



Big Data

Mohammad Amin Nikbakht

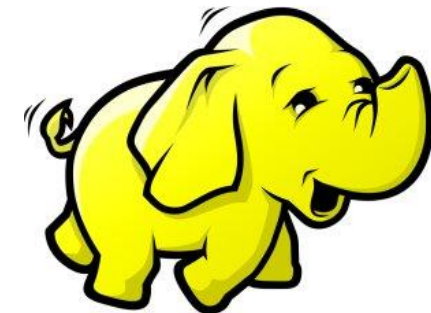
Electrical & Computer Department



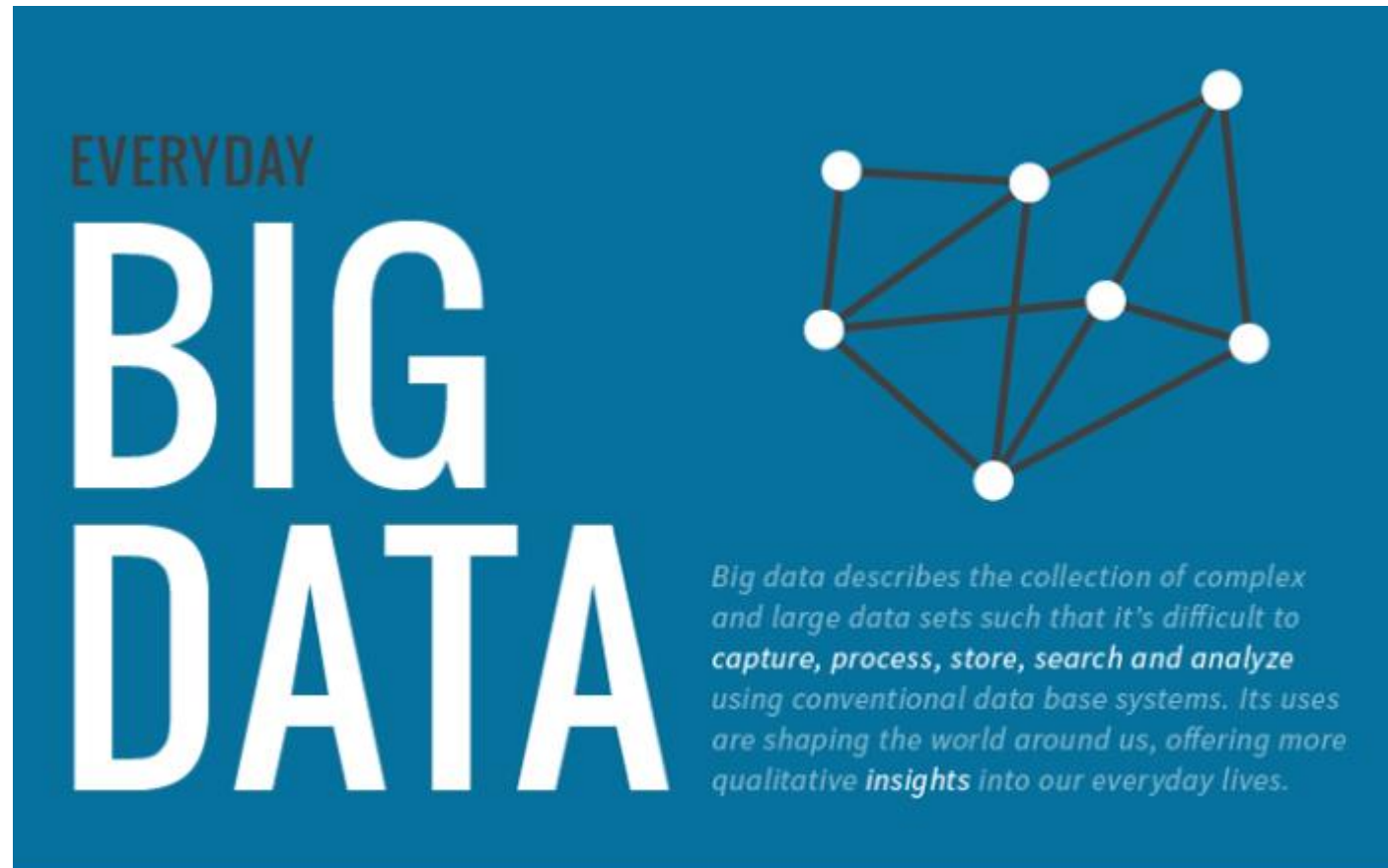
Isfahan University Of Technology

Agenda

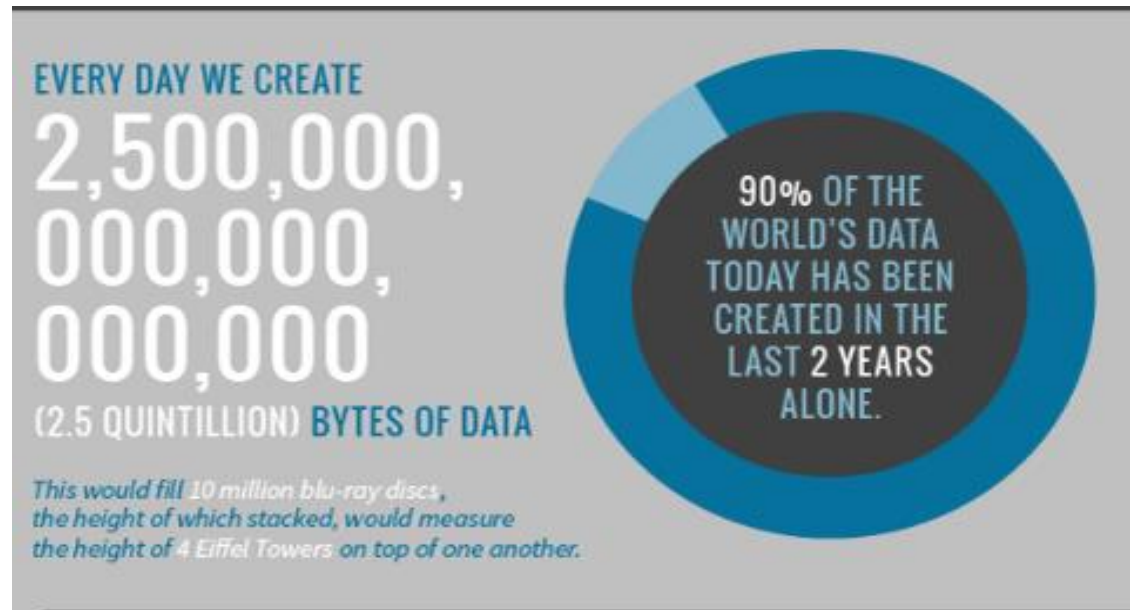
- What is Big Data?
- What is Hadoop?



