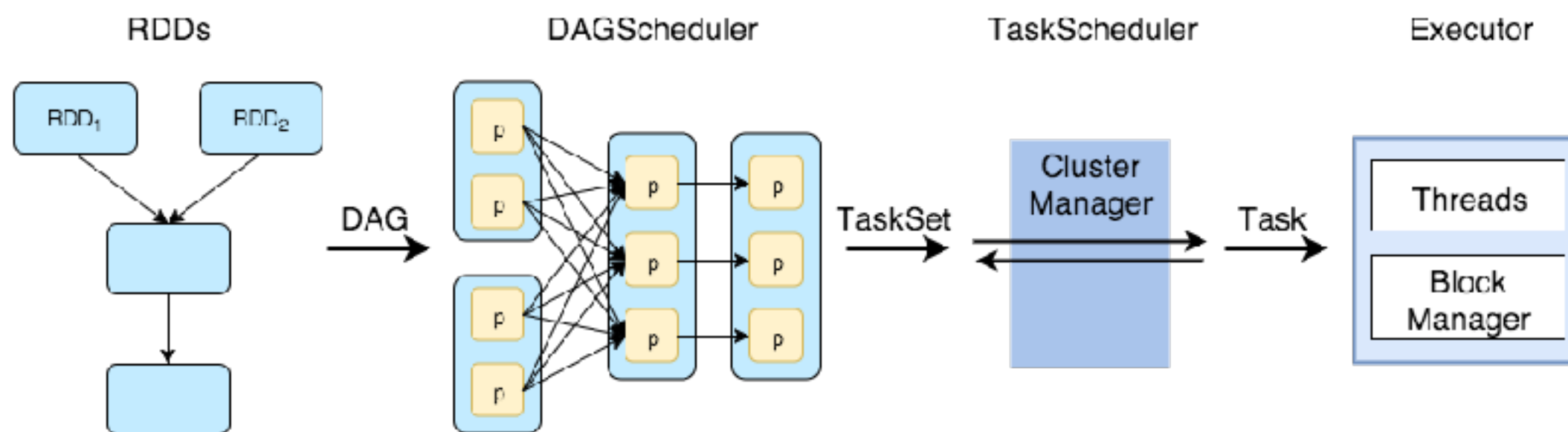


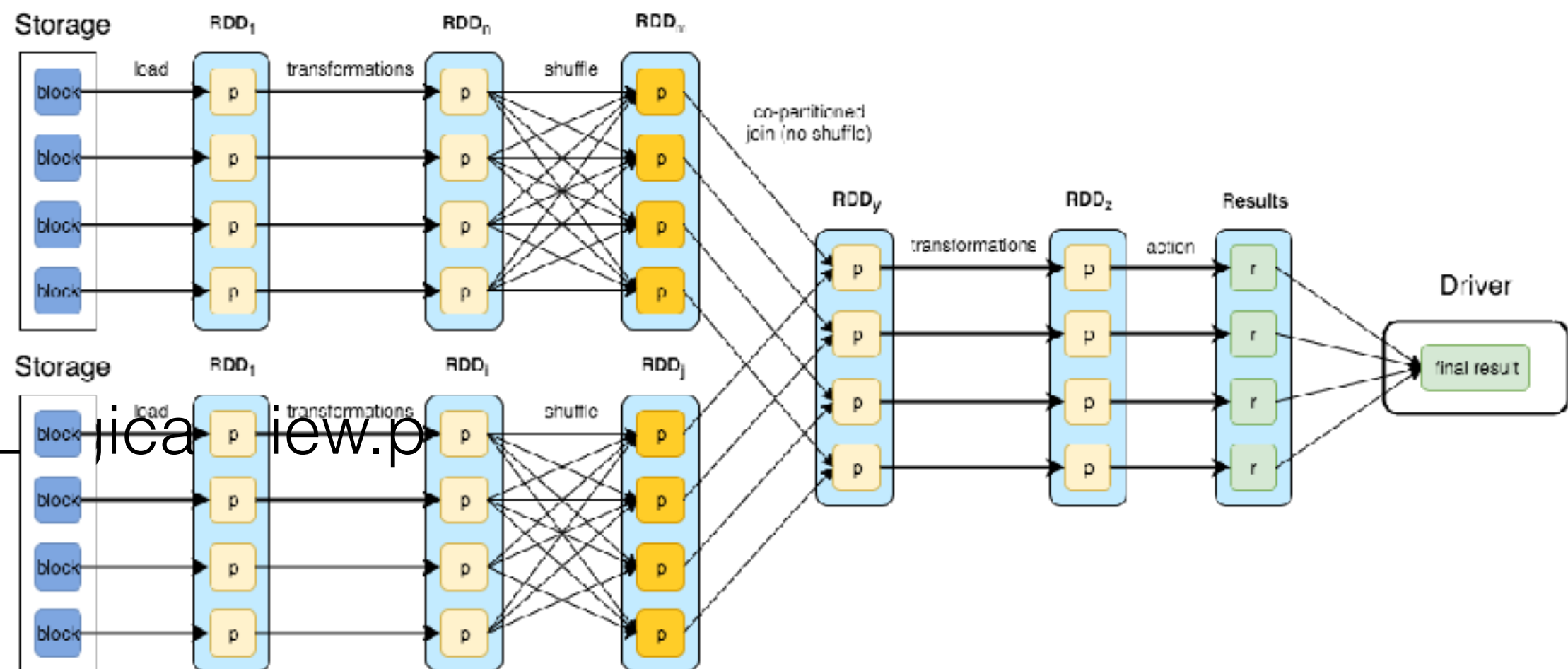
Apache Spark Internals

