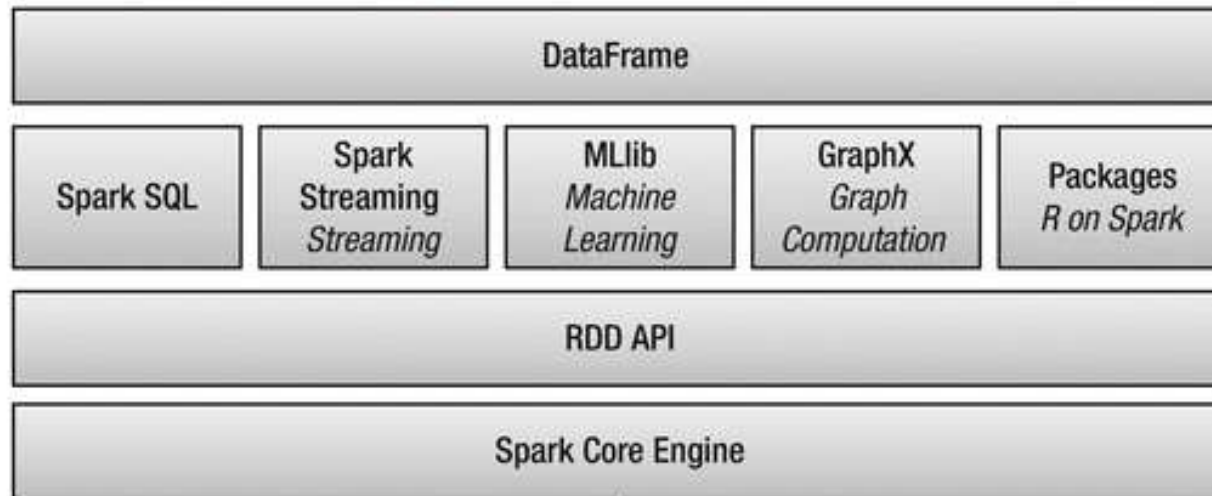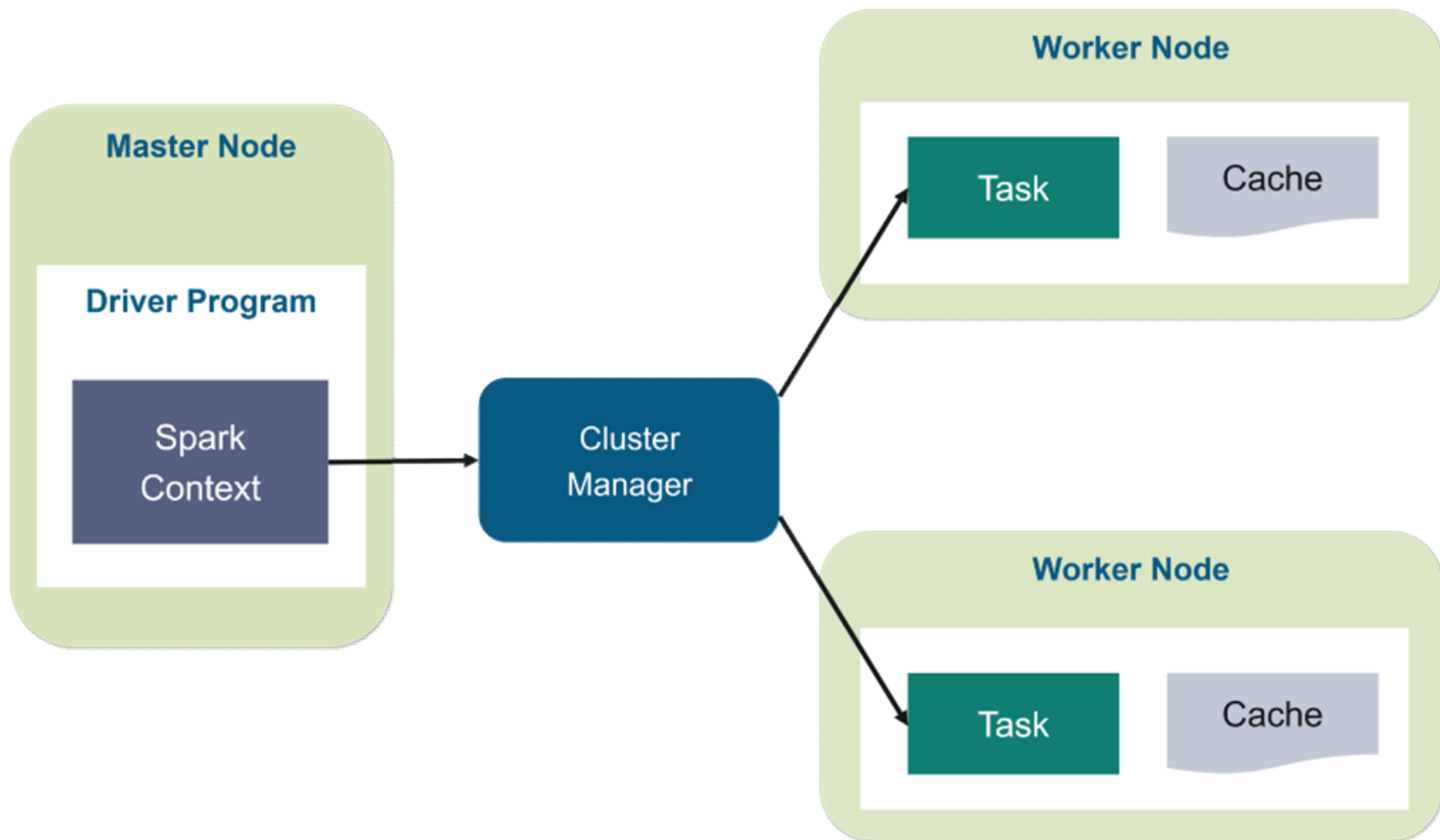