

DS221 | Aug, 2019

# Data Structures, Algorithms & Data Science Platforms

Yogesh Simmhan

[simmhan@cds.iisc.ac.in](mailto:simmhan@cds.iisc.ac.in)

©Department of Computational and Data Science, IISc, 2016

This work is licensed under a [Creative Commons Attribution 4.0 International License](https://creativecommons.org/licenses/by/4.0/)

Copyright for external content used with attribution is retained by their original authors



# L5: Big Data Platforms

Spark, Storm, Giraph

Slide Credits:

- [https://stanford.edu/~rezab/sparkclass/slides/itas\\_workshop.pdf](https://stanford.edu/~rezab/sparkclass/slides/itas_workshop.pdf)
- <https://www.slideshare.net/deanchen11/scala-bay-spark-talk>
- <https://databricks-training.s3.amazonaws.com/slides/advanced-spark-training.pdf>
- Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing, M. Zaharia, et al., NSDI 2012
- <http://spark.apache.org/docs/latest/programming-guide.html>

# What is Big Data?

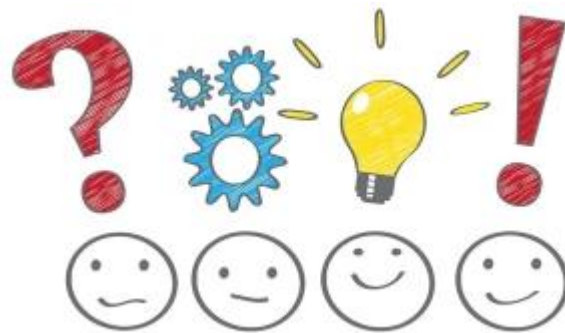




# So...What is Big Data?

Data whose characteristics exceeds the capabilities of conventional *algorithms, **systems** and techniques* to derive useful value.

<https://www.oreilly.com/ideas/what-is-big-data>

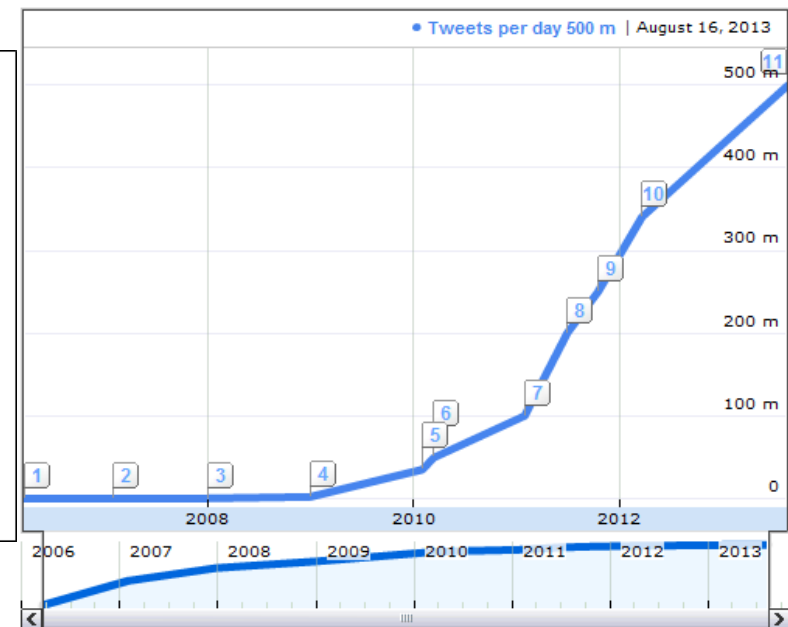
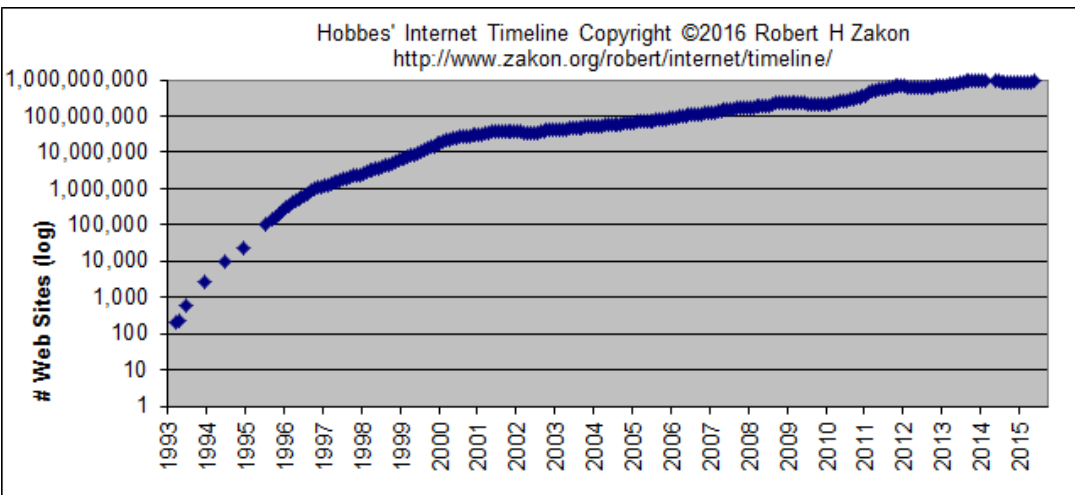




