



# Spark

1. Most important topic
2. Restart
3. Interview ✓  
+  
Work

~1 month work + 1 year

(Storage) HDFS → distributed storage, though cloud based gaining popularity

(resource manager) YARN → resource manager (Docker, Kubernetes, Mesos).

Challenging MR → waiting & performance → Spark

Some Common Questions.

Q:- Do I need to know hadoop before spark? → No

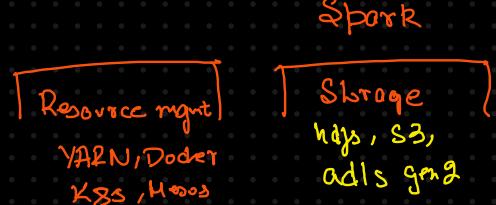
Q:- Do I need to install hadoop first? → No, standalone

Q: Prerequisites before learning spark?

→ Python, Java, Scala

→ SQL

→ hadoop essentials ⇒ dist. storage & processing



Q: Can you run spark on your laptop? → Yes

Q: Programming language?

Q: How is spark diff. from other databases? ⇒ spark is not a DB

Q: What storage can you use with spark? ⇒ hdfs, s3, adls gen2, Kafka

Q: Is spark free? ⇒ Yes

Android → One UI

Q: What is diff b/w spark & hbase?

Q: Can you use spark on cloud?

## MR and its limitation

1. Performance → high latency  
disk i/o
2. Complex & hard to write ∵ hard to write boilerplate code  
versatile
3. Only batch processing is supported ∵ not stream processing
4. Rigid Sequence :- Everything has to be thought of in terms of MR
5. No interactive mode/way to monitor jobs or run in ports?



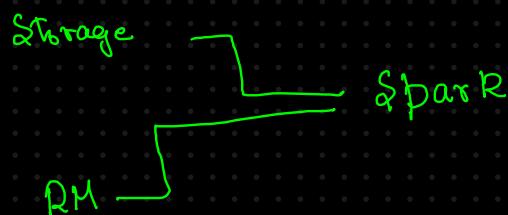
# Spark

Apache Spark is an open-source distributed computing system

Spark 3

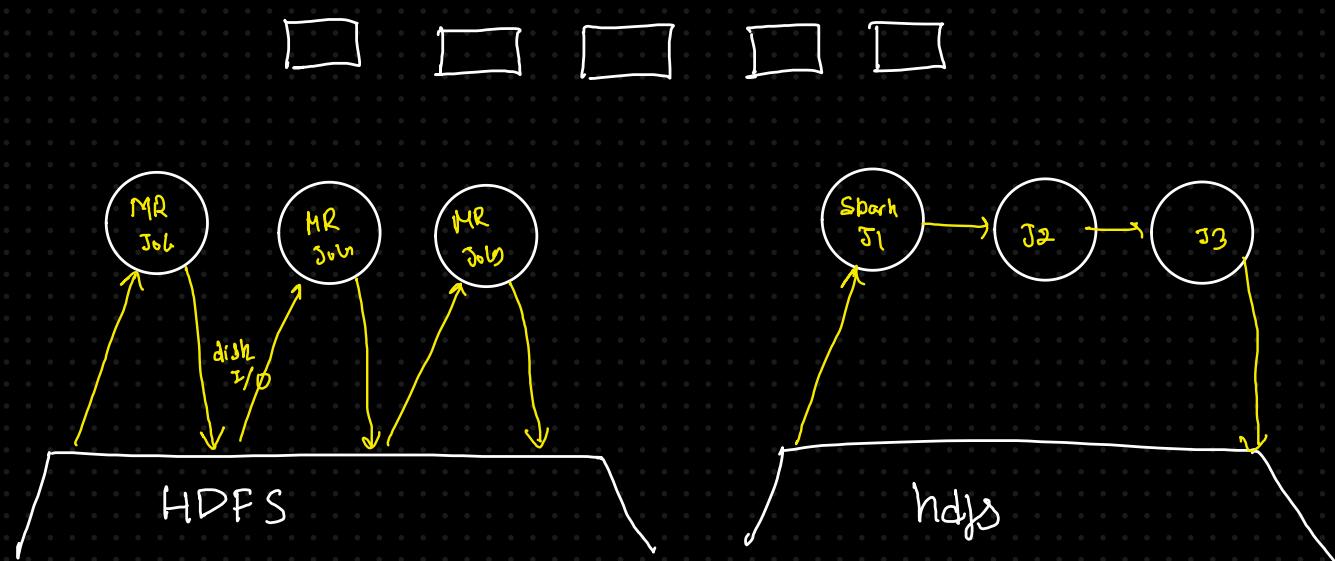
designed to process large datasets quickly.

