1



# Spark Introduction

2

# Agenda

- History of Spark
- Introduction
- Components of Stack
- Resilient Distributed Dataset – RDD

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

3

# History of Spark



2004
MapReduce paper

2010
Spark paper

| 2002 | | 2004 | | 2006 | | 2008 | | 2010 | | 2012 | | 2014 |

2002
MapReduce @ Google

2008
Hadoop Summit

2014
Apache Spark top-level

2006
Hadoop @Yahoo!

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

4

# History of Spark

circa 1979 – **Stanford**, **MIT**, **CMU**, etc.
 set/list operations in LISP, Prolog, etc., for parallel processing
**www-formal.stanford.edu/jmc/history/lisp/lisp.htm**

circa 2004 – **Google**
 *MapReduce: Simplified Data Processing on Large Clusters* Jeffrey Dean and Sanjay Ghemawat
**research.google.com/archive/mapreduce.html**

circa 2006 – **Apache**
*Hadoop*, originating from the Nutch Project  Doug Cutting
**research.yahoo.com/files/cutting.pdf**

circa 2008 – **Yahoo**
*web scale search indexing Hadoop Submit, HUG, etc.*
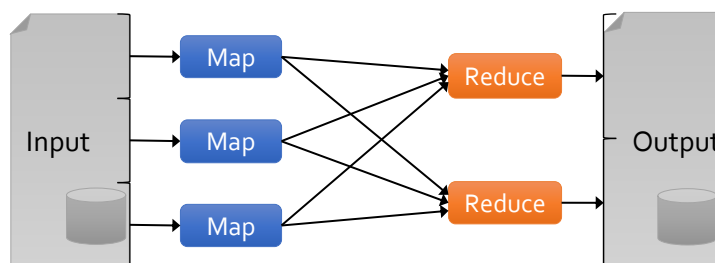 **developer.yahoo.com/hadoop/**

circa 2009 – **Amazon AWS**
 Elastic MapReduce
 Hadoop modified for EC2/S3, plus support for Hive, Pig, Cascading, etc.
**aws.amazon.com/elasticmapreduce/**

SOICT  VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

5

# MapReduce

Most current cluster programming models are based on *acyclic data flow* from stable storage to stable storage



SOICT  VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

6

3

# MapReduce

- Acyclic data flow is inefficient for applications that repeatedly reuse a *working set* of data:
    - **Iterative** algorithms (machine learning, graphs)
    - **Interactive** data mining tools (R, Excel, Python)

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

7

# Data Processing Improvement Goals

- **Low latency (interactive) queries on historical data:** enable faster decisions
    - E.g., identify why a site is slow and fix it

- **Low latency queries on live data (streaming):** enable decisions on real-time data
    - E.g., detect & block worms in real-time (a worm may infect **1mil** hosts in **1.3sec**)

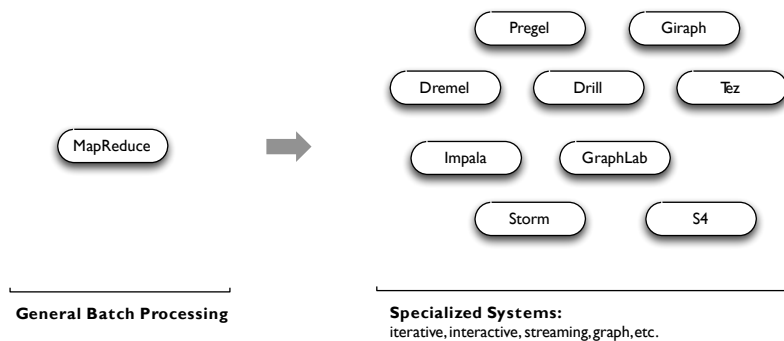- **Sophisticated data processing:** enable "better" decisions
    - E.g., anomaly detection, trend analysis

Therefore, people built specialized systems as workarounds…

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

8

# Specialized Systems

General Batch Processing — MapReduce

Specialized Systems:
iterative, interactive, streaming, graph, etc.

Pregel, Giraph, Dremel, Drill, Tez, Impala, GraphLab, Storm, S4

*The State of Spark, and Where We're Going Next*
**Matei Zaharia**
Spark Summit (2013)
youtu.be/ nU6vO2EJAb4

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

9

# Storage vs Processing Wars

**NoSQL battles**

Relational vs NoSQL

HBase vs Cassanrdra

Redis vs Memcached vs Riak

MongoDB vs CouchDB vs Couchbase

Neo4j vs Titan vs Giraph vs OrientDB

Solr vs Elasticsearch

**Compute battles**

MapReduce vs Spark

Spark Streaming vs Storm
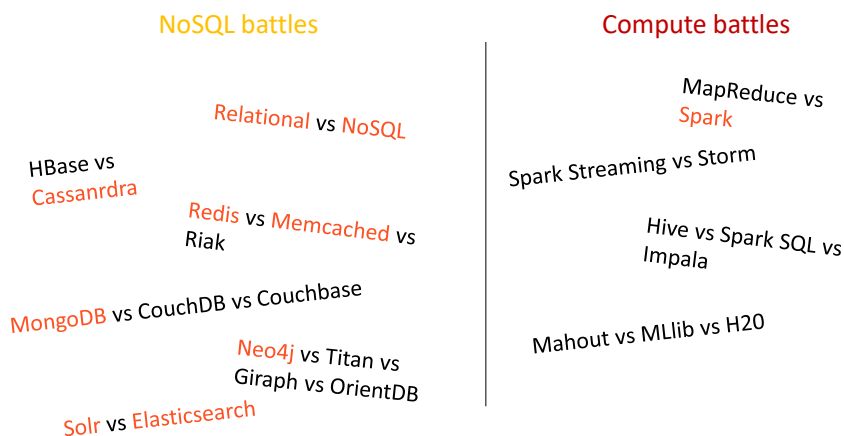
Hive vs Spark SQL vs Impala

