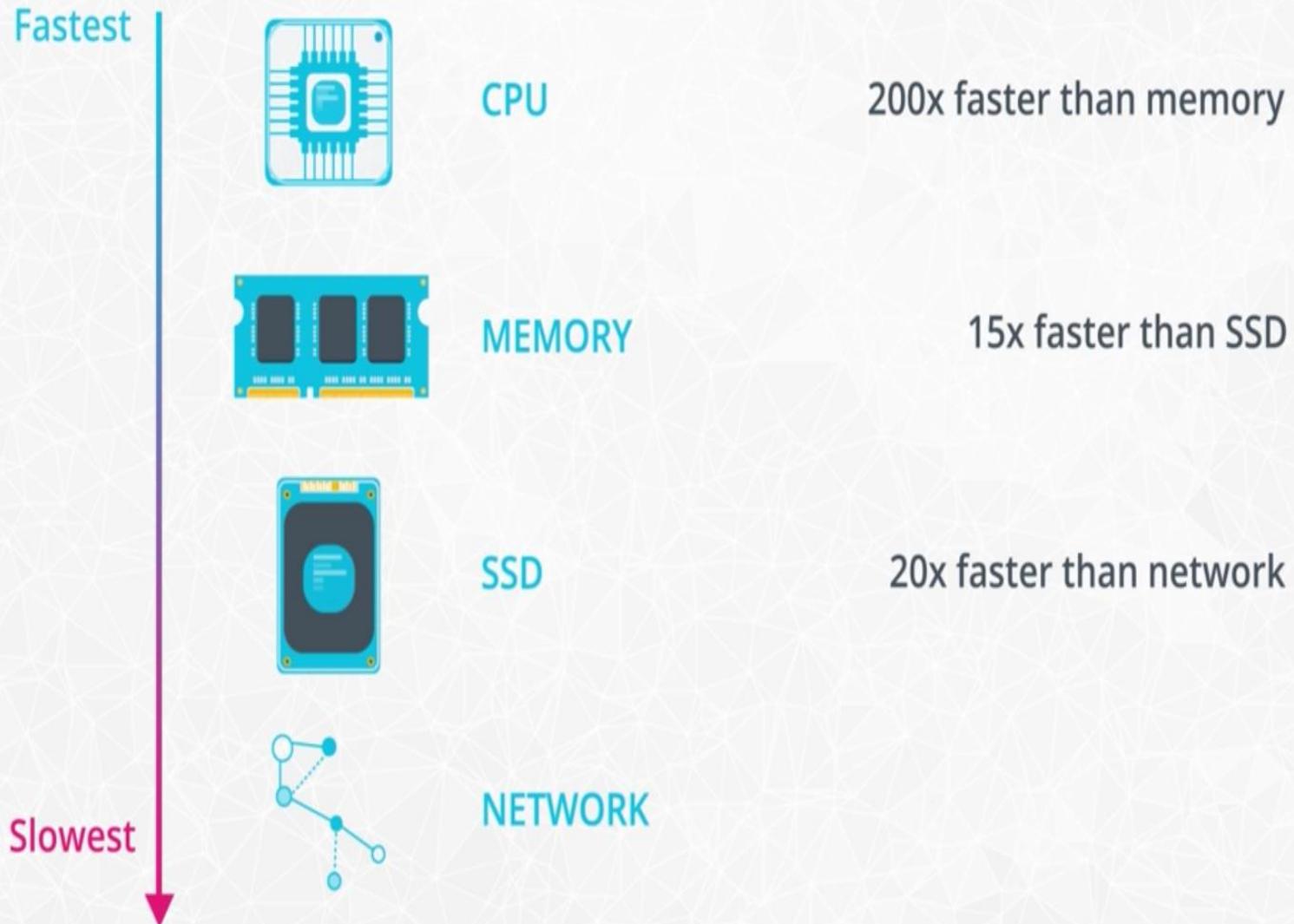


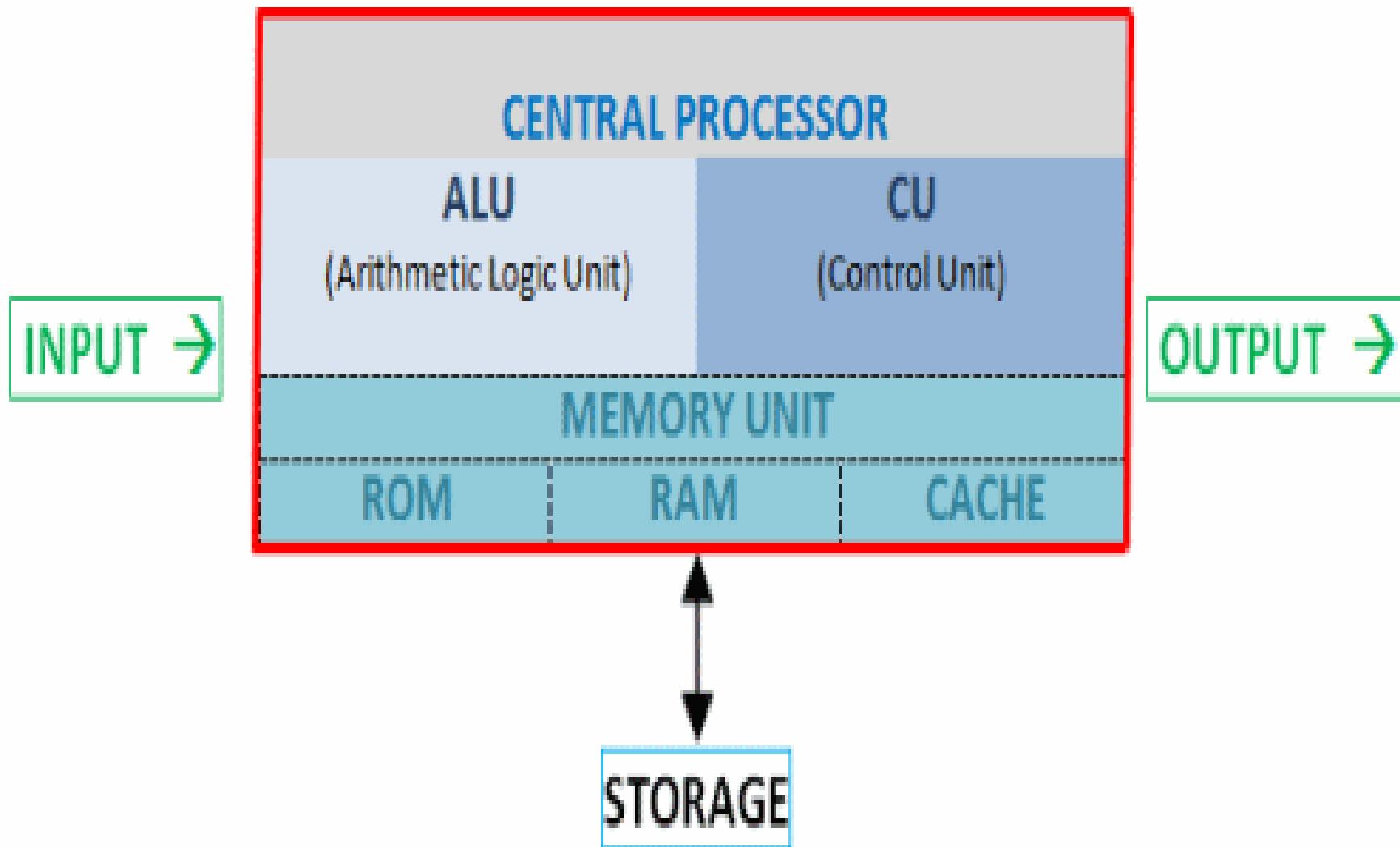
Prepared By : Ravi

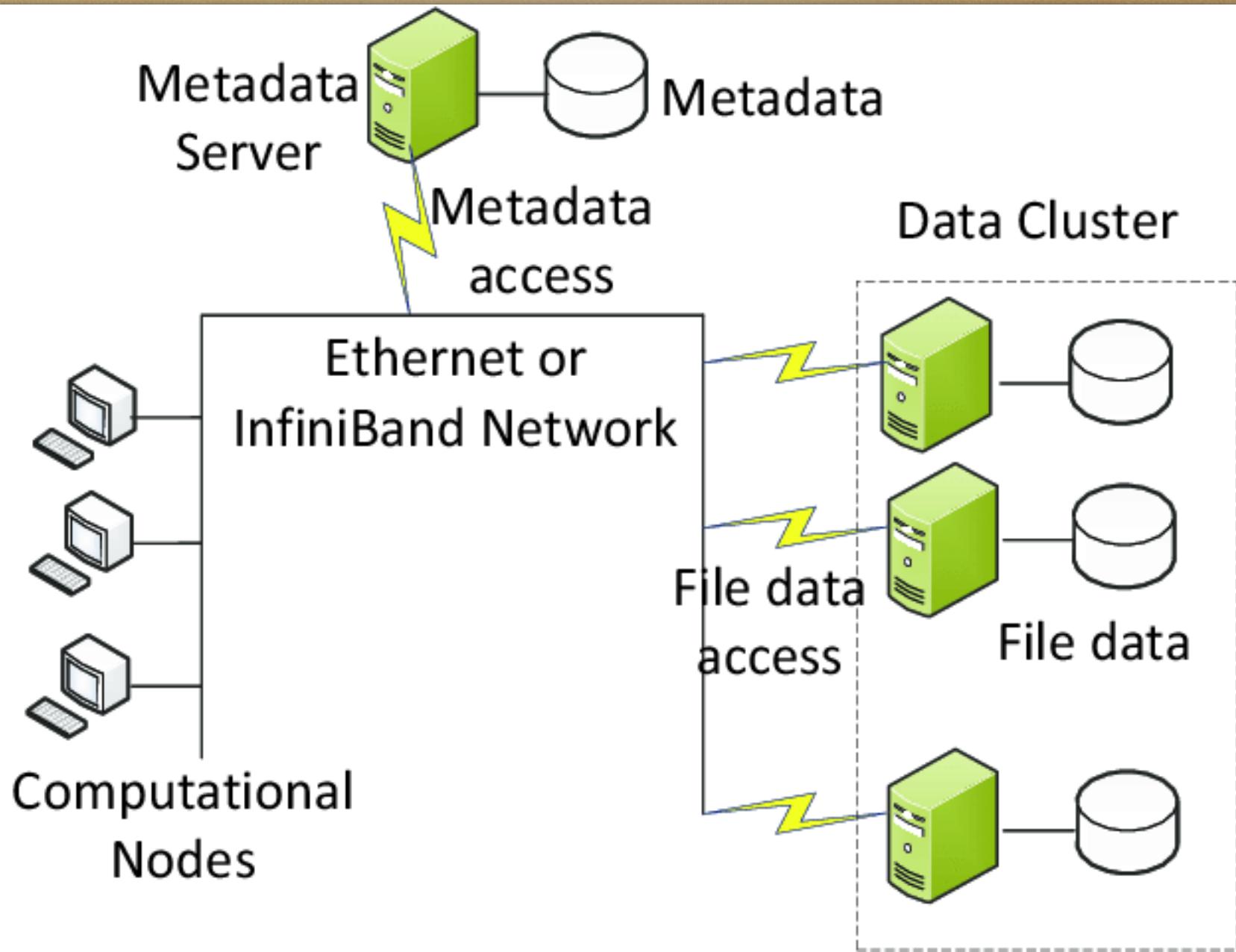


Source: [spark.apache.org](http://spark.apache.org) & [www.databricks.com](http://www.databricks.com)



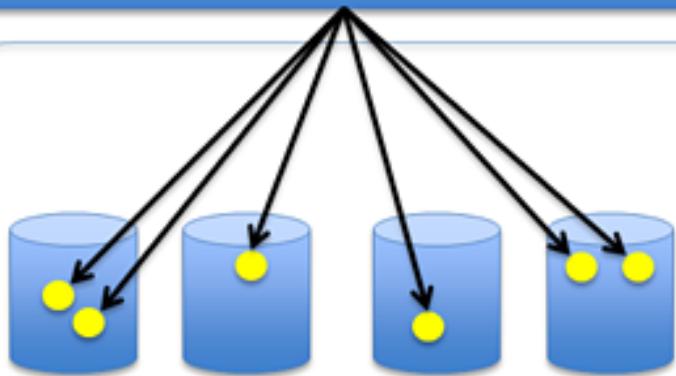
# Central Processing Unit



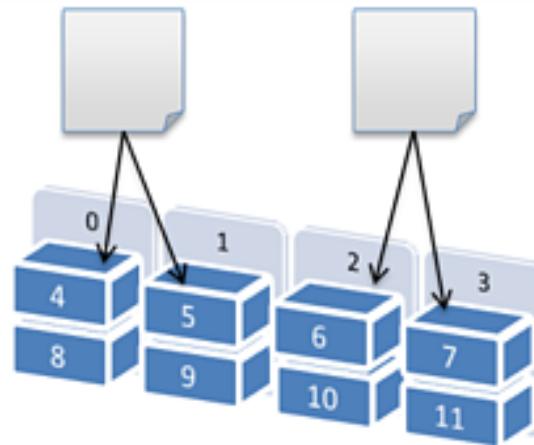


## HTTP(S) Interface

### Object Storage



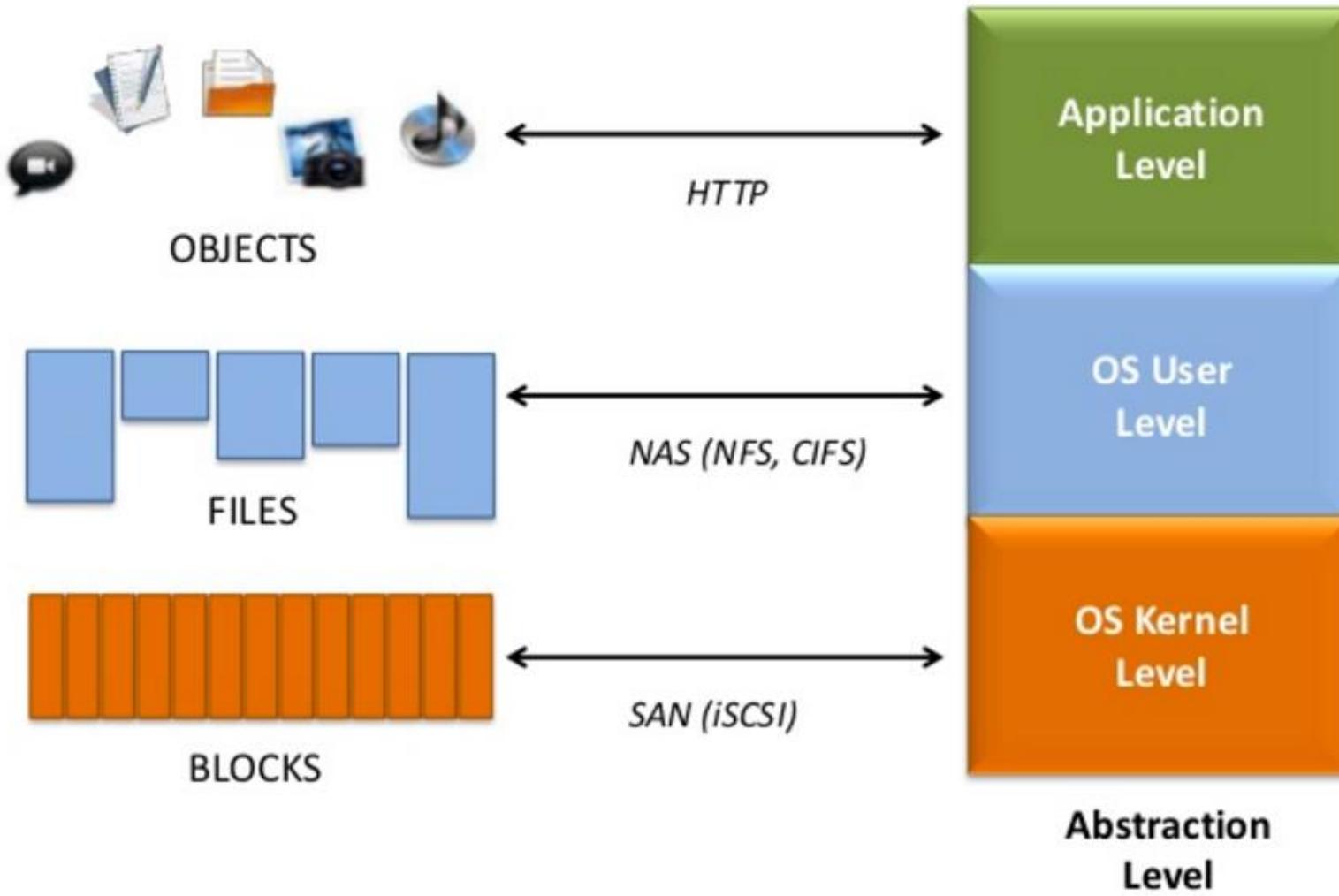
## Block Storage



- Store virtually unlimited files.
- Maintain file revisions.
- HTTP(S) based interface.
- Files are distributed in different physical nodes.

- File is split and stored in fixed sized blocks.
- Capacity can be increased by adding more nodes.
- Suitable for applications which require high IOPS, database, transactional data.

# Object vs. File vs. Block Storage





# What is Big Data?



# THE 3Vs OF BIG DATA

## VOLUME

- Amount of data generated
- Online & offline transactions
- In kilobytes or terabytes
- Saved in records, tables, files



## VELOCITY

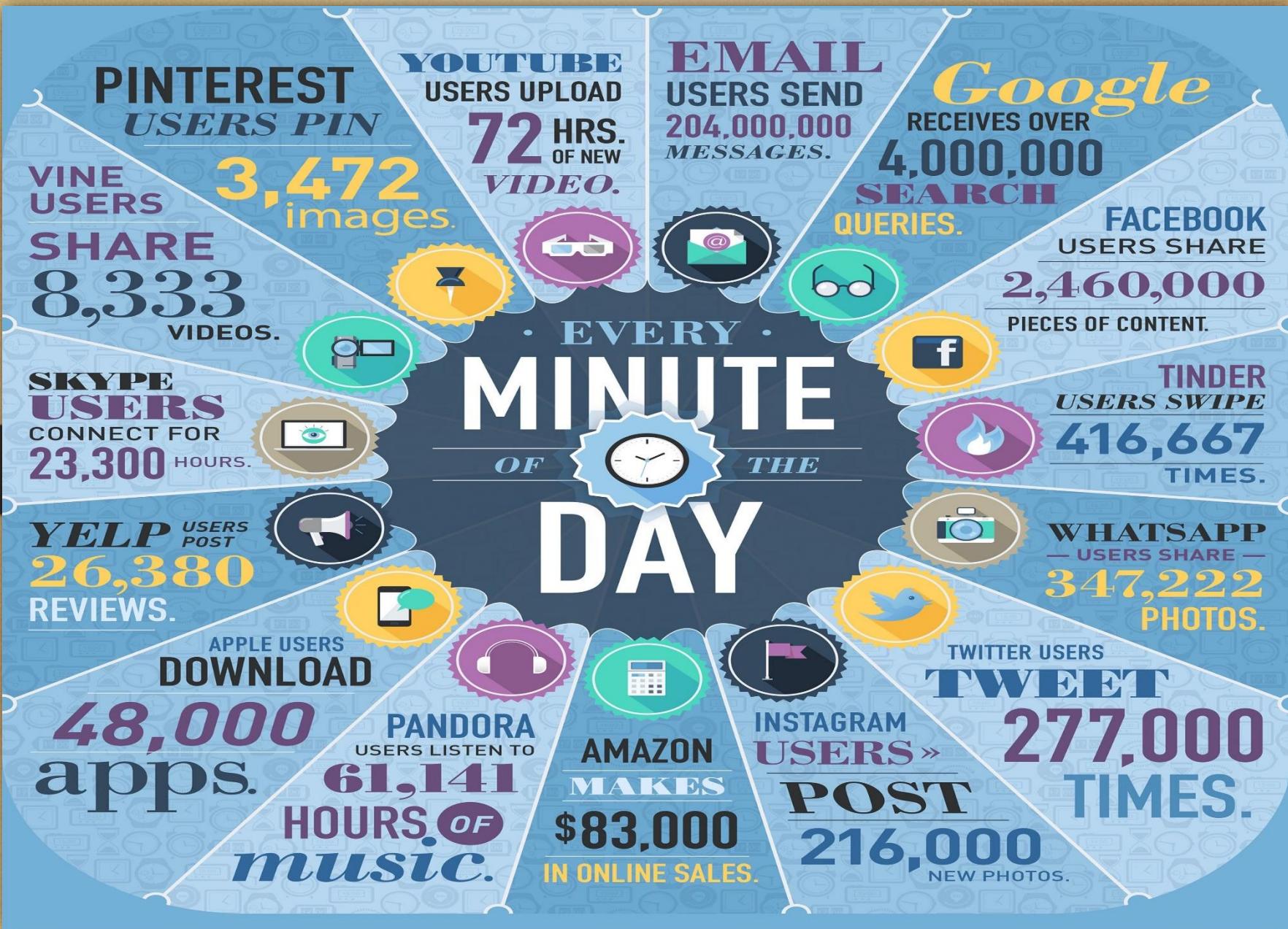
- Speed of generating data
- Generated in real-time
- Online and offline data
- In Streams, batch or bits



## VARIETY

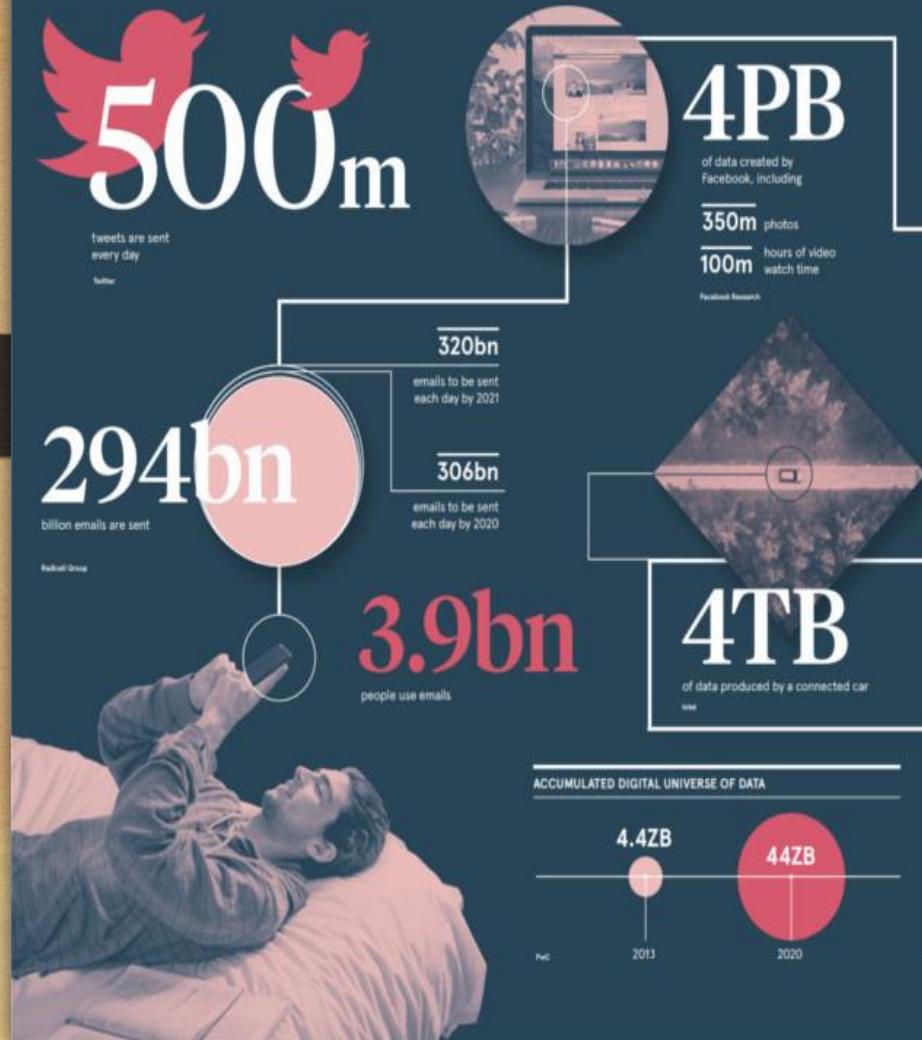
- Structured & unstructured
- Online images & videos
- Human generated - texts
- Machine generated - readings



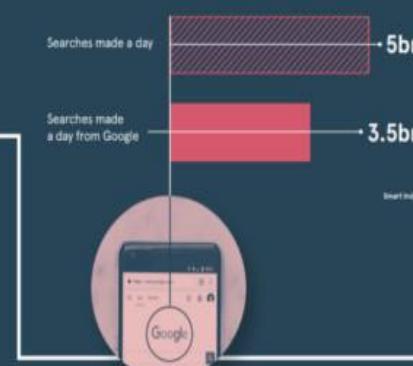


## A DAY IN DATA

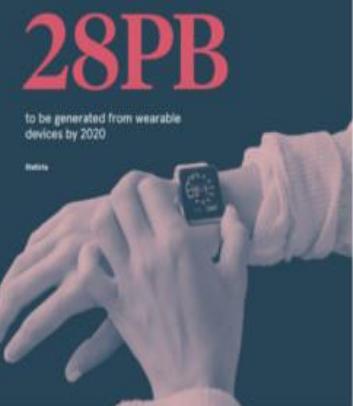
The exponential growth of data is undisputed, but the numbers behind this explosion - fuelled by internet of things and the use of connected devices - are hard to comprehend, particularly when looked at in the context of one day.



Unit	Value	Size
<b>B</b> bit	0 or 1	1/8 of a byte
<b>B</b> byte	8 bits	1 byte
<b>KB</b> kilobyte	1,000 bytes	1,000 bytes
<b>MB</b> megabyte	1,000 <sup>3</sup> bytes	1,000,000 bytes
<b>GB</b> gigabyte	1,000 <sup>6</sup> bytes	1,000,000,000 bytes
<b>TB</b> terabyte	1,000 <sup>12</sup> bytes	1,000,000,000,000 bytes
<b>PB</b> petabyte	1,000 <sup>15</sup> bytes	1,000,000,000,000,000 bytes
<b>EB</b> exabyte	1,000 <sup>18</sup> bytes	1,000,000,000,000,000,000 bytes
<b>ZB</b> zettabyte	1,000 <sup>21</sup> bytes	1,000,000,000,000,000,000,000 bytes
<b>YB</b> yottabyte	1,000 <sup>24</sup> bytes	1,000,000,000,000,000,000,000,000 bytes



of data will be created every day by 2025



# The six Vs of big data

Big data is a collection of data from various sources, often characterized by what's become known as the 3Vs: *volume, variety and velocity*. Over time, other Vs have been added to descriptions of big data:

VOLUME	VARIETY	VELOCITY	VERACITY	VALUE	VARIABILITY
The amount of data from myriad sources.	The types of data: structured, semi-structured, unstructured.	The speed at which big data is generated.	The degree to which big data can be trusted.	The business value of the data collected.	The ways in which the big data can be used and formatted.

**Velocity**

Frequency of data generation



300 hours

of video uploaded to YouTube every minute (estimate)



2,400,000 search queries

per minute on Google

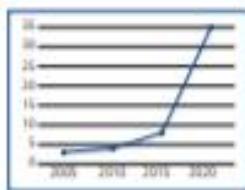


4,170,000 posts liked

on Facebook per minute

**Volume**

The growth of world data



1 terabyte holds the equivalent of roughly 210 single sided DVDs.

