

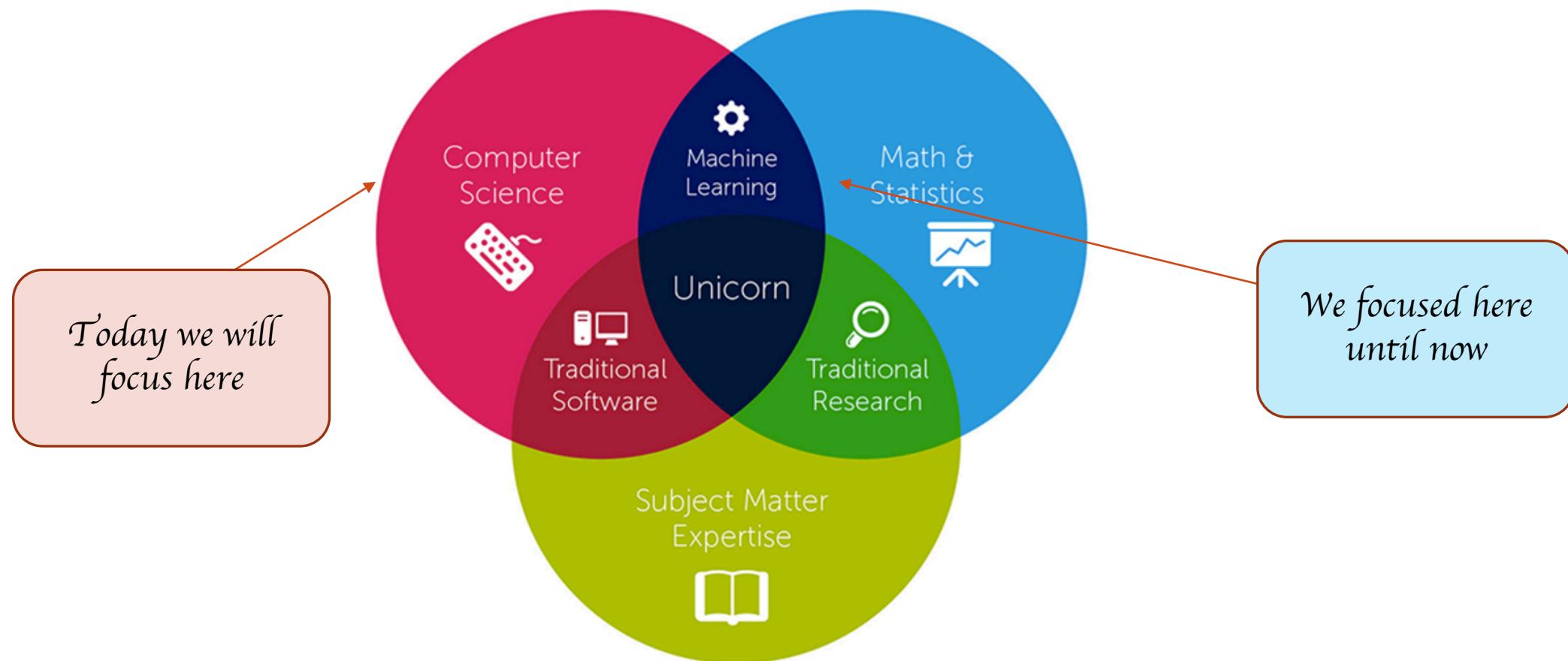


Introduction to Big Data

DSLA COURSE

ROHIT PADEBETTU

Who is a Data Scientist?



Time Gained in Flight

Summary

Model Details

Source Code

Select Airports

Flight origin

JFK

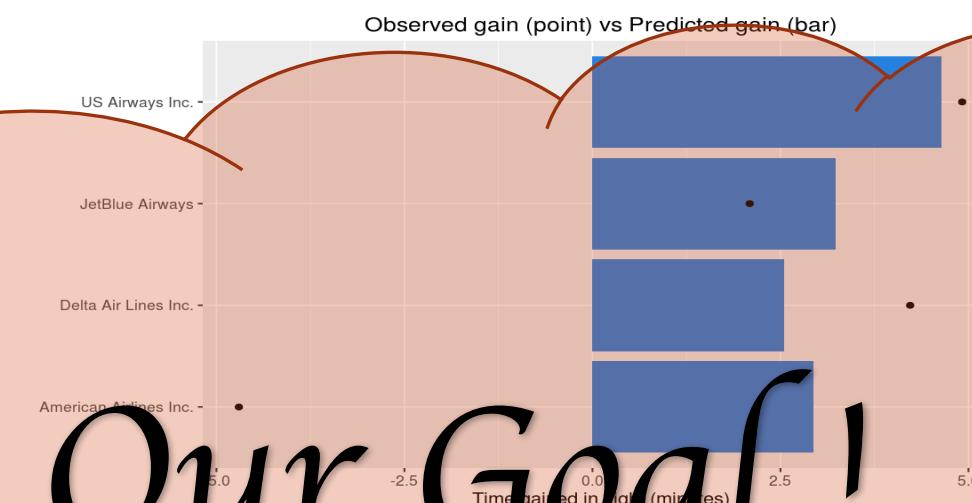
Flight destination

LAS

Background

Given that your flight was delayed by 15 minutes or more, what is the likelihood your airline carrier will make up time in route? Some of the most significant factors for making up time are flight distance and airline carrier. The data model behind this dashboard is based on flights from NYC airports in 2013.

Observed versus predicted time gain



Data details

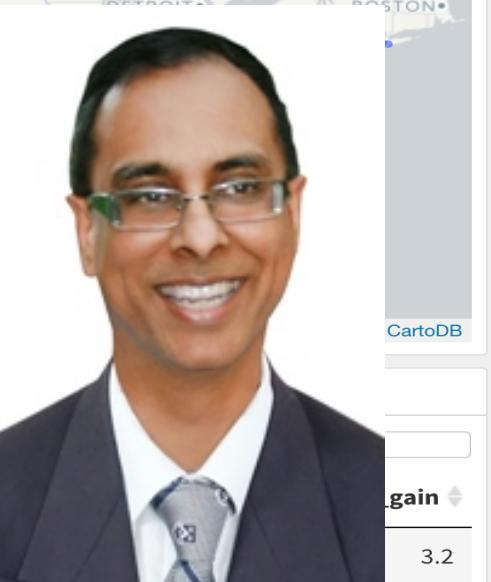
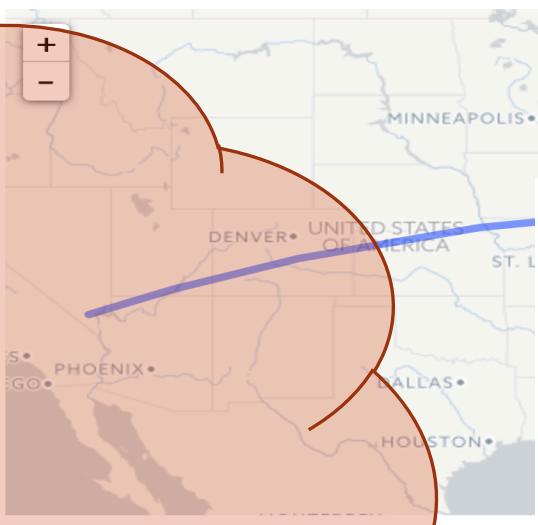
Show 10 entries

	airline	flights	distance	avg_dep_delay	avg_arr_delay	gain
1	JetBlue Airways	286	2248	61.7	59.6	3.2
2	American Airlines Inc.	97	2248	57.2	61.9	-4.7
3	Delta Air Lines Inc.	227	2248	58.7	54.4	4.2
4	US Airways Inc.	192	2248	57.4	52.5	4.9

Showing 1 to 4 of 4 entries

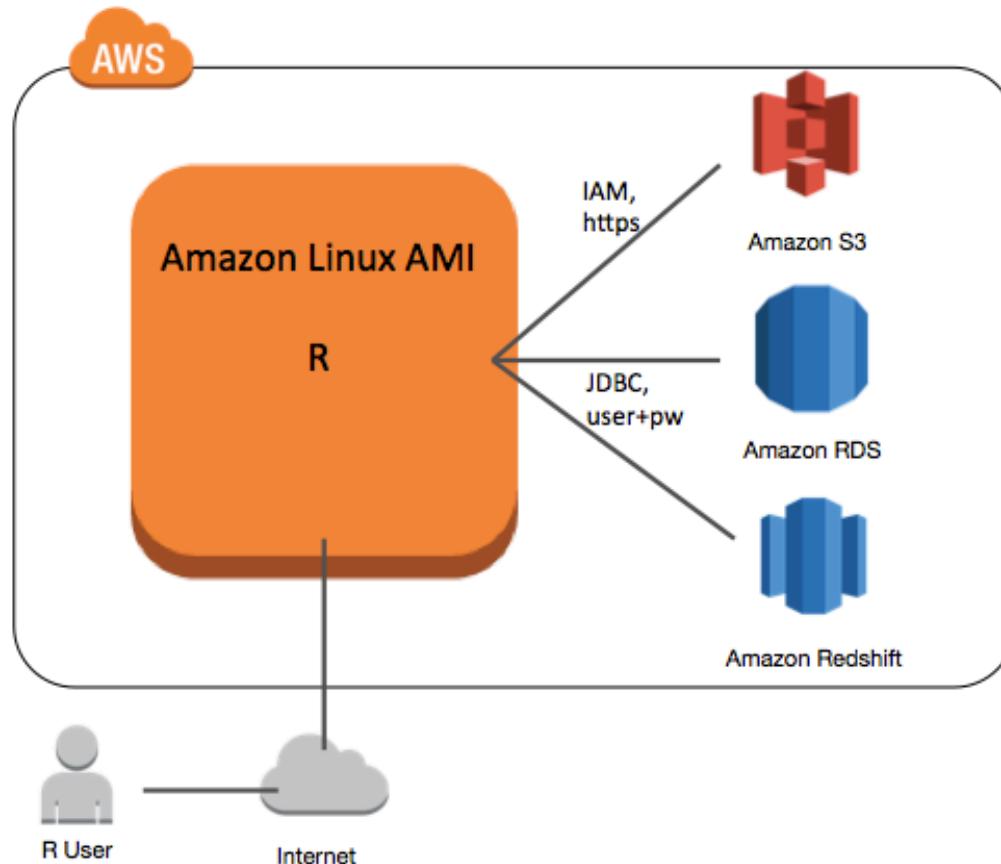
Previous 1 Next

Route



CartoDB

AWS Cloud Architecture



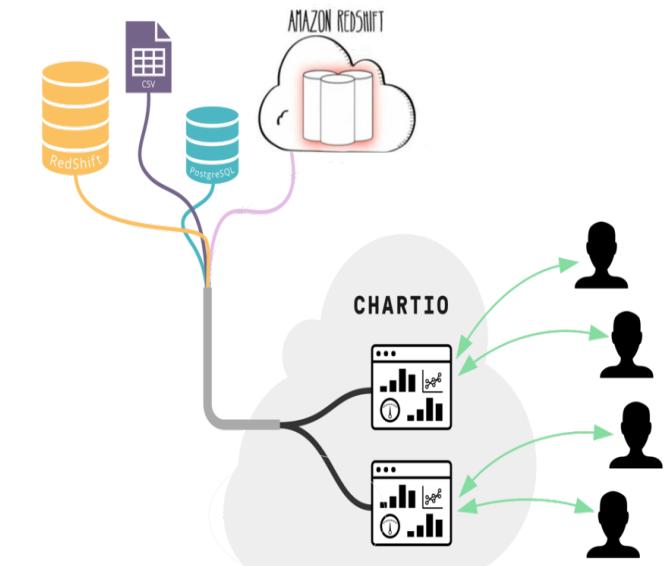
Yotta (Y)	10^{24}	1 septillion
Zetta (Z)	10^{21}	1 sextillion
Exa (E)	10^{18}	1 quintillion
Peta (P)	10^{15}	1 quadrillion
Tera (T)	10^{12}	1 trillion
Giga (G)	10^9	1 billion
Mega (M)	10^6	1 million
kilo (k)	10^3	1 thousand