Mahout vs MLlib vs H20
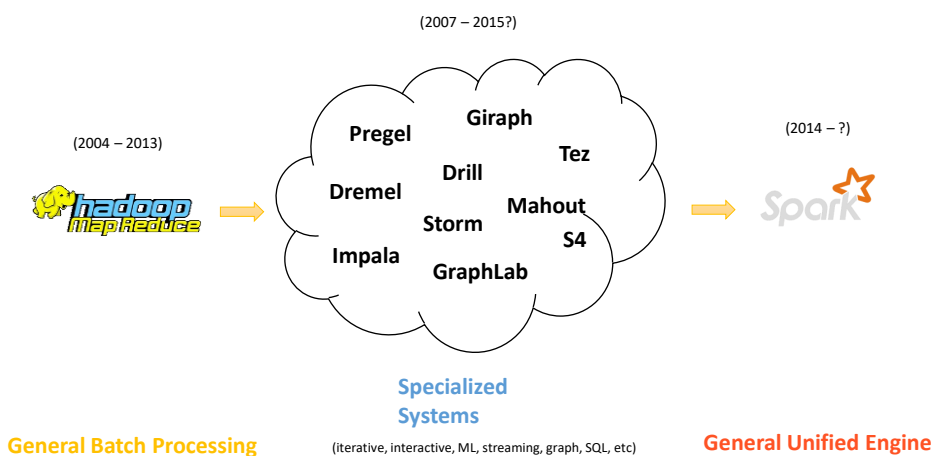
VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

10

# Storage vs Processing Wars

### NoSQL battles

Relational vs NoSQL

HBase vs Cassanrdra

Redis vs Memcached vs Riak

MongoDB vs CouchDB vs Couchbase

Neo4j vs Titan vs Giraph vs OrientDB

Solr vs Elasticsearch

### Compute battles

MapReduce vs Spark

Spark Streaming vs Storm

Hive vs Spark SQL vs Impala

Mahout vs MLlib vs H20

SOICT VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

11

# Specialized Systems

(2007 – 2015?)

(2004 – 2013)

(2014 – ?)

Pregel  Giraph  Tez
Drill
Dremel  Mahout
Storm  S4
Impala  GraphLab

**Specialized Systems**

**General Batch Processing**   (iterative, interactive, ML, streaming, graph, SQL, etc)   **General Unified Engine**

SOICT VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

12

**VS**

YARN — Mesos

Tachyon

SQL

MLlib

STORM — Streaming

13



OOZIE

10X – 100X
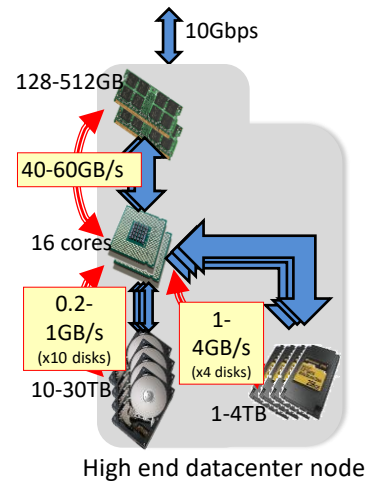
14

# Support Interactive and Streaming Comp.

- Aggressive use of *memory*
- Why?
  1. Memory transfer rates >> disk or SSDs
  2. Many datasets already fit into memory
     - Inputs of over 90% of jobs in Facebook, Yahoo!, and Bing clusters fit into memory
     - e.g., 1TB = 1 billion records @ 1KB each
  3. Memory density (still) grows with Moore's law
     - RAM/SSD hybrid memories at horizon



10Gbps

128-512GB

40-60GB/s

16 cores

0.2-1GB/s (x10 disks)

1-4GB/s (x4 disks)

10-30TB

1-4TB

High end datacenter node

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

15

# Support Interactive and Streaming Comp.
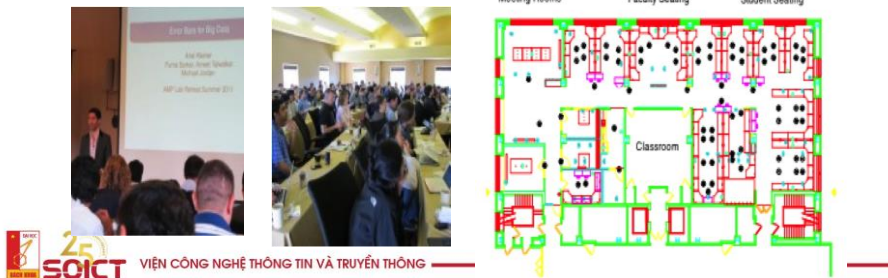
- Increase *parallelism*
- Why?
  - Reduce work per node → improve latency
- Techniques:
  - Low latency parallel scheduler that achieve high locality
  - Optimized parallel communication patterns (e.g., shuffle, broadcast)
  - Efficient recovery from failures and straggler mitigation



result

T

result

$T^{new}$ (< T)

