



# 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

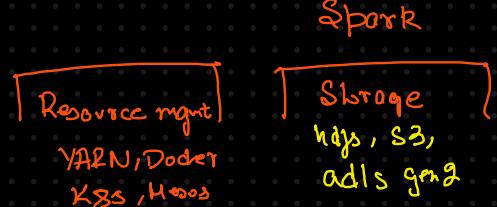
Q:- Do I need to install hadoop first? → No, standalone

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?



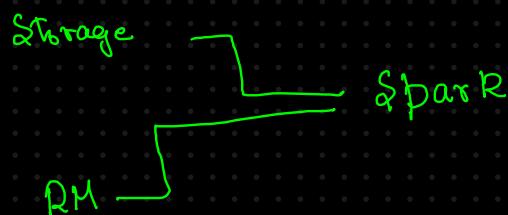
# Spark

Apache Spark is an open-source distributed computing system

Spark 3

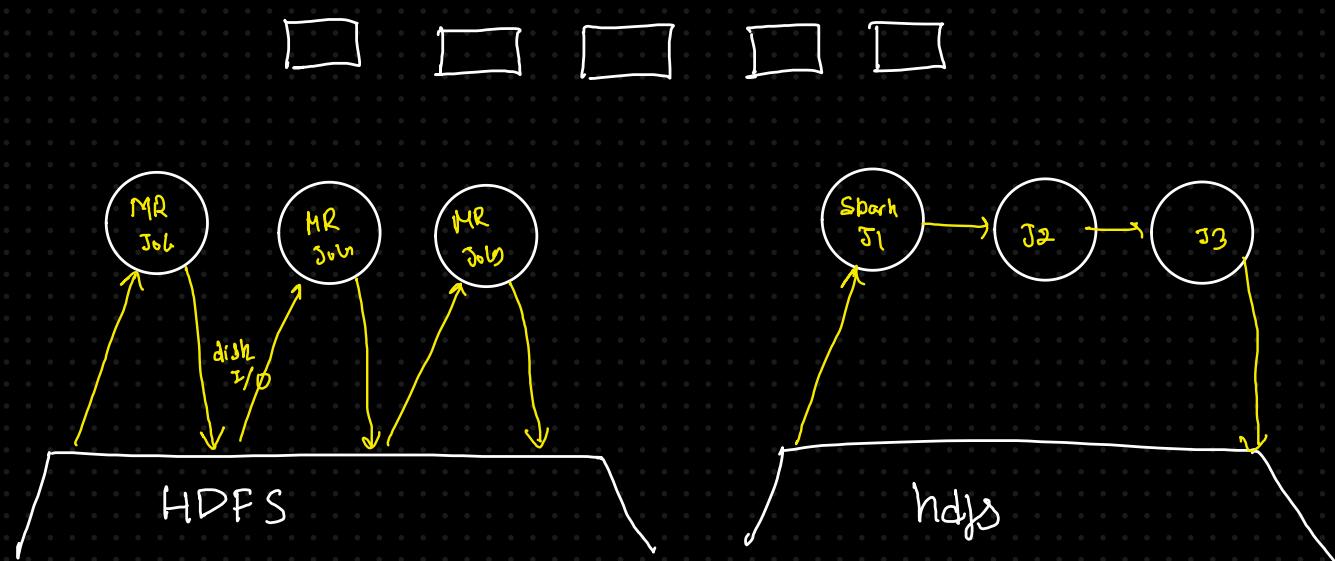
designed to process large datasets quickly.

