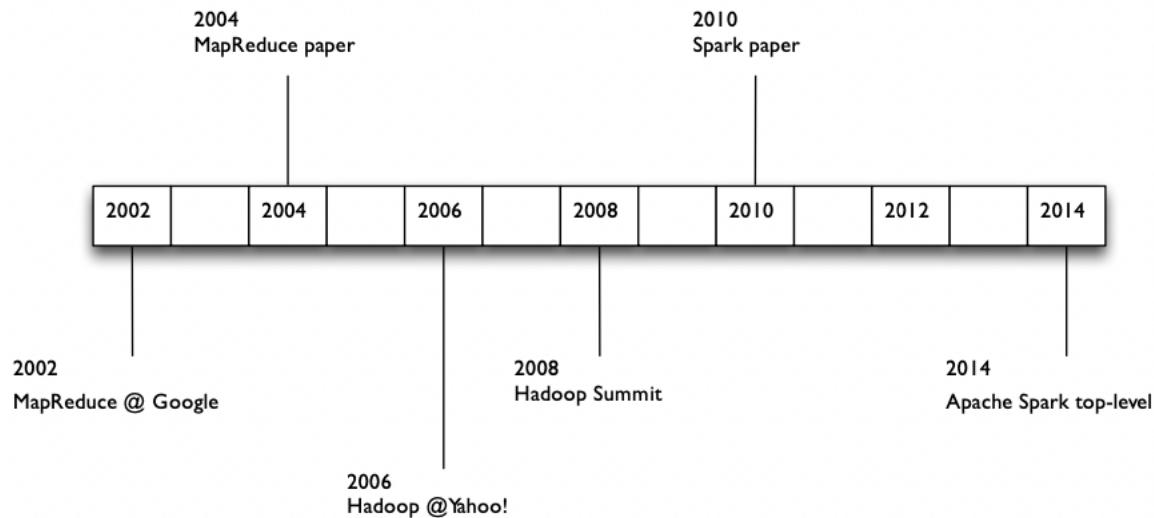


Spark Introduction

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

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

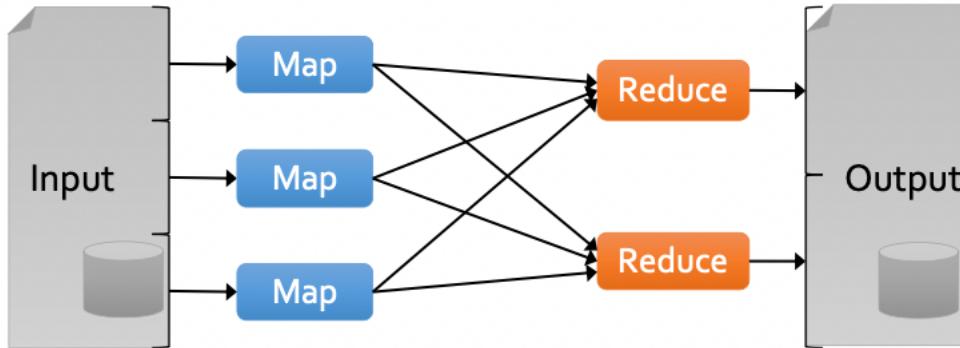
circa 2006 – **Apache**
Hadoop, originating from the Nutch Project Doug Cutting
research.yahoo.com/files/cutting.pdf

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/

MapReduce

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

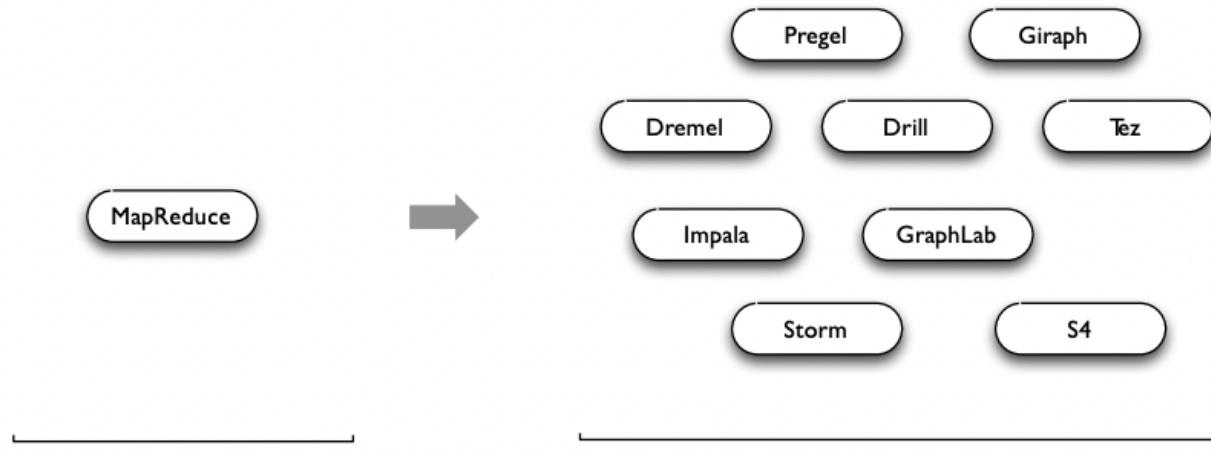


- Acyclic data flow is inefficient for applications that repeatedly reuse a **working set** of data:
 - o **Iterative** algorithms (machine learning, graphs).
 - o **Interactive** data mining tools (R, Excel, Python).