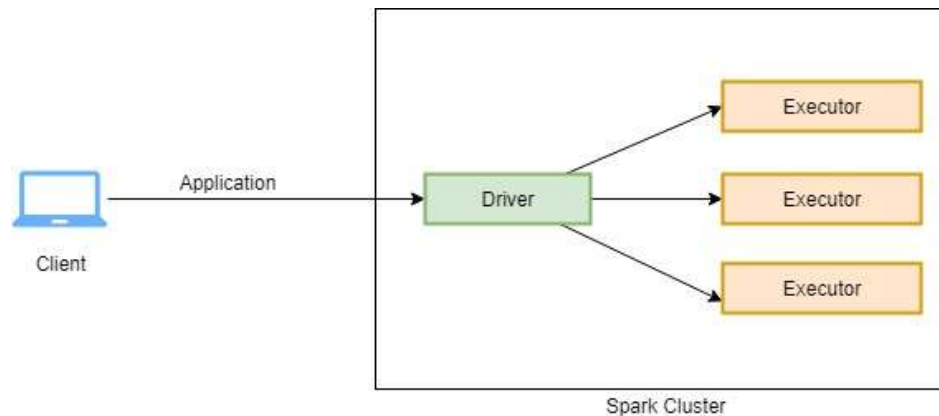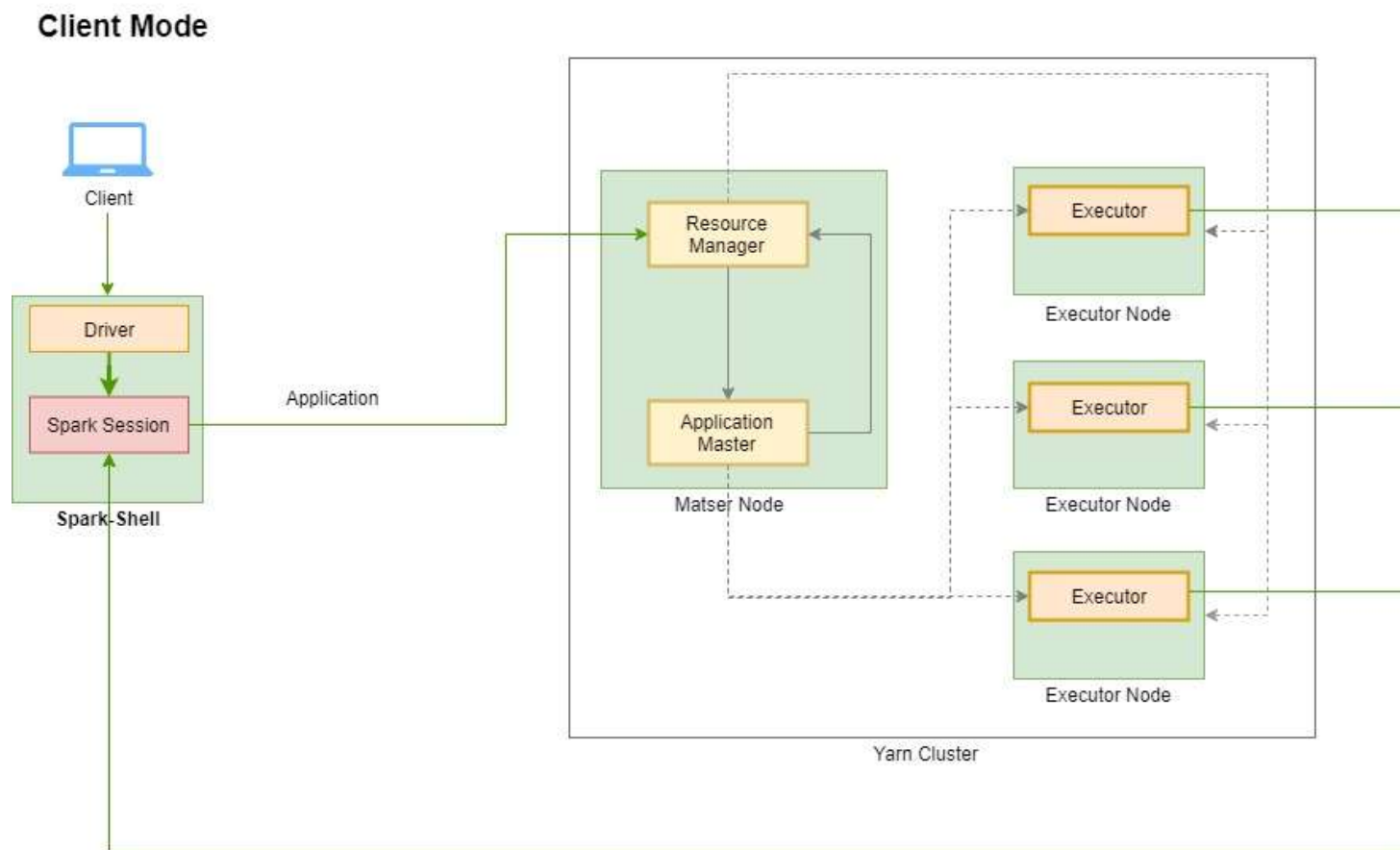