

INTRO to DATA SCIENCE

LECTURE 22: MAP-REDUCE

I. SETTING UP AMAZON WEB SERVICES (AWS)

II. KEEPING AN EYE ON COSTS

III. GENERATING KEY PAIRS

IV. LAUNCHING AN EC2 INSTANCE

0. ETHICS BY MONICA BULGER, DATA & SOCIETY

I. BIG DATA

II. MAP-REDUCE

III. LAUNCHING A SPARK CLUSTER

IV. PYSPARK EXERCISES ON AWS

V. PYTHON EXERCISES: MULTIPROCESSING, MRJOB

- **DESCRIBE WHY HADOOP AND MAP-REDUCE EXIST**
- **WRITE MAP-REDUCE JOBS IN PYTHON**
- **LAUNCH SPARK CLUSTER AND WRITE MAP-REDUCE JOBS**

INTRO TO DATA SCIENCE

I. BIG DATA

Big data is characterized by

Big data is characterized by

- ▶ *Volume*

Big data is characterized by

1 exabyte = 1K petabytes = 1M TB = 1B GB

► *Volume*

Volume

- *As of 2012, about 2.5 exabytes of data are created each day, and that number is doubling every 40 months or so*

Big data is characterized by

1 exabyte = 1K petabytes = 1M TB = 1B GB

▶ *Volume*

Volume

- ▶ *As of 2012, about 2.5 exabytes of data are created each day, and that number is doubling every 40 months or so*
- ▶ *Walmart collects 2.5+ PB of data every hour from its customers*

Big data is characterized by

- ▶ *Volume*
- ▶ *Velocity*

Velocity

- ▶ *Twitter's firehose sends 6,000 tweets a second on average*

Big data is characterized by

‣ *Volume*

‣ *Velocity*

Velocity

- *Twitter's firehose sends 6,000 tweets a second on average*
- *Twitter saw 140,000 tweets in a single second on August 3, 2013*

Big data is characterized by

▶ *Volume*

▶ *Velocity*

Velocity

- ▶ *Twitter's firehose sends 6,000 tweets a second on average*
- ▶ *Twitter saw 140,000 tweets in a single second on August 3, 2013*
- ▶ *Assuming a size of 3 kB per tweet, 6K tweets a second is ~1.5 TB a day*

Big data is characterized by

- ▶ *Volume*
- ▶ *Velocity*
- ▶ *Variety*

Variety

- ▶ *Cell phone data, texts, GPS signals, etc.*

*cellphones ubiquitous
early 2000s*

Big data is characterized by

- ▶ *Volume*
- ▶ *Velocity*
- ▶ *Variety*

Variety

- ▶ *Cell phone data, texts, GPS signals, etc.*
- ▶ *Messages, posts, images posted to social networks*

*cellphones ubiquitous
early 2000s*

*Facebook 2004
Twitter 2006*

Big data is characterized by

- ▶ *Volume*
- ▶ *Velocity*
- ▶ *Variety*

Variety

- ▶ *Cell phone data, texts, GPS signals, etc.*
- ▶ *Messages, posts, images posted to social networks*
- ▶ *Webpages, for desktop, mobile, tablet, etc.*

*cellphones ubiquitous
early 2000s*

*Facebook 2004
Twitter 2006*

*iPhone 2007
iPad 2010*

Big data is characterized by

- ▶ *Volume*
- ▶ *Velocity*
- ▶ *Variety*

Variety

- ▶ *Cell phone data, texts, GPS signals, etc.*
- ▶ *Messages, posts, images posted to social networks*
- ▶ *Webpages, for desktop, mobile, tablet, etc.*
- ▶ *Movies, music, etc.*

*cellphones ubiquitous
early 2000s*

*Facebook 2004
Twitter 2006*

*iPhone 2007
iPad 2010*

*In 2010, Netflix's streaming
surpassed its mailing business*

Big data is characterized by

- ▶ *Volume*
- ▶ *Velocity*
- ▶ *Variety*

Variety

- ▶ *Cell phone data, texts, GPS signals, etc.*
- ▶ *Messages, posts, images posted to social networks*
- ▶ *Webpages, for desktop, mobile, tablet, etc.*
- ▶ *Movies, music, etc.*
- ▶ *Internet of Things, readings from sensors, etc.*

