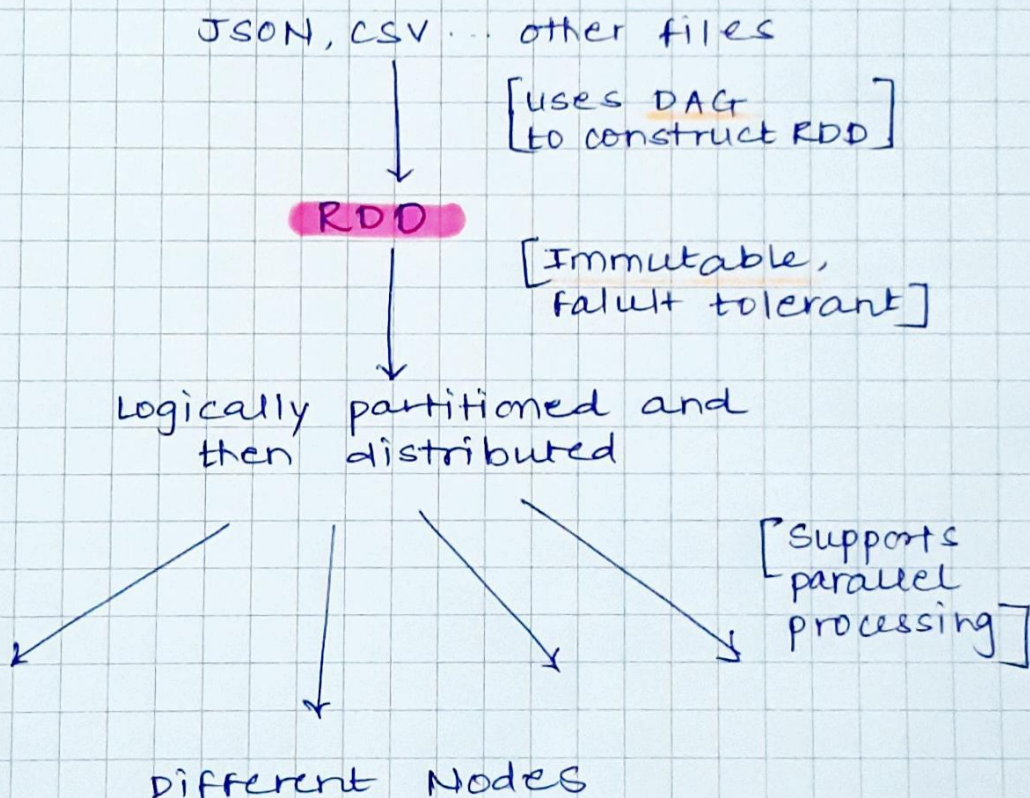


SPARK RDD

RDD : Resilient Distributed Dataset

- Fault Tolerant
- We can recompute the lost/missing/ corrupted data/partition using RDD.
- Every dataset is logically partitioned and distributed across servers.
- The data structure used for logical distribution is RDD.



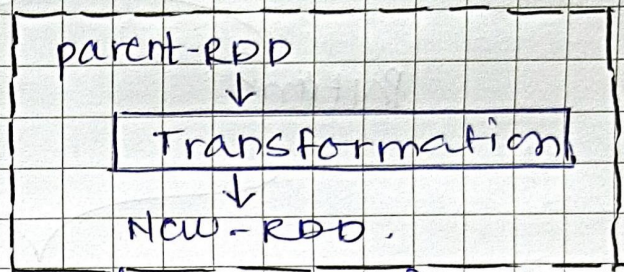
Note : DAG - Directed Acyclic Graph.