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

16

# Berkeley AMPLab

- ▪ "Launched" January 2011: 6 Year Plan
- 8 CS Faculty
- ~40 students
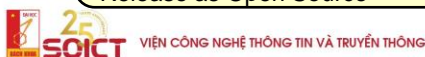- 3 software engineers
- Organized for collaboration:



17

# Berkeley AMPLab

- Funding:
  - DARPA XData, NSF CISE Expedition Grant
  - Industrial, founding sponsors amazon.com Google SAP
  - 18 other sponsors, including



**Goal:** Next Generation of Analytics Data Stack for Industry & Research:
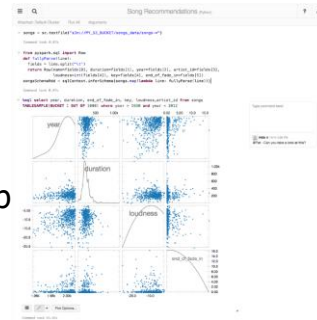• Berkeley Data Analytics Stack (BDAS)
• Release as Open Source

18

# Databricks



making big data simple



- Founded in late 2013
- by the creators of Apache Spark
- Original team from UC Berkeley AMPLab
- Raised $47 Million in 2 rounds

Databricks Cloud:
"A unified platform for building Big Data pipelines – from ETL to Exploration and Dashboards, to Advanced Analytics and Data Products."
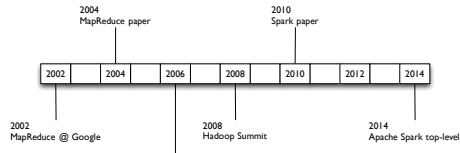
VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

19

---

The Databricks team contributed more than **75%** of the code
added to Spark in the 2014



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

20

# History of Spark

*Spark: Cluster Computing with Working Sets*
Matei Zaharia, Mosharaf Chowdhury,
Michael J. Franklin, Scott Shenker, Ion Stoica
USENIX HotCloud (2010)
people.csail.mit.edu/ matei/ papers/ 2010/ hotcloud_spark.pdf

*Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing*
Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave,
Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker,
Ion Stoica NSDI (2012)
usenix.org/ system/ files/ conference/ nsdi12/ nsdi12-final138.pdf

SOICT VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

21

---

# History of Spark

April 2012

SOICT VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

"We present Resilient Distributed Datasets (RDDs), a distributed memory abstraction that lets programmers perform in-memory computations on large clusters in a fault-tolerant manner.

RDDs are motivated by two types of applications that current computing frameworks handle inefficiently: iterative algorithms and interactive data mining tools.

In both cases, keeping data in memory can improve performance by an order of magnitude."
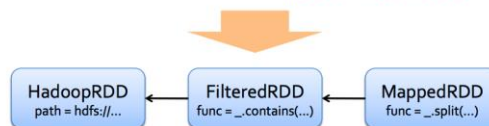
22

# History of Spark

## RDD Fault Tolerance

RDDs track the series of transformations used to build them (their *lineage*) to recompute lost data

E.g: `messages = textFile(...).filter(_.contains("error"))`
`.map(_.split('\t')(2))`

```
HadoopRDD              FilteredRDD              MappedRDD
path = hdfs://...      func = _.contains(...)   func = _.split(...)
```

*The State of Spark, and Where We're Going Next*
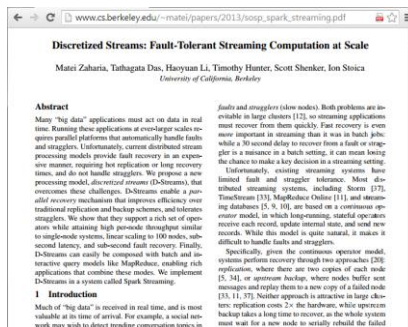**Matei Zaharia**
Spark Summit (2013)
youtu.be/ nU6yO2EJAb4

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

23

# History of Spark

## Spark STREAMING

Analyze real time streams of data in ½ second intervals

**Discretized Streams: Fault-Tolerant Streaming Computation at Scale**

Matei Zaharia, Tathagata Das, Haoyuan Li, Timothy Hunter, Scott Shenker, Ion Stoica
*University of California, Berkeley*

```
TwitterUtils.createStream(...)
    .filter(_.getText.contains("Sp
ark"))
    .countByWindow(Seconds(5))
```

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

24

12

# History of Spark

**Spark** SQL

Seemlessly mix SQL queries with Spark programs.

**Spark SQL: Relational Data Processing in Spark**

Michael Armbrust[1], Reynold S. Xin[1], Cheng Lian[1], Yin Huai[1], Davies Liu[1], Joseph K. Bradley[1], Xiangrui Meng[1], Tomer Kaftan[2], Michael J. Franklin[1][3], Ali Ghodsi[1], Matei Zaharia[1][4]

[1]Databricks Inc.   [2]MIT CSAIL   [3]AMPLab, UC Berkeley

**ABSTRACT**

Spark SQL is a new module in Apache Spark that integrates relational processing with Spark's functional programming API. Built on our experience with Shark, Spark SQL lets Spark programmers leverage the benefits of relational processing (e.g., declarative queries and optimized storage), and lets SQL users call complex analytics libraries in Spark (e.g., machine learning). Compared to previous systems, Spark SQL makes two main additions. First, it offers much tighter integration between relational and procedural processing, through a declarative DataFrame API that integrates with procedural Spark code. Second, it includes a highly extensible optimizer, Catalyst, built using features of the Scala programming language, that makes it easy to add composable rules, control code generation, and define extension points. Using Catalyst, we have built a variety of features (e.g., schema inference for JSON, machine learning types, and query federation to external databases) tailored for the complex needs of modern data analysis. We see Spark SQL as an evolution of both SQL-on-Spark and of Spark itself, offering richer APIs and optimizations while keeping the benefits of the Spark programming model.