# And, where does Big Data come from?

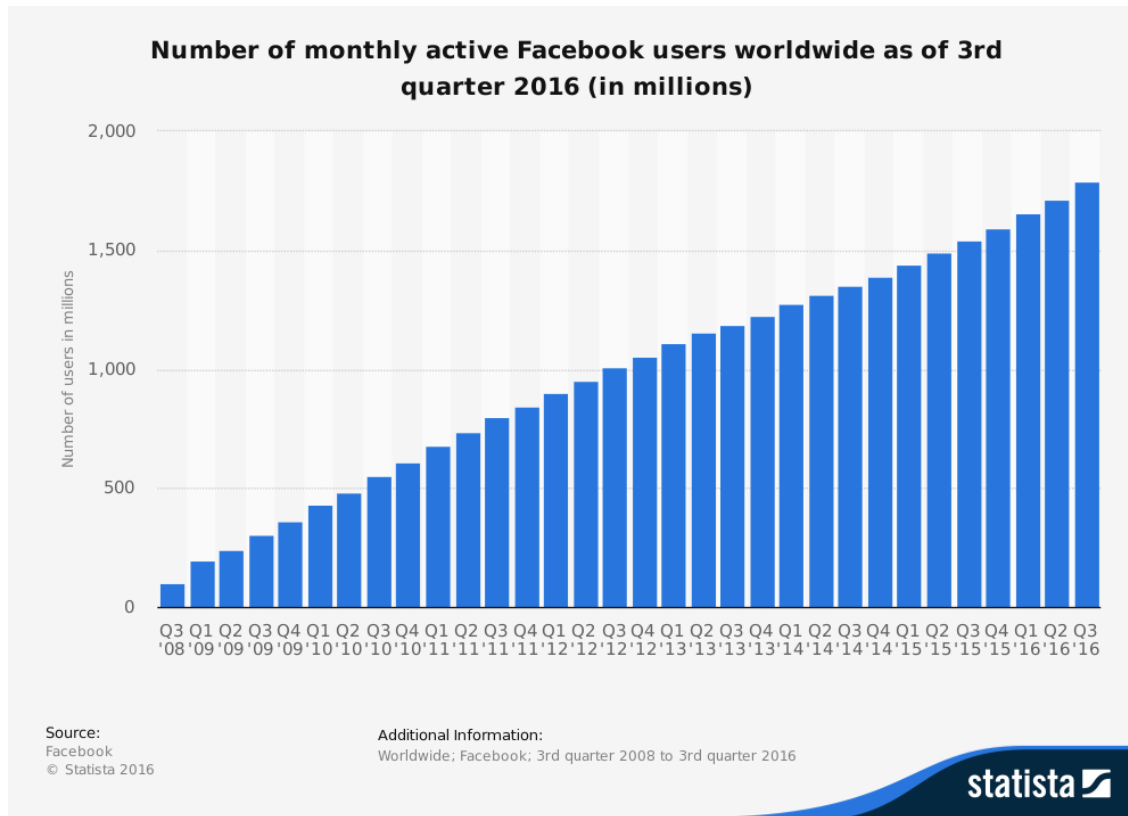
# Web & Social Media

## ■ Web search, Social Networks & Micro-blogs

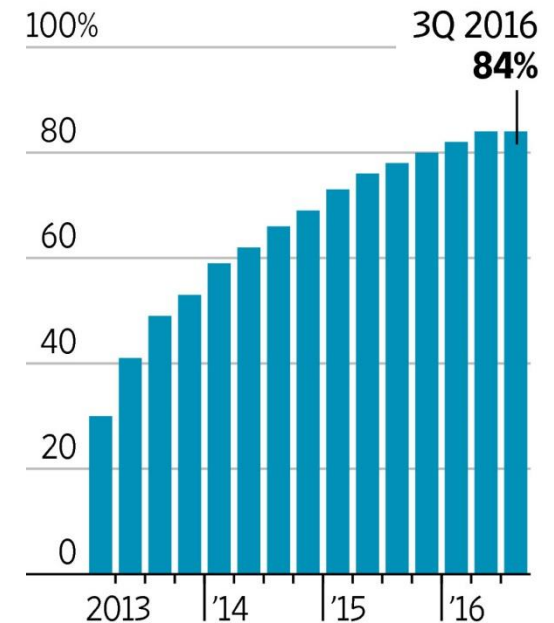


# Web & Social Media

## ■ Social Networks & Micro-blogs



## Facebook's mobile ad revenue as a share of total ad revenue



**1.79 billion monthly active users as of September 30, 2016**

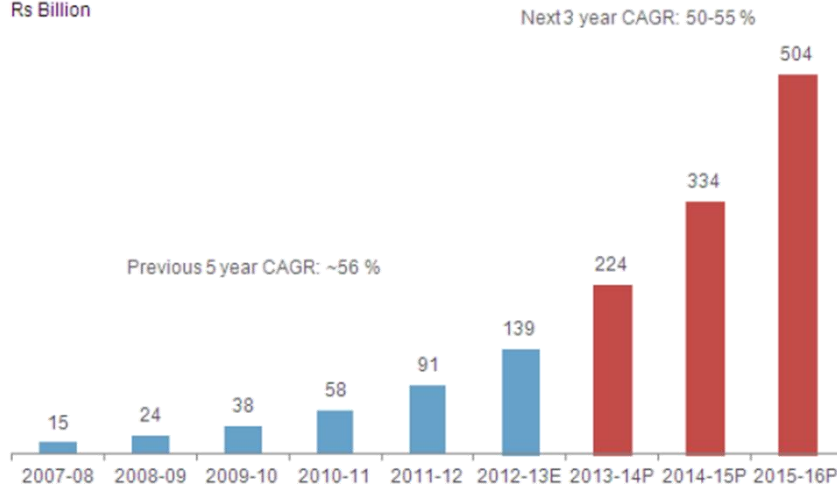


# Enterprises & Government

## ■ Online retail & eCommerce

Online retail market size and growth

Rs Billion

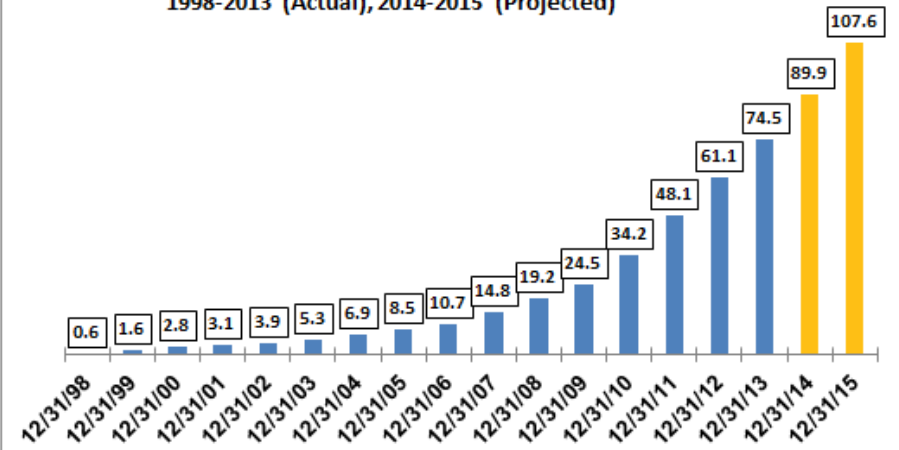


Source: CRISIL Research

<http://blogs.ft.com/beyond-brics/2014/02/28/online-retail-in-india-learning-to-evolve/>

Amazon Annual Revenue (\$ Billions)

1998-2013 (Actual), 2014-2015 (Projected)

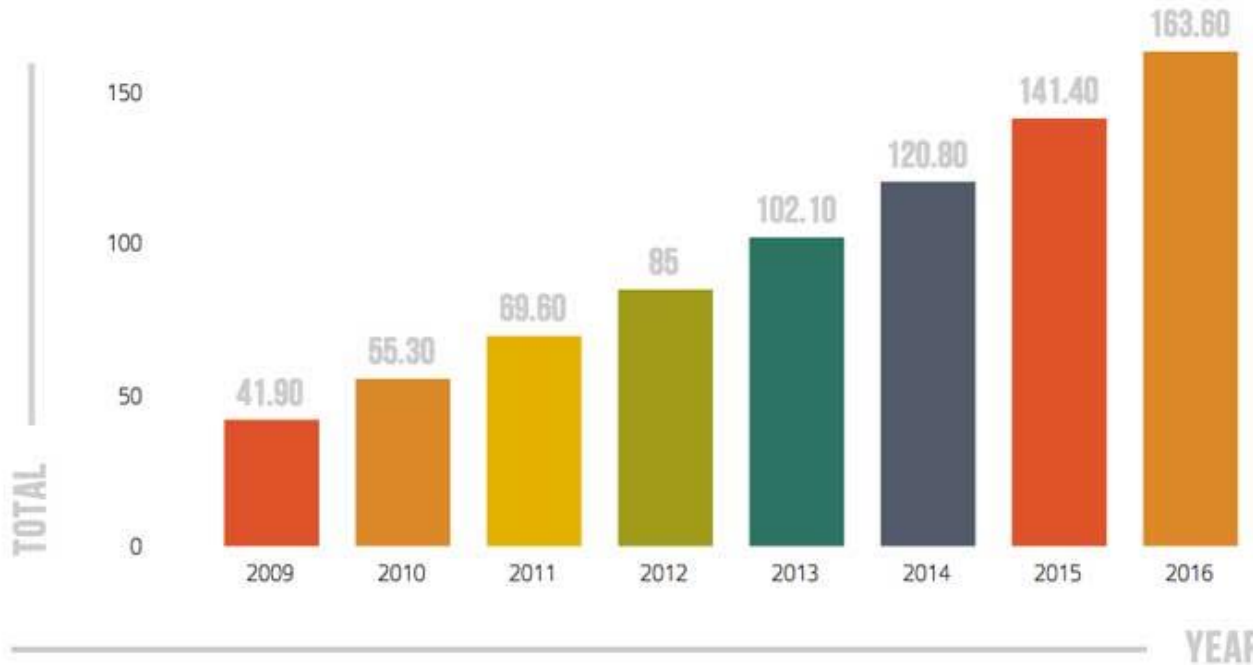


<http://www.peridotcapital.com/2014/04/amazon-sales-growth-projections-for-next-two-years-appear-overly-optimistic.html>

# Enterprises & Government: Finance

## ■ Mobile Transactions & FinTech

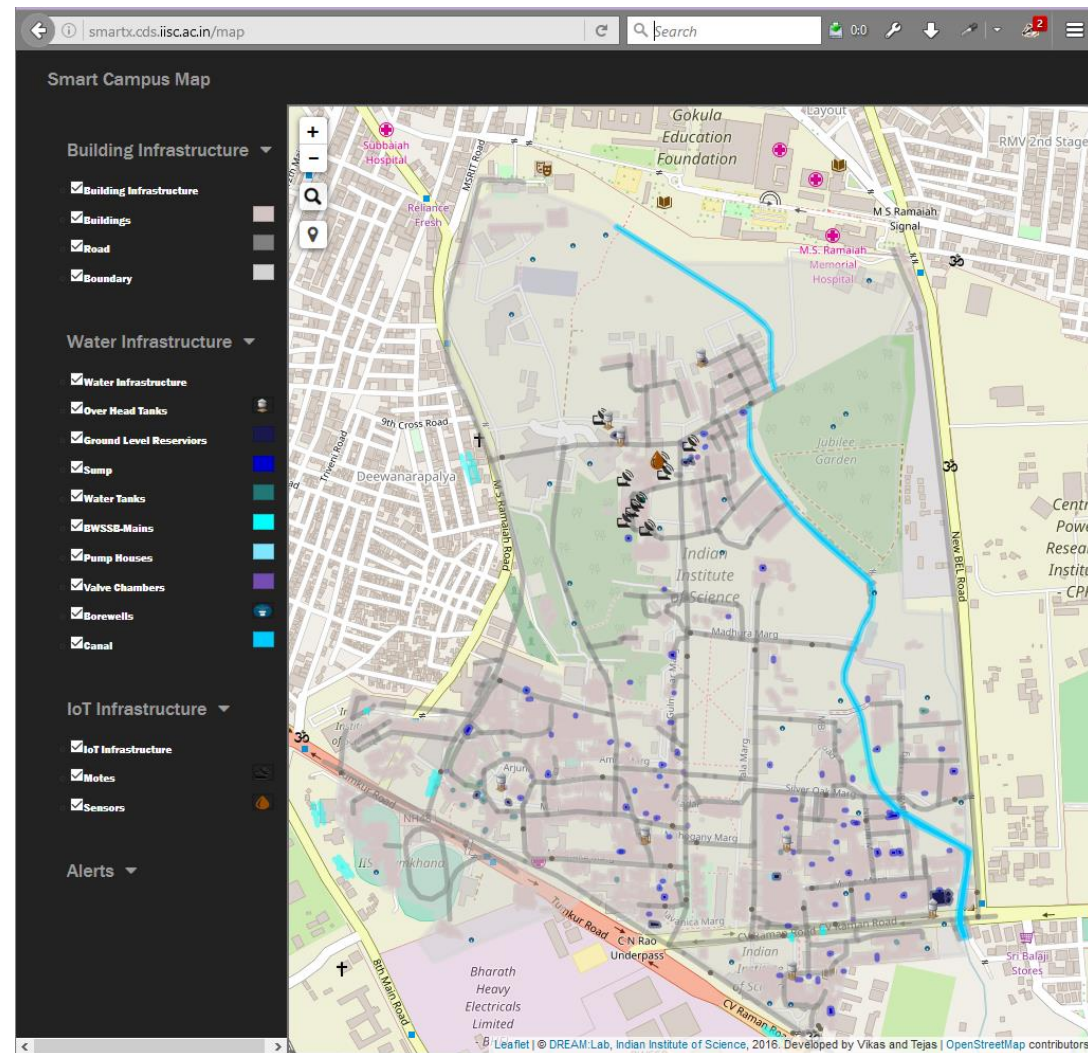
ASIA/PACIFIC (USERS IN MILLIONS)