Most modern databases can manage approximately $6T^3B$ of data

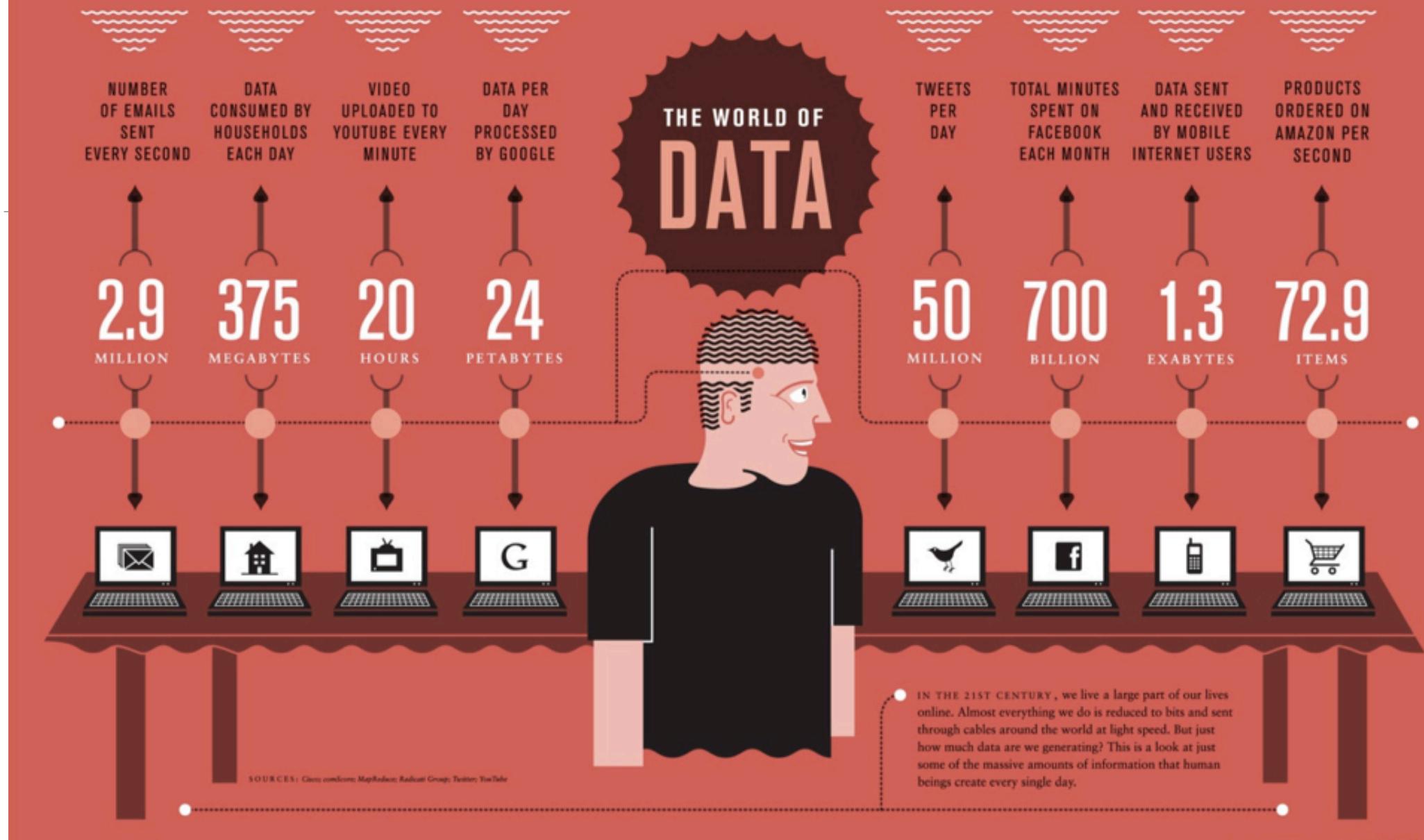
Data in all Shapes & Sizes



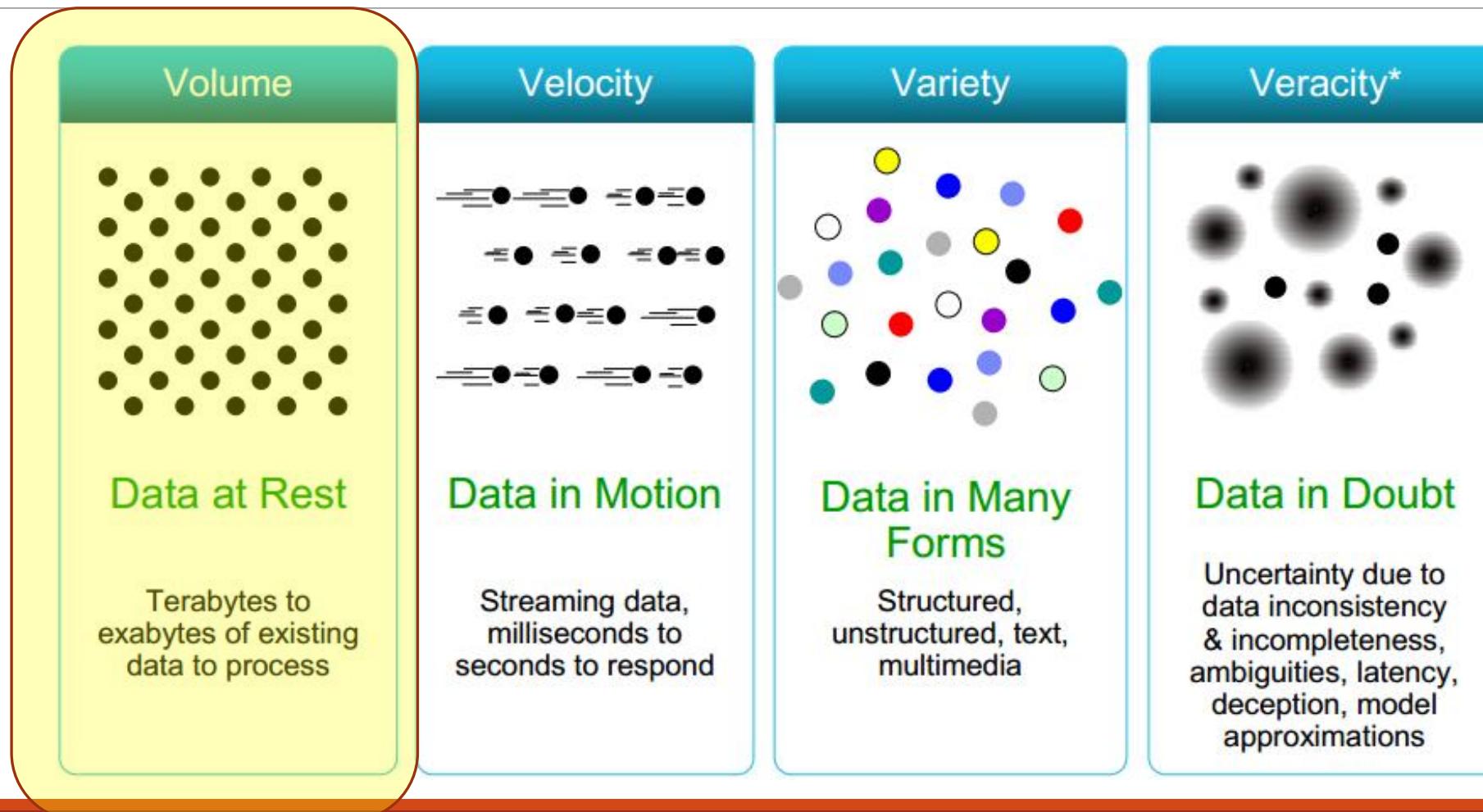
A Data Warehouse is a “*Relational Database*” that is used primarily for query and analysis rather than transaction processing.
(OLAP vs OLTP)