**Categories and Subject Descriptors**

H.2 [**Database Management**]: Systems

**Keywords**

Databases; Data Warehouse; Machine Learning; Spark; Hadoop

**1  Introduction**

Big data applications require a mix of processing techniques, data sources and storage formats. The earliest systems designed for these workloads, such as MapReduce, gave users a powerful, but

While the popularity of relational systems shows that users often prefer writing declarative queries, the relational approach is insufficient for many big data applications. First, users want to perform ETL, to and from various data sources that might be semi- or unstructured, requiring custom code. Second, users want to perform advanced analytics, such as machine learning and graph processing, that are challenging to express in relational systems. In practice, we have observed that most data pipelines would ideally be expressed with a combination of both relational queries and complex procedural algorithms. Unfortunately, these two classes of systems—relational and procedural—have until now remained largely disjoint, forcing users to choose one paradigm or the other.

This paper describes our effort to combine both models in Spark SQL, a major new component in Apache Spark [39]. Spark SQL builds on our earlier SQL-on-Spark effort, called Shark. Rather than forcing users to pick between a relational or a procedural API, however, Spark SQL lets users seamlessly intermix the two.

Spark SQL bridges the gap between the two models through two contributions. First, Spark SQL provides a DataFrame API that can perform relational operations on both external data sources and Spark's built-in distributed collections. This API is similar to the widely used data frame concept in R [32], but evaluates operations lazily so that it can perform relational optimizations. Second, to support the wide range of data sources and algorithms in big data, Spark SQL introduces a novel extensible optimizer called Catalyst. Catalyst makes it easy to add data sources, optimization rules, and data types for domains such as machine learning.

The DataFrame API offers rich relational/procedural integration within Spark programs. DataFrames are collections of structured records that can be manipulated using Spark's procedural API, or using new relational APIs that allow richer optimizations. They can

```
sqlCtx = new HiveContext(sc)
results = sqlCtx.sql(
  "SELECT * FROM people")
names = results.map(lambda p:
p.name)
```

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

25

---

# History of Spark

**Spark** GRAPHX

Analyze networks of nodes and edges using graph processing

**GraphX: A Resilient Distributed Graph System on Spark**

Reynold S. Xin, Joseph E. Gonzalez, Michael J. Franklin, Ion Stoica

AMPLab, EECS, UC Berkeley
{rxin, jegonzal, franklin, istoica}@cs.berkeley.edu

**ABSTRACT**

From social networks to targeted advertising, big graphs capture the structure in data and are central to recent advances in machine learning and data mining. Unfortunately, directly applying existing data-parallel tools to graph computation tasks can be cumbersome and inefficient. The need for intuitive, scalable tools for graph computation has lead to the development of new graph-parallel systems (e.g., Pregel, PowerGraph) which are designed to efficiently execute graph algorithms. Unfortunately, these new graph-parallel systems do not address the challenges of graph construction and transformation which are often just as problematic as the subsequent computation. Furthermore, existing graph-parallel systems provide limited fault-tolerance and support for interactive data mining.

We introduce GraphX, which combines the advantages of both data-parallel and graph-parallel systems by efficiently expressing graph computation within the Spark data-parallel framework. We leverage new ideas in distributed graph representation to efficiently distribute graphs as tabular data-structures. Similarly, we leverage advances in data-flow systems to exploit in-memory computation and fault-tolerance. We provide powerful new operations to simplify graph construction and transformation. Using these primitives we implement the PowerGraph and Pregel abstractions in less than 20 lines of code. Finally, by exploiting the Scala foundation of Spark, we enable users to interactively load, transform, and compute on massive graphs.

**1.  INTRODUCTION**

From social networks to advertising and the web, big graphs can be found in a wide range of important applications. By modeling the

and distributed systems. By abstracting away the challenges of large-scale distributed system design, these frameworks simplify the design, implementation, and application of new sophisticated graph algorithms to large-scale real-world graph problems.

While existing graph-parallel frameworks share many common properties, each presents a slightly different view of graph computation tailored to either the originating domain or a specific family of graph algorithms and applications. Unfortunately, because each framework relies on a separate runtime, it is difficult to compose these abstractions. Furthermore, while these frameworks address the challenges of graph computation, they do not address the challenges of data ETL (preprocessing and construction) or the process of interpreting and applying the results of computation. Finally, few frameworks have built-in support for interactive graph computation.

Alternatively, data-parallel systems like MapReduce and Spark [12] are designed for scalable data processing and are well suited to the task of graph construction (ETL). By exploiting data-parallelism, these systems are highly scalable and support a range of fault-tolerance strategies. More recent systems like Spark even enable interactive data processing. However, naively expressing graph computation and graph algorithms in these data-parallel abstractions can be challenging and typically leads to complex joins and excessive data movement that does not exploit the graph structure.