Since November 8, 2016, *Paytm* has surpassed its metrics -tripling *transactions per day* to 7.5 million

# Internet of Everything

- Personal Devices
  - ▶ Smart Phones, Fitbit
- Smart Appliances
- Smart Cities
  - ▶ Power, Water, Transportation, Environment
- Smart Retail
- Millions of sensor data streams





# Why is Big Data Difficult?



Sources: McKinsey Global Institute, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTec, QAS





**40 ZETTABYTES**

[ 43 TRILLION GIGABYTES ]

of data will be created by 2020, an increase of 300 times from 2005



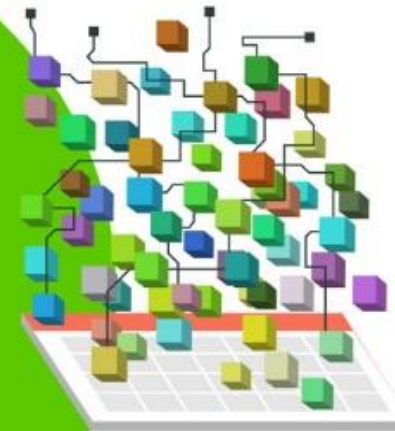
## Volume SCALE OF DATA

It's estimated that

**2.5 QUINTILLION BYTES**

[ 2.3 TRILLION GIGABYTES ]

of data are created each day



# The FOUR of Big Data

From traffic patterns and m  
history and medical reco  
stored, and analyzed to e  
and services that the wor  
But what exactly is big da  
massive amounts of data

As a leader in the sector  
break big data into four  
**Velocity, Variety and Veraci**

Depending on the industry  
data encompasses inform  
internal and external source

**6 BILLION  
PEOPLE**

have cell  
phones



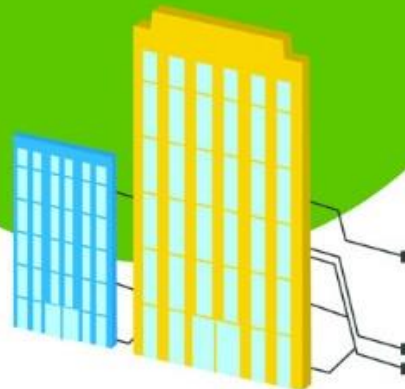
**WORLD POPULATION: 7 BILLION**

Most companies in the  
U.S. have at least

**100 TERABYTES**

[ 100,000 GIGABYTES ]

of data stored



The New York Stock Exchange  
captures

**1 TB OF TRADE**



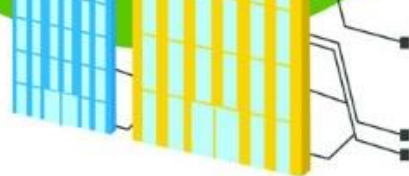
Modern cars have close to

**100 SENSORS**

that monitor items such as



**WORLD POPULATION: 7 BILLION**

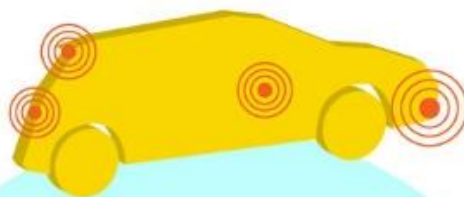


Most companies in the U.S. have at least  
**100 TERABYTES**  
[ 100,000 GIGABYTES ]  
of data stored

The New York Stock Exchange captures

**1 TB OF TRADE INFORMATION**

during each trading session



Modern cars have close to

**100 SENSORS**

that monitor items such as fuel level and tire pressure

# Velocity

## ANALYSIS OF STREAMING DATA

By 2016, it is projected there will be

**18.9 BILLION NETWORK CONNECTIONS**

– almost 2.5 connections per person on earth



and services that the v  
But what exactly is big  
massive amounts of dat

As a leader in the sec  
break big data into fo  
**Velocity, Variety and Ver**

Depending on the indu  
data encompasses inf  
internal and external sou  
social media, enterpris  
mobile devices. Comp  
adapt their products an  
customer needs, opt  
infrastructure, and find

By 2015

**4.4 MILLION IT JO**

will be created globally  
with 1.9 million in the





# R V's ig

and music downloads to web records, data is recorded, to enable the technology world relies on every day. data, and how can these data be used?

ctor, IBM data scientists our dimensions: **Volume**, **Velocity**, **Veracity**

stry and organization, big information from multiple sources such as transactions, use content, sensors and

As of 2011, the global size of data in healthcare was estimated to be

**150 EXABYTES**

[ 161 BILLION GIGABYTES ]

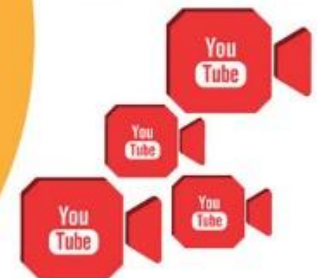


By 2014, it's anticipated there will be

**420 MILLION WEARABLE, WIRELESS HEALTH MONITORS**

**4 BILLION+ HOURS OF VIDEO**

are watched on YouTube each month



## Variety DIFFERENT FORMS OF DATA

**30 BILLION PIECES OF CONTENT**

are shared on Facebook every 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



Poor data quality costs the US economy around

**\$3.1 TRILLION A YEAR**



world relies on every day.  
data, and how can these  
be used?

or, IBM data scientists  
r dimensions: **Volume,**  
**Velocity**

ty and organization, big  
formation from multiple  
ces such as transactions,  
content, sensors and  
es can leverage data to  
services to better meet  
imize operations and  
ew sources of revenue.