*cellphones ubiquitous
early 2000s*

*Facebook 2004
Twitter 2006*

*iPhone 2007
iPad 2010*

*In 2010, Netflix's streaming
surpassed its mailing business*

Big data is characterized by

- ▶ *Volume*
- ▶ *Velocity*
- ▶ *Variety*

Variety

- ▶ *Cell phone data, texts, GPS signals, etc.*
- ▶ *Messages, posts, images posted to social networks*
- ▶ *Webpages, for desktop, mobile, tablet, etc.*
- ▶ *Movies, music, etc.*
- ▶ *Internet of Things, readings from sensors, etc.*
- ▶ *What's next?*

*cellphones ubiquitous
early 2000s*

*Facebook 2004
Twitter 2006*

*iPhone 2007
iPad 2010*

*In 2010, Netflix's streaming
surpassed its mailing business*

Big data is characterized by

- ▶ *Volume*
- ▶ *Velocity*
- ▶ *Variety*
- ▶ *Veracity*

Veracity

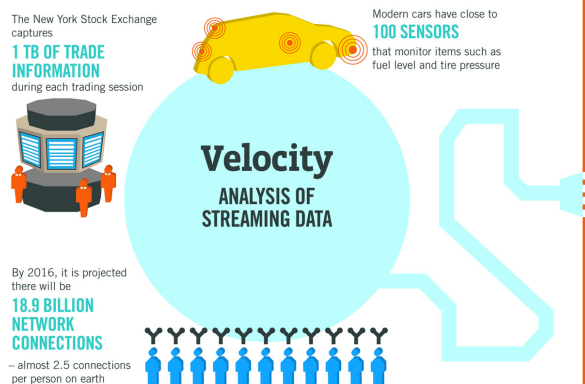
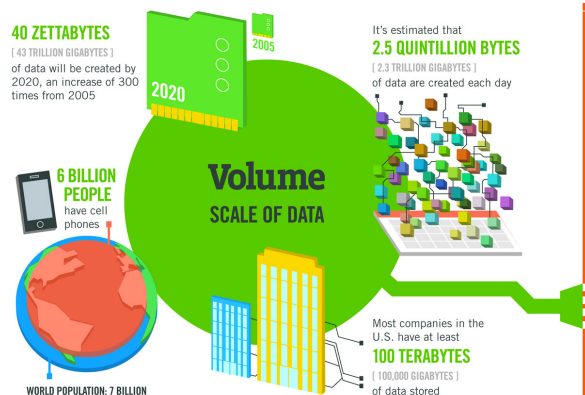
- ▶ *1 out of 3 business leaders don't trust the data they use to make decisions*

Big data is characterized by

- ▶ *Volume*
- ▶ *Velocity*
- ▶ *Variety*
- ▶ *Veracity*

Veracity

- ▶ *1 out of 3 business leaders don't trust the data they use to make decisions*
- ▶ *Poor data quality cost the US economy around \$3 trillion a year*



The FOUR V's of Big Data

From traffic patterns and music downloads to web history and medical records, data is recorded, stored, and analyzed to enable the technology and services that the world relies on every day. But what exactly is big data, and how can these massive amounts of data be used?

As a leader in the sector, IBM data scientists break big data into four dimensions: **Volume, Velocity, Variety and Veracity**

Depending on the industry and organization, big data encompasses information from multiple internal and external sources such as transactions, social media, enterprise content, sensors and mobile devices. Companies can leverage data to adapt their products and services to better meet customer needs, optimize operations and infrastructure, and find new sources of revenue.