**Variety**

Structured and unstructured data - types of Big Data



Web and social media



Machine to machine



Big transaction data



Biometric



Human-generated

**Veracity**

Establishing trust in data



1 in 3

business leaders don't trust the information they use



Uncertainty

due to inconsistency, ambiguity, latency and approximation

**Viability**

Relevance and feasibility

"Can we use mobile phone data to monitor cross-border tourism?"  
Hypothesis  
Validation to determine if the data will have a meaningful impact



Long-term rewards and better outcomes from hidden relationships in data

**Value**

Return on investment

**Costs**

Risk of simply creating Big Costs without creating the value

**Insights**

Sophisticated queries, counterintuitive insights and unique learning

## Operational Efficiency



### Use data to:

- Increase level of transparency
- Optimize resource consumption
- Improve process quality and performance

## Customer Experience



### Exploit data to:

- Increase customer loyalty and retention
- Perform precise customer segmentation and targeting
- Optimize customer interaction and service

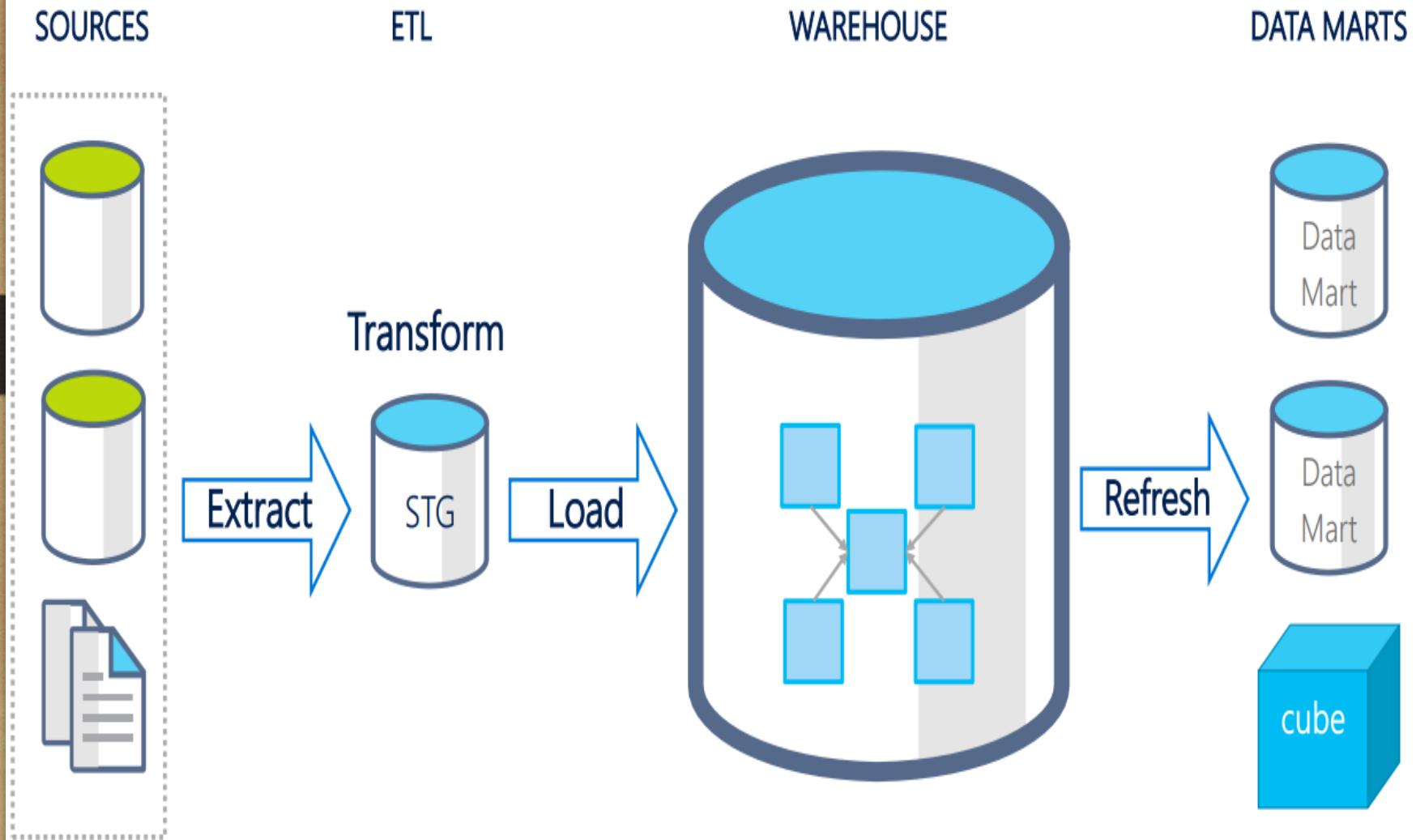
## New Business Models



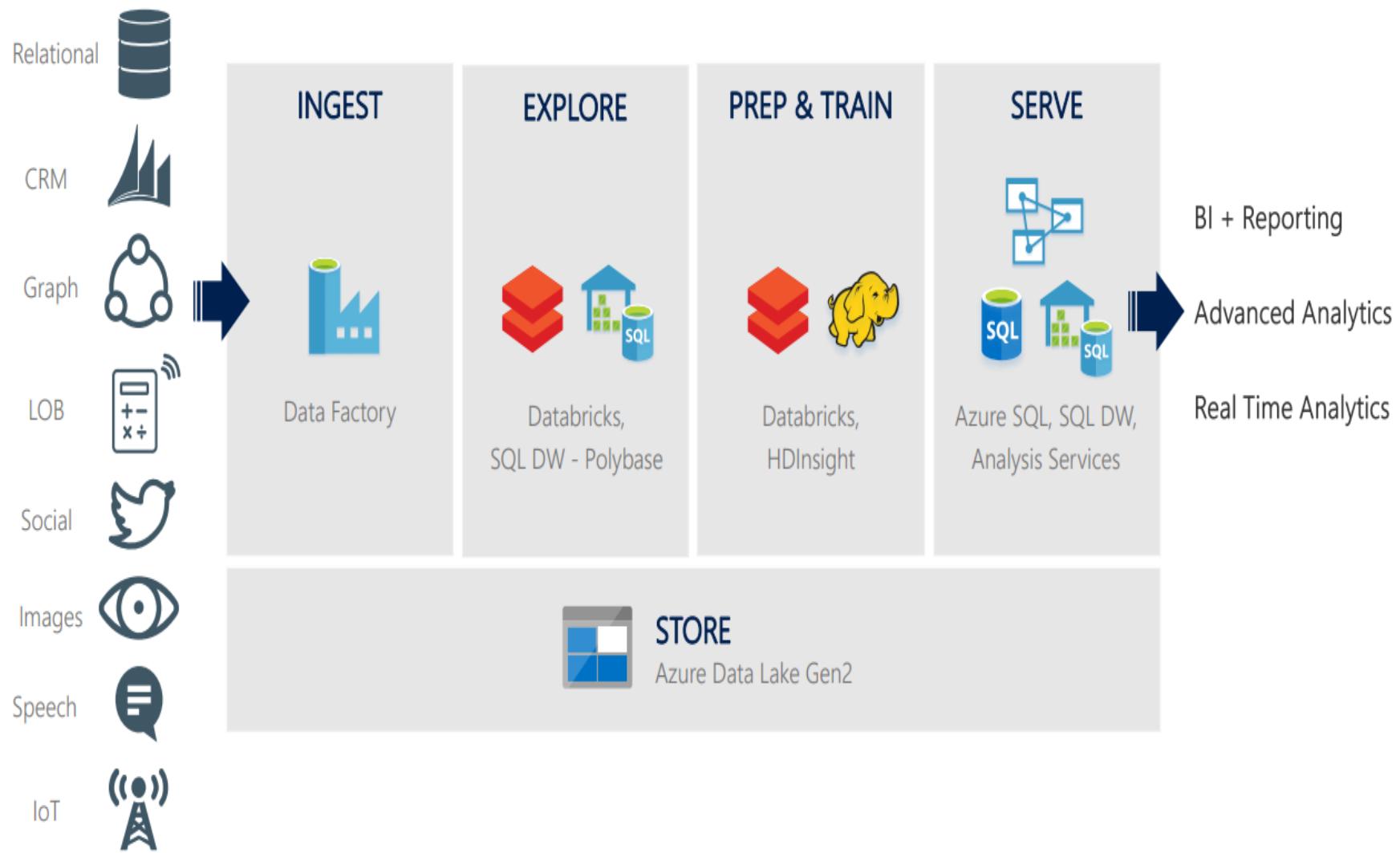
### Capitalize on data by:

- Expanding revenue streams from existing products
- Creating new revenue streams from entirely new (data) products

# Traditional Data Warehousing



# Modern Data Warehouse on Azure



# What is HDFS?



Hadoop  
Client



HDFS

# What is HDFS?



Hadoop  
Client

"I have a 200 TB file  
that I need to  
store."

"Wow - that is big data!  
I will need to distribute  
that across a cluster."



HDFS

# What is HDFS?



Hadoop  
Client

"I have a 200 TB file  
that I need to  
store."

"Wow - that is big data!  
I will need to distribute  
that across a cluster."

"Sounds risky! What  
happens if a drive  
fails?"



HDFS

# What is HDFS?



Hadoop  
Client

"I have a 200 TB file  
that I need to  
store."

"Wow - that is big data!  
I will need to distribute  
that across a cluster."

"Sounds risky! What  
happens if a drive  
fails?"



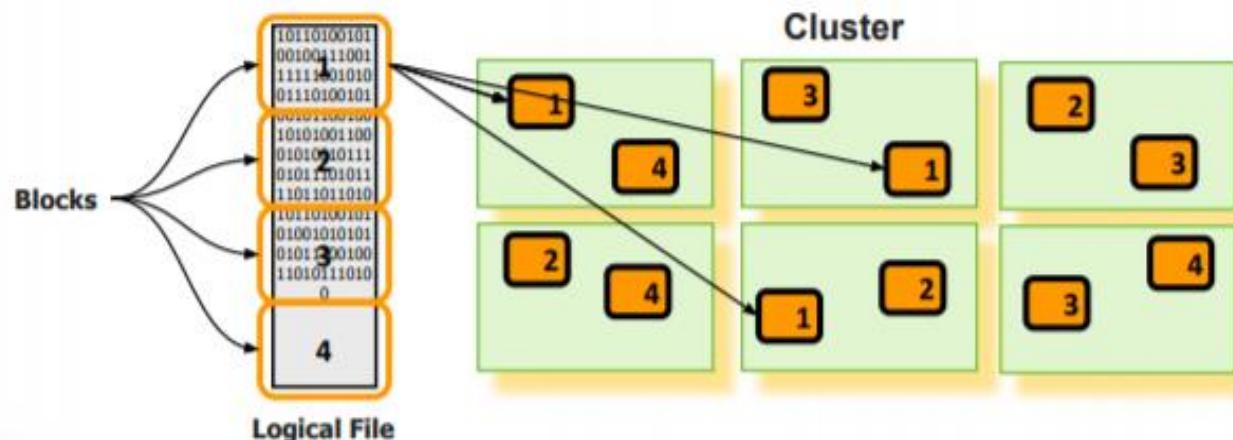
HDFS

"No need to worry! I  
am designed for  
failover."

# HDFS

## Key Ideas

- Write Once, Read Many times (WORM)
- Divide files into big blocks and distribute across the cluster
- Store multiple replicas of each block for reliability
- Programs can ask "where do the pieces of my file live?"



## History of Hadoop

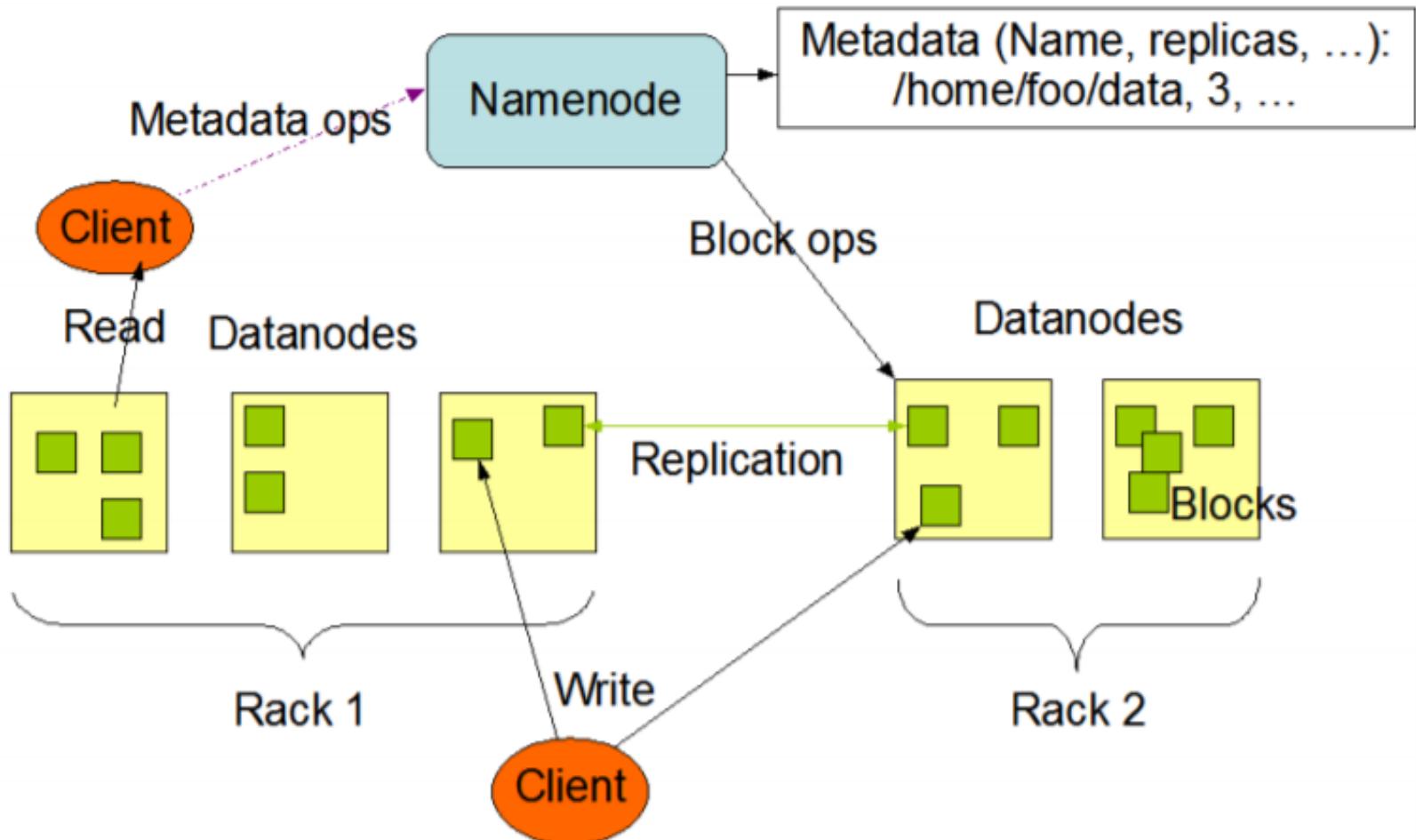
In 1990's Google had to come up with more data and to get the proper solution it has taken 13 years. In 2003, they had introduced GFS (Google File System) which is a technique to store huge data. In 2004, they have introduced MapReduce which is the best processing technique. They have published a "white paper" which has a description of GFS and MapReduce. Later, Yahoo which is the next best search engine after Google introduced HDFS in the year 2006-2007 and MapReduce was introduced in 2007-2008. They have taken the white paper which was given by Google and started implementing and came up with HDFS (Hadoop Distributed File System) and MapReduce. These are the two core components of Hadoop. Hadoop was then introduced by Doug Cutting in 2005.

# Map/Reduce History



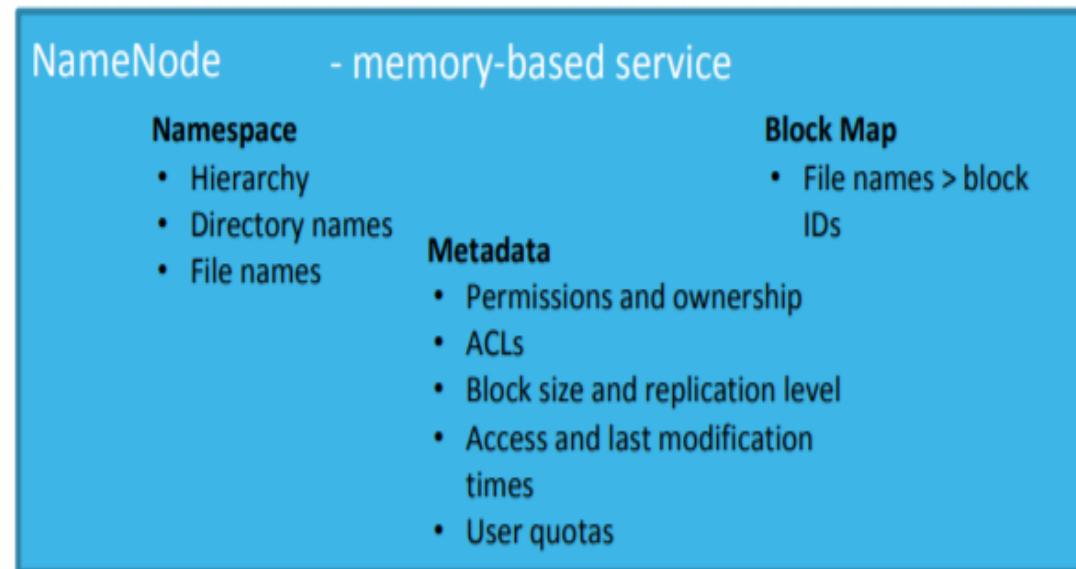
- **Invented by Google:**
  - Inspired by functional programming languages map and reduce functions
  - Seminal paper: Dean, Jeffrey & Ghemawat, Sanjay (OSDI 2004), "MapReduce: Simplified Data Processing on Large Clusters"
  - Used at Google to completely regenerate Google's index of the World Wide Web.
  - It replaced the old ad-hoc programs that updated the index and ran the various analysis
- **Uses:**
  - distributed pattern-based searching, distributed sorting, web link-graph reversal, term-vector per host, web access log stats, inverted index construction, document clustering, machine learning, statistical machine translation
- **Apache Hadoop:** Open source implementation matches Google's specifications
- **Amazon Elastic MapReduce** running on Amazon EC2

## HDFS Architecture

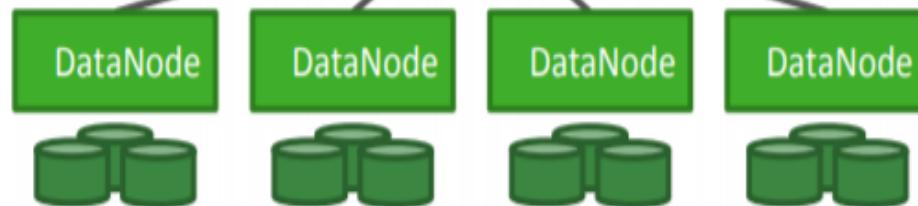


# HDFS Architecture

- The NameNode (master node) and DataNodes (worker nodes) are daemons running in a Java virtual machine.



**Block Storage**  
• Data blocks



## Working of HDFS

- ❖ Suppose a client is willing to put 150MB of data in a cluster and sends a request to the NameNode cluster as metadata. Metadata stores the data about the data given by the Client.
- ❖ 150MB of data is stored in a file with the file name as file.txt as shown in Figure2.
- ❖ The file is divided into 3 input splits a.txt, b.txt, c.txt of each 64MB block size (150MB / 64MB).

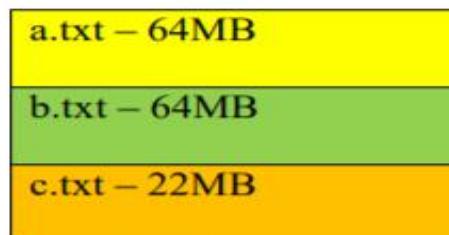
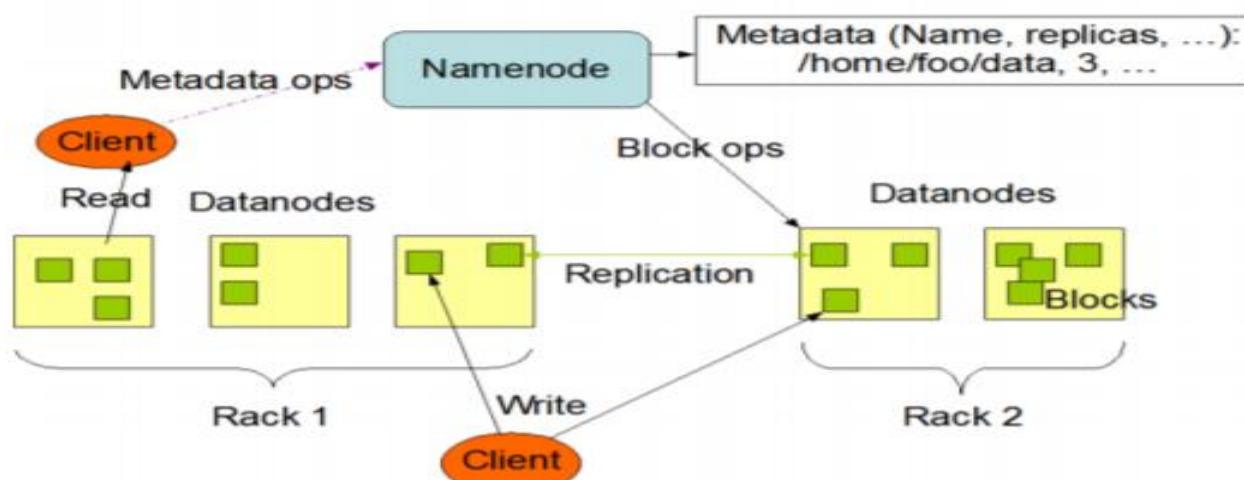


Figure 2. File.txt input splits

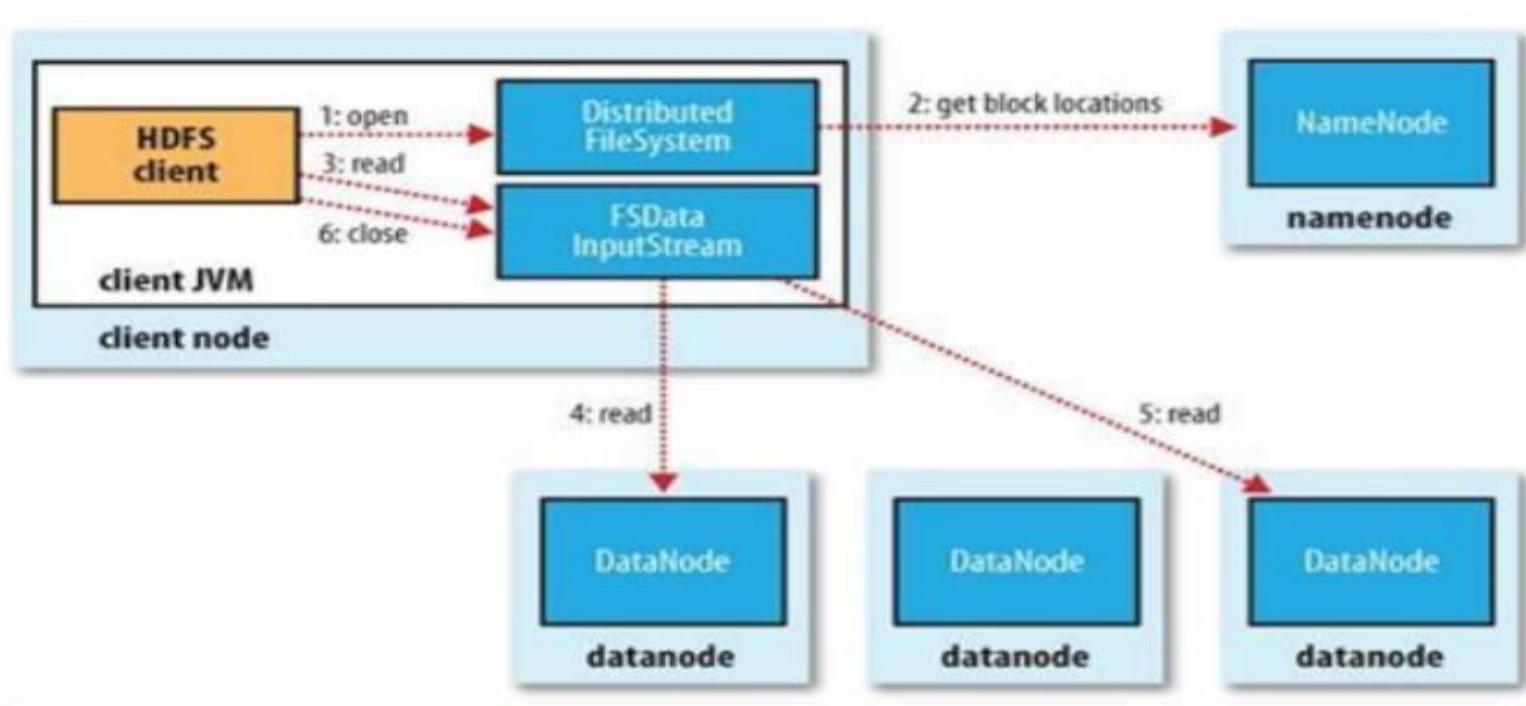
- ❖ NameNode responds to the client and requests to store 150MB data in the nodes which has space.
- ❖ Client store all the txt files in different DataNodes. However, all the files need not be in sequence order.

- ❖ DataNodes are commodity hardware which means if the system goes down the data doesn't lose since HDFS has been given 3 replications by default. Hence it has 2 more backup files for each text files stored in different DataNodes. Hence, the a.txt file occupies 450 MB (150 MB \* 3) of files in the whole cluster because of the replication. The same way other text files are also allocated to DataNodes with their corresponding replications. All the DataNodes which are SlaveNodes for that NameNode give proper block report and heartbeat to the NameNode. This acknowledgment gives the information of the condition of the DataNodes. Block report shows the DataNodes are still allocated with some size of block and heartbeat gives the status of the nodes. This is how the data is stored in HDFS.

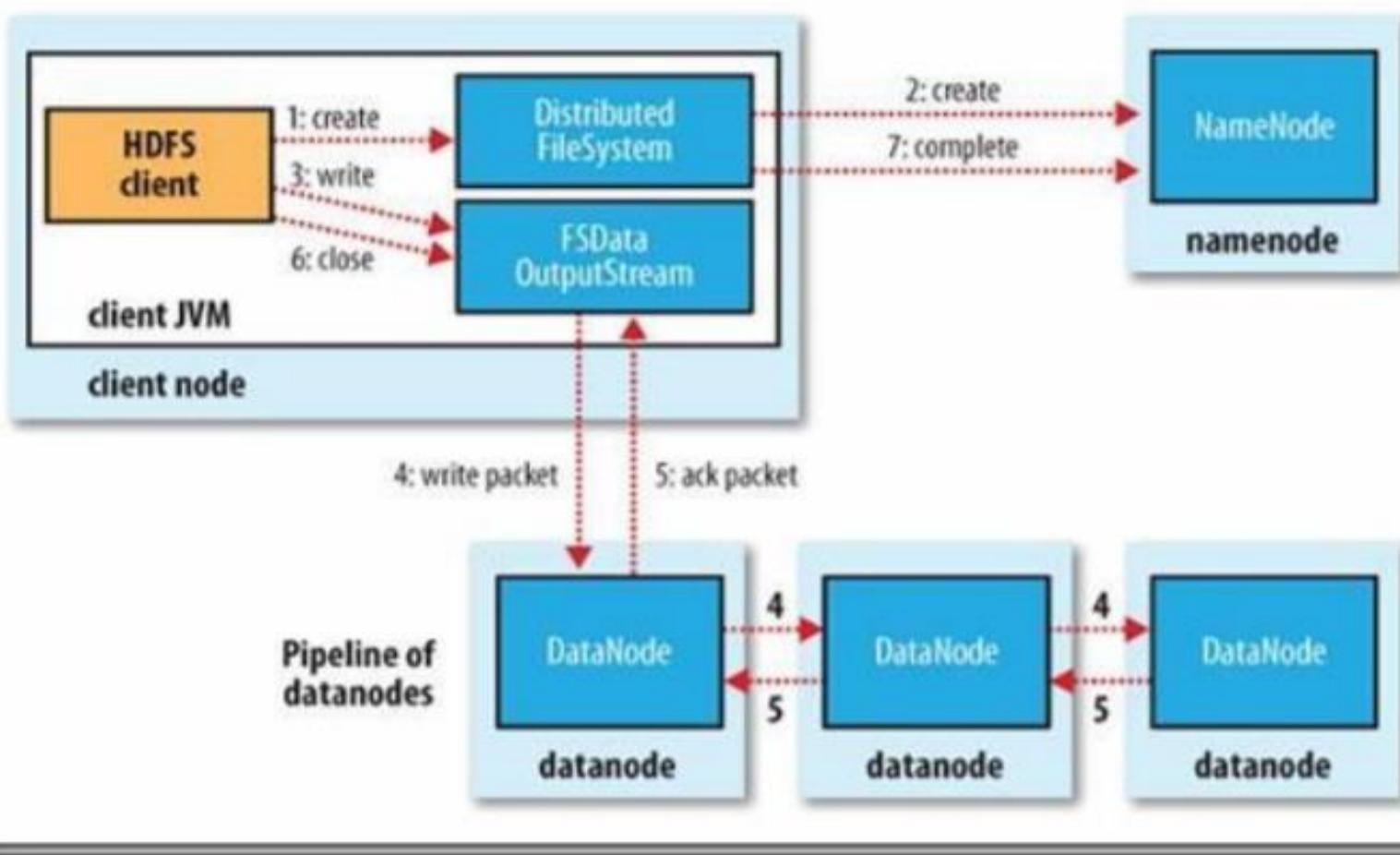


## Data Reading Process in HDFS

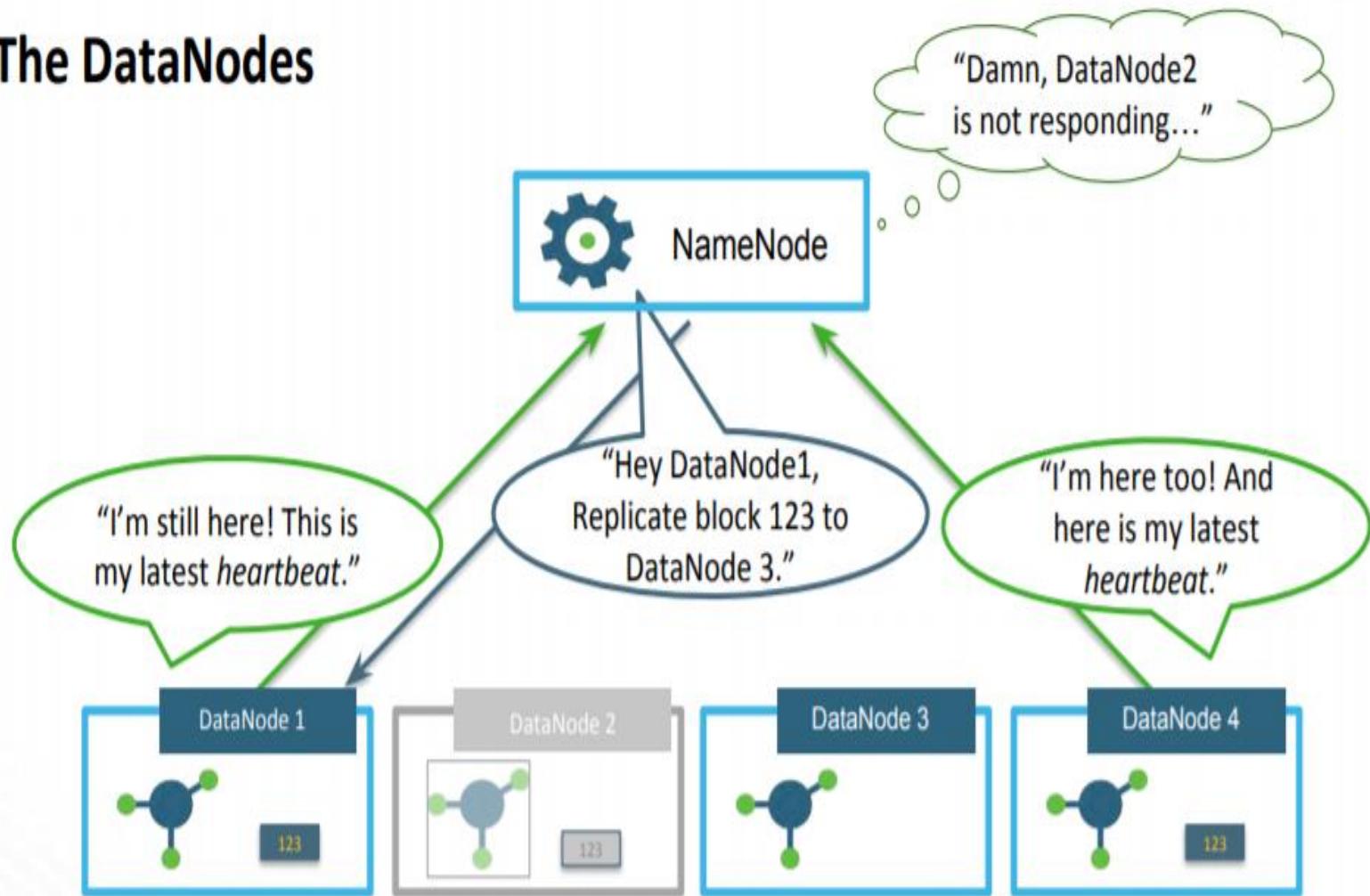
The data reading process in HDFS is not difficult. It is similar to the programming logic which has created the object, i.e., calling the method and performing the execution. The following section will introduce the reading processing of the HDFS.



