

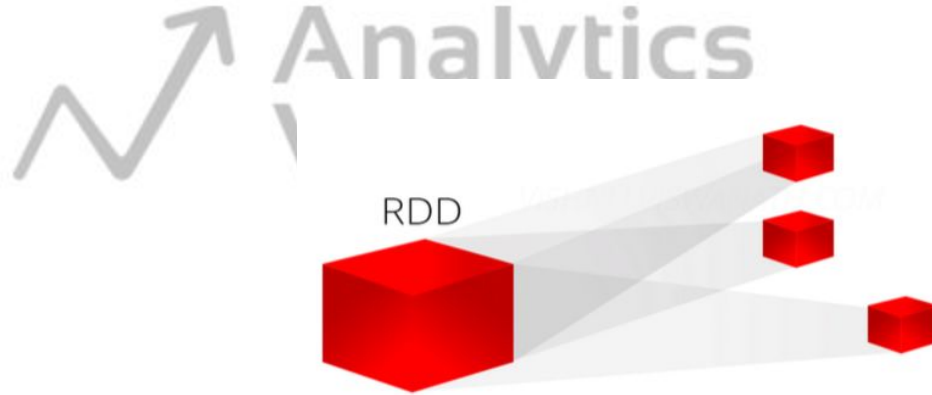


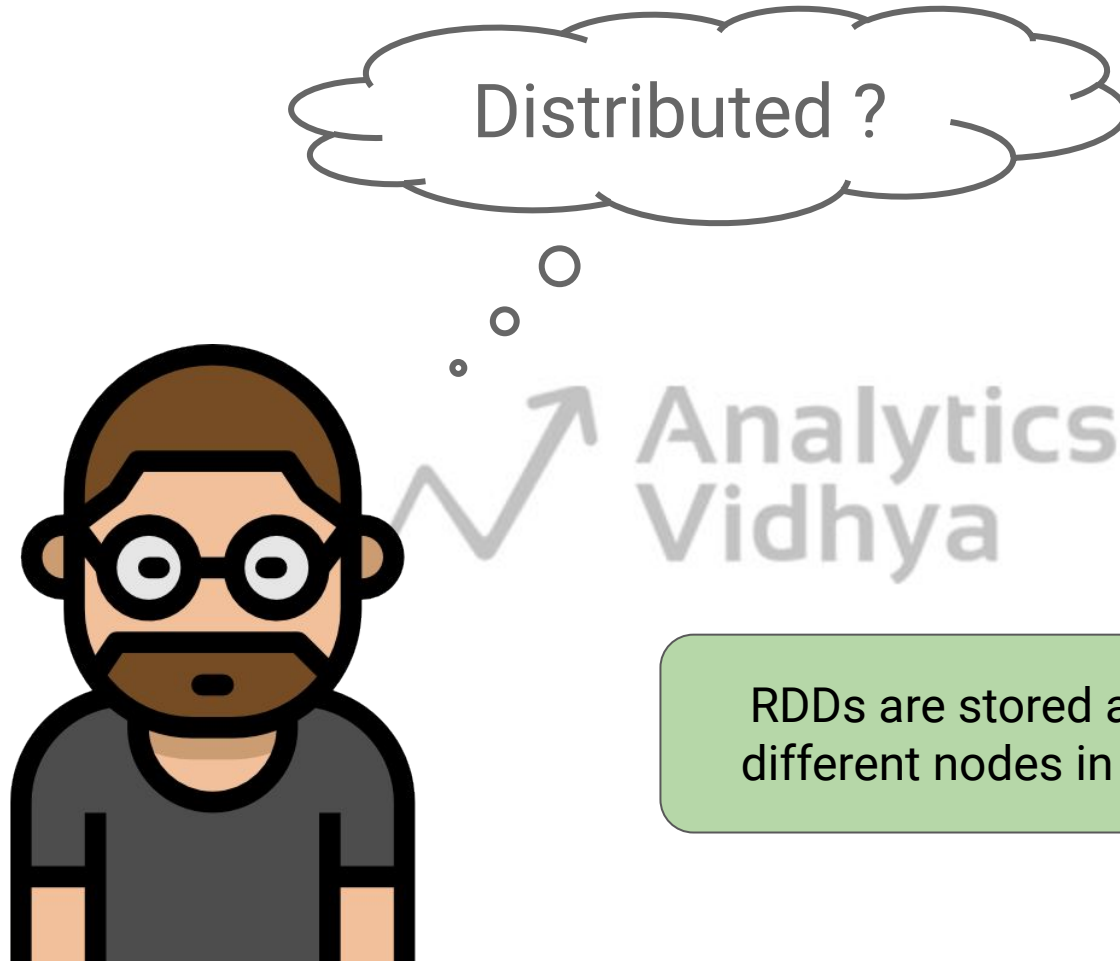
What are RDDs?