Data Processing Goals

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

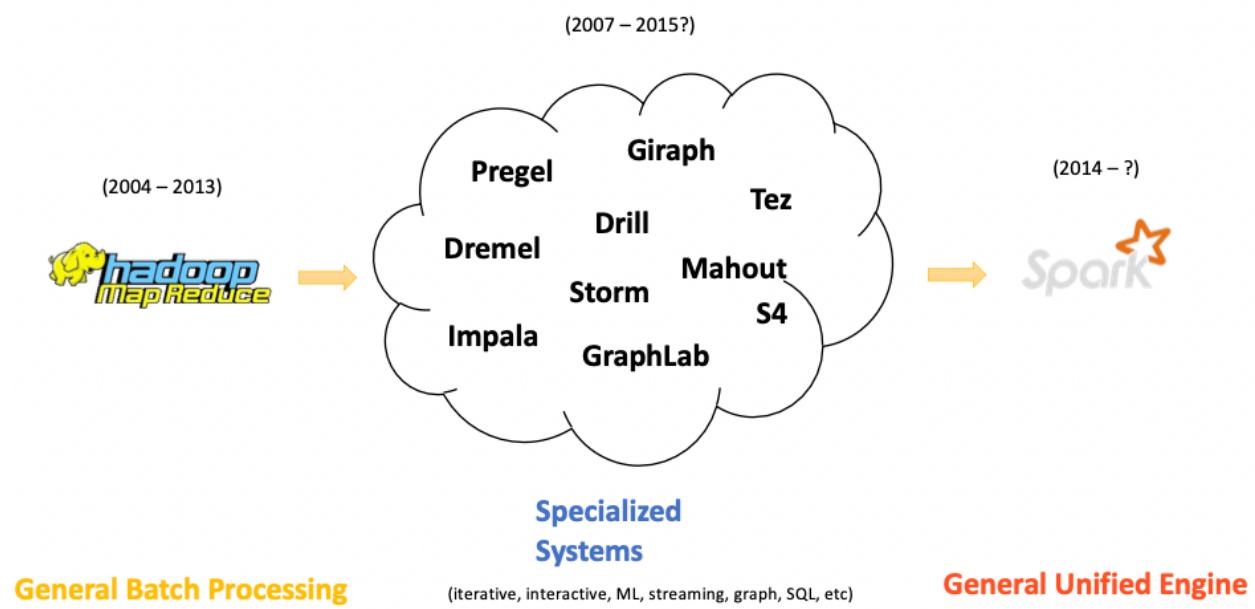
Specialized Systems

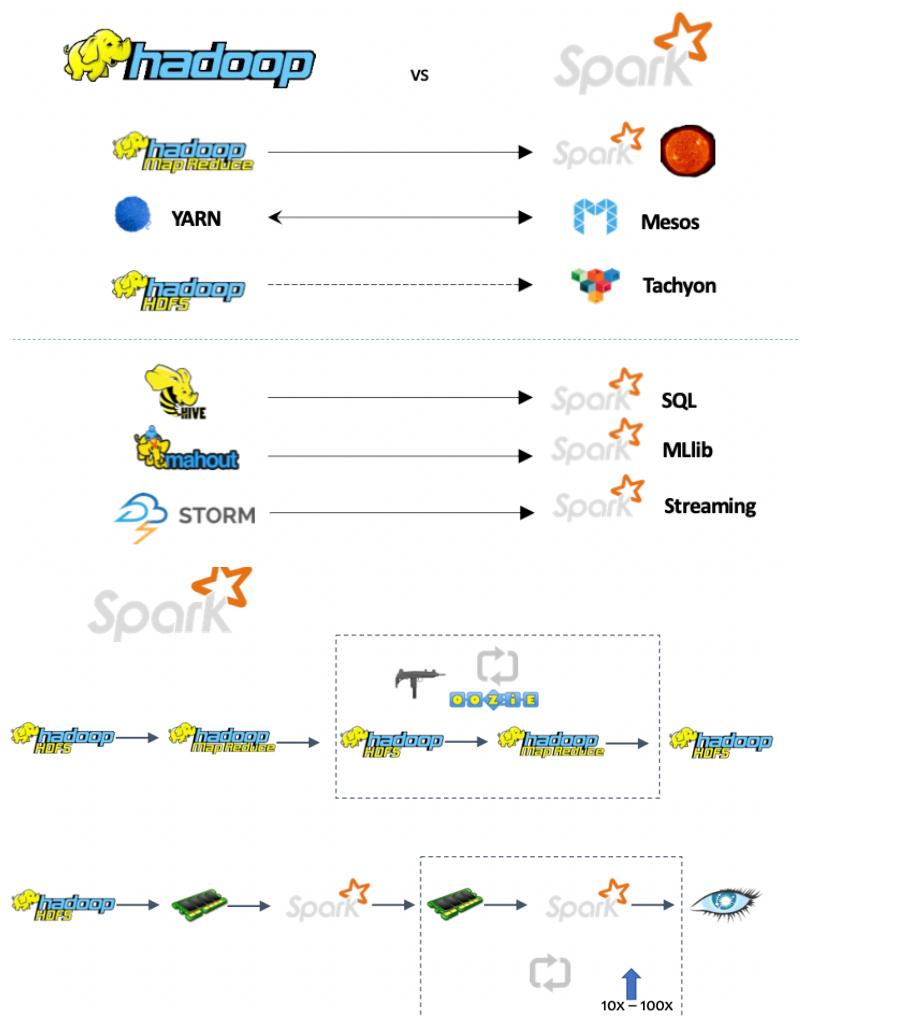


Storage vs Processing Wars



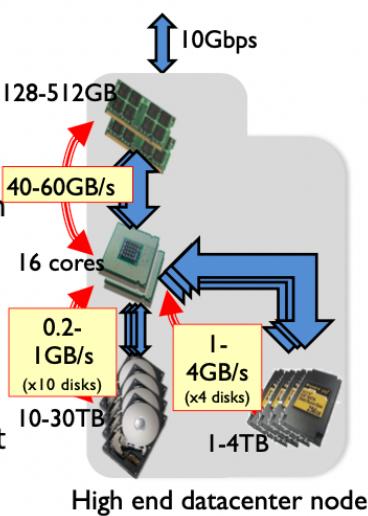
Specialized Systems



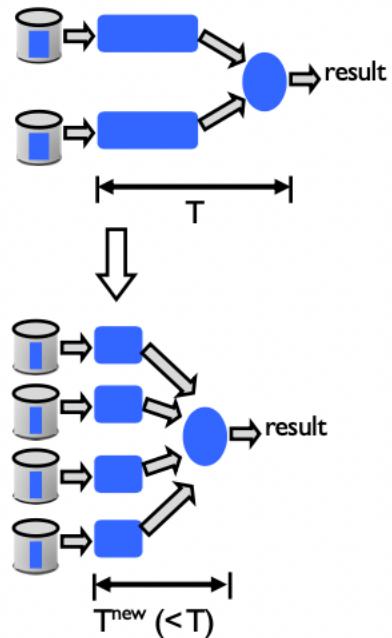


Support Interactive and Streaming Comp

- 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

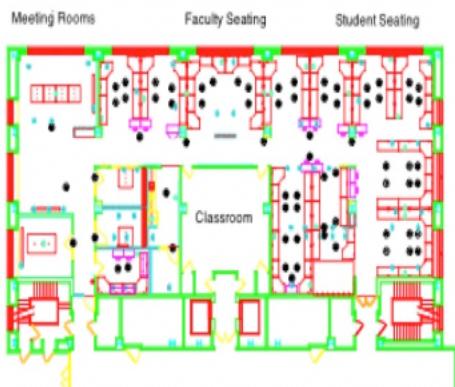


- Increase **parallelism**
- Why?
 - Reduce work per node → improve latency
- Techniques:
 - Low latency parallel **scheduler** that achieve high locality
 - Optimized **parallel communication patterns** (e.g., shuffle, broadcast)
 - Efficient **recovery** from failures and straggler mitigation



Berkeley AMP Lab

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



- Funding:
 - XData, CISE Expedition Grant



CISE Expedition Grant

- Industrial, founding sponsors
- 18 other sponsors, including



Goal: Next Generation of Analytics Data Stack for Industry &

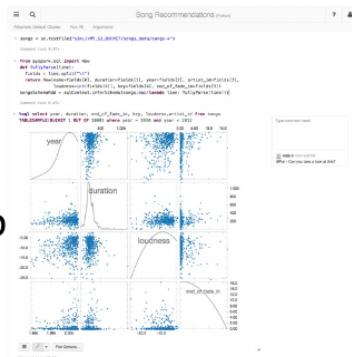
Research:

- Berkeley Data Analytics Stack (BDAS)
- Release as Open Source

Databricks



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



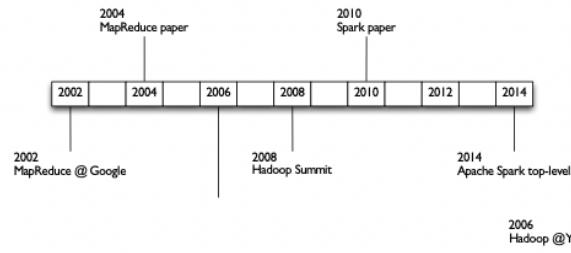
Databricks Cloud:
"A unified platform for building Big Data pipelines – from ETL to Exploration and Dashboards, to Advanced Analytics and Data Products."

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, Michael J. Franklin, Scott Shenker, Ion Stoica
University of California, Berkeley

Abstract

MapReduce and its variants have been highly successful in implementing large-scale data-intensive applications. However, these frameworks are built around an acyclic data flow model that is not well-suited for iterative or streaming computations. This paper proposes a working set abstraction that can support a working set of data across multiple parallel operations, allowing users to reuse data in memory and reduce disk access, as well as interactive data analysis tools.

We propose a new model of computation while retaining the availability and fault tolerance of MapReduce. To achieve these goals, Spark introduces an RDD abstraction that supports both iterative and streaming computations. An RDD is a read-only collection of objects partitioned across machines. Data is replicated by RDD, so if one partition is lost, Spark can recompute it by 10x in iterative machines, 100x in streaming jobs, and can be used to interactively analyze 100 GB datasets with sub-second response times.

