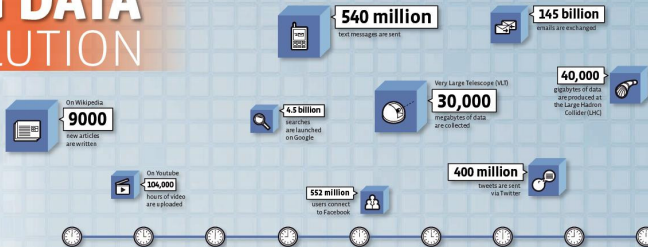




Développement logiciel pour le Cloud (TLC)

Introduction

The **BIG DATA** REVOLUTION



Source: CNRS magazine 2013

Big data

“Big data refers to data sets whose size is beyond the ability of typical database software tools to capture, store, manage and analyze.” — *The McKinsey Global Institute, 2011.*

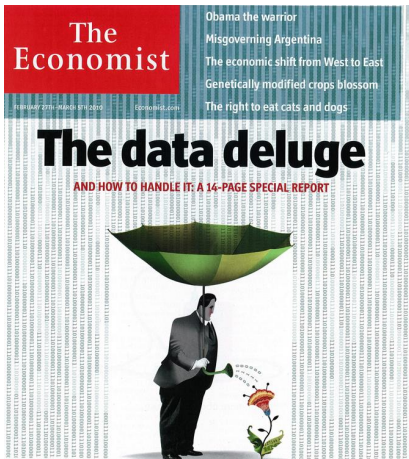
“Big data is the term for a collection of data sets so large and complex that it becomes difficult to process using on-hand database management tools or traditional data processing applications.” — *Wikipedia.*

How big is big data?

Earlier Berkeley studies estimated that by the end of 1999, the sum of human-produced information (including all audio, video recordings and text/books) was about 12 Exabytes of data (*1 exabyte = 1 million TB*).

Eric Schmidt: **Every 2 Days We Create As Much Information As We Did Up To 2003.**

<http://techcrunch.com/2010/08/04/schmidt-data/>

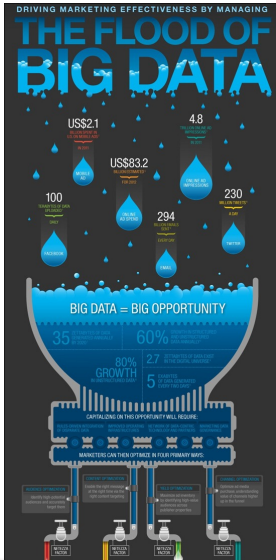


In 2010 the Digital Universe contained **1.2 zettabytes** (1 zettabyte = 1 *billion* TB)

In 2020 the Digital Universe will contain **35 zettabytes**.



Why do we want to analyze this data?



Big Data @ Work

Organizations in all industries are under increasing pressure to capitalize on data.

Healthcare

The average amount of data per hospital will increase from 147TB in 2013 to 467TB in 2015, driven by the exponential growth of medical images and electronic medical records.

With Big Data

Medical professionals can improve patient care and reduce costs by extracting relevant clinical information from vast amounts of data to better understand the past and predict future outcomes.

Customer Service

Today, 66% of consumers quit doing business with a company because of a bad customer experience, up from 57% four years ago.¹

With Big Data

Service representatives can use data to gain a more holistic view of their customers, understanding their likes and dislikes in real-time in order to resolve a problem or capture an happy client faster.

Insurance

Insurance companies and government agencies each gather flood data related to their own individual missions. But the kind, quality and volume of data compiled varies widely.²

With Big Data

An insurance or claims services provider can apply advanced analytics to data and detect fraud quickly, before funds are paid out.

Financial Services

Wall Street alone delivers 3 new research documents every minute. Dow Jones publishes upwards of 17,000 news items per day.³

With Big Data

Financial services professionals can better understand market changes through increased business insight from data, helping to anticipate performance gaps and more accurately assess investment alternatives.

Retail

\$140 billion in total sales are missed each year because retailers don't have the right products in stock to meet customer demand.⁴

With Big Data

Retailers can better understand their customers by analyzing sales trends and incorporating more accurate forecasting, ultimately increasing customer loyalty and revenue.

Communications

5 billion global subscribers in the tele industry are demanding unique and personalized offerings that match their individual lifestyles.⁵

With Big Data

Communications providers can use data to create a more personalized customer experience and avoid losing customers to competitors.