# What are RDDs?

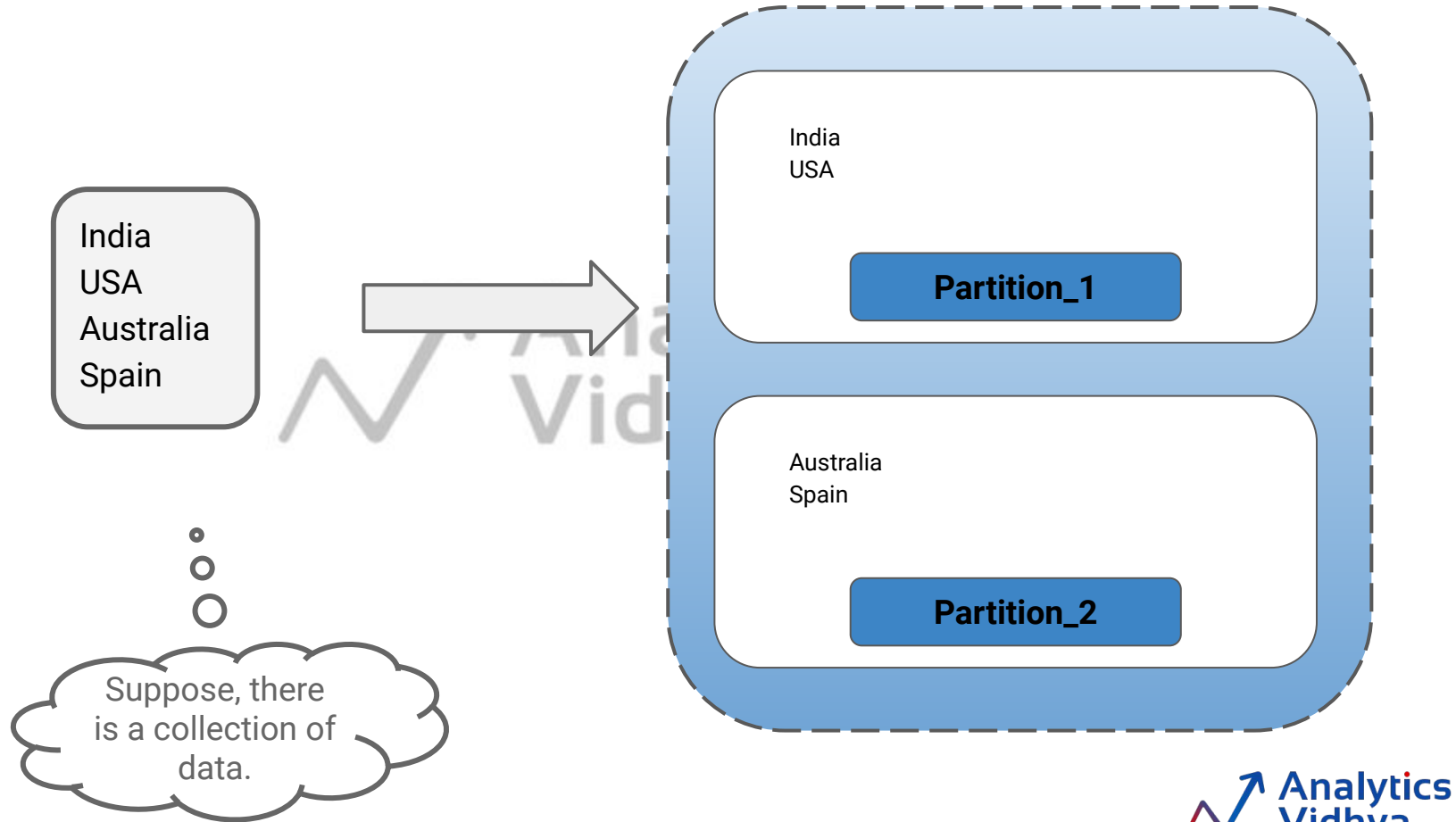
**RDDs (Resilient Distributed Dataset)** is the Spark's main abstraction and has the following features:

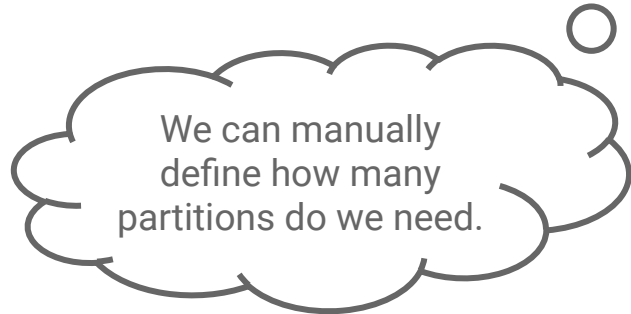
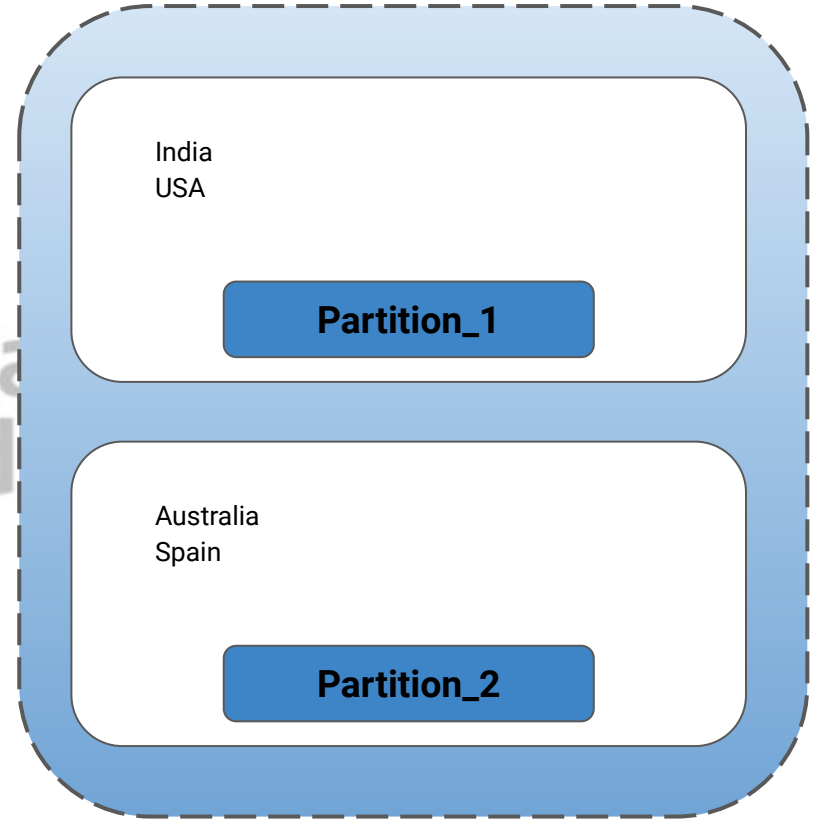
- Distributed
- Immutable
- Resilient





RDDs are stored and processed on different nodes in the Spark cluster.







Immutable ?

Once an RDD is created, it cannot be modified further. It can only be used to create a new RDD