1 Introduction

A new model of cluster computing has become widely popular, in which data-parallel computations are executed on clusters of commodity machines. Two main paradigms have emerged: MapReduce [10] and MapReduce/Merge [24]. MapReduce generalizes the types of data flows supported by MapReduce, while MapReduce/Merge generalizes the types of data flows supported by MapReduce. Both paradigms are implemented by providing a programming model where the user creates acyclic data flow graphs to pass input data through a set of map and reduce functions. In addition, they support iterative scheduling and re-use of results without user intervention. While MapReduce is well-suited for a wide class of applications, there are applications that can not be expressed efficiently as acyclic data flows. In fact, many machine learning and data mining applications require iterative data processing, which is what we see Hadoop users report that MapReduce is deficient:

- **Iterative jobs.** Many common machine learning algorithms (e.g., logistic regression) require iterative passes to optimize a parameter (e.g., through gradient descent). While each iteration can be expressed as a

MapReduce/Dryad job, each job must reload the data from disk, incurring a significant performance penalty.

• **Interactive analysis.** Hadoop is often used to run ad-hoc exploratory queries. These queries involve SQL interfaces such as Pig [21] and Hive [1]. Ideally, these queries can be run interactively, reading data from memory across a number of machines and query it rapidly. However, with Hadoop, each query incurs a significant performance overhead because it must run a separate MapReduce job and read data from disk.

This paper proposes a new model of computation, called **Spark**, which supports applications with working sets while providing similar scalability and fault tolerance.

The main abstraction in Spark is that of a **resilient distributed dataset** (RDD), which represents a collection of partitions of data that can be updated in place. A partition of data can be rebuilt if a partition is lost.

Users can explicitly update a partition of data in multiple MapReduce-like **parallel operations**. RDDs achieve fault tolerance through **replication**: if a partition of an RDD is lost, the RDD has enough information about how it was derived from other RDDs to be reconstructed. In contrast, MapReduce/Merge does not provide a general shared memory abstraction; they represent a coarse-grained abstraction on the one hand and a fine-grained abstraction on the other. In fact, MapReduce/Merge [24] is not well-suited for a variety of applications.

In this paper, we present a new abstraction for high-level programming language for the Java VM, and expose a simple API and a programming interface similar to MapReduce/Merge [24]. In fact, we have ported the MapReduce/Merge interface to our system. In addition, we found them well-suited for a variety of applications.

Our goal is to make iterative data processing easy, and early experiments with the system are encouraging. We show that Spark can outperform Hadoop by 10x in iterative computations, and can process 100 GB datasets with sub-second latency.

This paper is organized as follows. Section 2 describes

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

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
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 designed for applications that cannot be efficiently implemented using current computing frameworks. They provide abstractions for iterative and interactive data mining tools. In both cases, RDDs provide a restricted form of shared memory based on coarse-grained parallelism, and do not require data to be shared at shared state. However, we show that RDDs are expressive enough to capture a wide class of computations, including iterative operations (e.g., MapReduce) and interactive data mining (e.g., distributed file systems). Unlike MapReduce, RDDs do not capture. We have implemented RDDs in a system called Spark, and we evaluate through a variety of user applications and benchmarks.

1 Introduction

Current computing frameworks like MapReduce [10] and Dryad [11] have been widely adopted for large-scale data analytics. These systems let users write parallel computations using a set of high-level operators, without having to worry about low-level details.

Although current frameworks provide numerous abstractions for accessing a cluster's computational resources, they lack abstractions for leveraging distributed memory. This is particularly problematic for a class of emerging applications, those that reuse intermediate data results across multiple computations. Data reuse is common in iterative data mining, machine learning, and graph algorithms, including PageRank, K-means clustering, and logistic regression. Another compelling use case is interactive data mining, where a user runs multiple ad-hoc queries on the same set of data. Until recently, in most current frameworks, the only way to reuse data between computations (e.g., between two MapReduce passes) is to use an external distributed file system, e.g., a distributed file system. This incurs substantial overheads due to data replication, disk I/O, and serialization

time, which dominate application execution times.

Recognizing this problem, researchers have developed specialized frameworks for some applications that require shared memory. For example, GraphLab [12] offers iterative graph computations that keep intermediate data in memory, while Hadoop [7] offers an iterative MapReduce interface. However, these frameworks do not support iterative computations (e.g., keeping a sequence of MapReduce steps), and perform data sharing implicitly for these patterns. They do not provide abstractions for iterative data mining, nor do they support performing data mining on large datasets in memory and running ad-hoc queries across them.

In this paper, we present a new abstraction called **resilient distributed datasets** (RDDs) that enables users to work with data stored in a broad set of applications. RDDs are fault-tolerant, parallel data structures that let users explicitly persist intermediate results in memory, control their lifetime, and query them interactively, and manipulate them using a set of operators.

The main challenge in designing RDDs is defining a programming interface that can provide fault tolerance and interactivity. Existing abstractions for memory storage on clusters, such as distributed shared memory [34], key-value stores [25], databases, and Piccolo [27], offer an interface for users to store data in memory and query it (e.g., cells in a table). With this interface, the only ways to provide fault tolerance are to replicate the data across machines or to log updates across machines. Both approaches are slow, and are not suitable for iterative workflows, as they require copying large amounts of data over the cluster network, whose bandwidth is far lower than that of local memory.