Information gathered by IBM

1. The McKinsey Quarterly, "The Power of Personalized Marketing," October 2014

2. The McKinsey Quarterly, "The Power of Personalized Marketing," October 2014

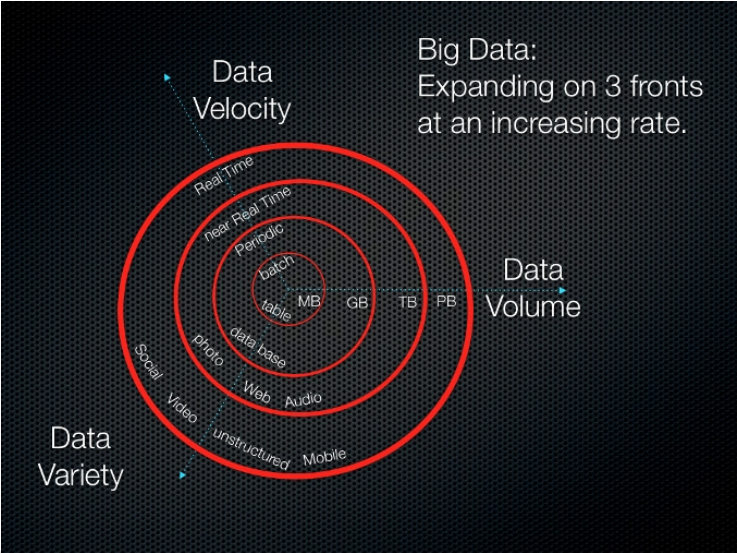
3. The McKinsey Quarterly, "The Power of Personalized Marketing," October 2014

4. The McKinsey Quarterly, "The Power of Personalized Marketing," October 2014

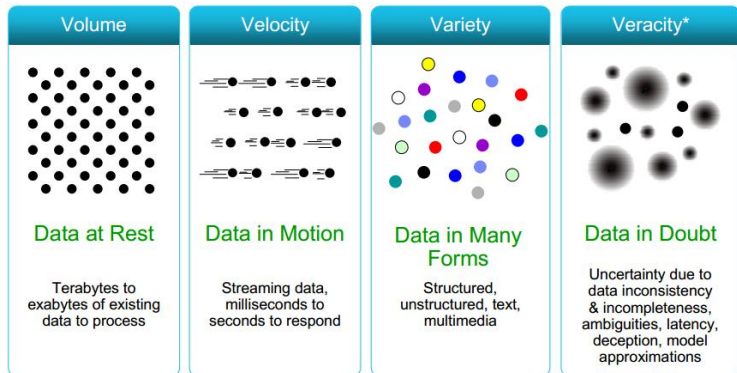
5. The McKinsey Quarterly, "The Power of Personalized Marketing," October 2014

© Copyright IBM Corporation 2015. All rights reserved.

Big data challenges: the “three V’s”



The “three V’s” are becoming the “four V’s”



Big data == big **privacy** concerns...

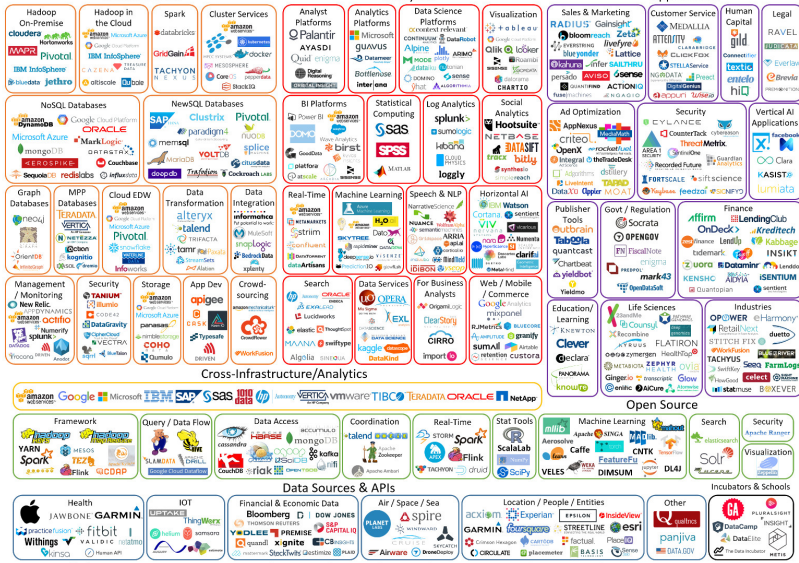


Big Data Landscape 2016 (Version 3.0)

Infrastructure

Analytics

Applications



Last Updated 3/23/2016

© Matt Turck (@mattturck), Jim Hao (@jimhao), & FirstMark Capital (@firstmarkcap)

FIRSTMARK

MapReduce

MapReduce was introduced by Google in 2004:

- ▶ Big data at that time: 20+ billion web pages \times 20 kB = 400+ TB
- ▶ One computer can read 30-35 MB/sec from disk
 - ⇒ 4 months to read the Web
 - ⇒ \sim 1000 hard drives just to store the web

