

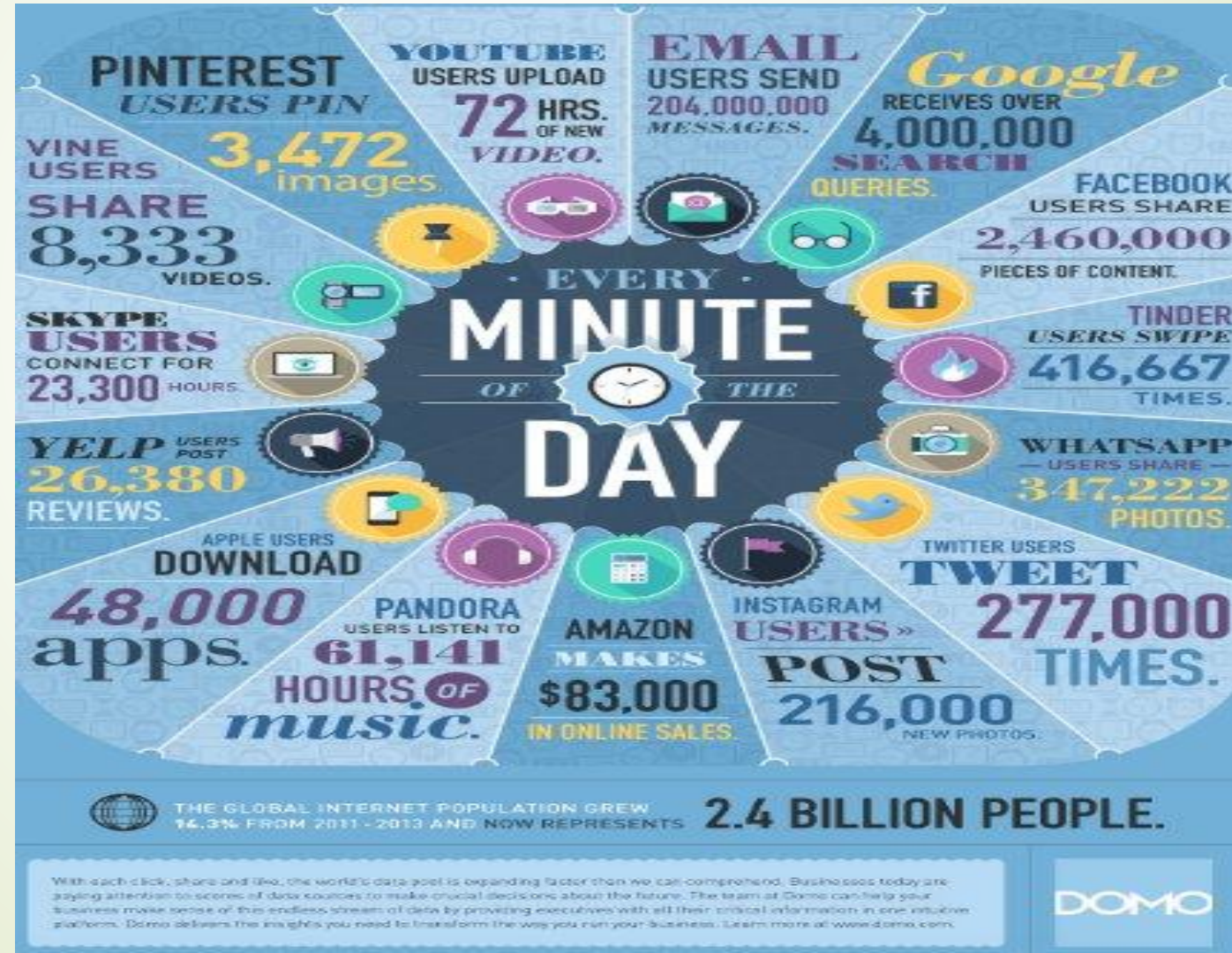


Data Processing with Apache Spark

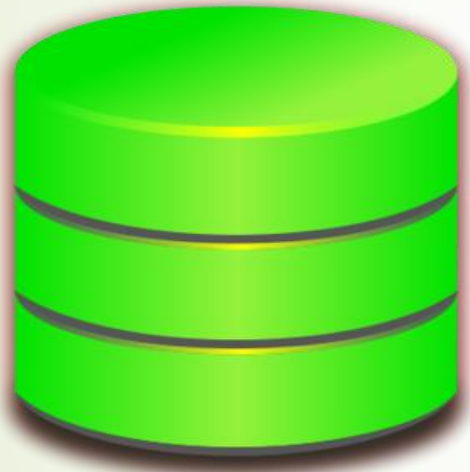
Harsha Puthalapattu

@harshappt

Big Data



Big Data Problem



Processing...