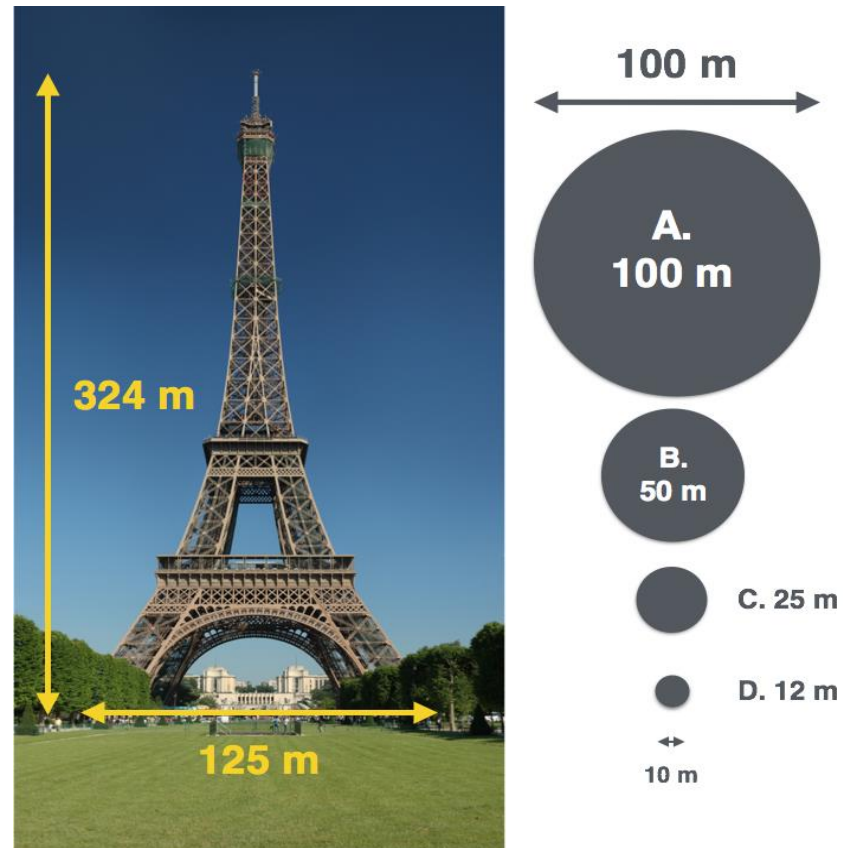
Big Data



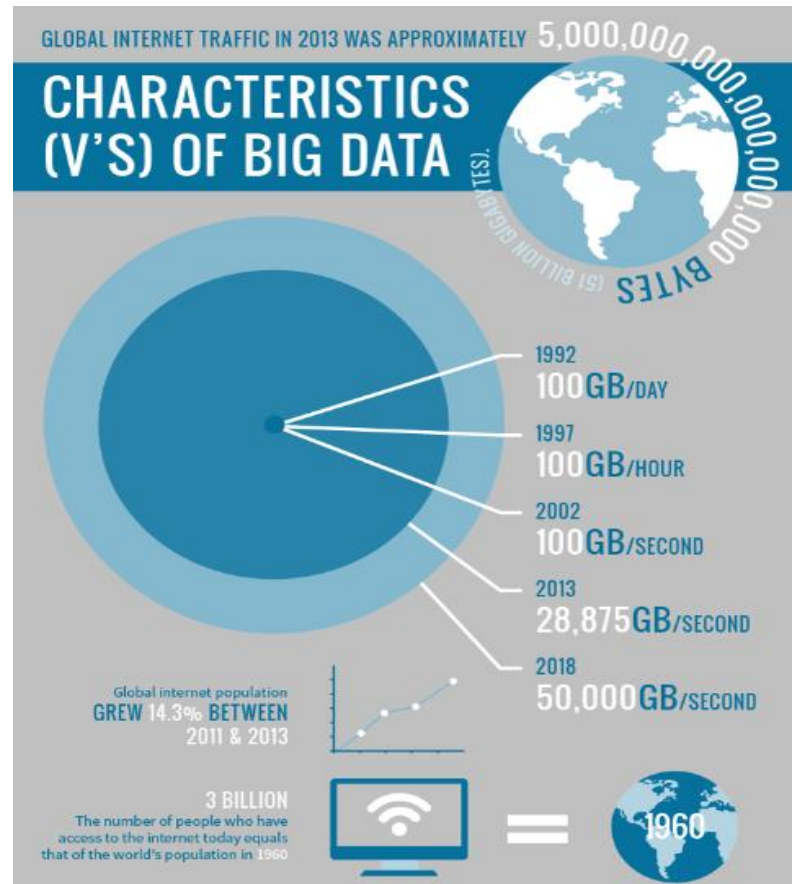
Big Data



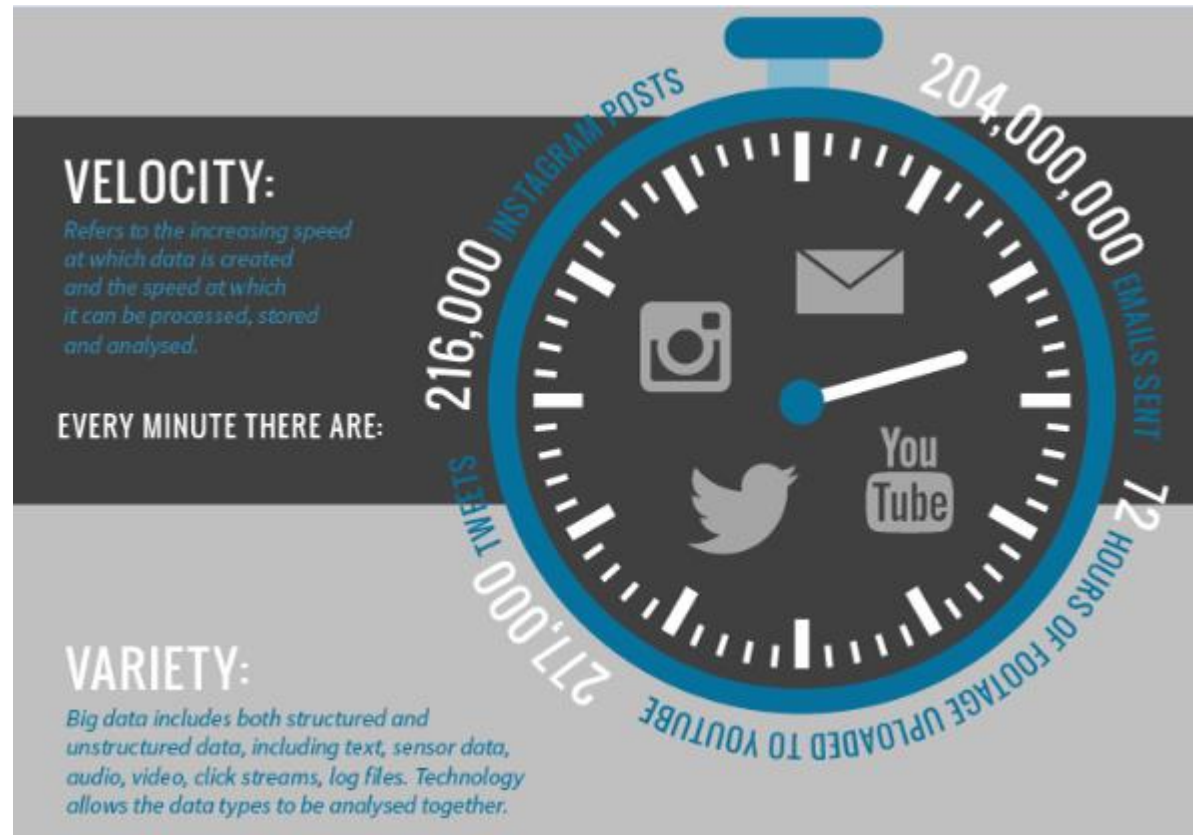
Big Data



Big Data



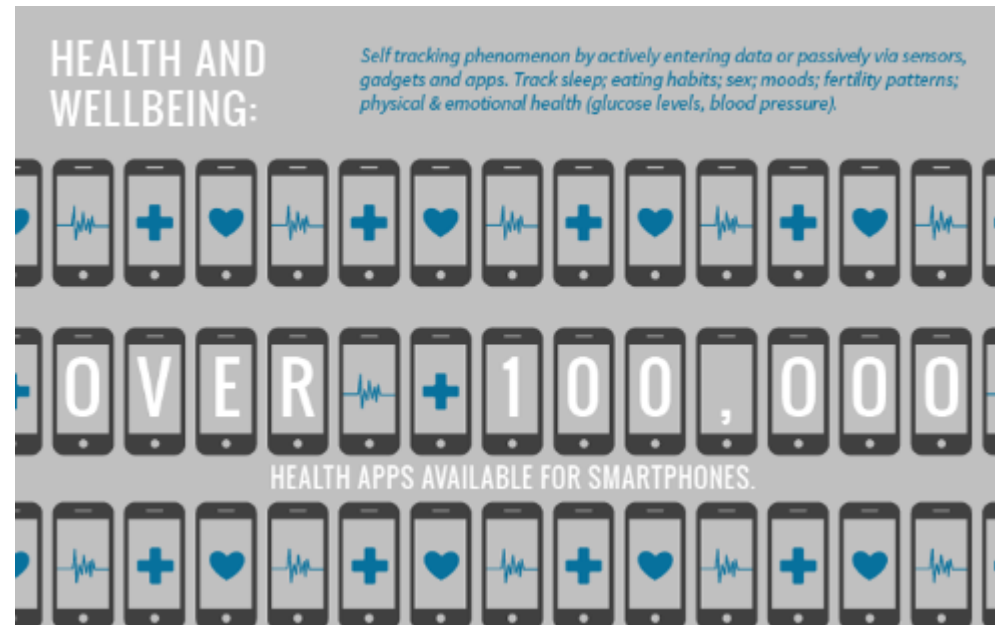
Big Data



Big Data



How does BIG DATA improve our daily lives ?



How does BIG DATA improve our daily lives ?



How does BIG DATA improve our daily lives ?

SAVING MONEY:



ONLINE:

- Big Data is being used to **match** market offers with consumers buying habits and individual needs
- **Relevant offers** from retailers you use and those who sell products that may be relevant to you
- Feedback gives opportunity to **engage** with businesses to ensure efficient service
- By saving money on their costs, businesses can pass these savings onto the consumer



AT HOME:

- Monitor and **reduce** energy usage

How does BIG DATA improve our daily lives ?



The infographic is divided into three sections, each with an icon and a list of benefits:

- TRAVEL:** (Icon: airplane)
 - Airlines have started to use customer data to improve customer service
 - Frequent fliers can soon expect the in-flight crew to know allergies; seat preferences; birthday; how they like their tea or coffee
- IN THE CAR:** (Icon: car key)
 - Monitor the condition of your car
 - Monitor mileage and fuel consumption
 - Insurance Telematics boxes can reduce insurance significantly
- SHOPPING:** (Icon: shopping cart)
 - Loyalty schemes enable shops to track what their customers are purchasing and tailor coupons accordingly
 - In-store location trackers interacting with smartphones as customers enter shops

How does BIG DATA improve our daily lives ?



URBAN TRANSPORT

USING:

- Real-time data capture
- Magnetic sensors installed in the road network
- GPS systems tracking public transport
- Social media monitoring systems

ENABLES TRANSPORT AGENCIES TO FACILITATE THE MANAGEMENT OF POTENTIAL TRAFFIC:

- by changing bus routes
- modifying traffic light sequences
- delivering information to drivers via mobile apps indicating approximate driving times and offering alternative routes

How does BIG DATA improve our daily lives ?



