

## IT4043E Lưu trữ và phân tích dữ liệu lớn

IT4043E

12/2022  
Thanh-Chung Dao Ph.D.

1

## Lecture Agenda

- W1: Spark introduction + Lab
- W2: Spark RDD + Lab
- W3: Spark Machine Learning + Lab
- W4: Spark on Blockchain Storage + Lab

2

2

1

## **Today's Agenda**

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

3

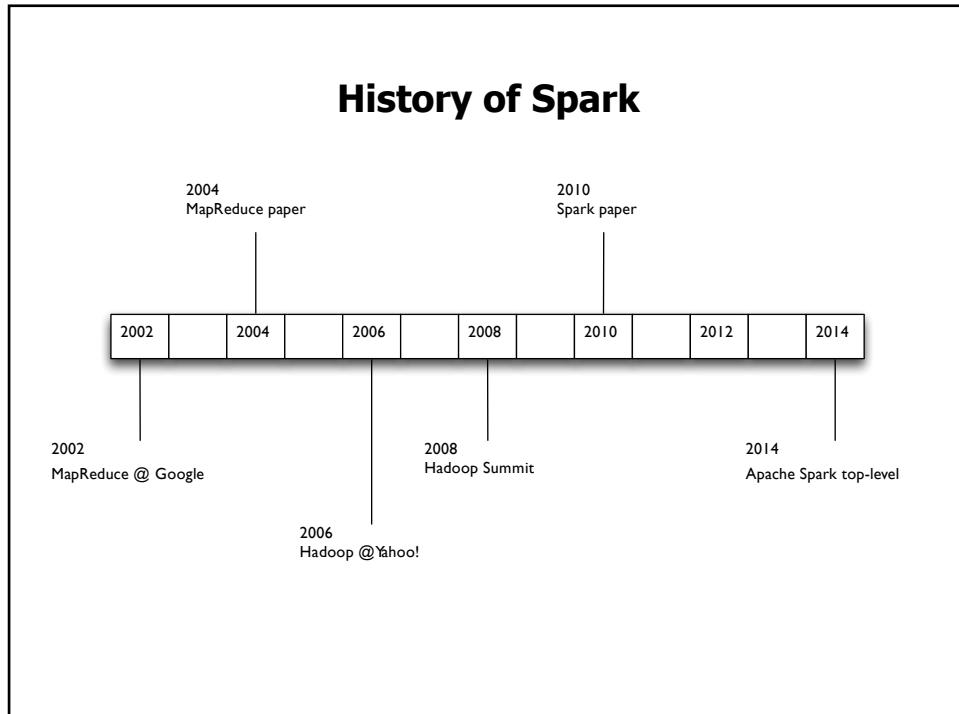
3

## HISTORY OF SPARK

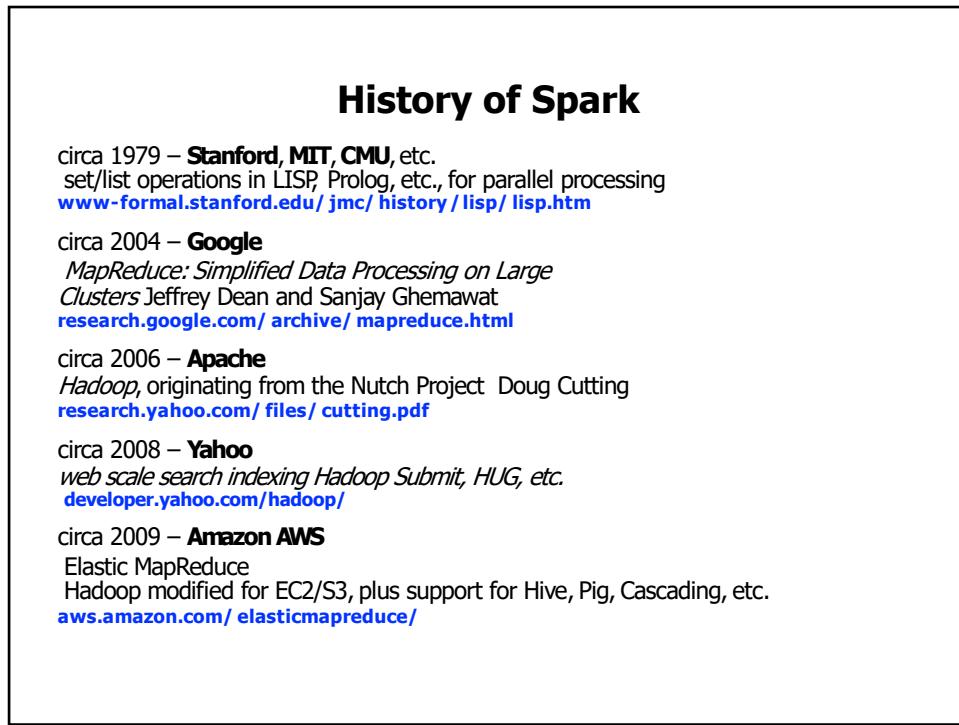
4

4

2



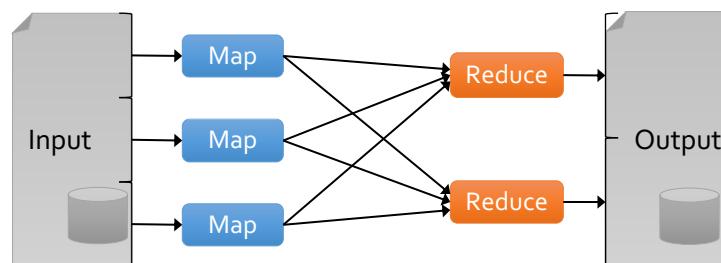
5



6

## MapReduce

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



7

## 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)

8

## Data Processing Goals

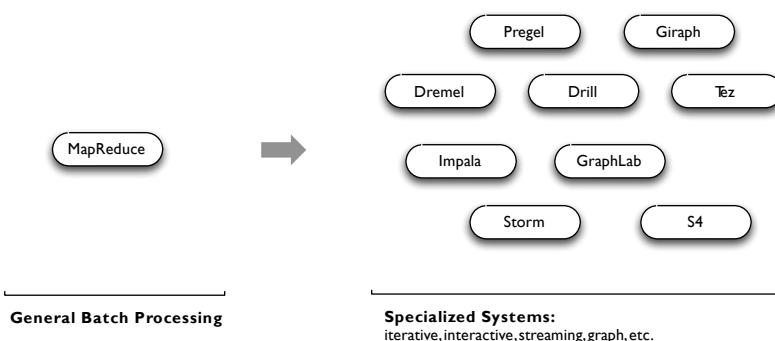