Amazon Redshift is branded as a Petabyte scale database



Characteristics of Data

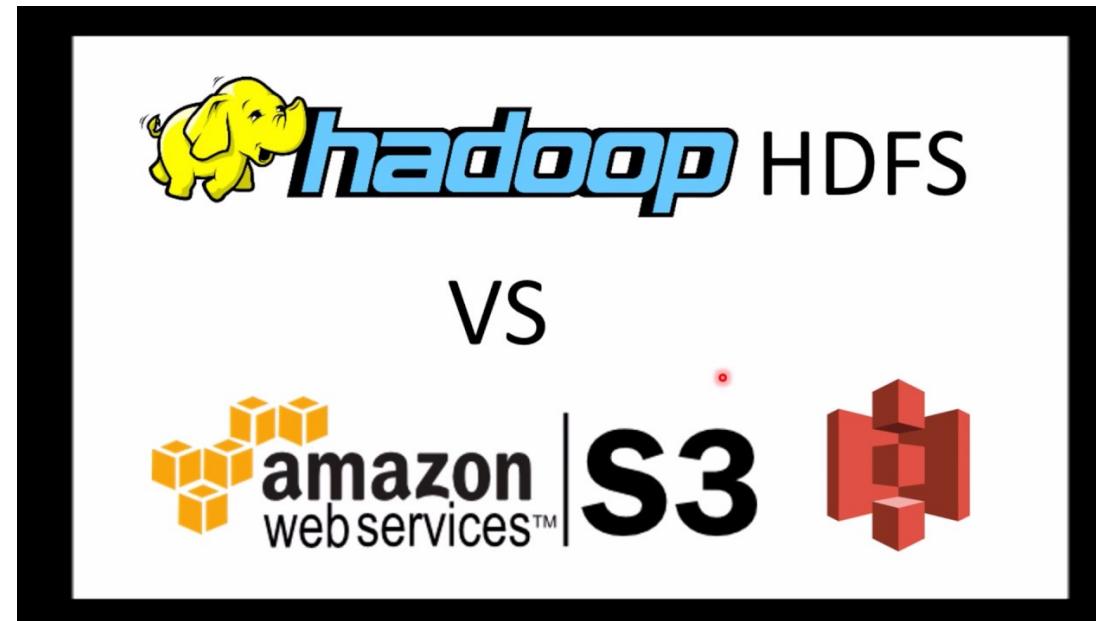


Big Data Will Scale To Exabytes

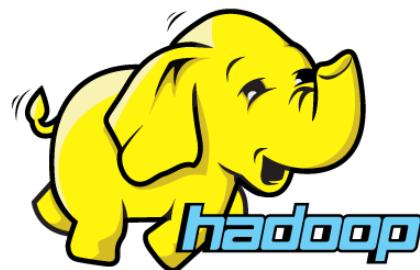


Store Big Data

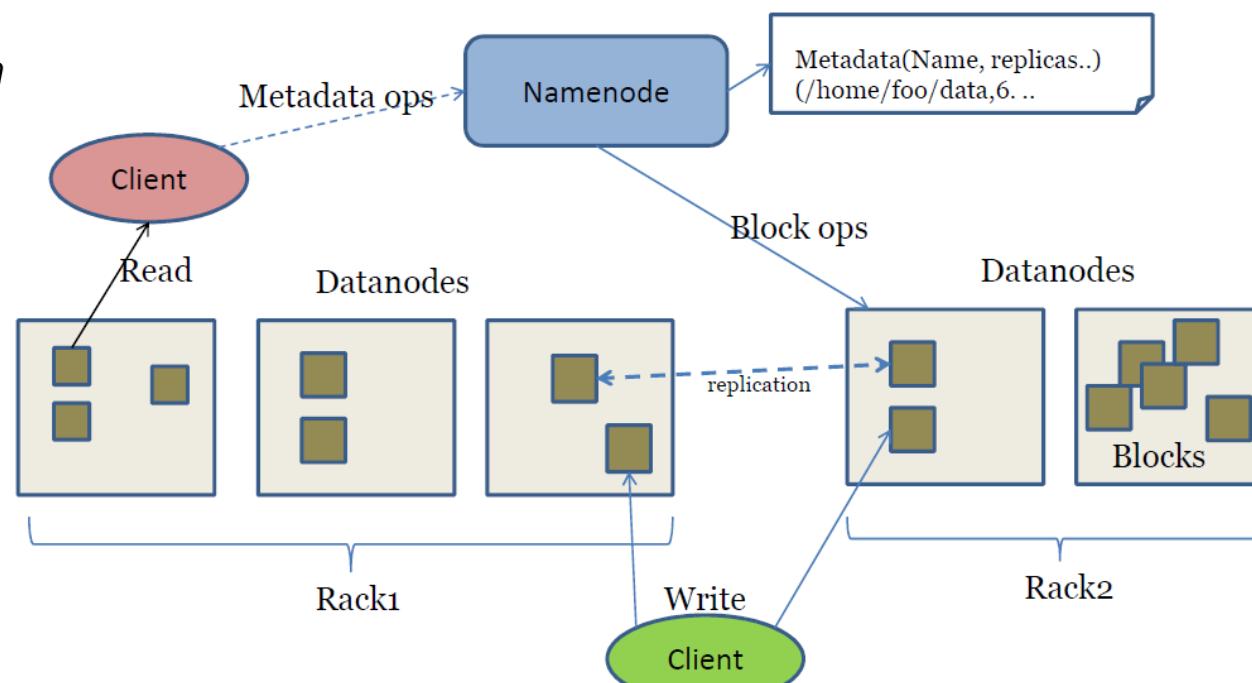
How ?