S  
to support big data,  
United States



**400 MILLION TWEETS**

are sent per day by about 200  
million monthly active users

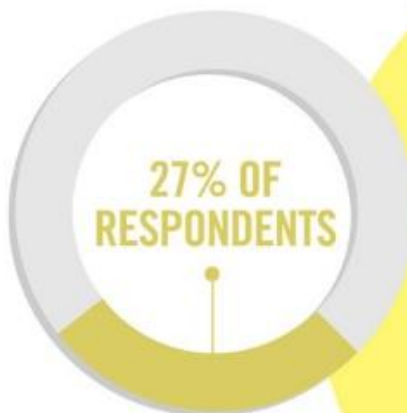
### 1 IN 3 BUSINESS LEADERS

don't trust the information  
they use to make decisions



Poor data quality costs the US  
economy around

**\$3.1 TRILLION A YEAR**

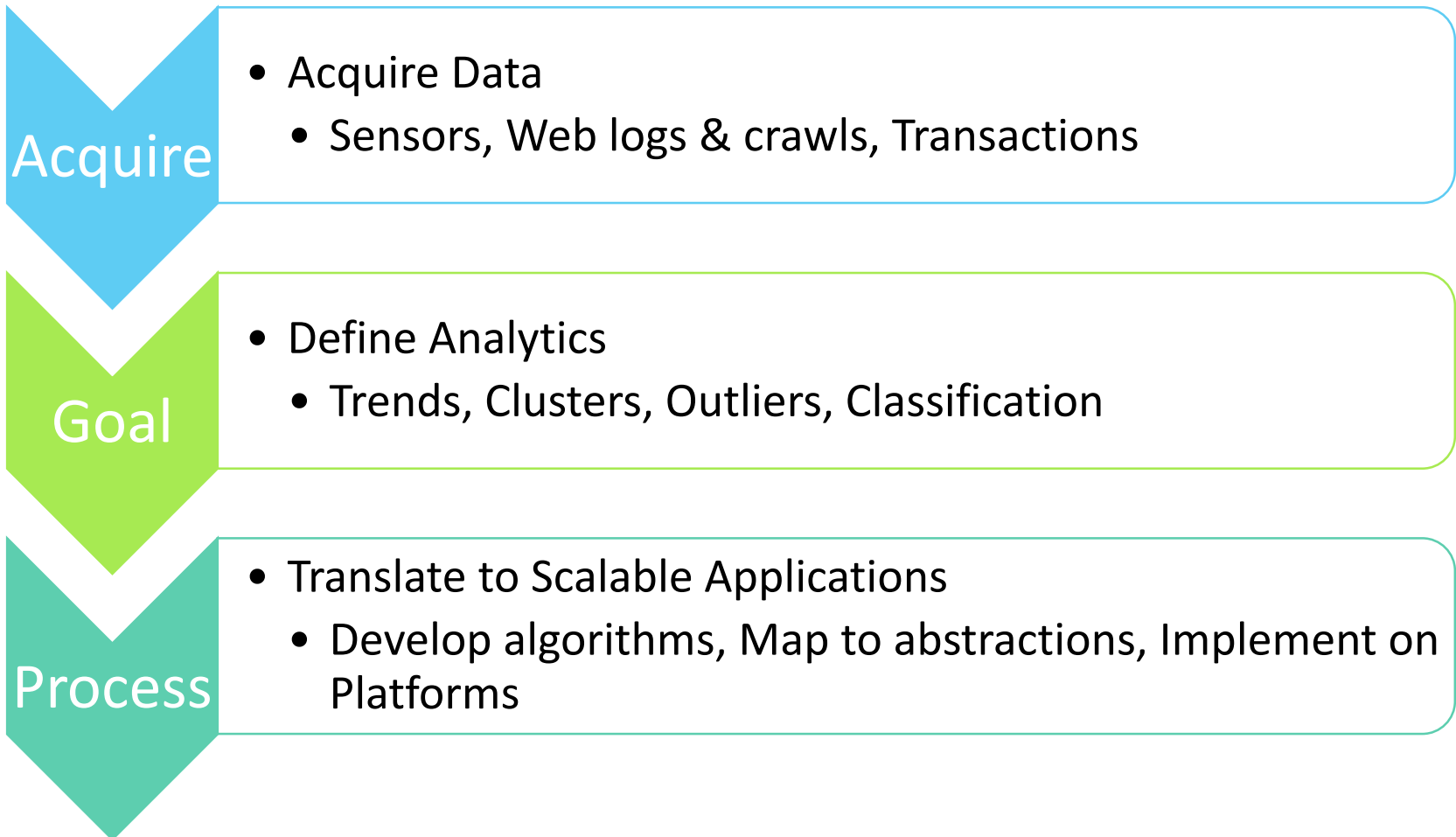


in one survey were unsure of  
how much of their data was  
inaccurate

## Veracity

### UNCERTAINTY OF DATA

# Data Analysis Lifecycle





# Data Platforms

- Acquire, manage, process Big Data
- At large scales
- To meet application needs

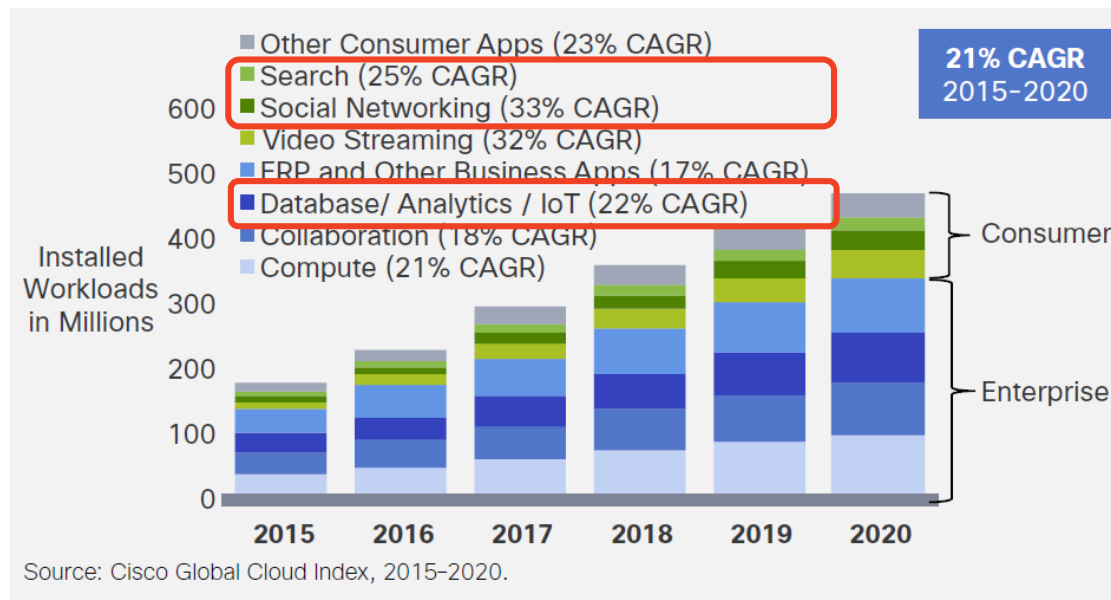


# Distributed Systems

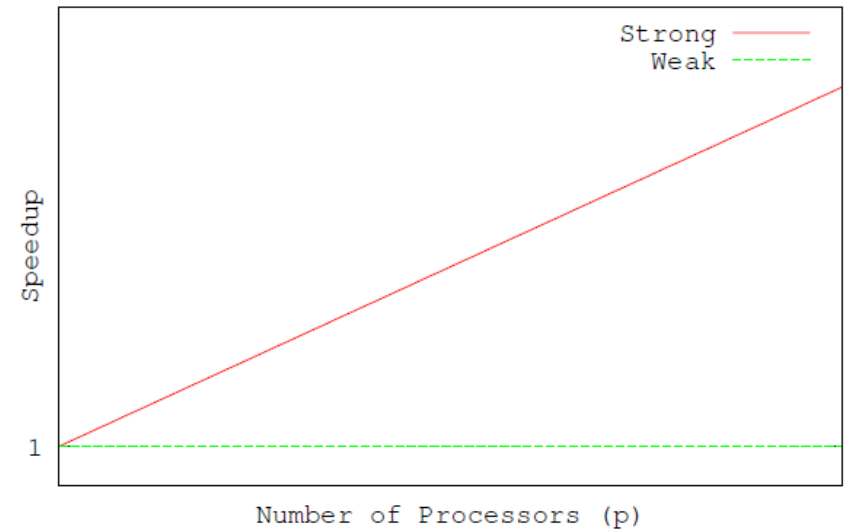
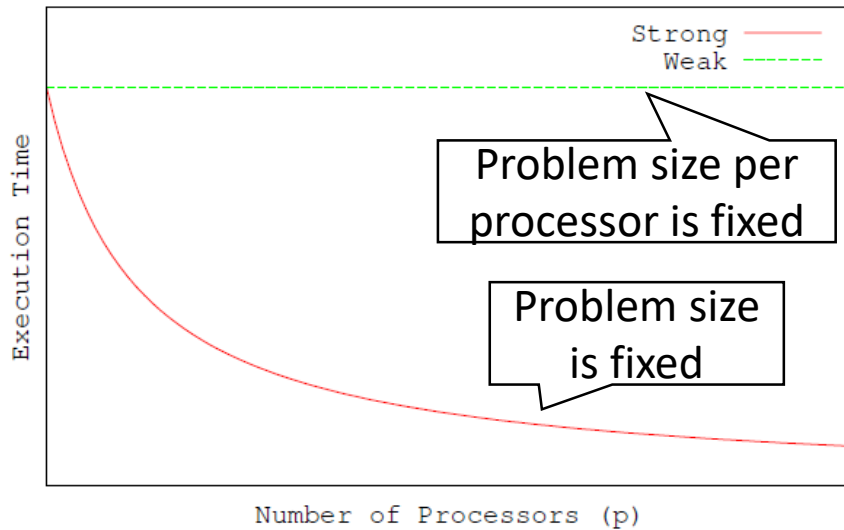
- Distributed Computing
  - ▶ Clusters of machines
  - ▶ Connected over network
- Distributed Storage
  - ▶ Disks attached to clusters of machines
  - ▶ Network Attached Storage
- *How can we make effective use of multiple machines?*
- **Commodity** clusters vs. **HPC** clusters
  - ▶ Commodity: Available off the shelf at large volumes
  - ▶ Lower Cost of Acquisition
  - ▶ Cost vs. Performance
    - Low disk bandwidth, and high network latency
    - CPU typically comparable (Xeon vs. i3/5/7)
    - Virtualization overhead on Cloud
- *How can we use many machines of modest capability?*

# Growth of Cloud Data Centers

**Figure 17.** Global Data Center Workloads by Applications



# Ideal Strong/Weak Scaling



# Scalability

- Strong vs. Weak Scaling
- **Strong Scaling:** How the performance varies with the # of processors for a *fixed total problem size*
- **Weak Scaling:** How the performance varies with the # of processors for a *fixed problem size per processor*
  - ▶ Big Data platforms are intended for “Weak Scaling”

# Ease of Programming

- Programming distributed systems is difficult
  - ▶ Divide a job into multiple tasks
  - ▶ Understand dependencies between tasks: Control, Data
  - ▶ Coordinate and synchronize execution of tasks
  - ▶ Pass information between tasks
  - ▶ Avoid race conditions, deadlocks
- Parallel and distributed programming models/languages/abstractions/platforms try to make these easy
  - ▶ E.g. Assembly programming vs. C++ programming
  - ▶ E.g. C++ programming vs. Matlab programming



