Lesson 1.1: Why Distributed Computing?

Slide 2: Learning Objectives

## Learning Objectives

Motivate the business need for processing big data

Identify key concepts related to distributed computing
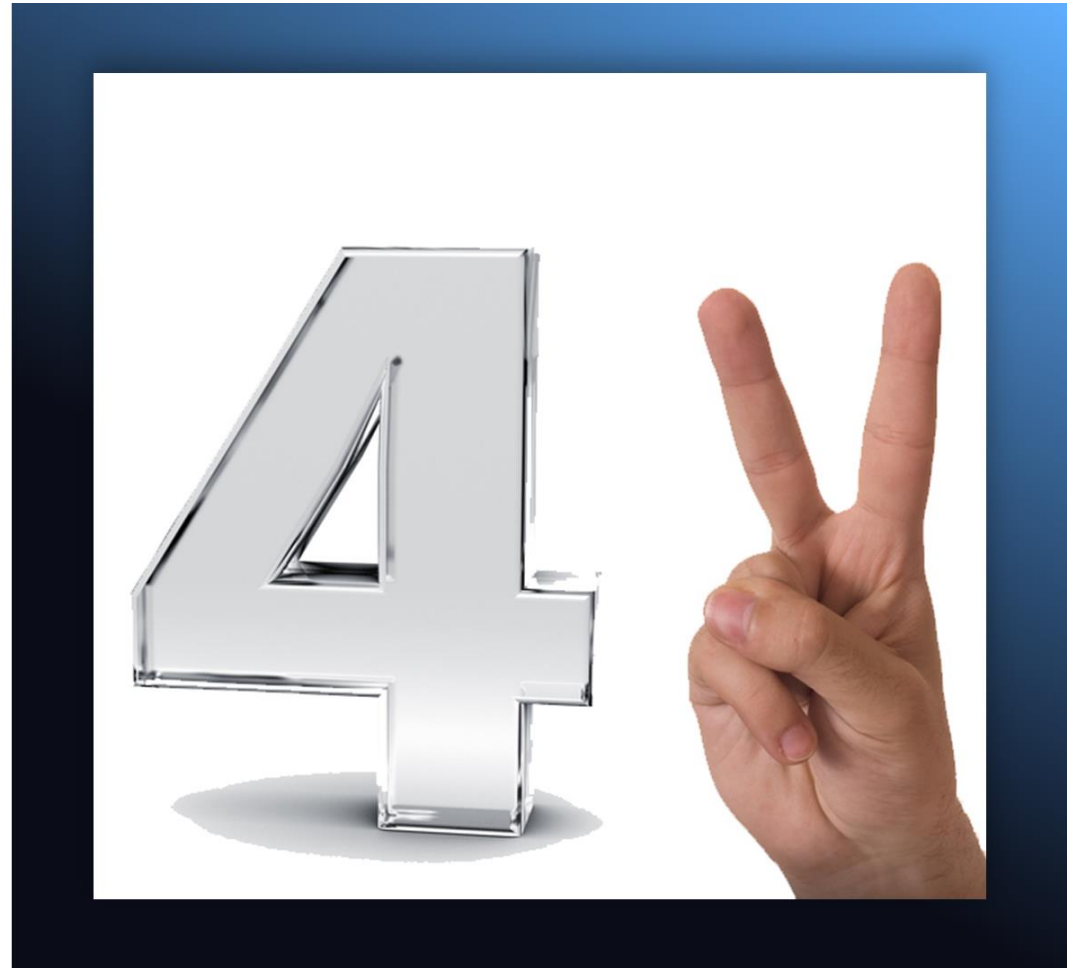
Slide 3: Qualities of Big Data