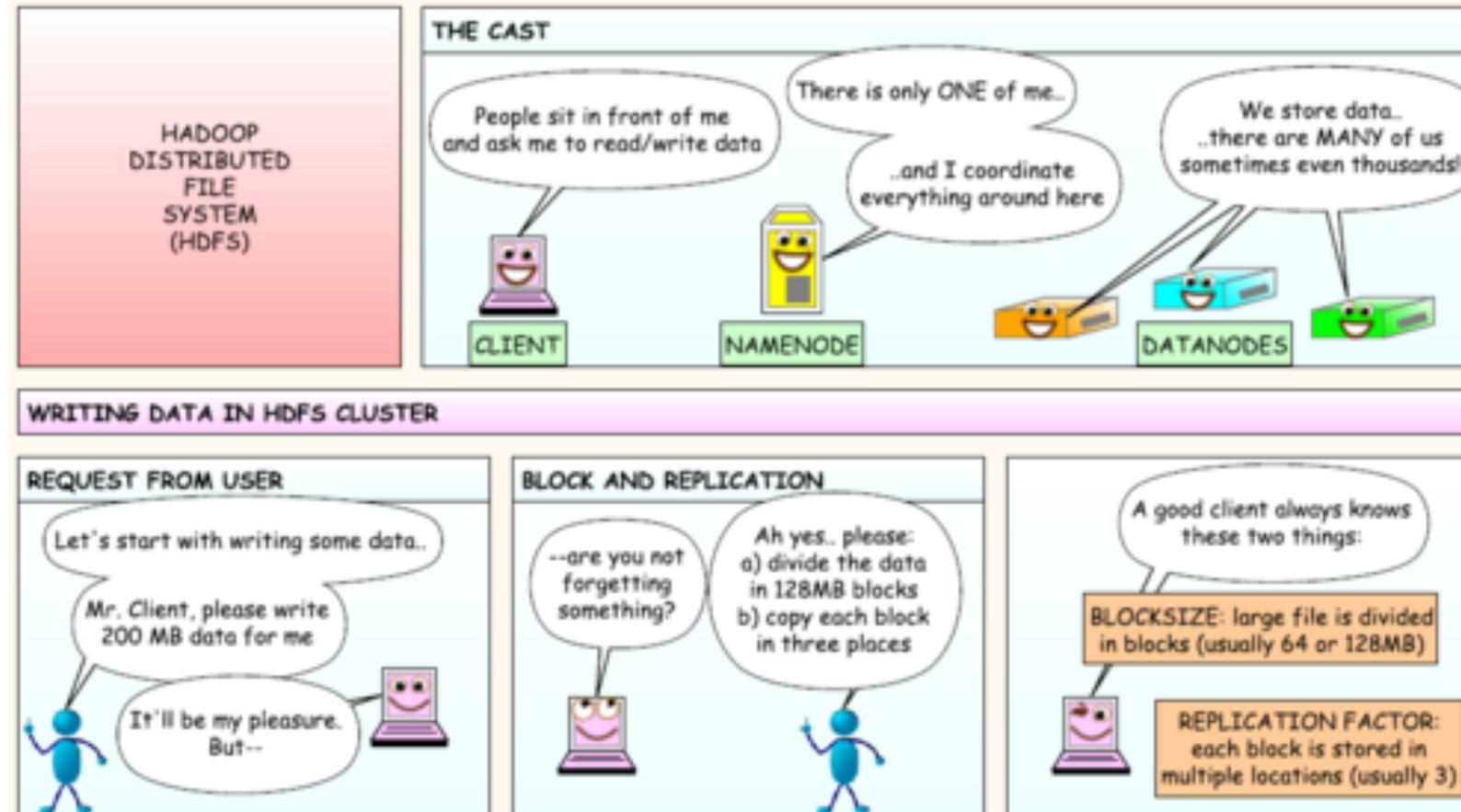
Distributed Storage



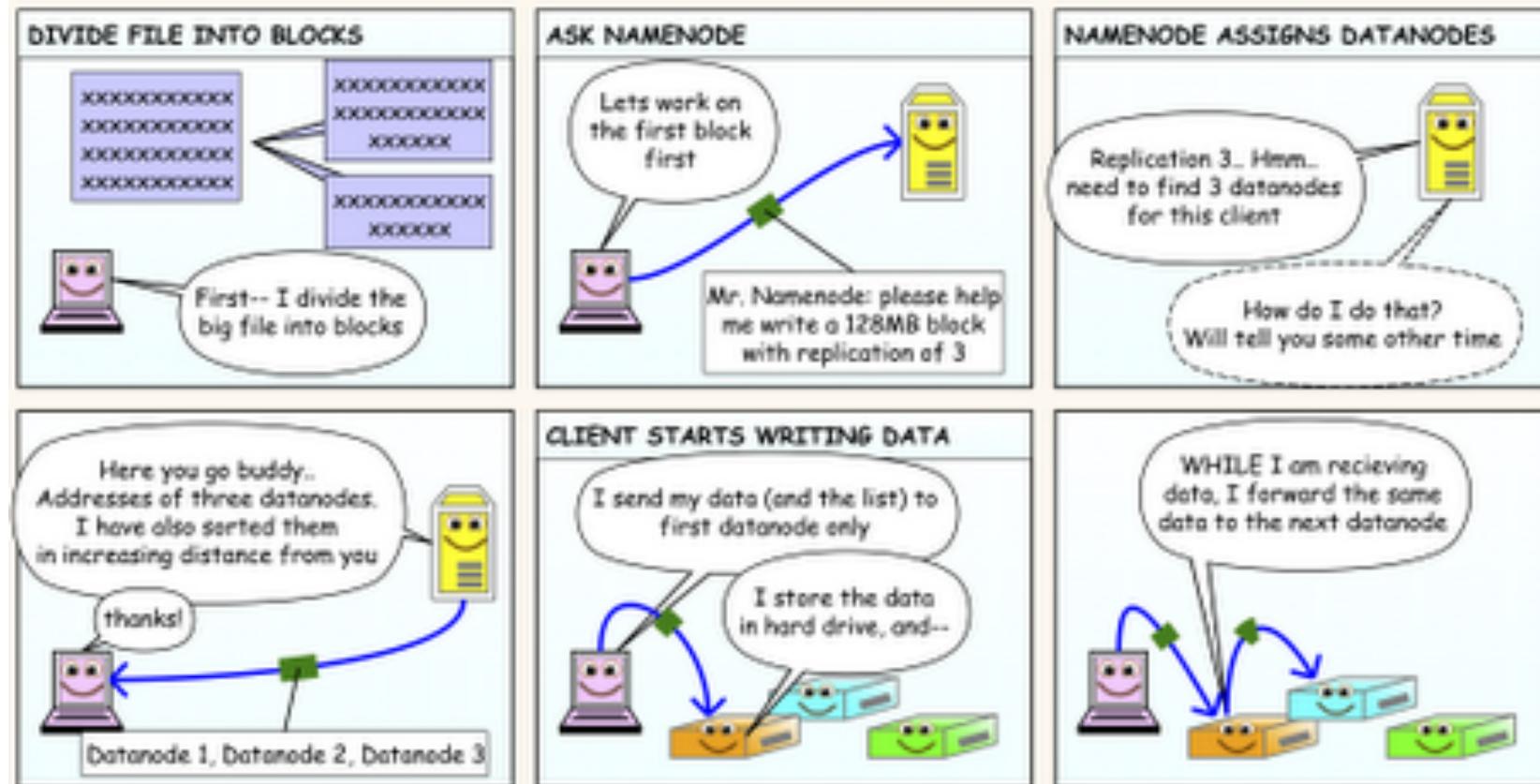
HDFS Architecture



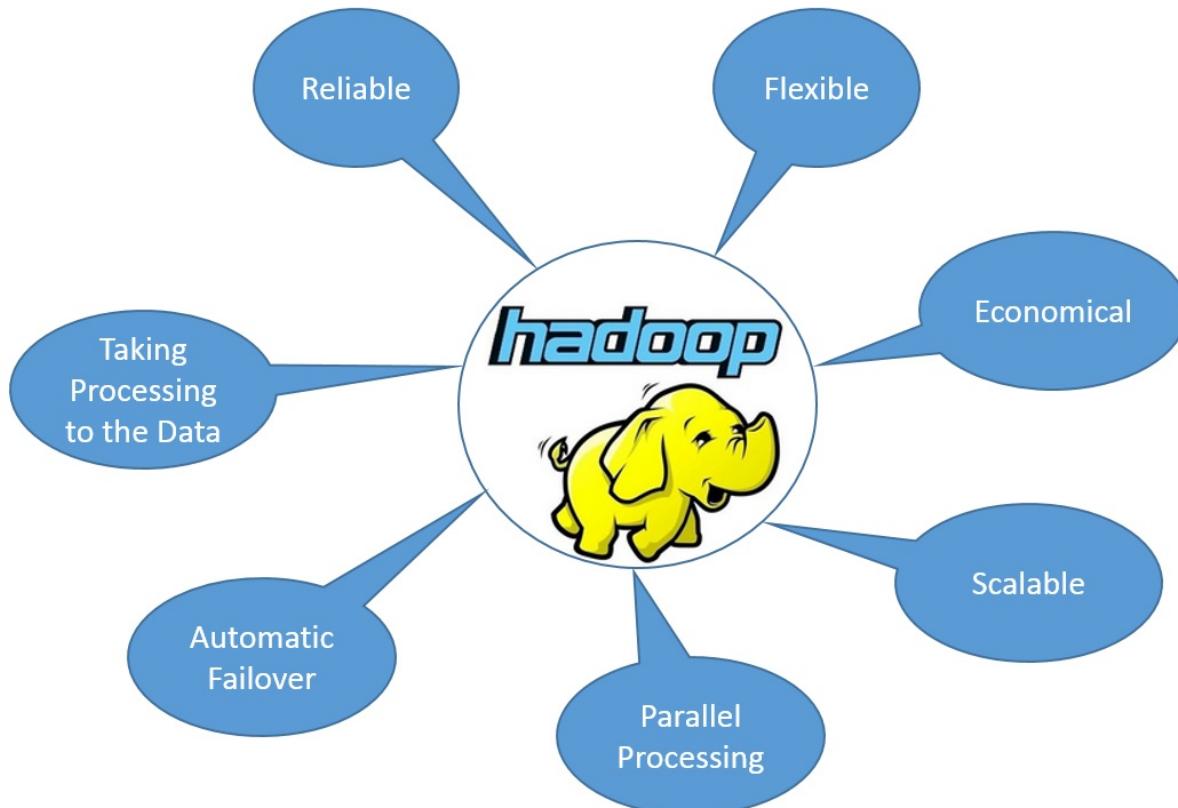
Distributed Storage



Distributed Storage



Hadoop Advantages



Access Distributed Data

How?

MapReduce: Simplified Data Processing on Large Clusters

Jeffrey Dean and Sanjay Ghemawat

jeff@google.com, sanjay@google.com

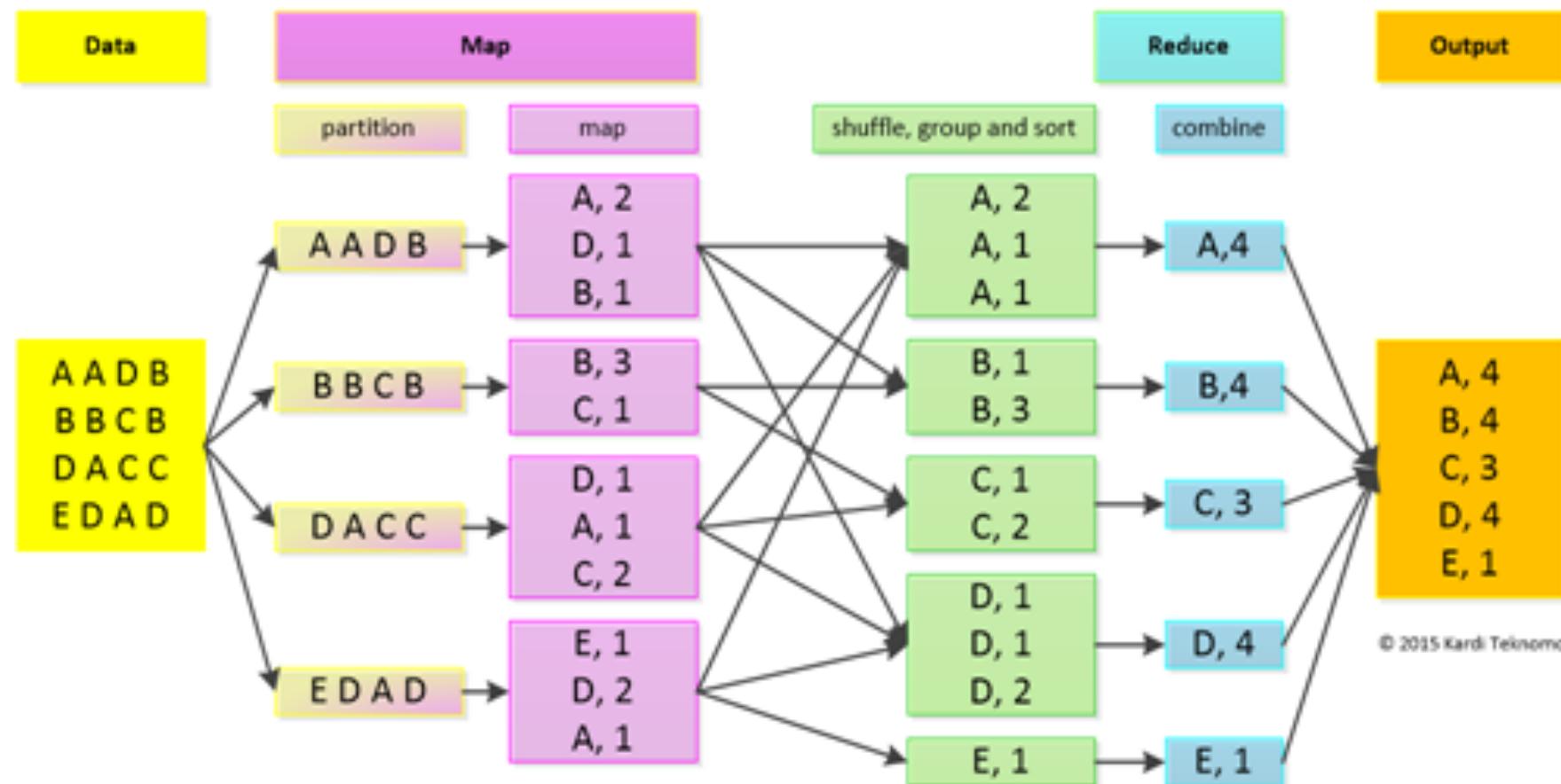
Google, Inc.

Abstract

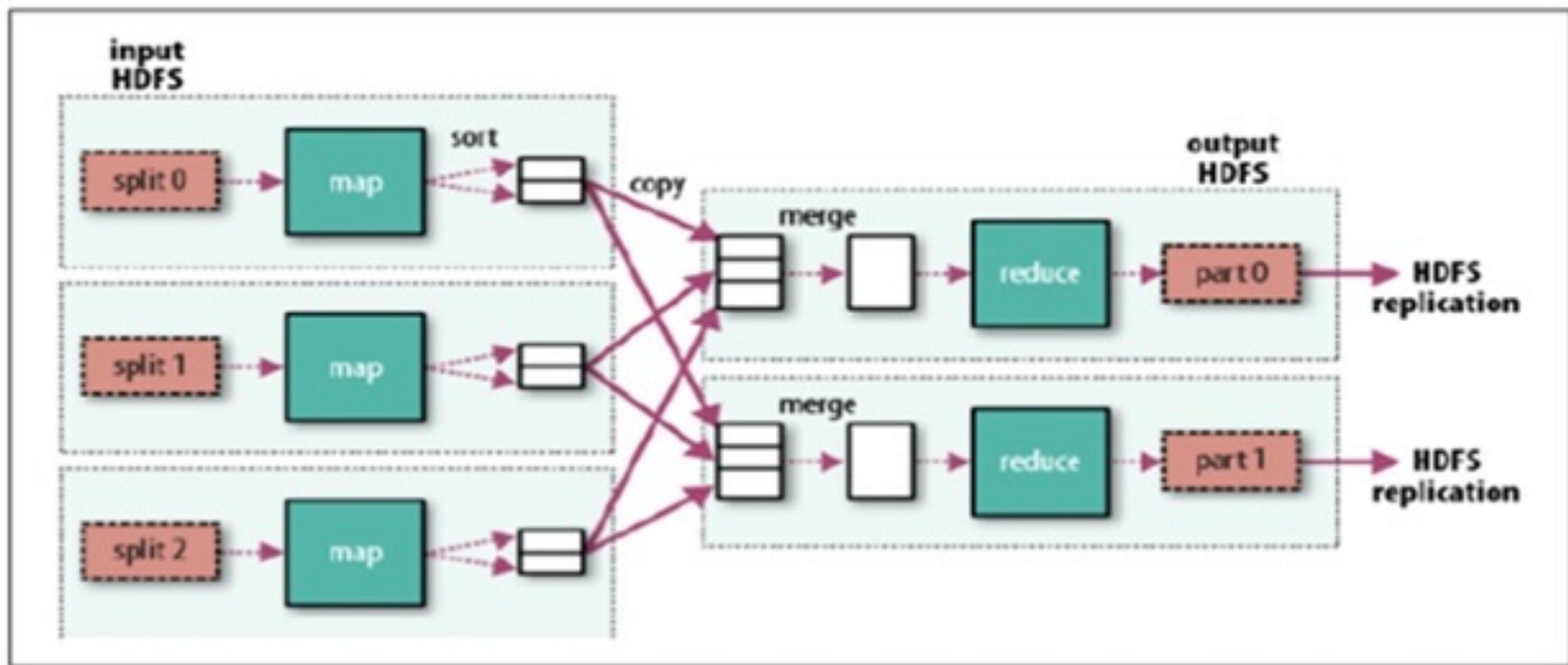
MapReduce is a programming model and an associated implementation for processing and generating large data sets. Users specify a *map* function that processes a key/value pair to generate a set of intermediate key/value

