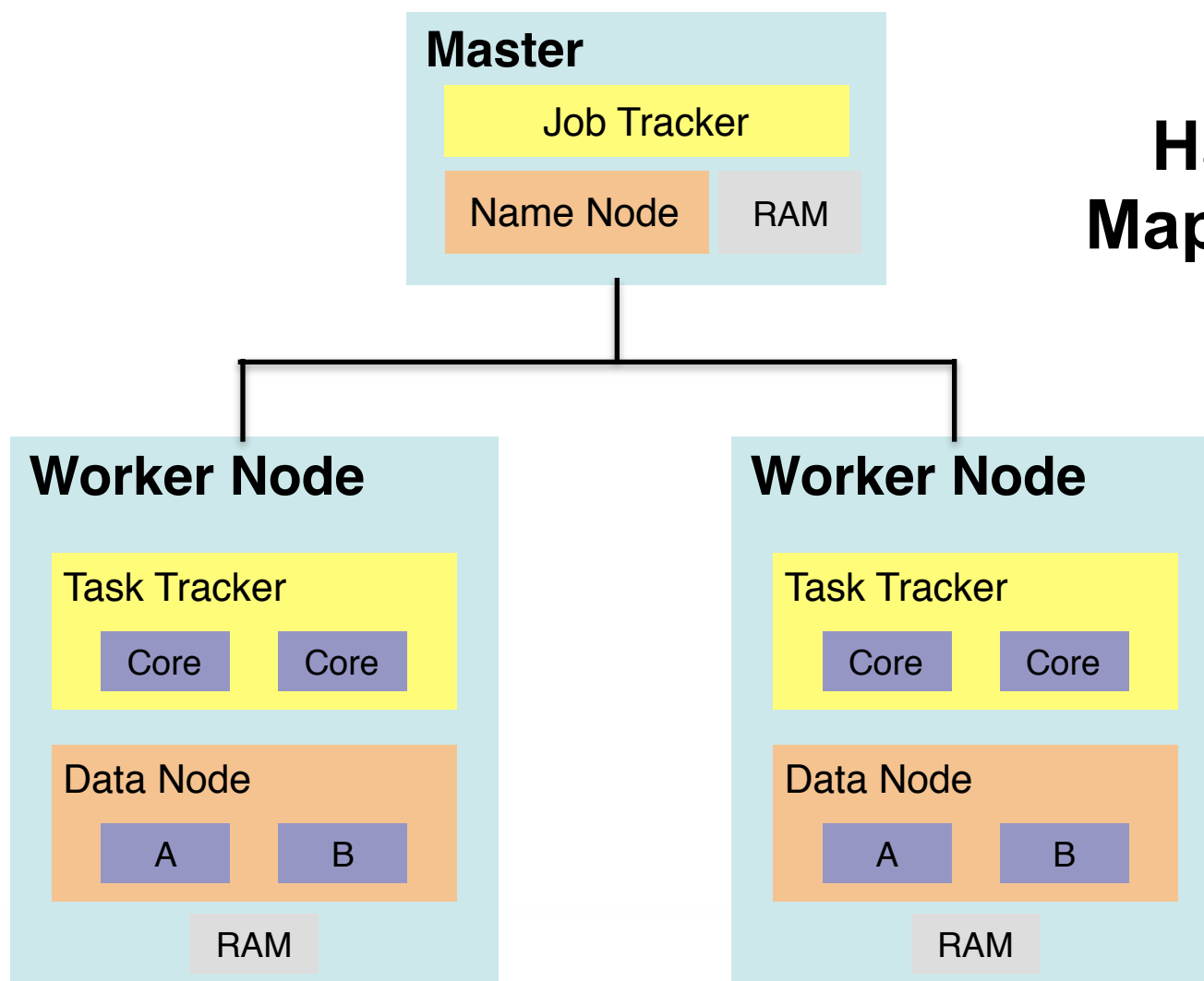


Lesson 4: Spark Internals

4.3 The Spark Execution Context

Hadoop MapReduce



Spark Execution Context

Laptop

Driver Program

SparkContext

Cluster Manager

Standalone
YARN
Mesos

Cluster

Worker Node

Executor

cache

Task

Task

Worker Node

Executor

cache

Task

Task



Terminology

<i>Term</i>	<i>Meaning</i>
Driver	<i>Process that contains the SparkContext</i>
Executor	<i>Process that executes one or more Spark tasks</i>
Master	<i>Process that manages applications across the cluster</i>
Worker	<i>Process that manages executors on a particular node</i>



Spark vs. Hadoop

- Spark only replaces MapReduce (**computation**)
- Still need a **data store**: HDFS, HBase, Hive, etc.
- Spark has a more **flexible/general** programming model
- Spark often faster for **iterative** computation