fast, general purpose & support batch & real time processing



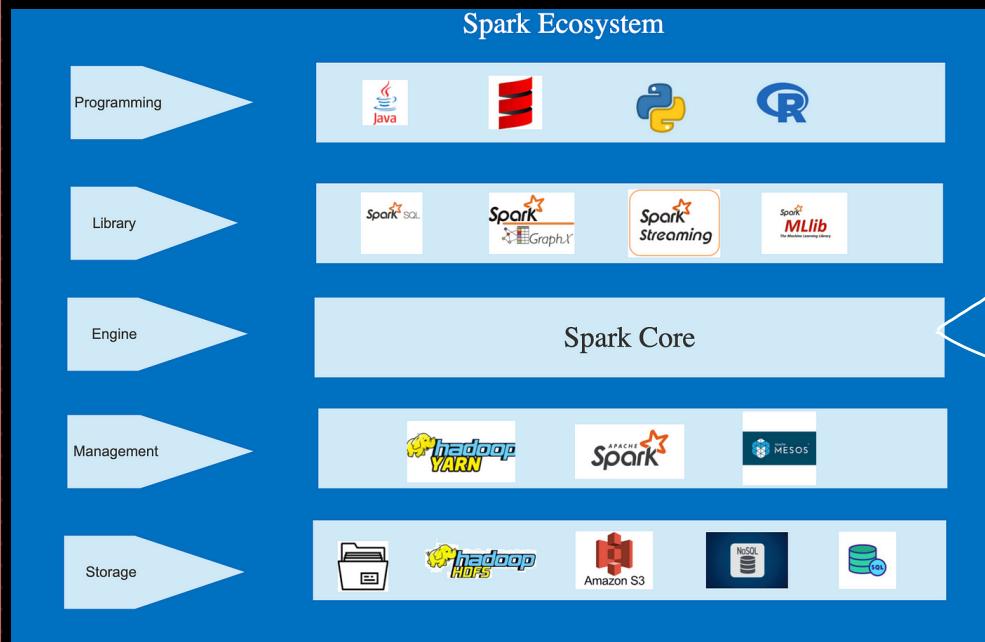
## Characteristic

1. In-memory Processing : in-memory processing faster than traditional disk-based system



2. Ease of use → Lang. support + life

3. Unified Framework →



developed to make life  
easy (clean & flexible)

↑  
Spark higher level

↓  
Spark core / low level  
RDD (most  
flexible)

## Features of Apache spark

1. Speed : in memory processing, DAG,

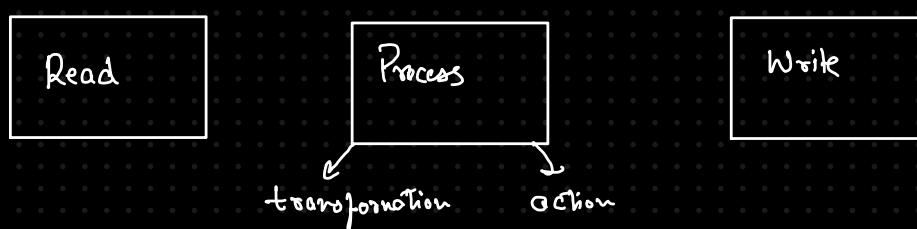
2. Scalable : one m/c → 1000 m/c

3. Fault tolerant : automatically recovers the data in case of failure

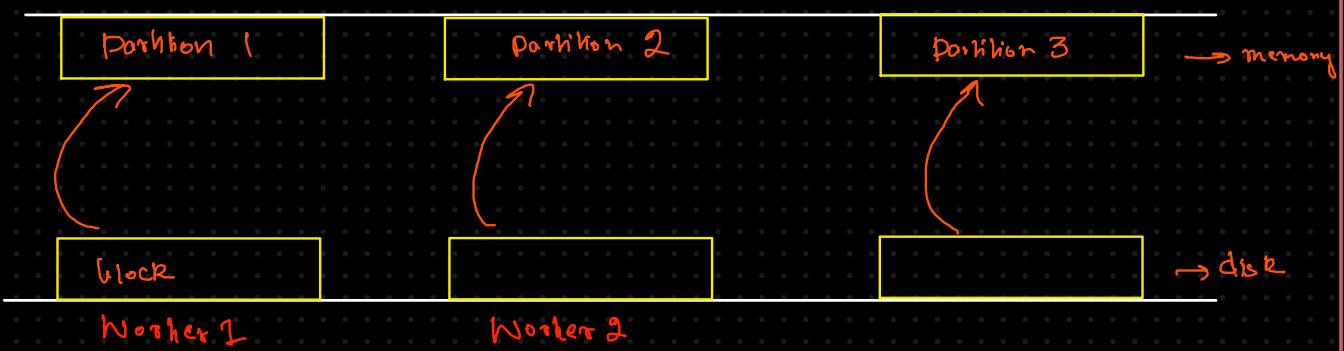
4. Multi-language / Polyglot: multi lang support

5. Unified framework

6. Integration with existing systems : hdfs , yarn, Kafka



## RDD & execution plan in Spark



data loaded from disk → memory

300mb ⇒ 3 blocks

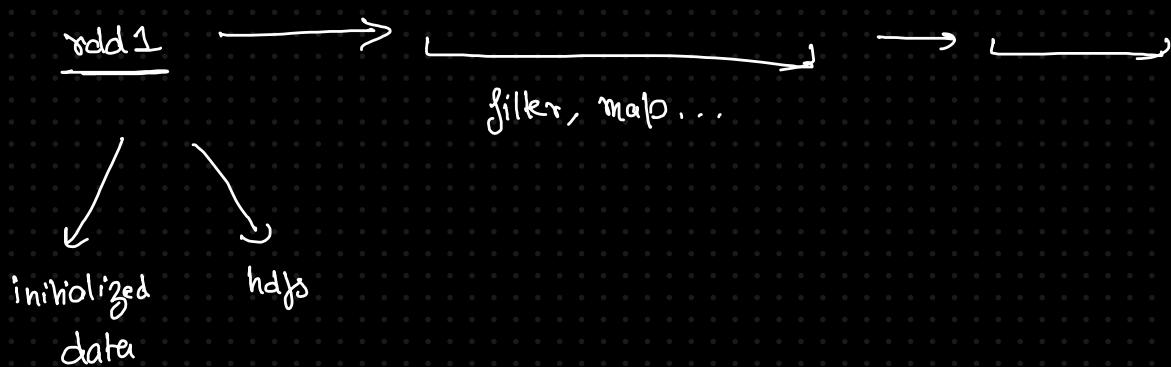
block on disk  
&  
position on memory

hdfs file → read in spark



processing → send it back to disk

add from some local data



Nothing happened in our code unless we ran collect()

