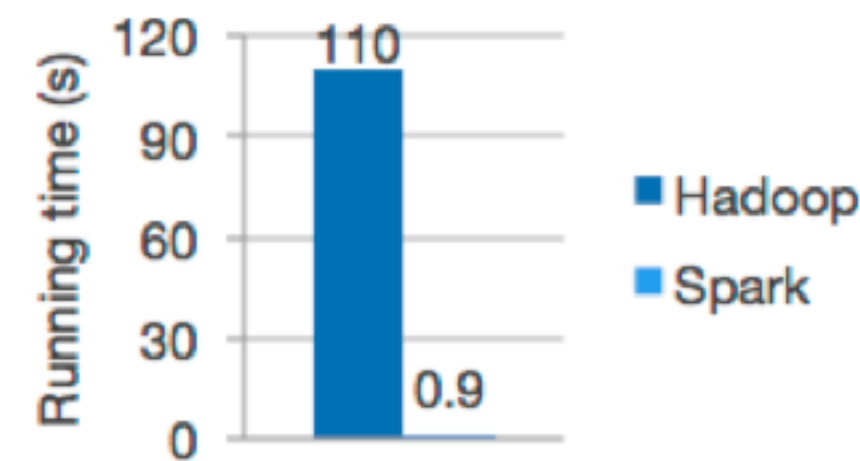


**Apache Spark**<sup>™</sup> is a unified analytics engine for large-scale data processing.

## Speed

Run workloads 100x faster.

Apache Spark achieves high performance for both batch and streaming data, using a state-of-the-art DAG scheduler, a query optimizer, and a physical execution engine.



Logistic regression in Hadoop and Spark

## Ease of Use

Write applications quickly in Java, Scala, Python, R, and SQL.

Spark offers over 80 high-level operators that make it easy to build parallel apps. And you can use it *interactively* from the Scala, Python, R, and SQL shells.

```
df = spark.read.json("logs.json")
df.where("age > 21")
  .select("name.first").show()
```

Spark's Python DataFrame API  
Read JSON files with automatic schema inference

### Latest News

Spark 2.4.3 released (May 08, 2019)

Spark 2.4.2 released (Apr 23, 2019)

Spark 2.4.1 released (Mar 31, 2019)

Spark 2.3.3 released (Feb 15, 2019)

[Archive](#)



[Download Spark](#)

### Built-in Libraries:

[SQL and DataFrames](#)

[Spark Streaming](#)

[MLlib \(machine learning\)](#)

[GraphX \(graph\)](#)

[Third-Party Projects](#)

	Hadoop MR Record	Spark Record	Spark 1 PB
Data Size	102.5 TB	100 TB	1000 TB
Elapsed Time	72 mins	23 mins	234 mins
# Nodes	2100	206	190