MapReduce

MapReduce was introduced by Google in 2004:

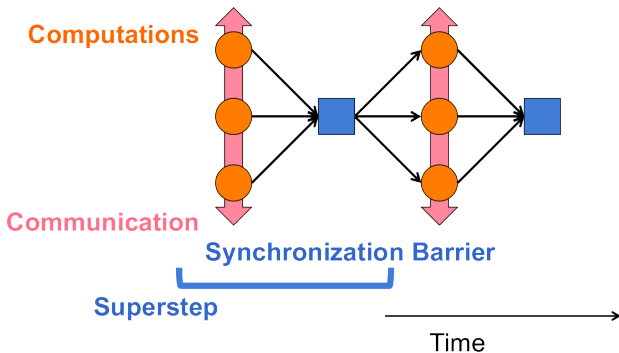
- ▶ Big data at that time: 20+ billion web pages \times 20 kB = 400+ TB
- ▶ One computer can read 30-35 MB/sec from disk
 - ⇒ 4 months to read the Web
 - ⇒ \sim 1000 hard drives just to store the web
- ▶ But they wanted to **process** the data! This requires much more computation, data, etc.

"Google Infrastructure for Massive Parallel Processing",

Walfredo Cirne, Presentation in the industrial track in CCGrid'2007.

The Bulk Synchronous Parallel model

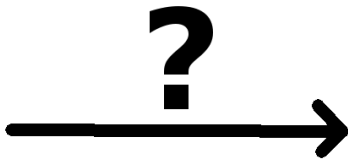
- ▶ Maximize I/O
- ▶ Minimize coordination



Parallelization is not so easy

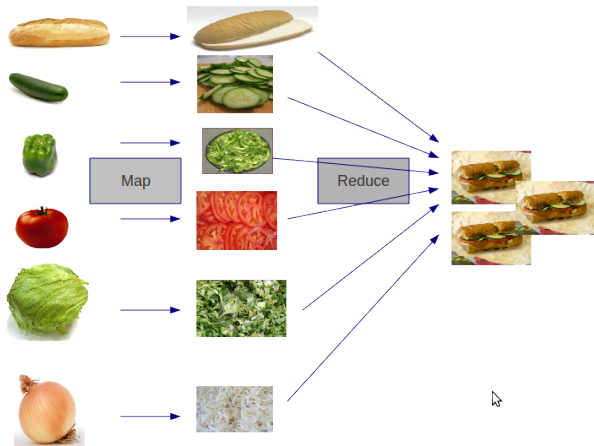
- ☺ “Easy” parallelization
 - ▶ Reading the Web on 1000 machines \Rightarrow less than 3 hours
- ☹ This requires lots of programming work
 - ▶ Communication & coordination
 - ▶ Debugging
 - ▶ Fault-tolerance
 - ▶ Management and monitoring
 - ▶ Optimization
- ☹ ☹ Repeat the same painful process for every problem you want to solve

Let's make sandwiches



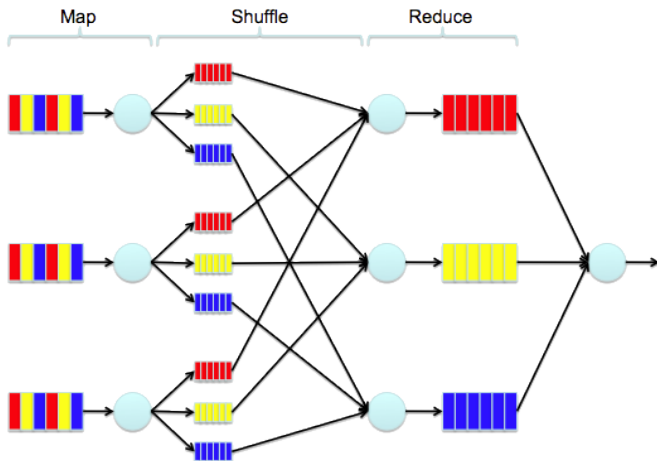
<https://twitter.com/tgrall>

Let's make sandwiches



<https://twitter.com/tgrall>

Back to MapReduce



<http://www.slideshare.net/lynnlangit/hadoop-mapreduce-fundamentals-21427224/>

Programmer must write two simple functions

- ▶ `map(key,value) → {key',value'}*`

The `map` function reads input data and produces intermediate tuples which are ready for the second phase

Programmer must write two simple functions

- ▶ $\text{map}(\text{key}, \text{value}) \rightarrow \langle \text{key}', \text{value}' \rangle^*$
The **map** function reads input data and produces intermediate tuples which are ready for the second phase
- ▶ $\text{reduce}(\text{key}', \langle \text{value}' \rangle^*) \rightarrow \langle \text{key}', \text{value}'' \rangle^*$
The **reduce** function takes all intermediate tuples with the same key, and produces output tuples

Example: word count

Let's take a (long) piece of text. Can we compute the number of occurrences of each word?

