

SPARK

→ fundamentals of

1. distributed processing
2. distributed storage

- most important topic

Interview

Work / Actual Job

→ Theory

→ Practicals

HDFS - majority getting replaced by cloud based platform

YARN - still preferred as a resource manager (mesos, docker, k8s)
but in DE system.

MR - writing & performance both issues are here

In hadoop ecosystem, spark is replacing MR

plugin → spark

Pre-requisite .

1. Python
2. SQL
3. Basics of big data & hadoop

We can separately learn spark
but it's better to learn with
hadoop.

Some common Questions about Spark

Q: Do I need to know Hadoop to learn Spark?

No, we can learn spark independently.

Having idea about distributed storage & distributed processing that is very helpful. YARN (Resource manager) is helpful.

Q. Do I need to install Hadoop to run Spark?

No. Spark can run in standalone mode with say YARN.

If we want HDFS then hadoop is required.

Q. What are the prerequisite to learn Spark?

Python, SQL + knowledge of distributed system is a plus.
or (Java, Scala)

Q. Can I use Spark on my laptop? Where all can we run it?

Yes, it is known as a local mode & good for learning or small testing

- ✓ 1. Google datalaboc cluster
- 2. Google colab
- 3. Local machine + ...
- 4. Datolrichs

Q. What Programming Language can I use with Spark?

Python Java Scala R

Q. How is spark different from traditional databases?

wrong question

Databases are for storing & querying data.

Spark is a data processing language.

Q. What storage system can Spark with?

HDFS, S3, ADLS Gen2, Cassandra, Kafka etc.

Q. Is spark free to use?

Yes, free to use. Apache Foundation

Android



Samsung, Motorola

Q. How is spark different from data warehousing tools like Hive?

Hive is a data warehousing tool

Spark is a general purpose processing engine

Q. Can I use spark on the cloud?

Yes, GCP, AWS, Azure, DataVicks

Limitations of Map Reduce

1. Performance : heavily dependent on disk → high latency
many I/O dependency operations on disk
2. Complexity & hard to write code: boilerplate code
verbose
3. Only batch processing supported: Real time processing not supported
so we cannot work on streaming data
4. Rigid Structure: filter
↳ map → shuffle → reduce
5. No interactive mode/way to monitor: JAR packaged & sent

Apache Spark

Apache spark is an open source, distributed computing system designed to process large scale datasets quickly.

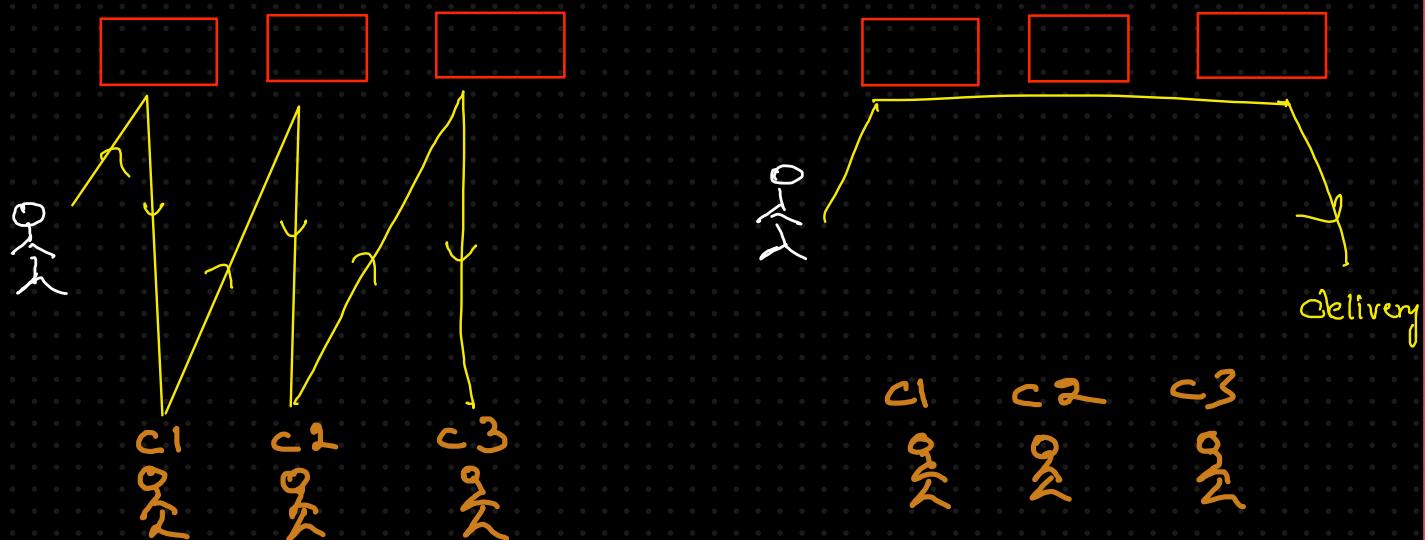
→ fast

→ general purpose → any storage and say resource

→ batch & real time processing

Characteristic

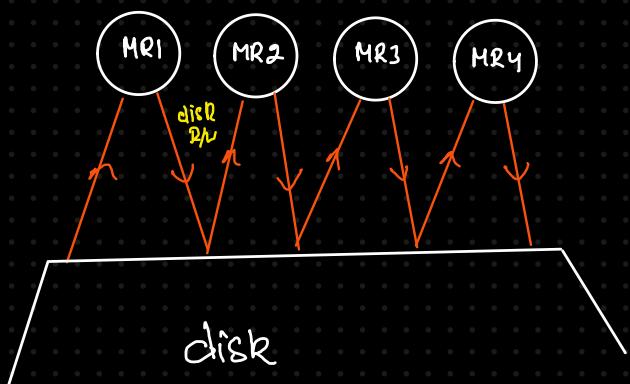
1. In-memory processing : spark processes data in memory as compared to traditional disk-based systems like MR.