To address these challenges we introduce GraphX, a graph computation system which runs in the Spark data-parallel framework. GraphX extends Spark's Resilient Distributed Dataset (RDD) abstraction to introduce the Resilient Distributed Graph (RDG), which associates records with vertices and edges in a graph and provides a collection of expressive computational primitives. Using these

```
graph = Graph(vertices, edges)
messages =
spark.textFile("hdfs://...")
graph2 =
graph.joinVertices(messages) {
  (id, vertex, msg) => ...
}
```
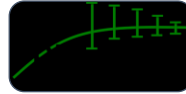
https://amplab.cs.berkeley.edu/wp-content/uploads/2013/05/grades-graphx_with_fonts.pdf

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

26

# History of Spark



SQL queries with Bounded Errors and Bounded Response Times

**BlinkDB: Queries with Bounded Errors and Bounded Response Times on Very Large Data**

Sameer Agarwal†, Barzan Mozafari*, Aurojit Panda†, Henry Milner†, Samuel Madden*, Ion Stoica*†

†University of California, Berkeley     *Massachusetts Institute of Technology     *Conviva Inc.
{sameerag, apanda, henrym, istoica}@cs.berkeley.edu, {barzan, madden}@csail.mit.edu

**Abstract**

In this paper, we present BlinkDB, a massively parallel, approximate query engine for running interactive SQL queries on large volumes of data. BlinkDB allows users to trade-off query accuracy for response time, enabling interactive queries over massive data by running queries on data samples and presenting results annotated with meaningful error bars. To achieve this, BlinkDB uses two key ideas: (1) an adaptive optimization framework that builds and maintains a set of multi-dimensional stratified samples from original data over time, and (2) a dynamic sample selection strategy that selects an appropriately sized sample based on a query's accuracy or response time requirements. We evaluate BlinkDB against the well-known TPC-H benchmarks and a real-world analytic workload derived from Conviva Inc., a company that manages video distribution over the Internet. Our experiments on a 100 node cluster show that BlinkDB can answer queries on up to 17 TBs of data in less than 2 seconds (over 200× faster than Hive), within an error of 2–10%.

**1.   Introduction**

Modern data analytics applications involve computing aggregates over a large number of records to *roll-up* web clicks,

cessing of large amounts of data by trading result accuracy for response time and space. These techniques include sampling [10, 14], sketches [12], and on-line aggregation [15]. To illustrate the utility of such techniques, consider the following simple query that computes the average SessionTime over all users originating in New York:

```
SELECT AVG(SessionTime)
FROM Sessions
WHERE City = 'New York'
```

Suppose the Sessions table contains 100 million tuples for New York, and cannot fit in memory. In that case, the above query may take a long time to execute, since disk reads are expensive, and such a query would need multiple disk accesses to stream through all the tuples. Suppose we instead executed the same query on a sample containing only 10,000 New York tuples, such that the entire sample fits in memory. This would be orders of magnitude faster, while still providing an approximate result within a few percent of the actual value, an accuracy good enough for many practical purposes. Using sampling theory we could even provide confidence bounds on the accuracy of the answer [16].

Previously described approximation techniques make different trade-offs between efficiency and the generality of the

| SELECT avg(sessionTime) | SELECT avg(sessionTime) |
|---|---|
| FROM Table | FROM Table |
| WHERE city='San Francisco' | WHERE city='San Francisco' |
| WITHIN 2  SECONDS | ERROR 0.1 CONFIDENCE 95.0% |
| Queries with Time Bounds | Queries with Error Bounds |

https://www.cs.berkeley.edu/~sameerag/blinkdb_eurosys13.pdf

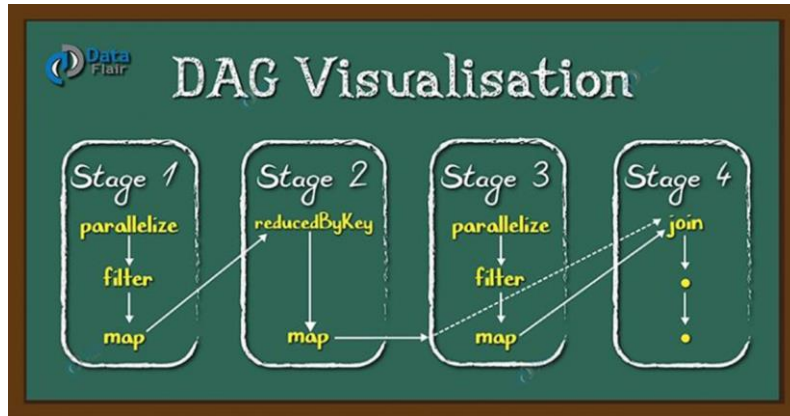VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

# History of Spark

• Unlike the various specialized systems, Spark's goal was to *generalize* MapReduce to support new apps within same engine

• Two reasonably small additions are enough to express the previous models:
  • *fast data sharing*
  • *general DAGs*

• This allows for an approach which is more efficient for the engine, and much simpler for the end users

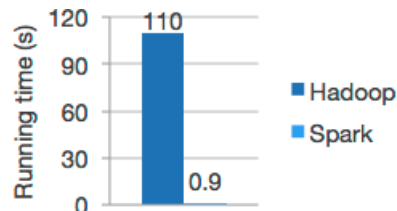VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

# Directed Acyclic Graph - DAG

29

# What is Apache Spark

- Spark is a unified analytics engine for large-scale data processing
- Speed: run workloads 100x faster
  - High performance for both batch and streaming data
  - Computations run in memory