- ▶ **Map function:** take a subset of the input, generate one intermediate tuple for every word in the text

```
def map(String input_key, String doc):  
    for each word w in doc:  
        EmitIntermediate(w, 1)
```

- ▶ **Shuffle operation:** all tuples with the same **key** are automatically sent to the same reducer
- ▶ **Reduce function:** count the occurrences we received for each word

```
def reduce(String output_key, Iterator output_vals):  
    int res = 0  
    for each v in output_vals:  
        res = res + v  
    Emit(res)
```

What makes MapReduce so great

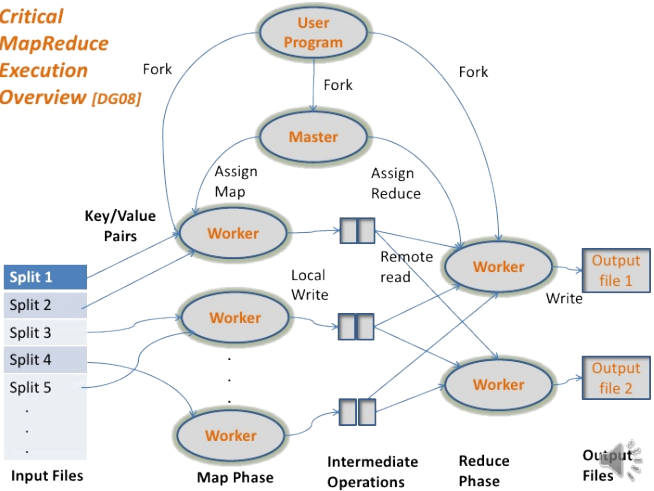
- ☺ `map()` functions run in parallel, creating different intermediate values from different input data sets
- ☺ `reduce()` functions also run in parallel, each working on a different output key
- ☺ All values are processed independently

What makes MapReduce so great

- 😊 **map()** functions run in parallel, creating different intermediate values from different input data sets
- 😊 **reduce()** functions also run in parallel, each working on a different output key
- 😊 All values are processed independently
- 😞 **Limitation:** the reduce phase cannot start until the map phase is totally finished

MapReduce architecture

Critical MapReduce Execution Overview [DG08]



MapReduce architecture

- ▶ **One master server, many worker servers**
 - ▶ Input data is split in chunks (~ 64 MB)
 - ▶ Tasks are assigned to workers dynamically
- ▶ The master assigns each **map task** to a free worker
 - ▶ Considers locality of data to worker when assigning task
 - ▶ Worker reads task input (often from local disk!)
 - ▶ Worker produces R local files containing intermediate key/value pairs
- ▶ The master assigns each **reduce task** to a free worker
 - ▶ Worker reads intermediate key/value pairs from map workers
 - ▶ Worker sorts & applies the user's Reduce function to produce the output


Fault tolerance

- ▶ If a worker fails:
 - ▶ The master will detect failure thanks to periodic heartbeats
 - ▶ The master re-executes the completed and in-progress map() tasks
 - ▶ The re-executes the in-progress reduce() tasks
- ▶ If the same input *always* makes the map() function crash:
 - ▶ The master will detect it and skip these values on re-execution

Fault tolerance

- ▶ If a worker fails:
 - ▶ The master will detect failure thanks to periodic heartbeats
 - ▶ The master re-executes the completed and in-progress map() tasks
 - ▶ The re-executes the in-progress reduce() tasks
- ▶ If the same input *always* makes the map() function crash:
 - ▶ The master will detect it and skip these values on re-execution
- ▶ If the master fails:

Fault tolerance

- ▶ If a worker fails:
 - ▶ The master will detect failure thanks to periodic heartbeats
 - ▶ The master re-executes the completed and in-progress map() tasks
 - ▶ The re-executes the in-progress reduce() tasks
- ▶ If the same input *always* makes the map() function crash:
 - ▶ The master will detect it and skip these values on re-execution
- ▶ If the master fails:


MapReduce in the real world



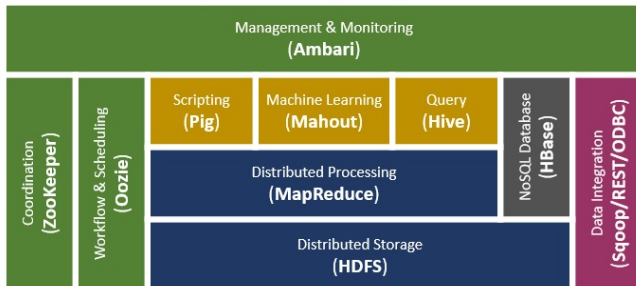
The reference open-source
implementation: Apache
Hadoop