2. Ease of Use : spark provides API for major lang. & easy to write code.

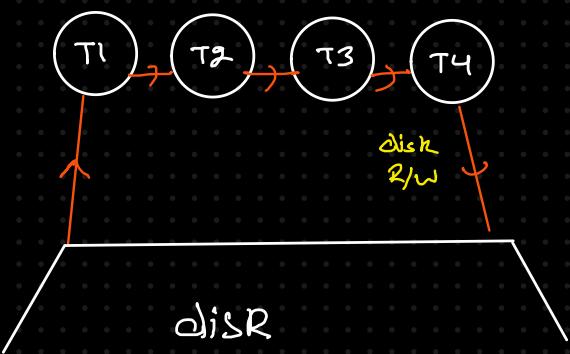
3. Unified Framework:

- Batch processing
- Realtime processing
- ML
- Graph Processing



Map Reduce

5 MR \rightarrow 10 Read & Write



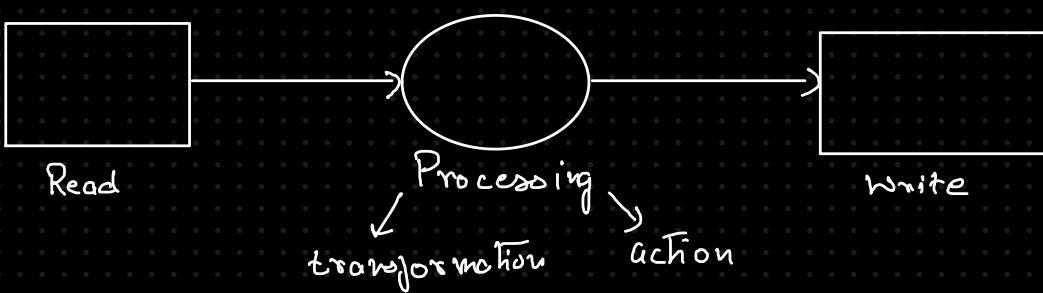
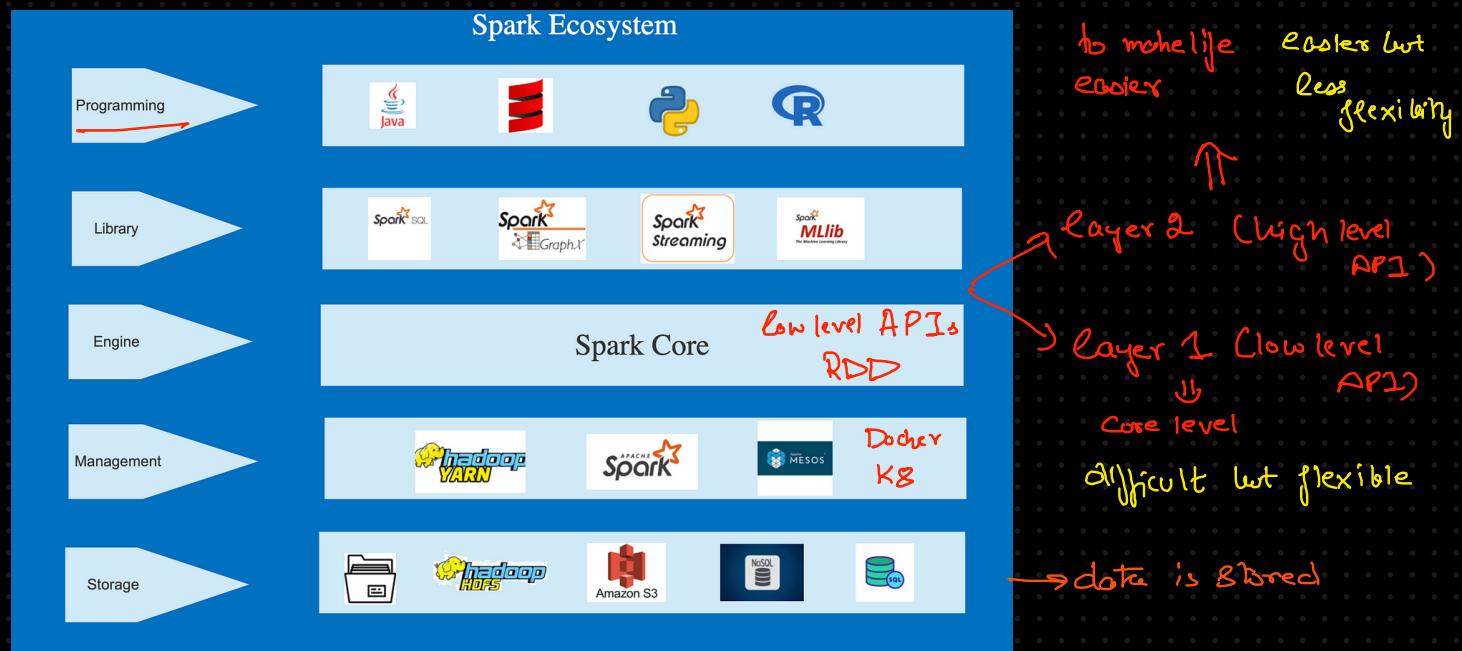
Spark