**WASTE
MANAGEMENT**

**DELIVERING EFFICIENT, EFFECTIVE AND ENVIRONMENTALLY
RESPONSIBLE WASTE COLLECTION AND RECYCLING
SERVICES BY USING:**

- sensors in waste container to detect the filling level
- historical data
- usage trends



Hadoop!

Ask bigger questions.

Timeline

- Dec 2004: Dean/Ghemawat (Google) MapReduce paper
- 2005: Doug Cutting and Mike Cafarella (Yahoo) create Hadoop, at first only to extend Nutch



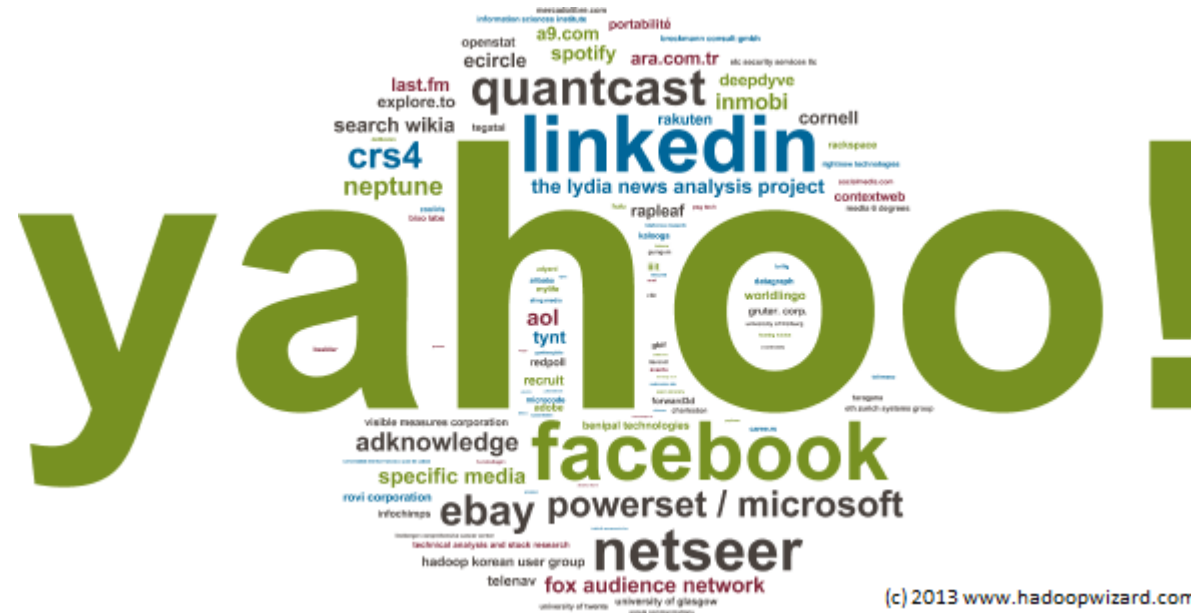
Timeline

- 2006: Yahoo runs Hadoop on 5-20 nodes
- March 2008: Cloudera founded
- July 2008: Hadoop wins TeraByte sort benchmark (1st time a Java program won this competition)
- April 2009: Amazon introduce “Elastic MapReduce” as a service on S3/EC2
- June 2011: Hortonworks founded

Timeline

- 27 dec 2011: Apache Hadoop release 1.0.0
- June 2012: Facebook claim “biggest Hadoop cluster”, totalling more than 100 PetaBytes in HDFS
- 2013: Yahoo runs Hadoop on 42,000 nodes, computing about 500,000 MapReduce jobs per day
- 15 oct 2013: Apache Hadoop release 2.2.0 (YARN)

Who uses Hadoop?



What is hadoop?

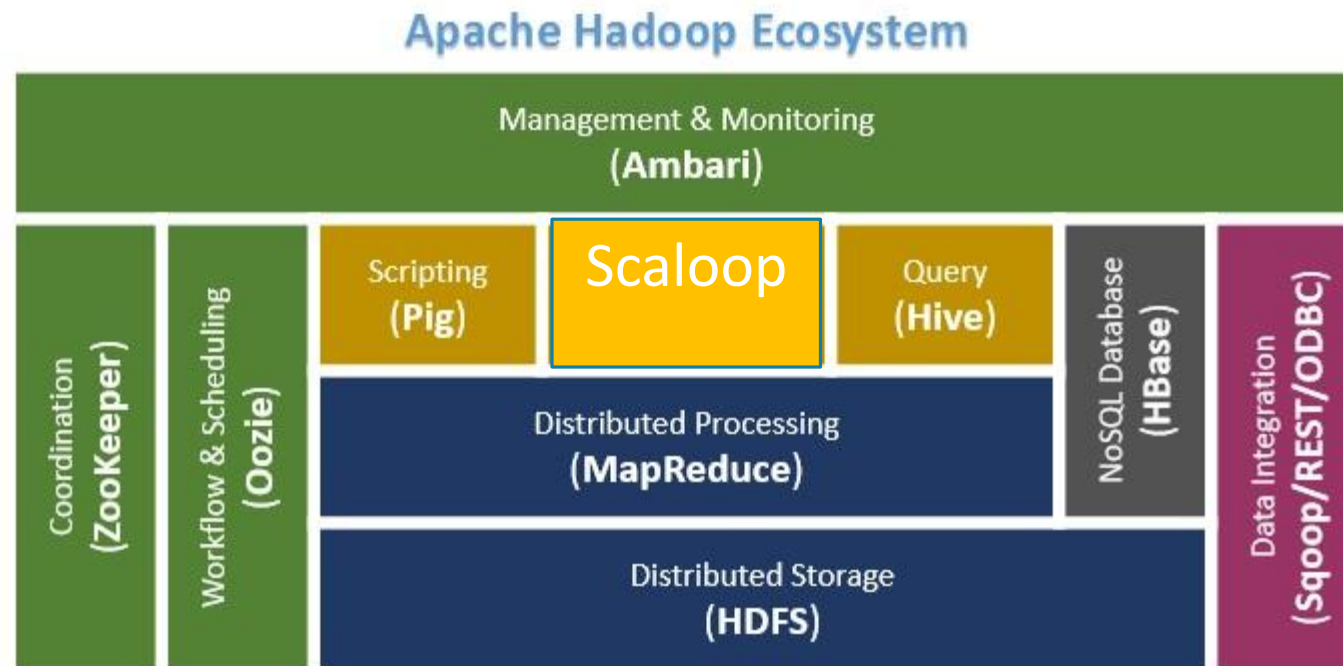


- **Open sourced, flexible and available** architecture for **large scale** computation and data processing on a network of **commodity hardware**
- **Inspired by** Google™