- **Low latency (interactive) queries on historical data:** enable faster decisions
  - E.g., identify why a site is slow and fix it
- **Low latency queries on live data (streaming):** enable decisions on real-time data
  - E.g., detect & block worms in real-time (a worm may infect **1mil** hosts in **1.3sec**)
- **Sophisticated data processing:** enable “better” decisions
  - E.g., anomaly detection, trend analysis

Therefore, people built specialized systems as workarounds...

9

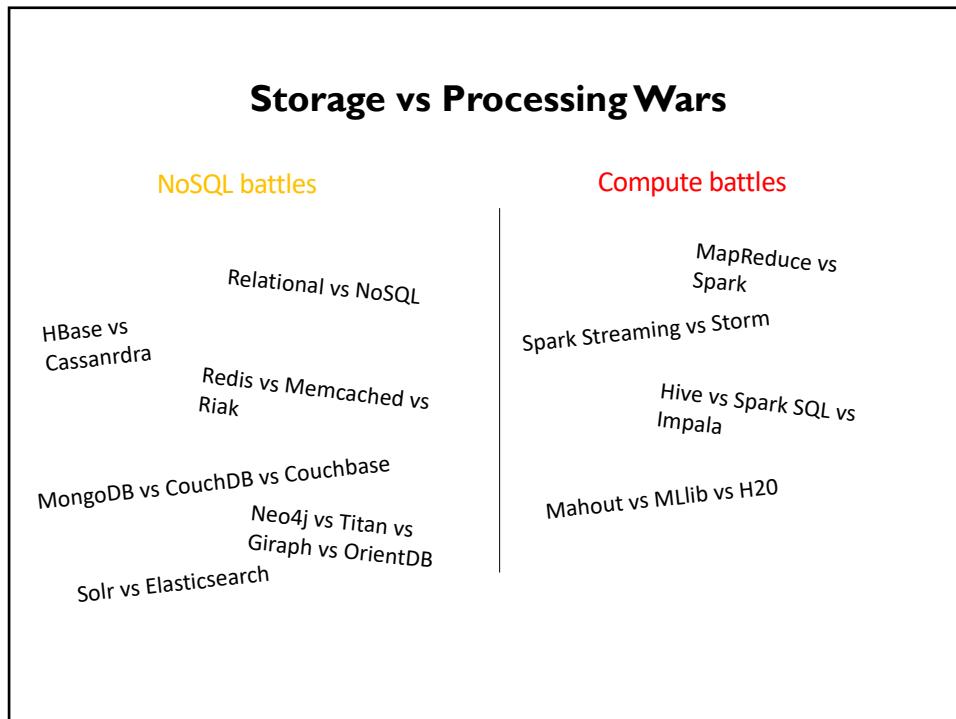
## Specialized Systems



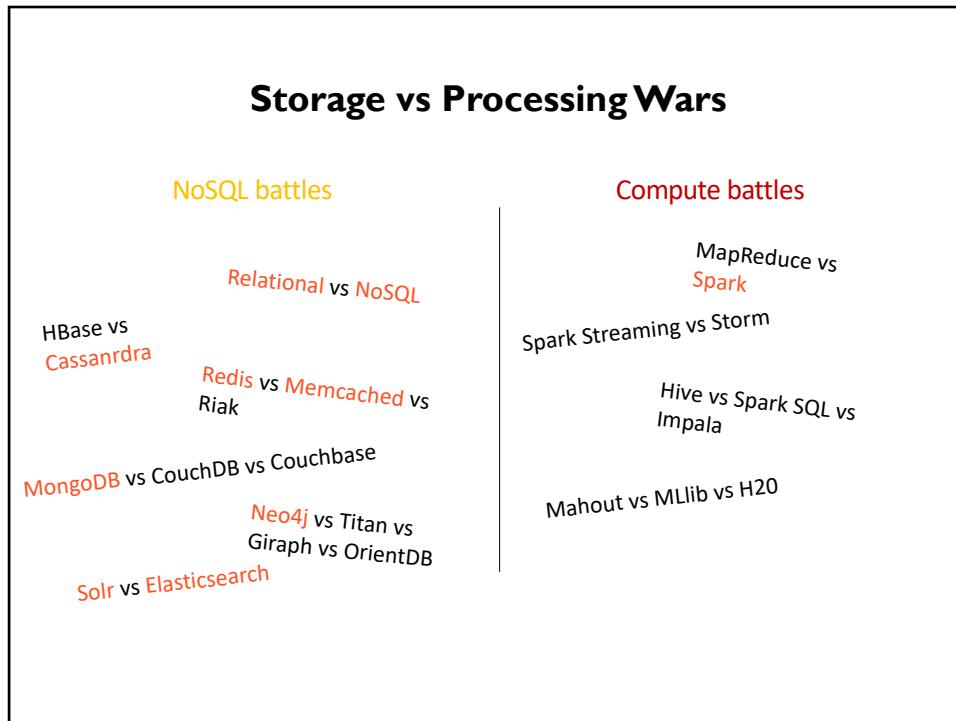
**Specialized Systems:**  
iterative, interactive, streaming, graph, etc.