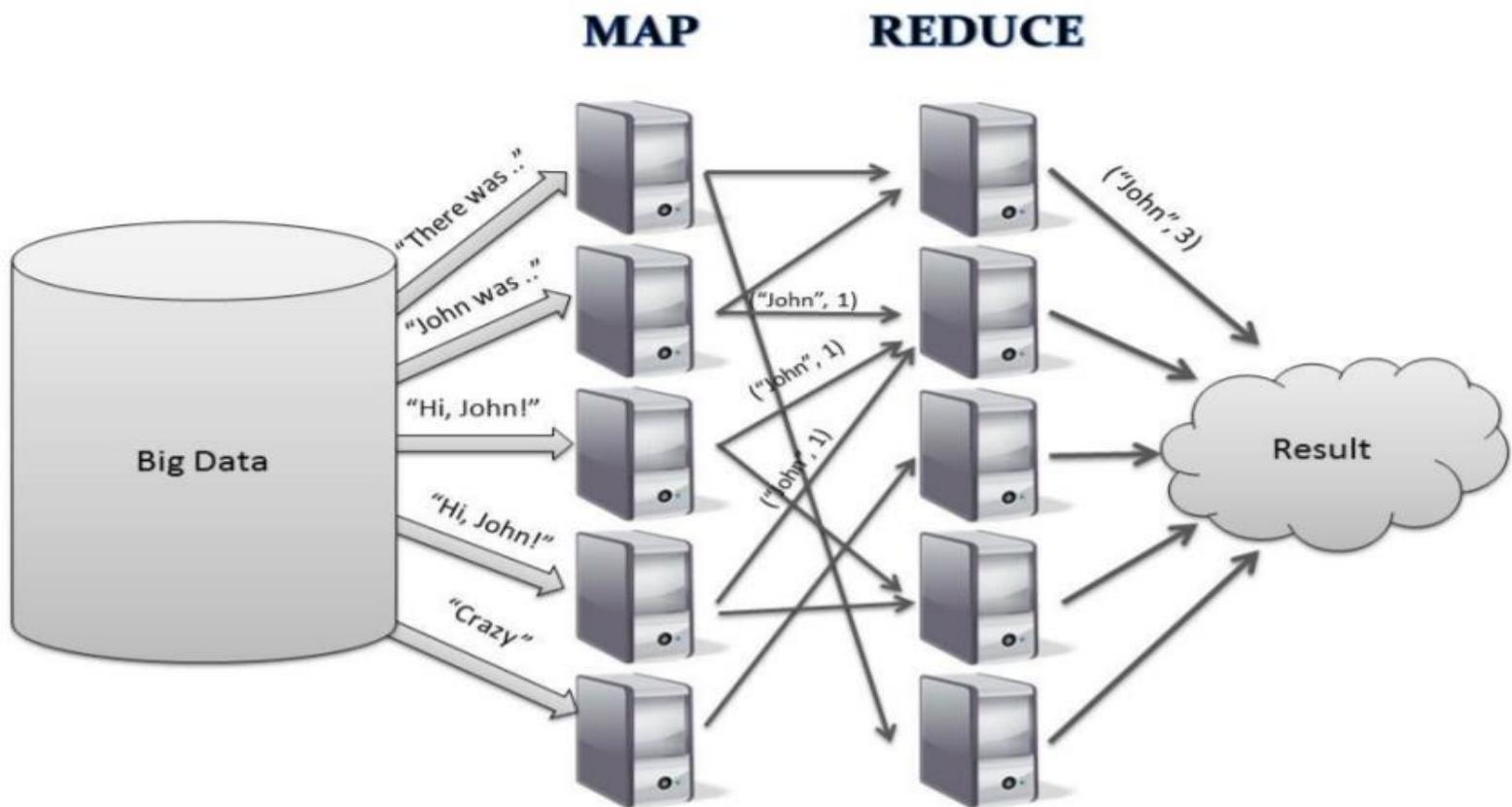
The structure of HDFS reading process is similar to the writing process. There are the following seven steps:

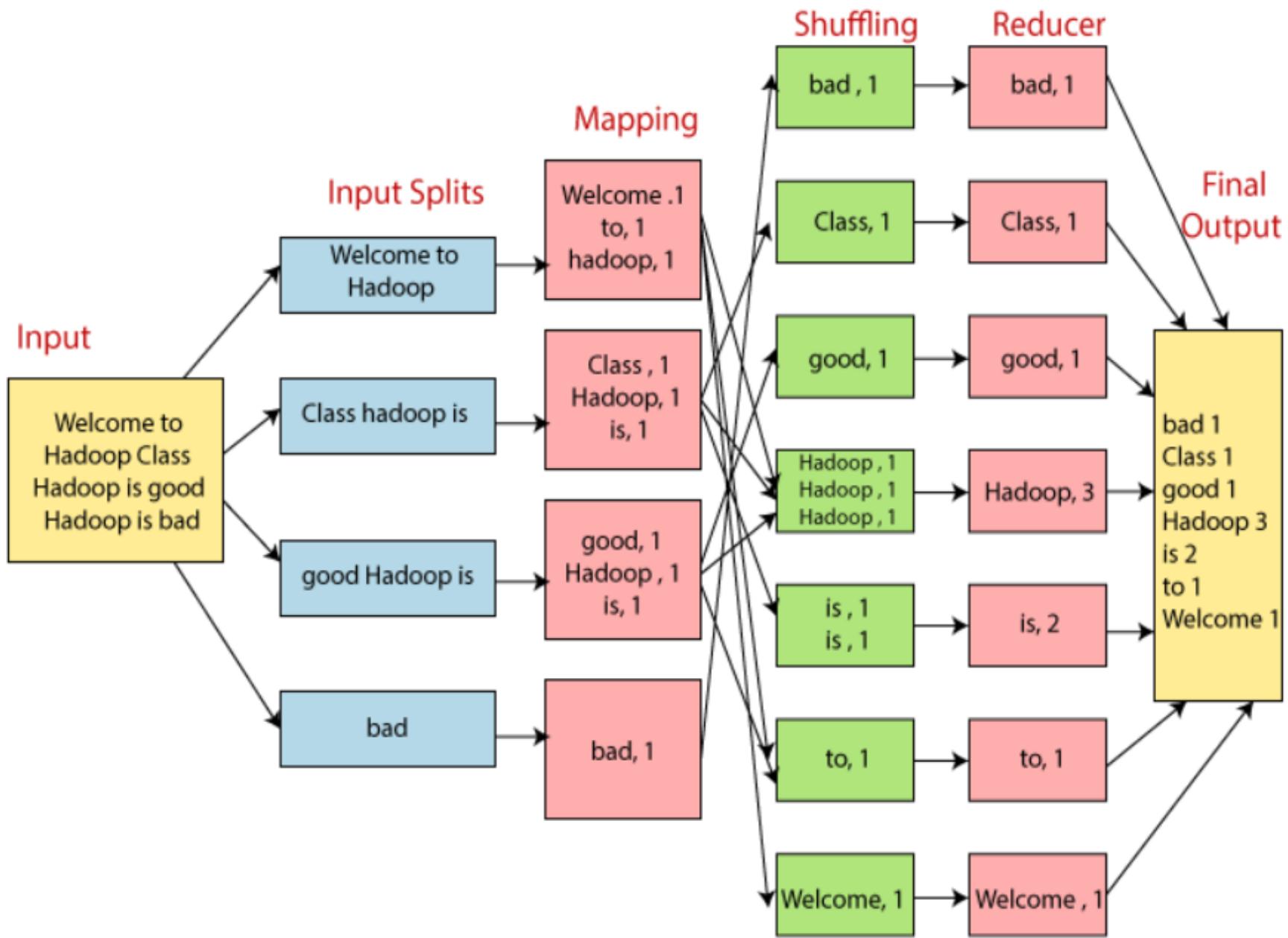


# The DataNodes



# Map/Reduce



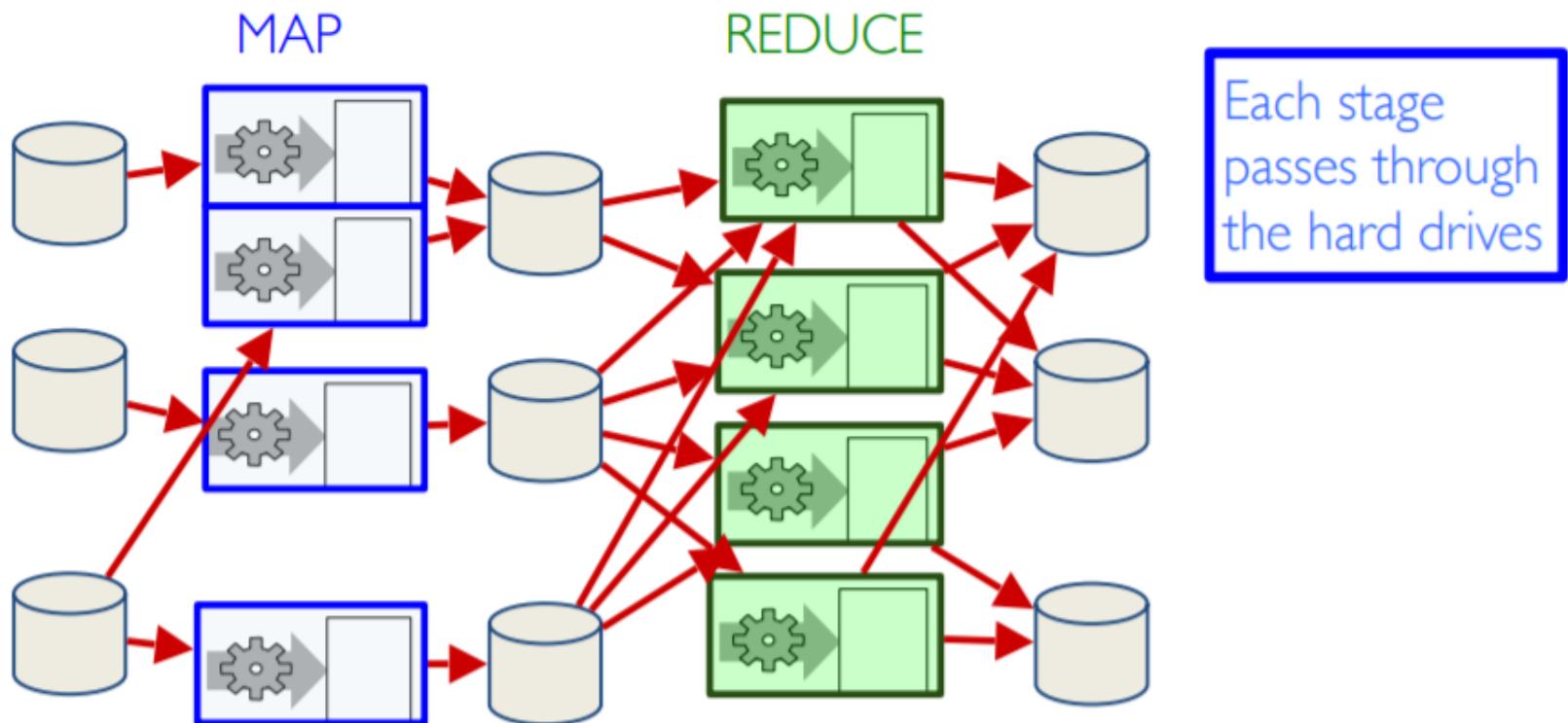


# Map/Reduce

- MapReduce is a framework for processing **parallelizable** problems across huge datasets using a large number of computers (nodes), collectively referred to as a **cluster**.
- **Map step:** Each worker node applies the "map()" function to the local data, and writes the output to a temporary storage. A master node orchestrates that for redundant copies of input data, only one is processed.
- **Shuffle step:** Worker nodes redistribute data based on the output keys (produced by the "map()" function), such that all data belonging to one key is located on the same worker node.
- **Reduce step:** Worker nodes now process each group of output data, per key, in parallel.

# Why Spark?

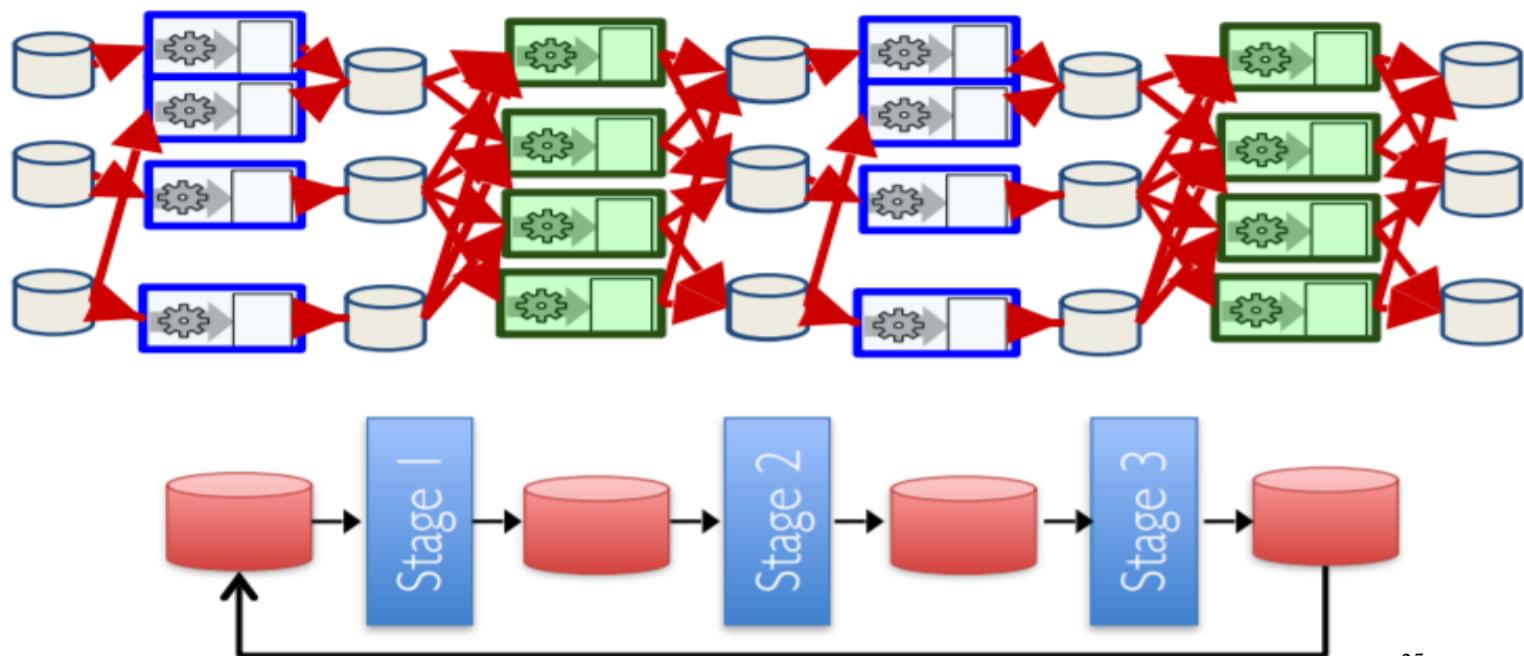
# MapReduce Execution



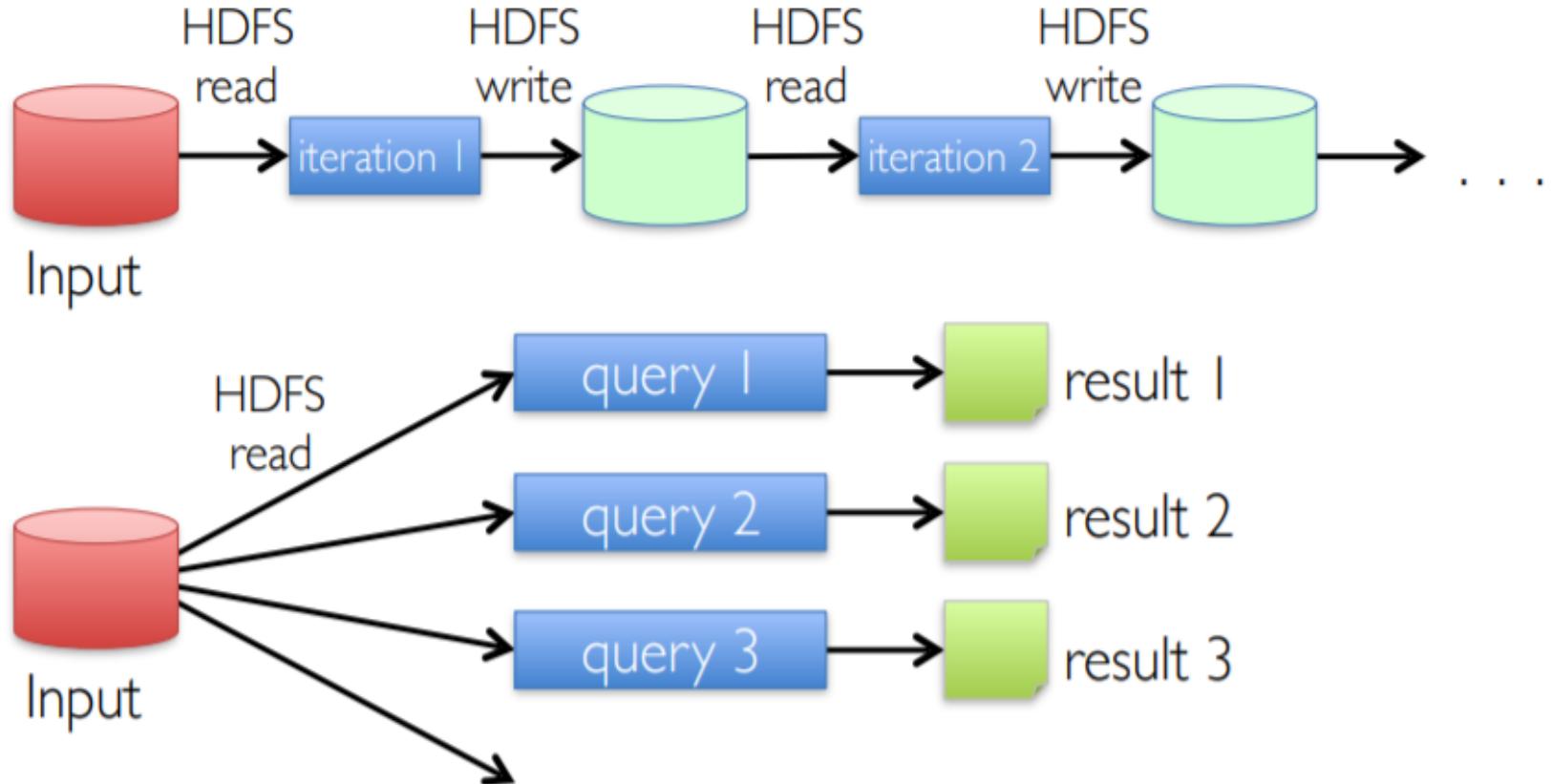
## Why Spark?

# Iterative Jobs

- Disk I/O for each repetition  
→ Slow when executing many small iterations

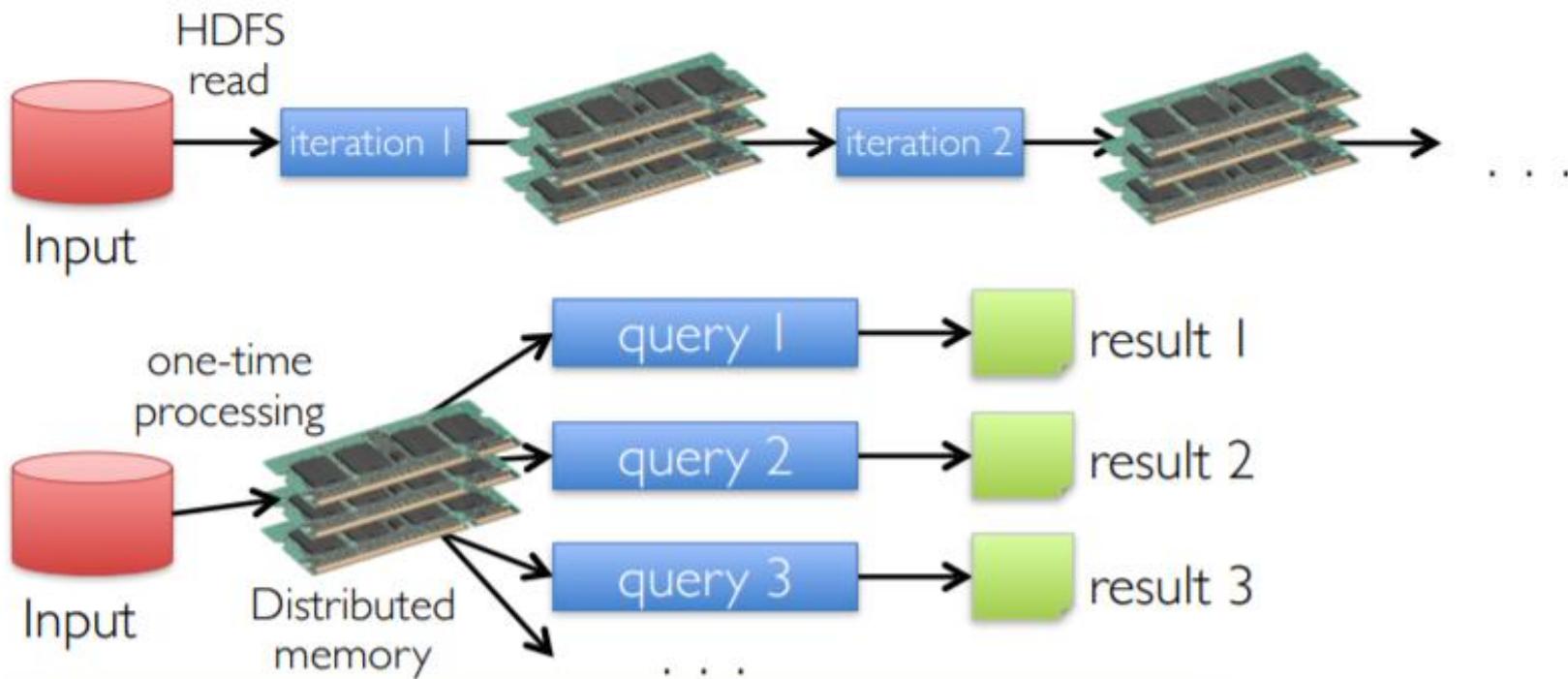


# Replace Disk with Memory

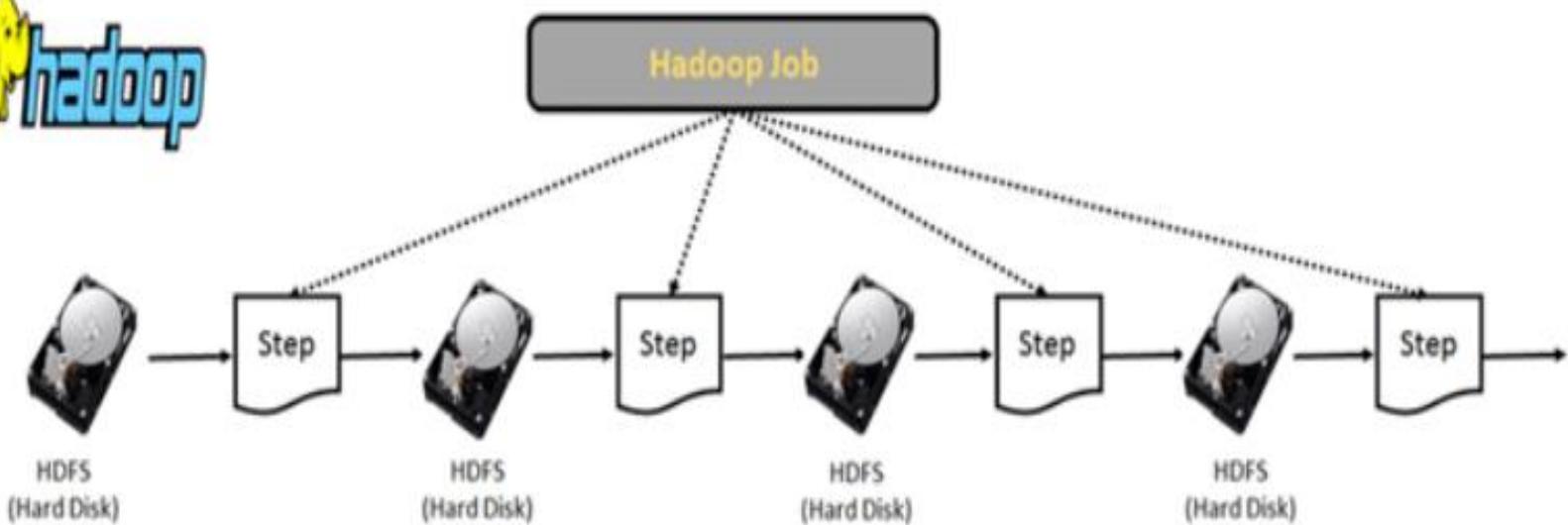
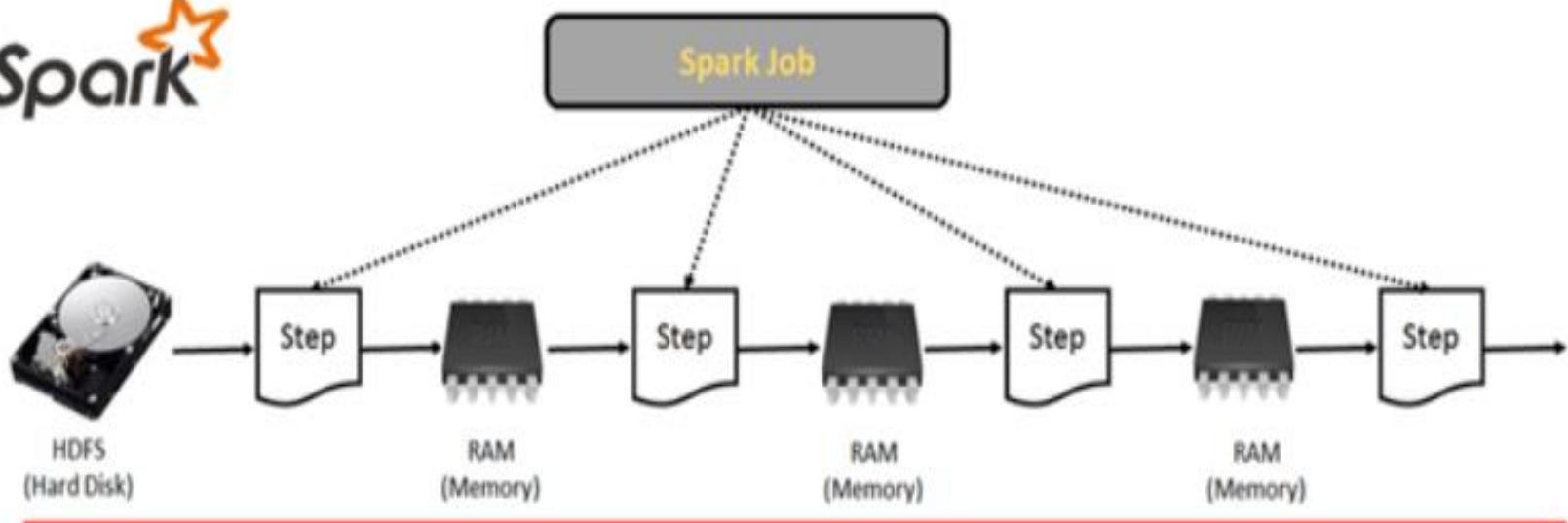