Ladle Patel

Job Flow



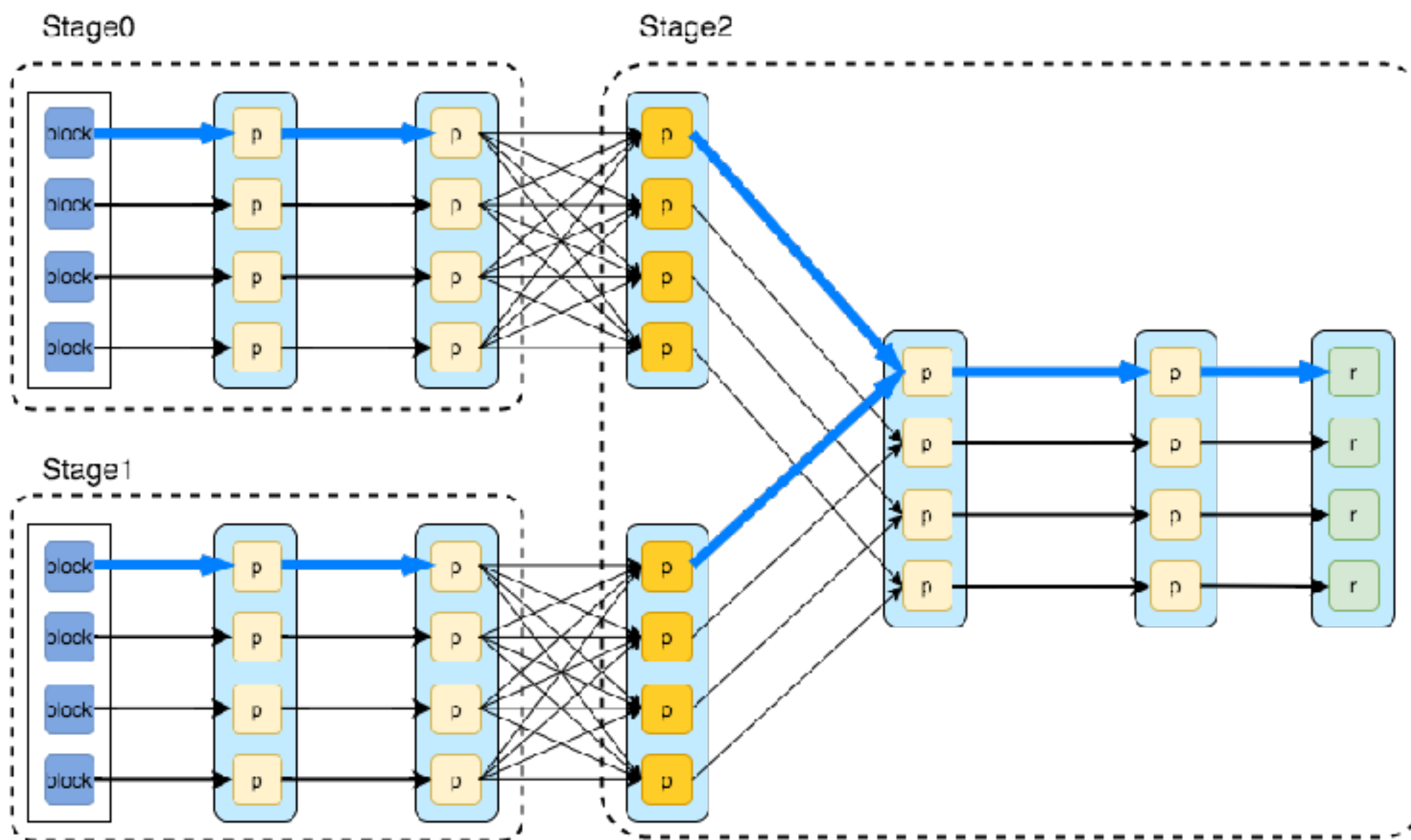
DAG(Logical plan)