Logistic regression in Hadoop and Spark

30

# What is Apache Spark

- Ease of Use: write applications quickly in Java, Scala, Python, R, SQL
  - Offer over 80 high-level operators
  - Use them interactively form Scala, Python, R, and SQL

```
df = spark.read.json("logs.json")
df.where("age > 21")
select("name.first").show()
```
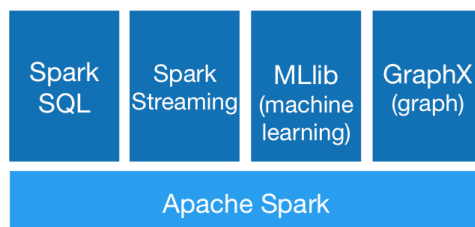
Spark's Python DataFrame API
Read JSON files with automatic schema inference

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

31

# What is Apache Spark

- Generality: combine SQL, Streaming, and complex analytics
  - Provide libraries including SQL and DataFrames, Spark Streaming, MLib, GraphX
  - Wide range of workloads e.g., batch applications, interactive algorithms, interactive queries, streaming

| Spark SQL | Spark Streaming | MLlib (machine learning) | GraphX (graph) |
|---|---|---|---|
| Apache Spark | | | |

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

32

# What is Apache Spark

- Run Everywhere:
  - run on Hadoop, Apache Mesos, Kubernetes, standalone or in the cloud.
  - access data in HDFS, Aluxio, Apache Cassandra, Apache Hbase, Apache Hive, etc.

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

33

# Comparison between Hadoop and Spark

| | hadoop MapReduce | Spark |
|---|---|---|
| **Strengths** | ▪ Can collect any data<br>▪ Limitless in size | ▪ Can work off any Hadoop collection<br>▪ Runs on Hadoop, or other clusters<br>▪ In-memory processing makes it very fast<br>▪ Supports Java, Scala, Python, and R*, and can be used with SQL. |
| **Used for** | ▪ Initial data ingestion<br>▪ Data curation<br>▪ Large-scale "boil the ocean" analytics<br>▪ Data archiving | ▪ Complex query processing of large amounts of data quickly<br>▪ Can handle ad hoc queries |
| **Limitations** | ▪ MapReduce is hard to program<br>▪ Disk-based batch nature limits speed, agility. | ▪ Limited only by processor speed, available memory, cores, and cluster size. |

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

34

# 100TB Daytona Sort Competition

|  | Hadoop MR Record | Spark Record | Spark 1 PB |
|---|---|---|---|
| Data Size | 102.5 TB | 100 TB | 1000 TB |
| Elapsed Time | 72 mins | 23 mins | 234 mins |
| # Nodes | 2100 | 206 | 190 |
| # Cores | 50400 physical | 6592 virtualized | 6080 virtualized |
| Cluster disk throughput | 3150 GB/s (est.) | 618 GB/s | 570 GB/s |
| Sort Benchmark Daytona Rules | Yes | Yes | No |
| Network | dedicated data center, 10Gbps | virtualized (EC2) 10Gbps network | virtualized (EC2) 10Gbps network |
| Sort rate | 1.42 TB/min | 4.27 TB/min | 4.27 TB/min |
| Sort rate/node | 0.67 GB/min | 20.7 GB/min | 22.5 GB/min |

Spark sorted the same data **3X faster** using **10X fewer machines** than Hadoop MR in 2013.

All the sorting took place on disk (HDFS) without using Spark's in-memory cache!

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

35

---



## Startup Crunches 100 Terabytes of Data in a Record 23 Minutes

BY KLINT FINLEY 10.13.14 | 2:36 PM | PERMALINK

Databricks demolishes big data benchmark to prove Spark is fast on disk, too

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

36

# Components of Stack

37

# The Spark stack

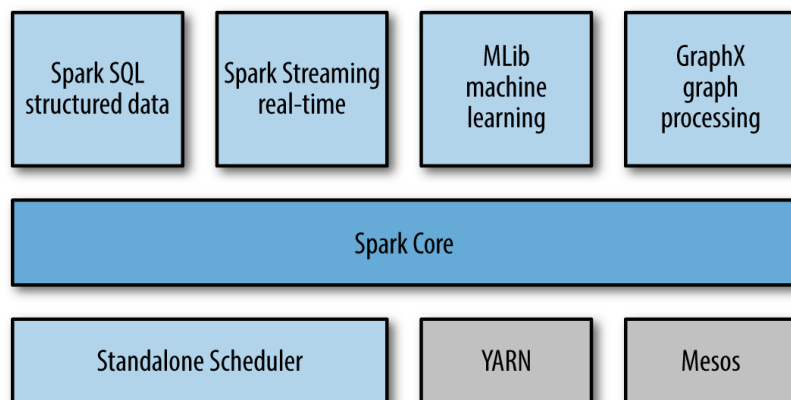| Spark SQL structured data | Spark Streaming real-time | MLib machine learning | GraphX graph processing |
|---|---|---|---|
| Spark Core | | | |
| Standalone Scheduler | | YARN | Mesos |

38

19

# The Spark stack

- Spark Core:
    - contain basic functionality of Spark including task scheduling, memory management, fault recovery, etc.
    - provide APIs for building and manipulating RDDs
- SparkSQL
    - allow querying structured data via SQL, Hive Query Language
    - allow combining SQL queries and data manipulations in Python, Java, Scala

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

39

# The Spark stack

- Spark Streaming: enables processing of live streams of data via APIs
- Mlib:
    - contain common machine language functionality
    - provide multiple types of algorithms: classification, regression, clustering, etc.
- GraphX:
    - library for manipulating graphs and performing graph-parallel computations
    - extend Spark RDD API

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

40

# The Spark stack

- Cluster Managers
  - Hadoop Yarn
  - Apache Mesos, and
  - Standalone Schedular (simple manager in Spark).