Lazy evaluation

### Transformation

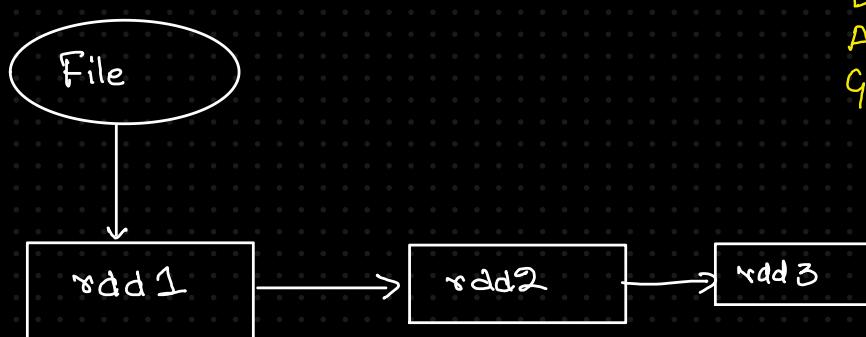
Op<sup>n</sup> that create a new RDD from existing one.

Lazy

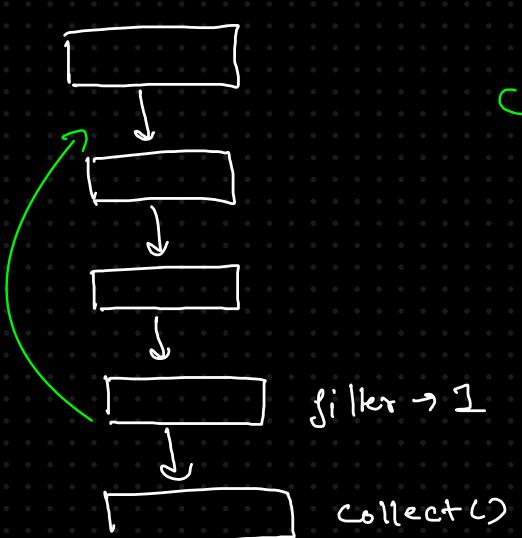
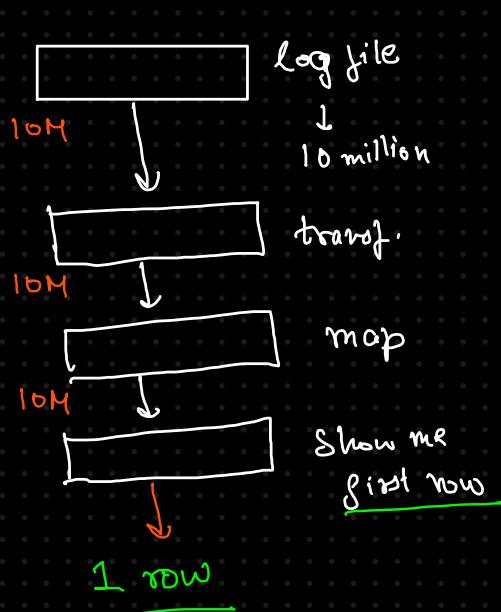
### Action

Op<sup>n</sup> that triggers the execution of all transformations & return result

Eager



Q: Why are transformation is lazy



## Transformations

map	join	union	distinct	repartition
mapPartitions	flatMap	intersection	pipe	coalesce
cartesian	cogroup	filter	sample	
sortByKey	groupByKey	reduceByKey	aggregateByKey	
mapPartitionsWithIndex			repartitionAndSortWithinPartitions	

## Actions

reduce	take	collect	takeSample	count
takeOrdered	countByKey	first	foreach	saveAsTextFile
saveAsSequenceFile		saveAsObjectFile		