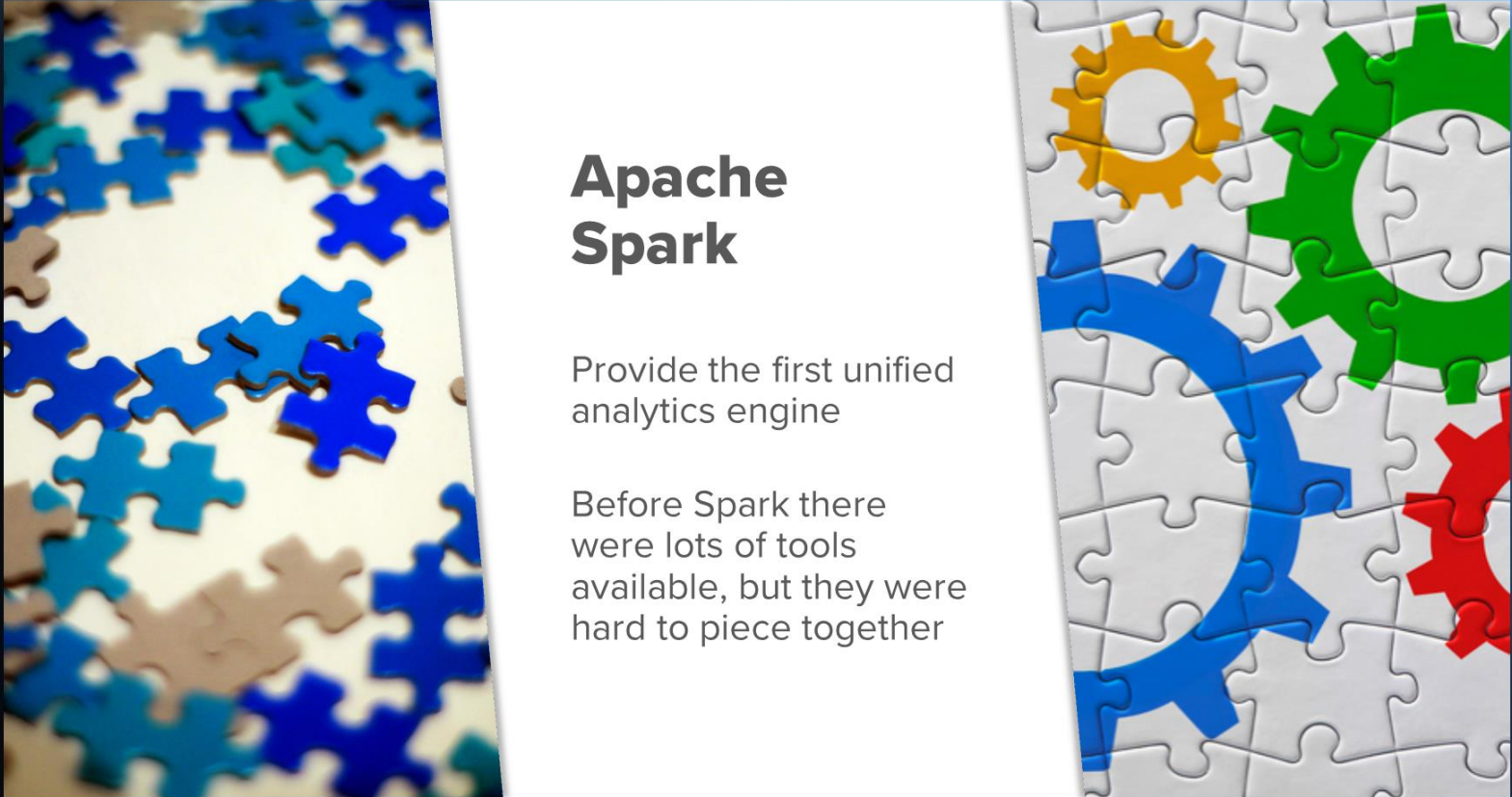## Qualities of Big Data

Volume

Velocity

Variety

Veracity

Slide 4: Big Data Defined

Slide 5: Apache Spark



**Apache Spark**

Provide the first unified analytics engine

Before Spark there were lots of tools available, but they were hard to piece together