# Availability, Failure

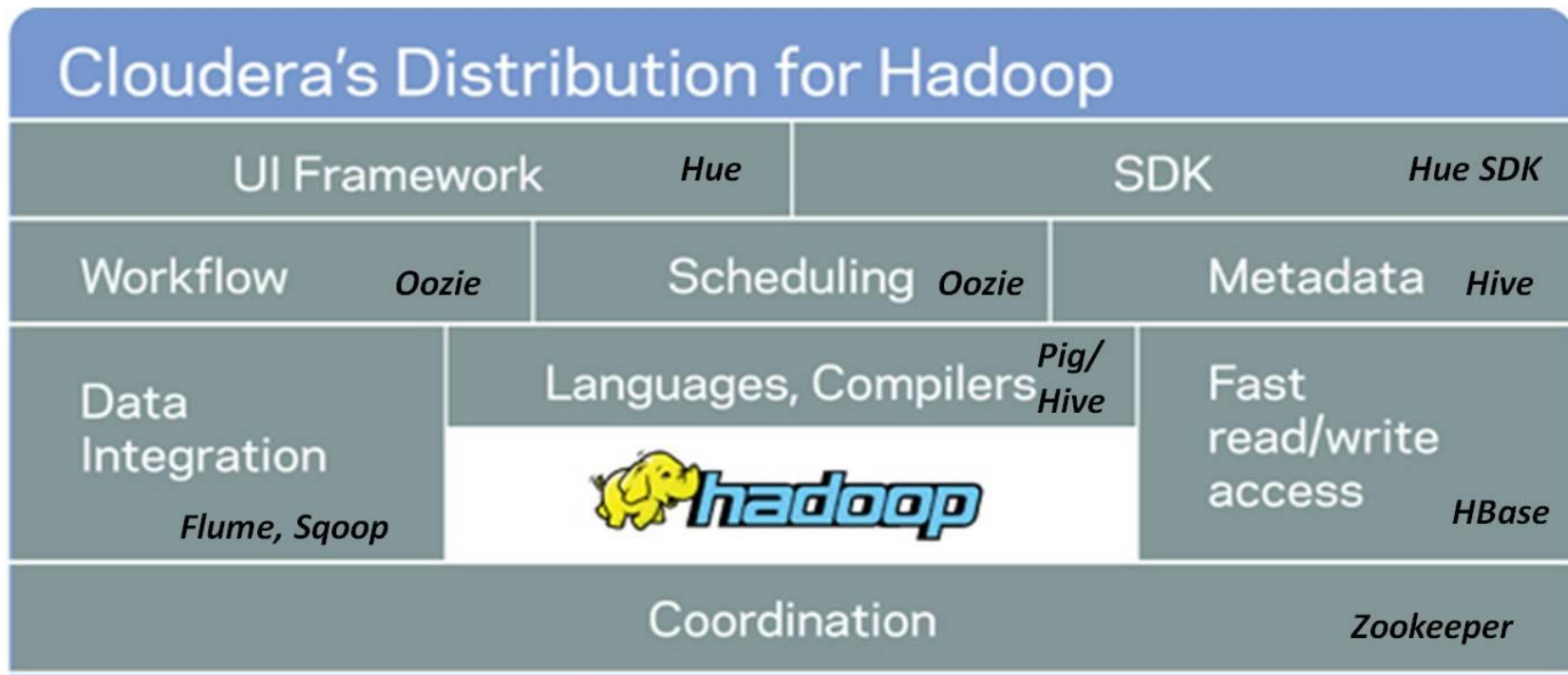
- Commodity clusters have lower reliability
  - ▶ Mass-produced
  - ▶ Cheaper materials
  - ▶ Smaller lifetime (~3 years)
- *How can applications easily deal with failures?*
- *How can we ensure availability in the presence of faults?*

# Early Technologies

- **MapReduce** is a distributed data-parallel programming model from Google
- MapReduce works best with a distributed file system, called **Google File System (GFS)**
- **Hadoop** is the open source framework implementation from Apache that can execute the MapReduce programming model
- **Hadoop Distributed File System (HDFS)** is the open source implementation of the GFS design
- **Elastic MapReduce (EMR)** is Amazon's PaaS