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

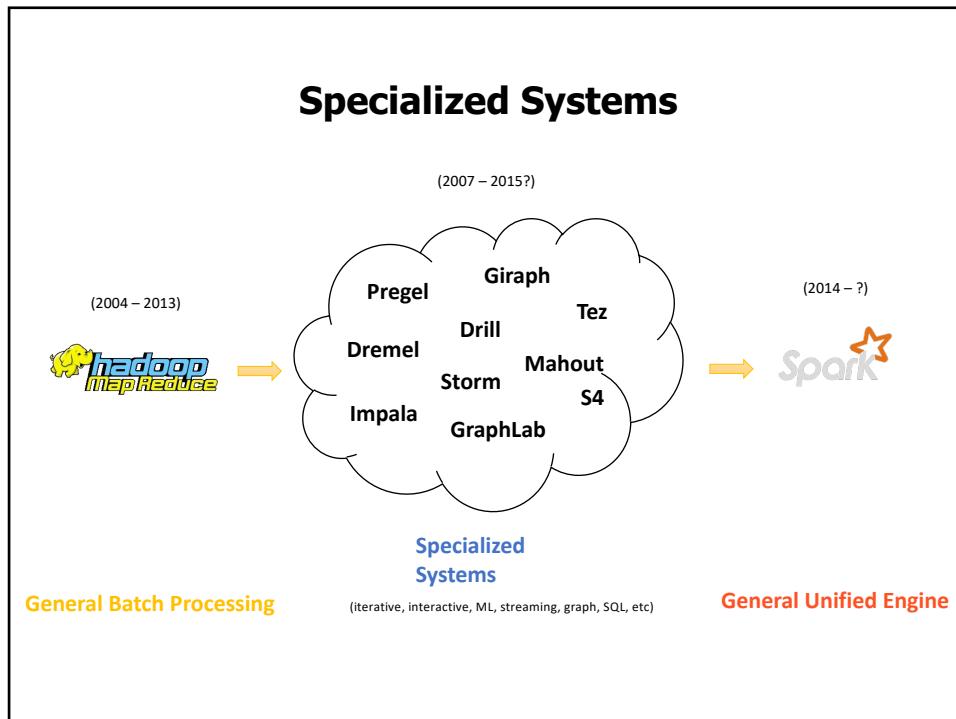
10



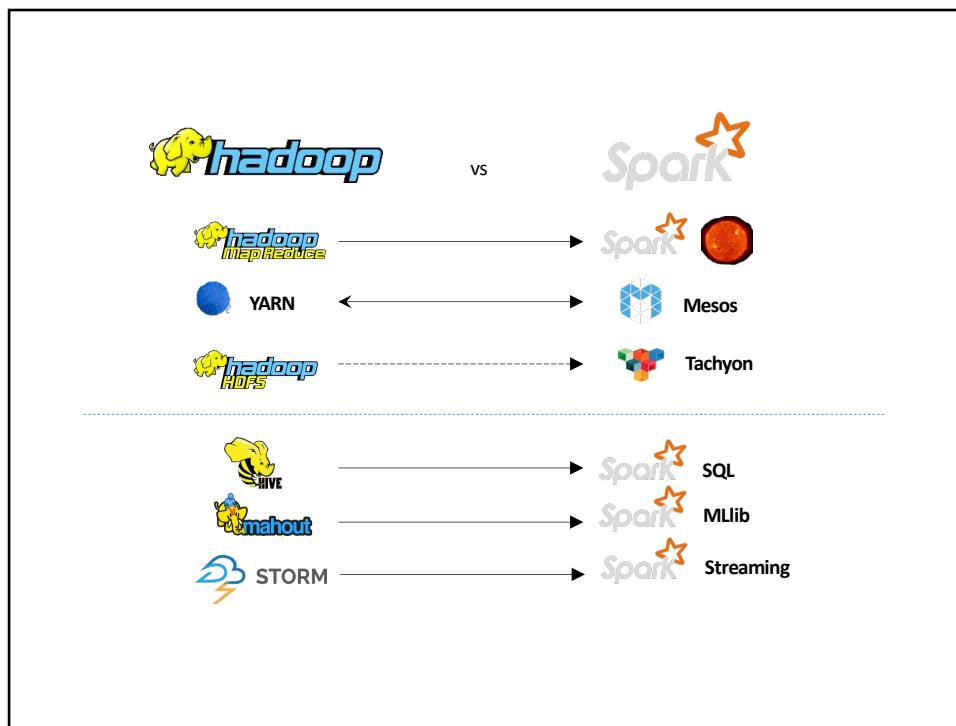
11



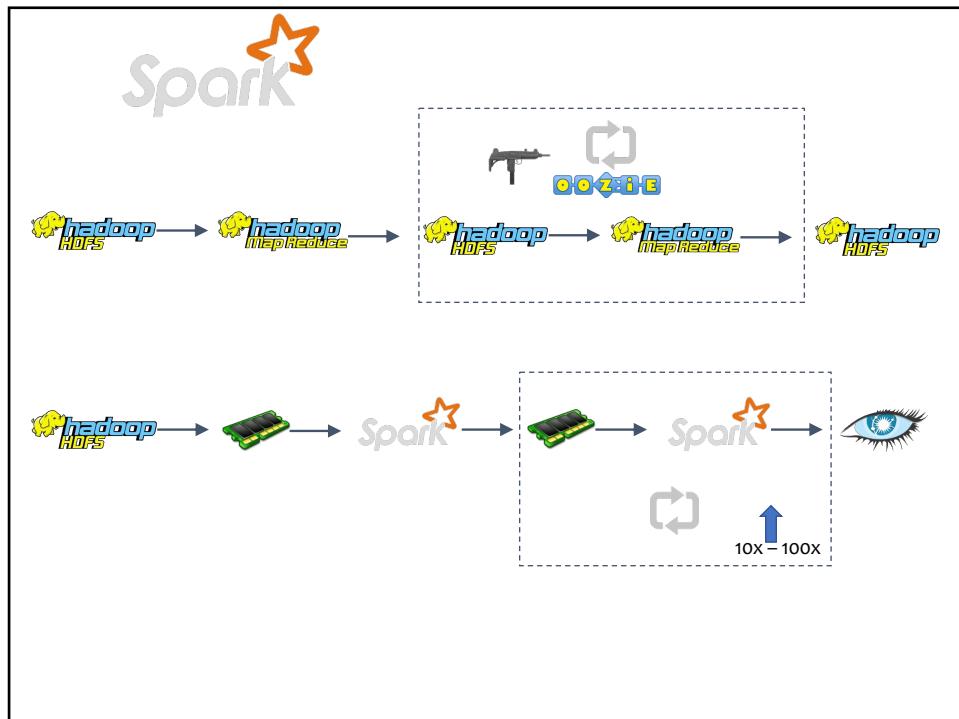
12



13



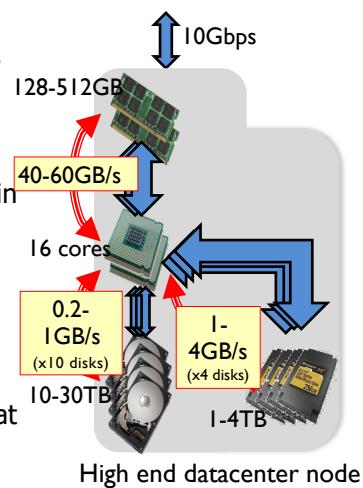
14



15

## Support Interactive and Streaming Comp.

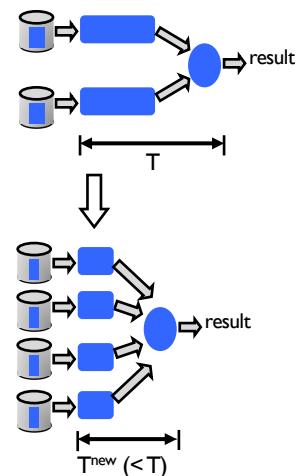
- Aggressive use of **memory**
- Why?
  1. Memory transfer rates  $\gg$  disk or SSDs
  2. Many datasets already fit into memory
    - Inputs of over 90% of jobs in Facebook, Yahoo!, and Bing clusters fit into memory
    - e.g., 1TB = 1 billion records @ 1KB each
  3. Memory density (still) grows with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



16

## 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



17

## Berkeley AMPLab

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



18

## Berkeley AMPLab

- Funding:
  - XData, DARPA, CISE Expedition Grant
  - NSF
- Industrial, founding sponsors
- 18 other sponsors, including

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

19

## Databricks

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

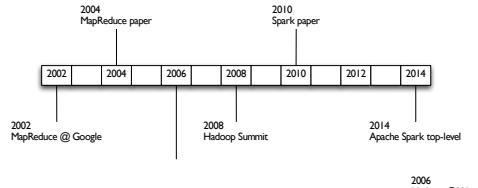
**Databricks Cloud:**  
 "A unified platform for building Big Data pipelines – from ETL to Exploration and Dashboards, to Advanced Analytics and Data Products."

20

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



# History of Spark



Spark: Cluster Computing with Working Sets  
Matei Zaharia, Mosharaf Chowdhury

Pratek Zanatta, Prashant Choudhary,  
Michael J. Franklin, Scott Shenker, Ion Stoica  
USENIX HotCloud (2010)

[http://CSAIL.mit.edu/~matei/papers/2010/hotcloud\\_spark.pdf](http://CSAIL.mit.edu/~matei/papers/2010/hotcloud_spark.pdf)

# Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing

*Archiver 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](http://usenix.org/system/files/conference/nsdi12/nsdi12-final138.pdf)

# History of Spark

**Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing**

Matei Zaharia, Michael Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, Ion Stoica  
*University of California, Berkeley*

**Abstract**  
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. Both cases benefit from keeping data in memory and performing computation in parallel. In order to implement these applications efficiently, RDDs provide a common abstraction for distributed memory and fault-tolerant transformations rather than fine-grained update semantics. This abstraction is general enough to capture a wide class of computations, including iterative algorithms and interactive data mining tools, such as PageRank and new applications that these models enable. We implemented RDDs in a distributed system called Spark, which we evaluate through a variety of benchmarks.

**1 Introduction**  
Cluster computing frameworks like MapReduce [10] and Dryad [11] have been widely adopted for large-scale data processing. They are well suited for batch processing of static data using a set of high-level operators, without having to worry about data placement or fault tolerance.