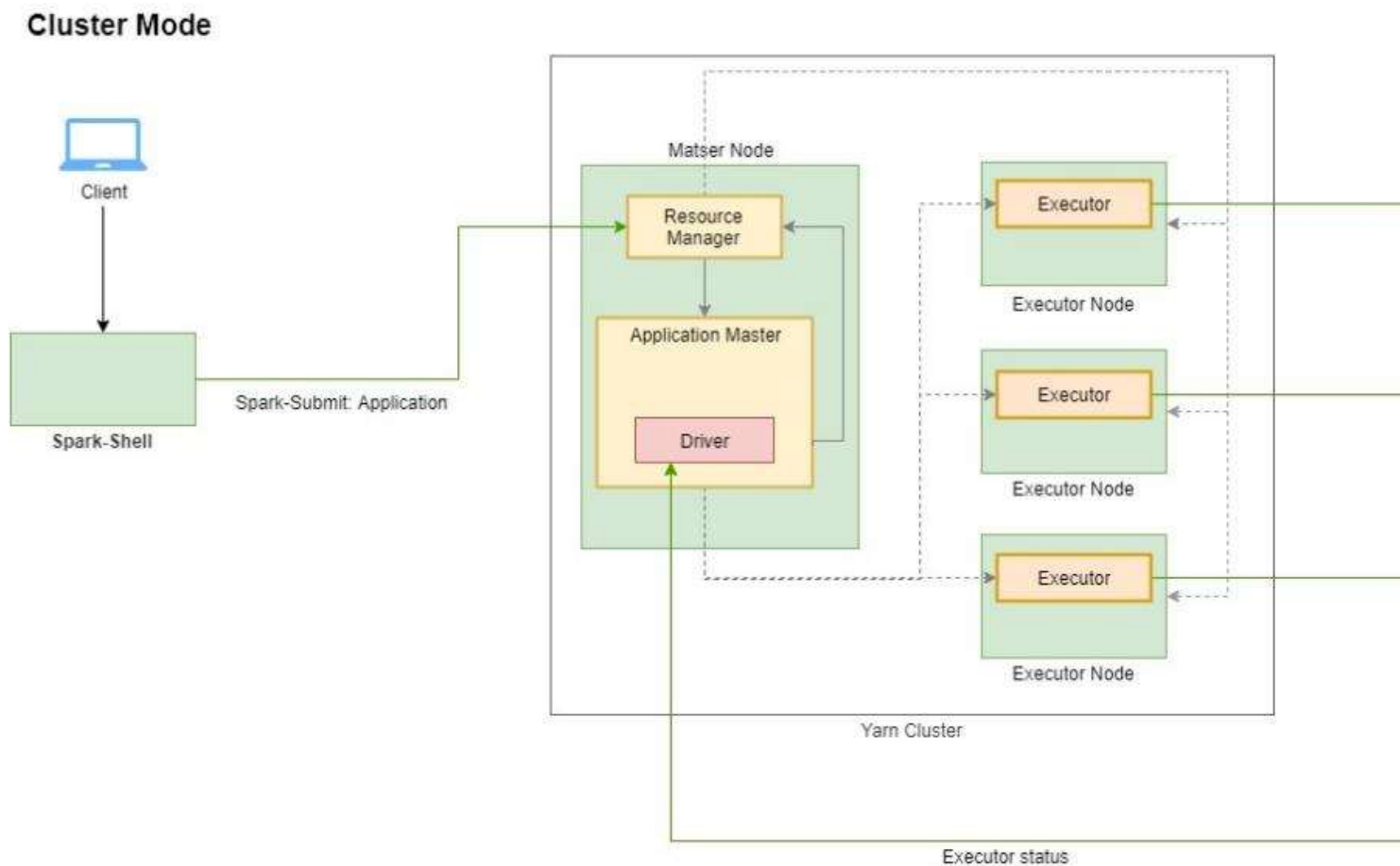# Apache Sparks Diagrams