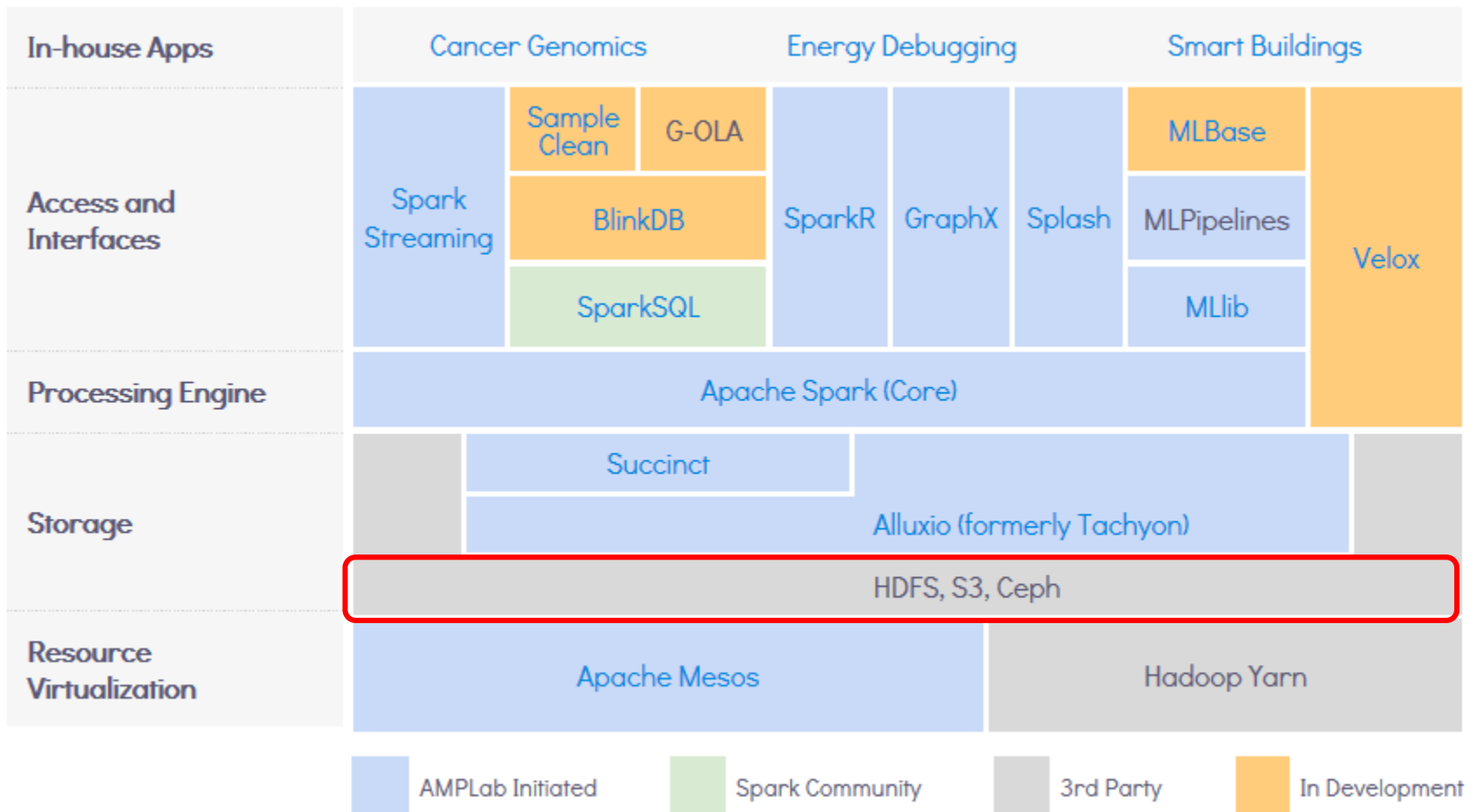
# Platforms...Think in terms of Stacks

## Cloudera



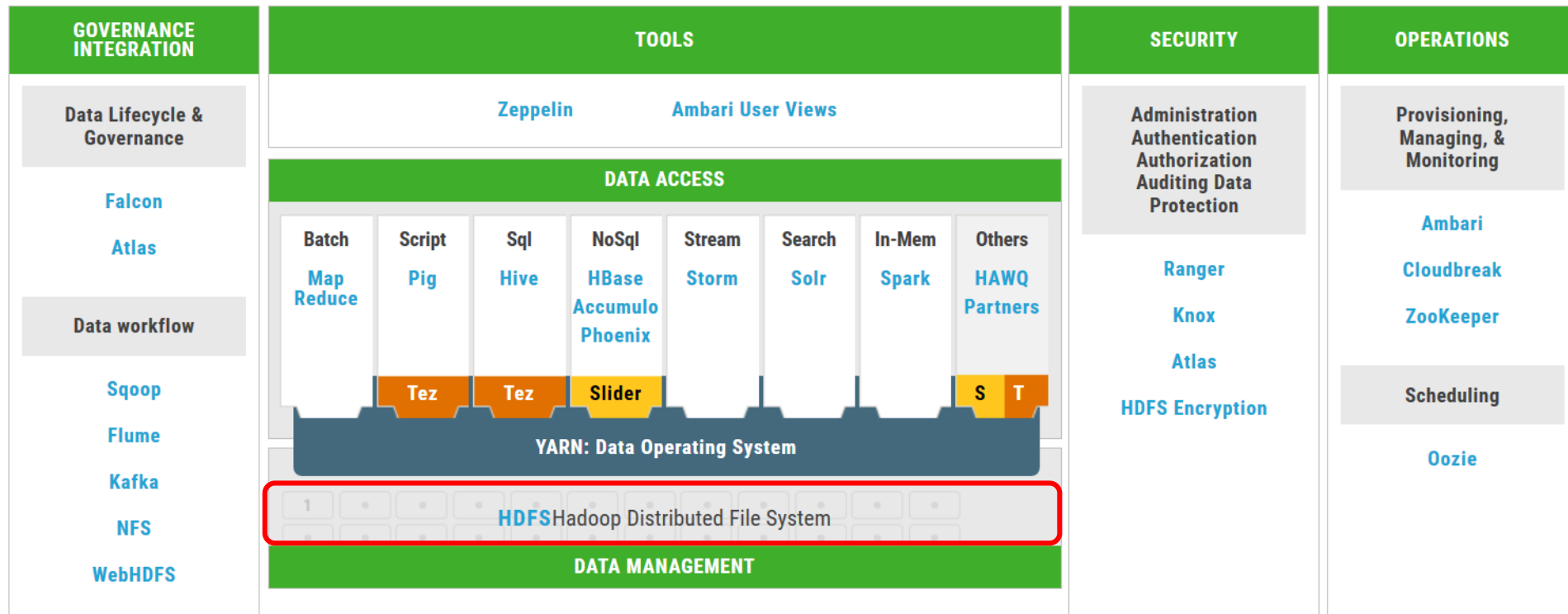
# Platforms...Think in terms of Stacks

## BDAS



# Platforms...Think in terms of Stacks

## HortonWorks





# Apache Spark

Slides & Additional Reading Courtesy

[https://stanford.edu/~rezab/sparkclass/slides/itas\\_workshop.pdf](https://stanford.edu/~rezab/sparkclass/slides/itas_workshop.pdf)

Resilient Distributed Datasets, Matei Zaharia

<http://spark.apache.org/docs/2.1.1/programming-guide.html>

<http://spark.apache.org/docs/latest/api/java/index.html>

<https://www.gitbook.com/book/jaceklaskowski/mastering-apache-spark/details>

**Apache Spark Internals**, Pietro Michiardi, Eurecom

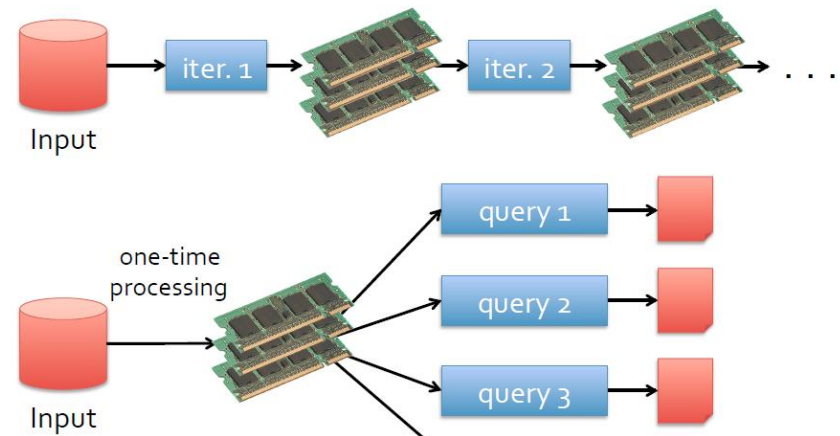
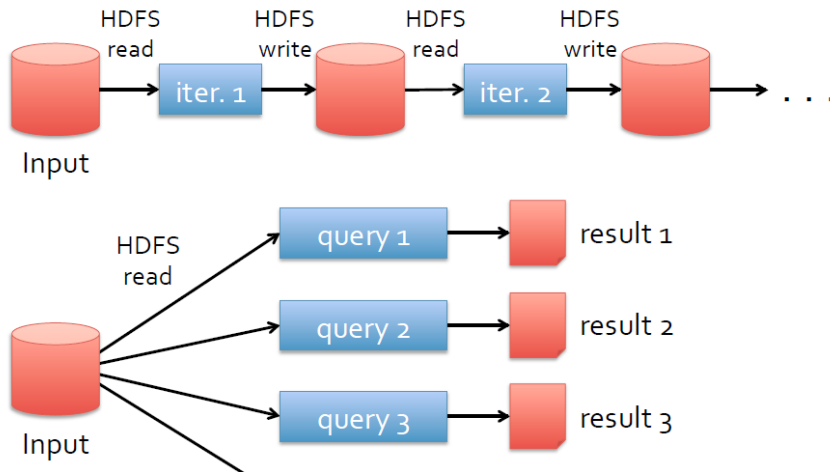


# Why Spark?

- Ease of language definition
  - ▶ Typing, dataflows,
  - ▶ But Pig, Hive, HBase, etc. give you that
- Better performance using “In memory” compute
  - ▶ Multiple stages part of same job
  - ▶ Lazy evaluation, caching/persistence

# In-memory computation

- Operate on data in (distributed) memory
  - ▶ Allows many operations to be performed locally
  - ▶ Write to disk only when data sharing required across workers
- This is unlike others like Hadoop Map/Reduce







# RDD: The Secret Sauce

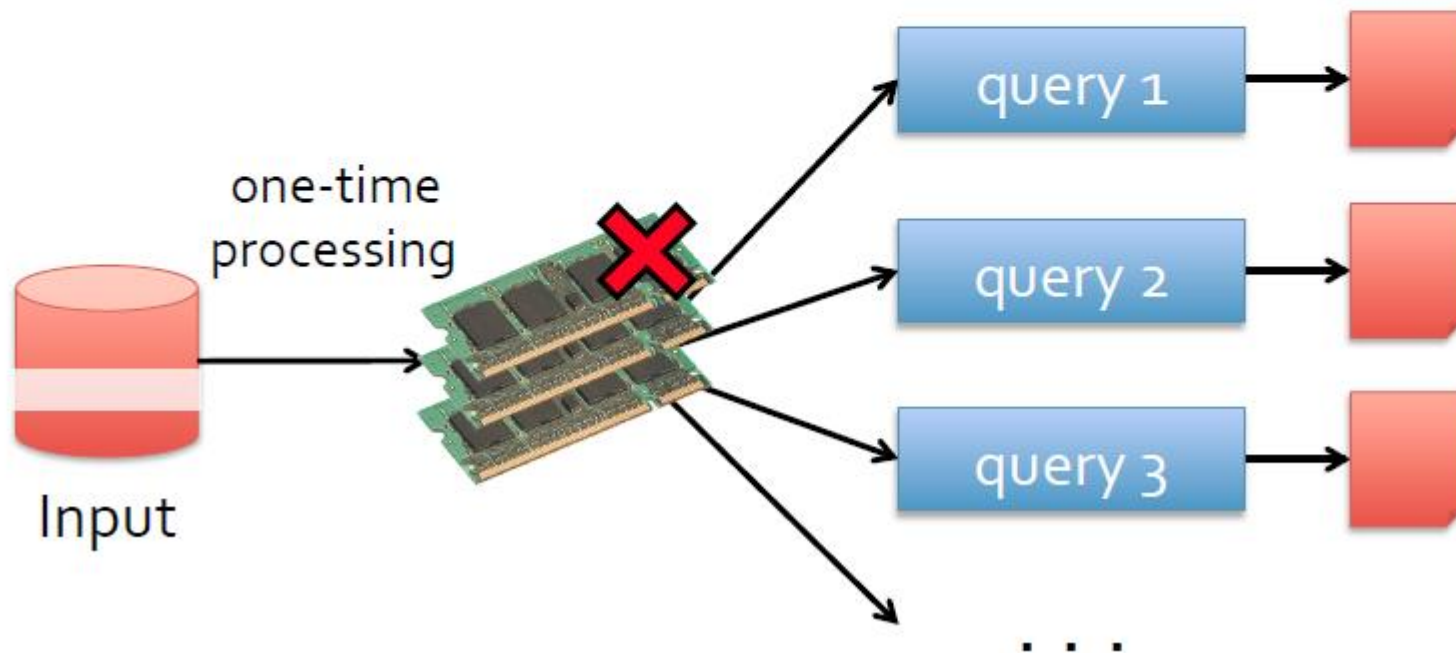
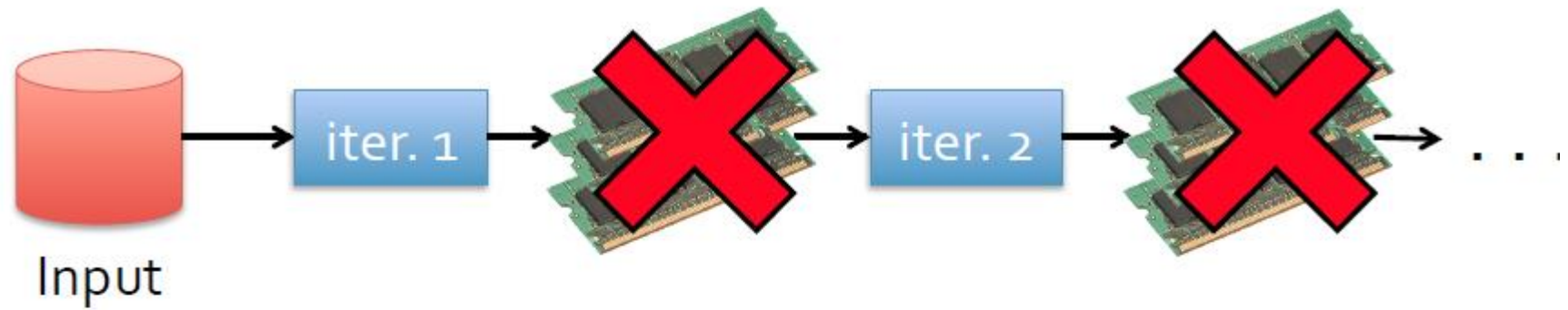
- RDD: Resilient Distributed Dataset