fast, general purpose & support batch & real time processing



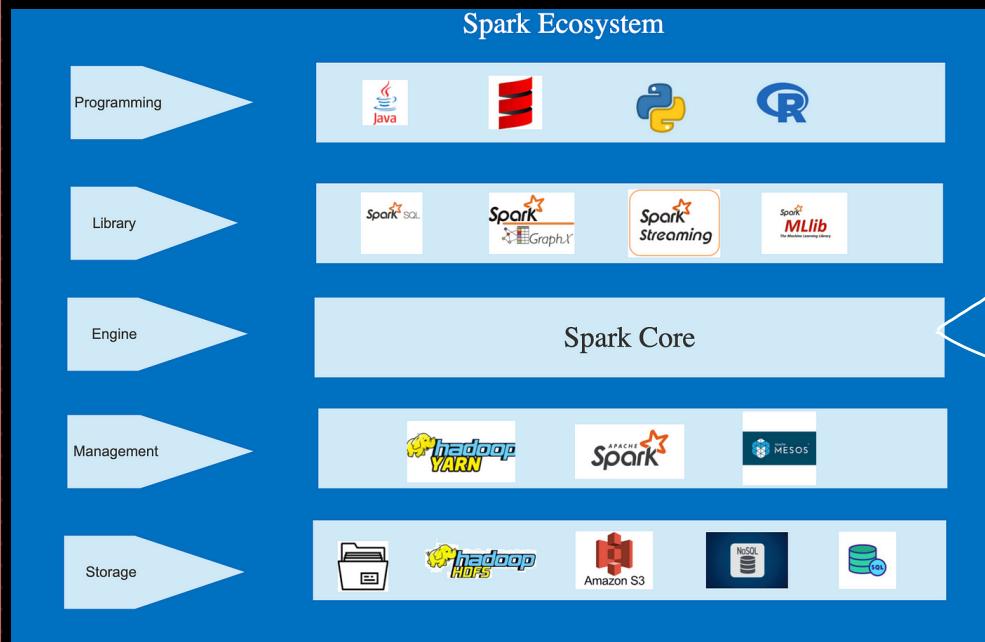
## Characteristic

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



2. Ease of use → Lang. support + life

3. Unified Framework →

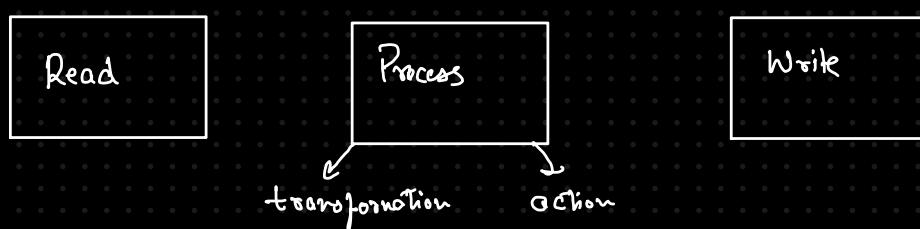


developed to make life  
easy (clean & flexible)  
↑  
Spark higher level

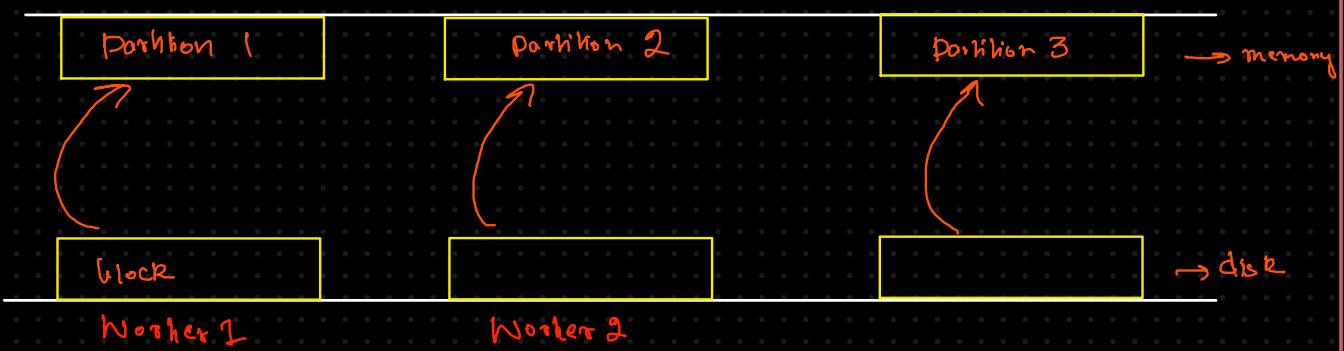
↓  
Spark core / low level  
RDD (most  
flexible)