5 Spark \rightarrow 2 read/write

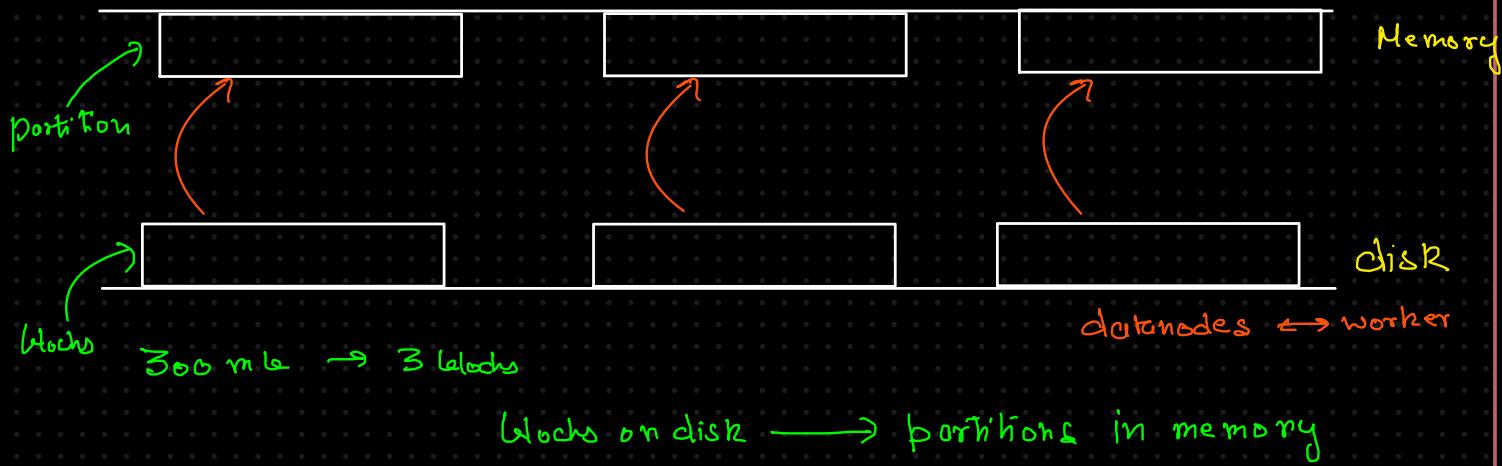
features :

1. Speed
2. Scalability \Rightarrow 1 m/c to 100 of m/c
3. Fault tolerance : DAG (direct acyclic Graph)
4. Polyglot : Multi lang. support
5. Unified Performance
6. Integration with Existing System

Spark Ecosystem



Execution Code in Spark

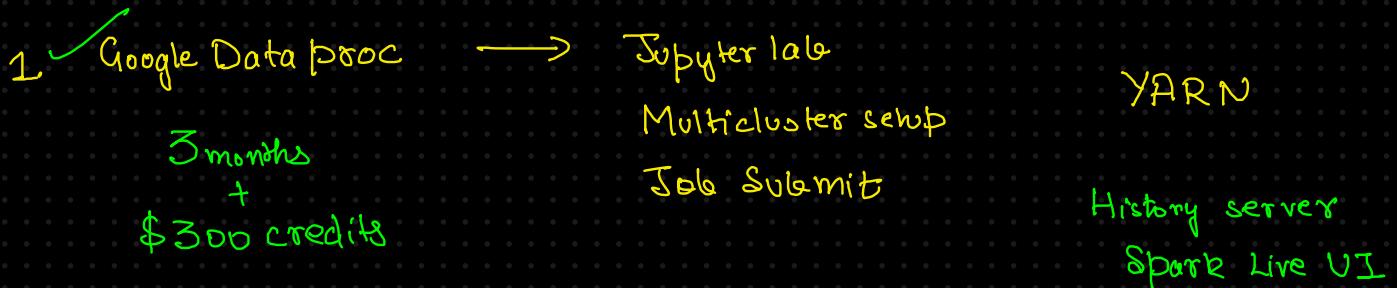


In partitions, the data is still distributed. \$ RDD is most basic unit



Word Count Program in Spark

Spark Code, Spark UI and History Server



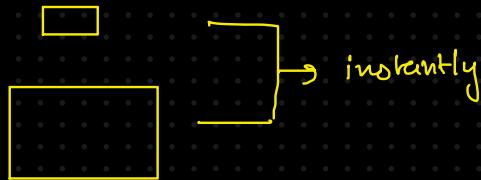
- use a new Id where you have never created your gcloud Id
- use a VISA/mastercard credit card

2. Datavincs easy to create free account
One machine 1gpu & 2cores
We can create notebooks

3. Google Colab We can access spark UI
run it online
free to use

Transformation vs Action

map filter flatmap



Execution
Transformation
dazily evaluated

Action

eagerly evaluated

Definition
Operation which creates
a new RDD from existing
one

an operation that triggers
execution of transform
& return results.