By 2015
4.4 MILLION IT JOBS
will be created globally to support big data, with 1.9 million in the United States

As of 2011, the global size of data in healthcare was estimated to be **150 EXABYTES**
[161 BILLION GIGABYTES]



30 BILLION PIECES OF CONTENT
are shared on Facebook every month



Variety DIFFERENT FORMS OF DATA

By 2014, it's anticipated there will be **420 MILLION WEARABLE, WIRELESS HEALTH MONITORS**

4 BILLION+ HOURS OF VIDEO
are watched on YouTube each month



400 MILLION TWEETS
are sent per day by about 200 million monthly active users



1 IN 3 BUSINESS LEADERS
don't trust the information they use to make decisions



Poor data quality costs the U.S. economy around **\$3.1 TRILLION A YEAR**



Veracity UNCERTAINTY OF DATA

27% OF RESPONDENTS

in one survey were unsure of how much of their data was inaccurate

II. MAP-REDUCE

How do we process big data?

*One approach would be to get
a huge supercomputer*

- expensive*
- difficult to maintain*
- scalability is bounded*

How do we process big data?

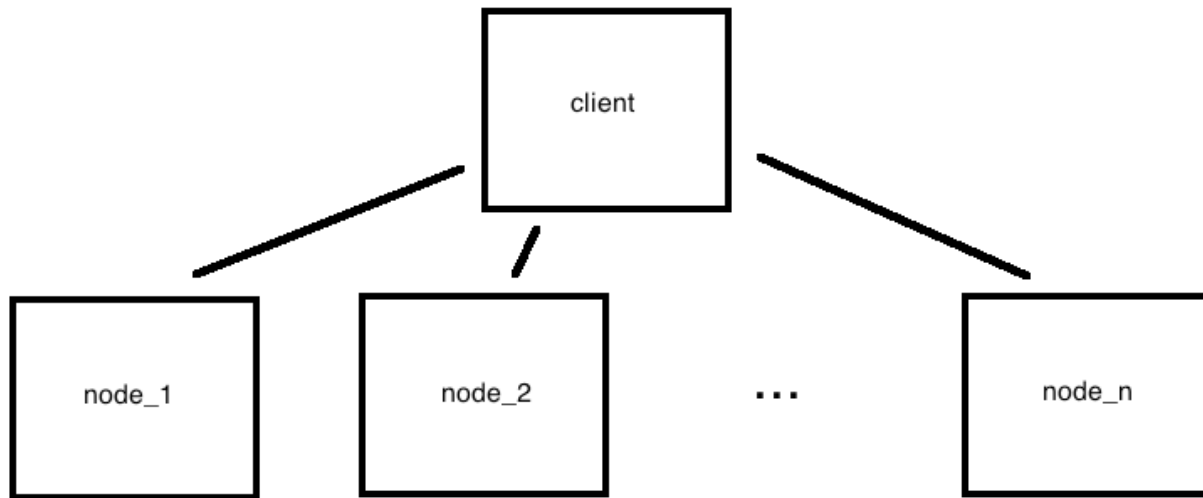
One approach would be to get a huge supercomputer