All the good clouds provide
Hadoop (or similar)
under a PaaS model

Hadoop's stack

Apache Hadoop Ecosystem



<http://bit.do/cSi9J>

Limitations of MapReduce

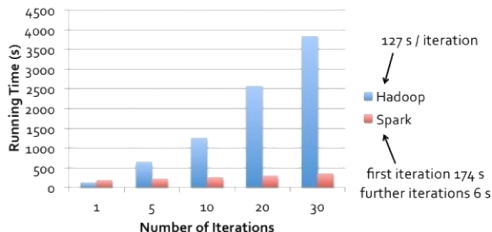
- ▶ **Rigid programming model**
 - ▶ Programs must be developed as a pair of map/reduce functions
 - ▶ Complex programs may be designed as a succession of iterative map/reduce steps
 - ▶ By default, **all intermediate result are stored in the HDFS file system**
 - ▶ Replicated, fault-tolerant, etc.
 - ▶ But there are *lots* of intermediate results in a map/reduce/map/reduce/map/reduce program!
- ⇒ Not-so-great performance

Spark

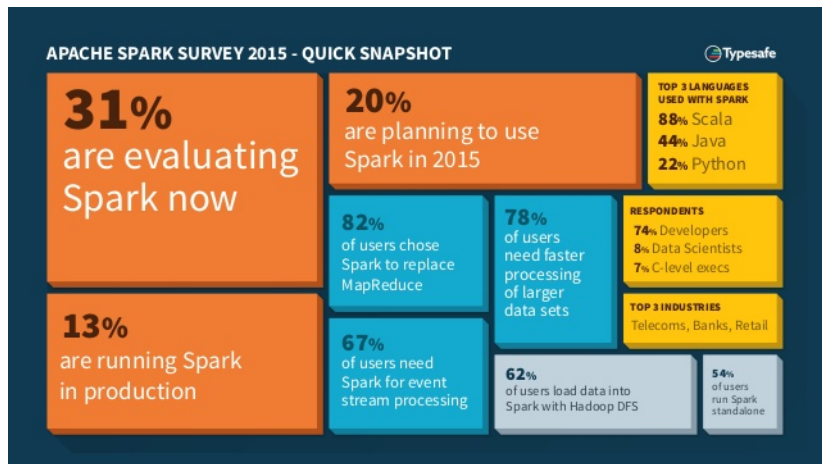
- ▶ Let's simplify application development: write “normal” code, let the system figure out how to execute it efficiently
 - ▶ Let's use **main memory** for all intermediate data
- ⇒ Major performance improvements



Logistic Regression Performance

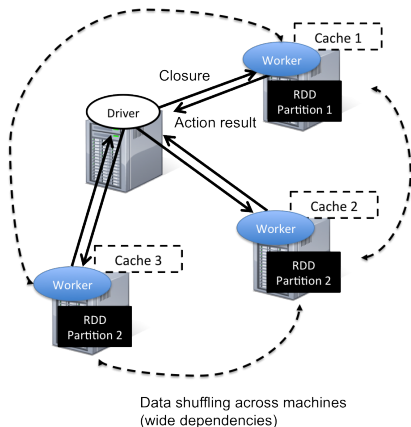


Survey within the big data developer community (2015)



Spark's architecture

- ▶ **One driver node**
~ master
 - ▶ orchestrates computation, assigns work
- ▶ **Many worker nodes**
 - ▶ execute tasks, report to driver node



<http://horicky.blogspot.fr/2013/12/spark-low-latency-massively-parallel.html>

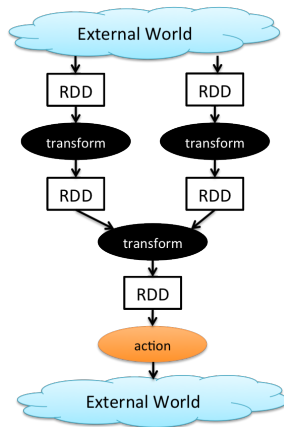
Spark RDDs

- ▶ RDD = **Resilient Distributed Dataset**
- ▶ Conceptually an array (or a map) of entries
 - ▶ Entries might be strings, numbers, maps, pairs, ...
 - ▶ Transparently partitioned / distributed by Spark
 - ▶ Transparently resilient (either by recomputation or storage)
- ▶ Creation:
 - ▶ Read from a local or distributed file system
 - ▶ Or produced by another Spark computation

Spark applications

A Spark application is composed of **transformations** and **actions**:

- ▶ **Transformations** specify how to produce an RDD from another RDD
 - ▶ But the system does not execute them immediately
- ▶ **Actions** trigger an actual computation
 - ▶ The system explores the graph of dependencies, and produces a directed acyclic graph of necessary transformations
 - ▶ Optimizes computations to be done
 - ▶ Distributes and organizes work



Spark transformations

Transformations

The following table lists some of the common transformations supported by Spark. Refer to the RDD API doc ([Scala](#), [Java](#), [Python](#), [R](#)) and pair RDD functions doc ([Scala](#), [Java](#)) for details.