- ▶ Immutable, partitioned collection of tuples
- ▶ Operated on by *deterministic transformations*
  - Object-oriented flavor
  - `RDD.operation()` → RDD

- Recovery by re-computation

- ▶ Maintains lineage of transformations
- ▶ Recompute missing partitions *if failure happens*
- ▶ Not possible/not automatic in Pig

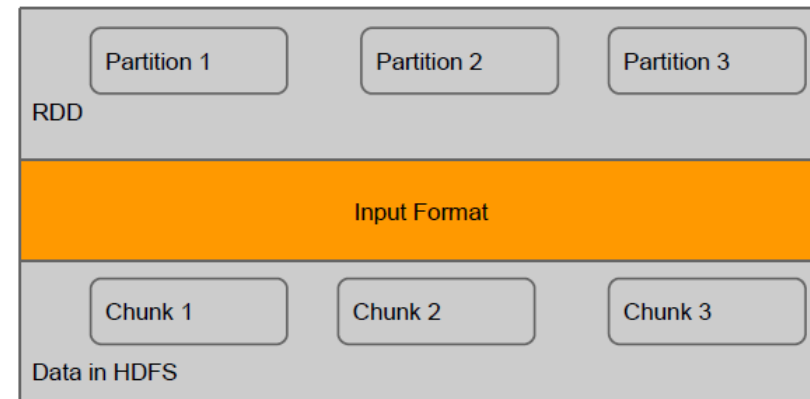
- Allows caching & persistence for reuse





# RDD Partitions

- RDD is internally a collection of partitions
  - Each partition holds a list of items
- Partitions may be present on a different machine
  - Partition is the *unit of execution*
  - Partition is the *unit of parallelism*
- They are immutable
  - Each transformation on an RDD generates a new RDD with different partitions
  - *Allows recovery of individual partitions*





# RDD Operations

Allows  
composability  
into Dataflows

<b>Transformations</b> (define a new RDD)	map filter sample groupByKey reduceByKey sortByKey	flatMap union join cogroup cross mapValues
<b>Actions</b> (return a result to driver program)	collect reduce count save lookupKey	



# A Sample Spark Program

- Movielens dataset, **movies.csv**

- ▶ movieId,title,genres

```
m = sc.textFile("hdfs:///ml/movies.csv").cache()
['movieId,title,genres']...
mcols = m.map(lambda l: l.split(",")).
mg = mcols.filter(lambda l: l[2] != 'genres')
['92363','Toy Story','cartoon|action|children']...
mgc = mg.map(lambda l: (len(l[2].split("|")), 1))
[3,['92363','Toy Story','cartoon|action|children']]...
maxgc = mgc.max()[0]
3
maxgcm = mgc.lookup(maxgc)
[3,['92363','Toy Story','cartoon|action|children']]...
```



# Creating RDD

- Load external data from distributed storage
- Create logical RDD on which you can operate
- Support for different input formats
  - ▶ HDFS files, Cassandra, Java serialized, directory, gzipped
- Can control the number of partitions in loaded RDD
  - ▶ Default depends on external DFS, e.g. 128MB on HDFS

```
m = sc.textFile("hdfs:///ml/movies.csv").cache()
```



# RDD Operations

## ■ Transformations

- ▶ From one RDD to one or more RDDs
- ▶ Lazy evaluation upon “action”...*use with care*
- ▶ Executed in a distributed manner

## ■ Actions

- ▶ Perform aggregations on RDD items
- ▶ Return single (or distributed) results to “driver” code
- ▶ **RDD.collect()** brings RDD partitions to single driver machine



# RDD and PairRDD

- RDD is logically a collection of items with a generic type
- PairRDD is a 2-tuple, like a “Map”, where each item in the collection is a <key,value> pair
  - ▶ But can have *duplicate keys*
- Transformation functions use RDD or PairRDD as input/output



# Transformations

Transformation	Meaning
<b>map</b> ( <i>func</i> )	Return a new distributed dataset formed by passing each element of the source through a function <i>func</i> .
<b>filter</b> ( <i>func</i> )	Return a new dataset formed by selecting those elements of the source on which <i>func</i> returns true.
<b>flatMap</b> ( <i>func</i> )	Similar to map, but each input item can be mapped to 0 or more output items (so <i>func</i> should return a Seq rather than a single item).
Implicit in PySpark	
<b>sample</b> ( <i>withReplacement</i> , <i>fraction</i> , <i>seed</i> )	Sample a fraction <i>fraction</i> of the data, with or without replacement, using a given random number generator <i>seed</i> .
<b>union</b> ( <i>otherDataset</i> )	Return a new dataset that contains the union of the elements in the source dataset and the argument.
<b>intersection</b> ( <i>otherDataset</i> )	Return a new RDD that contains the intersection of elements in the source dataset and the argument.
Also removes duplicates	
<b>distinct</b> ([ <i>numTasks</i> ])	Return a new dataset that contains the distinct elements of the source dataset.



# Transformations on PairRDD

**aggregateByKey**(zeroValue)(seqOp, combOp, [numTasks])

When called on a dataset of (K, V) pairs, returns a dataset of (K, U) pairs where the values for each key are aggregated using the given combine functions and a neutral "zero" value. Allows an aggregated value type that is different than the input value type, while avoiding unnecessary allocations. Like in `groupByKey`, the number of reduce tasks is configurable through an optional second argument.

**join**(otherDataset, [numTasks])

When called on datasets of type (K, V) and (K, W), returns a dataset of (K, (V, W)) pairs with all pairs of elements for each key. Outer joins are supported through `leftouterJoin`, `rightouterJoin`, and `fullouterJoin`.

# Aggregation: Average number of ratings given by users

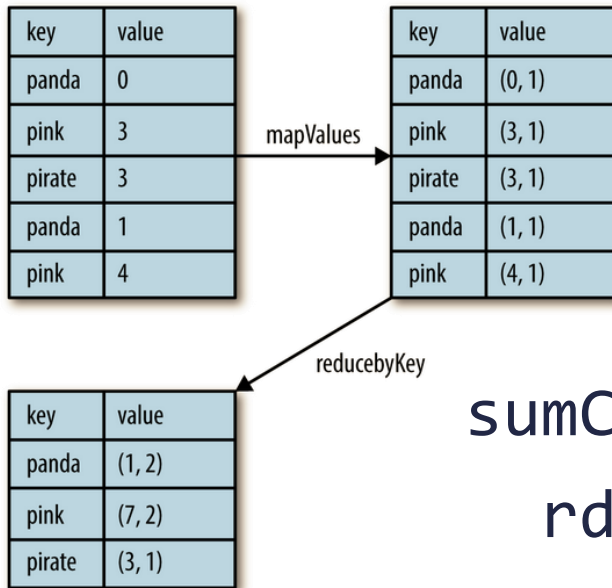
```
[userId,movieId,rating,timestamp]
rv = r.map(Lambda l: l.split(",")[2])
rfv = rv.filter(Lambda l:
                l != 'rating')
[rating]...
rvs = rfv.reduce(Lambda a, b:
                 float(a) + float(b))
rvc = rfv.count()
print rvs/rvc
```

**Action****Action**

# Actions

<b>reduce(func)</b>	Aggregate the elements of the dataset using a function <i>func</i> (which takes two arguments and returns one). The function should be commutative and associative so that it can be computed correctly in parallel.
<b>collect()</b>	Return all the elements of the dataset as an array at the driver program. This is usually useful after a filter or other operation that returns a sufficiently small subset of the data.
<b>count()</b>	Return the number of elements in the dataset.
<b>countByKey()</b>	Only available on RDDs of type (K, V). Returns a hashmap of (K, Int) pairs with the count of each key.
<b>first()</b>	Return the first element of the dataset (similar to take(1)).
<b>take(n)</b>	Return an array with the first <i>n</i> elements of the dataset.
<b>takeSample(withReplacement, num, [seed])</b>	Return an array with a random sample of <i>num</i> elements of the dataset, with or without replacement, optionally pre-specifying a random number generator seed.