# Resilient Distributed Dataset – RDD

- RDD Basics
- Creating RDDs
- RDD Operations
- Common Transformation and Actions
- Persistence (Caching)

# RDD Basics

- RDD:
  - Immutable distributed collection of objects
  - Split into multiple partitions => can be computed on different nodes
- All work in Spark is expressed as
  - creating new RDDs
  - transforming existing RDDs
  - calling actions on RDDs

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

43

# Example

Load error messages from a log into memory, then interactively search for various patterns

```
lines = spark.textFile("hdfs://...")
errors = lines.filter(_.startsWith("ERROR"))
messages = errors.map(_.split('\t')(2))
cachedMsgs = messages.cache()

cachedMsgs.filter(_.contains("foo")).count
cachedMsgs.filter(_.contains("bar")).count
. . .
```



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

44

# RDD Basics

- Two types of operations: transformations and actions
- Transformations: construct a new RDD from a previous one e.g., filter data
- Actions: compute a result base on an RDD e.g., count elements, get first element



45

# Transformations

- Create new RDDs from existing RDDs
- Lazy evaluation
    - See the whole chain of transformations
    - Compute just the data needed
- Persist contents:
    - persist an RDD in memory to reuse it in future
    - persist RDDs on disk is possible

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

46

# Typical works of a Spark program

1. Create some input RDDs form external data
2. Transform them to define new RDDs using transformations like filter()
3. Ask Spark to persist() any intermediate RDDs that will need to be reused
4. Launch actions such as count(), first() to kick off a parallel computation

**SOICT** VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

47

## Resilient Distributed Dataset – RDD

- RDD Basics
- Creating RDDs
- RDD Operations
- Common Transformation and Actions
- Persistence (Caching)

**SOICT** VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

48

# Two ways to create RDDs

1. Parallelizing a collection: uses parallelize()

• Python

```
lines = sc.parallelize(["pandas", "i like pandas"])
```

• Scala

```
val lines = sc.parallelize(List("pandas", "i like
pandas"))
```

• Java

```
JavaRDD<String> lines =
sc.parallelize(Arrays.asList("pandas", "i like
pandas"));
```

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

49

# Two ways to create RDDs

2. Loading data from external storage

• Python

```
lines = sc.textFile("/path/to/README.md")
```

• Scala

```
val lines = sc.textFile("/path/to/README.md")
```

• Java

```
JavaRDD<String> lines =
sc.textFile("/path/to/README.md");
```

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

50

# Resilient Distributed Dataset – RDD

- RDD Basics
- Creating RDDs
- RDD Operations
- Common Transformation and Actions
- Persistence (Caching)

51

# RDD Operations

- Two types of operations
    - Transformations: operations that return a new RDDs e.g., map(), filter()
    - Actions: operations that return a result to the driver program or write it to storage such as count(), first()
- Treated differently by Spark
    - Transformation: lazy evaluation
    - Action: execution at any time

52

# Transformation

- Example 1. Use filter()

- Python
```
inputRDD = sc.textFile("log.txt")
errorsRDD = inputRDD.filter(lambda x: "error" in x)
```

- Scala
```
val inputRDD = sc.textFile("log.txt")
val errorsRDD = inputRDD.filter(line =>
line.contains("error"))
```

- Java
```
JavaRDD<String> inputRDD = sc.textFile("log.txt");
JavaRDD<String> errorsRDD = inputRDD.filter(
new Function<String, Boolean>() {
public Boolean call(String x) {
    return x.contains("error"); }}
  });
```

53

# Transformation

- filter()
  - does not change the existing *inputRDD*
  - returns a pointer to an entirely new RDD
  - *inputRDD* still can be reused

- union()
```
errorsRDD = inputRDD.filter(lambda x: "error" in x)
warningsRDD=inputRDD.filter(lambda x: "warning" in x)
badLinesRDD = errorsRDD.union(warningsRDD)
```
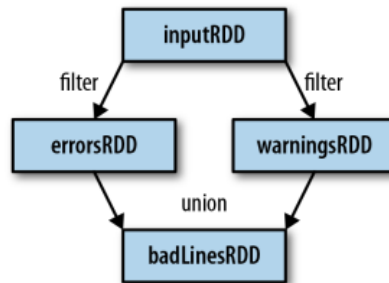
- transformations can operate on any number of input RDDs

54

# Transformation

- Spark keeps track dependencies between RDDs, called the lineage graph
- Allow recovering lost data

55

# Actions

- Example. count the number of errors
- Python

```python
print "Input had " + badLinesRDD.count() + " concerning lines"
print "Here are 10 examples:"
for line in badLinesRDD.take(10):
print line
```

- Scala

```scala
println("Input had " + badLinesRDD.count() + " concerning lines")
println("Here are 10 examples:")
badLinesRDD.take(10).foreach(println)
```

- Java

```java
System.out.println("Input had " + badLinesRDD.count() + " concerning
lines")
System.out.println("Here are 10 examples:")
for (String line: badLinesRDD.take(10)) {
System.out.println(line);
}
```

56

# Resilient Distributed Dataset – RDD

- RDD Basics
- Creating RDDs
- RDD Operations
- Common Transformation and Actions
- Persistence (Caching)