Spark RDD Features

<u>Features</u>	<u>Brief about feature.</u>
1. In-memory Computation	• RDDs are stored in-memory for desired timeline
2. Lazy Evaluation	• computation is performed only after action is triggered
3. Fault Tolerant	• In case of loss/failure the partition can be recomputed.
4. Immutable	• Once created, RDDs can't be manipulated.
5. Persistence	• Frequently used RDD's can be stored in in-memory.
6. Partitioning	• seamless distribution on nodes in the cluster.
7. Parallelism	• RDDs are processed parallelly.
8. Typed	• like <code>RDD[Int, long, String]</code>
9. No limitation	• Supports any number of RDDs.
10. Coarse-grained	• Coarse-grained: operations gets applied to entire RDD.
11. Location-stickiness	• Data is placed closer to task using placement reference.

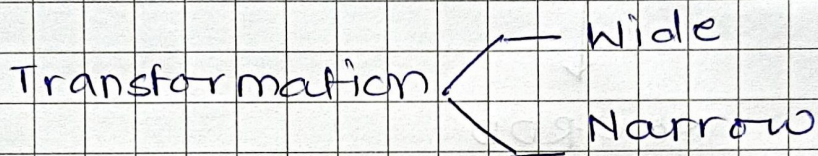
- # RDD operations
- 1. Transformation
 - 2. Action.

Transformation :