# Apache Spark Execution

- For every application submitted on spark cluster spark creates a dedicated Driver process and bunch of Executor processes.

- Driver process is responsible for analyzing, distributing, scheduling and monitoring of executor processes.

- Whereas the executor process is only responsible for running the task they were assigned by drivers and reporting the status back to the driver.

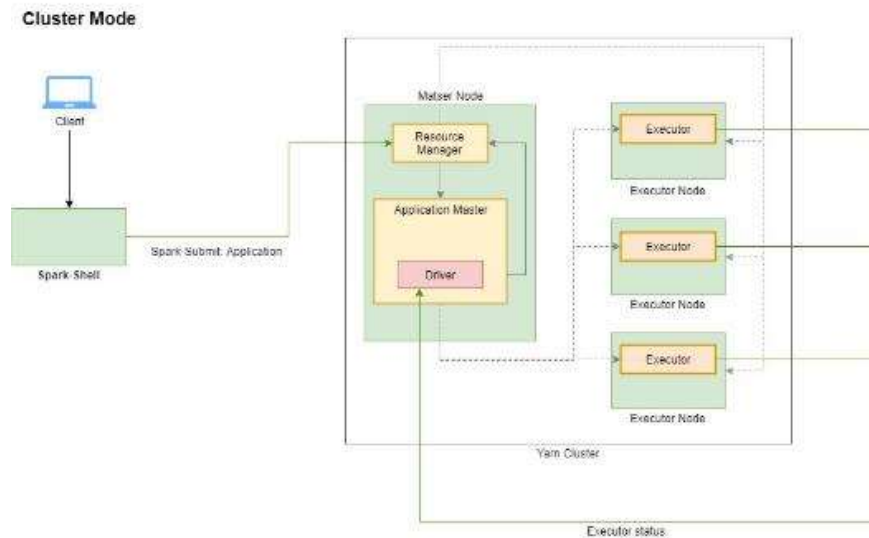# Apache Spark Execution – Client Mode

# Apache Spark Execution – Cluster Mode

**Cluster Mode**
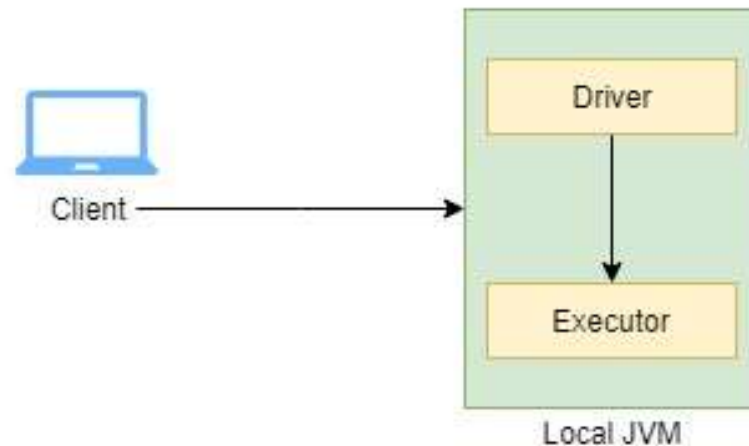