Username

Password

Big Data Problem



- ↑ Commodity hardware
- ↑ Cheaper
- ↑ Easy to add as data grows

- ↓ Non Reliable
- ↓ Performance Issues

- ↓ How to handle data serialization?
- ↓ How to ensure data integrity?
- ↓ How to recover failures?
- ↓ How to store and retrieve data?

Big Data Problem



IMPALA



- ↓ Map Reduce is Slow in large jobs
- ↓ Specialized processing requires specialized tools
- ↓ Data Engineers and Analysts need to learn multiple frameworks and tools



Spark – Fast and General purpose Cluster Computing System

- AMP Labs, UC Berkeley
- Paper published in 2010
- Apache top level project in 2014
- Evolved as ecosystem
- Primarily provides – Scheduling, Monitoring and Distributing Capabilities
- Unified API

Spark Components

Scala

Java

python

R

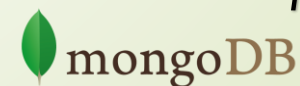
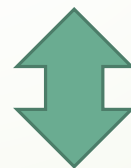
Data Frames
Data Sets

Streaming

MLlib
Machine Learning

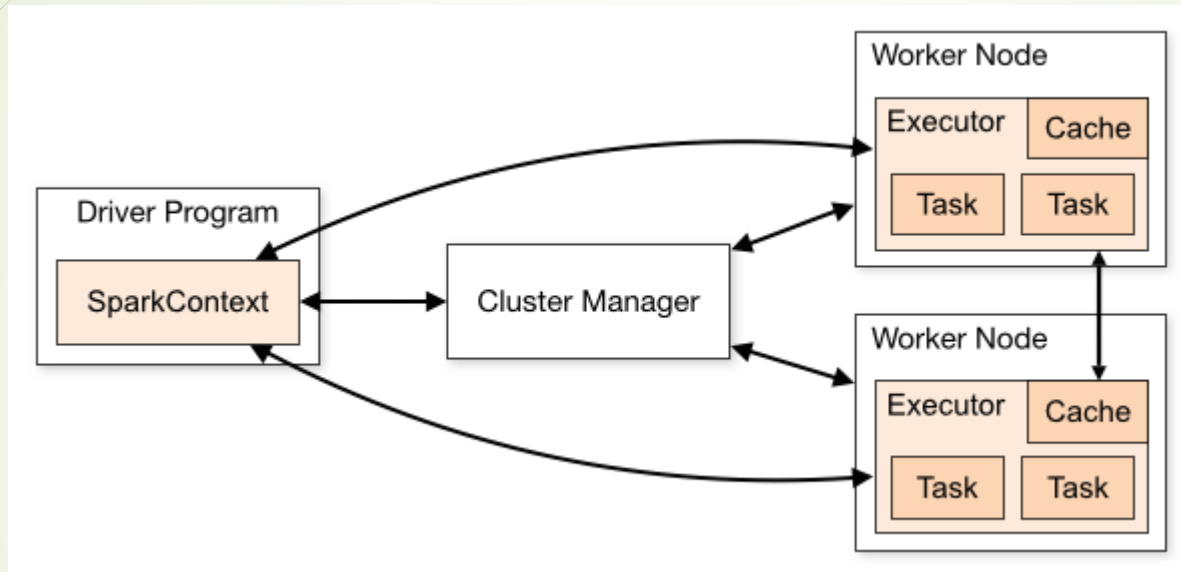
GraphX
Graph Computations

Spark Core
RDD API, Scheduling, Distributing, Monitoring



And
Many More

Spark Execution Architecture



- Driver is the process running the `main()` function of the application and creating the `SparkContext`
- Cluster Manager is an external service for acquiring resources on the cluster (e.g. standalone manager, Mesos, YARN)
- Worker Node is any node that can run application code in the cluster
- Executor is a process launched for an application on a worker node, that runs tasks and keeps data in memory or disk storage across them. Each application has its own executors



Demo

RDD and DataSet