Although current frameworks provide memory abstractions for distributed memory, they lack abstractions for leveraging distributed memory in a fault-tolerant manner. This is a common class of emerging applications that fit into iterative algorithms, such as machine learning and graph algorithms, and interactive data mining. Another compelling use case is interactive data mining, where a user runs multiple analyses on the same dataset. In current systems, the only way to reuse data is to write it to an external stable storage system (e.g., HDFS) and read it back in for each successive job. This leads to data replication, disk I/O, and serialization overheads, which can dominate application execution times.

Recognizing this problem, researchers have developed specialized frameworks for these applications. For example, GraphLab [21] is designed for iterative graph computation that keeps intermediate data in memory. Pregel [22] provides a distributed memory abstraction interface. However, these frameworks only support a subset of the operations available in MapReduce and MapReduce steps, and perform data sharing implicitly rather than explicitly. They also do not support the more general issue, e.g., to be able to use a user-defined function to read several datasets in parallel.

In this paper, we propose a new abstraction called *resilient distributed datasets* (RDDs) that enables efficient fault-tolerant, parallel data structures that let users explicitly define their lineage and dependencies, and support their partitioning to optimize data placement, and manage their lifetime.

The main challenge in designing RDDs is defining a programming model that is expressive enough to support efficiently. Existing abstractions for memory storage (such as shared memory [23], shared nothing [24], shared-value stores [25], databases, and Piccolo [26]) offer an interesting starting point, but they are not designed for fault tolerance. Instead, RDDs build on the well-known pattern of log-structured merge trees [27] to support fault tolerance by logging the transformations used to build a RDD. RDDs also support lineage tracing, which provides information about how a value was derived from an RDD. This requires copy-on-change propagation, however, and we discuss how to do it in §5.4.

**"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."**

April 2012

23

# 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)`

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

24

# History of Spark



Analyze real time streams of data in ½ second intervals

[www.cs.berkeley.edu/~matei/papers/2013/soosp\\_spark\\_streaming.pdf](http://www.cs.berkeley.edu/~matei/papers/2013/soosp_spark_streaming.pdf)

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

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

**Abstract**

Many “big data” applications must act on data in real time. Running these applications as a regular job requires parallel platforms that automatically handle faults and support fault recovery. Existing fault-tolerant processing models provide fault recovery in an expensive manner, requiring hot replication or long recovery times. In this paper, we introduce a new fault-tolerant processing model, *discretized streams* (DStreams), that overcomes these challenges. DStreams enable a ported of the MapReduce API to streaming, supports traditional replication and backup schemes, and tolerates stragglers. We show that they support a rich set of operators for stream processing, including joins, and scale from single-node systems, linear scaling to 100 nodes, without losing the guarantees of bounded fault recovery. Finally, DStreams can easily be converted into interactive query models like MapReduce, enabling rich applications to reuse these modes. We implement DStreams in a system called Spark Streaming.

**1 Introduction**

Much of “big data” applications must act on data in real time. For example, a social network may want to detect trending conversation topics in

**TwitterUtils.createStream(...)**  
**.filter(\_.getText.contains("Spark"))**  
**.countByWindow(Seconds(5))**

25

# History of Spark



Seemlessly mix SQL queries with Spark programs.

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

Michael Armbrust<sup>1</sup>, Reynold S. Xin<sup>1</sup>, Cheng Liang<sup>1</sup>, Yih Huai<sup>1</sup>, Davies Liu<sup>1</sup>, Joseph K. Bradley<sup>1</sup>, Kangru Meng<sup>1</sup>, Tomer Kaftan<sup>1</sup>, Michael J. Franklin<sup>2</sup>, Ali Ghodsi<sup>1</sup>, Matei Zaharia<sup>1</sup>  
<sup>1</sup>Databricks Inc.    <sup>2</sup>MIT CSAIL    <sup>3</sup>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 top of the distributed storage and computation infrastructure of Spark, Spark SQL provides a simple interface for users leveraging the benefits of relational processing (e.g., declarative grouping, joins, and window operations) while maintaining the analytics libraries in Spark (e.g., machine learning). Compared to previous systems, Spark SQL makes two main additions. First, it offers a declarative API for defining data sources and performing processing, through a declarative DataFrame API that integrates MapReduce, HDFS, and distributed databases. Second, it includes a rule-based optimizer, Catalyst, built using features of the Scala programming language, that performs both relational and procedural optimizations, generate code, and define execution plans. Using Catalyst, we have built a system that can process data in various formats, including learning types, and query distributed or external data sources tailored for the needs of big data analysis. We see Spark SQL as an evolution of both SQL-on-Hadoop and of SQL-like engines, richer APIs and optimizations while keeping the benefits of the functional API.

**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, due to the variety of data formats. The earlier systems designed for these workloads, such as MapReduce, gave users a powerful, but

While the declarative systems show that users often prefer writing declarative queries, the relational approach is inefficient for many big data queries. First, users often perform complex joins that are hard to express in a declarative, structured, requiring custom code. Second, users want to perform complex operations that are hard to express in a declarative language that are challenging to express in relational systems. In practice, users often want to combine multiple operations that are expressed with a combination of both relational queries and complex procedural logic. In addition, users often want to mix both relational and procedural—here we now renamed largely disjoint, forcing users to choose one paradigm or the other.