SOICT  VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

57

# RDD Basics

| Transformations | Actions |
|---|---|
| map<br>flatMap<br>filter<br>sample<br>union<br>groupByKey<br>reduceByKey<br>join<br>cache | reduce<br>collect<br>count<br>save<br>lookupKey<br>… |

SOICT  VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

58

# Transformations

| transformation | description |
|---|---|
| **map(**_func_**)** | return a new distributed dataset formed by passing each element of the source through a function _func_ |
| **filter(**_func_**)** | return a new dataset formed by selecting those elements of the source on which _func_ returns true |
| **flatMap(**_func_**)** | similar to map, but each input item can be mapped to 0 or more output items (so _func_ should return a Seq rather than a single item) |
| **sample(**_withReplacement_, _fraction_, _seed_**)** | sample a fraction _fraction_ of the data, with or without replacement, using a given random number generator _seed_ |
| **union(**_otherDataset_**)** | return a new dataset that contains the union of the elements in the source dataset and the argument |
| **distinct([**_numTasks_**]))** | return a new dataset that contains the distinct elements of the source dataset |

59

# Transformations

| transformation | description |
|---|---|
| **groupByKey([**_numTasks_**])** | when called on a dataset of (K, V) pairs, returns a dataset of (K, Seq[V]) pairs |
| **reduceByKey(**_func_, [_numTasks_]**)** | when called on a dataset of (K, V) pairs, returns a dataset of (K, V) pairs where the values for each key are aggregated using the given reduce function |
| **sortByKey([**_ascending_], [_numTasks_]**)** | when called on a dataset of (K, V) pairs where K implements Ordered, returns a dataset of (K, V) pairs sorted by keys in ascending or descending order, as specified in the boolean ascending 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 |
| **cogroup(**_otherDataset_, [_numTasks_]**)** | when called on datasets of type (K, V) and (K, W), returns a dataset of (K, Seq[V], Seq[W]) tuples – also called groupWith |
| **cartesian(**_otherDataset_**)** | when called on datasets of types T and U, returns a dataset of (T, U) pairs (all pairs of elements) |

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

60

30

3/26/2021

# Actions

| action | description |
|---|---|
| **reduce(**_func_**)** | aggregate the elements of the dataset using a function _func_ (which takes two arguments and returns one), and should also be commutative and associative so that it can be computed correctly in parallel |
| **collect()** | return all the elements of the dataset as an array at the driver program – usually useful after a filter or other operation that returns a sufficiently small subset of the data |
| **count()** | return the number of elements in the dataset |
| **first()** | return the first element of the dataset – similar to _take(1)_ |
| **take(**_n_**)** | return an array with the first _n_ elements of the dataset – currently not executed in parallel, instead the driver program computes all the elements |
| **takeSample(**_withReplacement, fraction, seed_**)** | return an array with a random sample of _num_ elements of the dataset, with or without replacement, using the given random number generator seed |

SOICT  VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

61

# Actions

| action | description |
|---|---|
| **saveAsTextFile(**_path_**)** | write the elements of the dataset as a text file (or set of text files) in a given directory in the local filesystem, HDFS or any other Hadoop-supported file system. Spark will call toString on each element to convert it to a line of text in the file |
| **saveAsSequenceFile(**_path_**)** | write the elements of the dataset as a Hadoop SequenceFile in a given path in the local filesystem, HDFS or any other Hadoop-supported file system. Only available on RDDs of key-value pairs that either implement Hadoop's Writable interface or are implicitly convertible to Writable (Spark includes conversions for basic types like Int, Double, String, etc). |
| **countByKey()** | only available on RDDs of type (K, V). Returns a `Map` of (K, Int) pairs with the count of each key |
| **foreach(**_func_**)** | run a function _func_ on each element of the dataset – usually done for side effects such as updating an accumulator variable or interacting with external storage systems |

62

31

# Resilient Distributed Dataset – RDD

- RDD Basics
- Creating RDDs
- RDD Operations
- Common Transformation and Actions
- Persistence (Caching)

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

63

# Persistence levels

| Level | Space used | CPU time | In memory | On disk | Comments |
|-------|-----------|----------|-----------|---------|----------|
| MEMORY_ONLY | High | Low | Y | N | |
| MEMORY_ONLY_SER | Low | High | Y | N | |
| MEMORY_AND_DISK | High | Medium | Some | Some | Spills to disk if there is too much data to fit in memory. |
| MEMORY_AND_DISK_SER | Low | High | Some | Some | Spills to disk if there is too much data to fit in memory. Stores serialized representation in memory. |
| DISK_ONLY | Low | High | N | Y | |

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

64

# Persistence

• Example

```
val result = input.map(x => x * x)
result.persist(StorageLevel.DISK_ONLY)
println(result.count())
println(result.collect().mkString(","))
```

SOICT   VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

65

---

**Books:**

• Holden Karau, Andy Konwinski, Patrick Wendell & Matei Zaharia. Learning Spark. Oreilly
• TutorialsPoint. Spark Core Programming

Acknowledgement and References

**Slides:**

• Paco Nathan. Intro to Apache Spark
• Harold Liu. Berkely Data Analytics Stack
• DataBricks. Intro to Spark Development

SOICT   VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

66

# Q&A

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

67



68

Thank you for your attentions!

soict.hust.edu.vn/  fb.com/groups/soict

69



ĐẠI HỌC BÁCH KHOA HÀ NỘI
VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

70