# Replace Disk with Memory

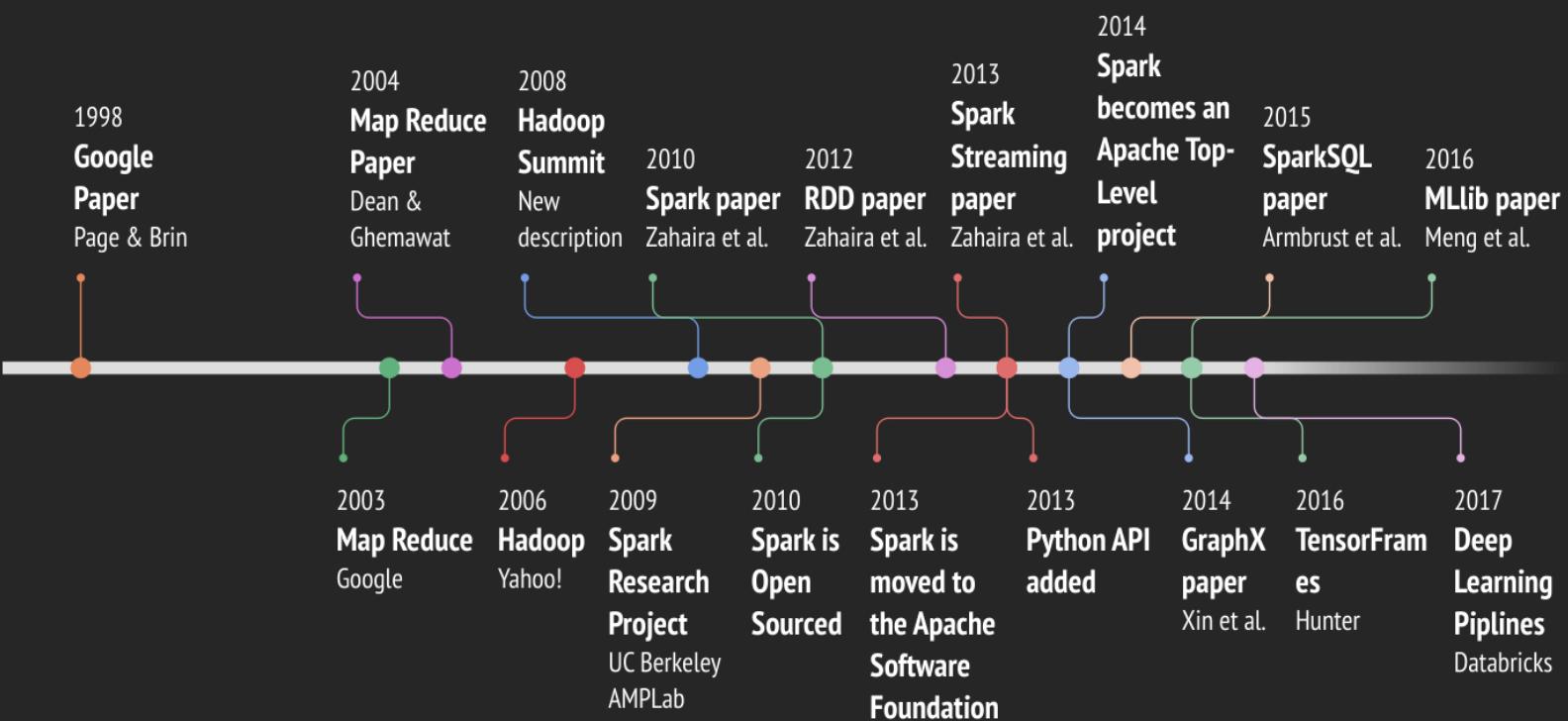
## In-Memory Data Sharing

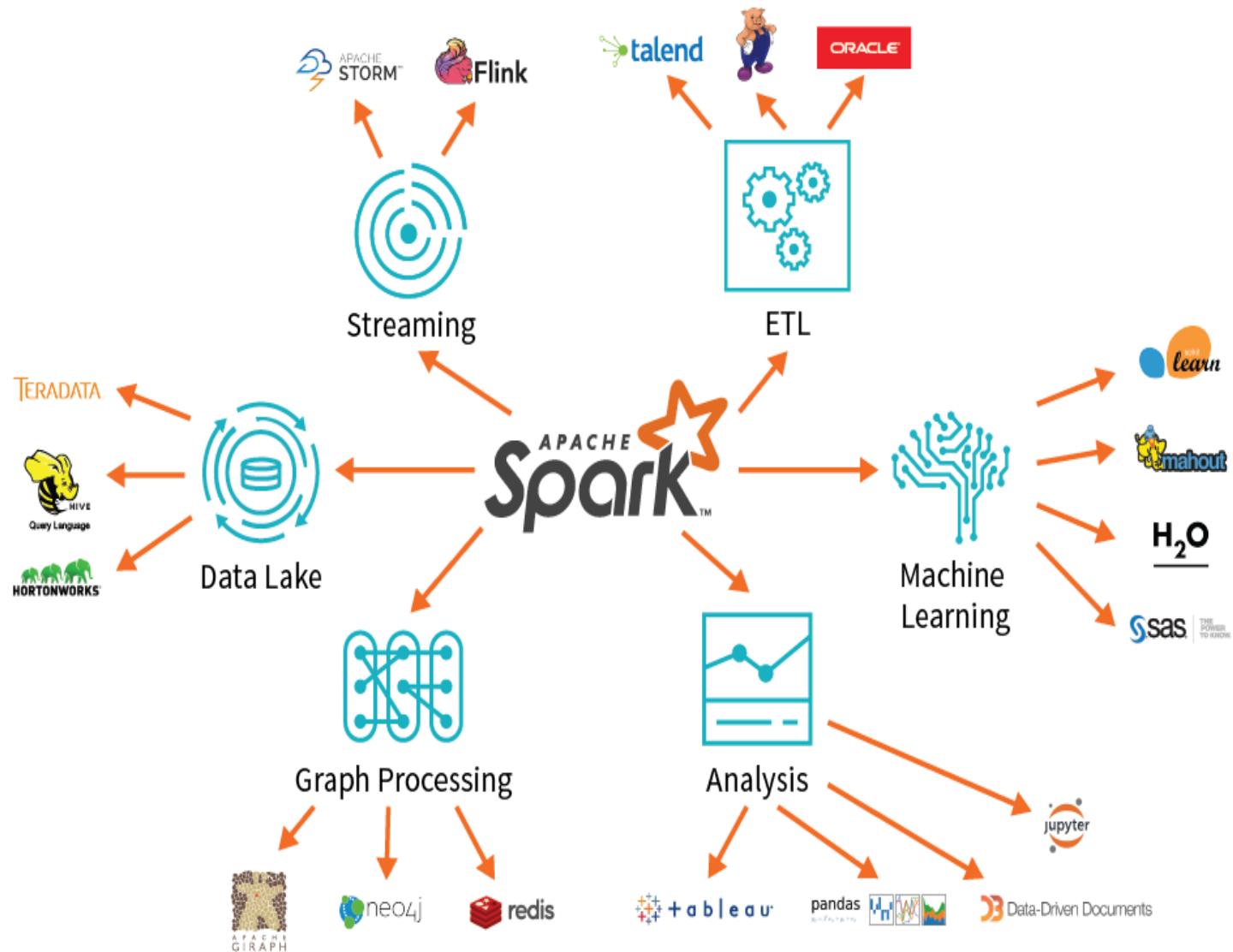


10-100x faster than network and disk



# Apache Spark Timeline

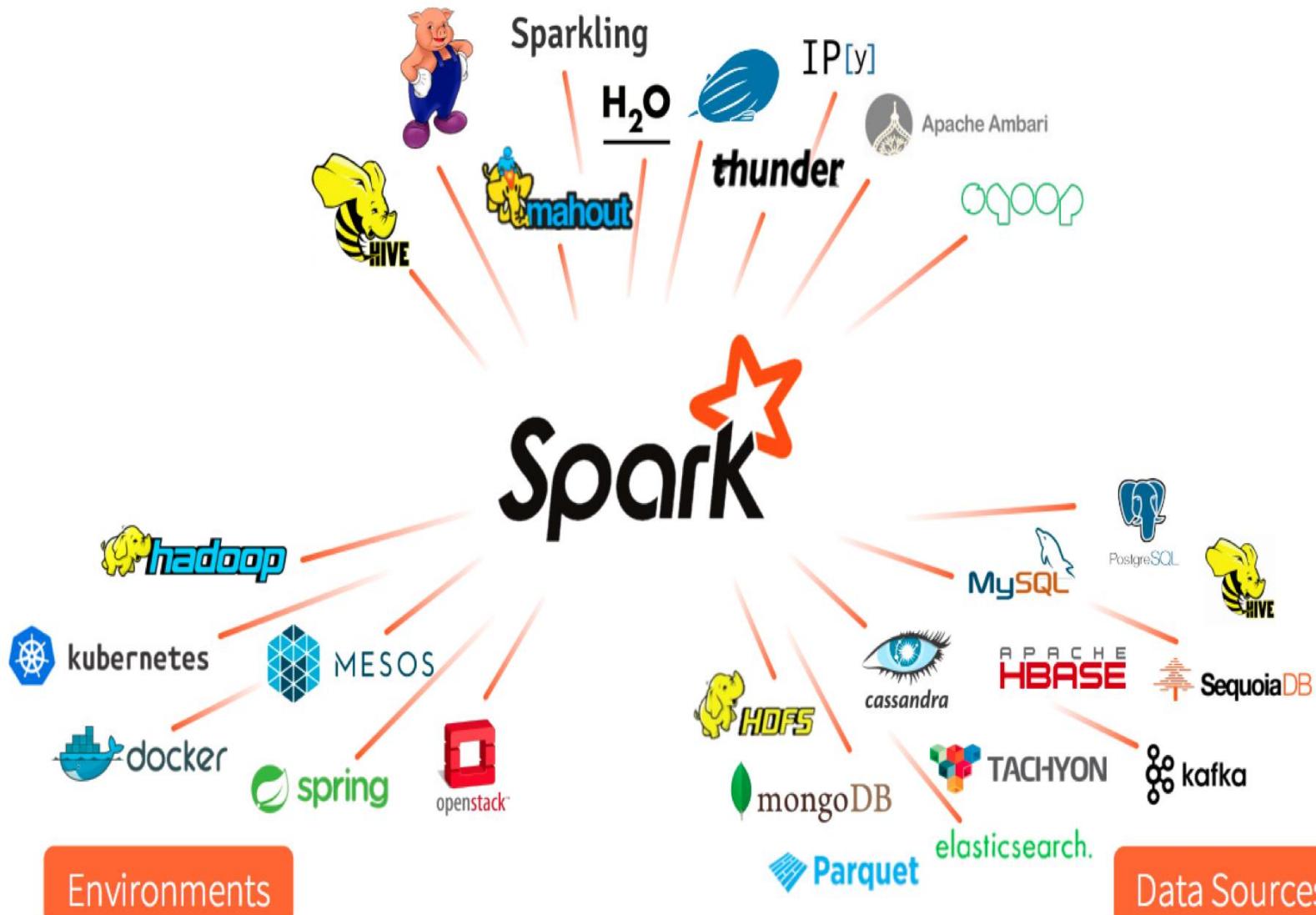




# In-Memory Processing

- Many datasets fit in memory (of a cluster)
  - Memory is fast and avoid disk I/O
- Spark distributed execution engine



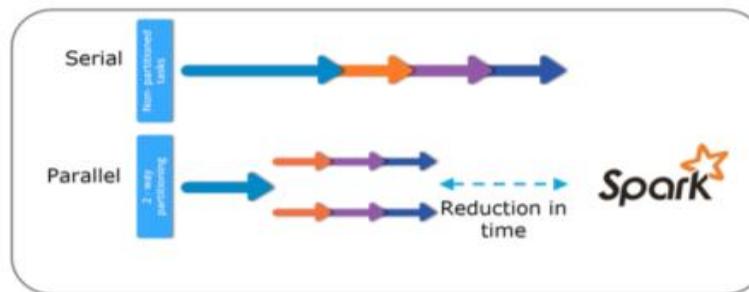


## Introduction to Apache Spark

- Fast, expressive cluster computing system compatible with Apache Hadoop
- It is much faster and much easier than Hadoop MapReduce to use due to its rich APIs
- Large community
- Goes far beyond batch applications to support a variety of workloads:
  - including interactive queries, streaming, machine learning, and graph processing



**Figure:** Real Time Processing In Spark



**Figure:** Data Parallelism In Spark

# Apache Spark Applications



## Uses Cases



Uses Spark Streaming to provide the best-in-class movie streaming and recommendation tool to its users.



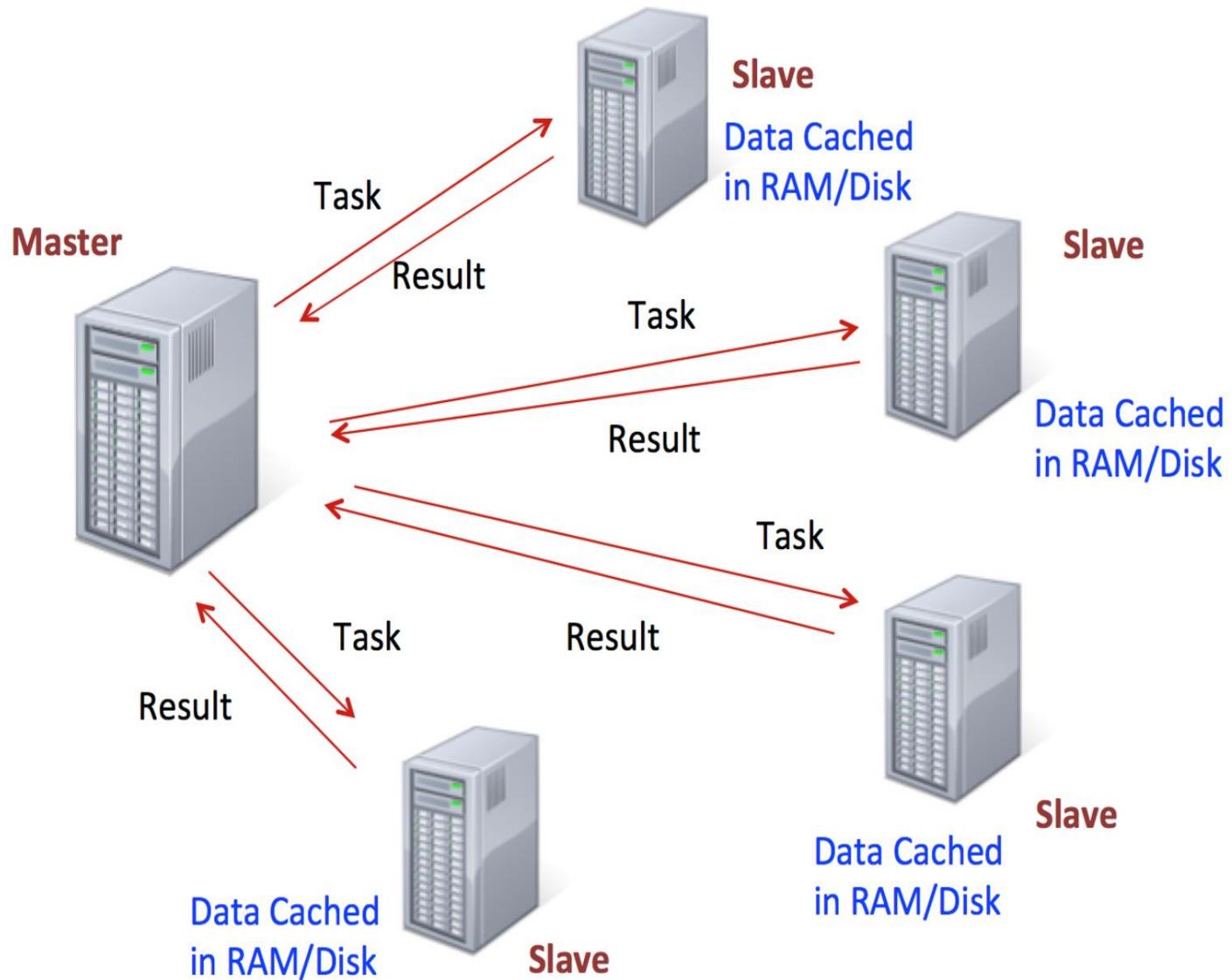
Uses Spark to collect TBs of raw and unstructured data every day from its users to convert it into structured data. This makes it ready for further complex analytics.



Feeds real-time data into Spark via Spark Streaming to get instant insights on how users are engaging with Pins globally. This makes Pinterest's recommendations (i.e. to show Pins) to be accurate.

	Data Warehouse	Hadoop M/R	APACHE Spark
Separate Compute & Storage	✗	✓	✓
More than SQL (i.e ML)	✗	✓	✓
Open Source at Scale	✗	✓	✓
SQL & Optimization	✓	✗	✓
Data Model & Catalog	✓	✗	✓
ACID Transactions	✓	✗	 DELTA LAKE

3.0



HADOOP MAPREDUCE	SPARK
The data is stored in the disc.	The data is stored in-memory.
Computing is based on the disc.	Computing relies on RAM.
Fault tolerance is done through replication.	Fault tolerance is done through RDD.
Hard to work with real-time data.	Easy to work with real-time data.
Less costly in comparison to spark.	More costly.
For batch processing only.	Supports interactive query.

# What is Apache Spark?

Apache Spark is a parallel processing framework and unified analytics engine that supports in-memory processing to boost the performance of big-data analytic and machine learning.

## What is Apache Spark used for?

Spark is used for many types of data processing. It supports ETL, interactive queries (SQL), advanced analytics (e.g. machine learning) and structured streaming over large datasets.

Spark integrates with many storage systems (e.g. HDFS, Cassandra, MySQL, HBase, MongoDB, S3). Spark is also pluggable, with dozens of applications, data sources, and environments



# Spark Architecture & Components



Apache Spark is a powerful open-source processing engine built around speed, ease of use, and sophisticated analytics.

