# Big Data Technologies

Trainer: Mr. Nilesh Ghule.

# Zookeeper

Distributed coordination & Sync service.

Singlenode kafka Cluster

F — Quorum Peer
F — Quorum Peer
Leader — Quorum Peer
F — Quorum Peer
F — Quorum Peer

Zookeeper data (tree)

Kafka
Quorum Peer

Kafka
HBase
etc...

brokers

ids
seq
topics

0

iot   onsiot   Consumer offset

Spark Core
① RDD & DAG execution
② Cluster (Master/Worker)

RDD: Resilient Distributed Dataset
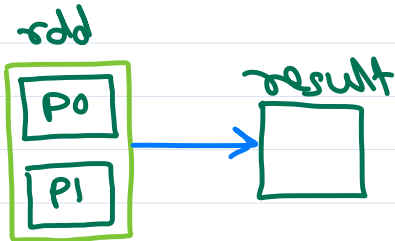↓ Can recover (Lineage)
↓ Split into Partitions on multiple worker nodes
↓ in-mem

↳ Logical entity
Physical entity

RDD operations
① narrow transformations :
② wide transformations :
③ actions :

rdd
P0
P1
→ result
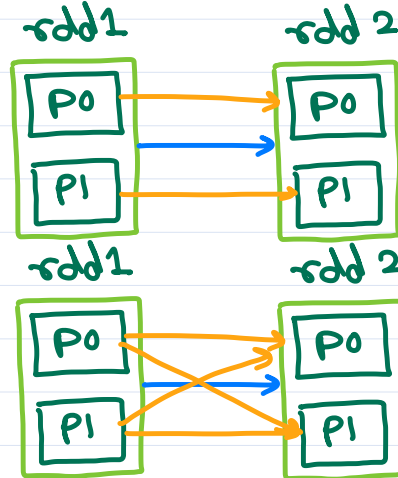
RDD characteristics
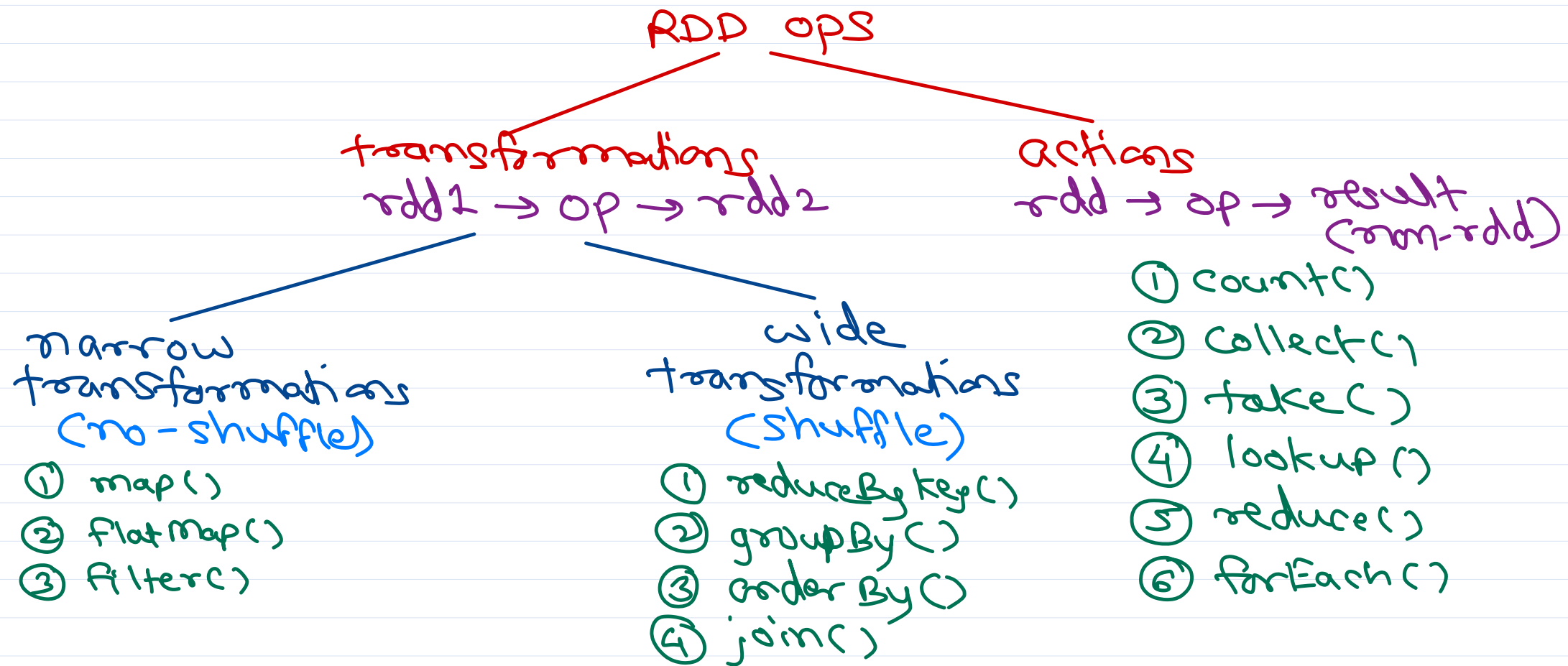① immutable
② resilient
③ lazy evaluated

RDD.programming
= Functional programming

Basics of fn prog:
① immutability
② location transparency

rdd1    rdd 2
P0  →  P0
P1  →  P1

rdd1    rdd 2
P0  →  P0
P1  →  P1

# RDD operations

RDD OPS

transformations
rdd1 → op → rdd2

narrow
transformations
(no-shuffle)
① map()
② flatMap()
③ filter()

wide
transformations
(shuffle)
① reduceByKey()
② groupBy()
③ orderBy()
④ join()

actions
rdd → op → result
(non-rdd)
① count()
② collect()
③ take()
④ lookup()
⑤ reduce()
⑥ forEach()

# Spark RDD - WordCount application in Scala



Directed Acyclic Graph: V=RDD, E=Operations

**Stage 0** | **Stage 1**

textFile

reduceByKey

Group By Key and Reduce all Values. e.g. Sum of all 1s.

wordcounts
RDD[(String,Int)] - KV RDD
ShuffleRDD.

map

capwordcounts
RDD[(String,Int)] - KV RDD

file
RDD[String]

collect()

result
Array[(String,Int)]

lines
RDD[String]

map

flatMap

Key   Value

words
RDD[String]

map

word1s
RDD[(String,Int)] <- Key-Value RDD
Each element is Tuple of 2 Elems.

---

```scala
val file=sc.textFile("/home/nilesh/
setup/bigdata/spark-3.3.1-bin-hadoop3/
LICENSE")

val lines=file.map(line=>line.
toLowerCase())

val words=lines.flatMap(line=>line.
split("[^a-z]"))

val word1s=words.map(word=>(word,1))

val wordcounts=word1s.reduceByKey((a,x)
=>a+x)

val capwordcounts=wordcounts.map(wc=>
(wc._1.toUpperCase(), wc._2))

val result=capwordcounts.collect()

capwordcounts.toDebugString
```

transforms

action.

# Thank you!

*Nilesh Ghule <nilesh@sunbeaminfo.com>*