

CSE 4/587

Data Intensive Computing

Dr. Eric Mikida

epmikida@buffalo.edu

208 Capen Hall

Day 20
Apache Spark

Announcements and Feedback

- Midterms have been returned
- Phase 1 grades have been released

References

- **Advanced Analytics with Spark** by S. Ryza, U. Laserson, S. Owen and J. Wills
- **Apache Spark documentation**
 - <http://spark.apache.org/>
 - <http://spark.apache.org/docs/latest/programming-guide.html>
- **Pyspark**
 - <http://spark.apache.org/docs/latest/api/python/pyspark.html>
- **Resilient Distributed Dataset: A Fault-tolerant Abstraction for in-Memory Cluster Computing.** M. Zaharia et al.
 - <https://www.usenix.org/system/files/conference/nsdi12/nsdi12-final138.pdf>

Challenges

Data cleaning: Majority of the work that goes into analyses lies in pre-processing data

- Munging, fusing, mushing and cleansing
- We need computational methods to clean data and data pipeline certainly should include an important step of “data cleaning” and “feature engineering”.
- Choosing from many features, the relevant features.
- Designing a math model from a 2D array (Ex: page rank)

Challenges

Data cleaning: Majority of the work that goes into analyses lies in pre-processing data

- Munging, fusing, musing and cleansing
- We need computational methods to clean data and data pipeline certainly should include an important step of “data cleaning” and “feature engineering”.
- Choosing from many features, the relevant features.
- Designing a math model from a 2D array (Ex: page rank)

Need to avoid delays in repeated reading of data

Challenges

Iteration: Iteration is a fundamental part of data science.

- Modeling and analysis require typically multiple passes over the same data
- Machine learning algorithms and statistical procedures like stochastic gradient and expected maximization involve repeated scans to reach convergence
- Choosing the right features, picking the right algorithms, running the right significance tests, finding the right hyperparameters: **all require experimentation**

Challenges

Iteration: Iteration is a fundamental part of data science.

- Modeling and analysis require typically multiple passes over the same data
- Machine learning algorithms and statistical procedures like stochastic gradient and expected maximization involve repeated scans to reach convergence
- Choosing the right features, picking the right algorithms, running the right significance tests, finding the right hyperparameters: **all require experimentation**

Need to avoid delays in repeated reading of data

Challenges

Information updates: The results of data analysis presented and **the application becomes part of the production system...**

- This system must frequently or in real time update itself driven by the availability of new data; ie fraud detection system.

Challenges

Information updates: The results of data analysis presented and **the application becomes part of the production system...**

- This system must frequently or in real time update itself driven by the availability of new data; ie fraud detection system.

How about the existing approaches?

- C++, Java are not good for EDA
- R is slow for large data sets and does not integrate well with production stacks
- Read-Evaluate-Print-Loop (REPL) are good for interaction but not work production

Challenges

Information updates: The results of data analysis presented and **the application becomes part of the production system...**

- This system must frequently or in real time update itself driven by the availability of new data; ie fraud detection system.

How about the existing approaches?

- C++, Java are not good for EDA
- R is slow for large data sets and does not integrate well with production stacks
- Read-Evaluate-Print-Loop (REPL) are good for interaction but not work production

Want a framework that makes modeling easy, but also fits well in production systems

The Apache Hadoop Stack



Hadoop User Experience (HUE)