Spark SQL +  
DataFrames

Streaming

MLlib  
*Machine  
Learning*

GraphX  
*Graph  
Computation*

Spark Core API

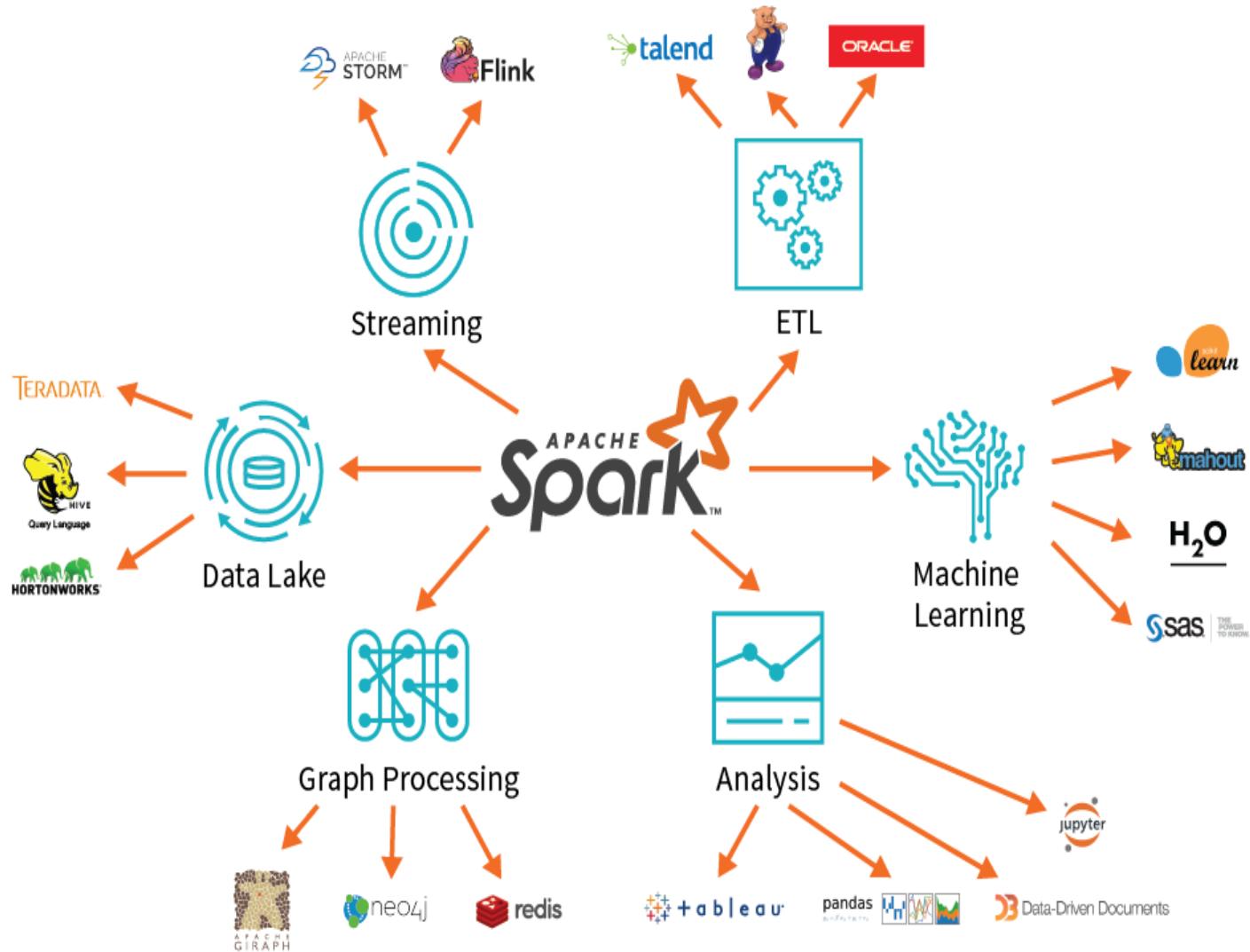
R

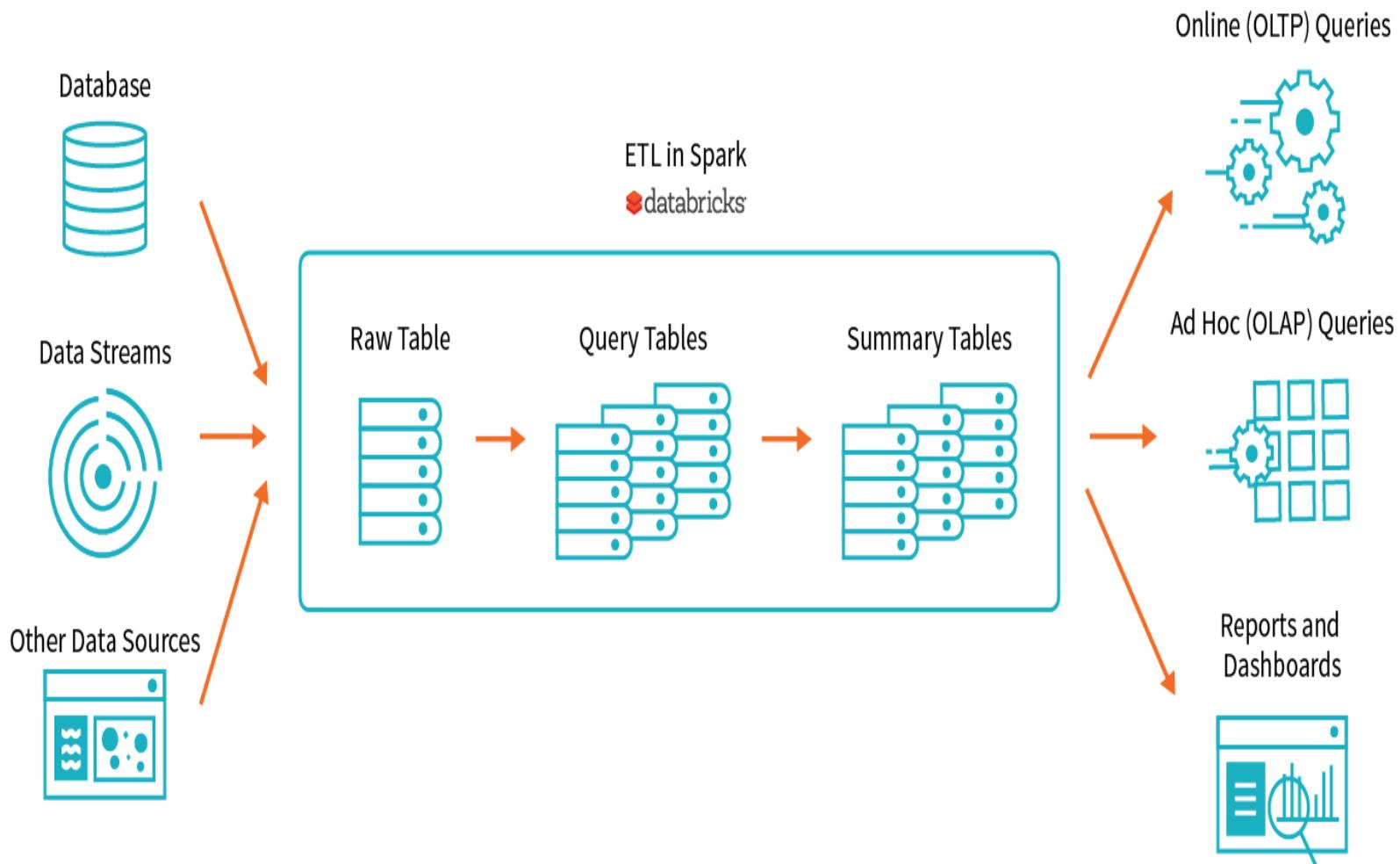
SQL

Python

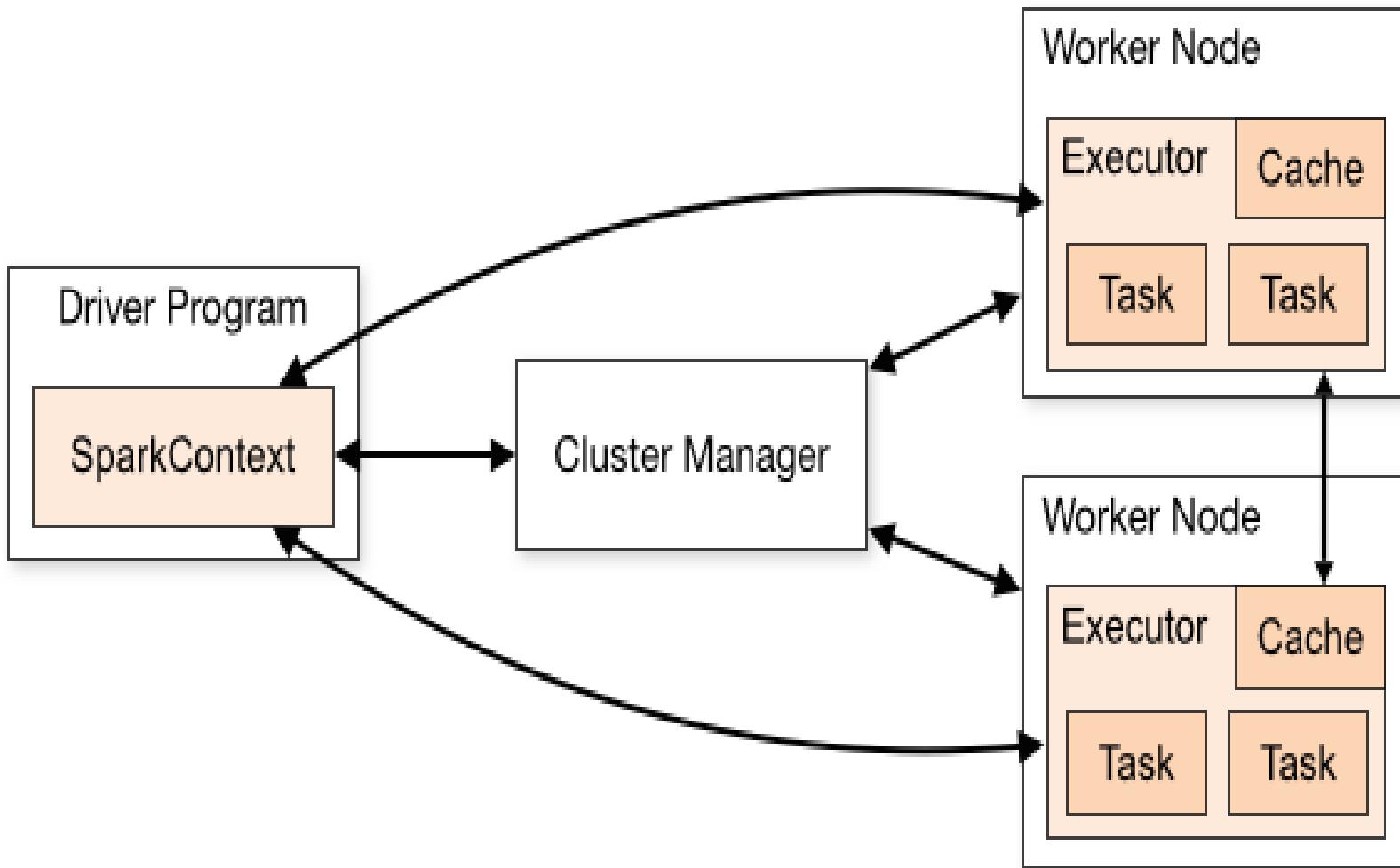
Scala

Java

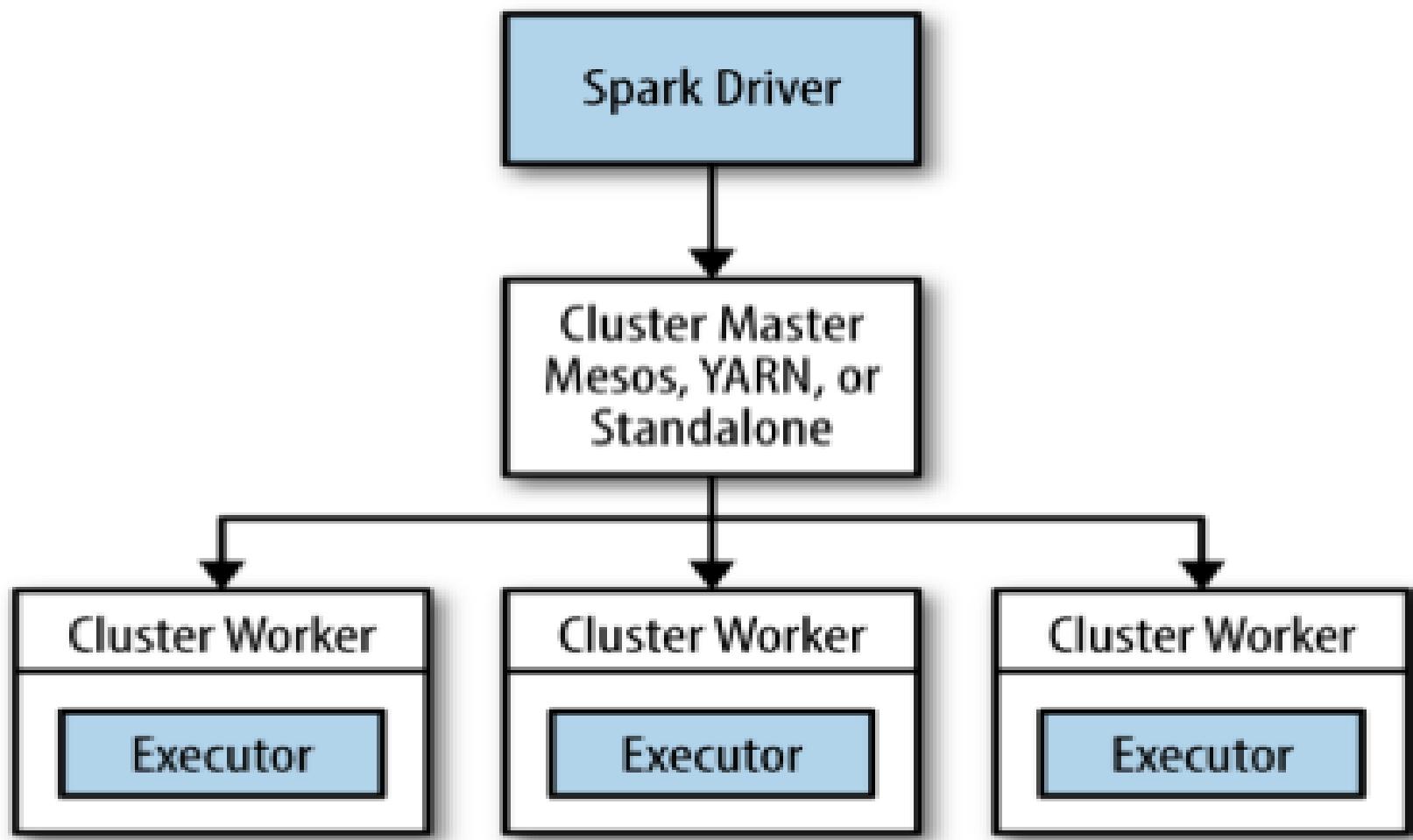


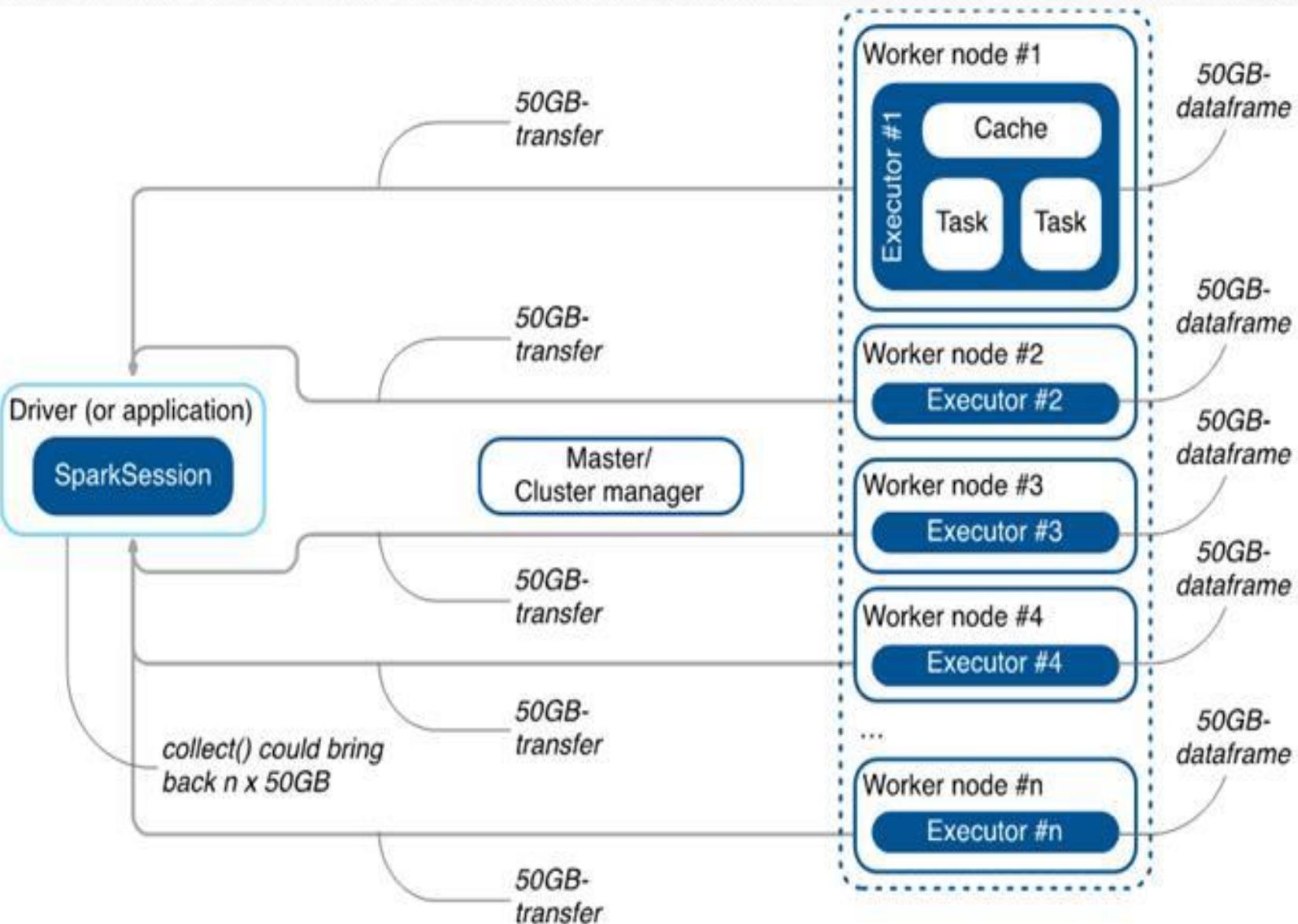


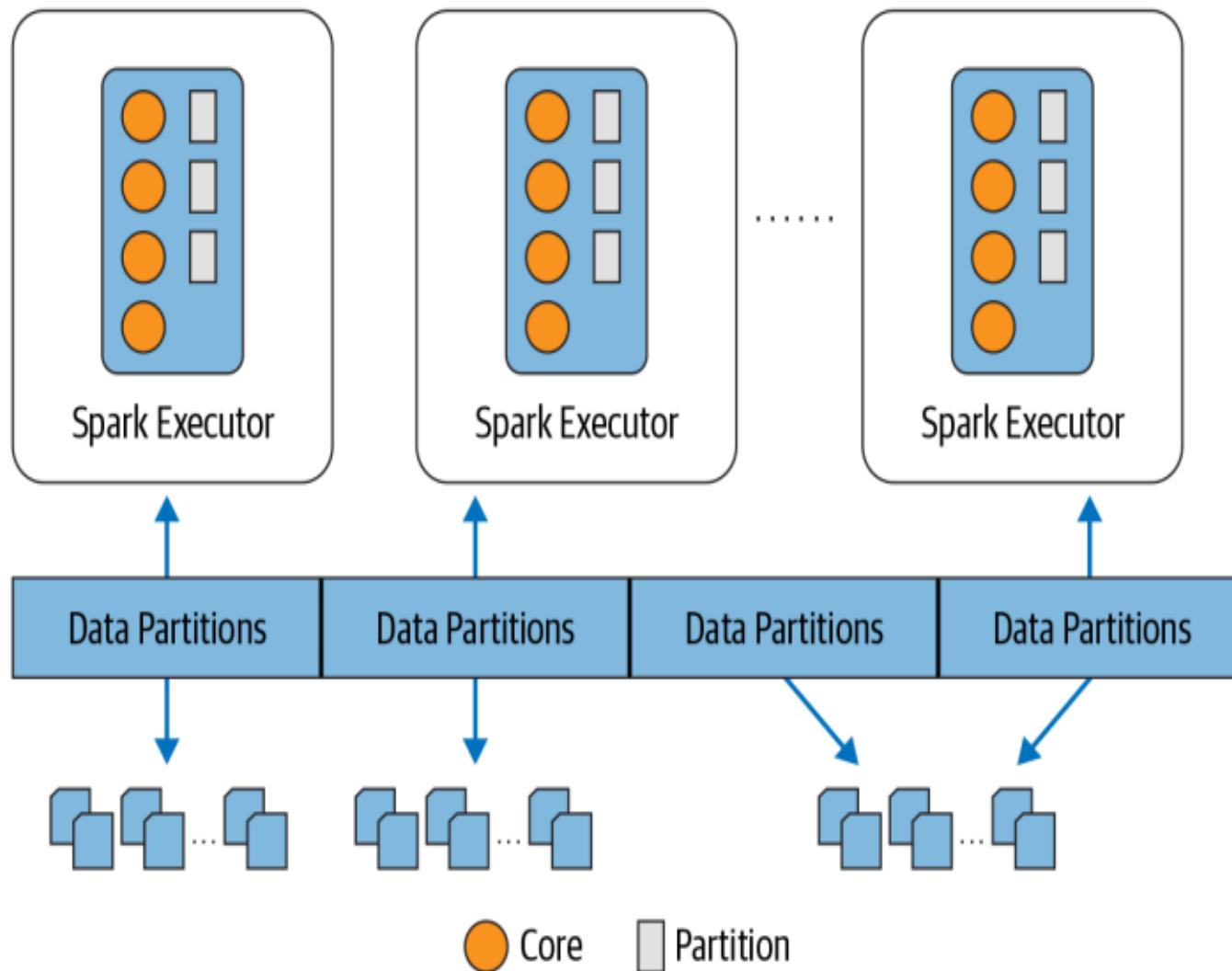
# Spark Components Architecture



In the cluster mode, Spark uses the master/slave architecture. The central coordinator is called the driver, and the driver communicates with the scattered workers. The scattered workers are also called executors. These drivers and executors are collectively called Spark Application.



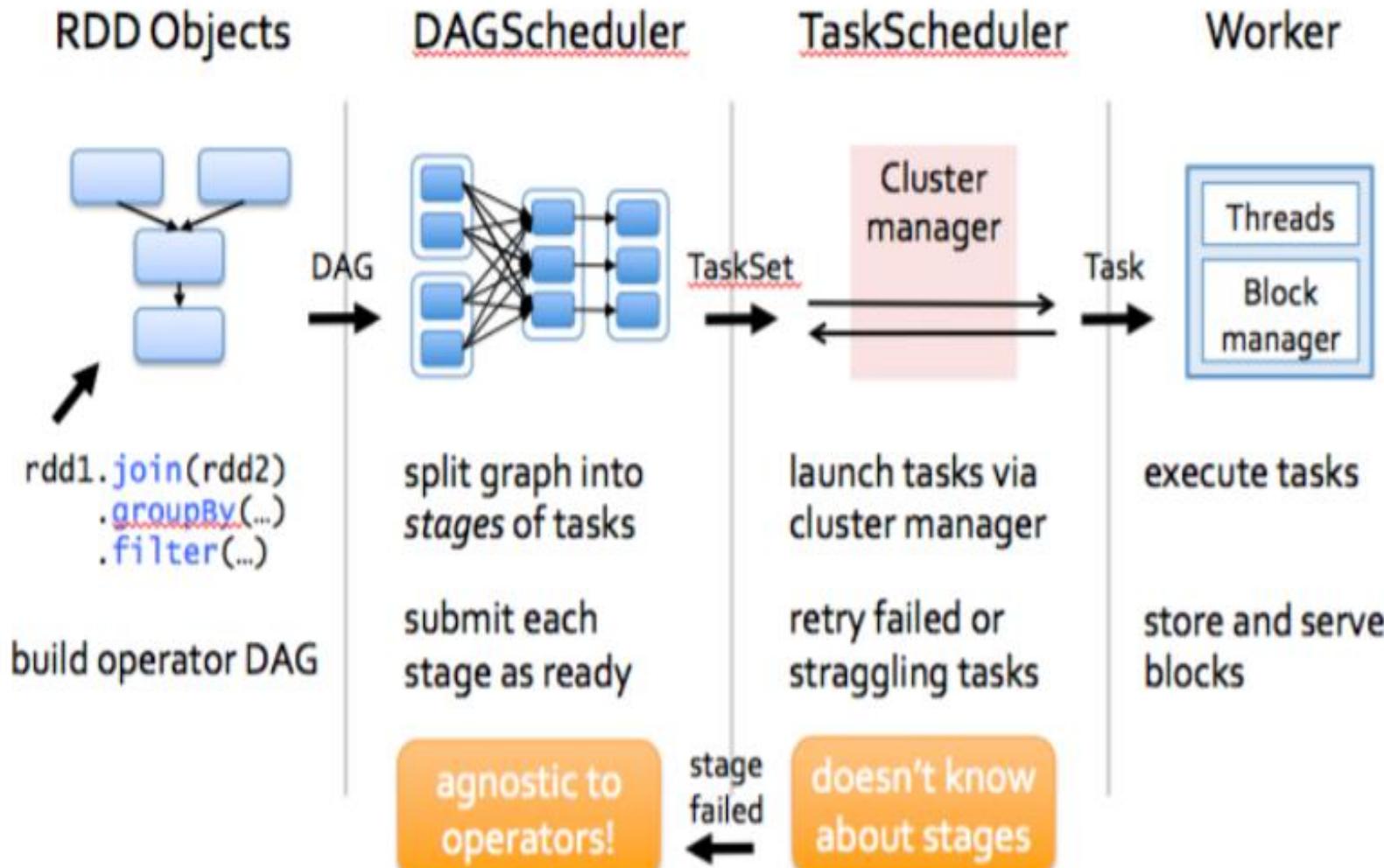




## Spark Core

- Contains the basic functionality for
  - task scheduling,
  - memory management,
  - fault recovery,
  - interacting with storage systems,
  - and more.
- Defines the Resilient Distributed Data sets (RDDs)
  - main Spark programming abstraction.

# SPARK EXECUTION



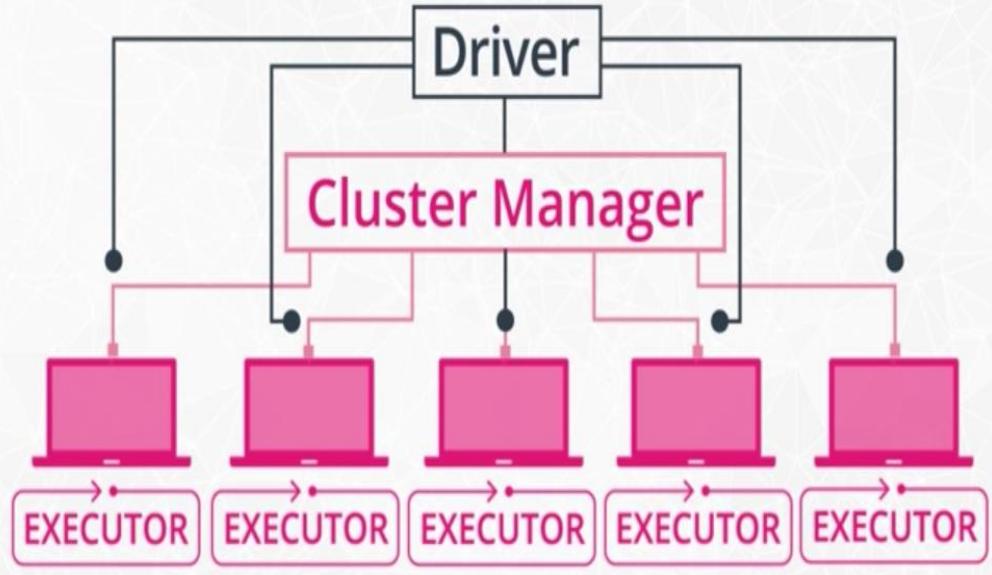
## SPARK EXECUTION PROCESS STEPS

- Job: A piece of code which reads some input from HDFS or local, performs some computation on the data and writes some output data.
- Stages: Jobs are divided into stages. Stages are classified as a Map or reduce stages (Its easier to understand if you have worked on Hadoop and want to correlate). Stages are divided based on computational boundaries, all computations (operators) cannot be Updated in a single Stage. It happens over many stages.
- Tasks: Each stage has some tasks, one task per partition. One task is executed on one partition of data on one executor (machine).
- DAG: DAG stands for Directed Acyclic Graph, in the present context its a DAG of operators.
- Executor: The process responsible for executing a task.
- Master: The machine on which the Driver program runs
- Slave: The machine on which the Executor program runs

## SPARK MODES



LOCAL MODE



## CLUSTER MODES

Standalone

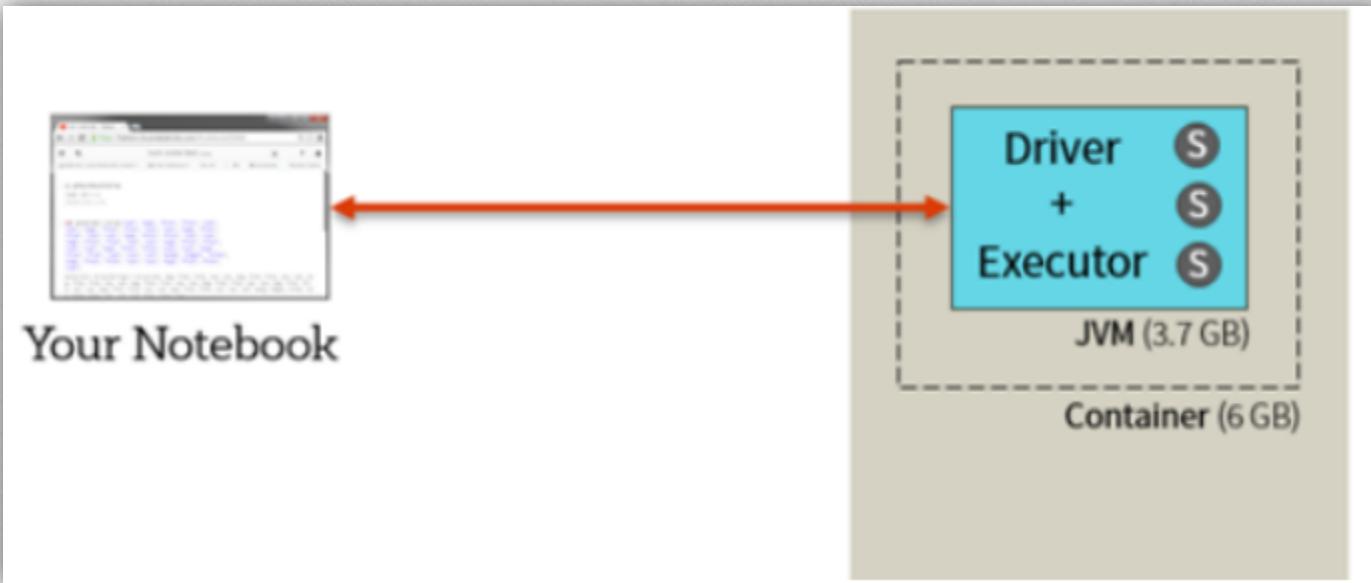
YARN

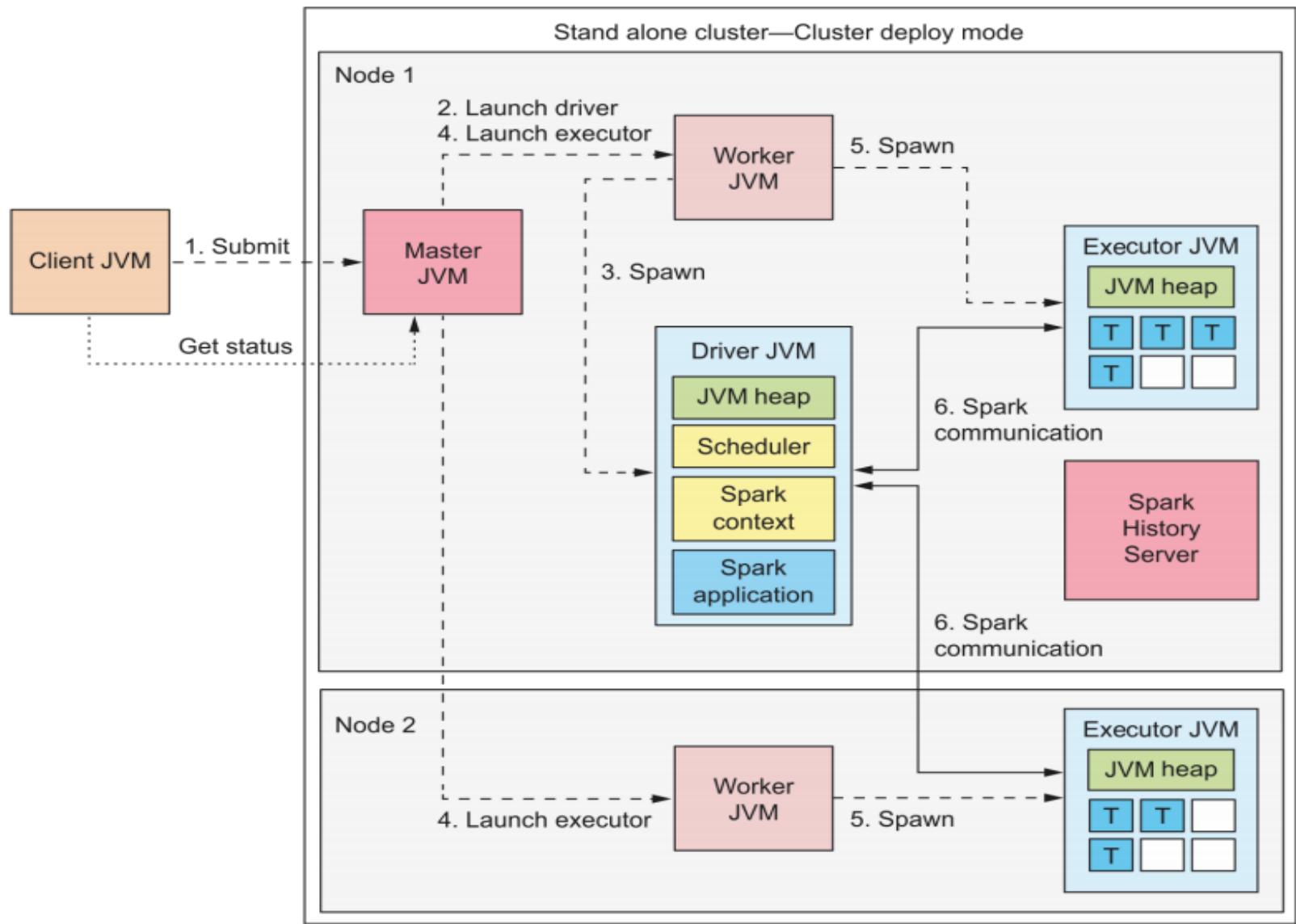
Mesos

# SPARK Cluster Execution Types

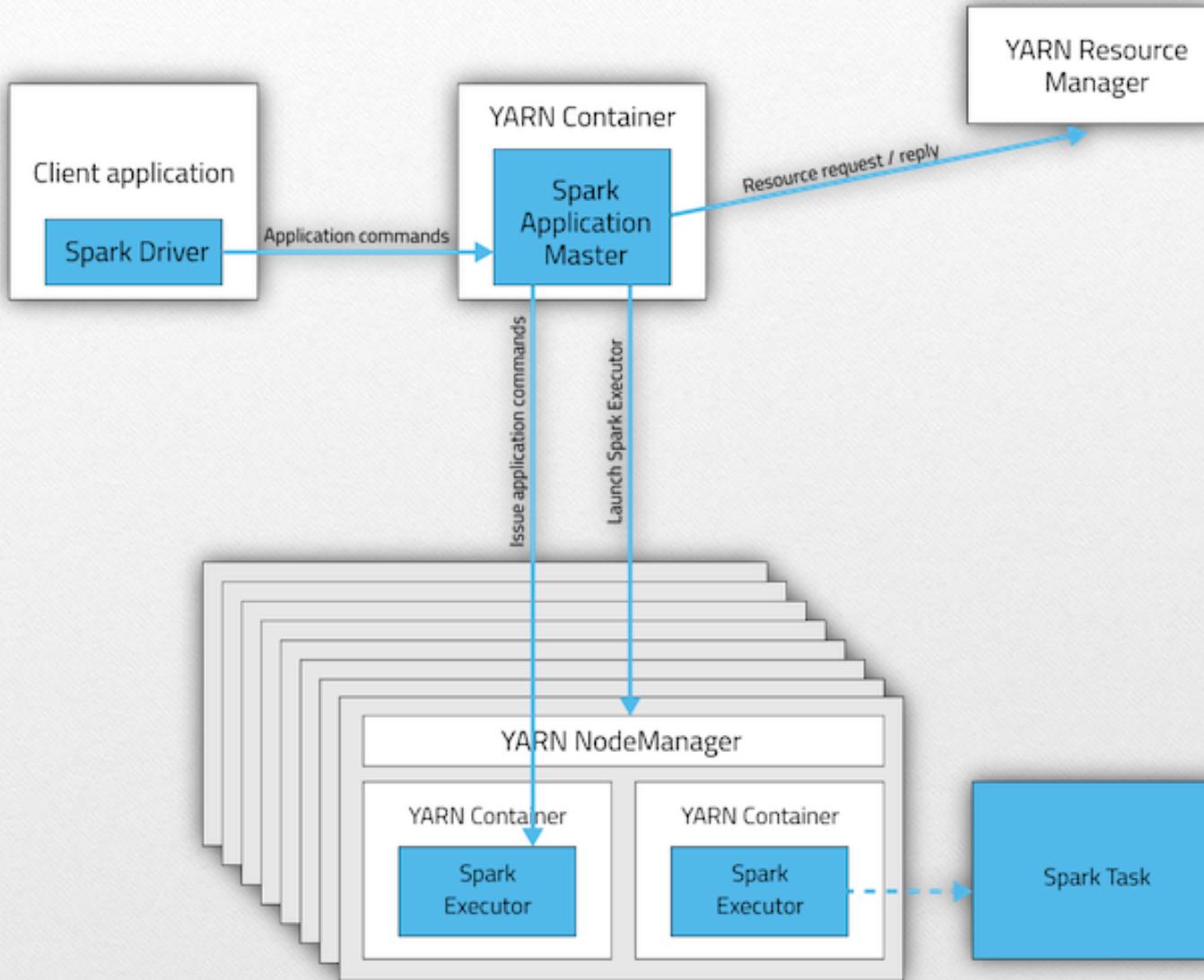
Mode	Spark driver	Spark executor	Cluster manager
Local	Runs on a single JVM, like a laptop or single node	Runs on the same JVM as the driver	Runs on the same host
Standalone	Can run on any node in the cluster	Each node in the cluster will launch its own executor JVM	Can be allocated arbitrarily to any host in the cluster
YARN (client)	Runs on a client, not part of the cluster	YARN's NodeManager's container	YARN's Resource Manager works with YARN's Application Master to allocate the containers on NodeManagers for executors
YARN (cluster)	Runs with the YARN Application Master	Same as YARN client mode	Same as YARN client mode
Kubernetes	Runs in a Kubernetes pod	Each worker runs within its own pod	Kubernetes Master

## Local Mode

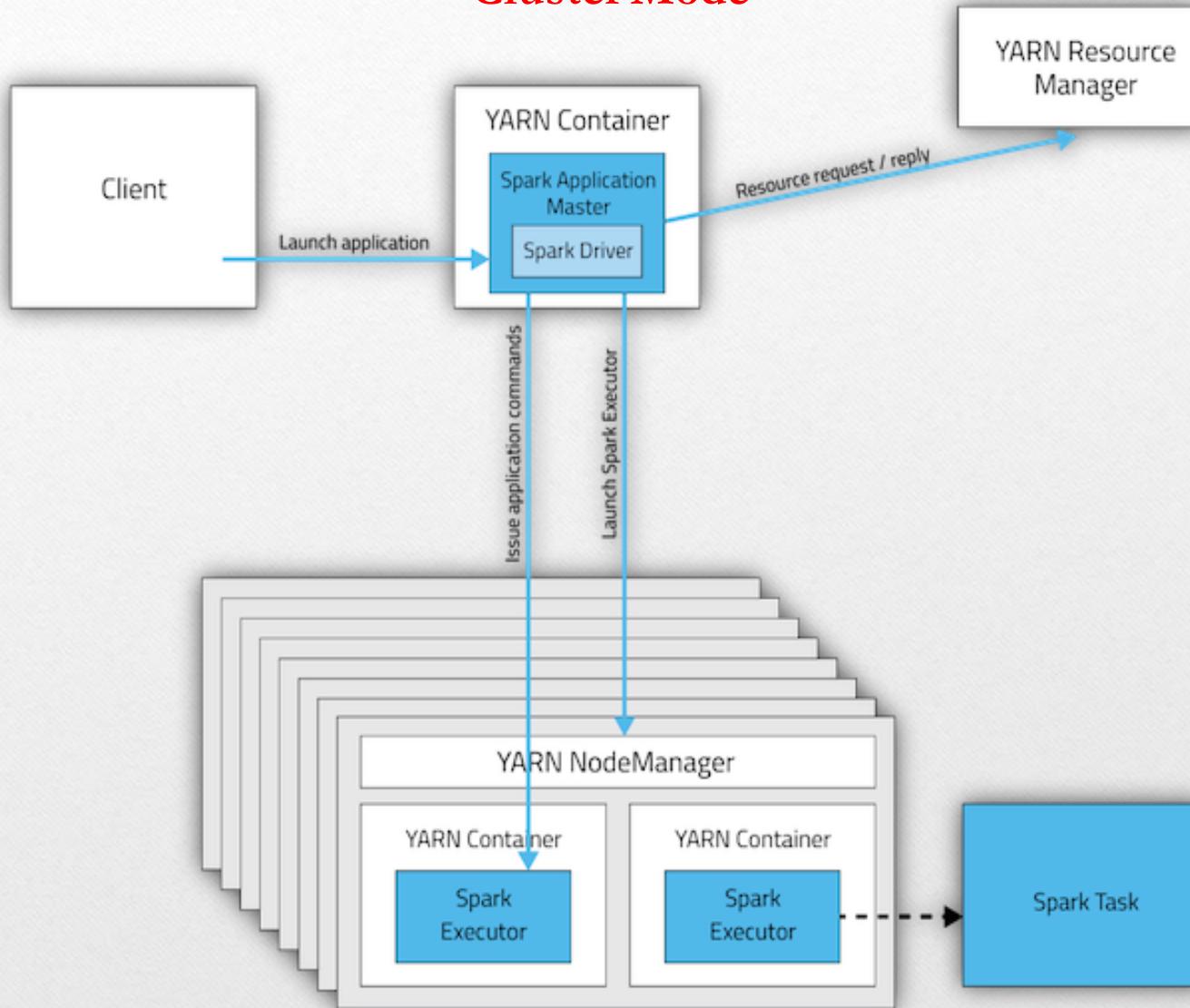




## Client Mode



## Cluster Mode



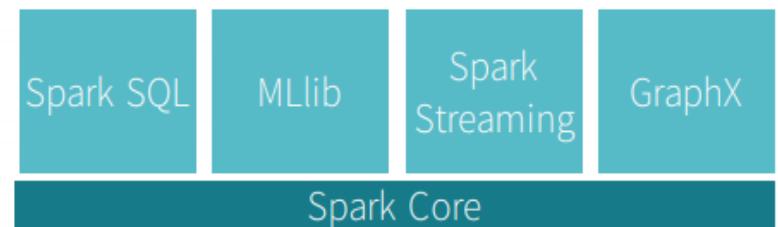
	YARN Cluster	YARN Client	Spark Standalone
<b>Driver runs in:</b>	Application Master	Client	Client
<b>Who requests resources?</b>	Application Master	Application Master	Client
<b>Who starts executor processes?</b>	YARN NodeManager	YARN NodeManager	Spark Slave
<b>Persistent services</b>	YARN ResourceManager and NodeManagers	YARN ResourceManager and NodeManagers	Spark Master and Workers
<b>Supports Spark Shell?</b>	No	Yes	Yes

In **yarn-cluster mode**, the driver runs in the Application Master. This means that the same process is responsible for both driving the application and requesting resources from YARN, and this process runs inside a YARN container. The client that starts the app doesn't need to stick around for its entire lifetime.

