

Apache Spark

Framework de computo distribuido

MLlib

Machine Learning

Streaming

Real-time analytics

SQL

Interactive Queries

GraphX

Graph processing

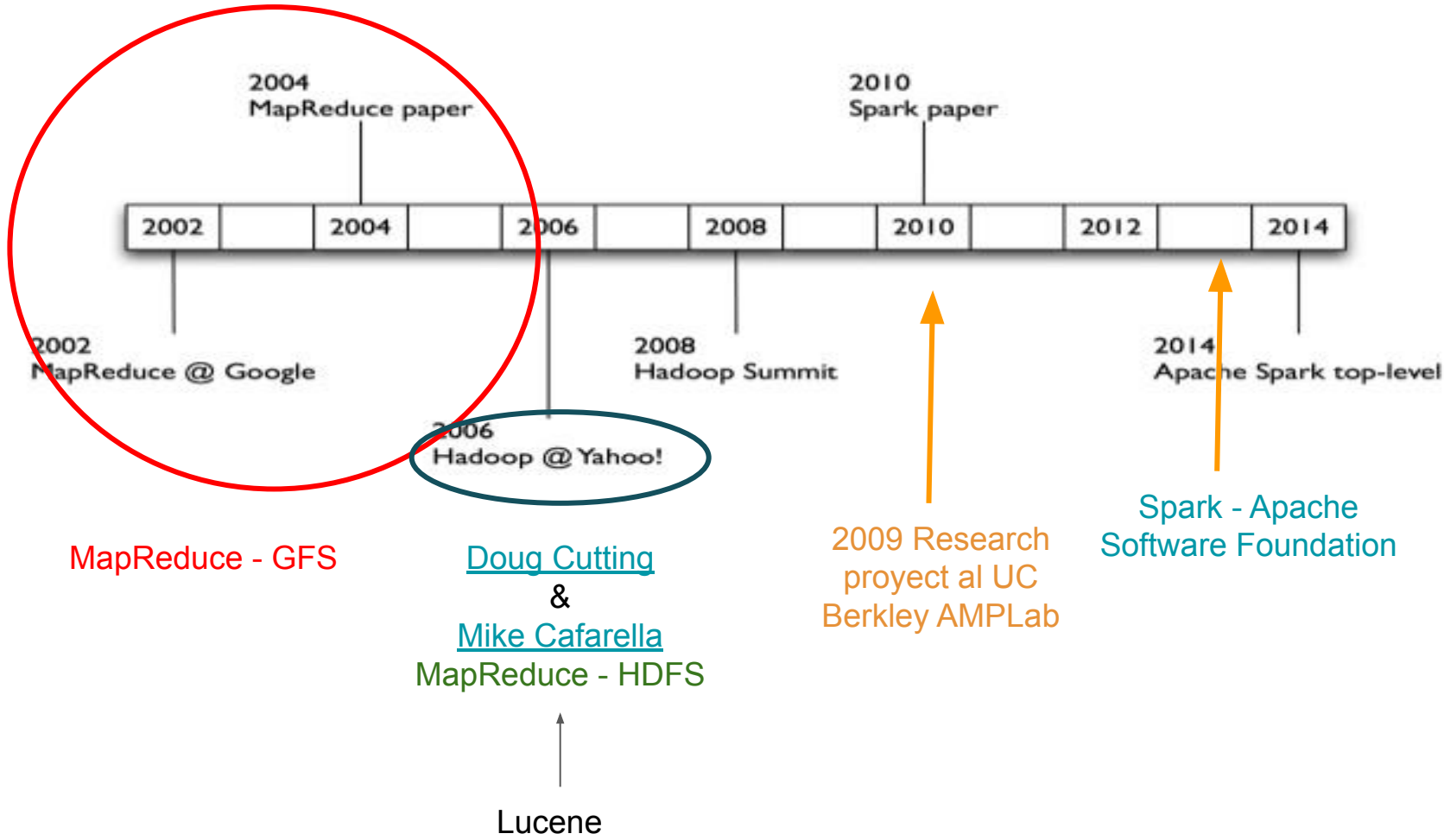
APACHE
Spark  **Core**

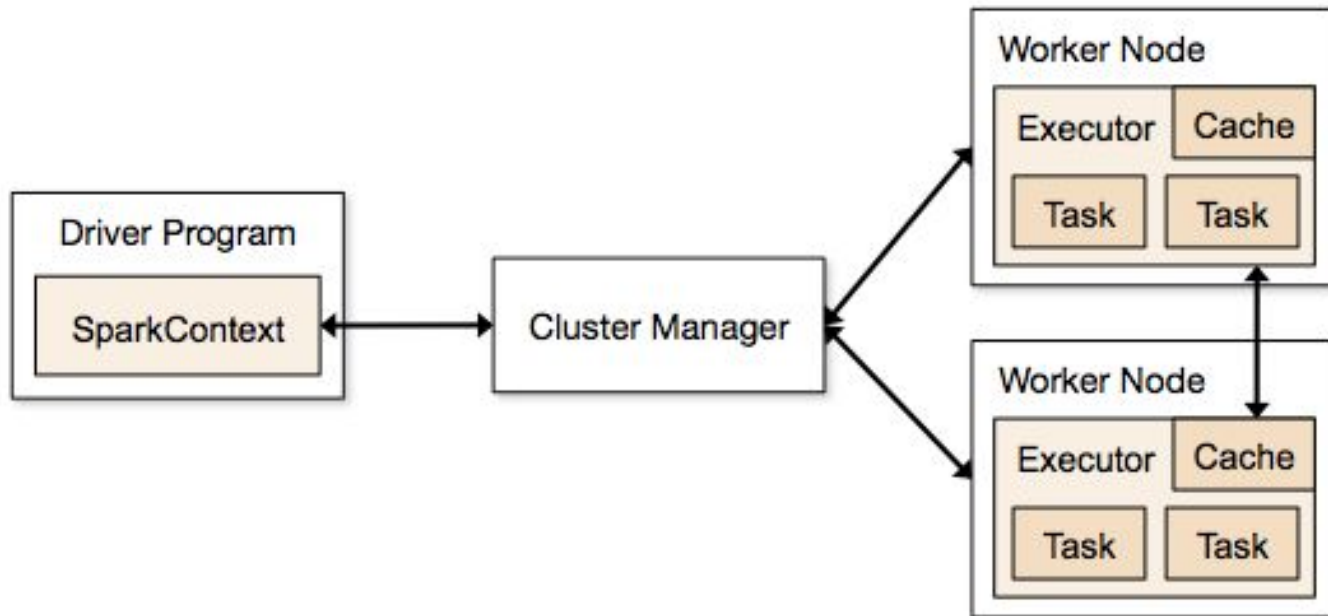
R

Python

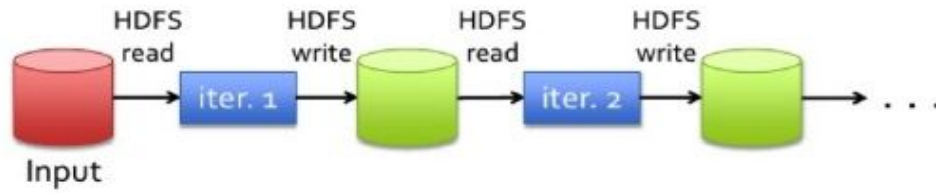
Scala

Java





- Master
- RDD
- Operaciones Lazy



<https://www.analyticsvidhya.com/wp-content/uploads/2016/09/big1.png>





Referencias

[¿por qué usar PySpark para grandes conjuntos de datos que exceden la memoria de la máquina de un solo nodo?](#)

[Comprehensive Introduction to Apache Spark, RDDs & Dataframes \(using PySpark\)](#)

[How to use a Machine Learning Model to Make Predictions on Streaming Data using PySpark](#)