Slide 6: Spark is Multilingual and Supports Many Languages
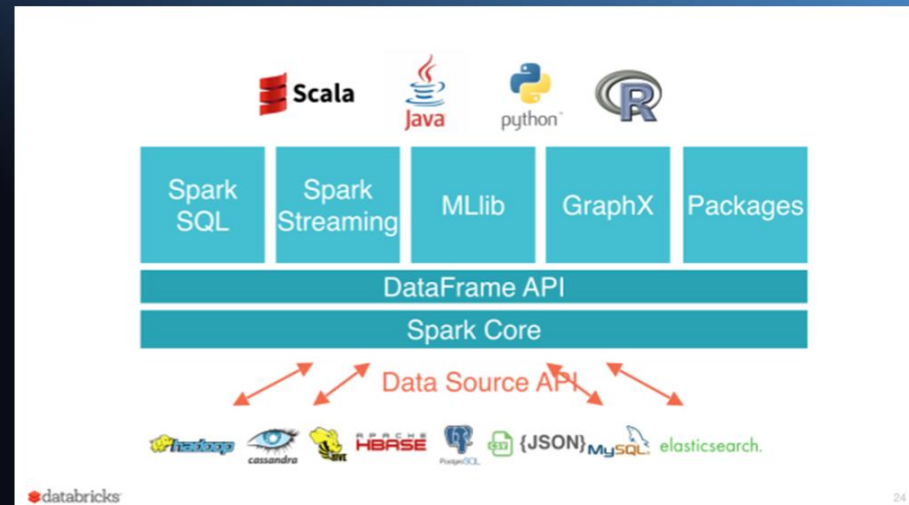
## Spark Supports Many Languages

SQL

Python

Scala

Java

R

Slide 7: Why Spark is Popular

## Why Spark is Popular

Reads & processes data from many sources

Works with many file types

Solves many data problems faced by analysts

Slide 8: Apache Spark: Origin Story

Slide 9: Apache Spark: Origin Story

Slide 10: Apache Spark: Origin Story

Slide 11: Let's Count Some M&Ms

Slide 12: Measure by Weight?