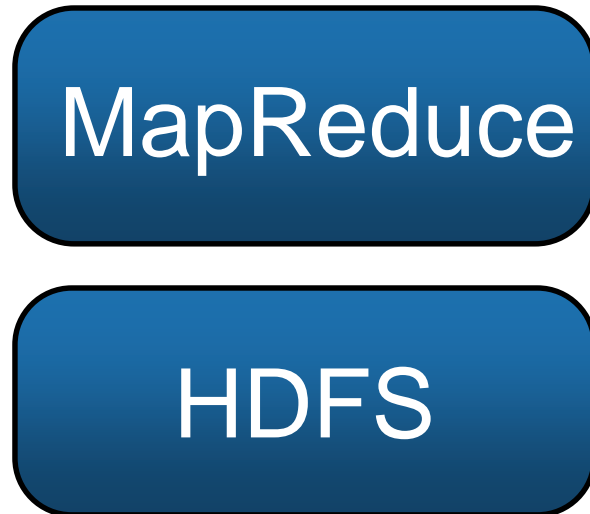
Hadoop amazing name!



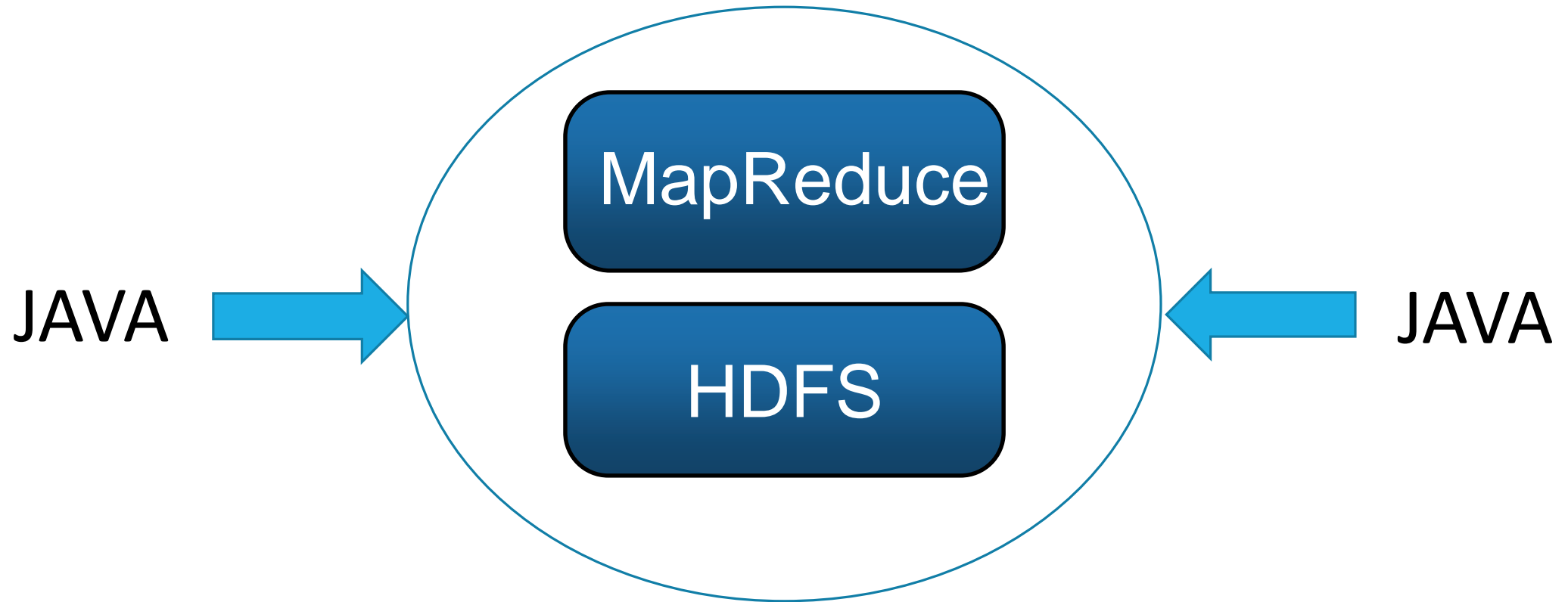
Hadoop Ecosystem!