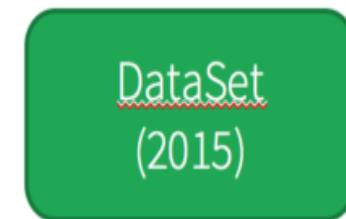
The **yarn-cluster mode**, is not well suited to using Spark interactively. Spark applications that require user input, like spark-shell and PySpark, need the Spark driver to run inside the client process that initiates the Spark application. In yarn-client mode, the Application Master is merely present to request executor containers from YARN. The client communicates with those containers to schedule work after they start:

# Spark Libraries on top of RDDs

- SQL (Spark SQL)
  - Full Hive SQL support with UDF, UDAFs, etc
  - how: Internally keep RDDs of row objects (or RDD of column segments)
- Machine Learning (MLlib)
  - Library of machine learning algorithms
  - how: Cache an RDD, repeatedly iterate it
- Streaming (Spark Streaming)
  - Streaming of real-time data
  - how: Series of RDDs, each containing seconds of real-time data
- Graph Processing (GraphX)
  - Iterative computation on graphs (e.g. social network)
  - how: RDD of Tuple<Vertex, Edge, Vertex> and perform self joins



# History of Spark APIs



Distribute collection  
of JVM objects

Functional Operators (map,  
filter, etc.)

Distribute collection  
of Row objects

Expression-based operations  
and UDFs

Internally rows, externally  
JVM objects

Almost the “Best of both  
worlds”: type safe + fast

Logical plans and optimizer

Fast/efficient internal  
representations

But slower than DF  
Not as good for interactive  
analysis, especially Python

# Spark API's

## 1) Resilient Distributed Dataset (RDD) => (Spark1.0)

An RDD stands for Resilient Distributed Datasets. It is Read-only partition collection of records. RDD is the fundamental data structure of Spark. It allows a programmer to perform in-memory computations on large clusters in a fault-tolerant manner. Thus, speed up the task

**RDD Limitations:** **Handling structured data** ->Unlike Dataframe and datasets, RDDs don't infer the schema of the ingested data and requires the user to specify it.

## 2) DataFrame => (Spark1.3)

DataFrame data organized into named columns. For example a table in a relational database. It is an immutable distributed collection of data. DataFrame in Spark allows developers to impose a structure onto a distributed collection of data, allowing higher-level abstraction.

**Dataframe Limitations:** **Compile-time type safety**

## 3) Dataset => (Spark1.6)

It is an extension to Dataframe API, the latest abstraction which tries to provide best of both RDD and Dataframe. Datasets API provides compile time safety which was not available in Data frames.

**DataSet Provides best of both RDD and Dataframe**

**RDD** (functional programming, type safe),

**DataFrame** (relational model, Query optimization , Tungsten execution, sorting and shuffling)

The Dataset API is available in Scala and Java

# Differences Between RDD & DataFrame & DataSet

## RDD Features:-

- 1) Distributed collection
- 2) Immutable
- 3) Fault tolerant
- 4) Lazy evaluations
- 5) Functional transformations => Transformations and Actions
- 6) Data processing formats => structured as well as unstructured data
- 7) Programming Languages supported => Java, Scala, Python and R.

## DataFrame Features:

- 1) Distributed collection of Row Object
- 2) Data Processing
- 3) Optimization using catalyst optimizer
- 4) Hive Compatibility
- 5) Tungsten
- 6) Programming Languages supported => Java, Scala, Python and R.

## DataSet Features:

- 1) Provides best of both RDD and Dataframe
- 2) Encoders
- 3) Programming Languages supported => Java, Scala
- 4) Type Safety

## Disadvantages Of RDD & DATAFRAME & DATASET

### Disadvantages of RDDs

If you choose to work with RDD you will have to optimize each and every RDD. In addition, unlike Datasets and DataFrames, RDDs don't infer the schema of the data ingested therefore you will have to specify it.

### Disadvantages of DataFrames

The main drawback of DataFrame API is that it does not support compile time safely, as a result, the user is limited in case the structure of the data is not known.

### Disadvantages of DataSets

The main disadvantage of datasets is that they require typecasting into strings.

# Spark Context

- A Spark program first creates a **SparkContext** object
  - » Tells Spark how and where to access a cluster
  - » pySpark shell and Databricks Cloud automatically create the **sc** variable
  - » [iPython](#) and programs must use a constructor to create a new **SparkContext**
- Use **SparkContext** to create RDDs

## Creating RDD And Working on RDD

Spark Operations =

+

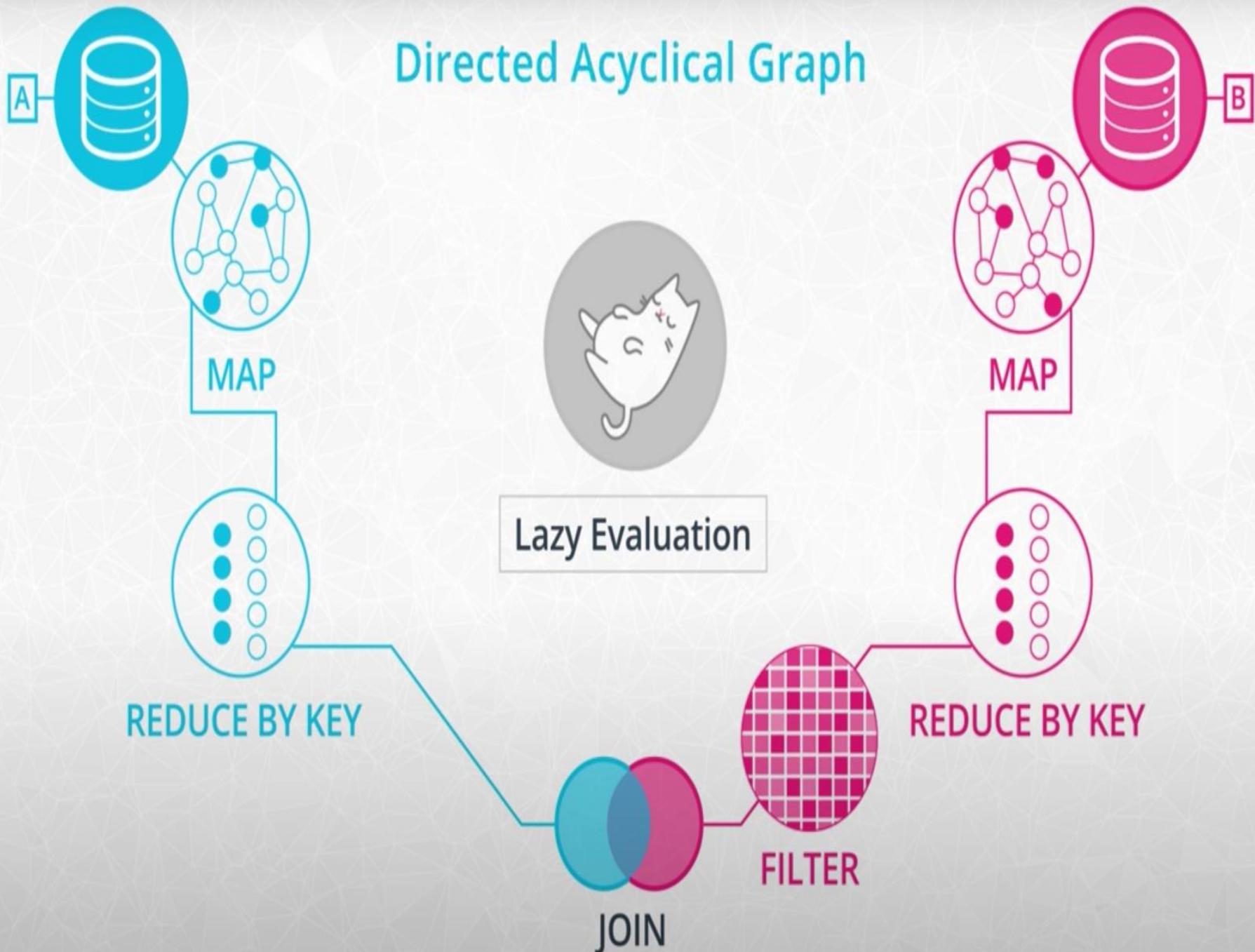


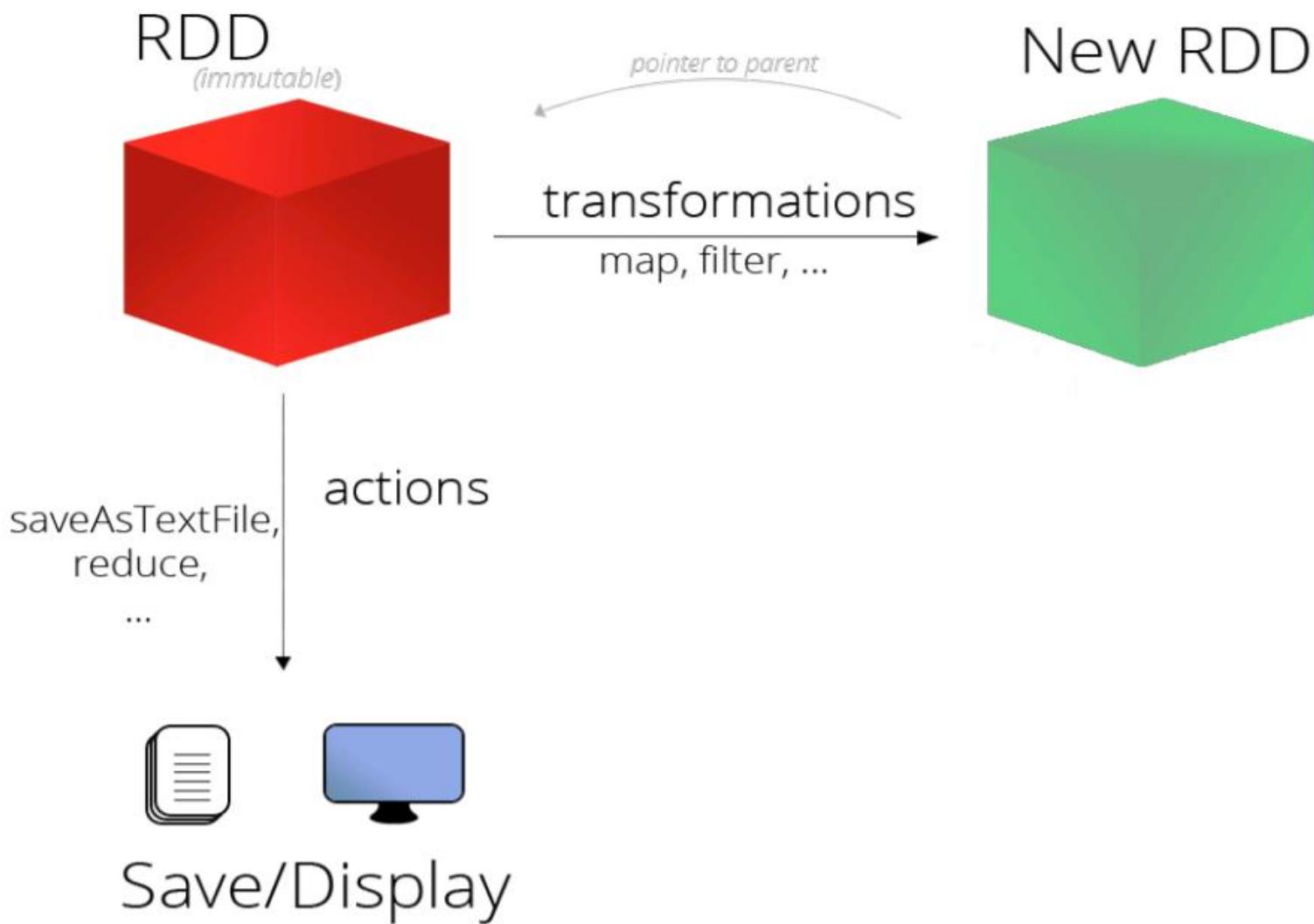
ACTIONS

# Operations on RDDs

- Transformations: lazy execution
  - Map, filter, intersection, groupByKey, zipWithIndex ...
- Actions: trigger execution of transformations
  - Collect, count, reduce, saveAsTextFile ...
- Caching
  - Avoid multiple executions of transformations when re-using RDD

# Directed Acyclical Graph





# An RDD can be created 2 ways

- Parallelize a collection

```
# Parallelize in Python  
wordsRDD = sc.parallelize(["fish", "cats", "dogs"])
```

- Take an existing in-memory collection and pass it to SparkContext's parallelize method
- Not generally used outside of prototyping and testing since it requires entire dataset in memory on one machine

- Read from File

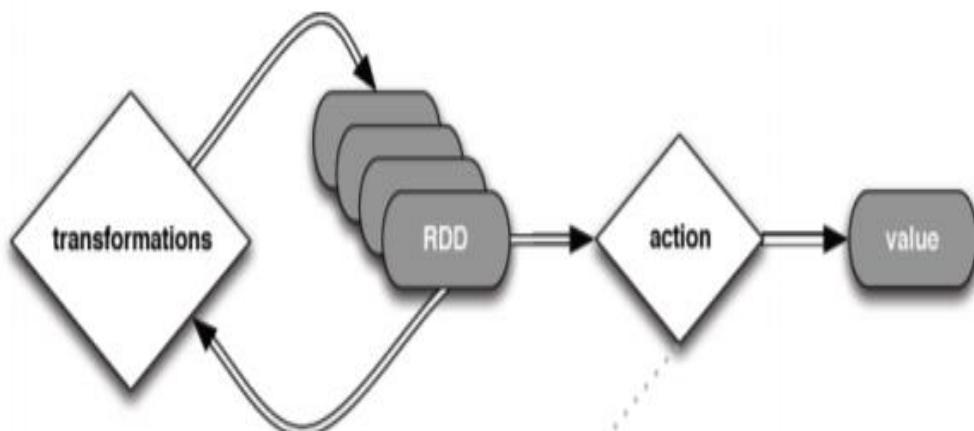
```
# Read a local txt file in Python  
linesRDD = sc.textFile("/path/to/README.md")
```

- There are other methods to read data from HDFS, C\*, S3, HBase, etc.

# Spark Core

## RDD – Resilient Distributed Dataset

- A primary abstraction in Spark – a fault-tolerant collection of elements that can be operated on in parallel.
- Two Types
  - Parallelized Scala collections
  - Hadoop datasets
- Transformations and Actions can be performed on RDDs.



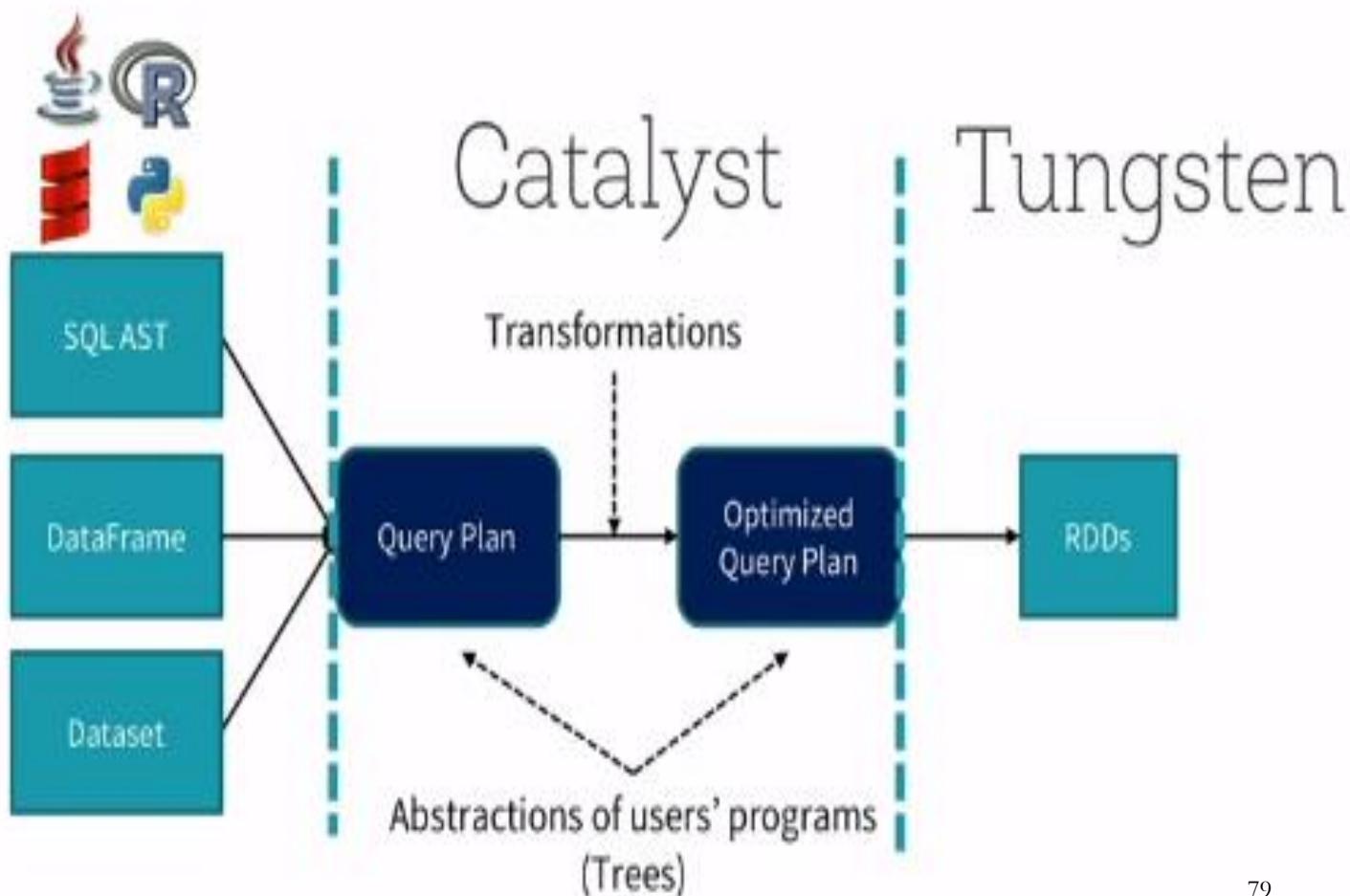
### Transformations

- Operate on an RDD and return a new RDD.
- Are Lazily Evaluated

### Actions

- Return a value after running a computation on a RDD.
- The DAG is evaluated only when an action takes place.

# Spark SQL Overview



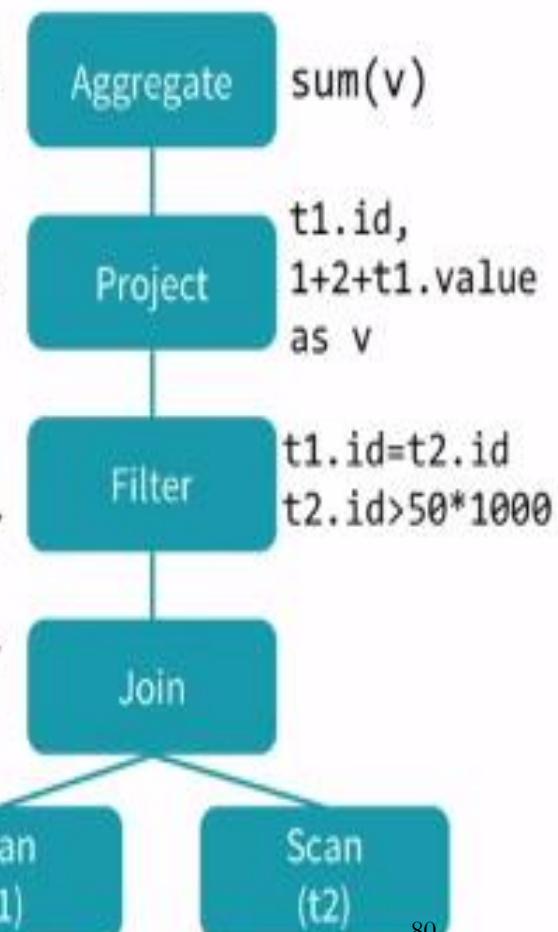
# Trees: Abstractions of Users' Programs

## Query Plan

```

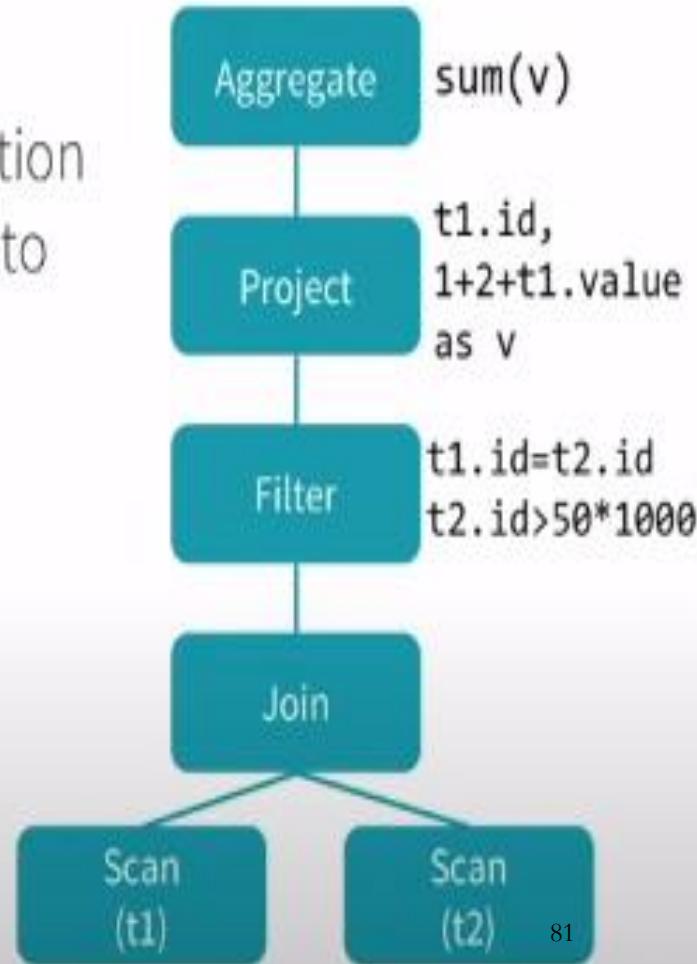
SELECT sum(v)
FROM (
    SELECT
        t1.id,
        1 + 2 + t1.value AS v
    FROM t1 JOIN t2
    WHERE
        t1.id = t2.id AND
        t2.id > 50 * 1000) tmp

```

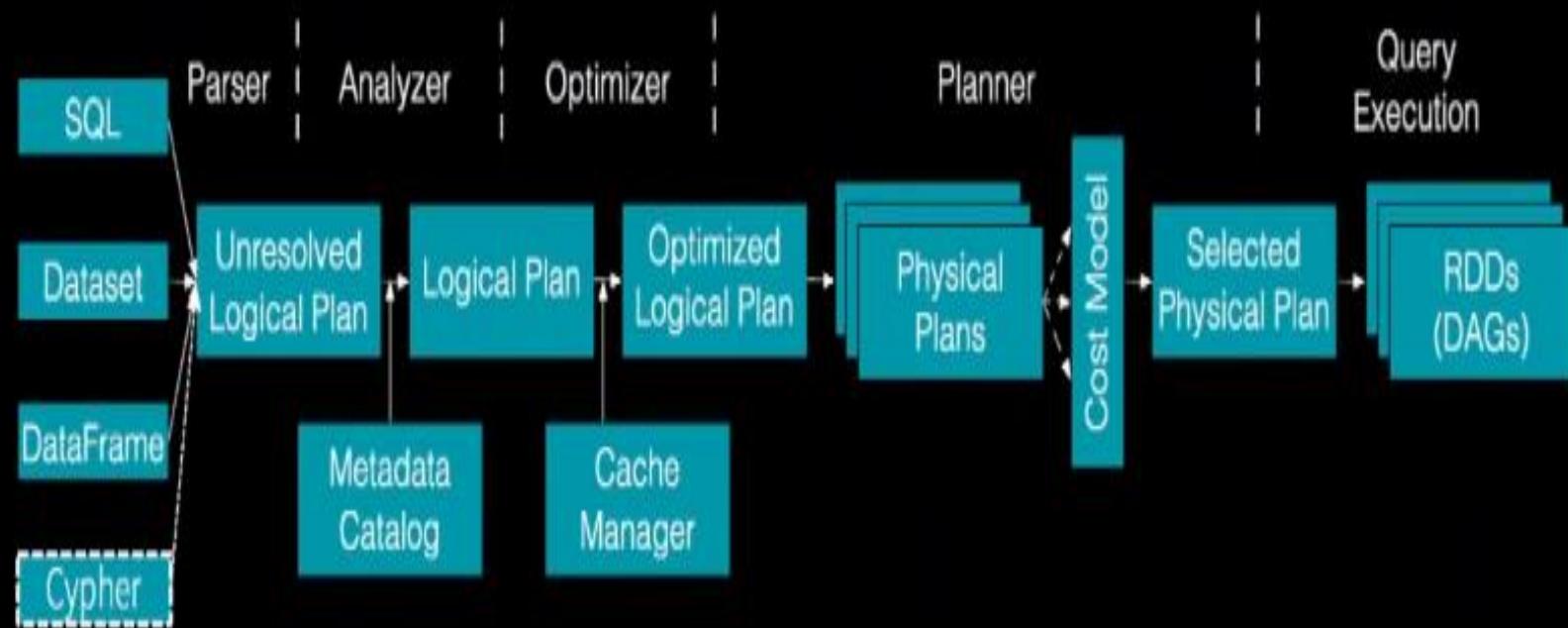


# Logical Plan

- A Logical Plan describes computation on datasets **without** defining how to conduct the computation



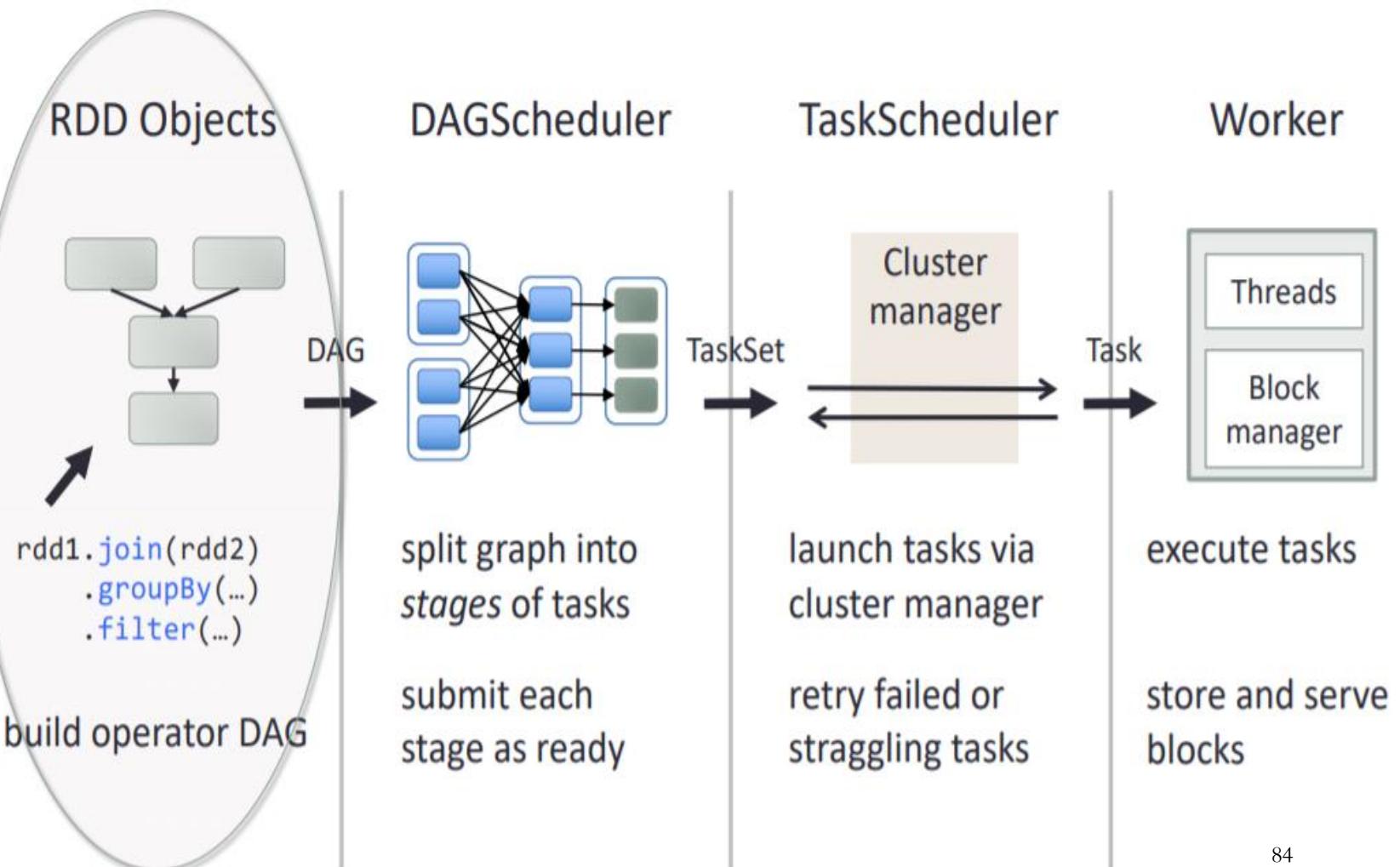
# From declarative queries to RDDs



## Programming with RDDs

- All work is expressed as either:
  - creating new RDDs
  - transforming existing RDDs
  - calling operations on RDDs to compute a result.
- Distributes the data contained in RDDs across the nodes (executors) in the cluster and parallelizes the operations.
- Each RDD is split into multiple **partitions**, which can be computed on different nodes of the cluster.

