# Distributed Computing with Spark

02807 Computational Tools for Data Science

# Plan for rest of the course

2019-11-12: Introduction to Spark

2019-11-19: Spark under the hood

2019-11-26: Guest lecture and project work

2019-12-03: Guest lecture and project work

2019-12-05: Hand-in project #3 and mandatory assignment #3

# Talk Announcement

Title: **Working with Data at Warp Speed**
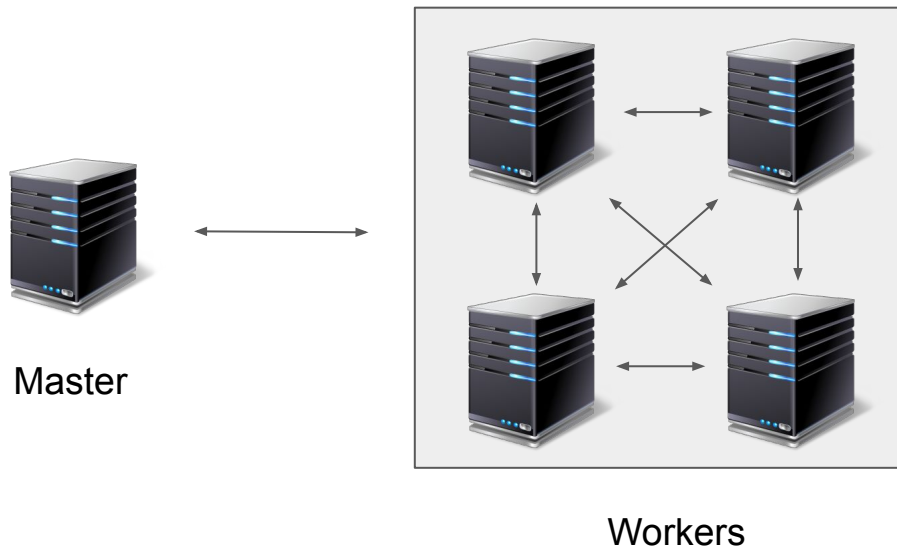
Speaker: **Brian Bell** (datarobot.com)

Date: **2019-11-26**

# Today

- Distributed computing
- Spark
  - Spark dataframes
  - PySpark
  - Spark UI

# Distributed Computing

Computation is distributed among multiple computers (nodes)



Master

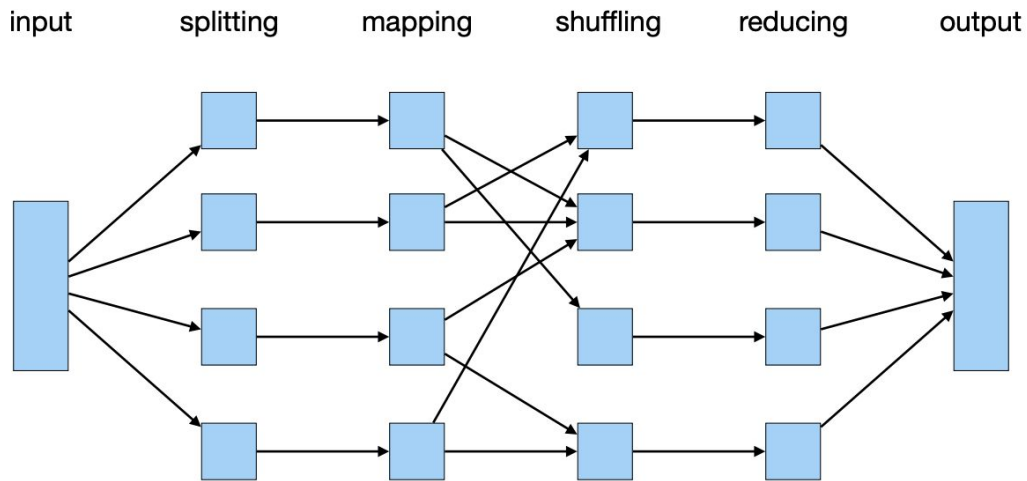Workers

# Distributed Computing

Pros:

- Speeding up computation
- Spreading risk of crashing

Cons:

- Exchanging data between workers is expensive (bottleneck: network)
- Synchronization is expensive (bottleneck: disk)
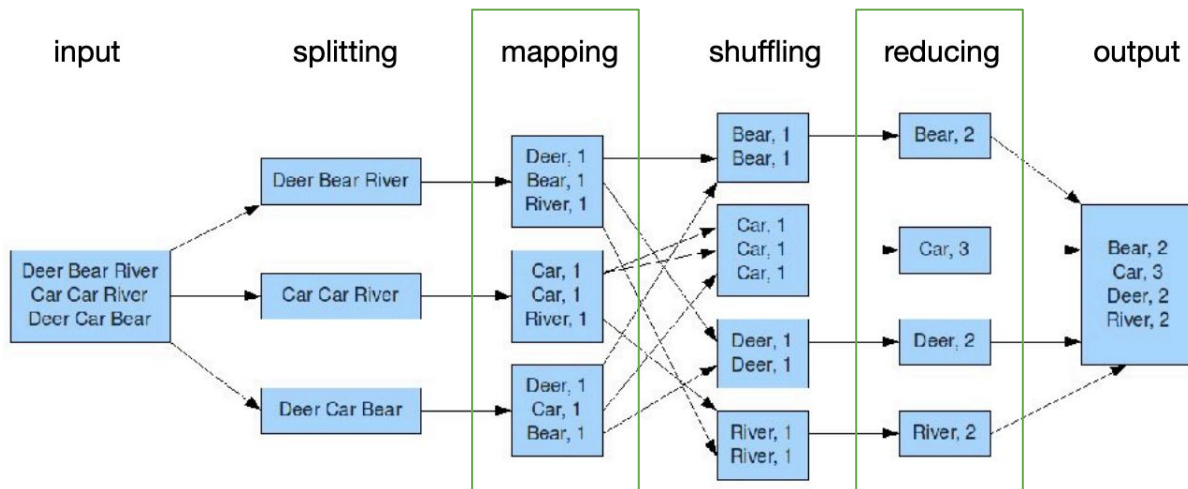- Writing distributed algorithms is hard