The DataFrames API is built on top of the RDD API in Spark SQL, a major new component in Apache Spark [19]. Spark SQL is built on top of the Catalyst optimizer [20] and the DataFrame API. While forcing users to pick between a relational or a procedural API, the DataFrame API offers a rich relational/procedural integration that allows users to write complex queries that manipulate 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)
```

26

# History of Spark



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, jgonzal, franklin, istoica}@cs.berkeley.edu

**ABSTRACT**  
 From social networks to targeted advertising, big graphs capture the complex web of data and are central to recent advances in machine learning and data mining. Unfortunately, directly applying existing large-scale distributed systems to graphs is inefficient and inelegant. The need for intuitive, scalable tools for graph computation has led to the development of new graph-parallel systems that can efficiently execute graph algorithms. Unfortunately, these new graph-parallel systems are often just as problematic as the subsequent components they replace. They are typically built on top of distributed data-parallel and support for interactive data mining is limited.

We present GraphX, a distributed graph system built on top of both data-parallel and graph-parallel systems by efficiently expressing graph computation within the Spark data-parallel framework. We leverage the strengths of both systems to support distributed graphs as tabular data-structures. Similarly, we leverage advances in distributed data-parallel systems to support distributed and fault-tolerance. We provide powerful new operations to simplify graph computation, and support for distributed graph processing. We implement the PowerGraph and Pregel abstractions in less than 20 lines of code, and by expressing the local formulation of Pregel, we enable users to interactively test, tune, and compare their 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 systems design, GraphX greatly simplifies the development of graph processing, and opens up a range of new sophisticated graph algorithms to large-scale real-world graph problems.

Most distributed systems for graph processing share many common properties; each presents a slightly different view of graph computation, and each is designed to support a different set of graph algorithms and applications. Unfortunately, because each framework is built from scratch, it is difficult to reuse or compose these abstractions. Furthermore, while these frameworks address the challenges of graph computation, they do not address the challenges of distributed graph processing, such as the process of interpreting and applying the results of computation. Finally, few frameworks support distributed graph processing and fault-tolerance.

Alternatively, data-parallel systems like MapReduce and Apache Hadoop have been shown to be well-suited to distributed graph computation [1]. We introduce GraphX, a distributed graph system built on top of the Spark data-parallel framework. GraphX provides a distributed graph abstraction, and a distributed abstraction to the task of graph computation (ETL). By exploiting distributed data-parallel primitives, GraphX can support distributed graphs as tabular data-structures. More recent systems like GraphX and PowerGraph have shown promise [2, 3]. However, merely expressing graph computation and graph patterns in these data-parallel frameworks can be challenging and typically leads to complex and error-prone code. In addition, these systems do not support the graph structure.

To address these challenges we introduce GraphX, a graph computation system which runs in the Spark data-parallel framework. GraphX provides a distributed graph abstraction, and a distributed abstraction to the task of graph computation (ETL). By exploiting distributed data-parallel primitives, GraphX can support distributed graphs as tabular data-structures. More recent systems like GraphX and PowerGraph have shown promise [2, 3]. However, merely expressing graph computation and graph patterns in these data-parallel frameworks can be challenging and typically leads to complex and error-prone code. In addition, these systems do not support the graph structure.

The main challenge we introduce GraphX, a graph computation system which runs in the Spark data-parallel framework. GraphX provides a distributed graph abstraction, and a distributed abstraction to the task of graph computation (ETL). By exploiting distributed data-parallel primitives, GraphX can support distributed graphs as tabular data-structures. More recent systems like GraphX and PowerGraph have shown promise [2, 3]. However, merely expressing graph computation and graph patterns in these data-parallel frameworks can be challenging and typically leads to complex and error-prone code. In addition, these systems do not support the graph structure.

graph = Graph(vertices, edges)  
 messages =  
 spark.textFile("hdfs://...")  
 graph2 =  
 graph.joinVertices(messages) {  
 (id, vertex, msg) => ...  
}

[https://AMPLab.berkeley.edu/wp-content/uploads/2013/05/grades-graphx\\_with\\_fonts.pdf](https://AMPLab.berkeley.edu/wp-content/uploads/2013/05/grades-graphx_with_fonts.pdf)

27

# 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  
 {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 returning results with bounded errors and bounded response times. To achieve this, BlinkDB uses two key ideas: (1) an adaptive optimization framework that builds and maintains a set of multi-stratified and multi-partitioned samples that reflects an appropriately sized sample based on a query's accuracy or response time requirements, and (2) a hybrid execution plan that selects an appropriate query plan based on the sample's size and quality. Our experiments on 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 200x faster than Hive), within an error of  $\pm 0\%$ .

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

processing of large amounts of data by trading off accuracy for response time. These techniques include sampling [10, 14], sketches [1], and online aggregation [15]. To illustrate the utility of such techniques, consider the following simple query that computes the average SessionTime over all sessions 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 slow. One way to speed up the query is to read disk and to stream through all the tuples. Suppose we instead executed the same query on a sample containing only 10,000 tuples. In this case, each tuple is read sequentially, one at a time. This would be orders of magnitude faster, while still providing an approximate result within a few percent of the actual value. Using this approach, we can 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

**Queries with Time Bounds**

```
SELECT avg(sessionTime)
FROM Table
WHERE city='San Francisco'
WITHIN 2 SECONDS
```

**Queries with Error Bounds**