# Apache Spark Execution – Cluster Mode

- Spark client submits the packed application to yarn request manager.
- Yarn request manager then creates the AM container. The Spark-Driver is also created in AM container.
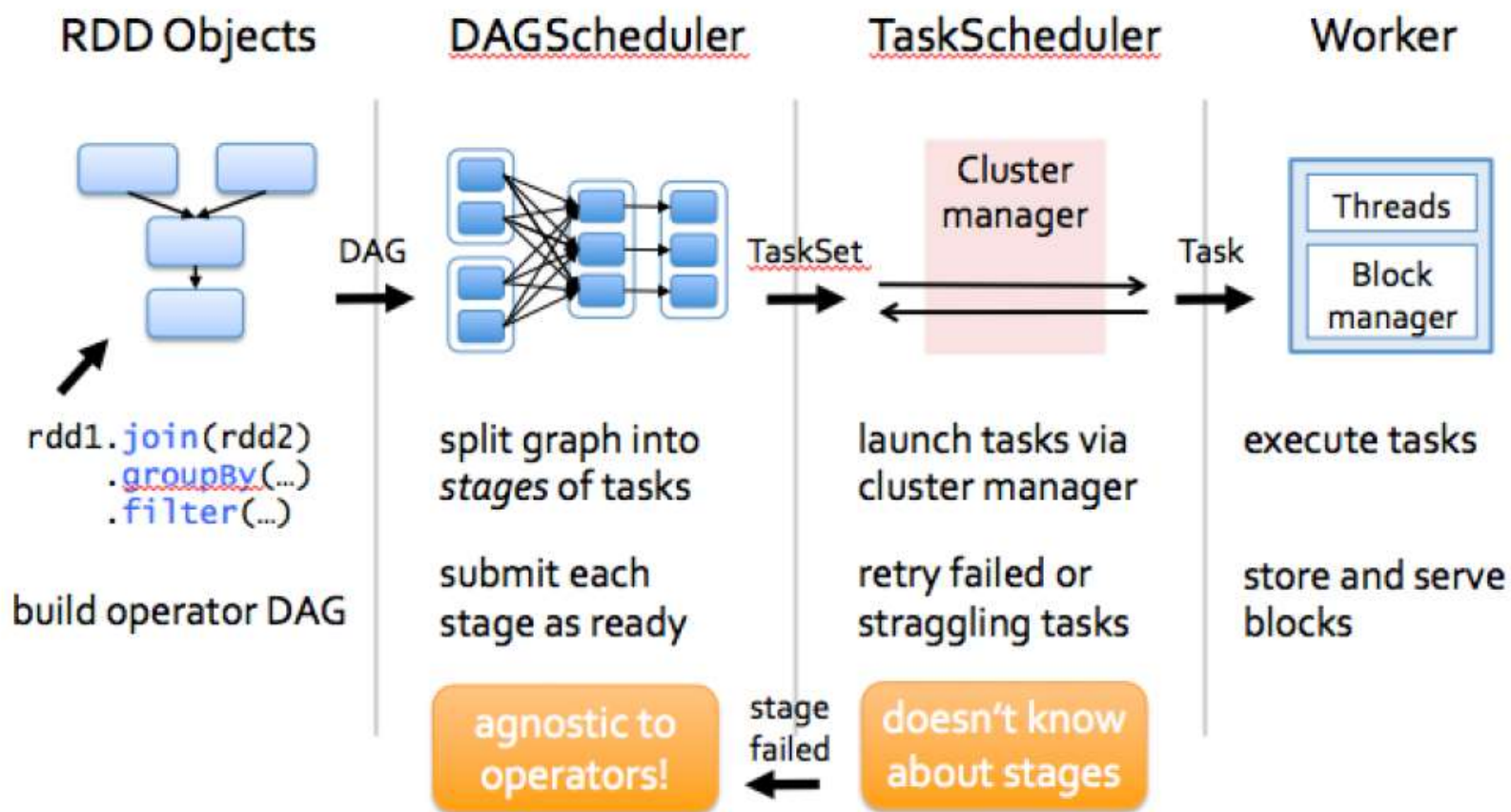- Rest of the flow is same as mentioned in the client mode.
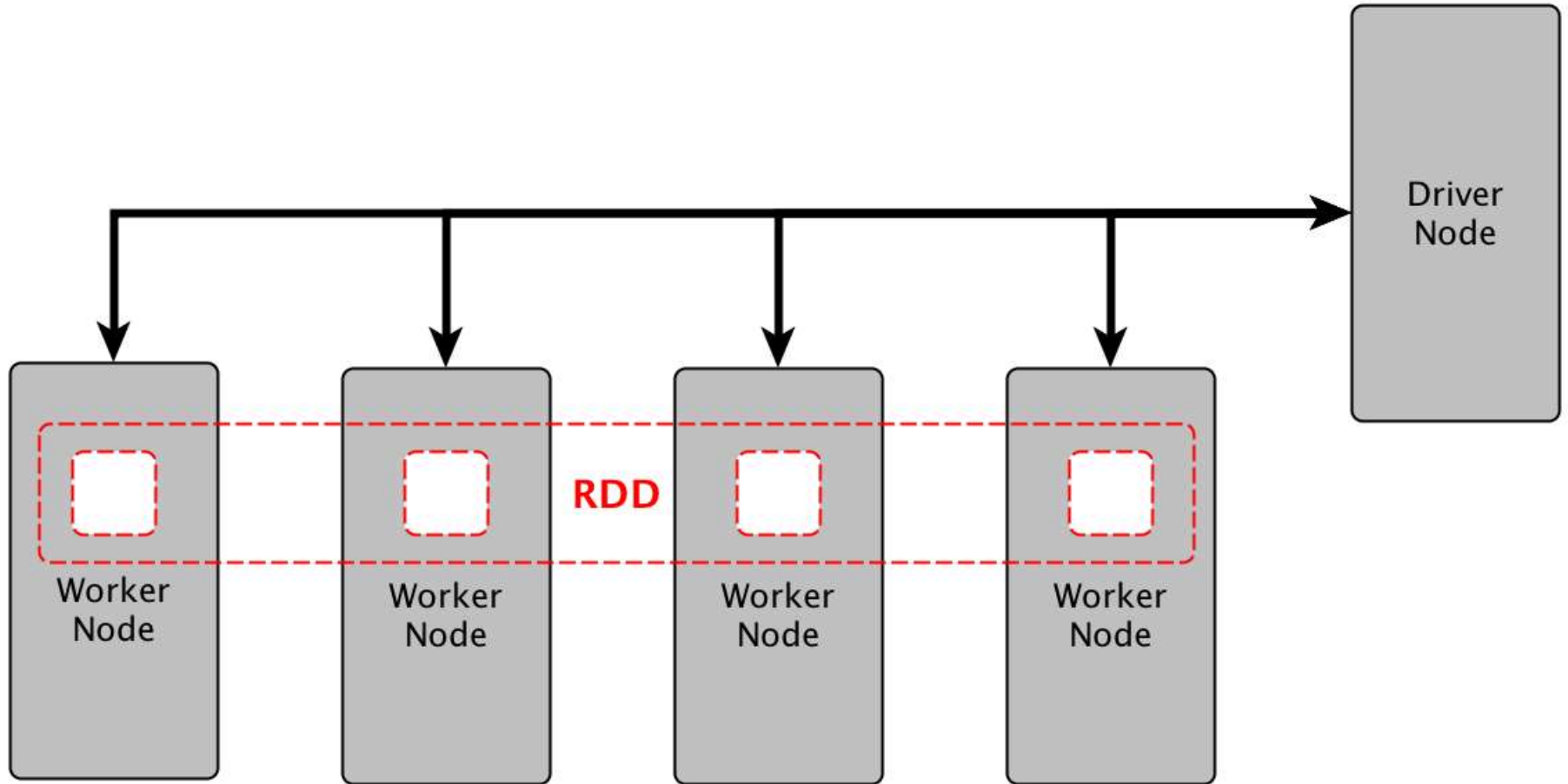
# Apache Spark Execution – Local Mode

- There is another mode in which spark can be run locally without any cluster requirement.

- This mode is suitable for scenarios when we do not have enough resources to create cluster.