RDD

1, 12, 4, 6, 8  
3, 14, 5, 1, 3

Partition\_1

1, 3, 4, 4, 5  
2, 4, 1, 1, 3

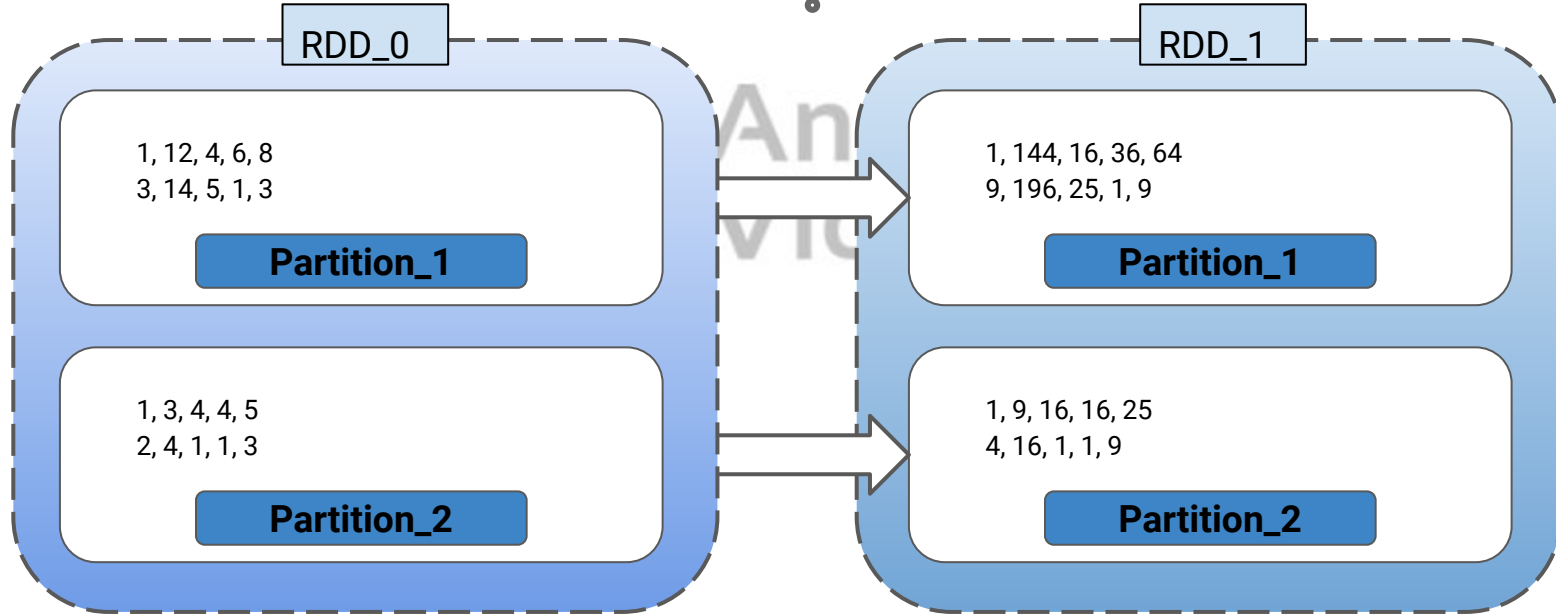
Partition\_2

Now, if you want to do any operations on a RDD

square of each number present in the data

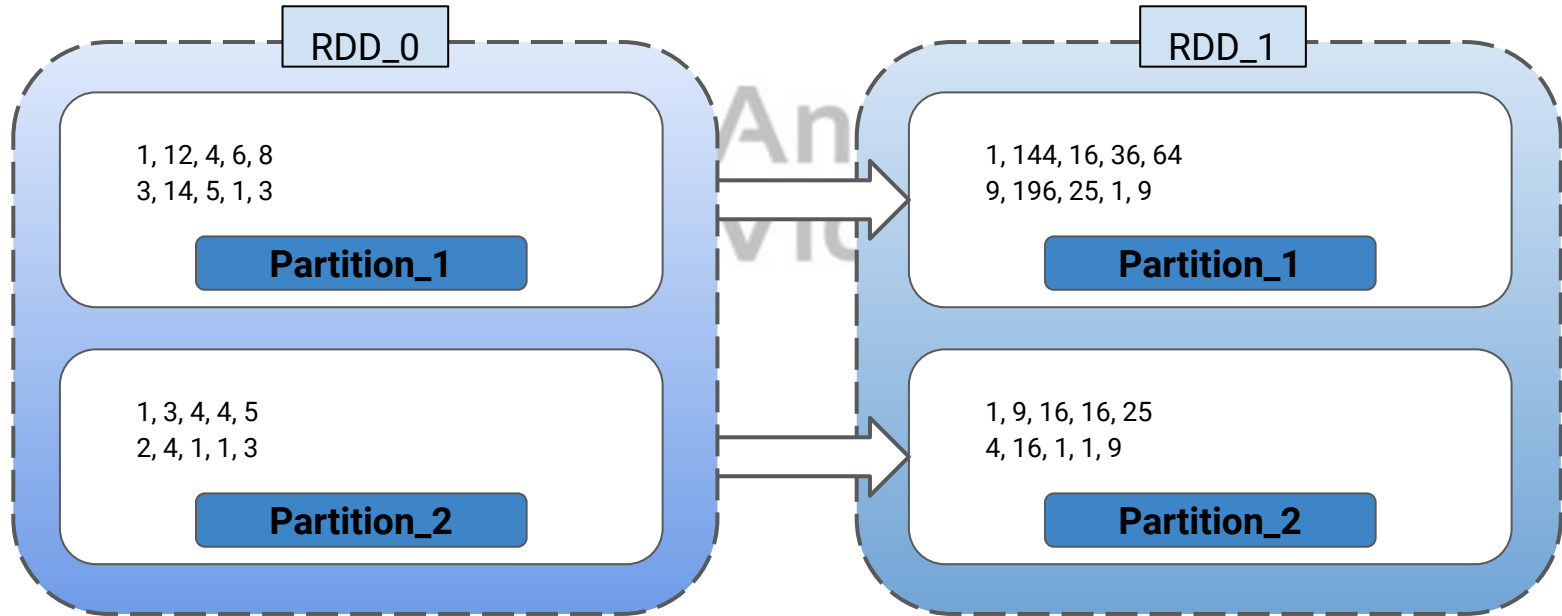
It will **not update** within the same RDD rather it will **create another one**

Square  
Operation on  
RDD





This process is called  
**Transformation.**





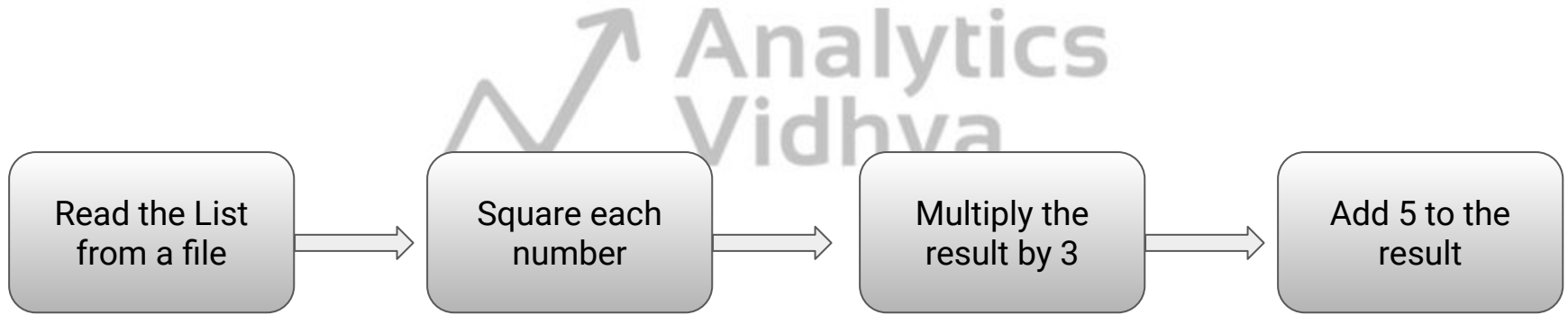
Resilient ?