```
SELECT avg(sessionTime)
FROM Table
WHERE city='San Francisco'
ERROR 0.1 CONFIDENCE 95.0%
```

[https://www.cs.berkeley.edu/~sameerag/blinkdb\\_eurosys13.pdf](https://www.cs.berkeley.edu/~sameerag/blinkdb_eurosys13.pdf)

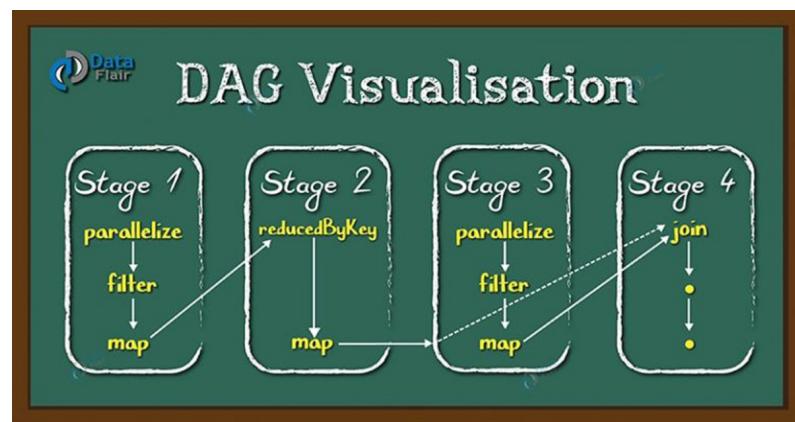
28

## 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

29

## Directed Acyclic Graph - DAG



30

30

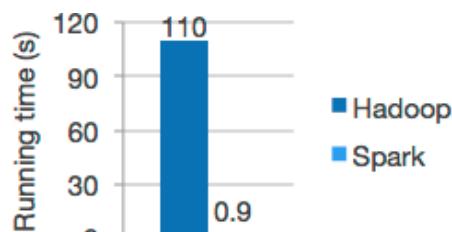
# SPARK INTRODUCTION

31

31

## 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

32

32

## 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 from Scala, Python, R, and SQL

```
df = spark.read.json("logs.json") df.  
where("age > 21")  
select("name.first").show()
```

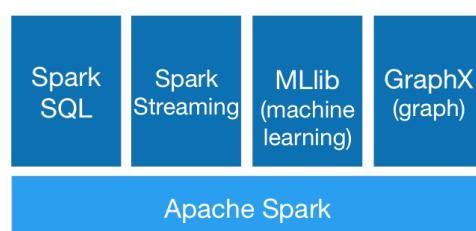
Spark's Python DataFrame API  
Read JSON files with automatic schema inference

33

33

## What is Apache Spark

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



34

34

## 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.



35

35

## Comparison between Hadoop and Spark

	Hadoop MapReduce	Spark
Strengths	<ul style="list-style-type: none"> <li>Can collect any data</li> <li>Limitless in size</li> </ul>	<ul style="list-style-type: none"> <li>Can work off any Hadoop collection</li> <li>Runs on Hadoop, or other clusters</li> <li>In-memory processing makes it very fast</li> <li>Supports Java, Scala, Python, and R*, and can be used with SQL.</li> </ul>
Used for	<ul style="list-style-type: none"> <li>Initial data ingestion</li> <li>Data curation</li> <li>Large-scale “boil the ocean” analytics</li> <li>Data archiving</li> </ul>	<ul style="list-style-type: none"> <li>Complex query processing of large amounts of data quickly</li> <li>Can handle ad hoc queries</li> </ul>
Limitations	<ul style="list-style-type: none"> <li>MapReduce is hard to program</li> <li>Disk-based batch nature limits speed, agility.</li> </ul>	<ul style="list-style-type: none"> <li>Limited only by processor speed, available memory, cores, and cluster size.</li> </ul>

36

36

## 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!

37

**WIRED** GEAR SCIENCE ENTERTAINMENT BUSINESS SECURITY DESIGN OPINION MAGAZINE

ENTERPRISE | big data | databricks | google | Hadoop

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

BY KLTN FINLEY 10.13.14 | 2:36 PM | PERMALINK

Facebook Share (1.1k) Twitter Tweet (789) Google+ (+1) LinkedIn Share (75) Pinterest Pin It (565)

**GIGAOM** EVENTS RESEARCH SIGN IN SUBSCRIBE

Cloud Data Media Mobile Science & Energy Social & Web Podcasts

Gigaom Research. Get unlimited market intelligence from over 200 i

**MUST READS**

- Google launches Contributor, a crowdfunding tool for publishers
- Net neutrality looks doomed in Europe before it even gets started
- Pure tech products that designers have fallen in love with

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

by Derrick Harris Oct 10, 2014 - 1:49 PM PST

Comment

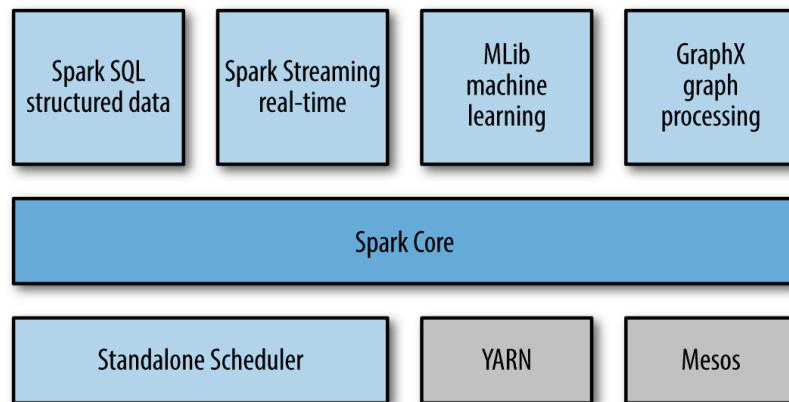
38

# Components of Stack

39

39

## The Spark stack



40

40

## 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

41

41

## 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

42

42

## The Spark stack

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

43

43

## Resilient Distributed Dataset – RDD

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

44

44

## 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

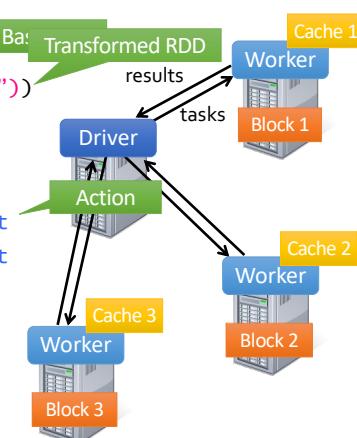
45

45

## Example

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

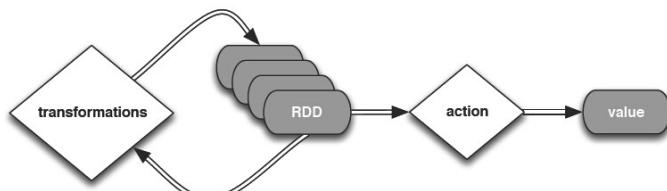
```
lines = spark.textFile("hdfs://...")  
errors = lines.filter(_.startswith("ERROR"))  
messages = errors.map(_.split('\t')(2))  
cachedMsgs = messages.cache()  
  