given day, etc. Most such computations are conceptually straightforward. However, the input data is usually large and the computations have to be distributed across hundreds or thousands of machines in order to finish in a reasonable amount of time. The issues of how to parallelize the computation, distribute the data, and handle

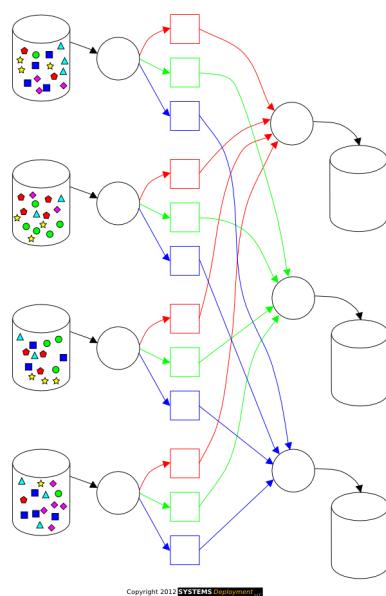
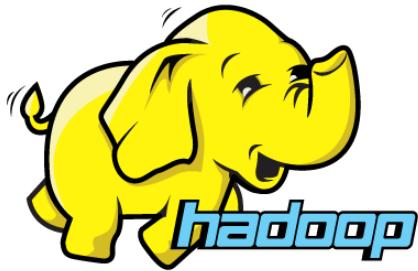
Map Reduce



Map Reduce on HDFS



Hadoop + Map Reduce



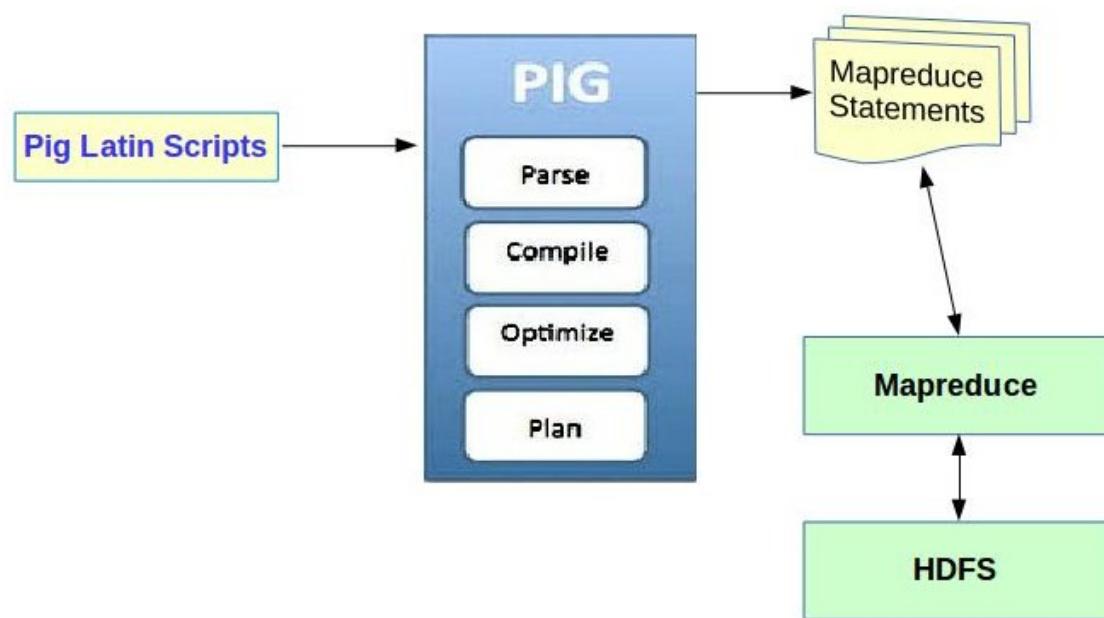
	RDBMS	Hadoop
Data size	Gigabytes (<i>terabytes</i>)	Petabytes (<i>hexabytes</i>)
Access	Interactive and batch	Batch
Updates	Read and write many times	Write once, read many times
Structure	Static schema	Dynamic schema
Integrity	High (ACID)	Low
Scaling	Nonlinear	Linear



Converting common and complex algorithms into MapReduce Framework involves many steps.
It is **very slow and cumbersome**

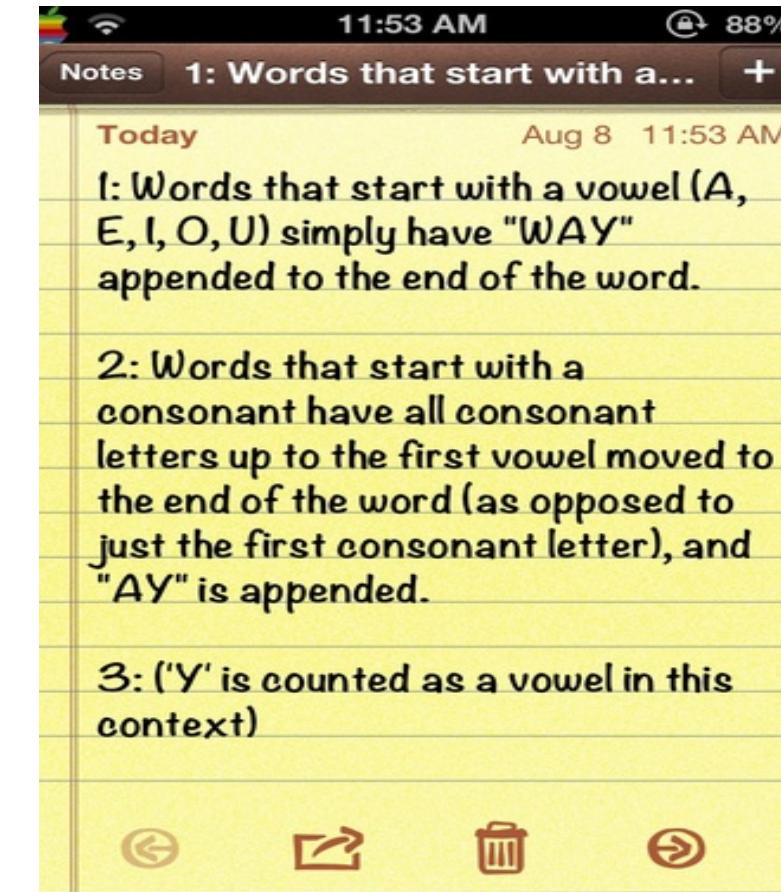
Pig

Apache Pig is a platform for analyzing large data sets that consists of a high-level language for expressing data analysis programs, coupled with infrastructure for evaluating these programs.



Pig Latin

*Apache Pig's Language layer is called **Pig Latin**.*



 A screenshot of an iPhone displaying a note titled "1: Words that start with a...". The note contains three numbered rules about Pig Latin:

- 1: Words that start with a vowel (A, E, I, O, U) simply have "WAY" appended to the end of the word.
- 2: Words that start with a consonant have all consonant letters up to the first vowel moved to the end of the word (as opposed to just the first consonant letter), and "AY" is appended.
- 3: ('Y' is counted as a vowel in this context)

 The note is timestamped "Today Aug 8 11:53 AM" and shows the phone's status bar at the top with the time "11:53 AM" and battery level "88%".