- expensive*
- difficult to maintain*
- scalability is bounded*

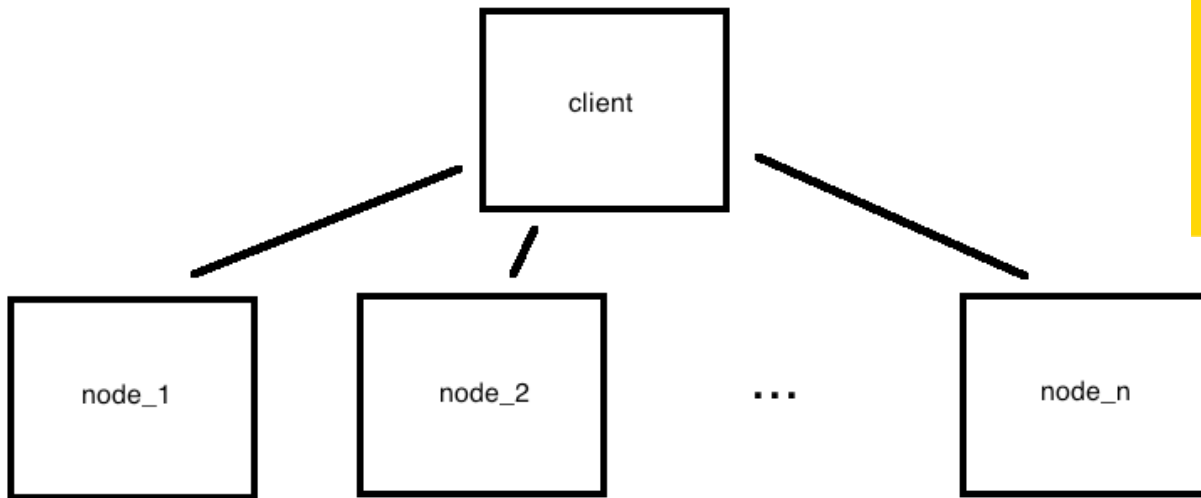
Another approach is to get a bunch of regular (commodity) machines

- cheaper*
- easier to maintain*
- scalability is unbounded*
(just add another node)

We can visualize this horizontal cluster architecture as a single client-multiple server relationship



We can visualize this horizontal cluster architecture as a single client-multiple server relationship

**NOTE**

A horizontally distributed system also has better *fault tolerance* than a single machine.

There are two ways to process data in a distributed architecture:

1) move data to code (& processing power)

2) move code to data

There are two ways to process data in a distributed architecture:

1) move data to code (& processing power)

- SETI

2) move code to data

- map-reduce → less overhead (network traffic, disk I/O)

“Computing nodes are the same as storage nodes.”

MapReduce *is method for parallel computing, introduced by Google.*

*Frequently when people say “map-reduce” they’re referring to Hadoop, but it is really a **paradigm** found in many other places as well
numpy, NoSQL, Cloudera, MapR, GFS, ...*

The MapReduce framework consists of two phases

- a **mapper**, which performs some local operation*
- a **reducer**, which aggregates the results that are sent back*

The *MapReduce* framework consists of two phases

- *a **mapper**, which performs some local operation*
- *a **reducer**, which aggregates the results that are sent back*

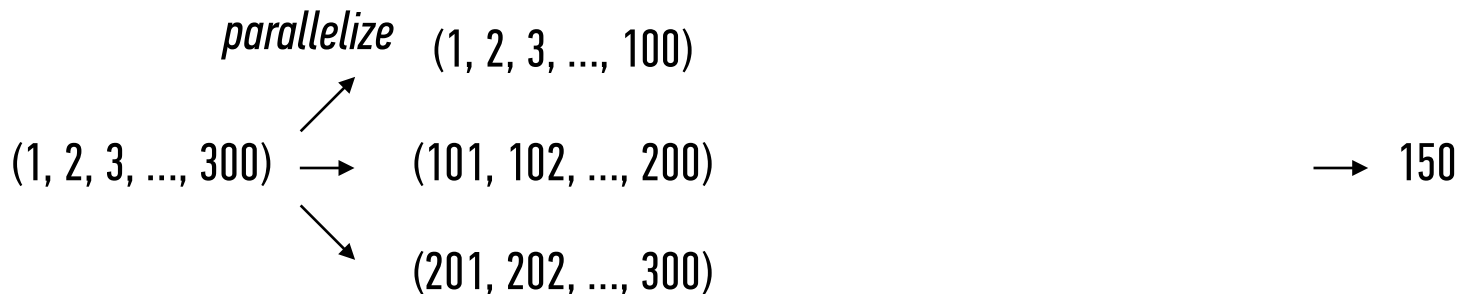
Example

(1, 2, 3, ..., 300) → *averaging* → 150

The *MapReduce* framework consists of two phases

- *a **mapper**, which performs some local operation*
- *a **reducer**, which aggregates the results that are sent back*