# Samples: Per-key average



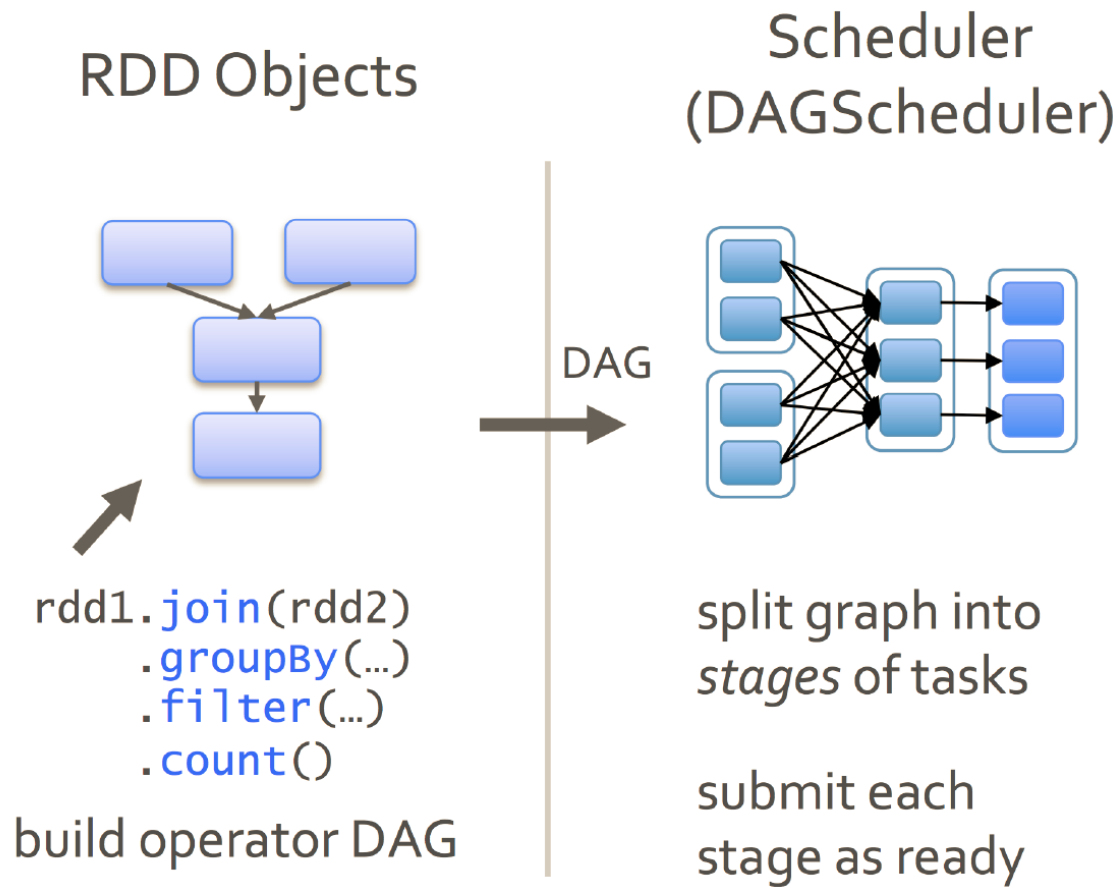
```
sumCount =  
  rdd.mapValues(x -> (x,1)).  
    reduceByKey((x, y) ->  
      (x[0]+y[0], x[1]+y[1]))
```



# RDD Persistence & Caching

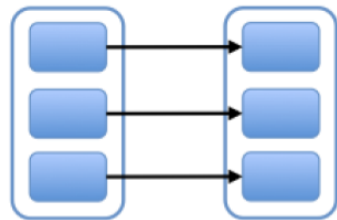
- RDDs can be reused in a dataflow
  - Branch, iteration
- But it will be re-evaluated each time it is reused!
- Explicitly persist RDD to reuse output of a dataflow path multiple times
- Multiple storage levels for persistence
  - Disk or memory
  - Serialized or object form in memory
  - Partial spill-to-disk possible
  - *Cache* indicates “persist” to memory

# Distributed Execution

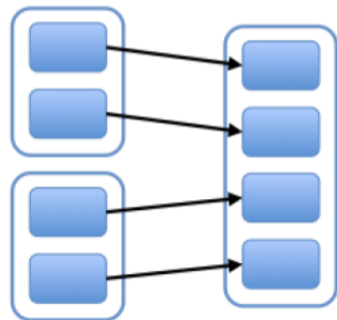


# Execution Dependency

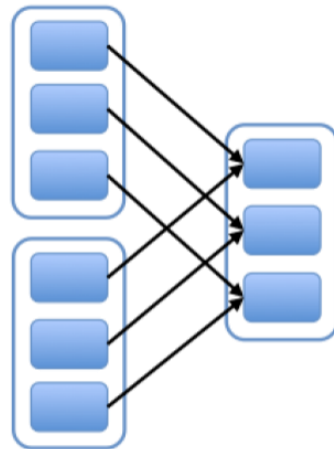
**NARROW DEPENDENCY:** Each partition of the parent RDD is used by at most one partition of the child RDD. Task can be executed locally and we don't have to shuffle.



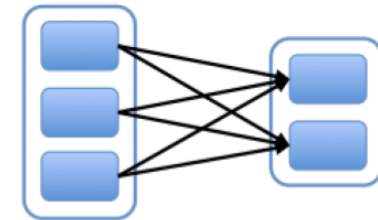
map, filter



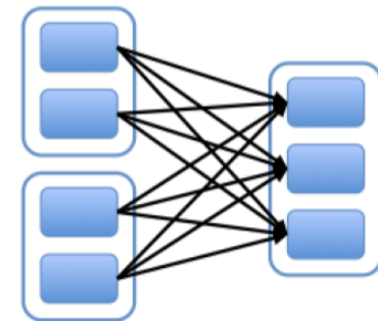
union



join with  
co-partitioned  
inputs



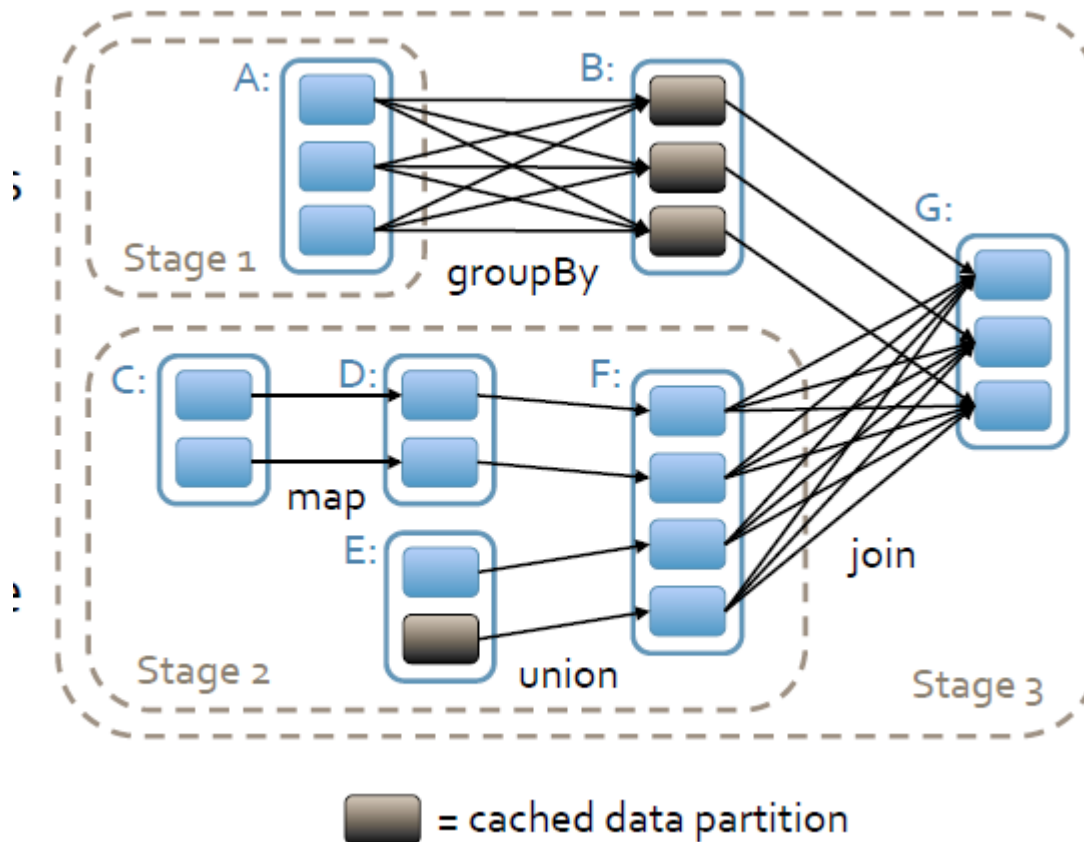
groupByKey



join with inputs not  
co-partitioned



# Lazy Execution



# From DAG to RDD lineage