Pig Latin

Pig Latin

```

countryrs = load '/user/gharriso/PIG_COUNTRIES' AS
  (country_id, country_name , country_subregion , region);

customers= load '/user/gharriso/PIG_CUSTOMERS' AS
  (cust_id,first_name, last_name, gender, yob, marital, postcode,city,country_id);

asianCountryrs = filter countryrs by region matches 'Asia';

joined = join customers by country_id, asianCountryrs by country_id;

grouped = group joined by country_name;

agged = foreach grouped generate group, COUNT(joined.customers::cust_id);

morethan500cust = filter agged by $1 > 500;

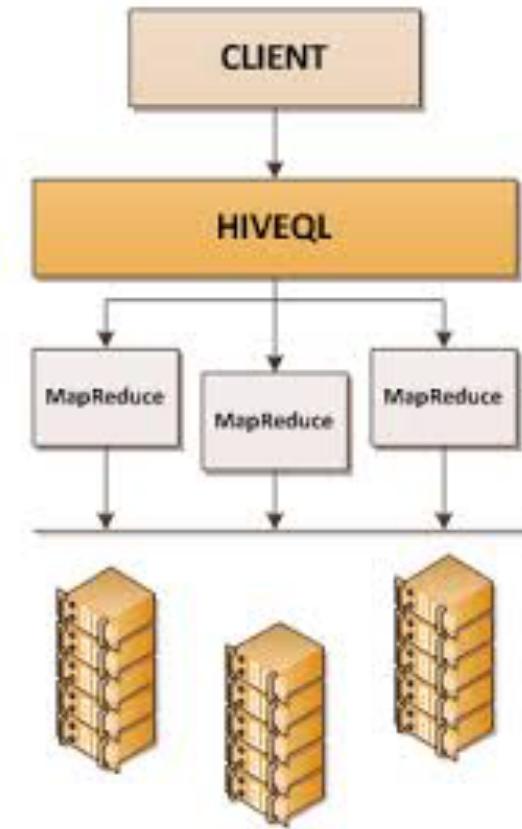
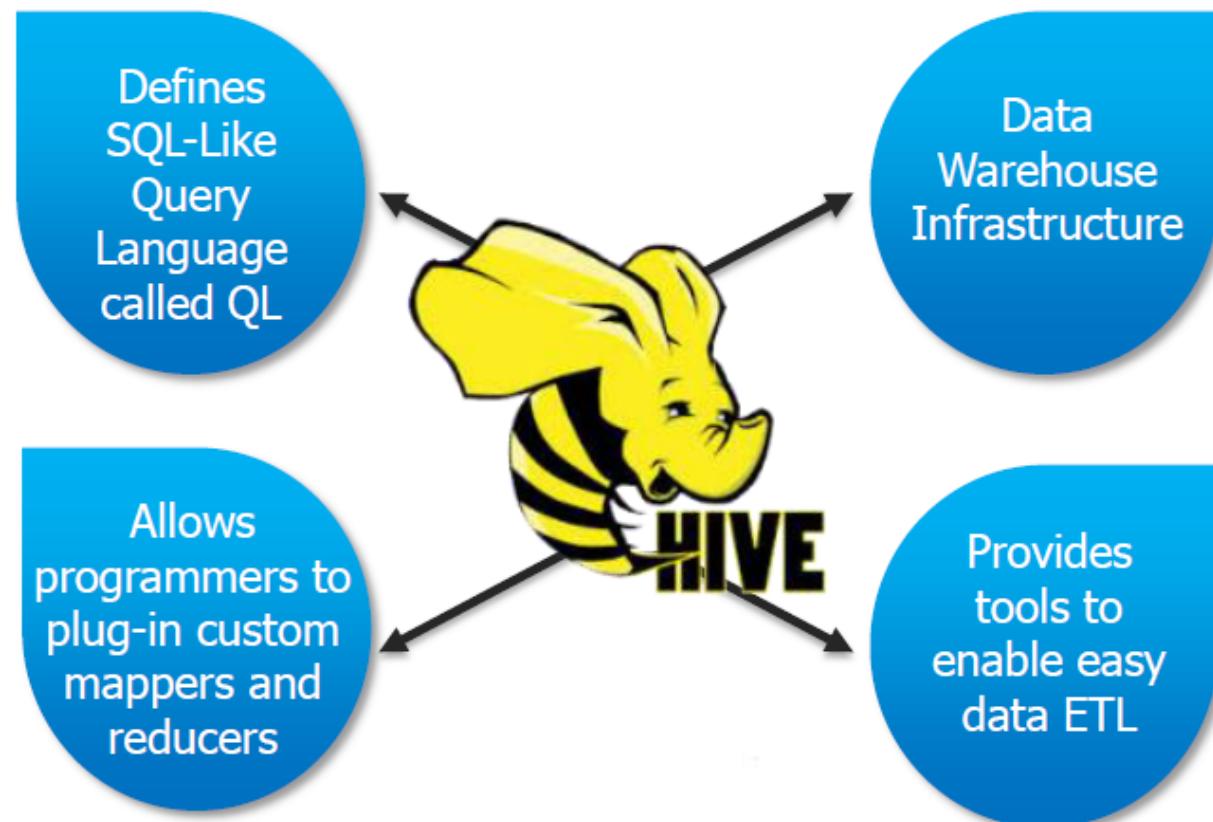
ordered =order morethan500cust by $1 desc;

dump ordered;

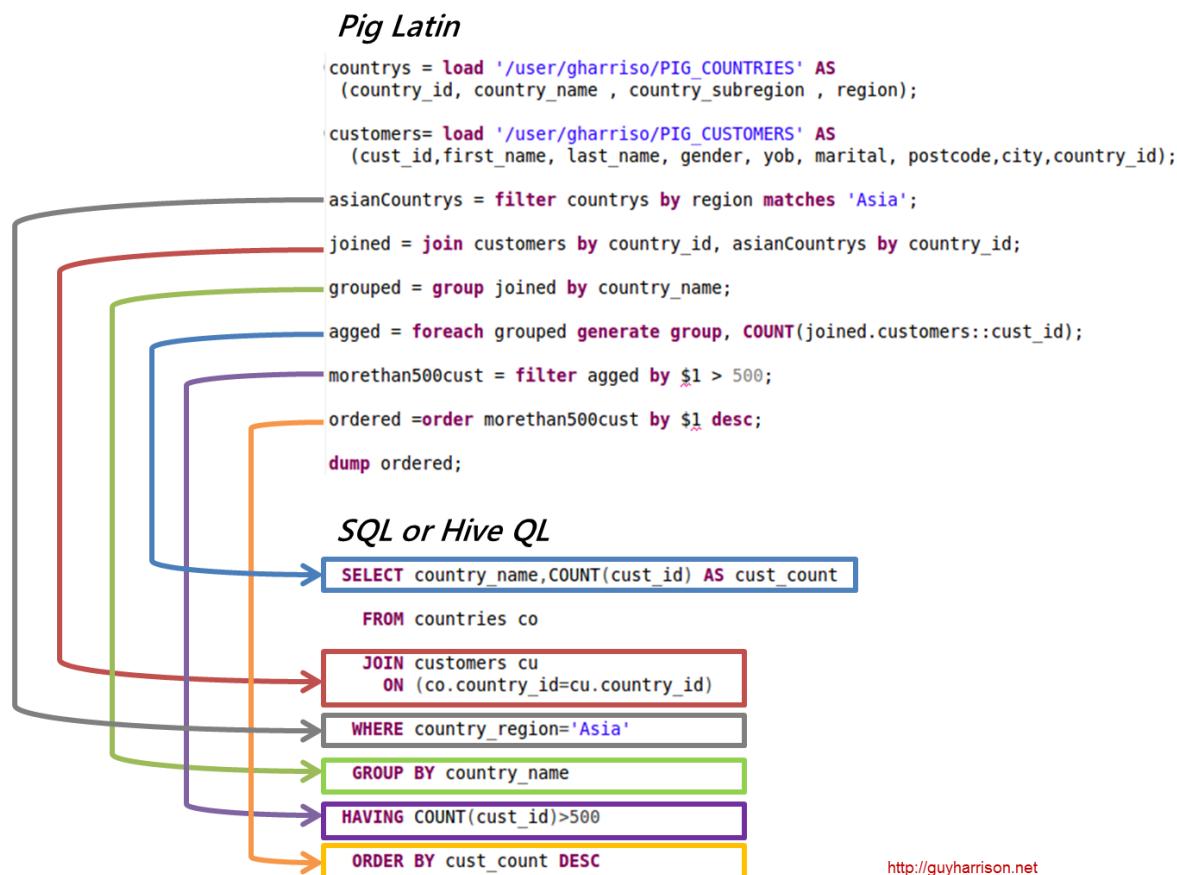
```

- *No Java knowledge required*
- *Easy & Fast to write compared to Map Reduce*
- *Slow to run compared to Map Reduce because of translation*
- *New language to Master, but not intuitive like SQL*

HIVE



HIVE Query Language



```

hive> desc emp;
OK
id                      int
name                     string
salary                   int
Time taken: 0.483 seconds, Fetched: 3 row(s)
hive> select * from emp;
OK
1      amit     80000
2      rakesh   35000
3      vishal   30000
4      rahul    50000
5      sameer   20000
1      ravi     30000
1      savi     26000
2      vikas    30000

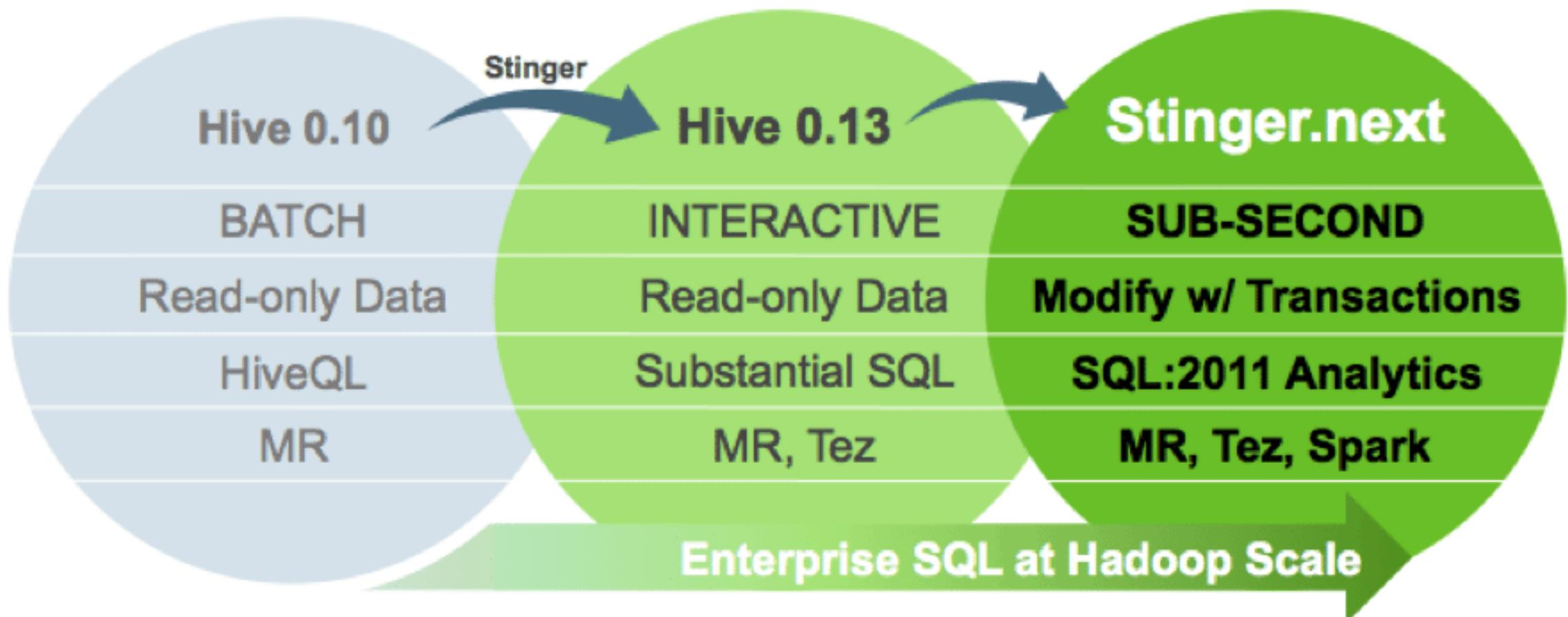
```