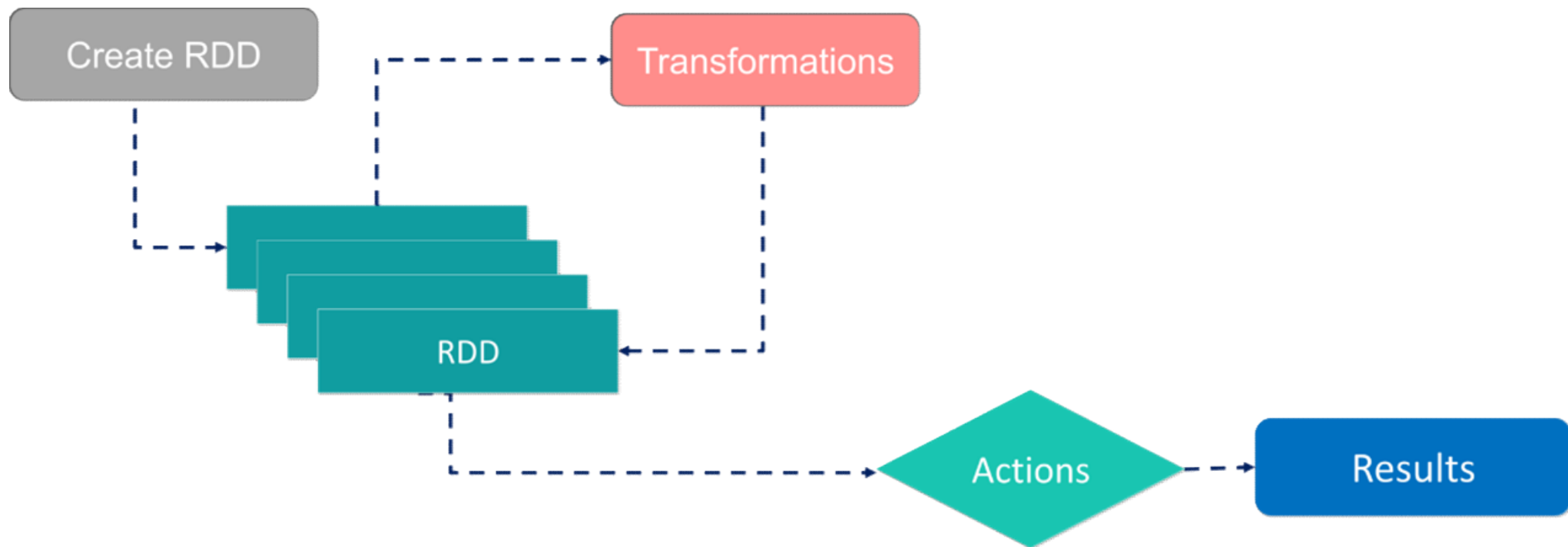- But in this mode you get only one executor and both the Driver and Excuter runs in the same JVM.

**Local Mode**

Client

Driver

Executor

Local JVM

# How Sparks work?



| RDD Objects | DAGScheduler | TaskScheduler | Worker |
|---|---|---|---|

DAG → TaskSet → Cluster manager → Task

```
rdd1.join(rdd2)
    .groupBy(...)
    .filter(...)
```

build operator DAG

split graph into *stages* of tasks

submit each stage as ready

launch tasks via cluster manager

retry failed or straggling tasks

execute tasks

store and serve blocks

Threads

Block manager

agnostic to operators!

stage failed

doesn't know about stages
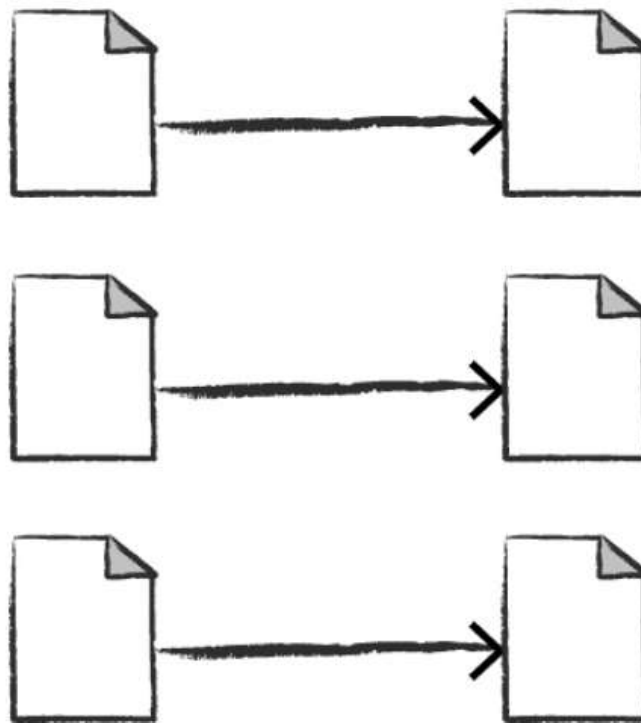
# Narrow transformations
## 1 to 1

Figure 2-4. A narrow dependency
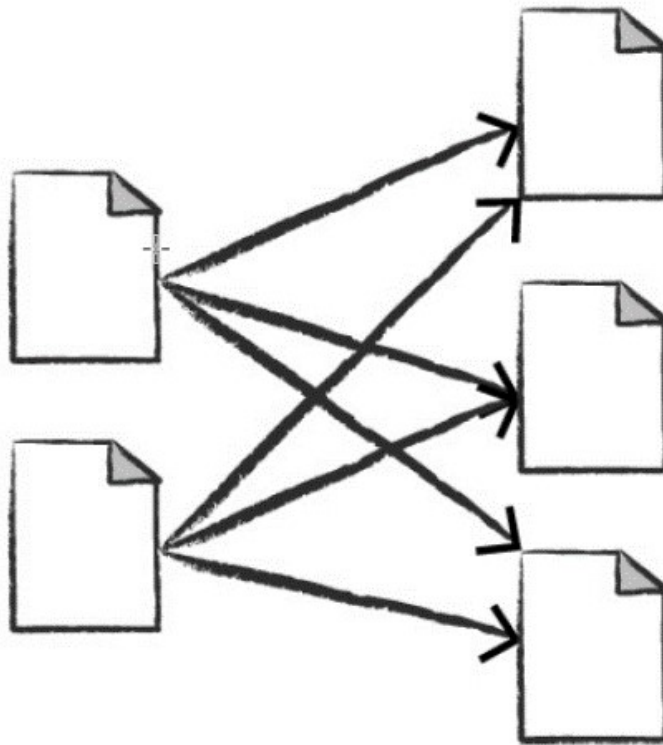
Wide transformations
(shuffles) 1 to N