In contrast to these systems, RDDs provide an interface based on coarse-grained transformations (e.g., map, filter, join) that apply the same operation to many partitions of data. This allows RDDs to provide fault tolerance by logging the transformations used to build a dataset (its lineage) rather than the actual data.¹ If a partition of an RDD is lost, the RDD can reconstruct itself based on the lineage information and the current state of the cluster.

¹Updating the data in some RDDs may be useful when a lineage chain grows large, however, and we discuss how to do it in §V.

"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.

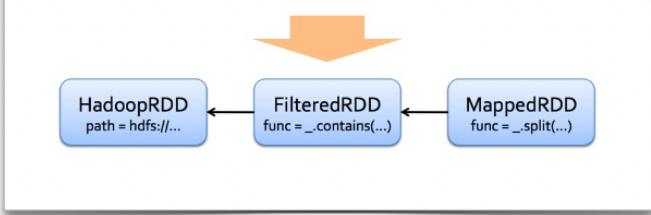
In both cases, keeping data in memory can improve performance by an order of magnitude."

April 2012

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/nU6vOZEJAb4



Analyze real time streams of data in ½ second intervals

www.cs.berkeley.edu/~matei/papers/2013/sosp_spark_streaming.pdf

Discretized Streams: Fault-Tolerant Streaming Computation at Scale
Matei Zaharia, Tohaga 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. These include applications at over-large scale that require parallel platforms that can handle faults and stragglers. Unfortunately, current distributed stream processing models provide fault recovery in an expensive manner, requiring hot replication or long recovery times. We propose a new model for distributed stream processing, *discretized streams* (DStreams), that overcomes these challenges. DStreams enable a partitioned state abstraction that augments traditional replication and backup schemes, and isolates stragglers. We show that streamlining high-per-node throughput amounts to a 10x improvement in both sub-second latency and sub-second fault recovery. Finally, DStreams can easily be composed with batch and interactive processing, enabling interesting mixed applications that combine these modes. We implement DStreams in a system called Spark Streaming.

1 Introduction
Much of “big data” is received in real time, and is most valuable at its time of arrival. For example, a social network may wish to detect trending conversation topics in

faults and stragglers (slow nodes). Both problems are inevitable in large clusters [17, 20]. Real-time streaming applications must deal with these issues. Fault recovery is even more important in streaming than it is in batch jobs, while a 30-second delay to recover from a fault or straggler is a catastrophe in a streaming setting, it can mean losing data for many key decisions, such as advertising targeting. Unfortunately, existing streaming systems have limited fault and straggler tolerance. Most streaming systems use a simple approach to handling faults [17], TimeStreams [33], MapReduce Online [11], and streaming databases [5, 9, 10], are based on a continuous operator model, in which long-running, stateful operators process data as it arrives, emitting results for new records. While this model is quite natural, it makes it difficult to handle faults and stragglers.

Quiescent gives us a new operator model, which performs recovery through two approaches [20]: replication, where there are two copies of each node [5, 34], or snapshot backup, where nodes buffer sent data and periodically checkpoint to a quiescent node [35, 11, 37]. Neither approach is attractive in large clusters: replication costs 2 × the hardware, while upzram backup takes a long time to recover, as the whole system must wait for a new node to serially rebroadcast the failed

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



Seamlessly mix SQL queries with Spark programs.

Spark SQL: Relational Data Processing in Spark

Michael Armbrust¹, Reynold S. Xin¹, Cheng Lian¹, Yin Huai¹, Davies Liu¹, Joseph K. Bradley¹, Xiangrui Meng¹, Tomer Kattan¹, Michael J. Franklin¹, Ali Ghodsi¹, Matei Zaharia^{1*}
¹Databricks Inc. *MIT CSAIL ¹AMPLab, 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 Spark, Spark SQL lets Spark programmers write declarative queries using the Scala programming language, reuse existing data sources and storage formats, generate, and define extension points. Using Catalyst, we have built a variety of features (e.g., schema inference for JSON, machine learning types, and query optimization) to better serve the needs for the complex needs of modern data analysis. We see Spark SQL as an evolution of SQL-on-Hadoop 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

Database; 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, care little about a workflow, but

while the popularity of relational systems shows that users often prefer writing declarative queries, the relational approach is inefficient for many big data applications. First, users want to perform complex operations that are hard to express in a declarative, structured, requiring custom code. Second, users want to perform advanced analytics, such as machine learning and graph processing, which are hard to express in a declarative, structured language. Finally, we have observed that most data pipelines would ideally be expressed with a combination of both declarative queries and complex procedural logic. Unfortunately, these two paradigms—relational and procedural—have until now remained largely disjoint, forcing users to choose one or the other.

This paper describes our effort to combine both models in Spark SQL, a major new component in Apache Spark [19]. Spark SQL provides a rich set of operators for both external data sources and internal data structures. It follows the well-known data frame API, the widely used data frame concept in R [32], but evaluates operations locally, rather than using distributed execution. This allows it to support the wide range of data sources and algorithms in big data. Spark SQL introduces a novel extensible optimizer called Catalyst. Catalyst is a hybrid optimizer that combines both relational rules, and data types for domains such as machine learning.