MapReduce vs Pig vs HIVE

Hadoop MapReduce Vs Pig Vs Hive

Hadoop MapReduce	Pig	Hive
Compiled Language	Scripting Language	SQL like query Language
Lower Level of Abstraction	Higher Level of Abstraction	Higher Level of Abstraction
More lines of Code	Comparatively less lines of Code than MapReduce	Comparatively less lines of Code than MapReduce and Apache Pig
More Development Effort is involved	Development Effort is less Code Efficiency is relatively less	Development Effort is less Code Efficiency is relatively less
Code Efficiency is high when compared to Pig and Hive	Code Efficiency is relatively less	Code Efficiency is relatively less

HIVE



Analysis

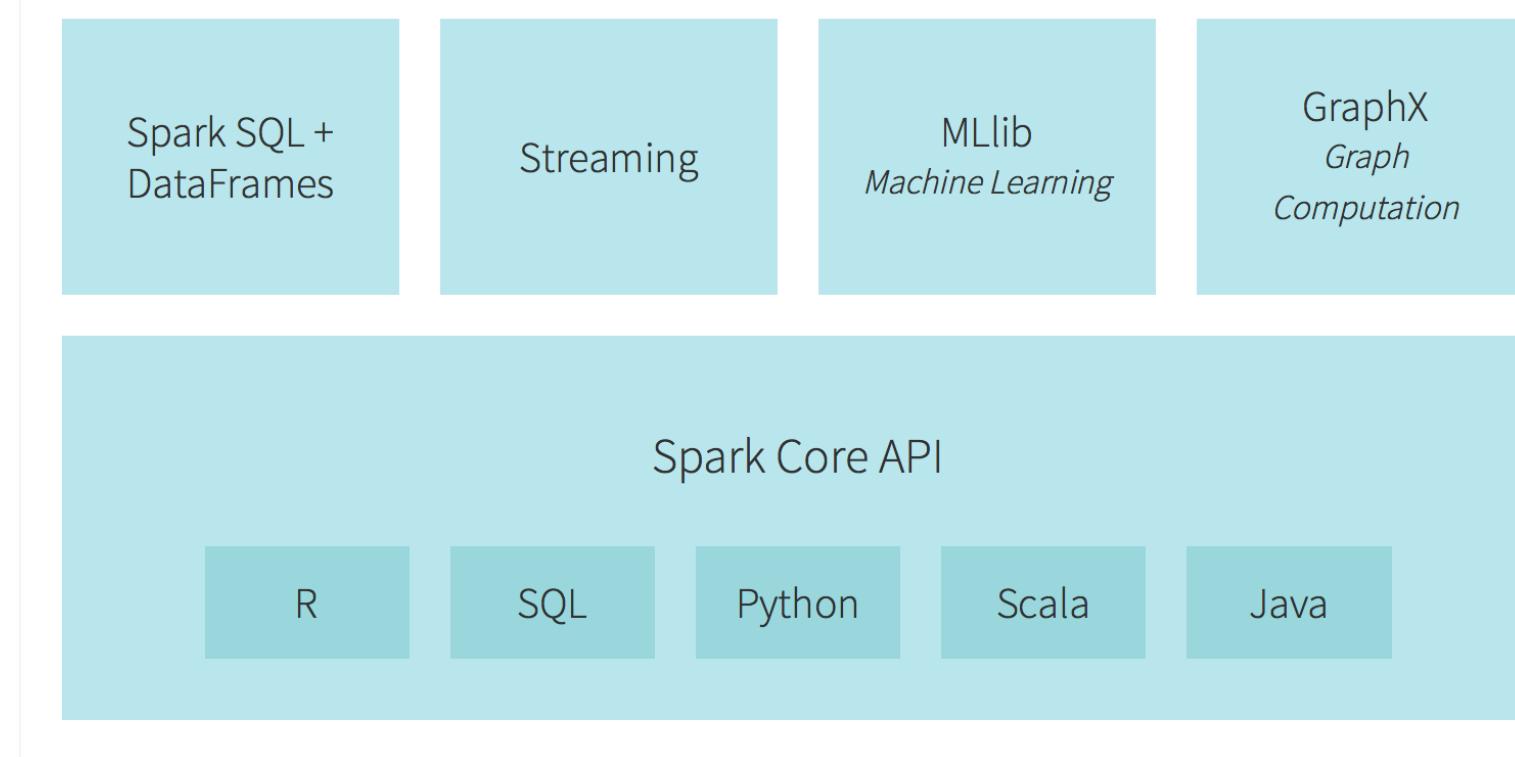


Batch Processing

How ?

*Real Time Processing
Interactive Analysis
Machine Learning
Graph Processing*

Apache Spark



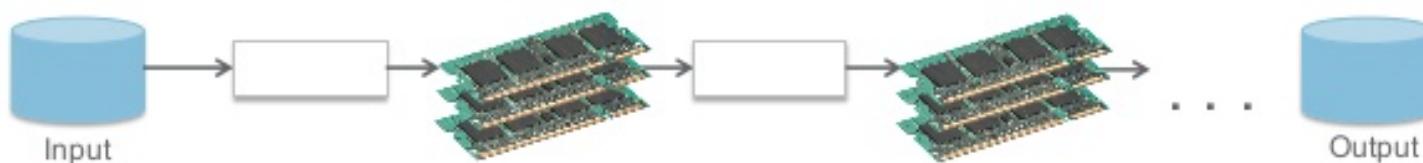
Apache Spark

Apache Spark utilizes in-memory caching and optimized execution for fast performance

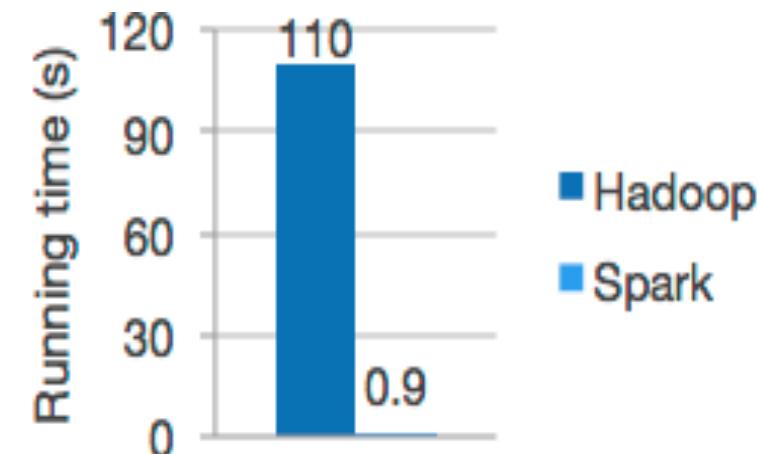
Hadoop MapReduce: Data Sharing on Disk



Spark: Speed up processing by using Memory instead of Disks



Logistic Regression in Hadoop and Spark

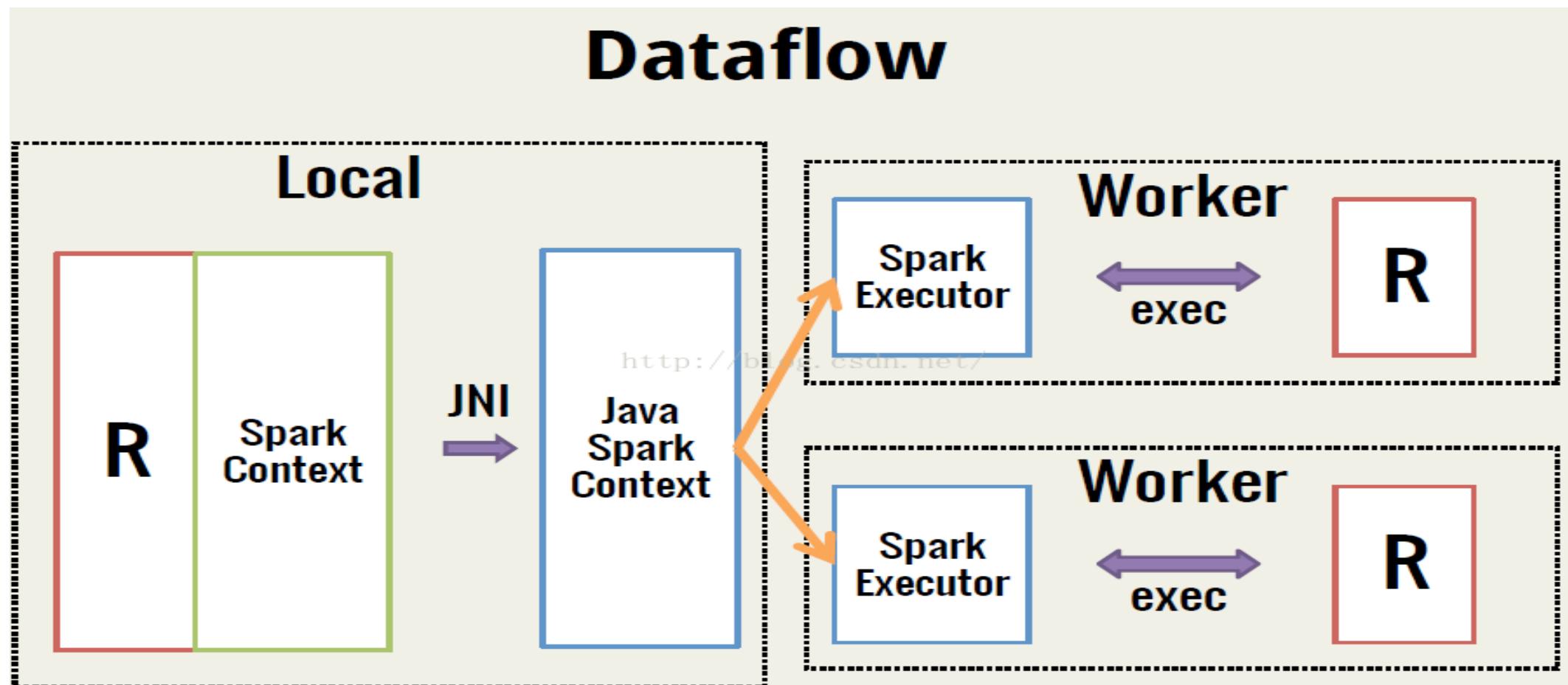


Apache Spark vs Hadoop MR

Parameters	Spark	Hadoop
Data Storage	Spark stores data in-memory.	Hadoop stores data on disk.
Fault tolerance	Spark's data storage model, resilient distributed datasets (RDD) guarantees fault tolerance.	It uses replication to achieve fault tolerance.
Line of code	Apache Spark is project of 20,000 Line of code.	Hadoop 2.0 has 1,20,000 Line of code
Speed	It is Faster due to In-memory computation.	It is relatively slower than Spark.
OS Support	<ul style="list-style-type: none"> • Linux • Windows • Mac OS 	<ul style="list-style-type: none"> • Linux
High level language	<ul style="list-style-type: none"> • Scala • Python • Java • R 	<ul style="list-style-type: none"> • Java
Streaming data	Spark can be used to process as well as modify real-time data with Spark streaming.	With Hadoop Map-Reduce one can process batch of stored data.
Machine Learning	Spark has its own set of Machine learning libraries (MLib).	Hadoop requires interface with other Machine learning library. Eg: Apache Mahout.