Transformation	Meaning
<code>map(func)</code>	Return a new distributed dataset formed by passing each element of the source through a function <i>func</i> .
<code>filter(func)</code>	Return a new dataset formed by selecting those elements of the source on which <i>func</i> returns true.
<code>flatMap(func)</code>	Similar to <code>map</code> , but each input item can be mapped to 0 or more output items (so <i>func</i> should return a <code>Seq</code> rather than a single item).
<code>mapPartitions(func)</code>	Similar to <code>map</code> , but runs separately on each partition (block) of the RDD, so <i>func</i> must be of type <code>Iterator<T> => Iterator<U></code> when running on an RDD of type <code>T</code> .
<code>mapPartitionsWithIndex(func)</code>	Similar to <code>mapPartitions</code> , but also provides <i>func</i> with an integer value representing the index of the partition, so <i>func</i> must be of type <code>(Int, Iterator<T>) => Iterator<U></code> when running on an RDD of type <code>T</code> .

etc...

<https://spark.apache.org/docs/latest/programming-guide.html>

Spark actions

Actions

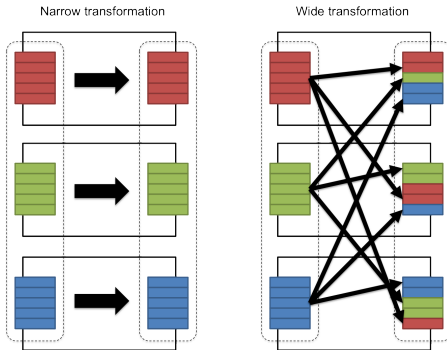
The following table lists some of the common actions supported by Spark. Refer to the RDD API doc ([Scala](#), [Java](#), [Python](#), [R](#)) and pair RDD functions doc ([Scala](#), [Java](#)) for details.

Action	Meaning
<code>reduce(func)</code>	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.
<code>collect()</code>	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.
<code>count()</code>	Return the number of elements in the dataset.
<code>first()</code>	Return the first element of the dataset (similar to <code>take(1)</code>).
<code>take(n)</code>	Return an array with the first <i>n</i> elements of the dataset.

etc...

<https://spark.apache.org/docs/latest/programming-guide.html>

Wide and Narrow Transformations



- ▶ Wide transformations require shuffling
 - ▶ e.g., `reduceByKey(...)`
 - ▶ Network costs, higher latency

The word count example in Spark

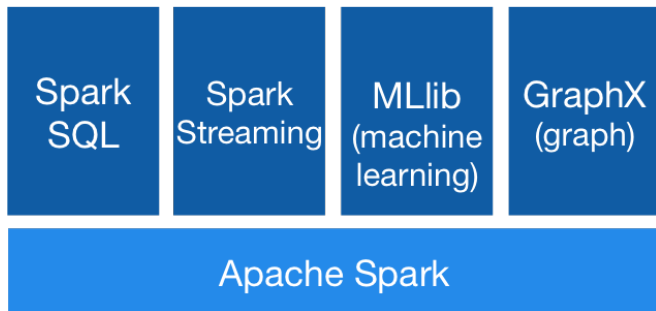
Transformations

```
val textFile = sc.textFile("hdfs://...")
val counts = textFile.flatMap(line => line.split(" "))
                    .map(word => (word, 1))
                    .reduceByKey(_ + _)
counts.saveAsTextFile("hdfs://...")
```

Action

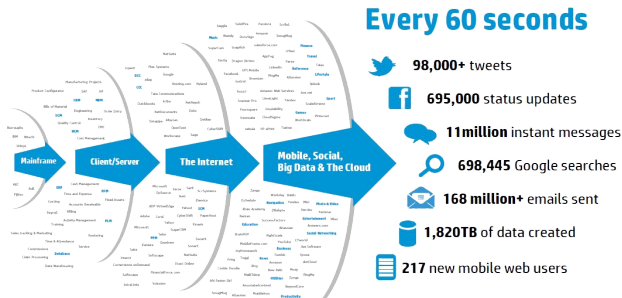
- ▶ The developer writes **simple, sequential code** using the Spark transformations and actions
- ▶ Spark **automatically parallelizes the code**, distributes it across many nodes, and coordinates the distributed execution