cachedMsgs.filter(_.contains("foo")).count  
cachedMsgs.filter(_.contains("bar")).count  
...
```



46

## 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



47

47

## 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

48

48

## Typical works of a Spark program

1. Create some input RDDs from 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

49

49

## Resilient Distributed Dataset – RDD

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

50

50

## 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(NSArray.asList("pandas", "i
like pandas"));
```

51

51

## 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");
```

52

52

## Resilient Distributed Dataset – RDD

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

53

53

## 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

54

54

## 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");
        }
    });

```

55

55

## 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)
```

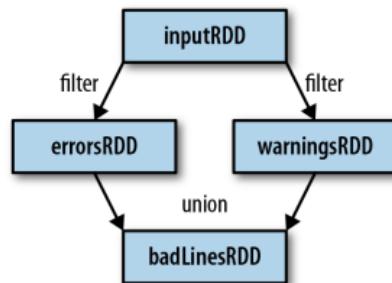
- transformations can operate on any number of input RDDs

56

56

## Transformation

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



57

57

## Actions

- Example. count the number of errors

- Python

```

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

- Scala

```

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

- 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);
}
  
```

58

58

<h2 style="color: #4682B4; font-weight: bold;">Resilient Distributed Dataset – RDD</h2>	<ul style="list-style-type: none"> <li>• RDD Basics</li> <li>• Creating RDDs</li> <li>• RDD Operations</li> <li>• Common Transformation and Actions</li> <li>• Persistence (Caching)</li> </ul>
59	

59

RDD Basics	
<b>Transformations</b>	<b>Actions</b>
map flatMap filter sample union groupByKey reduceByKey join cache ...	reduce collect count save lookupKey ...

Email: info@pti.edu.vn | Website: pti.edu.vn

60

## Transformations

<i>transformation</i>	<i>description</i>
<b>map(func)</b>	return a new distributed dataset formed by passing each element of the source through a function <i>func</i>
<b>filter(func)</b>	return a new dataset formed by selecting those elements of the source on which <i>func</i> returns true
<b>flatMap(func)</b>	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)
<b>sample(withReplacement, fraction, seed)</b>	sample a fraction <i>fraction</i> of the data, with or without replacement, using a given random number generator <i>seed</i>
<b>union(otherDataset)</b>	return a new dataset that contains the union of the elements in the source dataset and the argument
<b>distinct([numTasks]))</b>	return a new dataset that contains the distinct elements of the source dataset

61

## Transformations

<i>transformation</i>	<i>description</i>
<b>groupByKey([numTasks])</b>	when called on a dataset of (K, V) pairs, returns a dataset of (K, Seq[V]) pairs
<b>reduceByKey(func, [numTasks])</b>	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
<b>sortByKey([ascending], [numTasks])</b>	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
<b>join(otherDataset, [numTasks])</b>	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
<b>cogroup(otherDataset, [numTasks])</b>	when called on datasets of type (K, V) and (K, W), returns a dataset of (K, Seq[V], Seq[W]) tuples – also called groupWith
<b>cartesian(otherDataset)</b>	when called on datasets of types T and U, returns a dataset of (T, U) pairs (all pairs of elements)

62

## Actions

action	description
<code>reduce(func)</code>	aggregate the elements of the dataset using a function <code>func</code> (which takes two arguments and returns one), and should also 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 – 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 <code>n</code> elements of the dataset – currently not executed in parallel, instead the driver program computes all the elements
<code>takeSample(withReplacement, fraction, seed)</code>	return an array with a random sample of <code>num</code> elements of the dataset, with or without replacement, using the given random number generator seed

63

## Actions

action	description
<code>saveAsTextFile(path)</code>	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 <code>toString</code> on each element to convert it to a line of text in the file
<code>saveAsSequenceFile(path)</code>	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).
<code>countByKey()</code>	only available on RDDs of type (K, V). Returns a 'Map' of (K, Int) pairs with the count of each key
<code>foreach(func)</code>	run a function <code>func</code> on each element of the dataset – usually done for side effects such as updating an accumulator variable or interacting with external storage systems

64

## Resilient Distributed Dataset – RDD

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

65

65

## 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	

66

66

## Persistence

- Example

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

67

67

## Acknowledgement and References

### Books:

- Holden Karau, Andy Konwinski, Patrick Wendell & Matei Zaharia. Learning Spark. O'reilly
- TutorialsPoint. Spark Core Programming

### Slides:

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

68

68