Analytics  
Vidhya

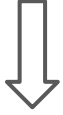
RDDs have the ability to continue the operations even if there is a failure.

# Let us Assume?

- We have to build the following operations on a list of numbers.



[ 1, 4, 7, 5, 3]  
[ 2, 3, 3, 6, 1]



RDD\_0

[ 1, 4, 7, 5, 3]

Partition\_1

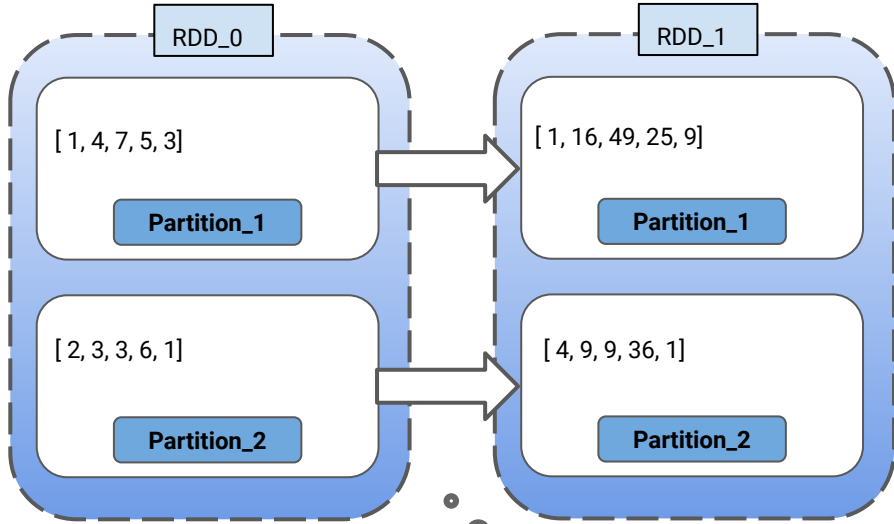
[ 2, 3, 3, 6, 1]

Partition\_2

data is partitioned  
on 2 nodes.