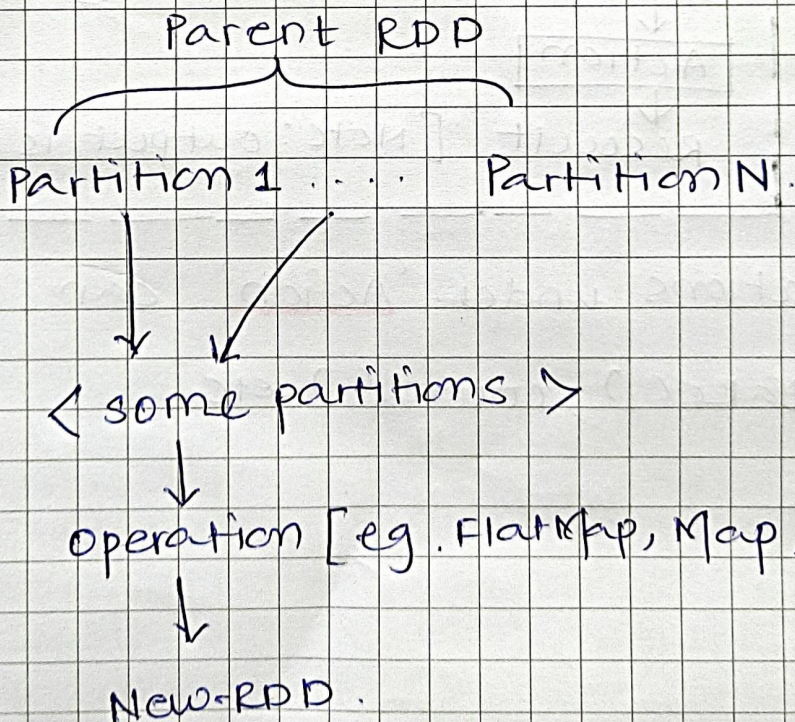


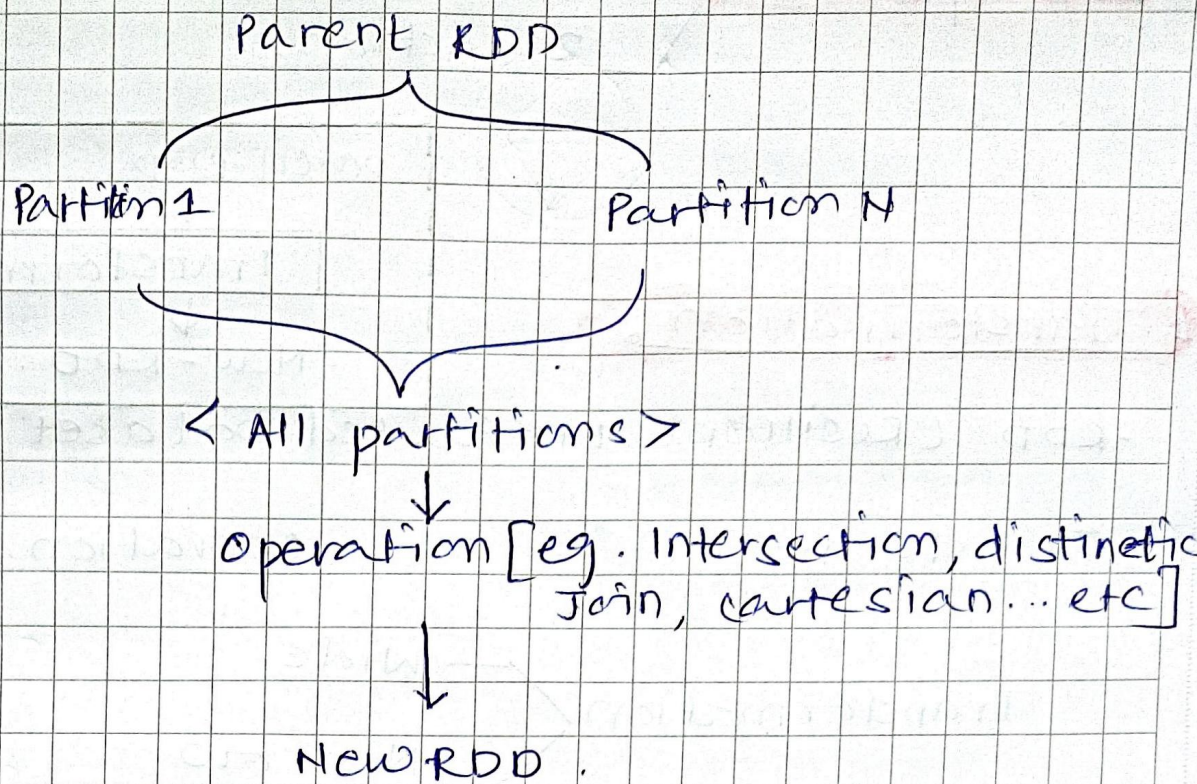
- RDD (Resilient Distributed Dataset)

- Transformation is lazy operation.

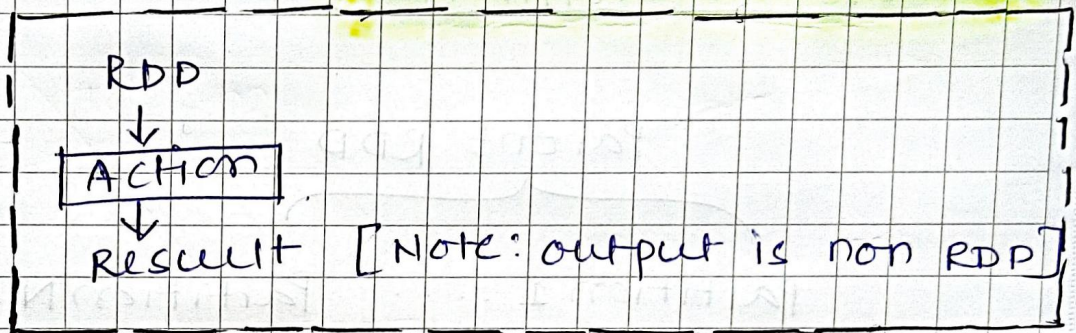


Narrow-Transformation





2. Action:



- Few operations under 'Action' can be `first()`, `take()`, `count()` etc.

Spark Map and FlatMap

- What is Map?

Takes '1' entity as an input and generates exactly 1 entity as a output.

- What is FlatMap?

takes '1' entity as a input and generates '0' or 'many' entities as a output.

In Spark, Map and FlatMap are transformations, hence input and output generated are RDDs.

Map

RDD



Map

programming-logic()



RDD

Size of input RDD and output RDD are same.

Flat-Map

RDD



FlatMap

Programming-logic()



0 or N-RDD

Size of input and output RDD is different.