Output
new RDD

final value / writes output
to storage

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		

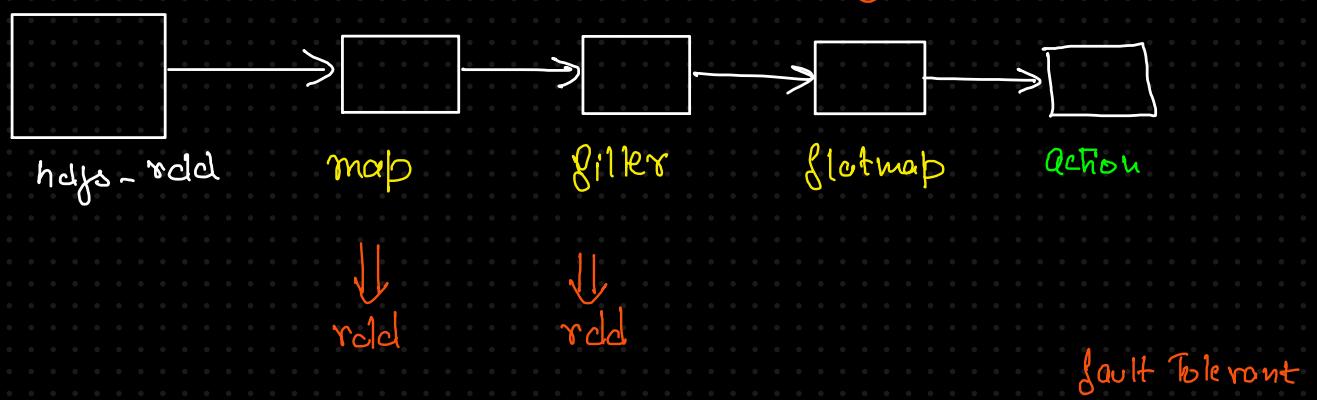
Occurs when an action
is called

Execution is immediate

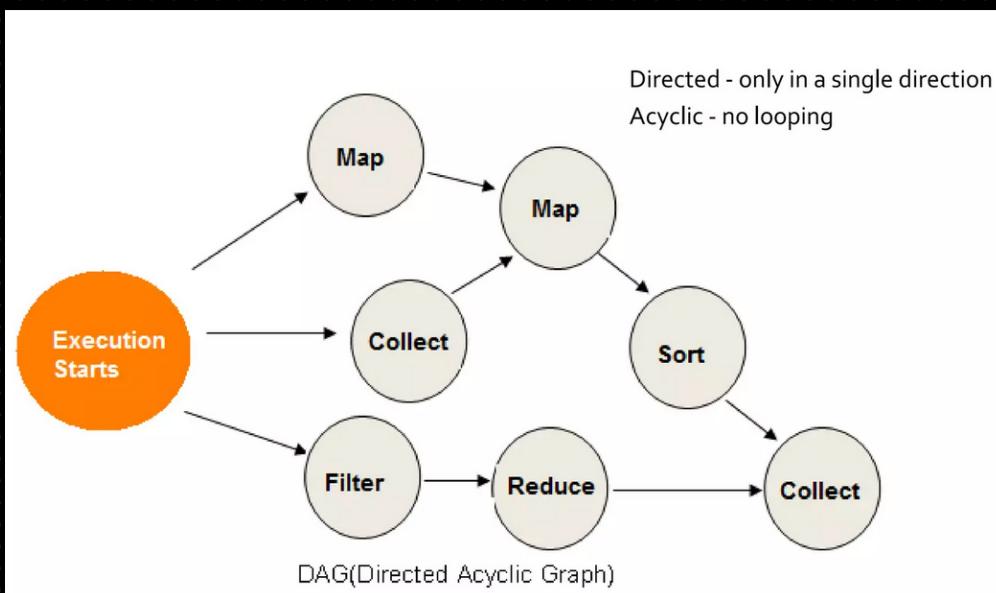
Represented as an intermediate
node in our DAG

Represented as a terminal
node.

Direct acyclic Graph (DAG)