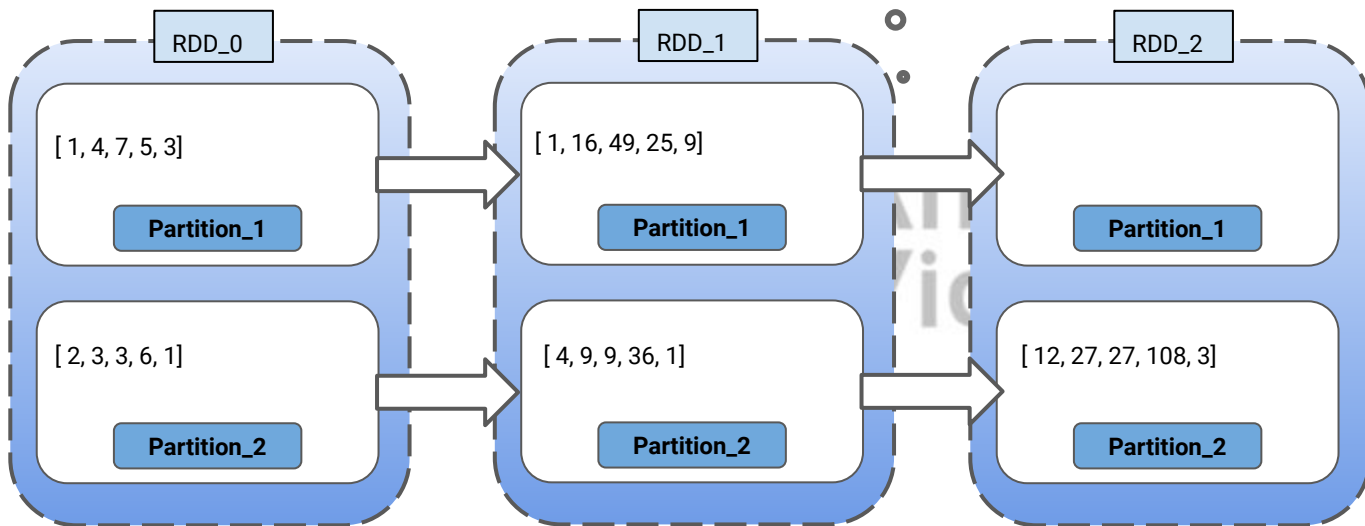


Analytics  
Vidhya

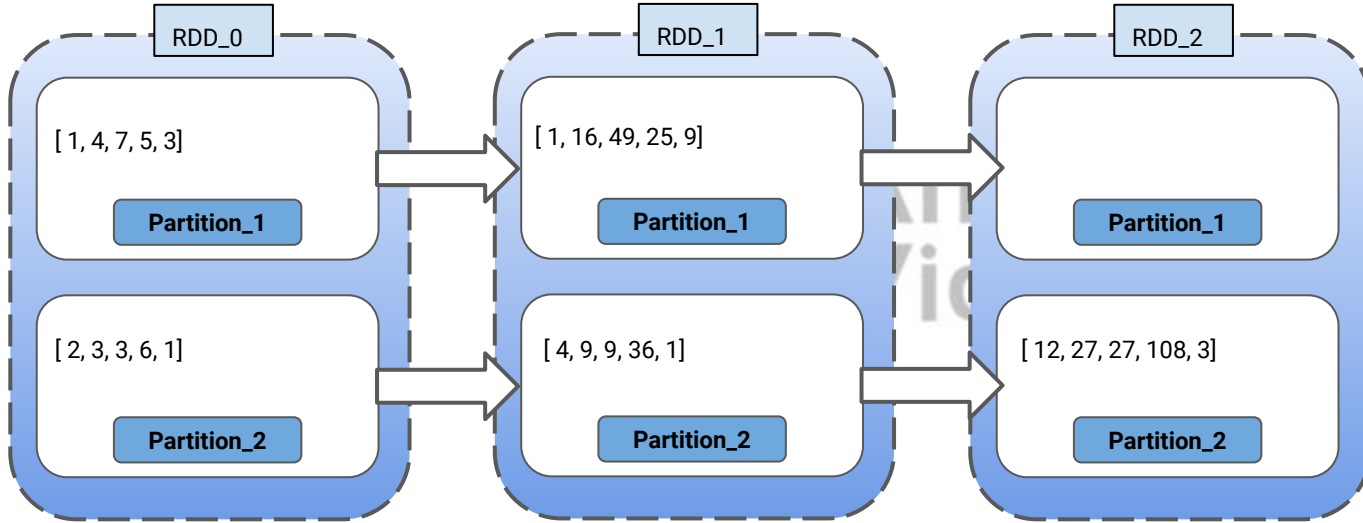


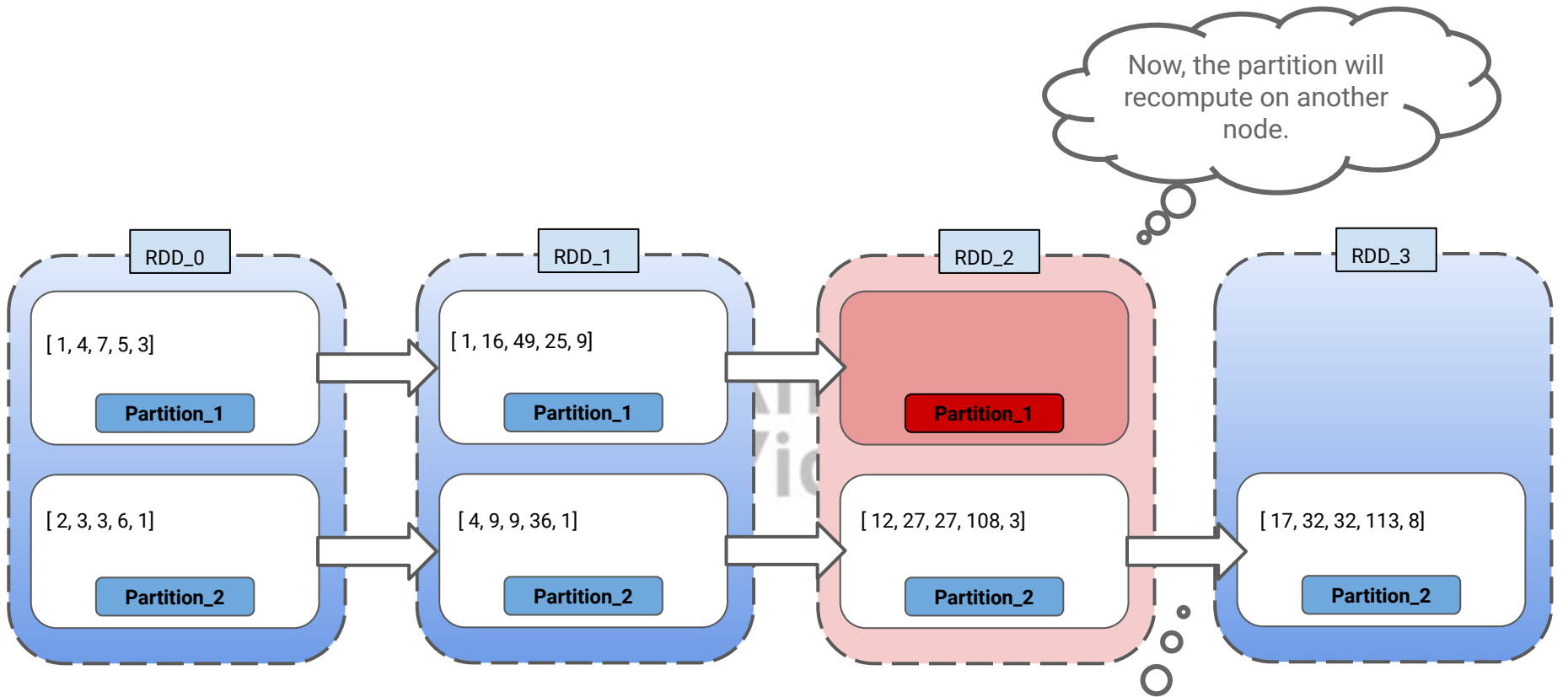
Now, in the first step we will