## 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

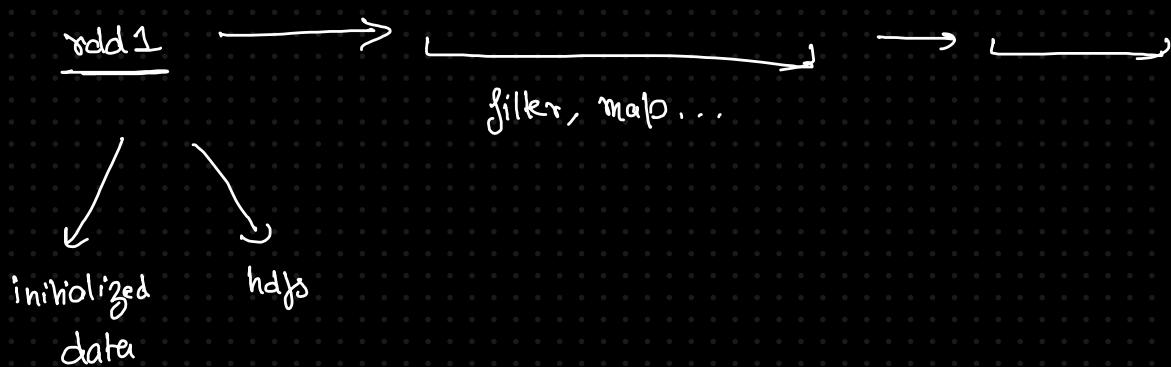
block on disk  
&  
position on memory

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

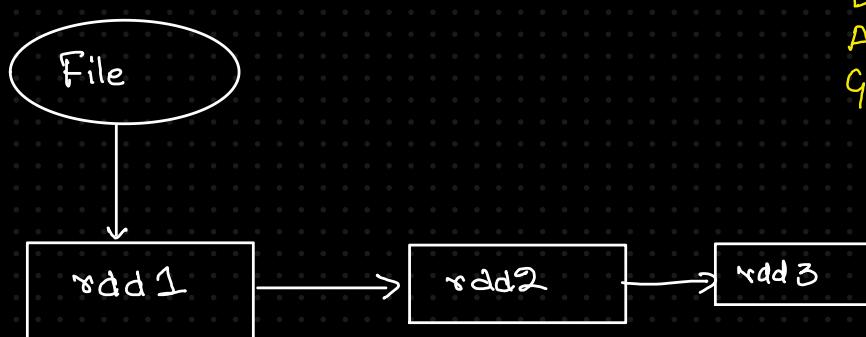
Op<sup>n</sup> that create a new RDD from existing one.

Lazy

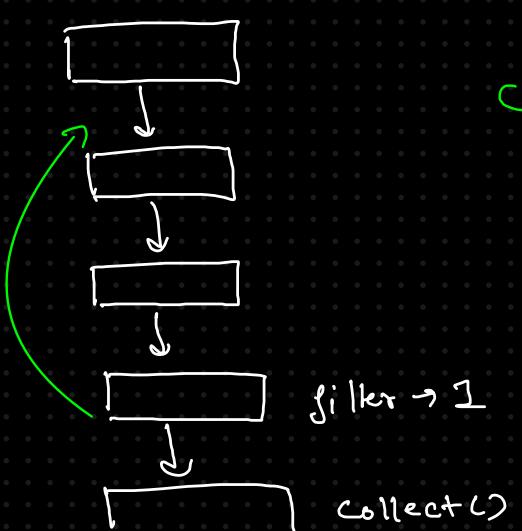
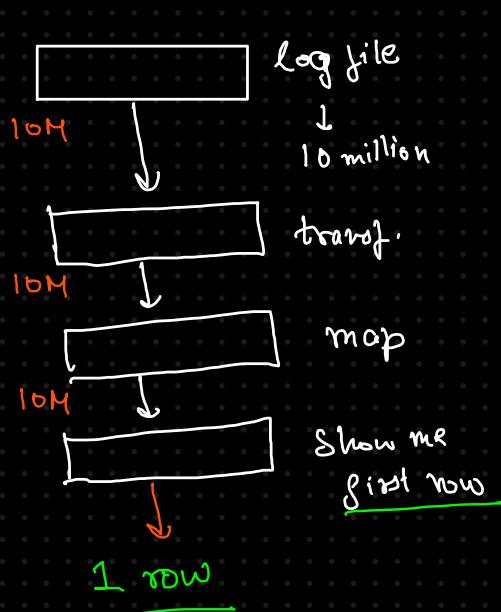
### Action