# Distributed Computing

- Writing distributed algorithms is hard, but...
- When processing data, the MapReduce model makes it simple
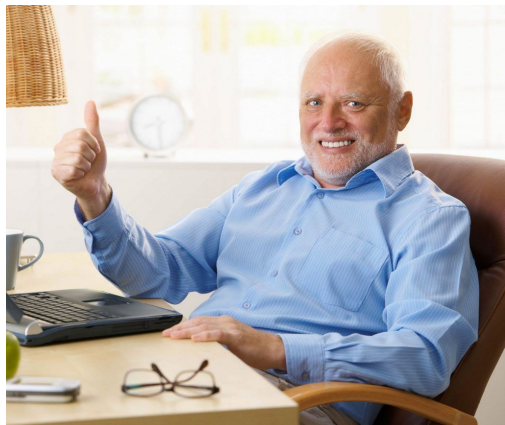
# MapReduce model example

- Input: Document of words (Deer Bear River Car Car River Deer Car Bear)
- Output: The frequency of each word



- Frameworks like Spark enables you to provide map and reduce functions - everything else is taken care of by the framework
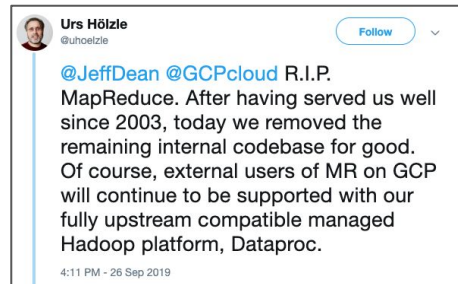
# *But I don't have cluster!*

- Modern laptops are single node, multi-core systems
- We want full utilization of our local resources
- Programs written in Spark will run locally as well as on clusters without modification

# Introducing Apache Spark



- Distributed computing framework
- Invented at Berkeley
- Written in Scala
- Free
- Has a growing community
- Becoming de-facto standard in the industry
- Faster than its predecessors

> **Urs Hölzle**
> @uhoelzle
>
> @JeffDean @GCPcloud R.I.P.
> MapReduce. After having served us well
> since 2003, today we removed the
> remaining internal codebase for good.
> Of course, external users of MR on GCP
> will continue to be supported with our
> fully upstream compatible managed
> Hadoop platform, Dataproc.
>
> 4:11 PM - 26 Sep 2019

Databricks' Growth Draws $400 Million Series F Investment and $6.2 Billion Valuation

October 22, 2019 10:00 PT

# Spark feature highlights

- Built-in resilience
- Node awareness
- Supports many file formats/data sources
- Streaming
- Same code, many platforms
- High-level API

# Spark DataFrame

- High level/declarative SQL-like API
- Immutable and lazy evaluation of commands
- Automatic optimization
- No random access (direct access to a specific row)

# Spark DataFrame: High level/declarative SQL-like API

- You tell Spark *what* to do, not *how* to do it
- Example: counting males

What:

```
df.filter(f.col('sex') == 'male').count()
```

How:

```
males = 0
for person in people:
    if person['sex'] == 'male':
        males += 1
```

# Spark DataFrame: Immutability and lazy evaluation

- Immutability: Every operation produces a new dataframe
- Two kinds of operations: transformations and actions
  - A transformation creates a new dataframe, but nothing is actually computed (*lazy* operation)
  - An action triggers an actual computation
- Example:

```
df = spark.read.csv('people.csv')
filtered_df = df.filter(f.col('sex') == 'male')
number_of_males = filtered_df.count()
```

| Read data | | Filter | | Count |
|-----------|---|--------|---|-------|
| *Transformation* | | *Transformation* | | *Action* |

# Spark DataFrame: Automatic optimization

- Spark re-organizes your operations to optimize performance
- More about this next week!

# Spark DataFrame: No random access

- You cannot access row $i$ in a dataframe the same was in a Python list
- After the split, the master knows the size of the split, but doesn't know which rows that went into which splits
- You **have to** think in terms of the MapReduce model to leverage Spark

# Let's code

# Job analysis

- Compute hash(ip, domain) for each row and use value to decide splitting
- For each split, compute uniques per domain
- Compute hash(domain) to group domains in new splits
- Sum up the uniques