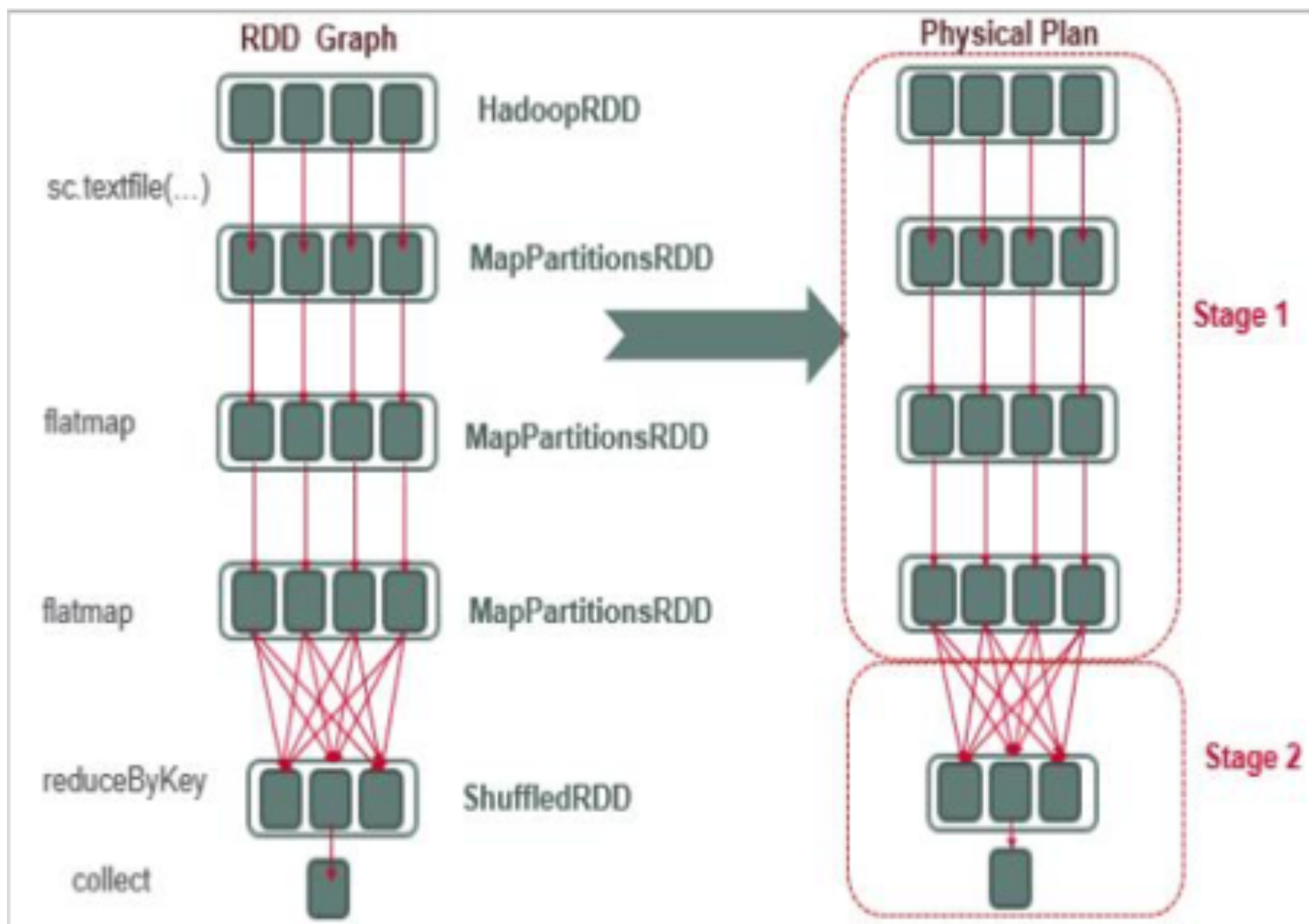
- Logical plan view.ppt

Splitting DAG into Stages

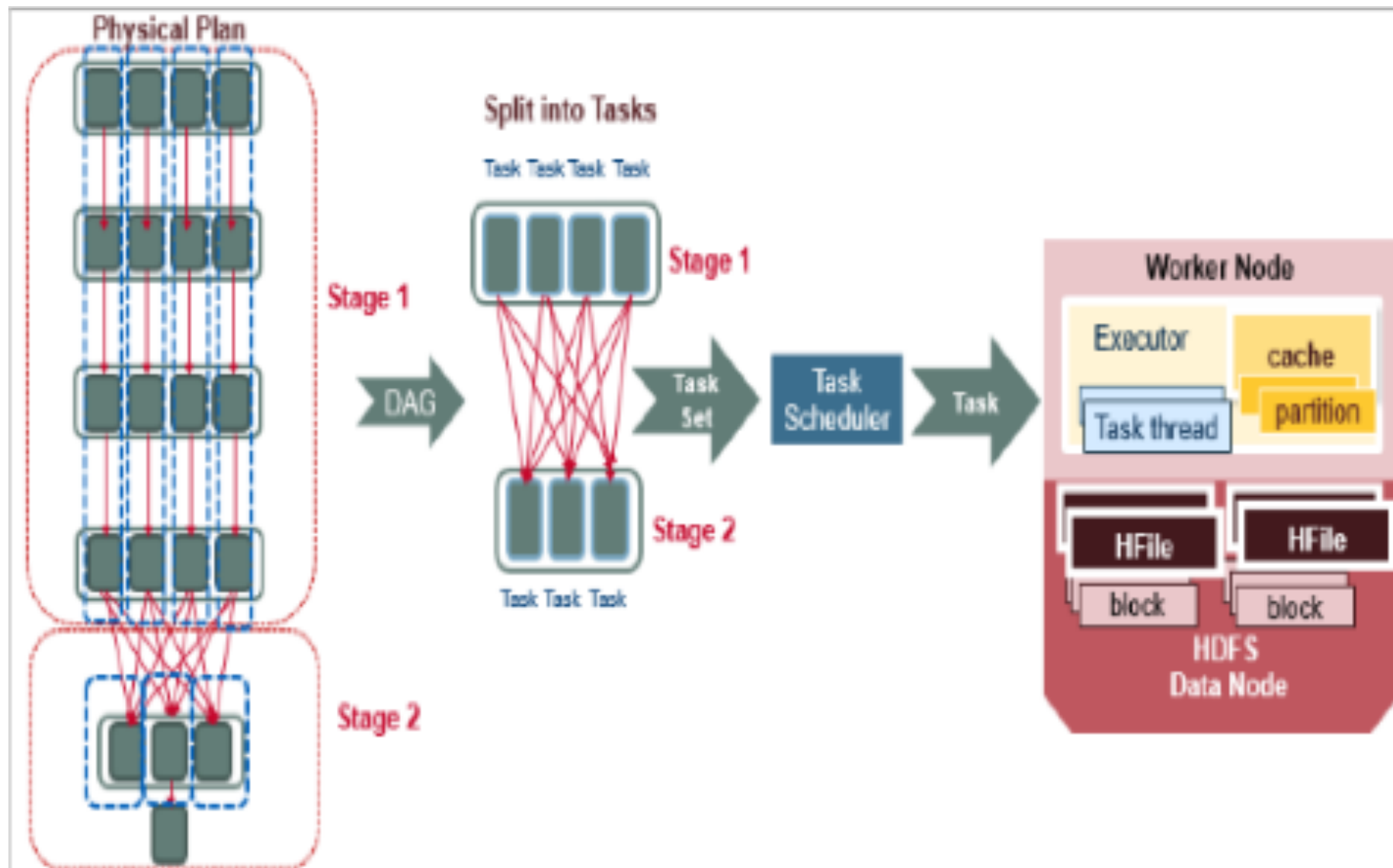
(Physical plan)



Logical and Physical plan



Logical and Physical plan Cont..



Directed Acyclic Graph

- Graph :Structure consisting of nodes, that are connected to each other with edges.
- Directed :The connections between the nodes (edges) have a direction: $A \rightarrow B$ is not the same as $B \rightarrow A$.
- Acyclic :Non-circular moving from node to node by following the edges, you will never encounter the same node for the second time.