## Job scheduling

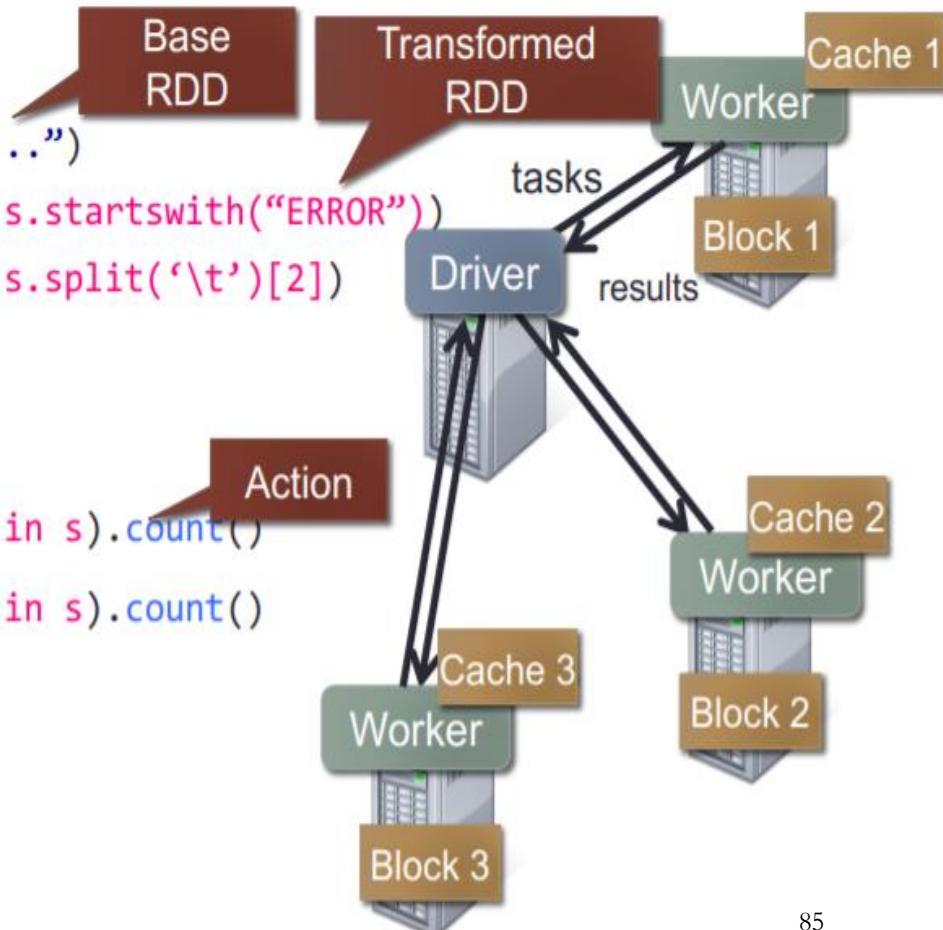


# Example: Mining Console Logs

Load error messages from a log into memory, then interactively search for patterns

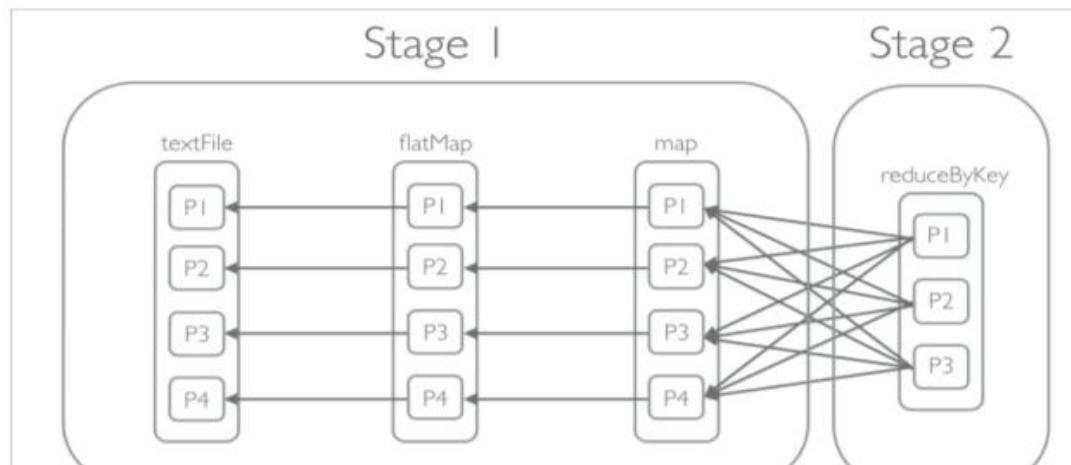
```
lines = spark.textFile("hdfs://...")  
errors = lines.filter(lambda s: s.startswith("ERROR"))  
messages = errors.map(lambda s: s.split('\t')[2])  
messages.cache()
```

```
messages.filter(lambda s: "foo" in s).count()  
messages.filter(lambda s: "bar" in s).count()  
...
```



# Partitions

- Each RDD is split up into a number of partitions
- Parallelism is determined by the number of partitions
  - `rdd.getNumberOfPartitions()`
- Why is important ?
  - Operations can be performed in parallel in each partition:
    - Textfile + flatMap + map operations can be performed in parallel in P1, P2, P3, P4
    - Operations that can run on the same partition are executed in stages



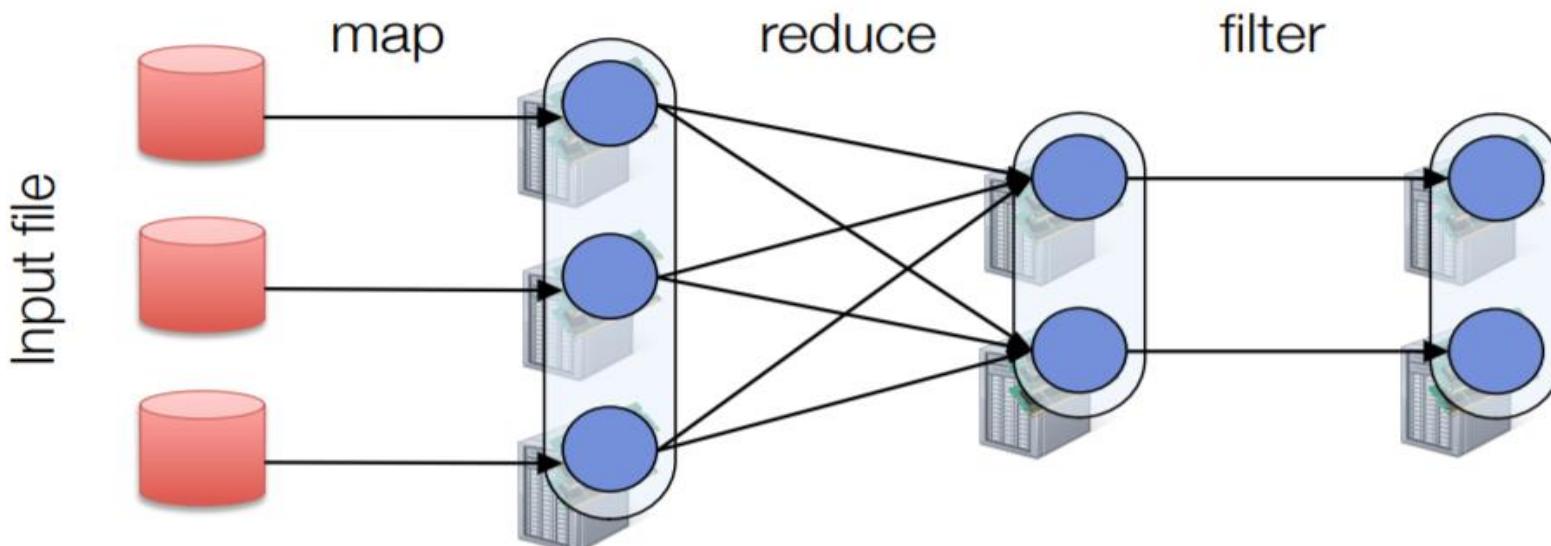
# RDD Fault Tolerance

- Hadoop conservative/pessimistic approach
  - Go to disk / stable storage (HDFS)
- Spark optimistic
  - Don't "waste" time writing to disk, re-compute in case of crash using lineage

# Fault Tolerance

RDDs track *lineage* info to rebuild lost data

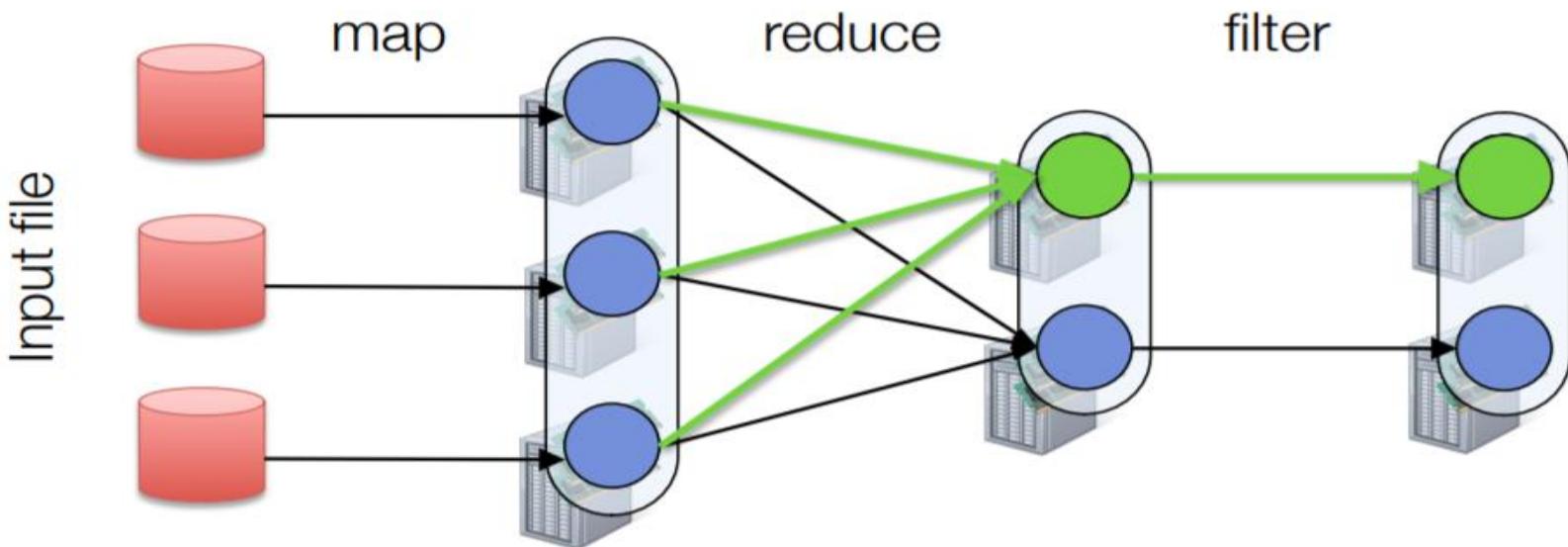
```
file.map(lambda rec: (rec.type, 1))  
    .reduceByKey(lambda x, y: x + y)  
    .filter(lambda (type, count): count > 10)
```



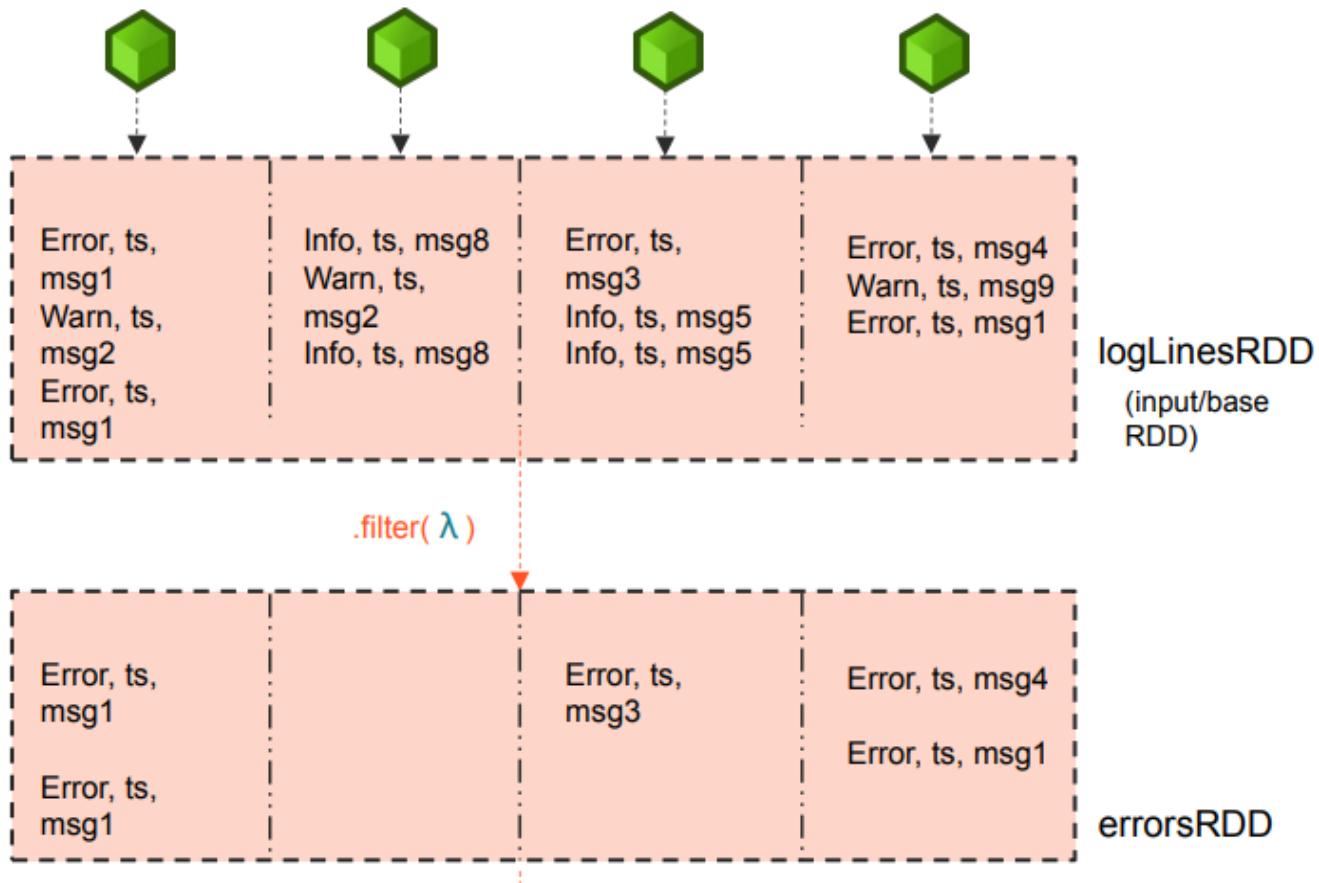
# Fault Tolerance

RDDs track *lineage* info to rebuild lost data

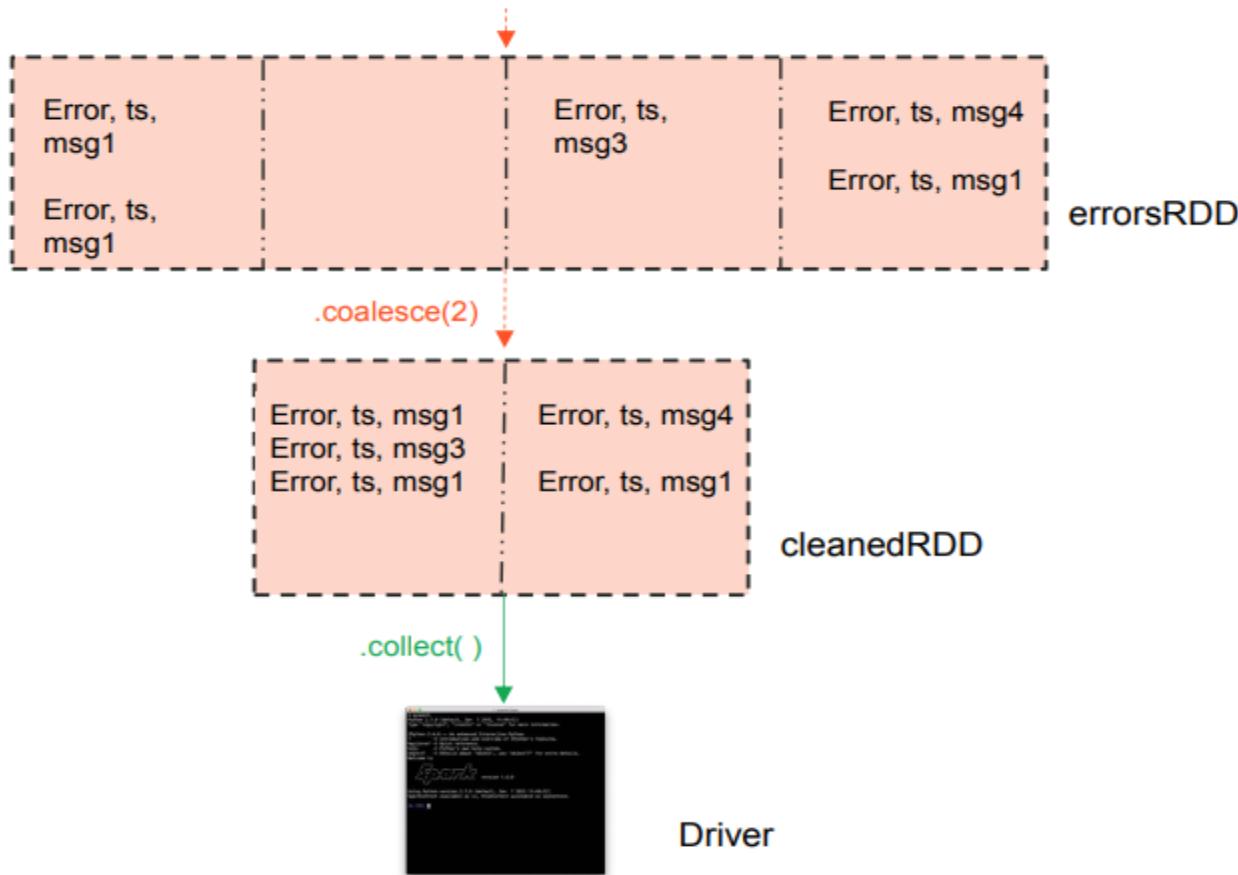
```
file.map(lambda rec: (rec.type, 1))  
    .reduceByKey(lambda x, y: x + y)  
    .filter(lambda (type, count): count > 10)
```



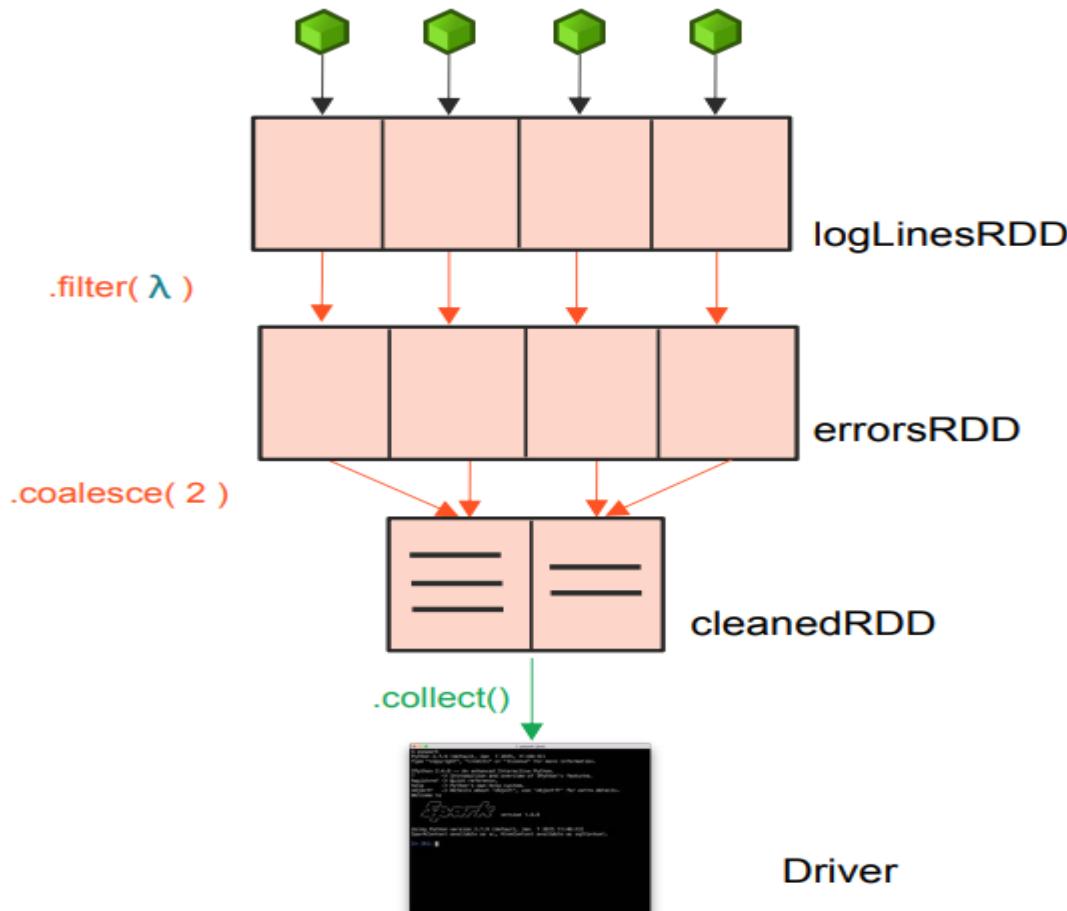
# RDD Example



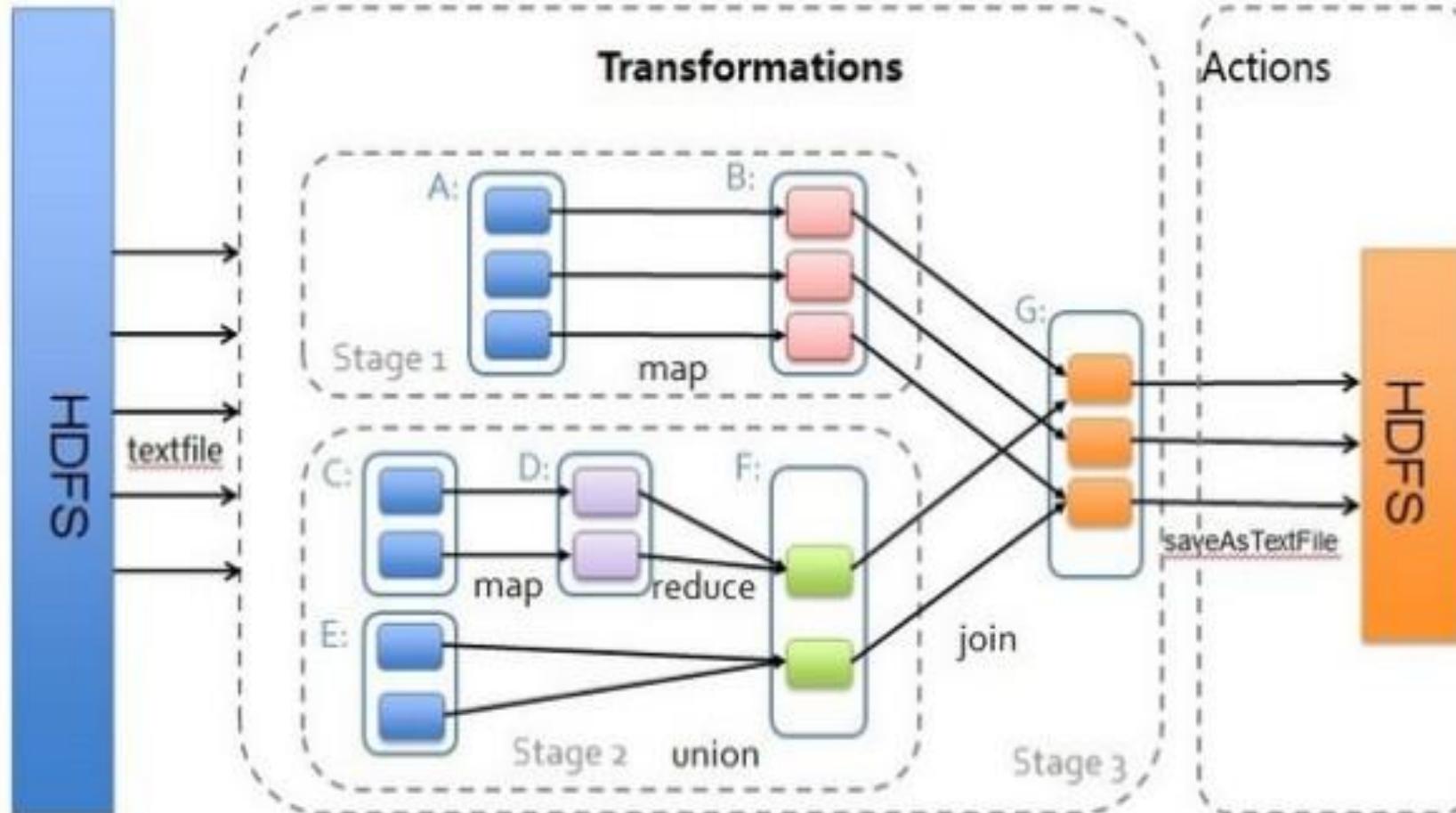
# RDD Example



# RDD Lineage



## Spark: Transformations & Actions



# Spark End to End operations





## ACTIONS



## General

- map
- filter
- flatMap
- mapPartitions
- mapPartitionsWithIndex
- groupBy
- sortBy

## Math / Statistical

- sample
- randomSplit

## Set Theory / Relational

- union
- intersection
- subtract
- distinct
- cartesian
- zip

## Data Structure / I/O

- keyBy
- zipWithIndex
- zipWithUniqueId
- zipPartitions
- coalesce
- repartition
- repartitionAndSortWithinPartitions
- pipe

- reduce
- collect
- aggregate
- fold
- first
- take
- foreach
- top
- treeAggregate
- treeReduce
- foreachPartition
- collectAsMap

- count
- takeSample
- max
- min
- sum
- histogram
- mean
- variance
- stdev
- sampleVariance
- countApprox
- countApproxDistinct