Hadoop core

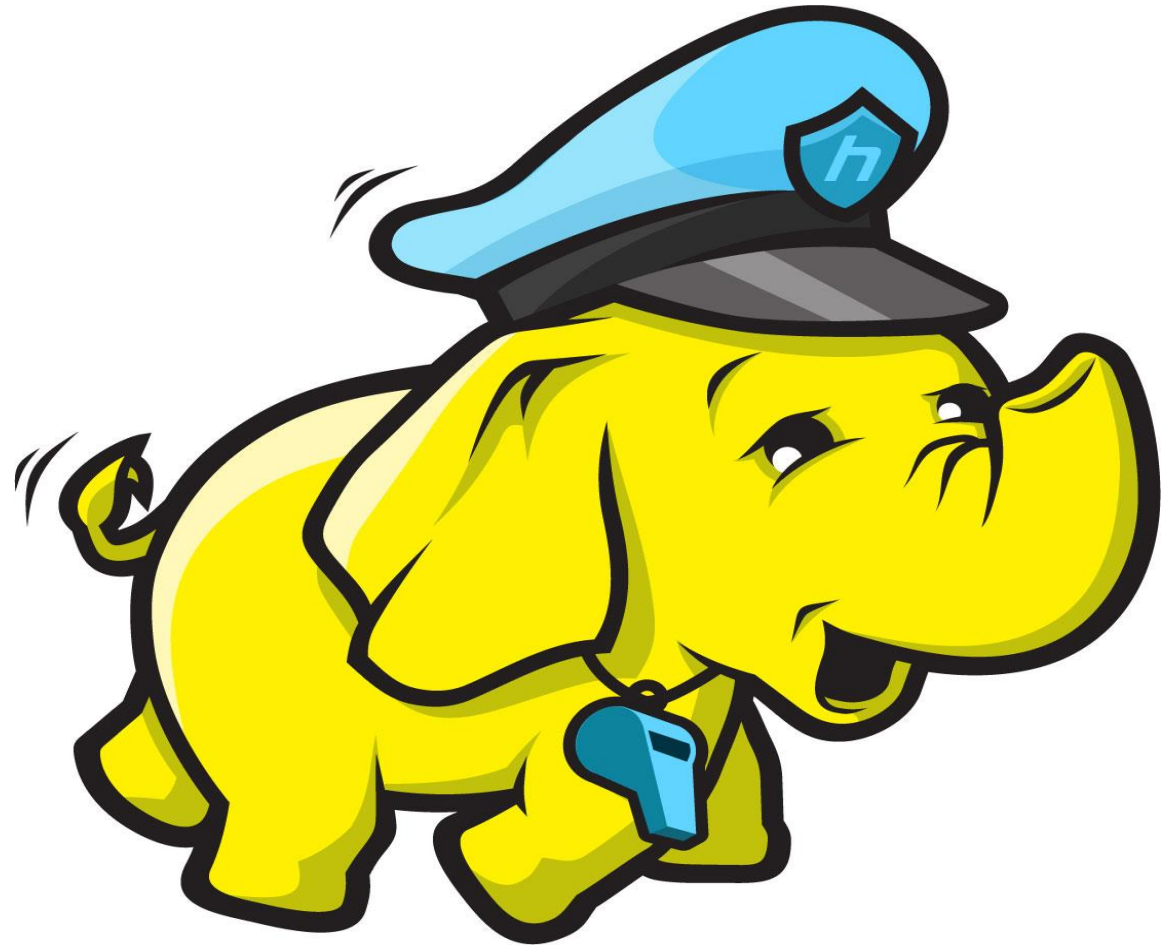


Hadoop core



What is HDFS?

- Hadoop Distributed File System
- Redundancy
- Fault Tolerant
- Scalable
- Self Healing
- Write Once, Read Many Times
- Java API
- Command Line Tool



Map Reduce



- Is a framework for performing high performance distributed data processing using the divide and aggregate programming paradigm.

Simple example

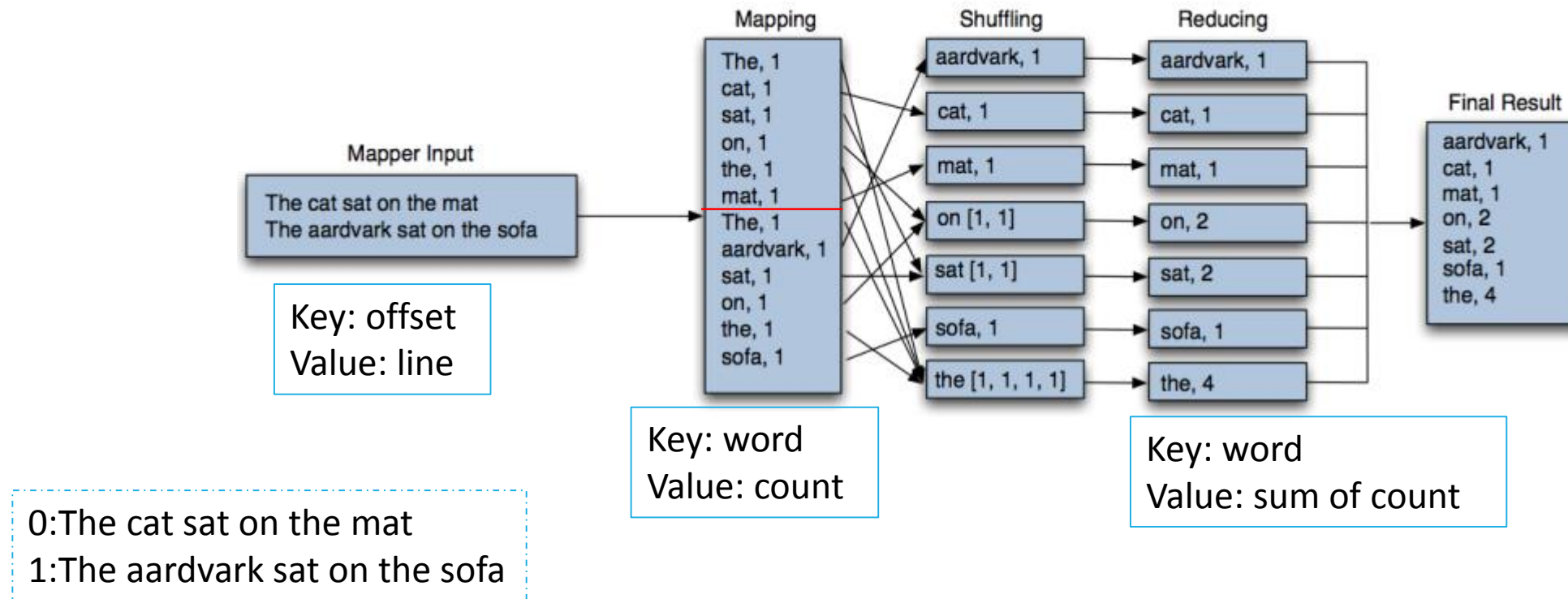
- Goal: count the number of books in the library.
- Map:
You count up shelf #1, I count up shelf #2
(The more people we get, the faster this part goes)
- Reduce:
We all get together and add up our individual counts.