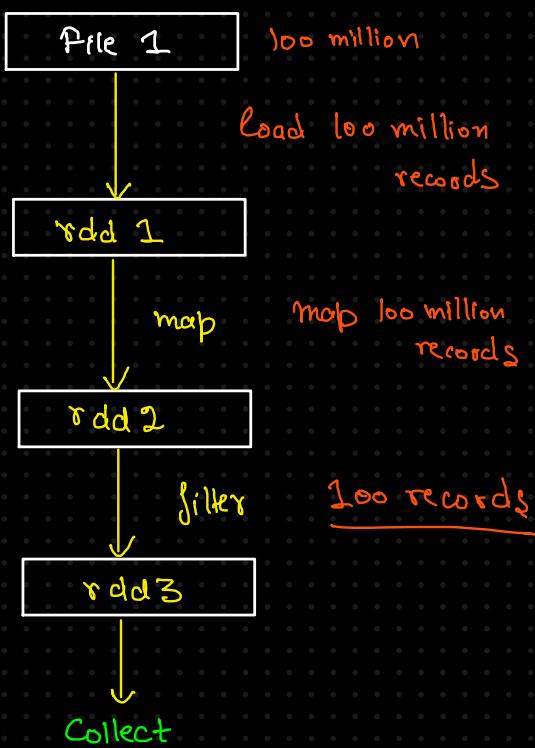


At each step a new RDD is getting created

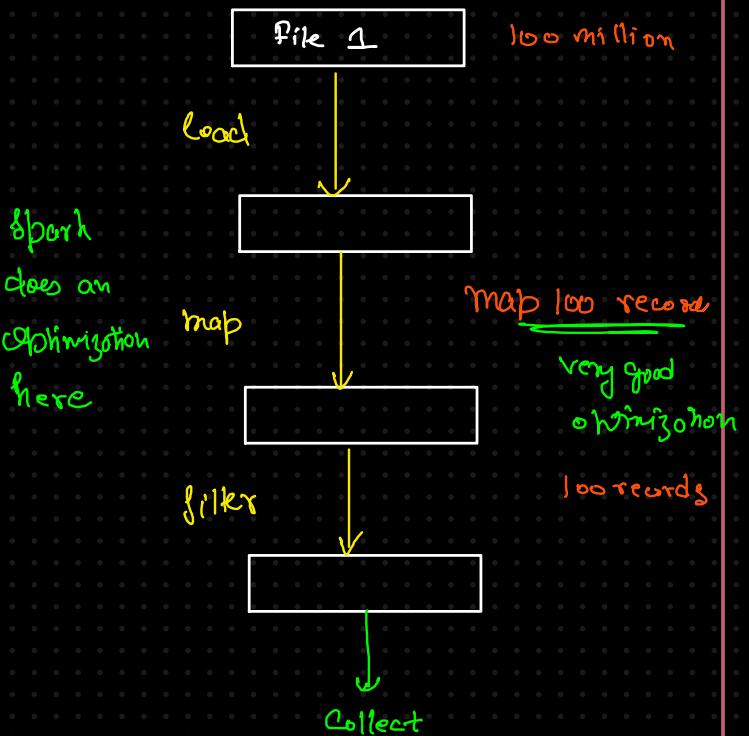


Why transformations are lazy?

not lazy



Lazy



Shopping ⇒ Discount

Cake ⇒ 1kg flour + 1kg chocolate + 1kg sugar → mix

↓
add some water

↓
bake

↓
5 people

eat/serve

Spark



low level APIs

R - resilient → resistant to failure

D - distributed →

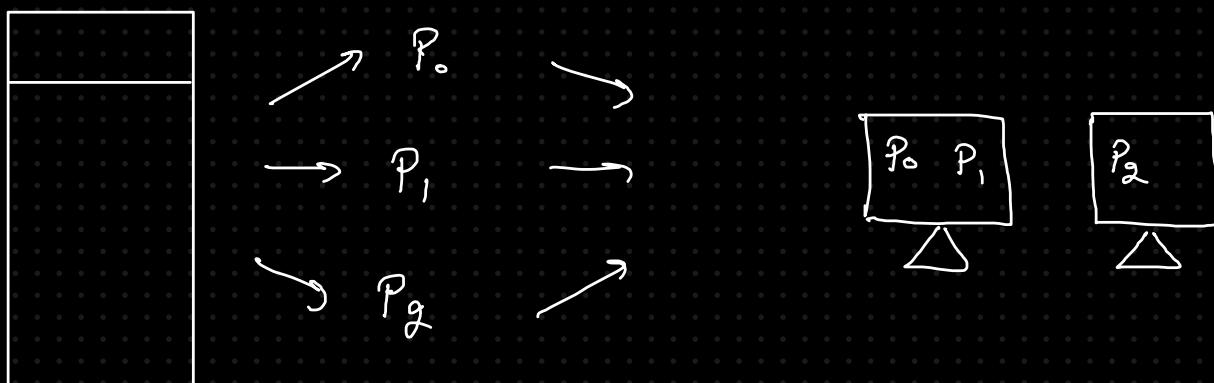
D - dataset

DAG

we need spark context to work with RDD

Properties

1. Immutability : Cannot be modified
2. Lazy evaluation: Computation not performed until an action is triggered
3. Partitioned :
4. Foundation data structure.

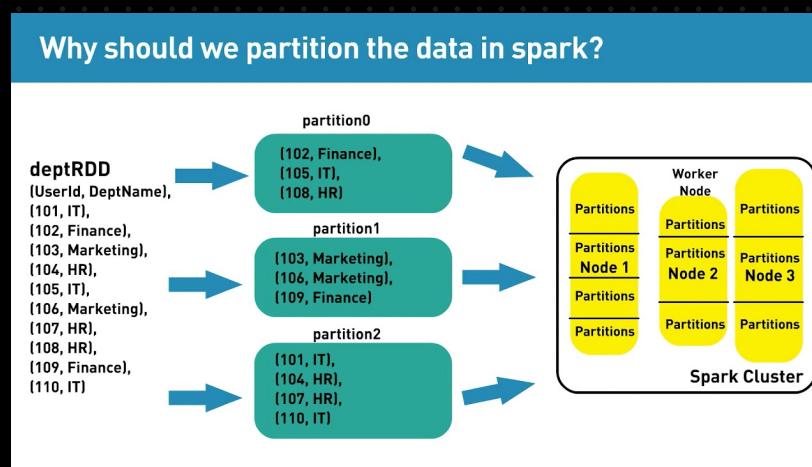


How Spark does data partitioning?

Reading the data / defining data will of course be there

Data → rdd

Partitioning in spark refers to dividing large datasets into smaller, manageable chunks called partitions.