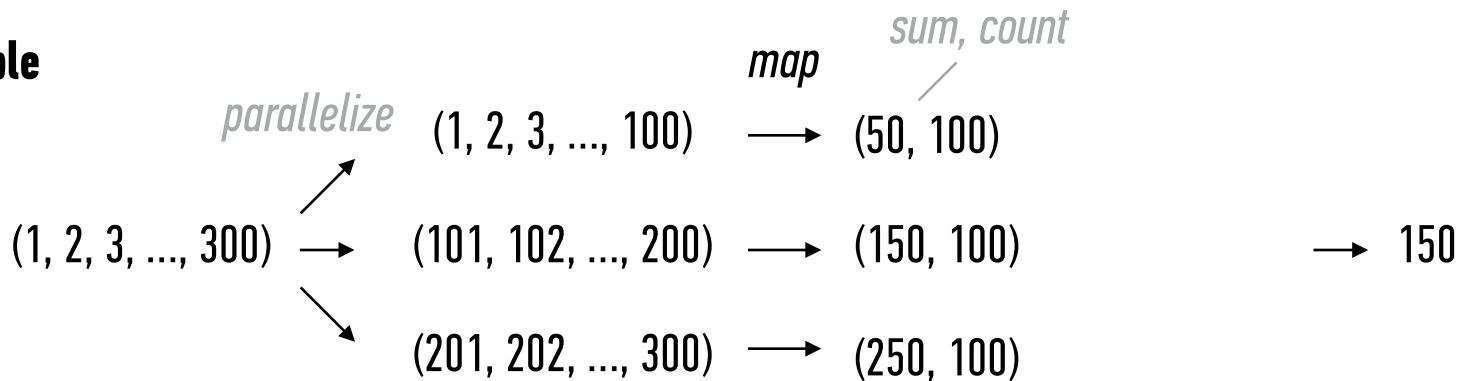
Example



The *MapReduce* framework consists of two phases

- **a mapper**, which performs some local operation
- **a reducer**, which aggregates the results that are sent back

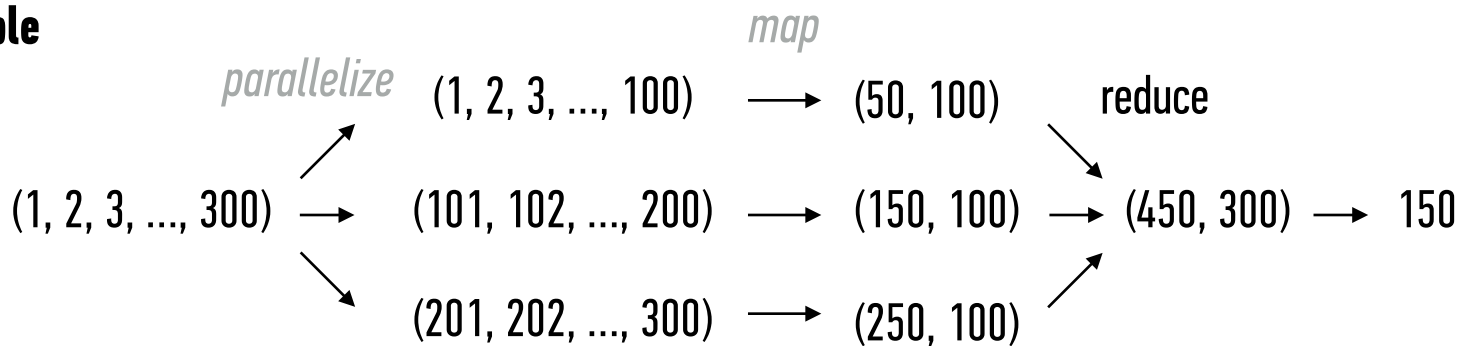
Example



The *MapReduce* framework consists of two phases

- *a **mapper**, which performs some local operation*
- *a **reducer**, which aggregates the results that are sent back*

Example



Example: Word Count

```
where  
where in  
where in the  
where in the world  
where in the world is  
where in the world is carmen  
where in the world is carmen sandiego
```