(Cf. http://www.chrisstucchio.com/blog/2011/mapreduce_explained.html)

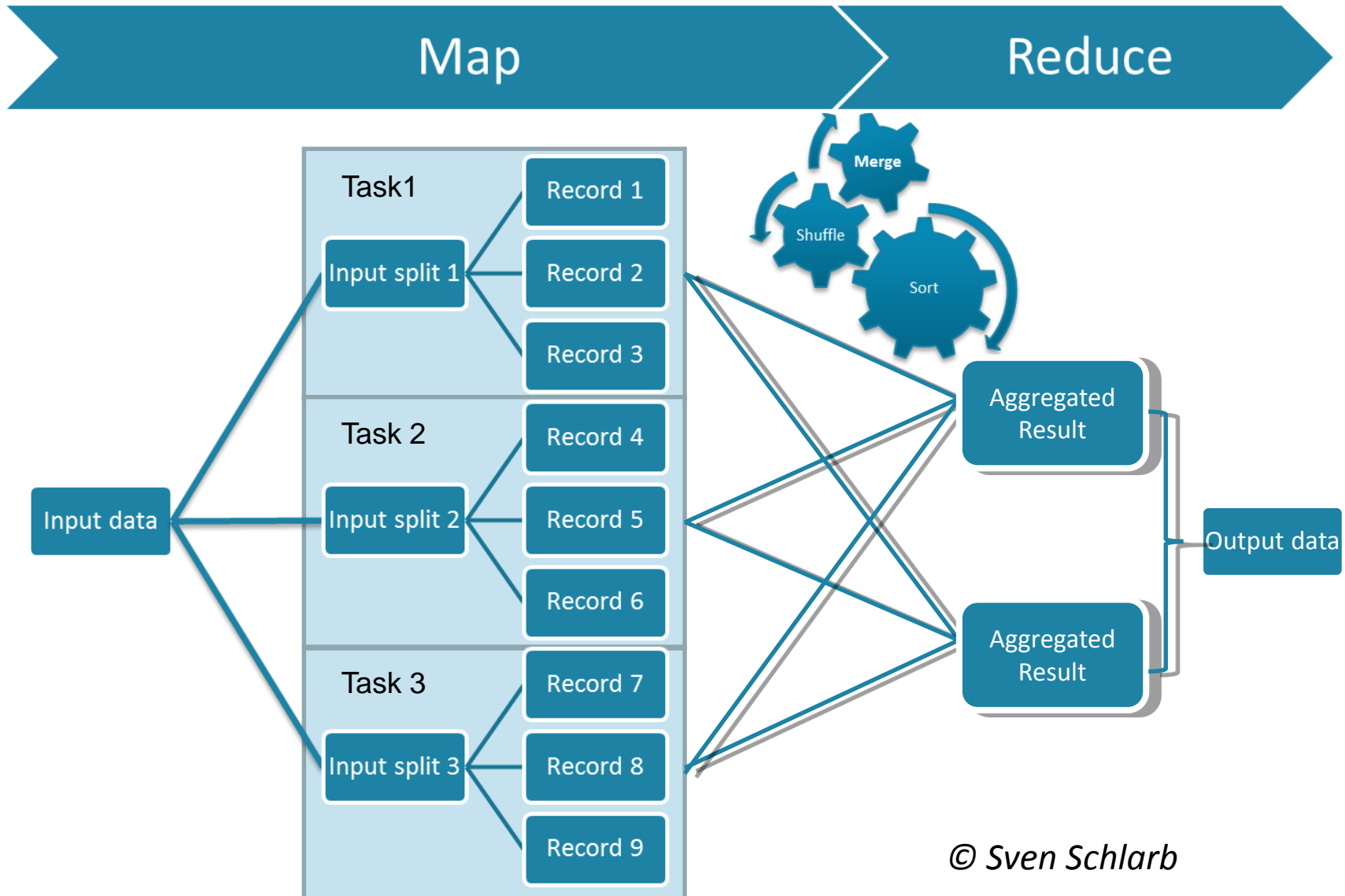


Routine example

The overall word count process

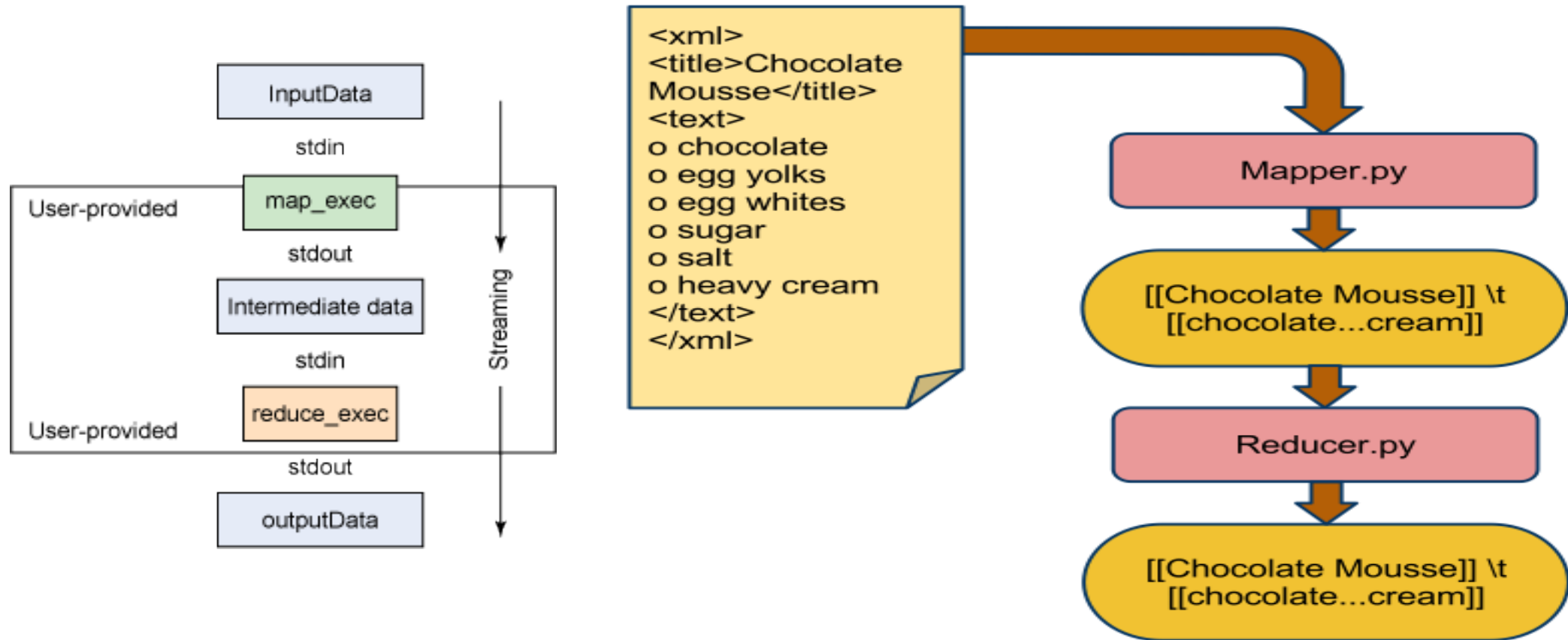


MapReduce in a nutshell



© Sven Schlarb

Hadoop streaming



Exercise

Goal : find
total sales of
each city!