square each number in the list.

In the next step, we will multiply the result by 3.



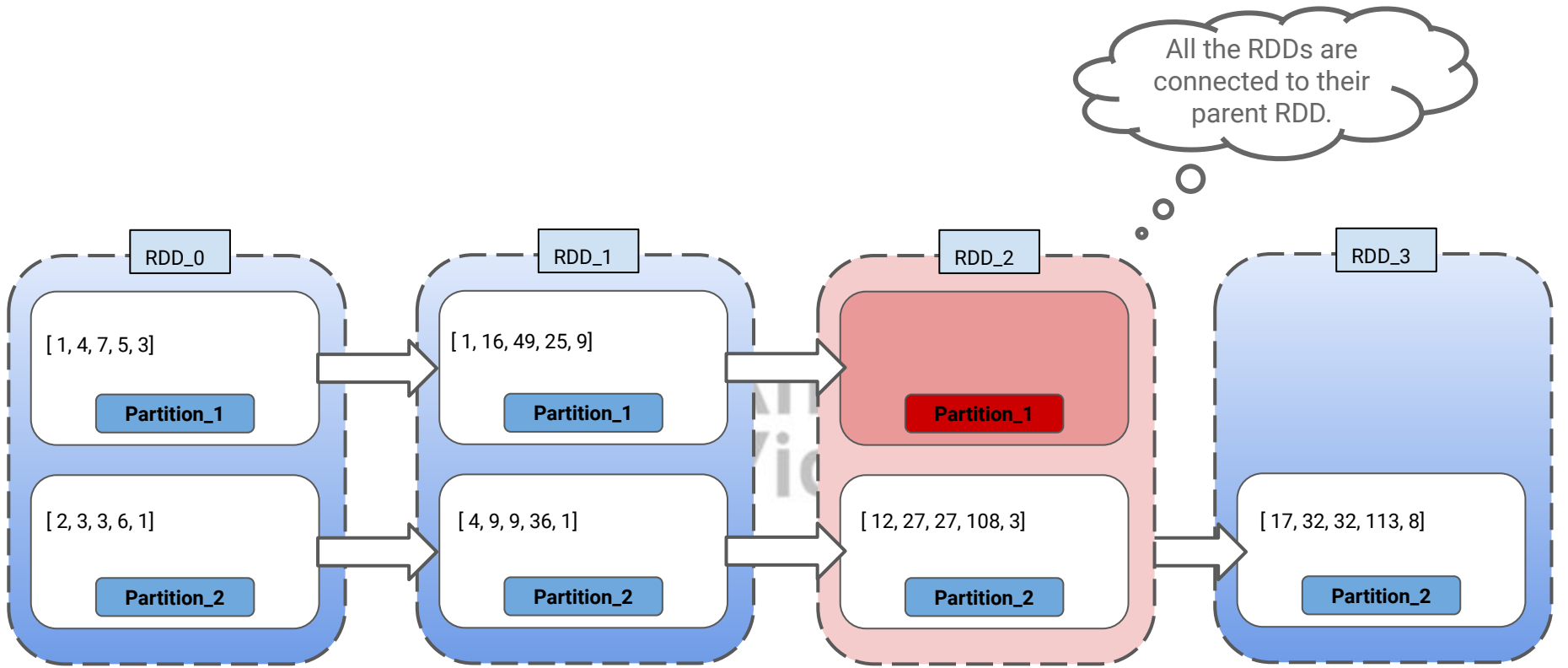
Now, due to some reasons one of the node fails.