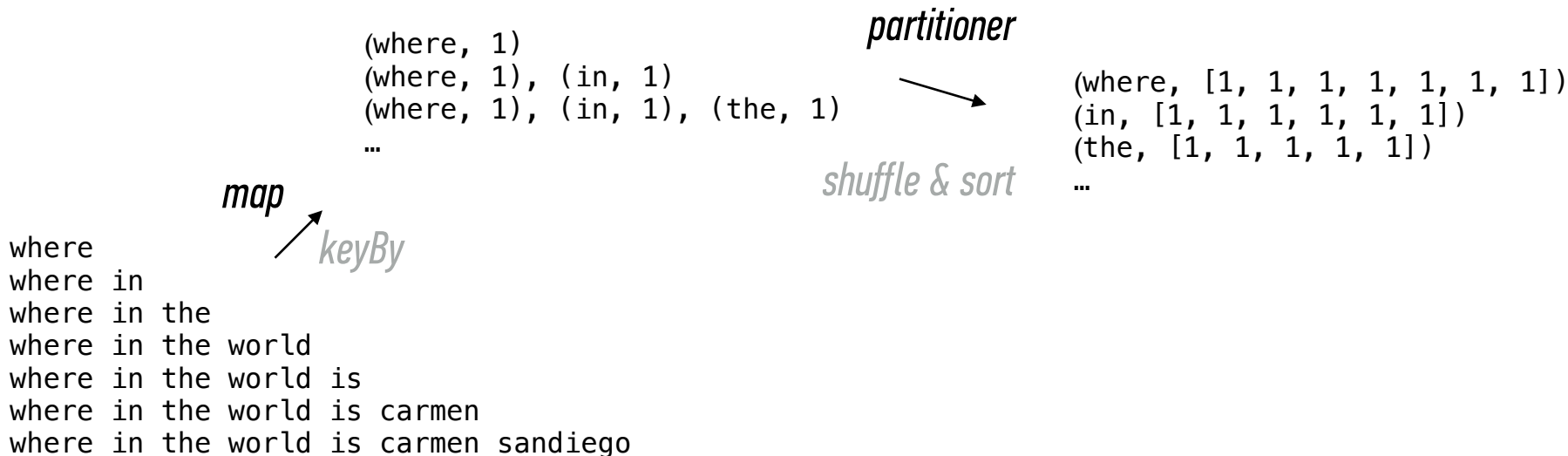
Example: Word Count

```
(where, 1)
(where, 1), (in, 1)
(where, 1), (in, 1), (the, 1)
...
```

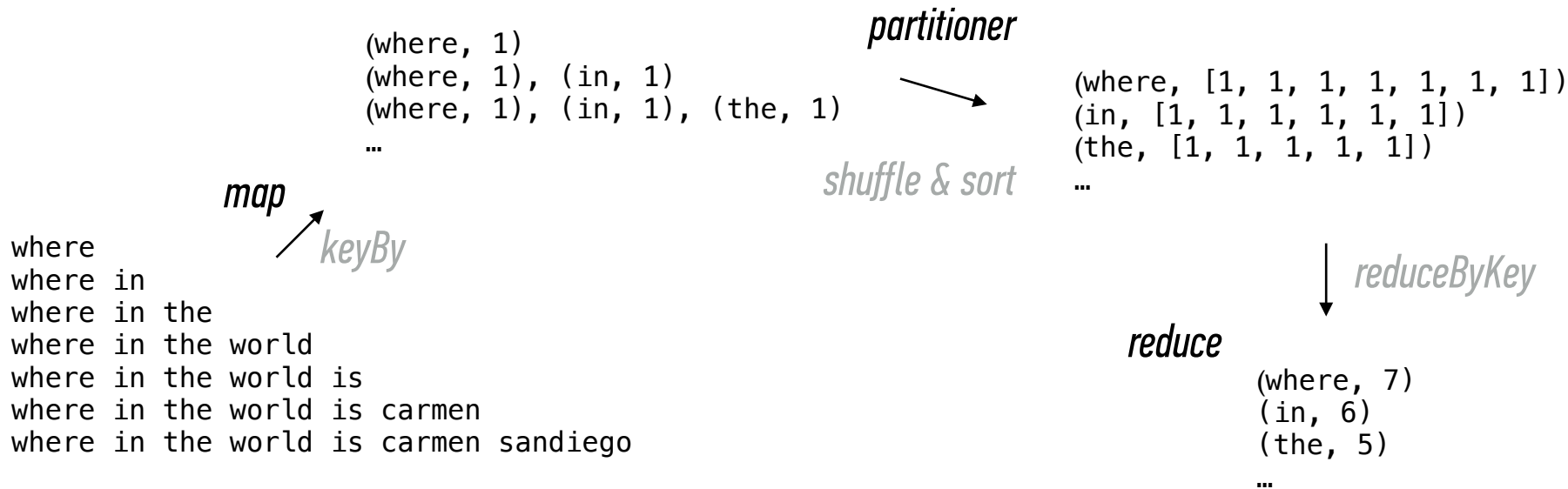
map ↗ *keyBy*

```
where
where in
where in the
where in the world
where in the world is
where in the world is carmen
where in the world is carmen sandiego
```