Aspect	Input Split (Hadoop)	Partition (Spark)
Nature	Logical division of input data for mappers	Logical division of RDDs/DataFrames for tasks
Creation	Defined by `InputFormat`, often based on HDFS blocks	Typically one per HDFS block but configurable
Control	Less dynamic control over split size	Dynamic control via transformations
Execution Model	Disk-based processing	In-memory processing
Performance Impact	Affects map task efficiency	Directly influences task scheduling and parallelism

How spark reads the data

1. Source of data → S3, hdfs, adls gen2, dB

Types of
data sources

2. Logical partitioning →

alec.txt → 1 glo
(hdfs)

8 blocks → 8 partitions

already
distributed
like s3, hdfs

read data
& create
partition
based on
configuration

3. Distributed reading / processing

Each partition is processed independently by a
executor in the cluster.

↓

fetch data from storage → in-memory

↓

decompose & action

8 blocks → 8 partitions

+ -

Spark parallelize and partitioning

sc.parallelize()

a function to create an RDD from an in-memory location

sc.textfile()

a function to read a file.

540mb file \Rightarrow 5 blocks \Rightarrow 5 partition

1mb file \Rightarrow 1 block \Rightarrow 2 partition

default minPartitions
 $\underline{= 2}$

for local file \Rightarrow n bytes \Rightarrow 4 partition

\swarrow default Parallelism
 $\underline{= 4}$

Properties are dependent on platform.

\downarrow
4 cores

Spark want to use all the cores in the machine for better parallelization.

```
# get num of Partition from rdd =>getNumPartitions()

# 544mb => 5 blocks => 5 partition
big_rdd_from_hdfs.getNumPartitions()

5

# 1 mb => 1 block --> 1 partition ( This is wrong)
small_rdd_from_hdfs.getNumPartitions()

2

print(f"Max Partition Bytes: {spark.conf.get('spark.sql.files.maxPartitionBytes')}")
Max Partition Bytes: 134217728b

spark.sparkContext.defaultMinPartitions

2

local_rdd.getNumPartitions()

4

spark.sparkContext.defaultParallelism

4
```

Understanding the test data & its Generation

Table	Column Name	Description
Customers	customer_id	Unique identifier for each customer.
	name	Name of the customer.
	city	City where the customer resides.
	state	State of the customer.
	country	Country of the customer (set to India).
	registration_date	Date when the customer registered.
Orders	is_active	Whether the customer is currently active or not.
	order_id	Unique identifier for each order.
	customer_id	Foreign key linking to the <code>customers</code> table (which customer placed the order).
	order_date	Date when the order was placed.
	total_amount	Total amount of the order.
Items	status	Current status of the order (e.g., "Pending", "Shipped", "Delivered", "Cancelled").
	item_id	Unique identifier for each item.
	order_id	Foreign key linking to the <code>orders</code> table (which order the item belongs to).
	item_name	Name of the item.
	category	Category of the item (e.g., "Electronics", "Clothing").
Payments	price	Price of the item.
	payment_id	Unique identifier for each payment.
	order_id	Foreign key linking to the <code>orders</code> table (which order the payment is for).
	payment_date	Date when the payment was made.
	amount	Amount paid.
Shipments	payment_method	Method of payment (e.g., "Credit Card", "Debit Card", "UPI").
	shipping_id	Unique identifier for each shipping.
	order_id	Foreign key linking to the <code>orders</code> table (which order the shipping is for).
	shipping_date	Date when the shipping was made.
	shipping_address	Shipping address where the order was delivered.
Shipments	shipping_method	Method used for shipping (e.g., "Standard", "Express").

$\text{deg} \Rightarrow \{ , , , , \}$

Customers - 1 mle
- 500 mle

E-commerce-related → 50+ tiles

1. We try to understand the data
 2. We are going to understand problem statement.
 3. We use/solve the problem on smaller data & then we scale it.

1 sec
↓
0.5 sec

for projects, we will use this data and some real world so we get full confidence.

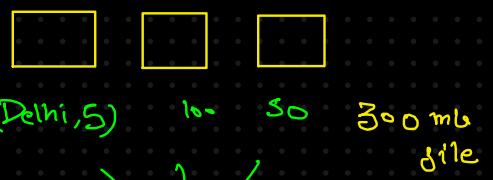
Spark RDD operations

Refer the notebook

```
# we can chain the functions together
customers_per_city = parsed_rdd.map(lambda row : (row[2],1)).reduceByKey(lambda x,y:x+y)
customers_per_city
PythonRDD[40] at RDD at PythonRDD.scala:53
customers_per_city.collect()
[('Bangalore', 3), ('Hyderabad', 1), ('Ahmedabad', 1), ('Delhi', 1)]
# CountByValue
cust_per_city = parsed_rdd.map(lambda row:row[2]).countByValue()
cust_per_city
defaultdict(int, {'Bangalore': 3, 'Hyderabad': 1, 'Ahmedabad': 1, 'Delhi': 1})
```

Transformation: Narrow vs Wide

Transformations in Spark are operations that create new RDD.