1	2012-01-01	09:00	San Jose	Men's Clothing	214.05	Amex
2	2012-01-01	09:00	Fort Worth	Women's Clothing	153.57	Visa
3	2012-01-01	09:00	San Diego	Music	66.08	Cash
4	2012-01-01	09:00	Pittsburgh	Pet Supplies	493.51	Discover
5	2012-01-01	09:00	Omaha	Children's Clothing	235.63	MasterCard
6	2012-01-01	09:00	Stockton	Men's Clothing	247.18	MasterCard
7	2012-01-01	09:00	Austin	Cameras	379.6	Visa
8	2012-01-01	09:00	New York	Consumer Electronics	296.8	Cash
9	2012-01-01	09:00	Corpus Christi	Toys	25.38	Discover
10	2012-01-01	09:00	Fort Worth	Toys	213.88	Visa
11	2012-01-01	09:00	Las Vegas	Video Games	53.26	Visa
12	2012-01-01	09:00	Newark	Video Games	39.75	Cash
13	2012-01-01	09:00	Austin	Cameras	469.63	MasterCard
14	2012-01-01	09:00	Greensboro	DVDs	290.82	MasterCard
15	2012-01-01	09:00	San Francisco	Music	260.65	Discover
16	2012-01-01	09:00	Lincoln Garden		136.9	Visa
17	2012-01-01	09:00	Buffalo	Women's Clothing	483.82	Visa
18	2012-01-01	09:00	San Jose	Women's Clothing	215.82	Cash

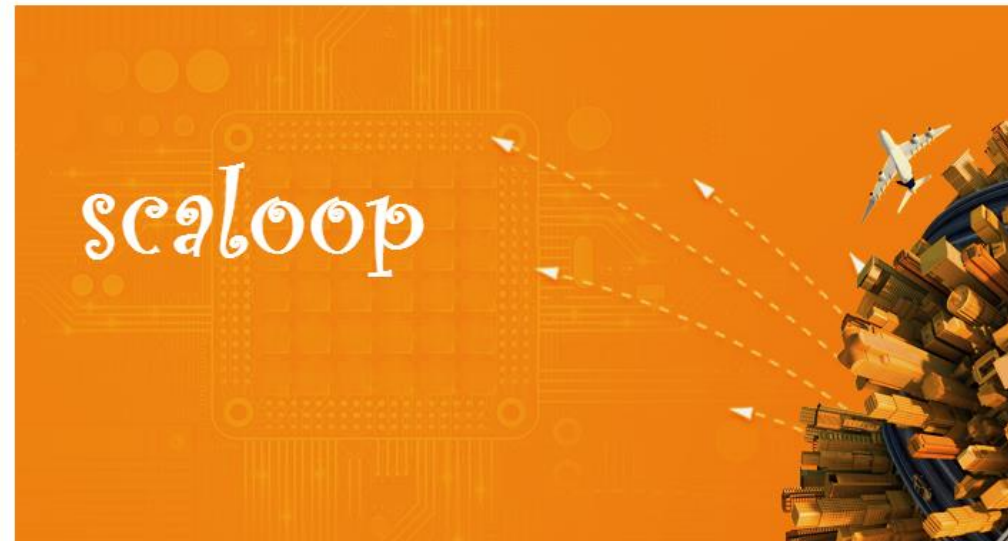
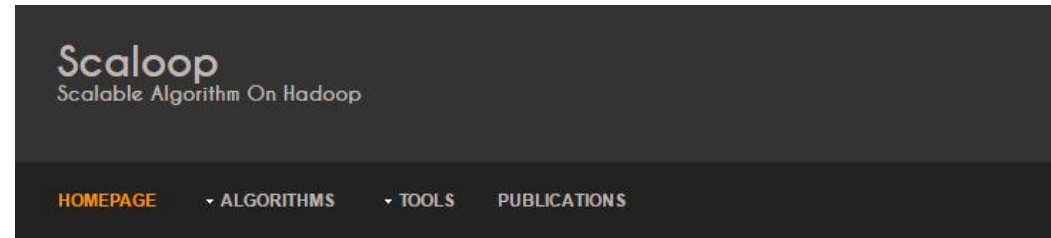
Mapper

```
1 #!/usr/bin/env python
2 import sys
3
4 def mapper():
5     for line in sys.stdin:
6         data = line.strip().split("\t")
7         if len(data)==6:
8             date,time,store,item,cost,payment = data
9             print "{0}\t{1}".format(store,cost)
10
11
12 mapper()
13
```


Reducer

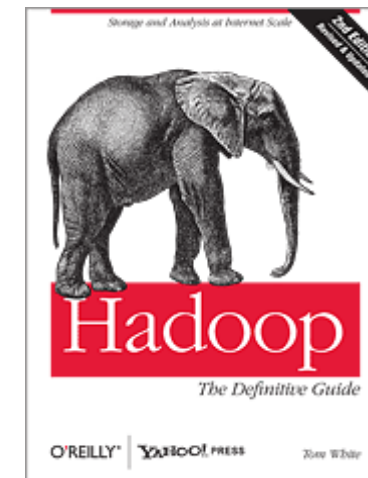
```
1 #!/usr/bin/env python
2 import sys
3
4 def reducer():
5     salesTotal = 0
6     oldKey = None
7     for line in sys.stdin:
8         data = line.strip().split("\t")
9         if len(data) != 2:
10             continue
11         thisKey, thisSale = data
12         if oldKey and oldKey != thisKey:
13             print "{0}\t{1}".format(oldKey, salesTotal)
14             salesTotal = 0
15         oldKey = thisKey
16         salesTotal += float(thisSale)
17
18
19 reducer()
```

Scalooop Story



Resources

- ❖ Tom White: Hadoop. The Definitive Guide
- ❖ Udacity Big Data course
- ❖ www.vcloudnews.com



Questions?