Spark's stack



<http://spark.apache.org/>

Limitations of MapReduce/Spark



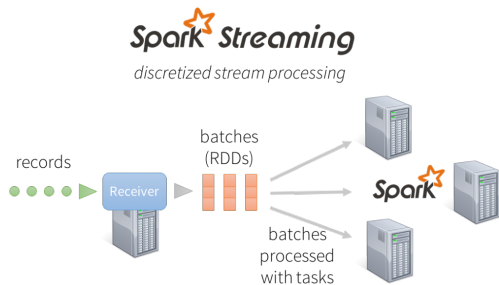
How can we keep up with the **velocity** of big data?

- ▶ Store incoming data (e.g., tweets)
 - ▶ One in a while: process the new data, produce new results
- ⇒ The results are always late!

☞ We need to be able to process incoming data **in real time**, not as a succession of batch jobs

Spark Streaming

- ▶ Spark Streaming relies on **micro-batches**
 - ▶ Ingest incoming real-time data from various sources
 - ▶ Generate a new “micro-batch” at fixed time intervals (e.g., 1 second)
 - ▶ Process each micro-batch as a separate Spark job



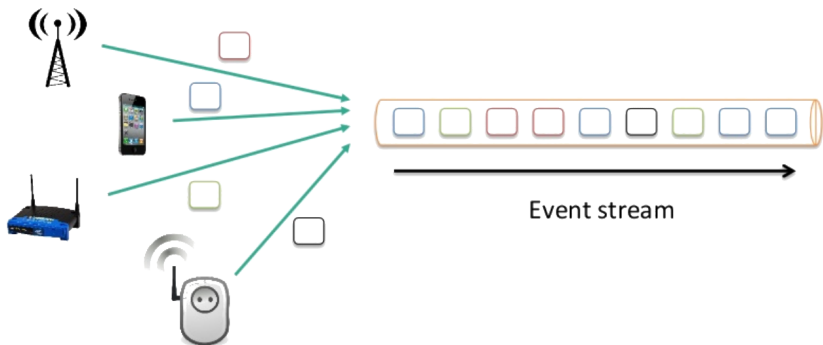
records processed in batches with short tasks
each batch is a RDD (partitioned dataset)

Limitations of micro-batches

- ▶ Data arriving out of order is hard to handle
 - ▶ How do you detect missing data, data gaps, correct out of time order data etc?
- ▶ Batch length restricts Window-based analytics
 - ▶ Large batches \Rightarrow poor responsiveness
 - ▶ Small batches \Rightarrow the system is obliged to work on very small window sizes
- ▶ Code is hard to write
 - ▶ As soon as you try to **update** existing results with each micro-batch

<http://bit.do/cSi4L>

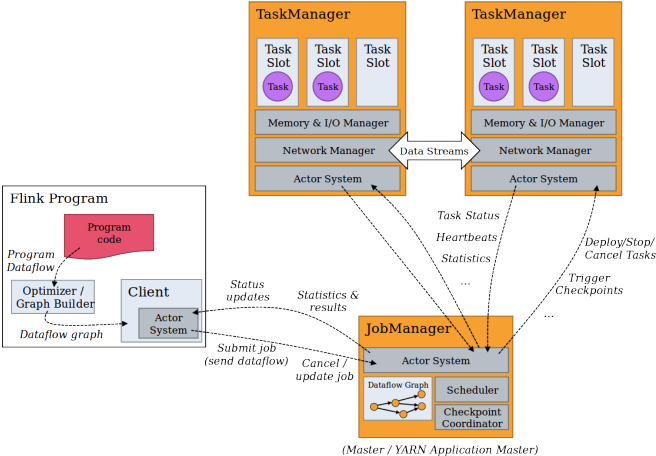
Many applications are fundamentally based on **streaming**



Apache Flink is a big-data framework based on a **distributed streaming dataflow** engine



Flink's architecture

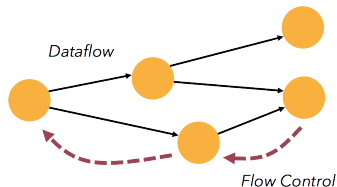


<https://ci.apache.org/projects/flink/flink-docs-release-1.1/concepts/concepts.html>

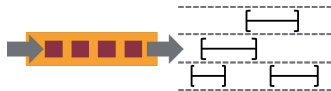
Streaming operators

Flink uses similar directed acyclic graphs (DAGs) of operators to Spark. But:

- ▶ In streaming mode, the DAG remains in place, and data **flows** along the DAG