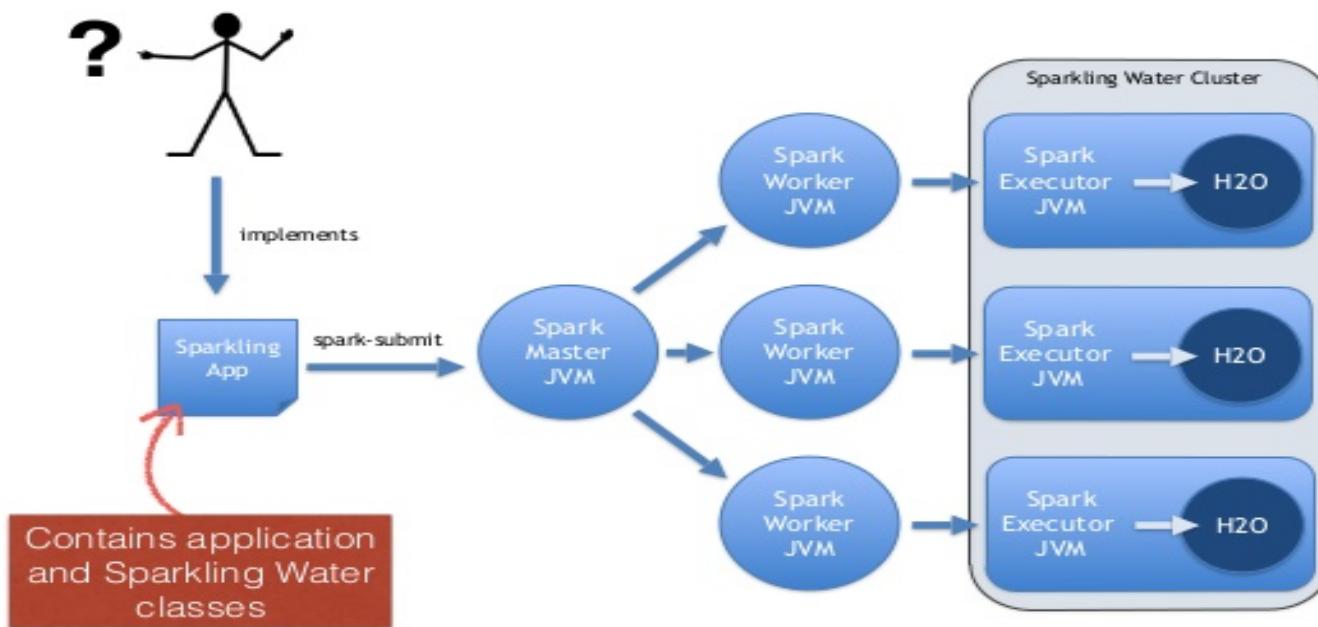
	Hadoop World Record	Spark	Spark
Data Size	102.5 TB	100 TB	1 PB
Elapsed Time	72 mins	23 mins	234 mins
# Nodes	2100	206	190
# Cores	50400	6592	6080
# Reducers	10,000	29,000	250,000
Rate	1.42 TB/min	4.27 TB/min	4.27 TB/min
Rate/node	0.67 GB/min	20.7 GB/min	22.5 GB/min
Sort Benchmark	Yes	Yes	No
Daytona Rules			
Environment	dedicated data center	EC2 (i2.0xlarge)	EC2 (i2.0xlarge)

Spark Dataflow in R



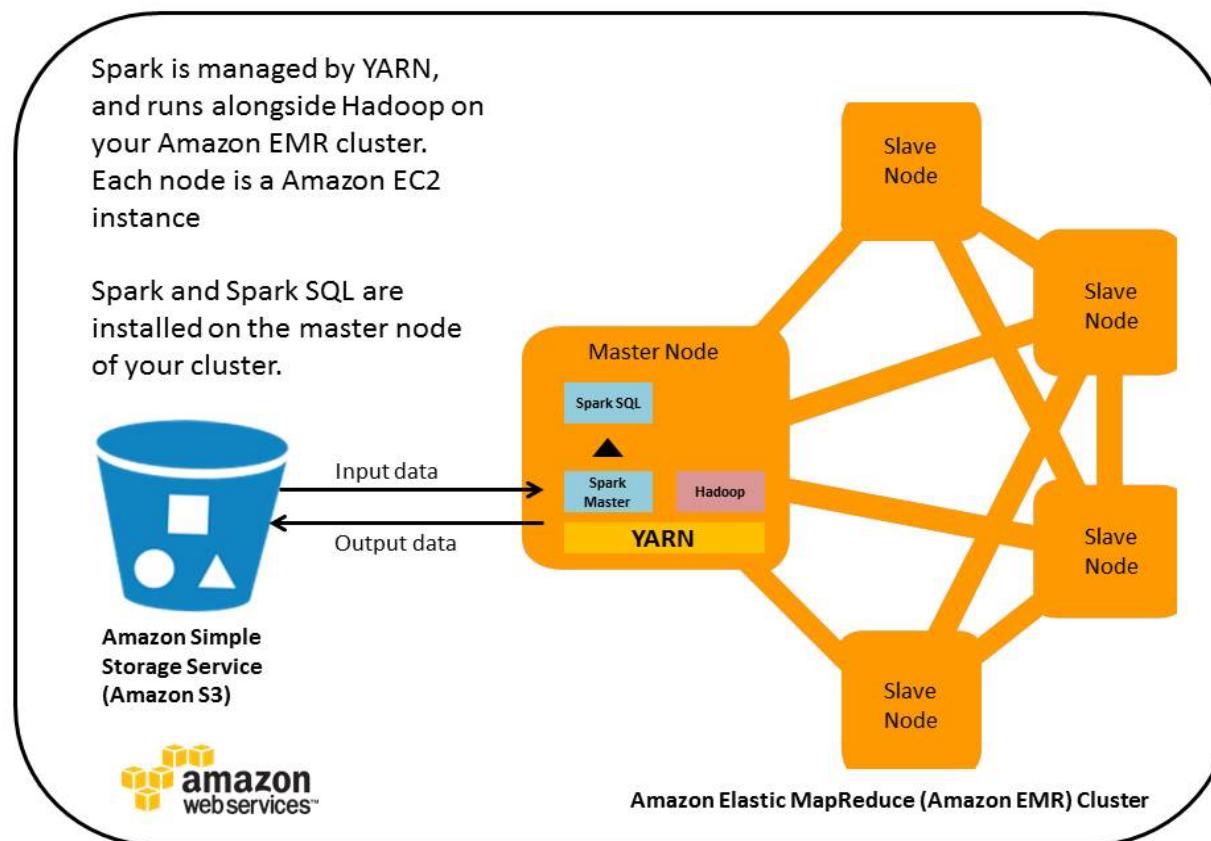
Spark + H2O = Sparkling Water

Sparkling Water Design



*Distributed Machine
Learning
on
Distributed Data Sets*

Amazon Elastic Map Reduce (EMR)



Amazon Elastic Map Reduce (EMR)

Create Cluster - Quick Options [Go to advanced options](#)

General Configuration

Cluster name: Rittman Mead PoC
 Logging
S3 folder: 
Launch mode: Cluster Step execution

Software configuration

Vendor: Amazon MapR
Release: emr-5.0.0
Applications: Spark: Spark 2.0.0 on Hadoop 2.7.2 YARN with Ganglia 3.7.2 and Zeppelin 0.6.1
 Core Hadoop: Hadoop 2.7.2 with Ganglia 3.7.2, Hive 2.1.0, Hue 3.10.0, Mahout 0.12.2, Pig 0.16.0, and Tez 0.8.4
 HBase: HBase 1.2.2 with Ganglia 3.7.2, Hadoop 2.7.2, Hive 2.1.0, Hue 3.10.0, Phoenix 4.7.0, and ZooKeeper 3.4.8
 Presto: Presto 0.150 with Hadoop 2.7.2 HDFS and Hive 2.1.0 Metastore

Hardware configuration

Instance type: m3.xlarge
Number of instances: 3 (1 master and 2 core nodes)

Security and access

EC2 key pair: Choose an option [Learn how to create an EC2 key pair.](#)
Permissions: Default Custom
Use default IAM roles. If roles are not present, they will be automatically created for you with managed policies for automatic policy updates.
EMR role: EMR_DefaultRole
EC2 instance profile: EMR_EC2_DefaultRole

[Cancel](#) [Create cluster](#)

Amazon Elastic Map Reduce (EMR)

Command Line Interface

*Launch a 1 Master 5 Worker node m3.2x large cluster
With Hadoop, Hive, Spark, RStudio, Shiny, sparklyr all installed*

```
aws emr create-cluster --applications Name=Hadoop Name=Spark Name=Hive Name=Pig Name=Tez Name=Ganglia
--release-label emr-5.2.0 --name "EMR 5.2.0 RStudio + sparklyr" --service-role EMR_DefaultRole \
--instance-groups InstanceGroupType=MASTER,InstanceCount=1,InstanceType=m3.2xlarge \
InstanceGroupType=CORE,InstanceCount=5,InstanceType=m3.2xlarge --bootstrap-actions \
Path=s3://aws-bigdata-blog/artifacts/aws-blog-emr-rstudio-sparklyr/rstudio_sparklyr_emr5.sh, \
Args=[ "--rstudio", "--sparkr", "--rexamples", "--plyrmr", "--rhdfs", "--sparklyr"], \
Name="Install RStudio" --ec2-attributes InstanceProfile=EMR_EC2_DefaultRole,KeyName=<Your Key> \
--configurations '[{"Classification": "spark", "Properties": {"maximizeResourceAllocation": "true"}}]' \
--region us-east-1
```

Big Data

Break Time

To give AWS time to spin up the cluster ...

Big Data

Student Presentations

While Flights data is being downloaded ...

Demo

Flights using Spark on Amazon EMR

Data



Have good rest of weekend!