R - resilient →

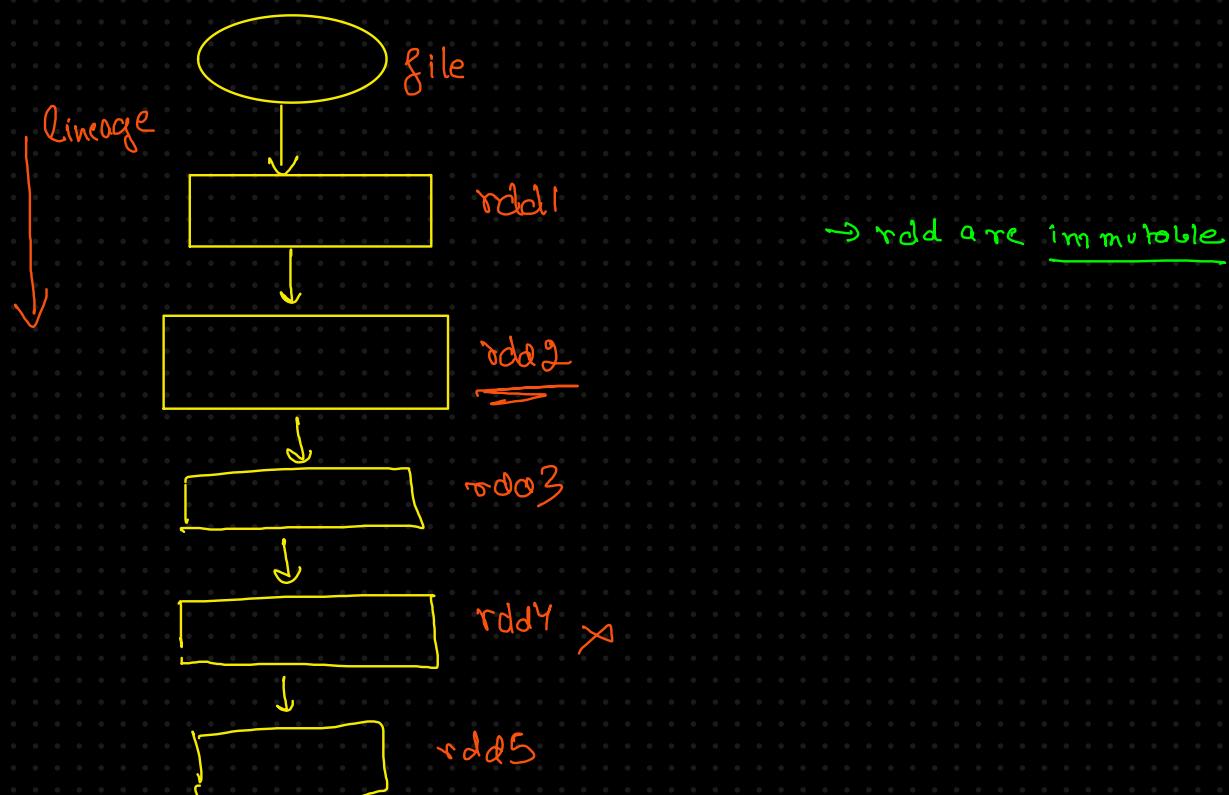
D - distributed ✓

D - dataset ✓

rdd



Resilient will mean that we can quickly recover

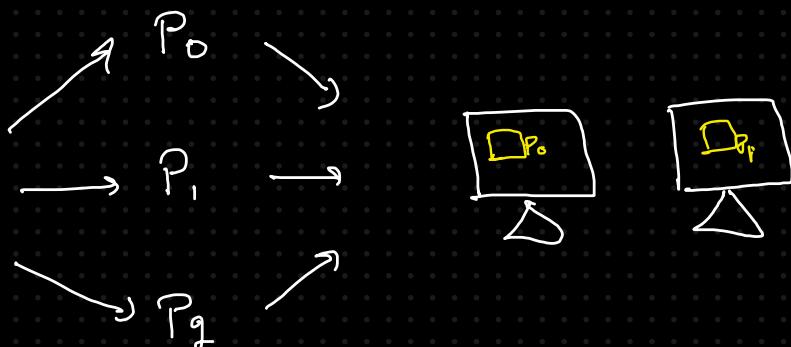
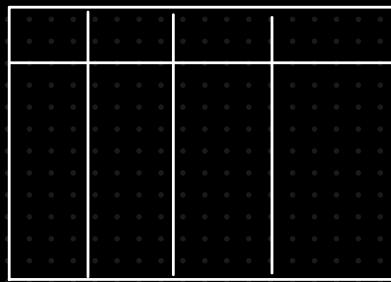


## Properties of rdd

1. Immutable
2. Lazy evaluated
3. Partitioned

## Characteristic

1. Fundamental data structure
2. Resilient
3. Fault Tolerance
4. Lazily evaluated



### 1. SparkContext (sc):

- It's the older entry point for Spark functionality
- Primarily works with RDDs (Resilient Distributed Datasets)
- Each JVM can only have one active SparkContext
- More low-level control over Spark configurations

### 2. SparkSession (spark):

- Introduced in Spark 2.0 as the unified entry point
- Encapsulates SparkContext, SQLContext, and HiveContext
- Provides access to DataFrame and Dataset APIs
- Can have multiple SparkSessions in the same application
- Recommended for modern Spark applications

### 3. "Normal" Spark:

- This usually refers to running Spark without specialized contexts
- Basic functionality without SQL or Hive support
- Limited compared to using SparkSession

Spark 2.0

Python & Spark → PySpark

1. GCP DataProc
2. Terminal of DataProc
3. Google Colab
4. Databricks

Python map reduce higher orders fn

# Reading & Partitioning Data in Spark

Data is read and divided into partition

What is partitioning in Spark?

> Logical chunks of data

Processing

large dataset is divided into smaller chunks  
called partition.

1. When already divided  
into blocks → hdfs, S3  
lake gen 2

2. read data & create partition  
based on the config.