Figure 2-7. Reading a CSV file into a DataFrame and converting it to a local array or list of rows
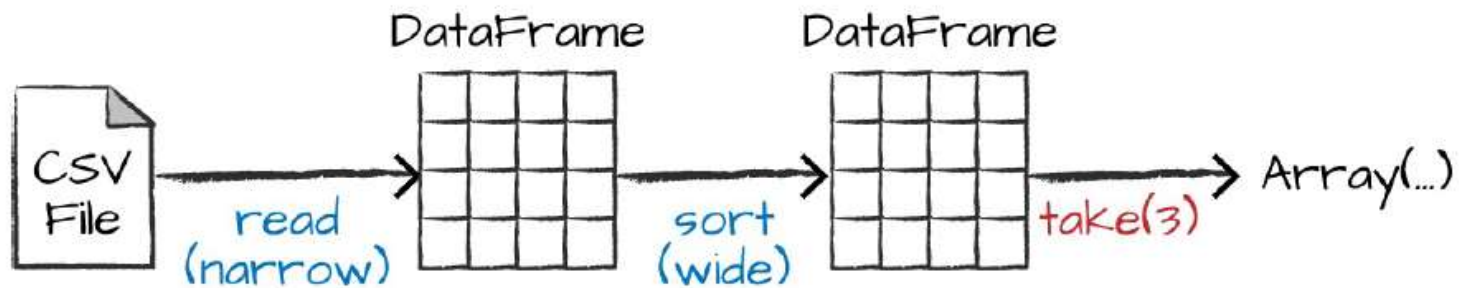


Figure 2-8. Reading, sorting, and collecting a DataFrame

# Catalyst Optimizer

- Spark SQL uses an optimizer called catalyst to optimize all the queries

- This optimizer makes queries run much faster

- An optimizer automatically finds out the most efficient plan to execute data operations specified in the user's program.

- logical plan — series of algebraic or language constructs, as for example: SELECT, GROUP BY or UNION keywords in SQL. It's usually represented as a tree.

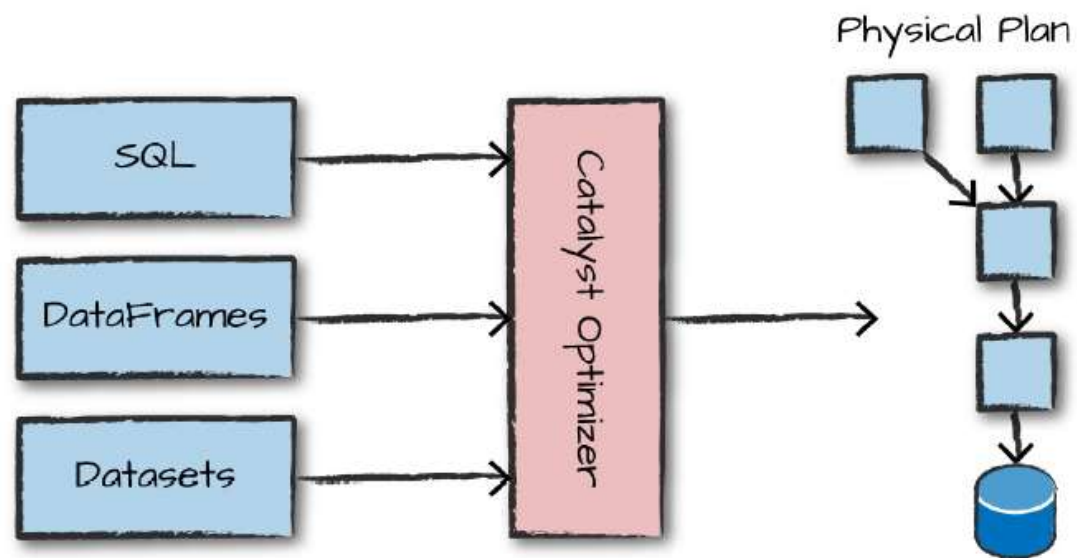- physical plan — Concerns low level operations.