Narrow

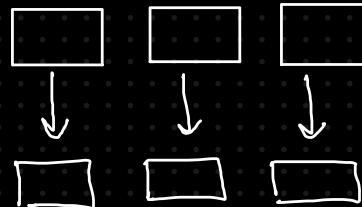
filters

Data movement

Data is not shuffled across partitions

Dependency

Each child partition depends on single parent partition



Performance

Faster as it avoids data shuffling

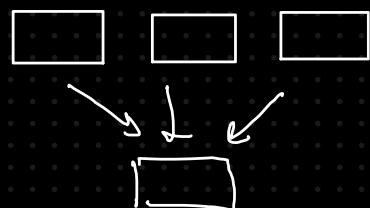
map, filter, flat map

Wide

reduce By Key

Data is shuffled across partitions, requiring network I/O

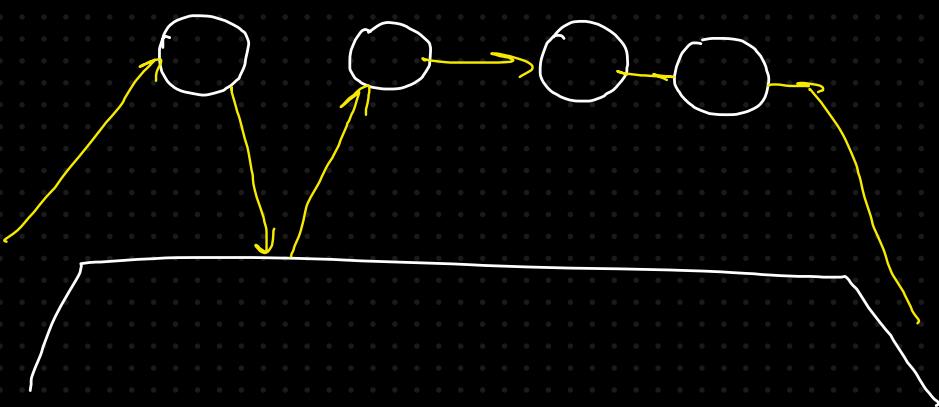
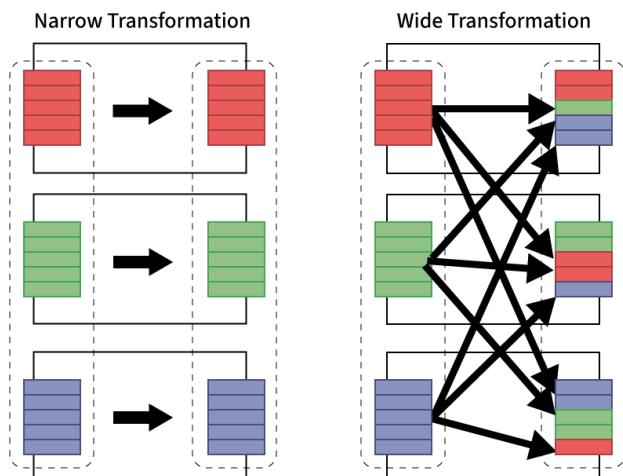
Each child partition depends on multiple parent partition



Slower due to network & disk I/O

reduce By Key , Group By Key

Wide & Narrow dependencies in Apache Spark



Final aim

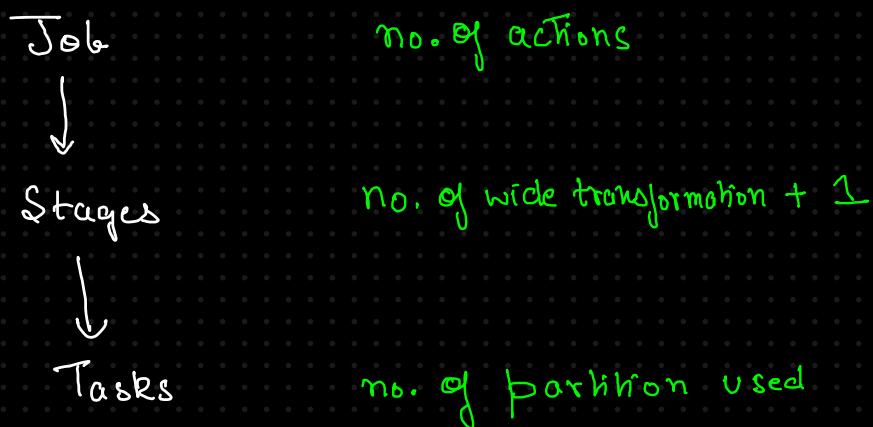
1. try to have very less wide transformation
2. narrow down the data so that data getting shuffled is less.

Bang, 1

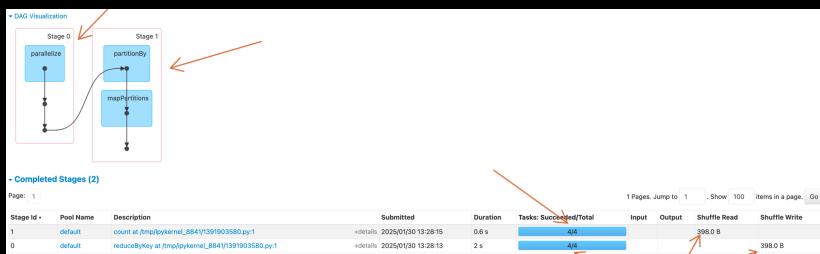
10k

(Bang, 10k)

Job, tasks and stages in Spark UI



The same can be optimized / diff based on platform



Reduce By Key and Group by Key

472 bytes

2 KB

674 bytes

9 KB

→ for Ime file

→ for Sample file

Reduce By Key

Output

Returns a key-value pair with aggregated values

Local aggregation is there.