Execution of one partition of data = one task  
= 1 CPU core required for

Execution  $\approx$  Mapper

(executor)  $\approx$  Reducer

performs a task

& it can have  
more than  
1 core as  
well

8 blocks  $\rightarrow$  8 partition  $\rightarrow$  8 tasks

( $\lg b$ )

inc.

dec

i will  
get more illism

less parallelism

## Spark parallelize

We can create an RDD using  
parallelize & then work on it

Subset of data



< 128 mb

$\Downarrow$   
2 partition  $\Rightarrow$  2 cores



8 cores

+ partition  $\Rightarrow$  better results

→ Read the data → Parallelize & partition, counting value

→ Transformations → map, reduce, filter, reduceByKey

Reduce By Key

transformation

Count key value

Action

Notebook → functions in spark

reduceByKey 5 keys  $\rightarrow$  5 values

map 100 rows  $\rightarrow$  100

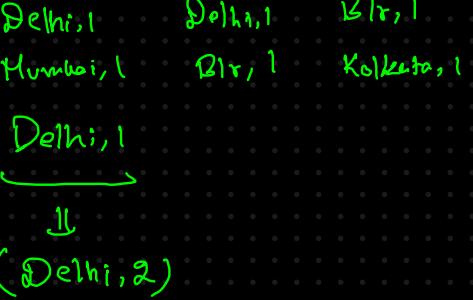
filter 100 rows  $\rightarrow$  0 - 100

This is still the most diff. way of writing code in spark  
 $\Rightarrow$  more easy way

• Save as text file (hdfs)

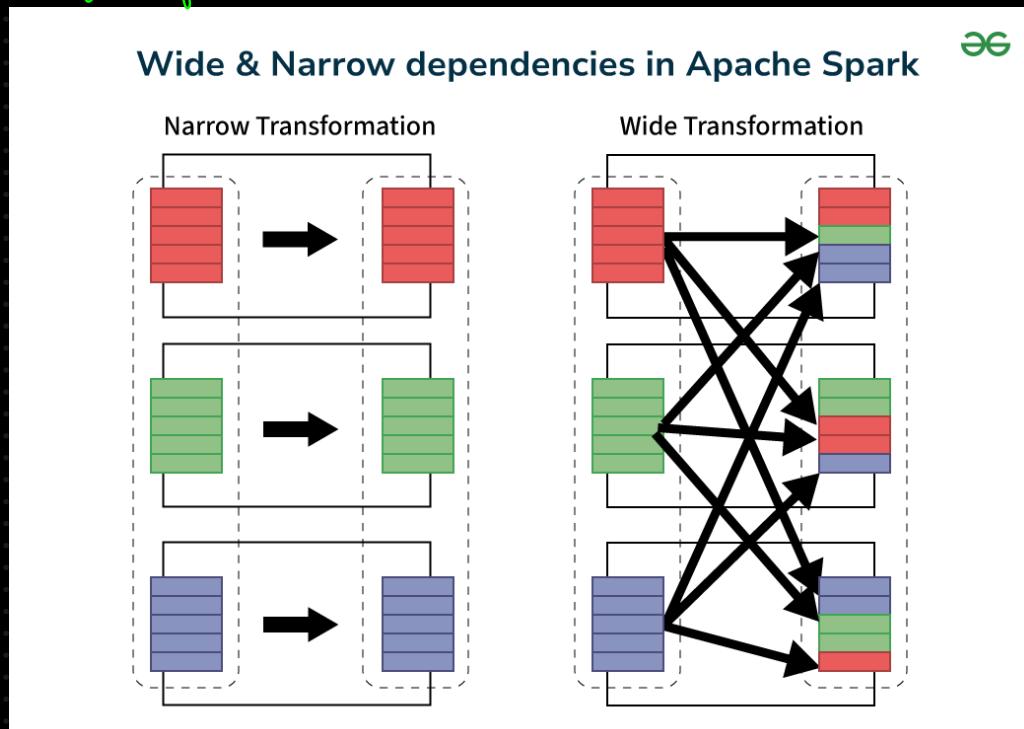
# Transformation in Spark : Narrow vs wide

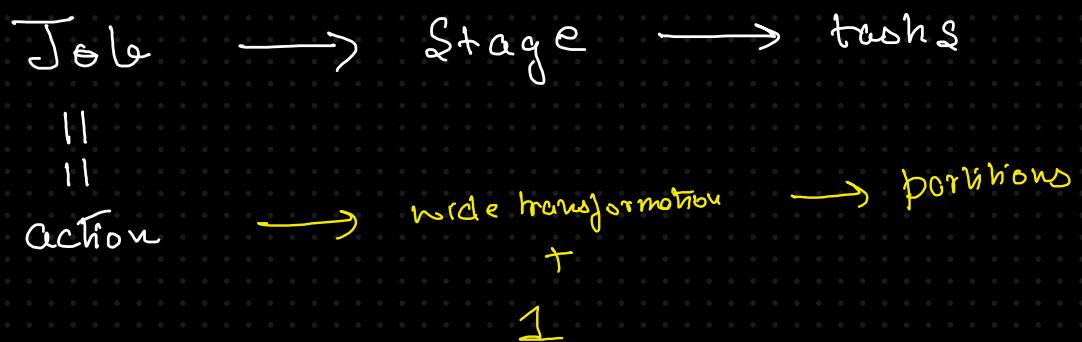
transformations are operation that create a new rdd from existing one

Narrow	Wide	Broad
data is not shuffled across partition	Data is shuffled across our partitions	
faster → no shuffle	Slower - as shuffling ⇒ I/O network	 Delhi,1      Delhi,1      Blr,1 Mumbai,1      Blr,1      Kolkata,1 Delhi,1 ↓ (Delhi,2)
map, filter, flatmap dependent on one parent position	map groupByKey reduceByKey SortByKey dependent on multiple position	

either

- try to have very less wide transformation
- narrow down the data so that data getting shuffled is less.