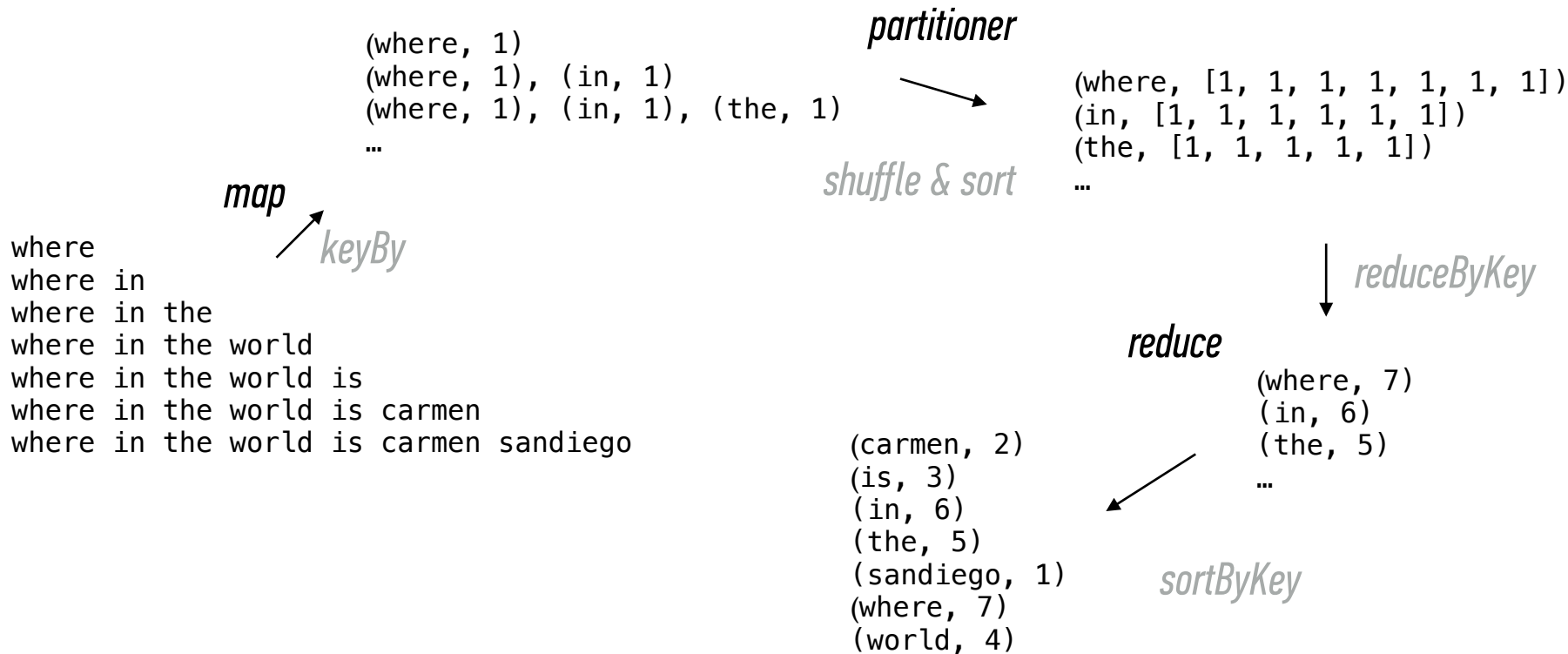
Example: Word Count



Example: Word Count



Example: Word Count



III. IMPLEMENTATION

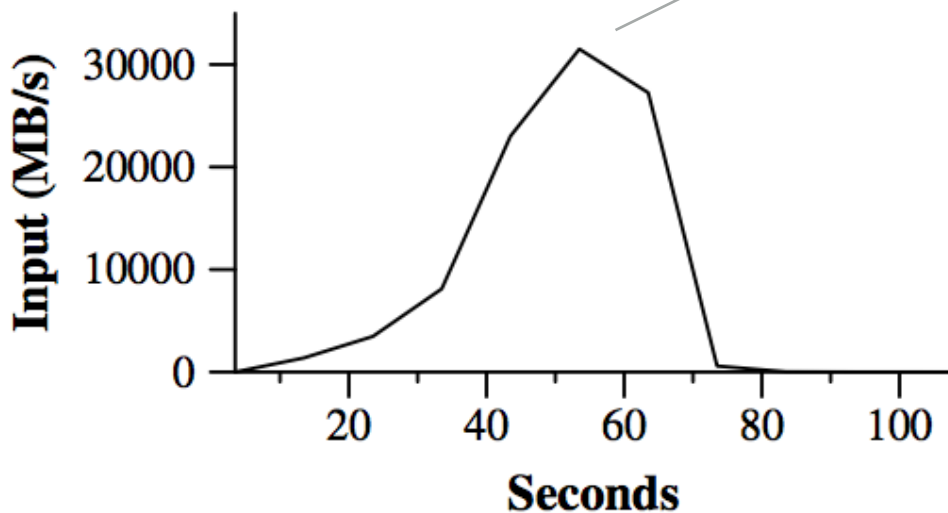
The map-reduce framework handles a lot of messy details for you:

- parallelization & distribution (eg, input splitting)*
- partitioning (shuffle/sort/redirect)*
- fault-tolerance (fact: tasks/nodes will fail!)*
- I/O scheduling*
- status and monitoring*

This (along with the functional semantics) allows you to focus on solving the problem instead of accounting & housekeeping details.

Data transfer rate over time

*Peaks at 30+ GB/s when
1764 workers have been assigned*



grep program

- 10^{10} 100-byte records
- searches for 3-char pattern (occurs ~90K times)
- input split into ~64MB pieces ($M = 15,000$)

The Google File System (GFS) was developed alongside map-reduce to serve as the native file system for this type of processing.

The Hadoop platform is bundled with an open-source implementation of this file system called HDFS.

If you use Amazon EMR, you can use their file system (Amazon S3) as well.

Hadoop is a popular open-source Java-based implementation of the map-reduce framework (including file storage for input/output).

You can download Hadoop and configure a set of machines to operate as a map-reduce cluster, or you can run it as a service via Amazon's Elastic Map-Reduce.

Hadoop is written in Java, but the Hadoop Streaming utility allows client code to be supplied as executables (eg, written in any language).

Doug Cutting, Hadoop's creator, named the framework after his child's stuffed toy elephant

IV. EXERCISES

- *Launch Spark cluster in AWS*
- *Write jobs in pyspark using AWS cluster*
- *multiprocessing example in ipython notebook*
- *mrjob implementation in python*

INTRO TO DATA SCIENCE

DISCUSSION