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 →

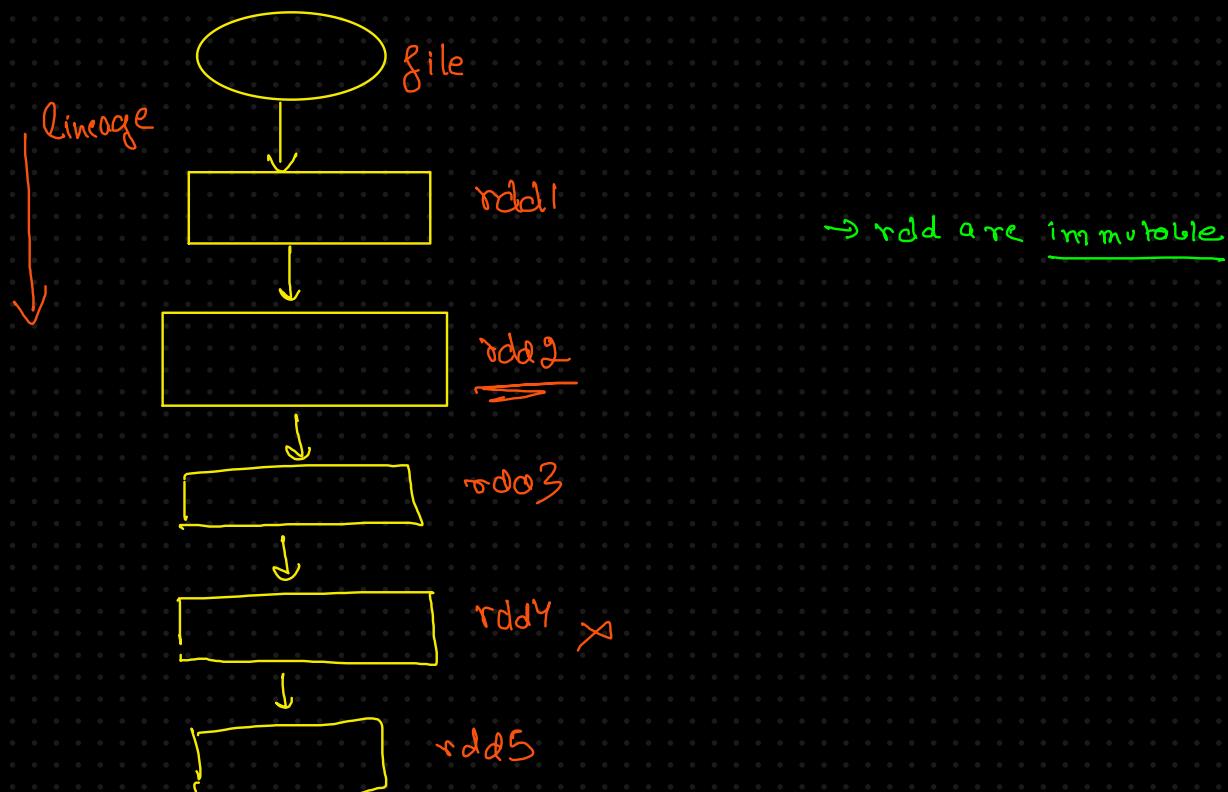
D - distributed ✓

D - dataset ✓

rdd



Resilient will mean that we can quickly recover

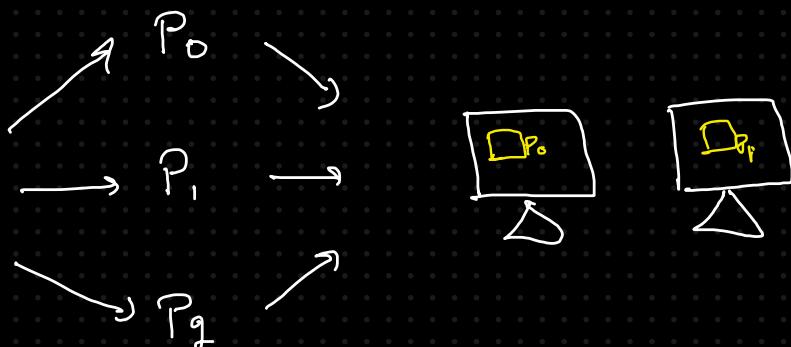
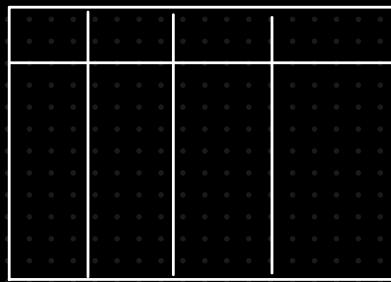


## Properties of rdd

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