When we do a wide transformation  $\Rightarrow$  a stage is created



→ homework

→ just run what I did & understand spark UI

→ be comfortable in pyspark

from next class, lots of code without major error.

# Reduce By Key & Group by Key

## Reduce by Key

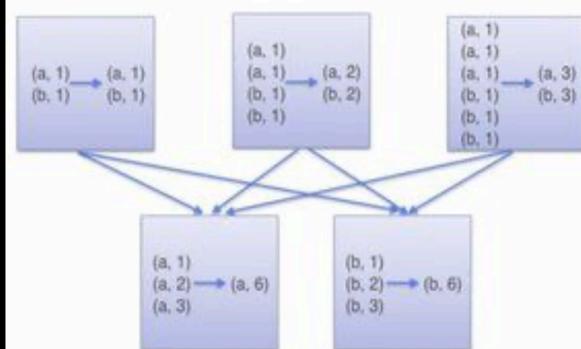
Definition part

Aggregating the value of keys during shuffle phase

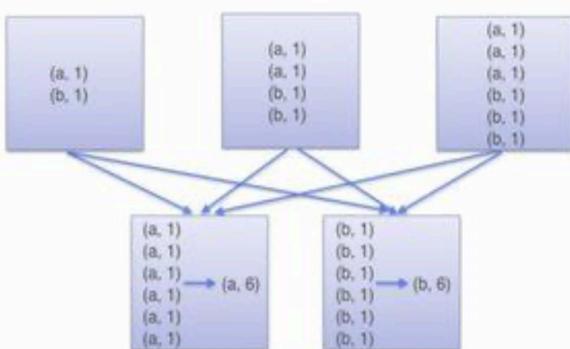
## Group by Key

Group all the key values and we can aggregate on top of them

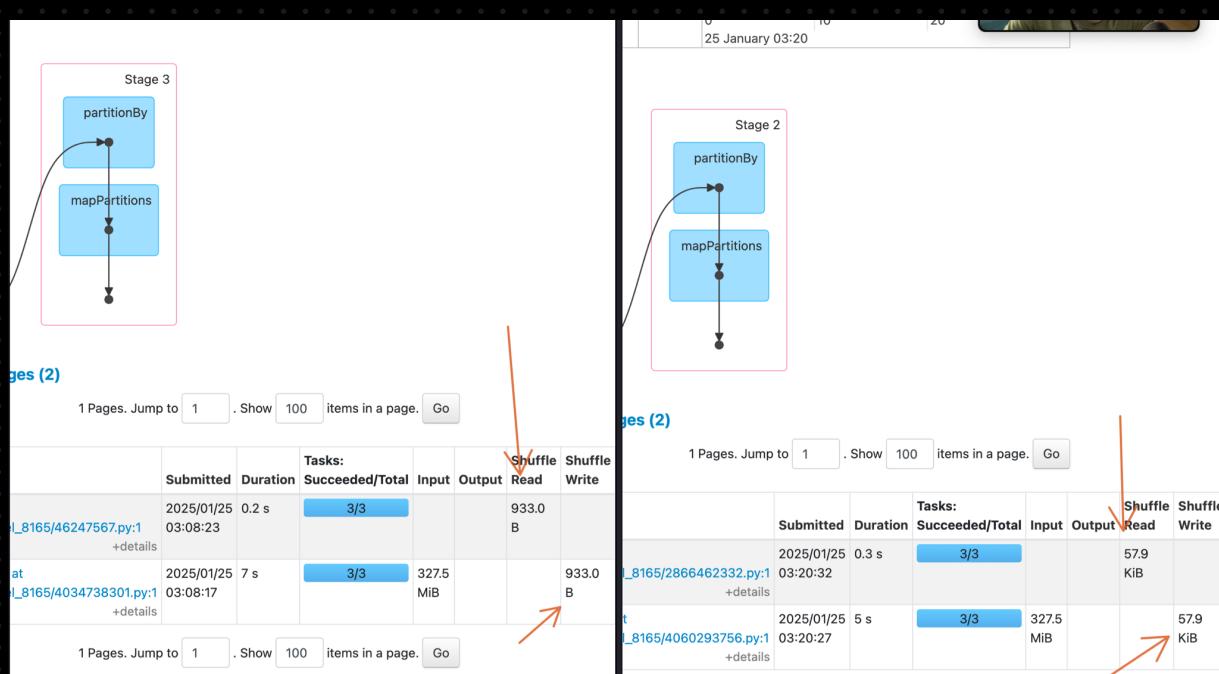
### ReduceByKey



### GroupByKey



For say median, reduceByKey will give wrong answer  
and



## Increasing and Decreasing No. of partition

↳ Repartition vs coalesce

Why would we need to increase or decrease no. of partition?