(City, 1)
(City, 1)
(Delhi, 1)



(Brs, 2)
(Delhi, 1)



minimize data transfer
by doing a local aggregation

less data for shuffling
+
I/O for disk/network

Group by Key

Returns a key value pair where values are collections

no local aggregation

(Brs, 1)
(Brs, 1)
(Delhi, 1)

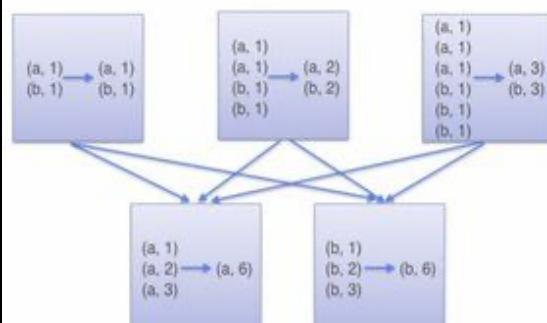


Group by key
(Brs, 1)
(Brs, 1)
(Delhi, 1)

Group all the values
based on key & aggregation
involves after shuffle.

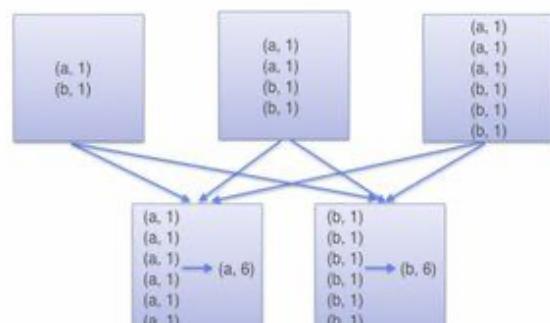
more data for
shuffling

ReduceByKey



max sum min

GroupByKey



Note: For non-commutative and non-associative property like say median, avg etc. we have to use group by key

➡ of group by key

1. chances of OOM err
2. Shuffling lots of data
3. not doing / restricting lism.

Reduce key key

	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
el_32253/266034041.py:1 +details	2025/01/30	2 s	5/5			2.1 KiB	
at el_32253/266034041.py:1 +details	2025/01/30	12 s	5/5	544.6 MiB		2.1 KiB	

Group by key

	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
el_32253/1858560816.py:1 +details	2025/01/30	0.3 s	5/5			96.1 KiB	
at el_32253/1858560816.py:1 +details	2025/01/30	10 s	5/5	544.6 MiB		96.1 KiB	

~49%

Increasing or Decreasing Number of Partition

No. of partitions affect the parallelism
 \approx no. of tasks

Increase the number of partitions

2 file \rightarrow 16 partitions

You have a great 30 node cluster

it's better to use all the nodes

\Rightarrow can be helpful to increase them

Cluster should be able to handle it.

Decreasing the number of partition

1 file \rightarrow 2000 partition

& say a 50 node cluster

\Rightarrow decreasing is better

say we have 100 partitions

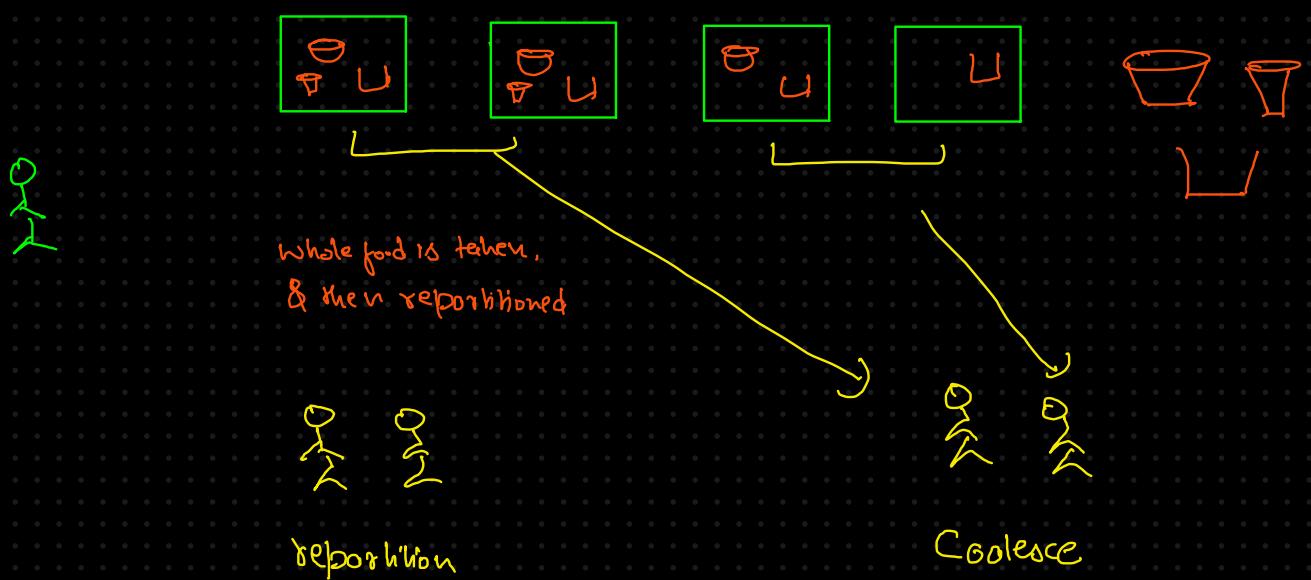
8 node cluster

100 partitions \Rightarrow filter



Avoid sparse partition but should aim for dense partition

Repartition vs Coalesce



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

we prefer coalesce when
we have to reduce
no. of partition

Optimization in Spark is a incremental process