# Catalyst Optimizer



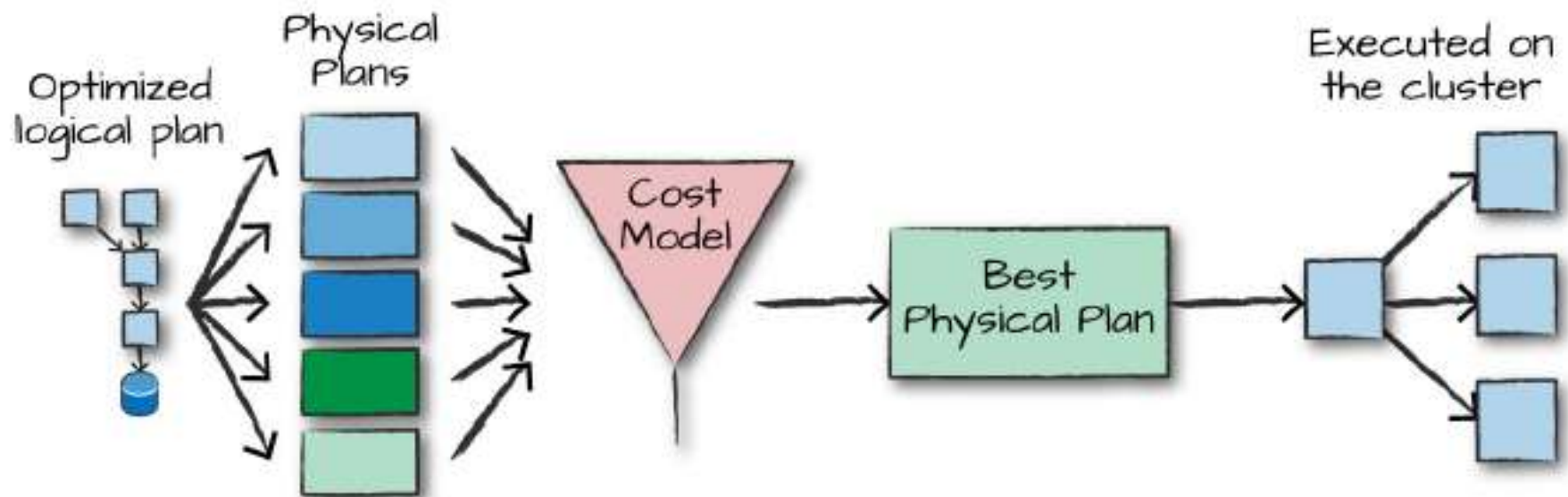Figure 4-1. The Catalyst Optimizer

# Catalyst Optimizer
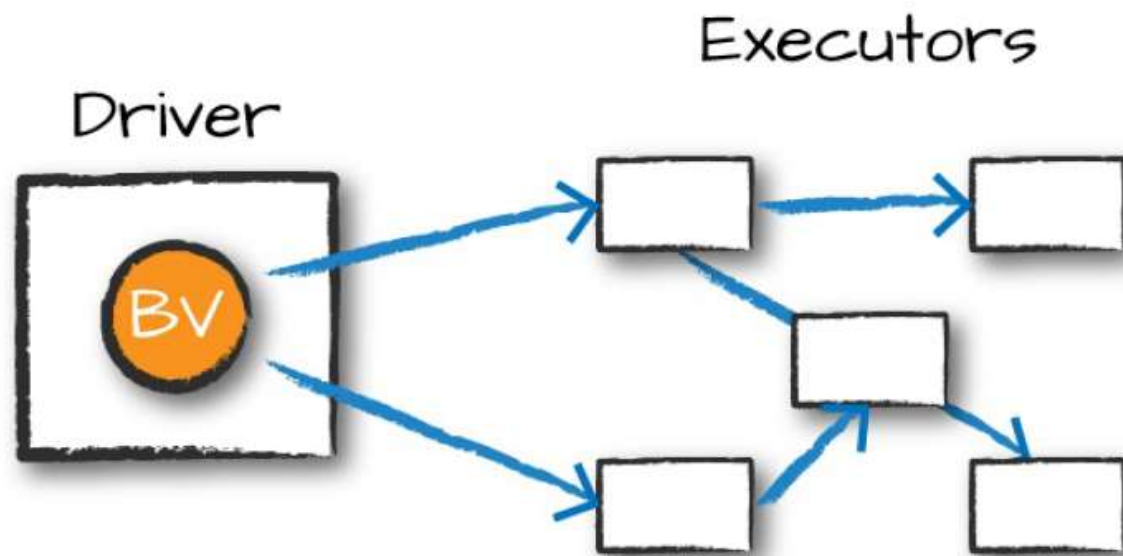


Figure 4-3. The physical planning process
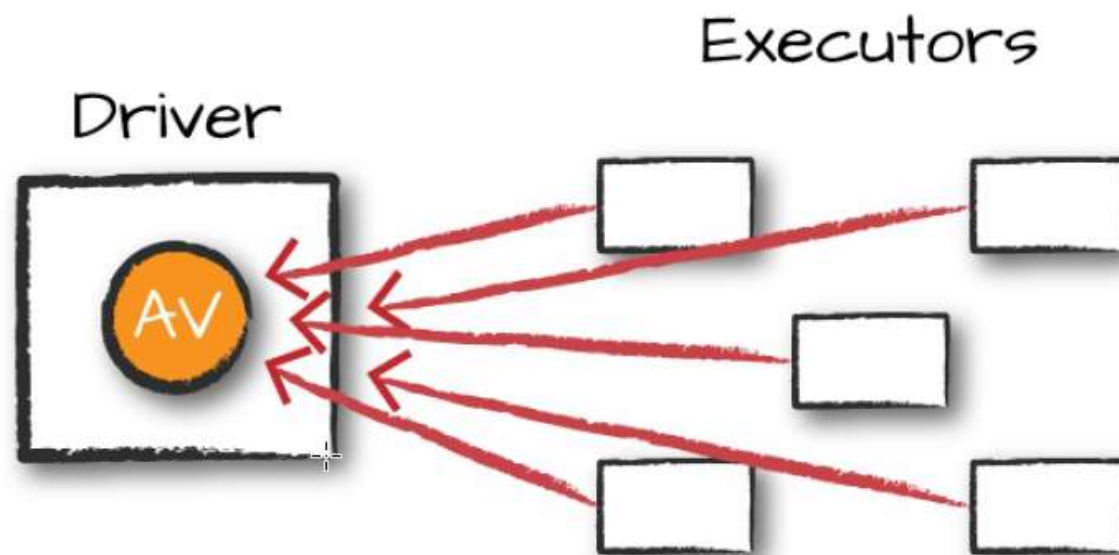
Executors

Driver



Figure 14-1. Broadcast variables

Figure 14-2. Accumulator variable

# Thanks