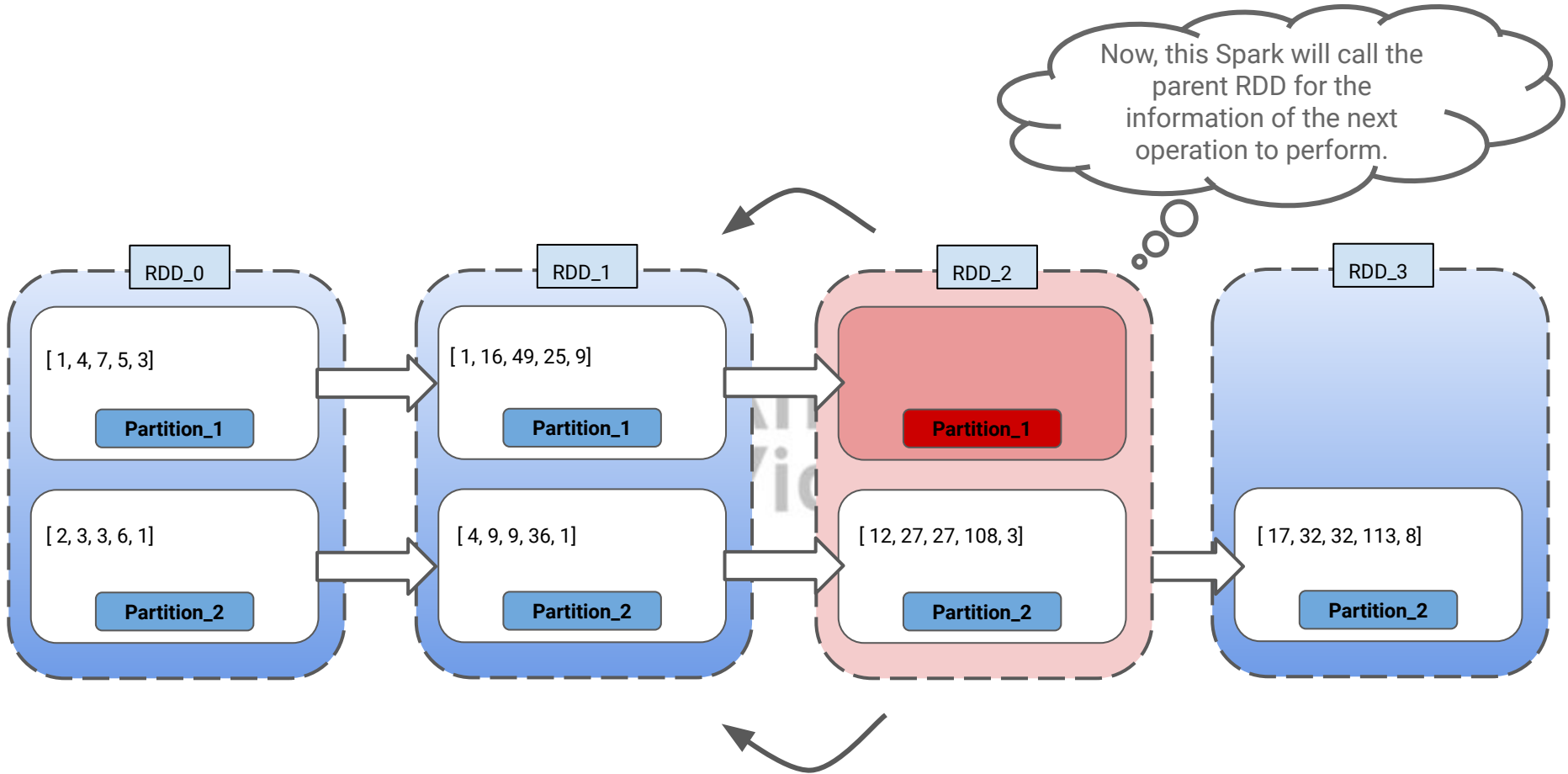


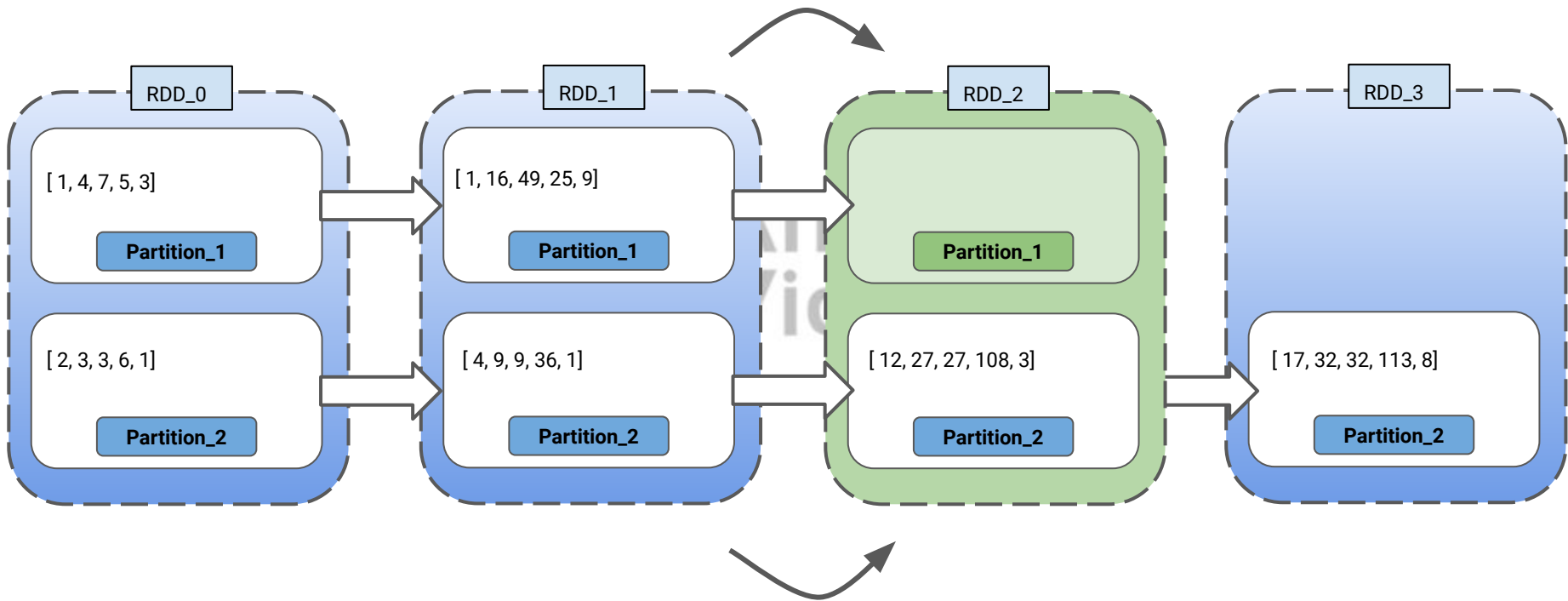


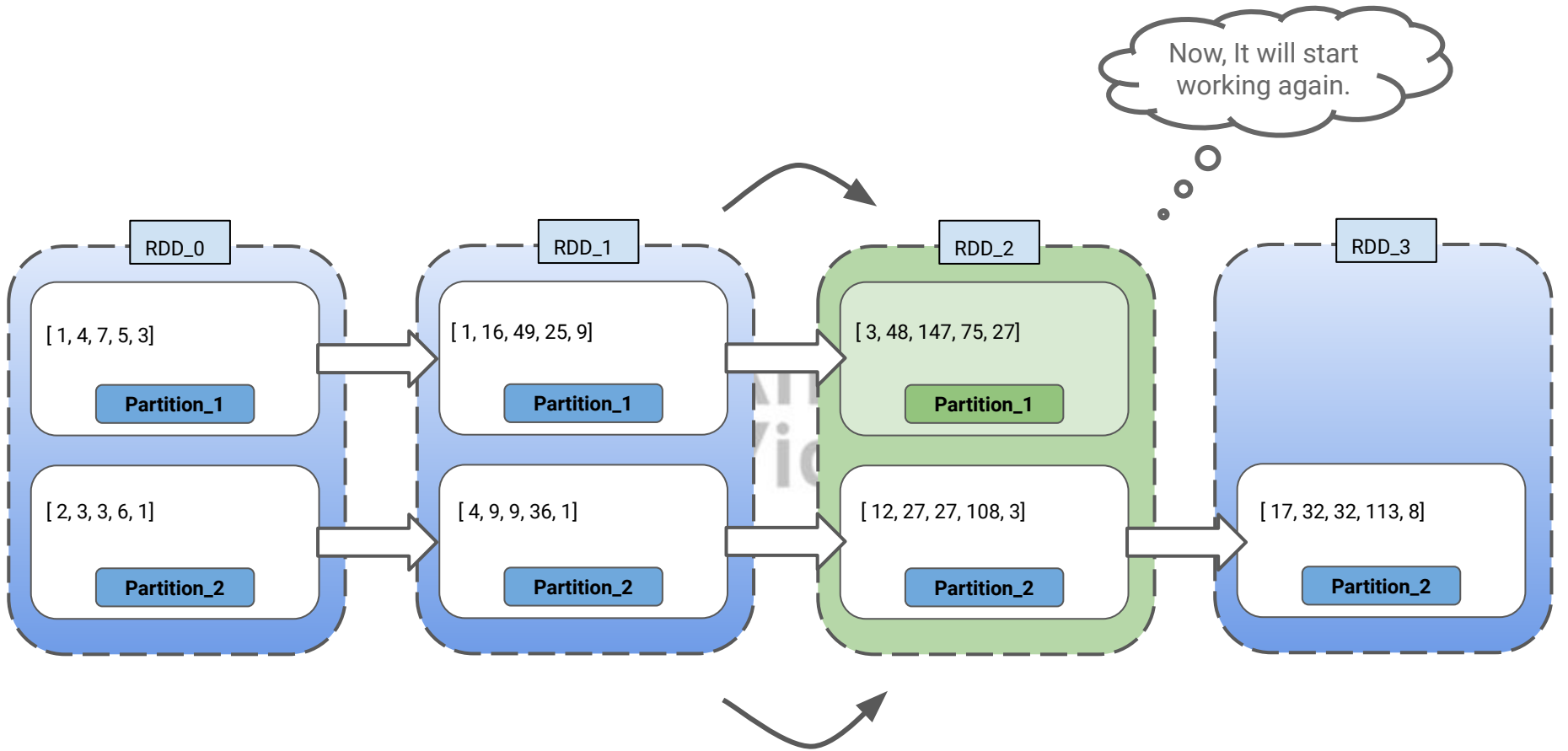
Meanwhile other partitions  
will continue their work for  
the next operation.