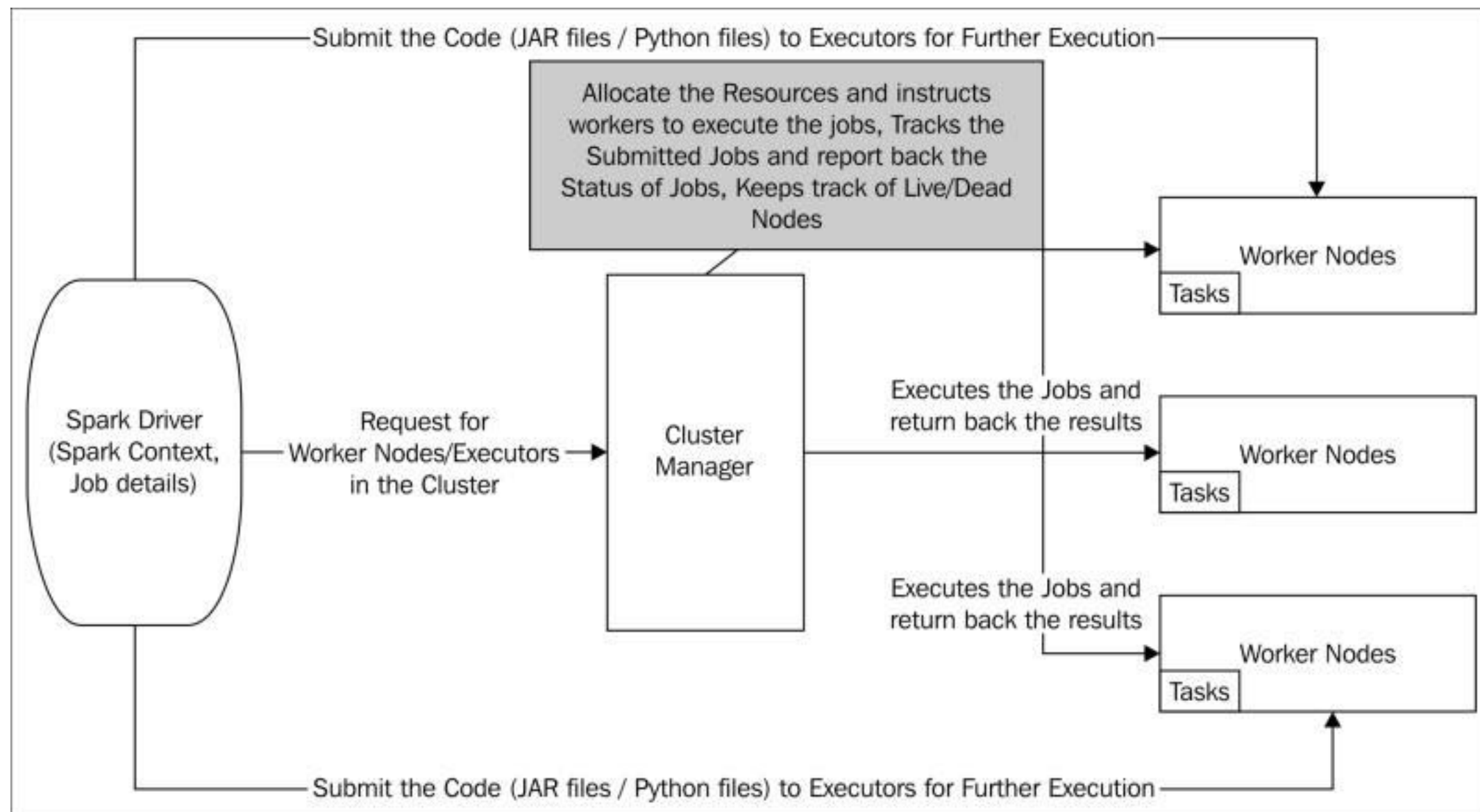
Narrow and Wide Transformation

- Narrow transformation : Transformation which doesn't require the data to be shuffled across the partitions. for example, Map, filter etc..
- Wide transformation : Transformation which requires the data to be shuffled for example, reduceByKey etc..

Terminologies

- DAGScheduler :Computes a DAG of stages for each job and submits them to TaskScheduler.
- TaskScheduler:Responsible for sending tasks to the cluster, running them, retrying if there are failures, and mitigating stragglers.
- BlockManager:Provides interfaces for putting and retrieving blocks both locally and remotely into various stores (memory, disk, and off-heap)
- Shuffle:The Transfer of Data between stages.

Jar Execution



For more details.

- <https://github.com/JerryLead/SparkInternals>
- <https://github.com/JerryLead/SparkLearning>



Thank You.