Slide 13: Involving Others



**Involving Others**

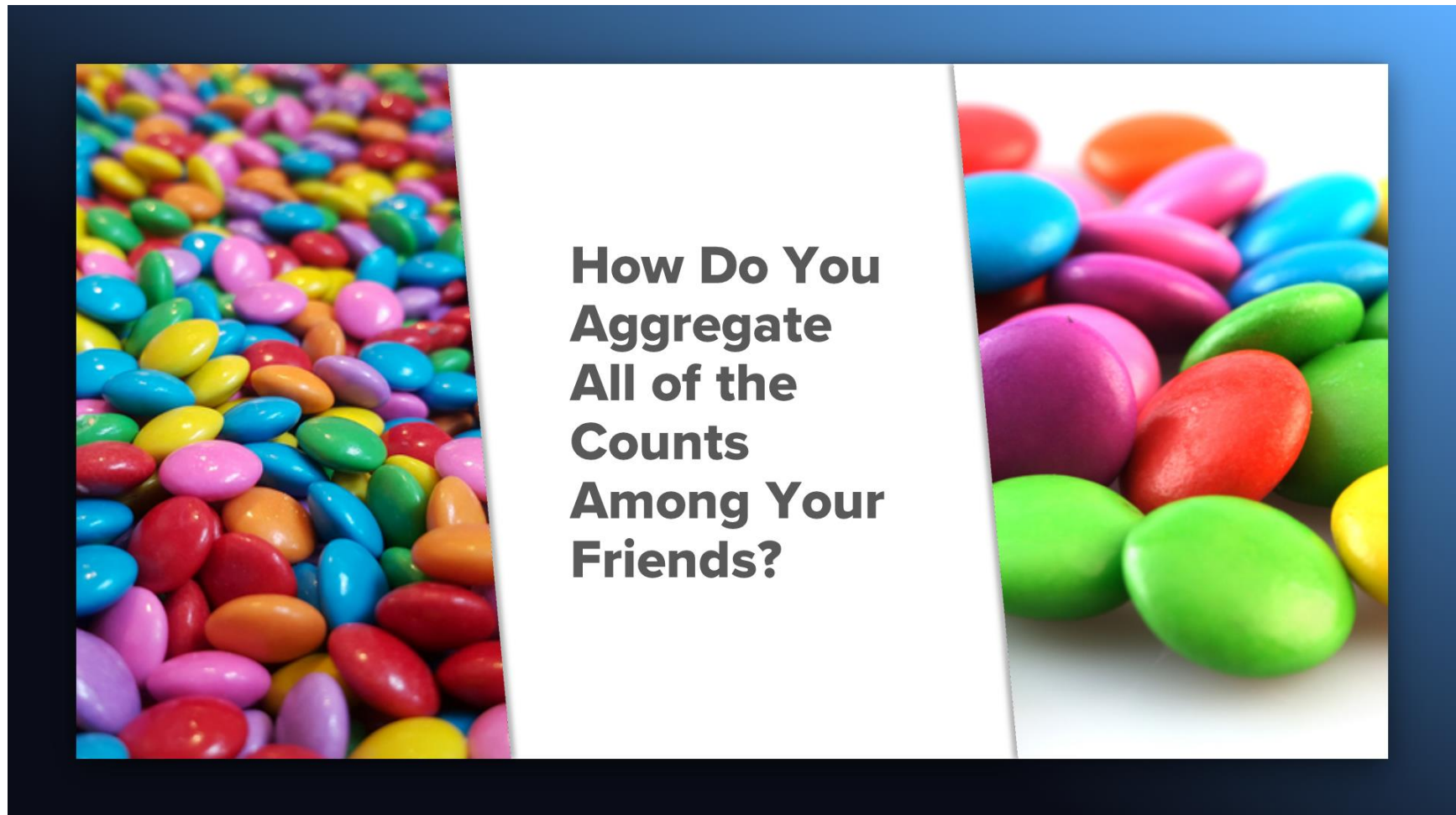Give everyone a handful to count
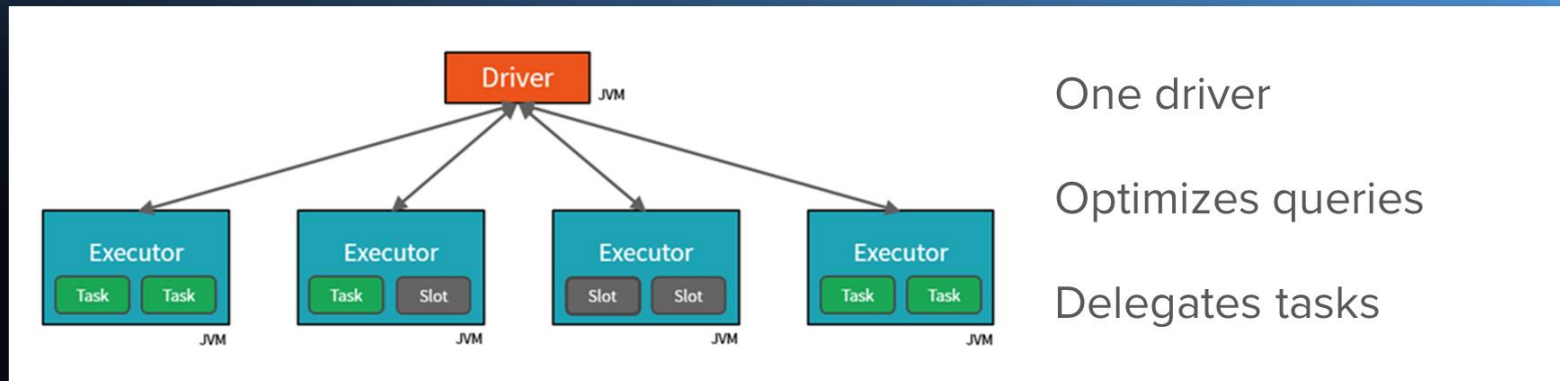
Slide 14: How Do You Aggregate All of the Counts Among Your Friends?

Slide 15: Drivers and Executors

Slide 16: Drivers and Executors



One driver

Optimizes queries

Delegates tasks

One or many executors

Perform actual queries

More is not always faster

Slide 17: Why More Computing Power Isn't Always Faster

Slide 18: Distributing Computation is Parallelism

Slide 19: Amdahl's Law

# Amdahl's Law

The amount of acceleration we would see from parallelizing a task is a function of what portion of the task can be completed in parallel

Slide 20: Linear Scalability



**Linear Scalability**

Dividing tasks across a cluster of machines

We see improvements up to thousands of machines

Spark is on par with best distributed computing solutions on the market

Slide 21: Scalability



**Scale out**

More data to process than on one machine

**Speed up**

More computer resources may speed up your query

Slide 22: Coming Up

## Coming Up

Learn about core Spark concepts