### Increase

1gb data → 8

but say we have  
12 nodes

We can increase the  
partitions here to  
use all nodes & get  
better I/O

Cluster should be able to  
handle

### Decrease

Say we have 1tb file →

8192

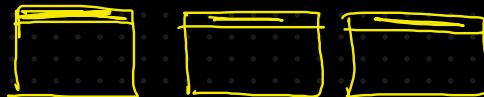
& Say we have  
50 nodes cluster

⇒ decrease its better



8192 partitions

↓ filters



Say 500 partition of 1mb  
each

spare partition

### Repartition

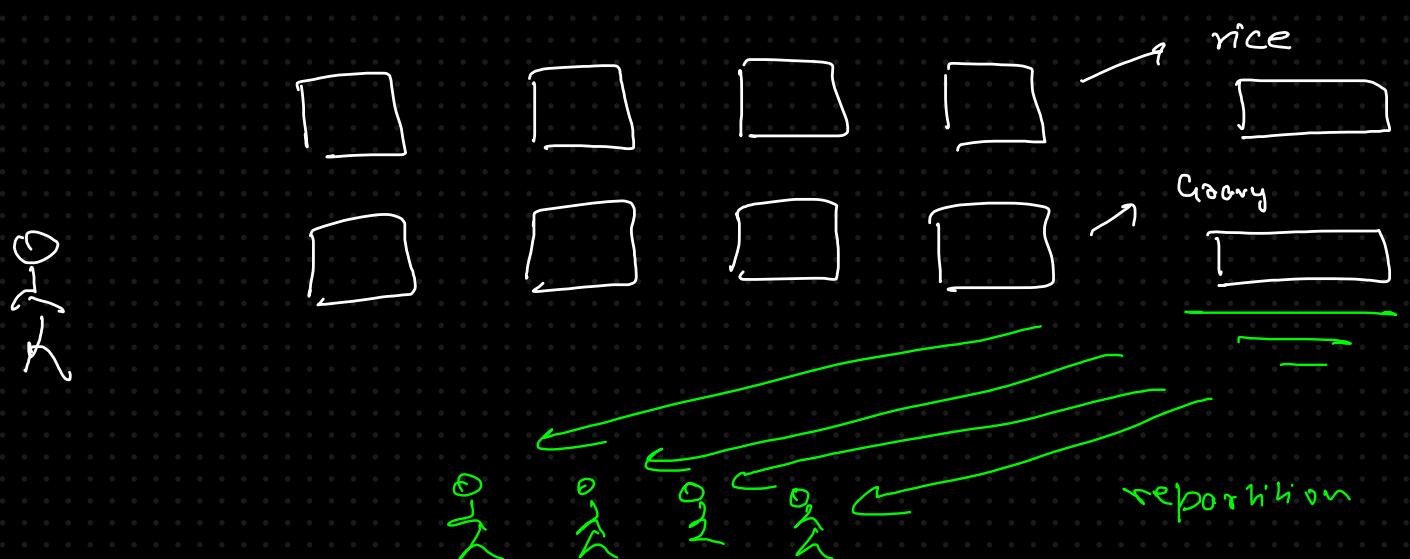
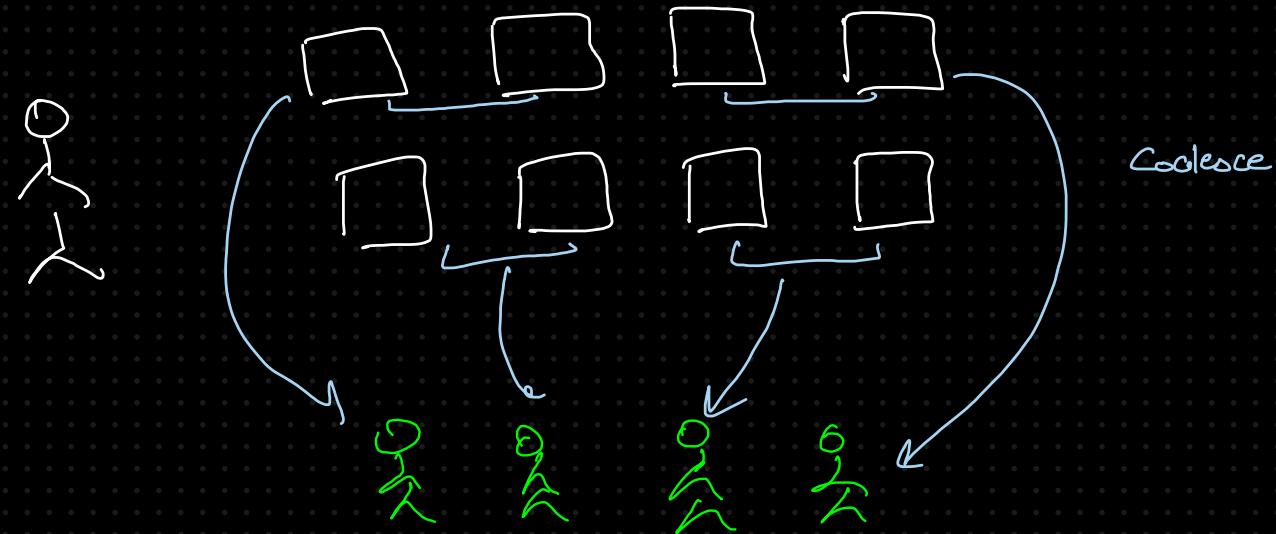
Can be used to increase or  
decrease no. of partition

It does a full shuffle &  
then repartitions the  
data

### Coalesce

decrease the number  
of partitions.

Tries to avoid shuffling  
while reducing no.  
of partition.



Aspect	Repartition	Coalesce
<b>Definition</b>	Shuffles data across nodes to <b>increase or decrease</b> partitions.	Combines partitions <b>without shuffle</b> , typically reducing their number.
<b>Use Case</b>	To balance load or increase parallelism.	To optimize performance by reducing partitions.
<b>Shuffling</b>	Yes, full shuffle across the cluster.	No shuffle; only adjacent partitions are merged.
<b>Performance</b>	Slower due to shuffling overhead.	Faster as it avoids shuffling.
<b>Scalability</b>	Used to increase partitions or even redistribute them.	Used only for decreasing partitions.
<b>When to Use</b>	After reading unbalanced data (e.g., uneven partitions).	When shrinking partitions after a filter operation.