The Dataframe API offers rich relational/procedural integration without forcing users to learn both paradigms. It also supports streaming 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)
```



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 structures in data that are central to several applications 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 distributed, scalable tools for graph computation has led to the development of new graph-processing 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. In addition, these graph-processing systems provide limited fault-tolerance and support for interactive queries.

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 store and process sparse graph data structures. Similarly, we advance 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, our flexible 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 these graph processing 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 systems. Furthermore, while these frameworks address the challenges of distributed computation, they do not address the challenges of graph construction. Finally, these frameworks lack a collection of data-ETL, preprocessing and composition operators for ingesting and applying the results of computation. Finally, few frameworks have built-in support for instance graph computation.

Alternatively, also-parallel systems like MapReduce and Spark have been developed for 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 systems can be challenging and typically result in complex joins and excessive data movement that does not exploit the graph structure.

To address these challenges, we introduce GraphX, a graph computation library built on top of the open-source framework Spark. 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



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 access 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)
FROM Table
WHERE city='San Francisco'
WITHIN 2 SECONDS
```

Queries with Time Bounds

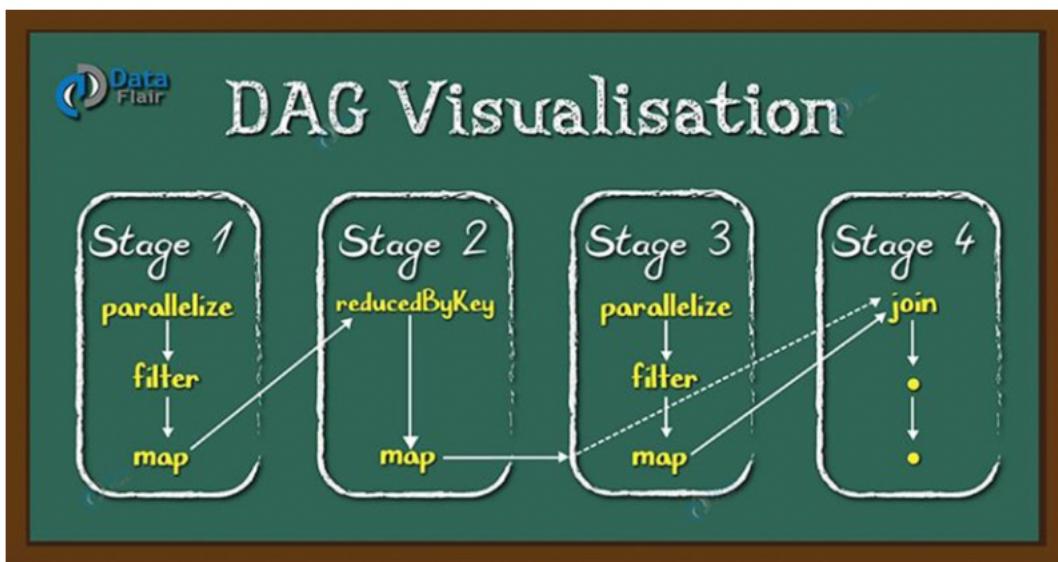
```
SELECT avg(sessionTime)
FROM Table
WHERE city='San Francisco'
ERROR 0.1 CONFIDENCE 95.0%
```