- takeOrdered

- saveAsTextFile
- saveAsSequenceFile
- saveAsObjectFile
- saveAsHadoopDataset
- saveAsHadoopFile
- saveAsNewAPIHadoopDataset
- saveAsNewAPIHadoopFile

## Essential Core &amp; Intermediate Spark Operations



= easy



= medium



## ACTIONS



# Essential Core & Intermediate PairRDD Operations

= easy      = medium

General	Math / Statistical	Set Theory / Relational	Data Structure
<ul style="list-style-type: none"> <li>• flatMapValues</li> <li>• groupByKey</li> <li>• reduceByKey</li> <li>• <b>reduceByKeyLocally</b></li> <li>• foldByKey</li> <li>• aggregateByKey</li> <li>• sortByKey</li> <li>• combineByKey</li> </ul>	<ul style="list-style-type: none"> <li>• sampleByKey</li> </ul>	<ul style="list-style-type: none"> <li>• cogroup (=groupWith)</li> <li>• join</li> <li>• subtractByKey</li> <li>• fullOuterJoin</li> <li>• leftOuterJoin</li> <li>• rightOuterJoin</li> </ul>	<ul style="list-style-type: none"> <li>• partitionBy</li> </ul>

---

<ul style="list-style-type: none"> <li>• keys</li> <li>• values</li> </ul>	<ul style="list-style-type: none"> <li>• countByKey</li> <li>• countByValue</li> <li>• countByValueApprox</li> <li>• countApproxDistinctByKey</li> <li>• countApproxDistinctByKey</li> <li>• countByKeyApprox</li> <li>• sampleByKeyExact</li> </ul>
--	--

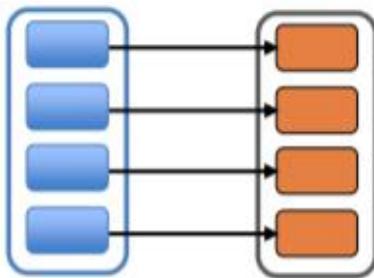


VS



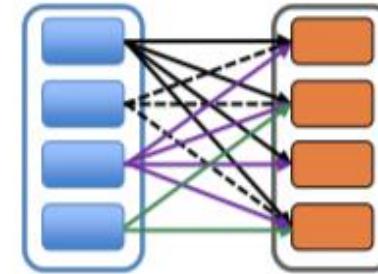
narrow

*each partition of the parent RDD is used by  
at most one partition of the child RDD*



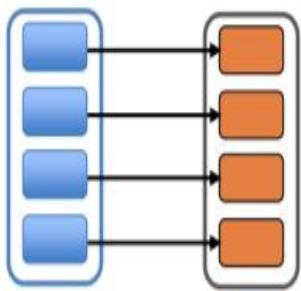
wide

*multiple child RDD partitions may depend  
on a single parent RDD partition*

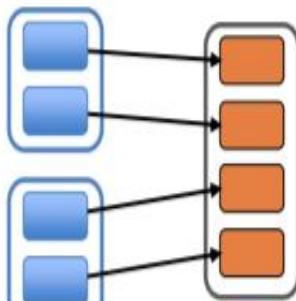


## narrow

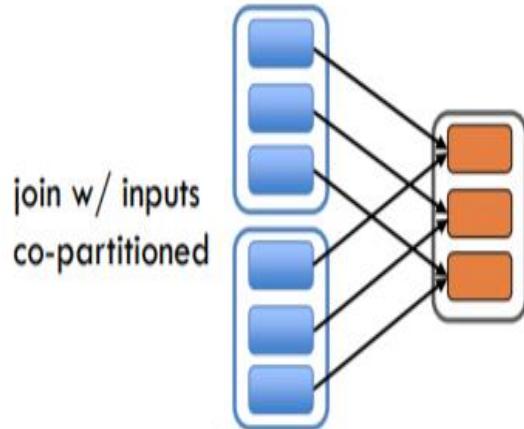
*each partition of the parent RDD is used by at most one partition of the child RDD*



map, filter



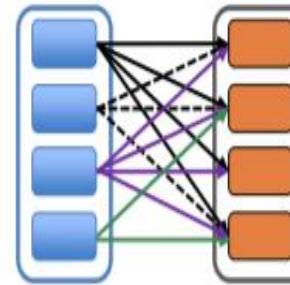
union



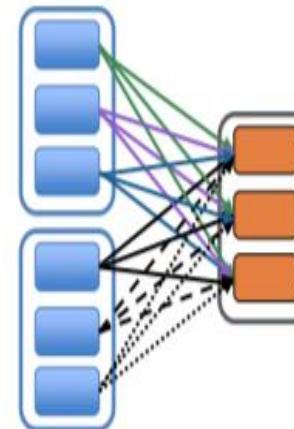
join w/ inputs co-partitioned

## wide

*multiple child RDD partitions may depend on a single parent RDD partition*



groupByKey



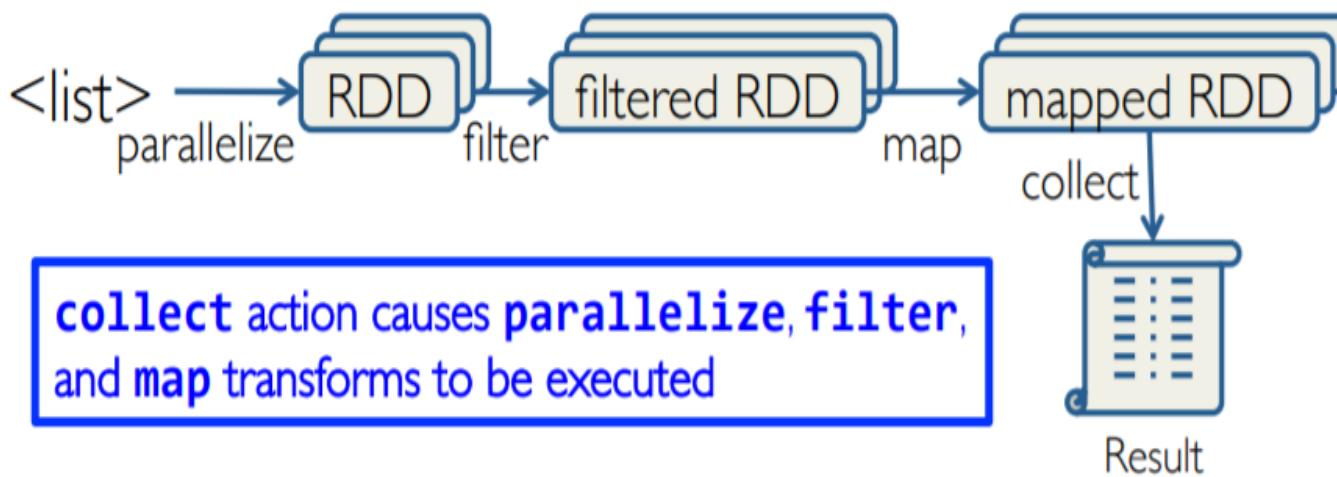
join w/ inputs not co-partitioned

# RDD Dependencies

- Narrow dependencies
  - allow for pipelined execution on one cluster node
  - easy fault recovery
- Wide dependencies
  - require data from all parent partitions to be available and to be shuffled across the nodes
  - a single failed node might cause a complete re-execution.

# Working with RDDs

- Create an RDD from a data source:  <list>
- Apply transformations to an RDD: map filter
- Apply actions to an RDD: collect count



# Creating an RDD

- Create RDDs from Python collections (lists)

```
>>> data = [1, 2, 3, 4, 5]  
  
>>> data  
  
[1, 2, 3, 4, 5]
```

No computation occurs with `sc.parallelize()`

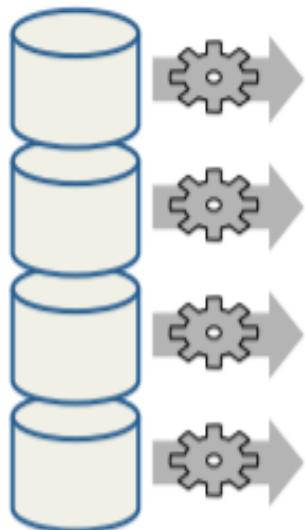
- Spark only records how to create the RDD with four partitions

```
>>> rDD = sc.parallelize(data, 4)  
  
>>> rDD
```

ParallelCollectionRDD[0] at parallelize at PythonRDD.scala:229

# Creating an RDD from a File

```
distFile = sc.textFile("...", 4)
```



- RDD distributed in 4 partitions
- Elements are lines of input
- *Lazy evaluation* means no execution happens now

## Creating RDDs

There are many ways to create RDD objects:

1. From list or arrays defined within the program
2. By reading from normal files
3. Reading from Hadoop HDFS
4. From the output of Hive queries
5. From the output of normal databases queries

**Help() function:** The help() function is used to display the documentation string and also facilitates you to see the help related to modules, keywords, attributes, etc.

**Dir() function:** The dir() function is used to display the defined methods.

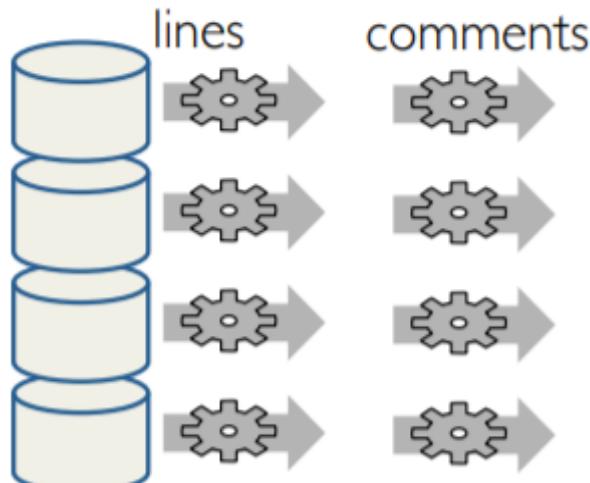
# Spark Transformations

- Create new datasets from an existing one
- Use *lazy evaluation*: results not computed right away – instead Spark remembers set of transformations applied to base dataset
  - » Spark optimizes the required calculations
  - » Spark recovers from failures and slow workers
- Think of this as a recipe for creating result

# Transforming an RDD

```
lines = sc.textFile("...", 4)
```

```
comments = lines.filter(isComment)
```



Lazy evaluation means  
nothing executes –  
Spark saves recipe for  
transforming source

## RDD - TRANSFORMATIONS

No.	RDD Action	Expecting Result
1	<code>map()</code>	Return a new distributed dataset formed by passing each element of the source through a function <i>func</i> .
2	<code>filter()</code>	Return a new dataset formed by selecting those elements of the source on which <i>func</i> returns true.
3	<code>flatMap()</code>	Similar to map, but each input item can be mapped to 0 or more output items (so <i>func</i> should return a Seq rather than a single item).
4	<code>mapPartitions()</code>	Similar to map, but runs separately on each partition (block) of the RDD
5	<code>mapPartitionsWithIndex()</code>	Similar to mapPartitions, but also provides <i>func</i> with an integer value representing the index of the partition
6	<code>sample()</code>	Sample a fraction <i>fraction</i> of the data, with or without replacement, using a given random number generator seed.
7	<code>union()</code>	Return a new dataset that contains the union of the elements in the source dataset and the argument.
8	<code>intersection()</code>	Return a new RDD that contains the intersection of elements in the source dataset and the argument.
9	<code>distinct()</code>	Return a new dataset that contains the distinct elements of the source dataset.
10	<code>groupByKey()</code>	When called on a dataset of (K, V) pairs, returns a dataset of (K, Iterable<V>) pairs

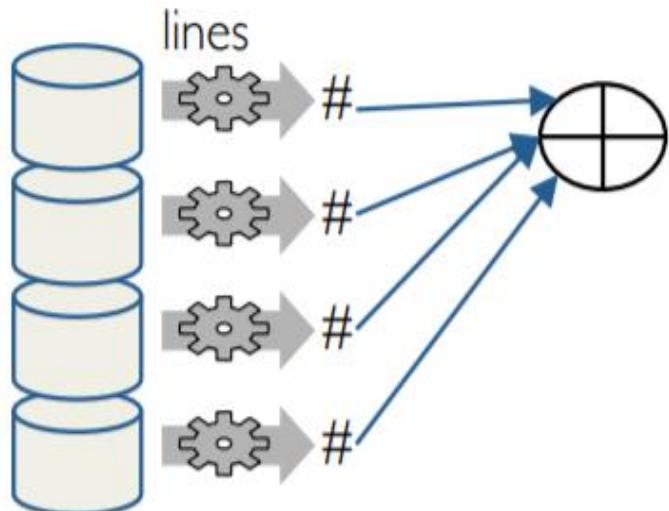
## RDD – TRANSFORMATIONS - JOINS

No	TRANSFORMATION	Expecting Result
1	<code>join()</code>	This takes in 2 Pair RDDs, and returns only matched records from both RDD's
2	<code>leftOuterJoin()</code>	This takes in 2 Pair RDDs, and its returns matched records from both RDD's and unmatched record from left RDD
3	<code>rightOuterJoin()</code>	This takes in 2 Pair RDDs, and its returns matched records from both RDD's and unmatched record from Right RDD
4	<code>fullOuterJoin()</code>	This takes in 2 Pair RDDs, and its returns matched records from both RDD's and unmatched records from both left and right RDD's
5	<code>cartesian()</code>	This takes in 2 pair RDD's data and it will apply cross product ( multiplication) of each key (record) to another RDD.

# Spark Programming Model

```
lines = sc.textFile("...", 4)
```

```
print lines.count()
```

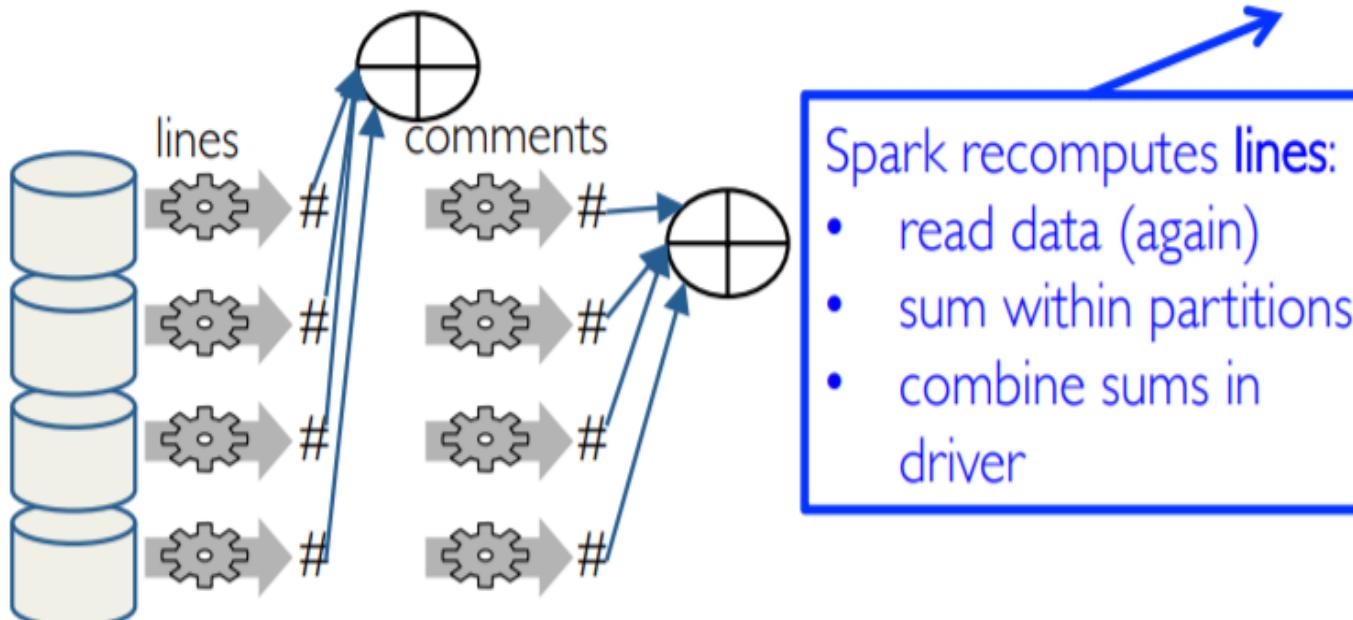


count() causes Spark to:

- read data
- sum within partitions
- combine sums in driver

# Spark Programming Model

```
lines = sc.textFile("...", 4)
comments = lines.filter(isComment)
print lines.count(), comments.count()
```



Spark recomputes `lines`:

- read data (again)
- sum within partitions
- combine sums in driver

## RDD ACTIONS LIST

No.	RDD Action	Expecting Result
1	<code>collect()</code>	Convert RDD to in-memory list
2	<code>take(3)</code>	First 3 elements of RDD
3	<code>top(3)</code>	Top 3 elements of RDD
4	<code>count()</code>	Find total no of values in RDD.
5	<code>min()</code>	Find minimum value from the RDD list
6	<code>max()</code>	Find maximum value from the RDD List
7	<code>sum()</code>	Find element sum (assumes numeric elements)
8	<code>mean()</code>	Find element mean (assumes numeric elements)
9	<code>stdev()</code>	Find element deviation (assumes numeric elements)
10	<code>takeSample(withReplacement=True, 3)</code>	Create sample of 3 elements with replacement

## RDD ACTIONS - LIST 2

No.	RDD Action	Expecting Result
11	<code>reduce()</code>	Reduce is a spark action that aggregates a data set (RDD) element using a function.
12	<code>countByKey()</code>	Count the number of elements for each key, and return the result to the master as a dictionary.
13	<code>CountByValue()</code>	Return the count of each unique value in this RDD as a dictionary of (value, count) pairs.
14	<code>fold()</code>	Aggregate the elements of each partition
15	<code>range()</code>	Create a new RDD of int containing elements from <i>start</i> to <i>end</i> (exclusive)
16	<code>variance()</code>	Compute the variance of this RDD's elements.
17	<code>sampleVariance()</code>	Compute the sample variance of this RDD's elements (which corrects for bias in estimating the variance by dividing by N-1 instead of N).
18	<code>saveAsTextFile()</code>	Save this RDD as a text file, using string representations of elements.
19	<code>saveAsPickleFile()</code>	Save this RDD as a SequenceFile of serialized objects
20	<code>Stats()</code>	Stats will give complete information count, min, max, stdev and mean

# Spark Program Lifecycle

1. Create RDDs from external data or parallelize a collection in your driver program
2. Lazily transform them into new RDDs
3. **cache()** some RDDs for reuse
4. Perform actions to execute parallel computation and produce results

# Review: Python `lambda` Functions

- Small anonymous functions (not bound to a name)  
`lambda a, b: a + b`
  - » returns the sum of its two arguments
- Can use lambda functions wherever function objects are required
- Restricted to a single expression

# Spark Caching & Persistence

**Spark caching** can be used to pull data sets into a cluster-wide in-memory cache. This is very useful for accessing repeated data, such as querying a small “hot” dataset or when running an iterative algorithm

There are two ways to persist RDDs in Spark:

1. cache()
2. persist()

There are some advantages of RDD caching and persistence mechanism in spark.

- Time efficient
- Cost efficient
- Lesser the execution time.

## STORAGE TYPES:

- 1) MEMORY\_ONLY
- 2) MEMORY\_AND\_DISK
- 3) DISK\_ONLY
- 4) MEMORY\_ONLY\_SER
- 5) MEMORY\_AND\_DISK\_SER

## RDD Persist different storage level

**MEMORY\_ONLY** Store RDD as deserialized Java objects in the JVM. If the RDD does not fit in memory, some partitions will not be cached and will be recomputed on the fly each time they're needed. This is the default level.

**MEMORY\_AND\_DISK** Store RDD as deserialized Java objects in the JVM. If the RDD does not fit

in memory, store the partitions that don't fit on disk, and read them from there when they're needed.

**MEMORY\_ONLY\_SER** (Java and Scala) Store RDD as serialized Java objects (one byte array per partition). This is generally more space-efficient than deserialized objects, especially when using a fast serializer, but more CPU-intensive to read.

**MEMORY\_AND\_DISK\_SER** (Java and Scala) Similar to MEMORY\_ONLY\_SER, but spill partitions that don't fit in memory to disk instead of recomputing them on the fly each time they're needed.

**DISK\_ONLY** Store the RDD partitions only on disk.

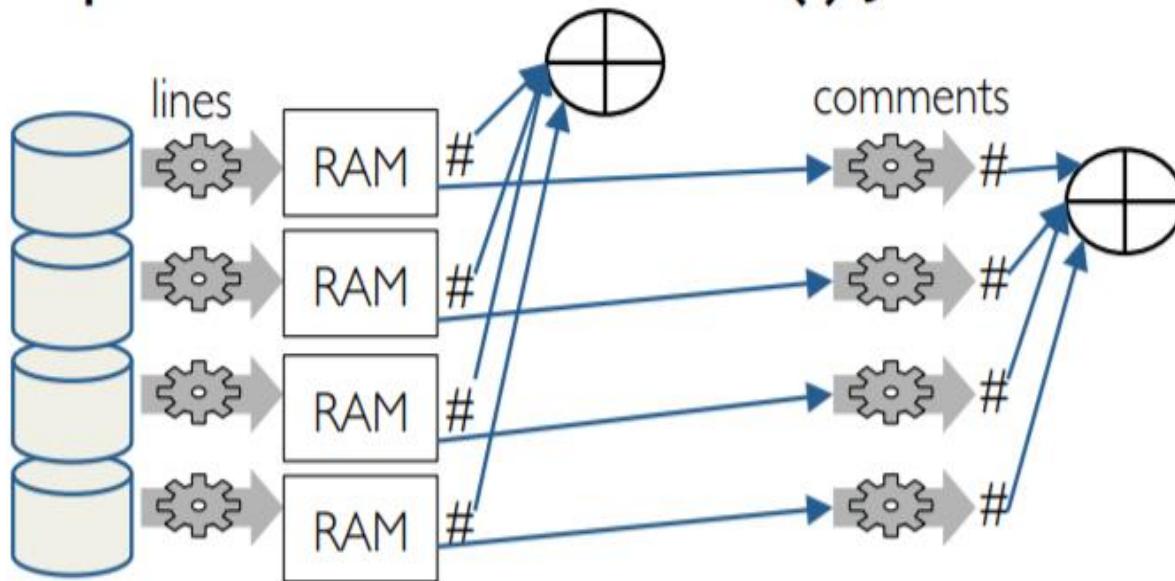
MEMORY\_ONLY\_2, MEMORY\_AND\_DISK\_2, etc. Same as the levels above, but replicate each partition on two cluster nodes.

**OFF\_HEAP** (experimental) Similar to MEMORY\_ONLY\_SER, but store the data in off-heap memory.

This requires off-heap memory to be enabled.

# Caching RDDs

```
lines = sc.textFile("...", 4)
lines.cache() # save, don't recompute!
comments = lines.filter(isComment)
print lines.count(), comments.count()
```



# Clear Spark Cache

## RDD Unpersist

Spark automatically monitors cache usage on each node and drops out old data partitions in a least-recently-used (**LRU**) fashion.

If you would like to manually remove an RDD instead of waiting for it to fall out of the cache,

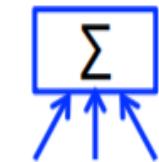
use the **RDD.unpersist()** method.



# pySpark Shared Variables

## Broadcast Variables

- » Efficiently send large, *read-only* value to all workers
- » Saved at workers for use in one or more Spark operations
- » Like sending a large, read-only lookup table to all the nodes



## Accumulators

- » Aggregate values from workers back to driver
- » Only driver can access value of accumulator
- » For tasks, accumulators are write-only
- » Use to count errors seen in RDD across workers

# Broadcast Variables

**Broadcast variables** allow the programmer to keep a read-only variable cached on each machine rather than shipping a copy of it with tasks.

They can be used, for example, to give every node a copy of a large input dataset in an efficient manner. Spark also attempts to distribute broadcast variables using efficient broadcast algorithms to reduce communication cost.

## Syntax:

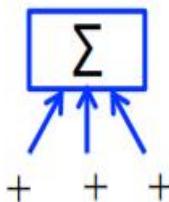
```
#Create Broadcast Variable using SparkContext.broadcast (LIST)
broadcast_v = sc.broadcast([1, 2, 3])
# Print the broadcast variable value using .value
broadcast_v.value
```

- **Broadcast variables** are created from a variable v by calling **SparkContext.broadcast(v)**.
- The broadcast variable is a wrapper around v, and its value can be accessed by calling the **.value** method.

`broadcast_v.unpersist()`

`unpersist(self, blocking=False)`

Delete cached copies of this broadcast on the executors. If the broadcast is used after this is called, it will need to be re-sent to each executor.

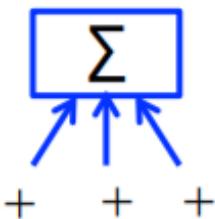


## Accumulators

- Variables that can only be “added” to by associative op
- Used to efficiently implement parallel counters and sums
- Only driver can read an accumulator’s value, not tasks

```
>>> accum = sc.accumulator(0)
>>> rdd = sc.parallelize([1, 2, 3, 4])
>>> def f(x):
...     global accum
...     accum += x

>>> rdd.foreach(f)
>>> accum.value
Value: 10
```



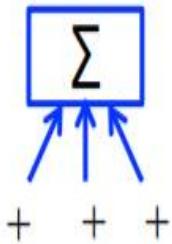
# Accumulators Example

- Counting empty lines

```
file = sc.textFile(inputFile)
# Create Accumulator[Int] initialized to 0
blankLines = sc.accumulator(0)

def extractCallSigns(line):
    global blankLines # Make the global variable accessible
    if (line == ""):
        blankLines += 1
    return line.split(" ")

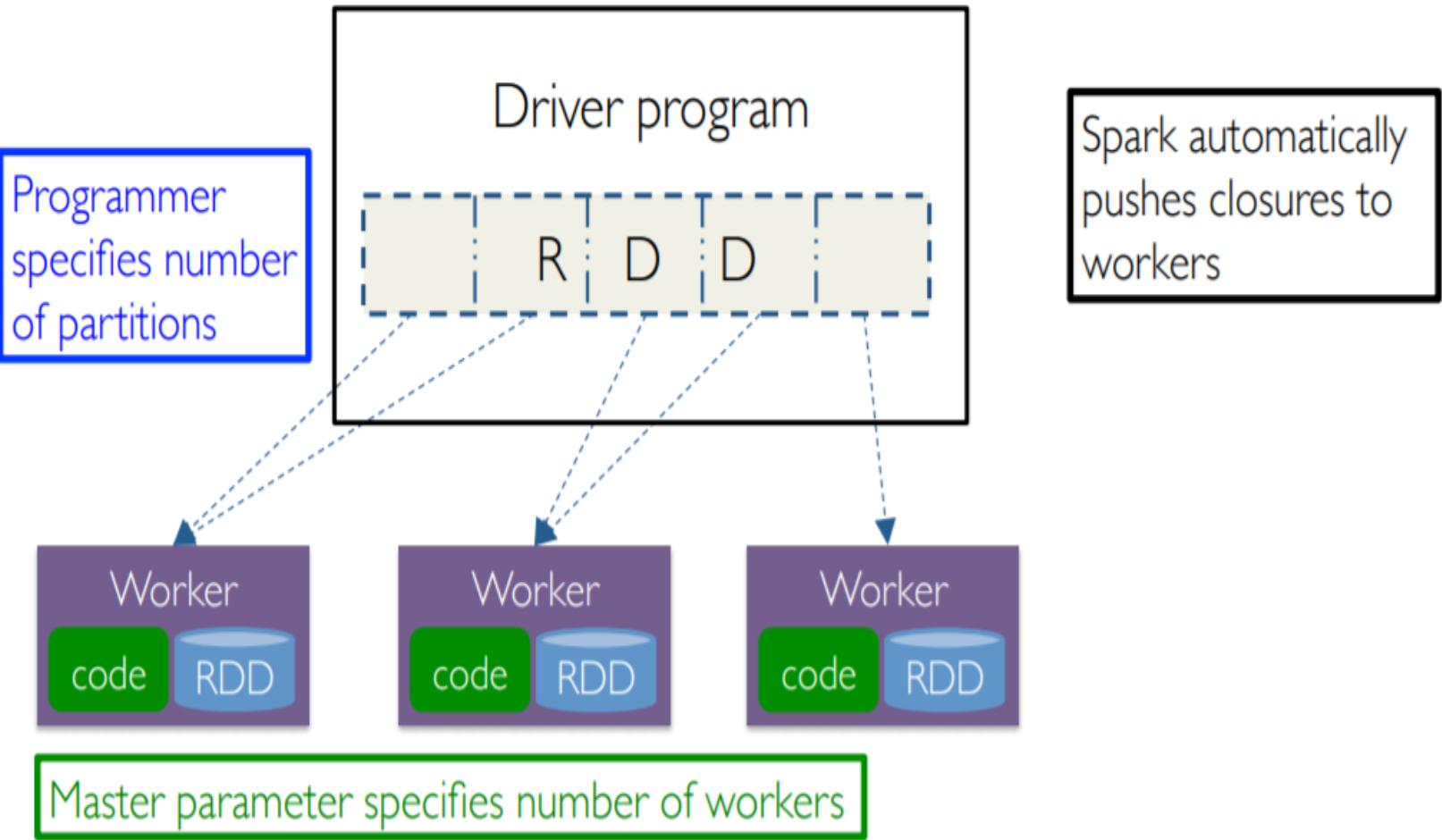
callSigns = file.flatMap(extractCallSigns)
print "Blank lines: %d" % blankLines.value
```



# Accumulators

- Tasks at workers cannot access accumulator's values
- Tasks see accumulators as write-only variables
- Accumulators can be used in actions or transformations:
  - » Actions: each task's update to accumulator is *applied only once*
  - » Transformations: *no guarantees* (use only for debugging)
- Types: integers, double, long, float
  - » See lab for example of custom type

# Summary



# THANK YOU