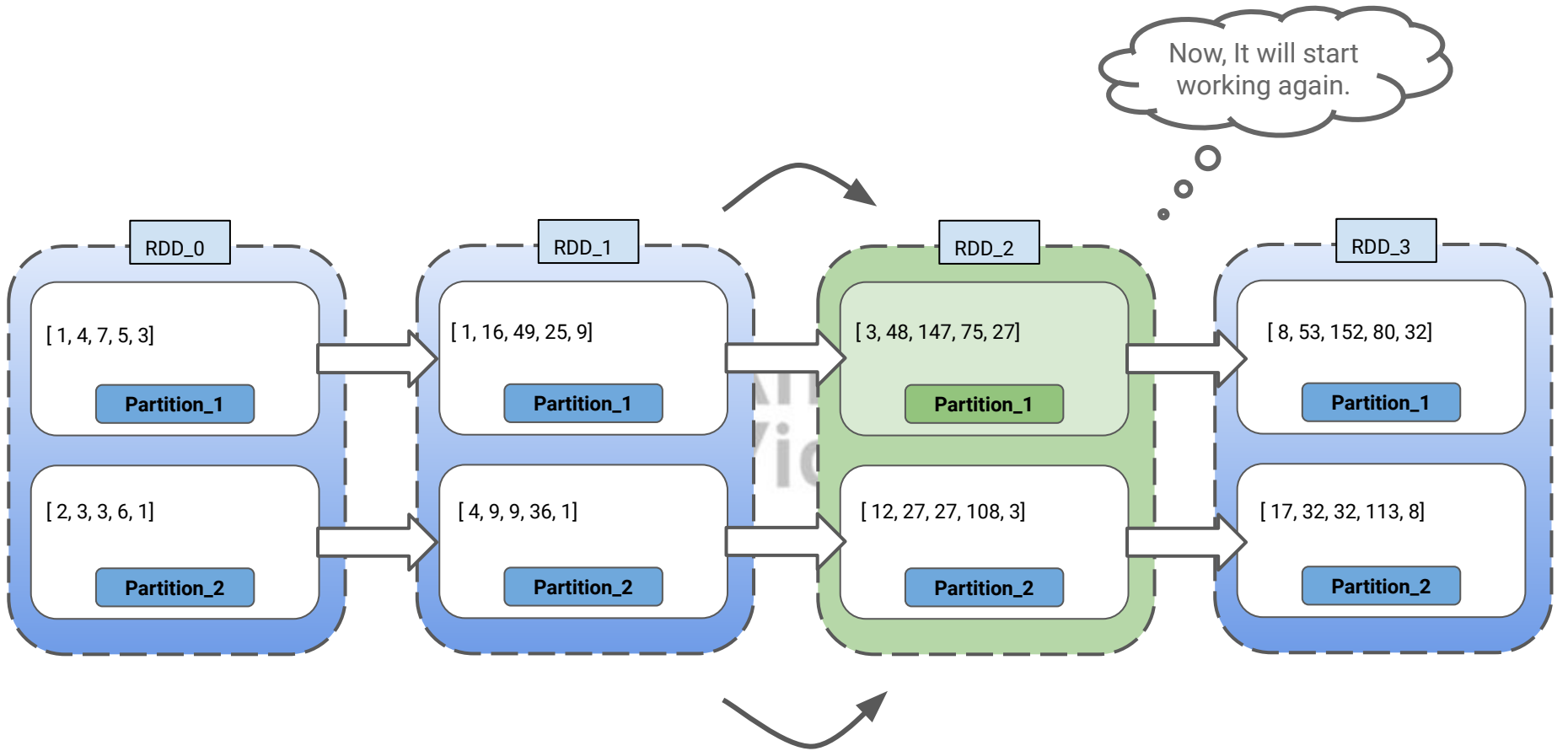


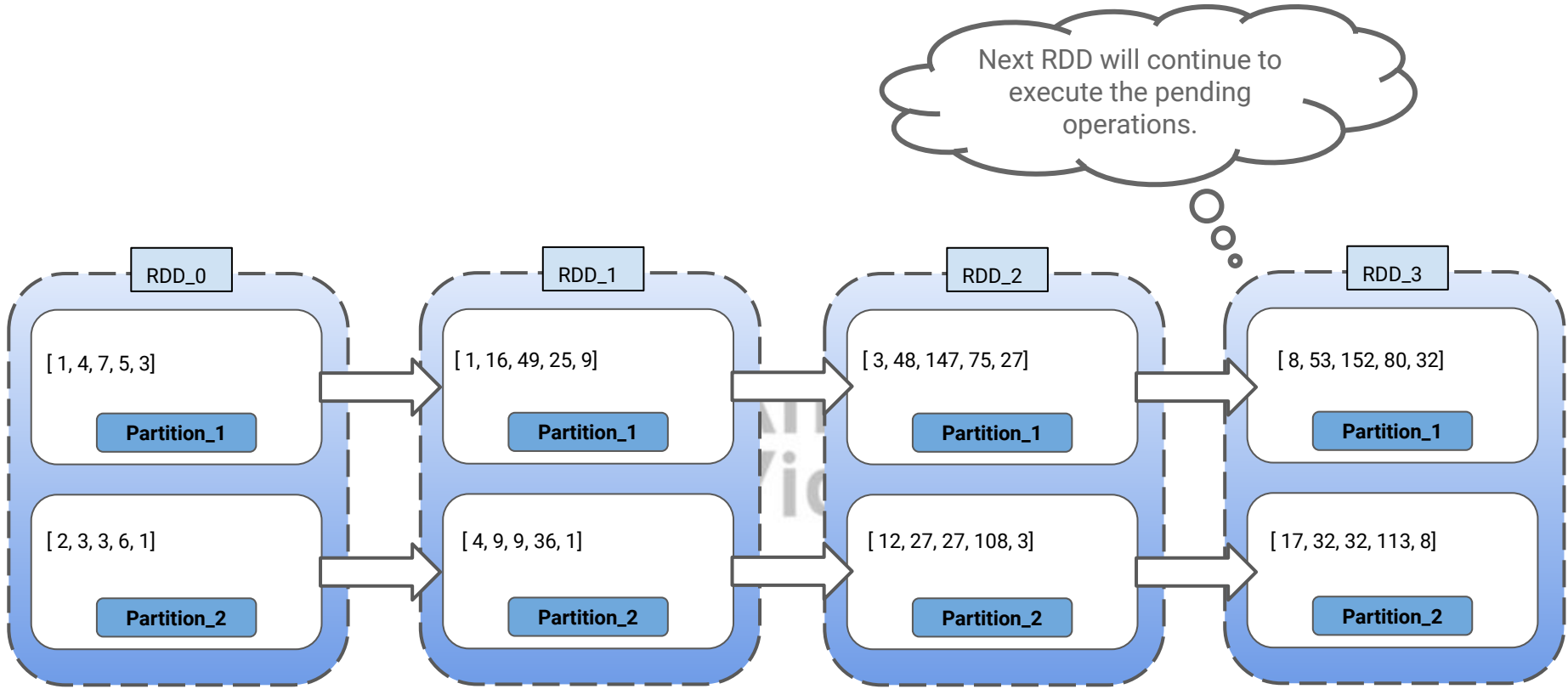














Thank You!!