Queries with Error Bounds

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

- 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

DAG

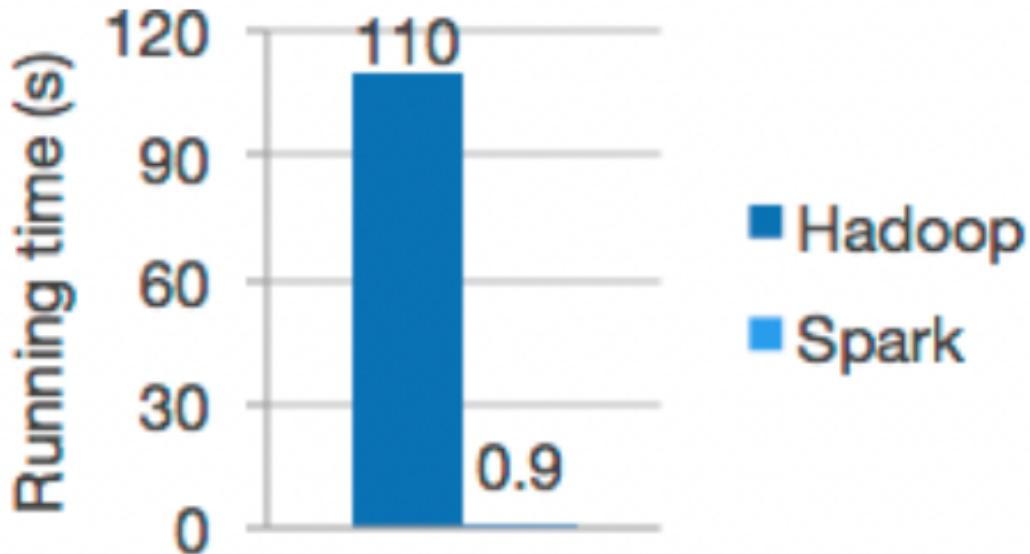


Spark Introduction

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

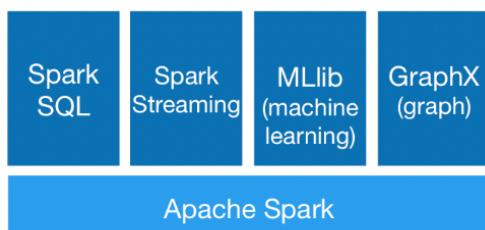
- **Ease of Use:** write applications quickly in Java, Scala, Python, R, SQL.
 - Offer over 80 high-level operators.
 - Use them interactively from Scala, Python, R and SQL.

```
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

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



- **Run Everywhere:**
 - o Run on Hadoop, Apache Mesos, Kubernetes, standalone or in the cloud.
 - o Access data in HDFS, Aluxio, Apache Cassandra, Apache HBase, Apache Hive, etc.

Comparison between Hadoop and Spark

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

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.**

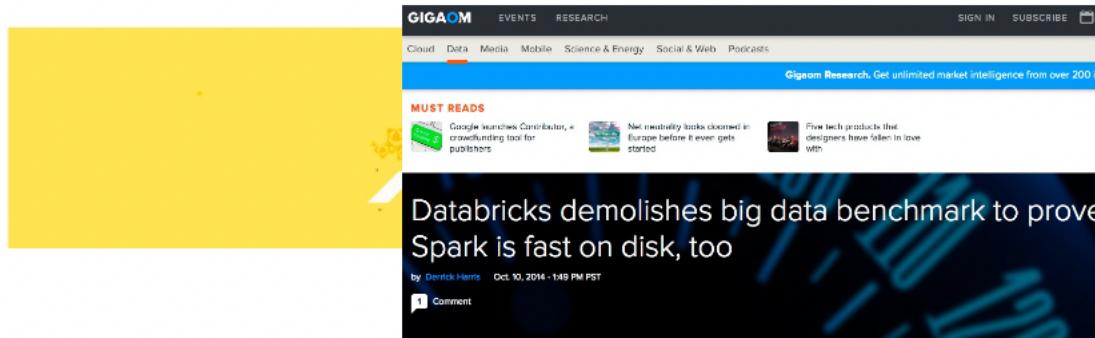
**All the sorting took place on disk (HDFS)
without using Spark's in-memory cache!**

WIRED GEAR SCIENCE ENTERTAINMENT BUSINESS SECURITY DESIGN OPINION MAGAZINE
 ENTERPRISE | big data | databricks | google | Hadoop

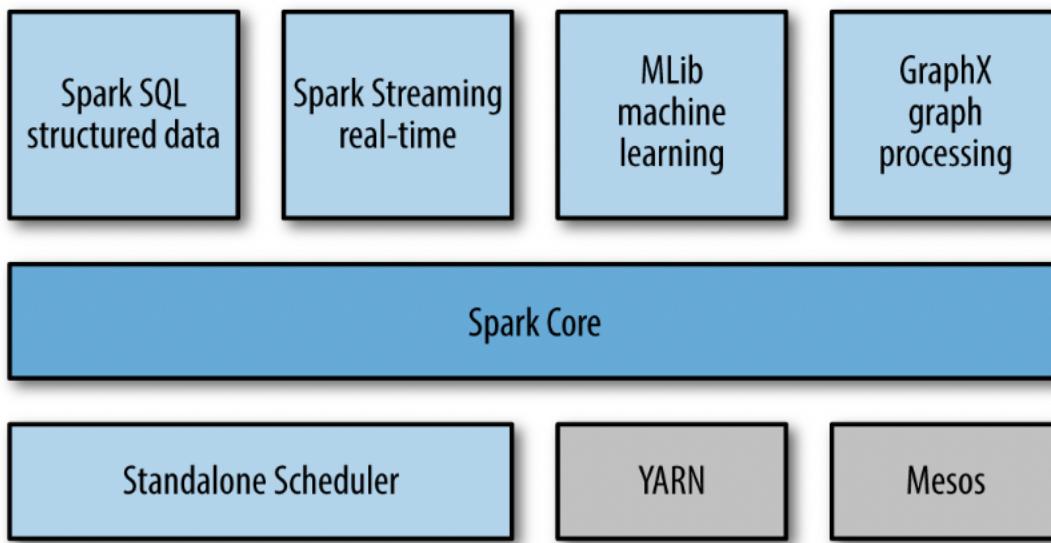
Startup Crunches 100 Terabytes of Data in a Record 23 Minutes

BY KLINT FINLEY 10.13.14 | 2:36 PM | PERMALINK

[Share](#) 1.1k [Tweet](#) 789 [+1](#) 75 [in Share](#) 565 [Pin it](#)



Components of Stack



- **Spark Core:**
 - o Contain the basic functionality of Spark including task scheduling, memory management, fault recovery, etc.
 - o Provide APIs for building and manipulating RDDs.
- **SparkSQL:**
 - o Allow querying structured data via SQL, Hive Query Language.
 - o Allow combining SQL queries and data manipulations in Python, Java, Scala.
- **SparkStreaming:** enables processing of live streams of data via APIs.
- **MLlib:**

- Contains common machine language functionality.
- Provides multiple types of algorithms: classification, regression clustering, etc.
- GraphX:
 - Library for manipulating graphs and performing graph-parallel computations.
 - Extend Spark RDD API.
- Cluster Managers:
 - Hadoop YARN.
 - Apache Mesos
 - Standalone Scheduler (simple manager in Spark).

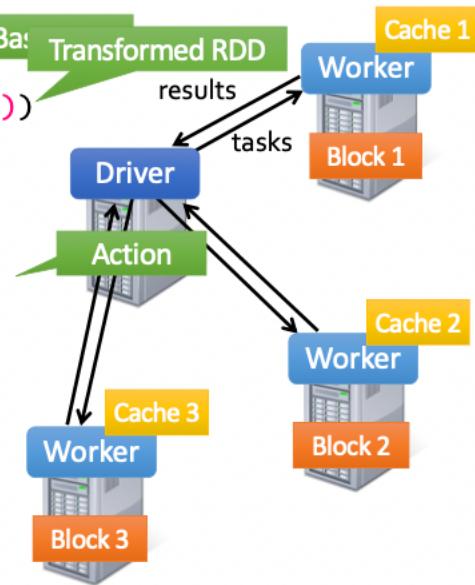
Resilient Distributed Dataset – RDD Basics

- 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.

Example

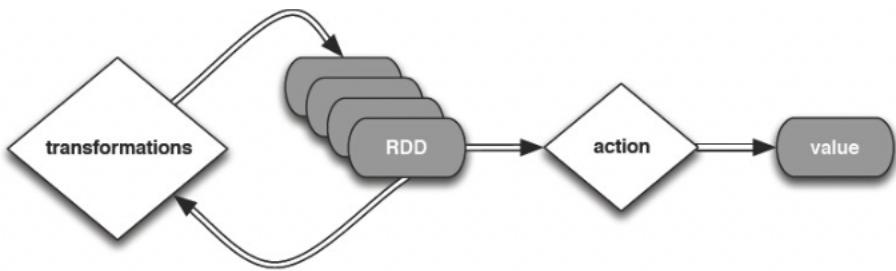
- Load error messages from a log in to memory, then interactively search for various patterns.