## Spark Higher Level APIs

RDDs not used majorly.

much more

than required

RDDs are core abstraction in Spark

but they lack schema & metadata

1. Verbose

2. Error-prone

3. Performance-based :

Very imp for basics &  
interview prep.

Not persistent

Schema & metadata

2. What are higher level APIs

Dataframes & Spark SQL

no dataset in  
PySpark

& provide schema aware abstraction  
for data processing

Schema awareness → Optimize Query Execution

SQL like programming

Integrate with other data sources (CSV, Parquet, etc)

## Dataframe

$df$  is a distributed collection of data organized into named columns, similar to a table in a dB

$\text{rdd}$  with some structure  
( schema & metadata)  
data

not persistent

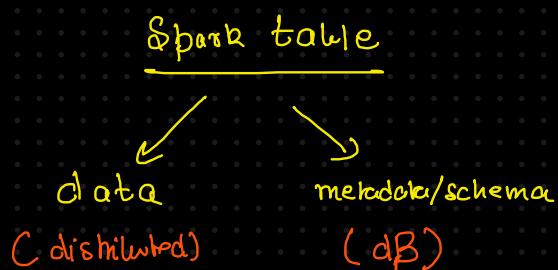
Feature	RDDs	DataFrames
Schema Support	No schema	Schema-aware
Ease of Use	Low	High
Performance	Less optimized	Catalyst-optimized

## Spark SQL

module for structured data processing.

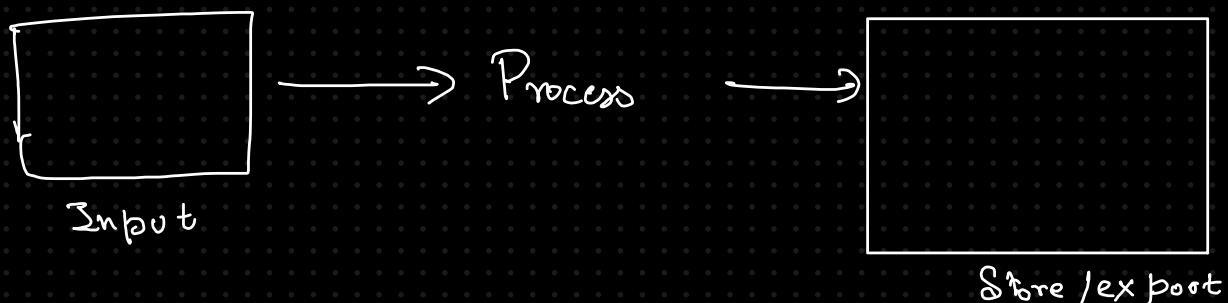
provides you an interface for running SQL Queries on  $df$

& table .

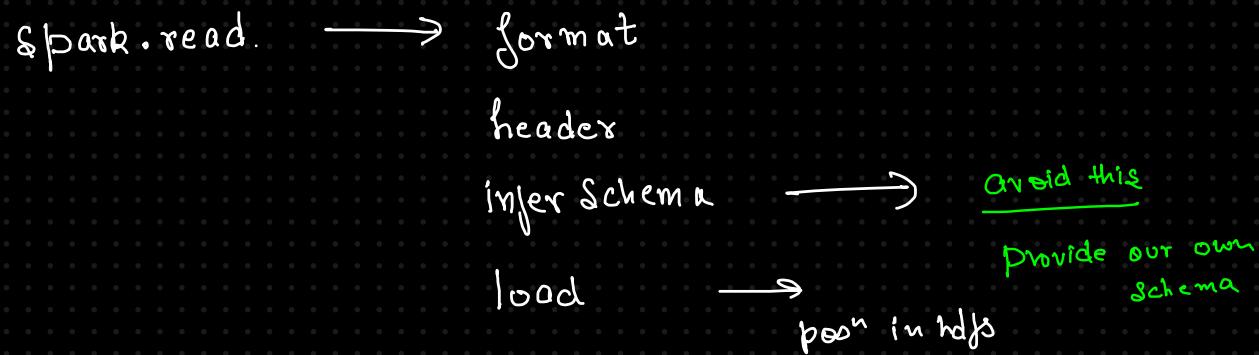


## SQL Queries

Spark  $\rightarrow$  process



## Spark DataFrame



`spark.read`  
action or transformation?

Scenario	Behavior	Explanation
Without <code>inferSchema</code> , with <code>header=True</code>	Eager Evaluation	Reads only the first line to determine column names. Triggers a lightweight job ( <code>collect</code> with <code>limit 1</code> ).
With <code>inferSchema=True</code>	Eager Evaluation	Scans the entire dataset to infer schema. Triggers two jobs (1 for column names, 1 for schema inference).
With Explicit Schema Definition	Lazy Evaluation	No upfront job is triggered. Schema validation occurs when an action is performed on the DataFrame.

# Schema enforcement in Spark df

Challenges with infer schema

Incorrect inference

Performance issues

Bad for Production

1. takes time ←

2. wrong

1. headers → -c1, -c2, ...

2. infer schema → ↗

## Schema Enforcement Techniques

1. Struct Type

Avoids incorrect inference

consistent schema

stricter checks → null

2. DDL string