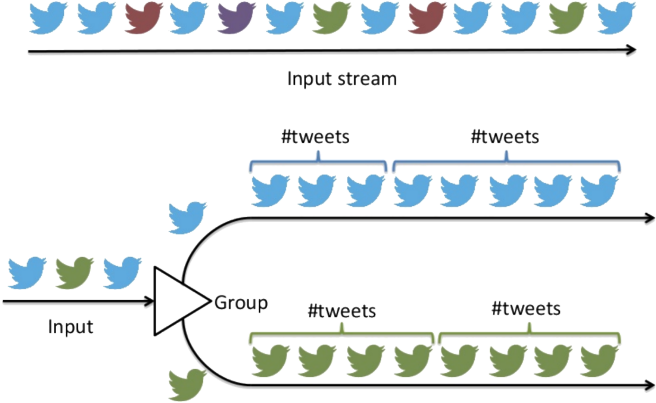


- ▶ Each operator works over a **window** of data items



<http://flink.apache.org/features.html>

Example: classify and count tweets

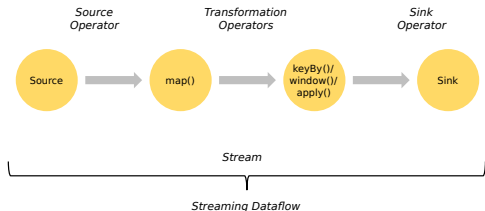


<http://www.slideshare.net/tillrohrmann/apache-flink-streaming-done-right-fosdem-2016>

Flink programs are compiled into an operator DAG

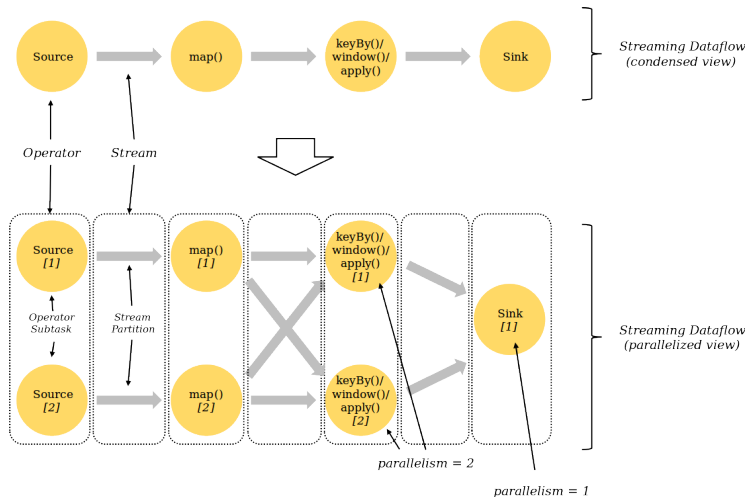
```
DataStream<String> lines = env.addSource(  
    new FlinkKafkaConsumer<>(...));  
DataStream<Event> events = lines.map((line) -> parse(line));  
DataStream<Statistics> stats = events  
    .keyBy("id")  
    .timeWindow(Time.seconds(10))  
    .apply(new MyWindowAggregationFunction());  
stats.addSink(new RollingSink(path));
```

Source
Transformation
Transformation
Sink



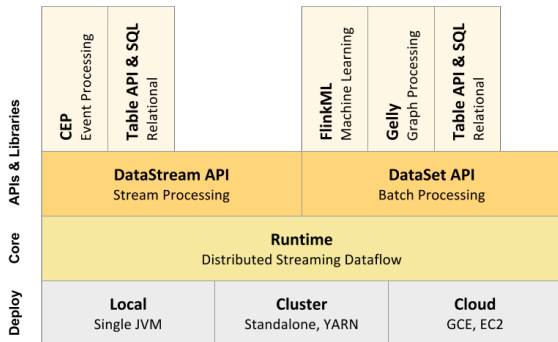
<https://ci.apache.org/projects/flink/flink-docs-release-1.1/concepts/concepts.html>

And each operator can be parallelized



<https://ci.apache.org/projects/flink/flink-docs-release-1.1/concepts/concepts.html>

Flink's stack



<https://ci.apache.org/projects/flink/flink-docs-release-1.1/>

Conclusion

- ▶ Processing big data is very difficult
 - ▶ Volume, Variety, Velocity
 - ▶ Parallel programming is hard!
- ▶ Cloud frameworks are being proposed to facilitate the developers' task
 1. **MapReduce** automatically parallelizes programs expressed as pairs of map/reduce functions
 2. **Spark** simplifies the development model by automatically compiling sequential code based on specific operators
 3. **Flink** extends Spark with data stream processing
- ▶ Many new frameworks are being proposed. Stay tuned for very fast progress in this exciting domain!