```
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  
...
```



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.



Transformations

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

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.

Resilient Distributed Dataset – Creating RDDs

1. Parallelizing a collection: uses `parallelize()`
- Python:
 - o `lines = sc.parallelize(["pandas", "i like pandas"])`
- Scala:
 - o `val lines = sc.parallelize(List("pandas", "i like pandas"))`
- Java:
 - o `JavaRDD<String> lines = sc.parallelize(Arrays.asList("pandas", "I like pandas"));`
2. Loading data from external storage
- Python:
 - o `lines = sc.textFile("/path/to/README.md")`
- Scala:
 - o `val lines = sc.textFile("/path/to/README.md")`
- Java:
 - o `JavaRDD<String> lines = sc.textFile("/path/to/README.md");`

Resilient Distributed Dataset – RDD Operations

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

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");
        }
    }
);
```

`filter()`

- Does not change the existing `inputRDD`.
- Returns a pointer to an entirely new RDD.
- `inputRDD` still can be reuse.

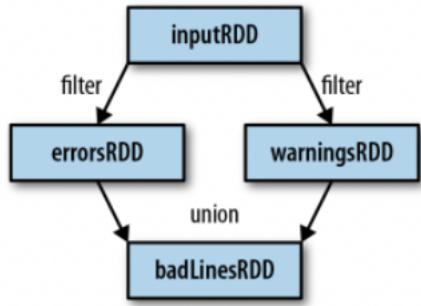
`union()`

```
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.

Spark keeps track dependencies between RDDs, called the lineage graph.

Allow recovering lost data:



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);
}

```

Resilient Distributed Dataset – Common Transformation and Actions

Transformations	Actions
map flatMap filter sample union groupByKey reduceByKey join cache ...	reduce collect count save lookupKey ...

Transformation

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

<i>transformation</i>	<i>description</i>
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 <i>ascending</i> 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)

Actions

action	description
reduce(func)	aggregate the elements of the dataset using a function <i>func</i> (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 <i>take(1)</i>
take(n)	return an array with the first <i>n</i> 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 <i>num</i> elements of the dataset, with or without replacement, using the given random number generator seed

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 <i>toString</i> 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 <i>func</i> on each element of the dataset – usually done for side effects such as updating an accumulator variable or interacting with external storage systems

Resilient Distributed Dataset – Persistence (Caching)

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	

Persistence

Example:

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

Table of Contents

<i>Spark Introduction</i>	1
<i>History of Spark</i>	1
<i>MapReduce</i>	2
<i>Data Processing Goals</i>	2
<i>Specialized Systems</i>	2
<i>Storage vs Processing Wars</i>	3
<i>Specialized Systems</i>	3
<i>Support Interactive and Streaming Comp</i>	4
<i>Berkeley AMP Lab</i>	5
<i>Databricks</i>	6
<i>History of Spark</i>	7
<i>DAG</i>	10
<i>Spark Introduction</i>	10
<i>What is Apache Spark</i>	10
<i>Comparison between Hadoop and Spark</i>	12
<i>100TB Daytona Sort Competition</i>	12
<i>Components of Stack</i>	13
<i>Resilient Distributed Dataset – RDD Basics</i>	14
<i>Example</i>	14
<i>RDD Basics</i>	14
<i>Transformations</i>	15
<i>Typical works of a Spark program</i>	15
<i>Resilient Distributed Dataset – Creating RDDs</i>	15
<i>Resilient Distributed Dataset – RDD Operations</i>	16
<i>Transformation</i>	16
<i>Actions</i>	17
<i>Resilient Distributed Dataset – Common Transformation and Actions</i>	17
<i>Transformation</i>	18
<i>Actions</i>	19
<i>Resilient Distributed Dataset – Persistence (Caching)</i>	20

<i>Persistence levels</i>	20
<i>Persistence</i>	20