Data Exchange



Sqoop



Zoo Keeper

Coordination

Pig Scripting



Hive SQL



Mahout ML



Oozie Workflow



APACHE HBASE

Hbase

Columnar data store

Flume



Log Control



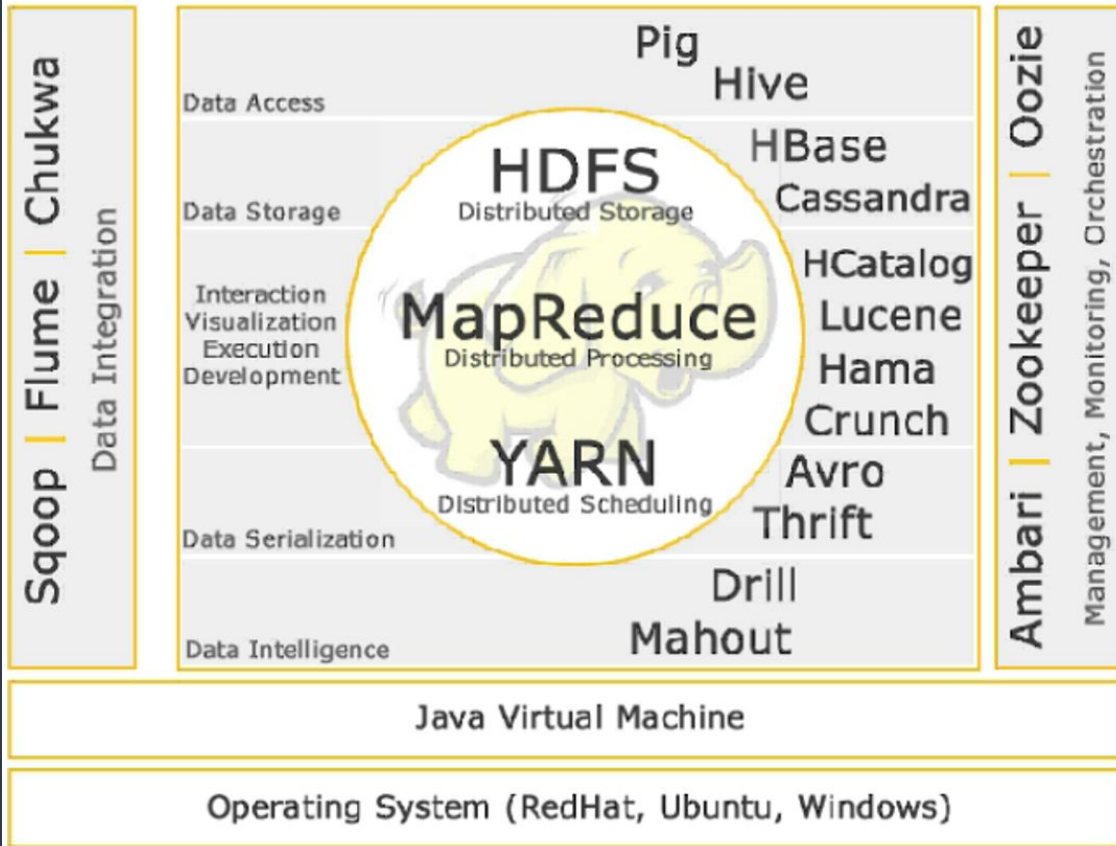
YARN/Map Reduce V2



Hadoop Distributed File System



Hadoop Technologies



Pig Abstraction Layer

Raw MapReduce is difficult to install and program (so why did we learn it?)

There are many models that simplify designing MapReduce applications:

- MRJob for python developers
- Elastic MapReduce (EMR) from Amazon AWS
- Hive from Apache via Facebook
- Pig from Apache via Yahoo
- And others...

Pig is a data flow language, conceptually closer to the way we solve problems...

Pig Data Flow Language (Pig Latin)

Pig data flow language describes a *directed acyclic graph* (DAG)

- Edges are data flows
- Nodes are operations.

There are no if statements or for loops in pig

- Procedural and object-oriented languages describe control flow; data flow is a side-effect.
- Pig focuses on data flow

Sample Pig Script

2	4	5
-2	3	4
3	5	6
-4	5	7
-7	4	6
3	4	5

```
A = LOAD 'data3' AS (x,y,z);
```

```
B = FILTER A BY x > 0;
```

```
C = GROUP B BY x;
```

```
D = FOREACH C GENERATE group,COUNT(B);
```

```
STORE D INTO 'p6out';
```

Pig Latin Constructs

- LOAD
- FILTER
- GROUP
- FOREACH
- GENERATE (apply some function from piggybank)
- STORE (DUMP for interactive debugging)

Word Count in Pig

```
input_lines = LOAD '/tmp/my-copy-of-all-pages-on-internet' AS (line:chararray);
words = FOREACH input_lines GENERATE FLATTEN(TOKENIZE(line)) AS word;
-- filter out any words that are just white spaces
filtered_words = FILTER words BY word MATCHES '\\w+';
-- create a group for each word
word_groups = GROUP filtered_words BY word;
-- count the entries in each group
word_count = FOREACH word_groups GENERATE COUNT(filtered_words) AS count,group AS word;
ordered_word_count = ORDER word_count BY count DESC;
STORE ordered_word_count INTO '/tmp/number-of-words-on-internet';
```

Apache Spark

Apache Spark is an open-source, distributed processing system commonly used for big data workloads.

- Utilizes in-memory caching
- Optimized execution for fast performance
- Supports general batch processing, streaming analytics, machine learning, graph databases, and ad hoc queries

Spark vs MapReduce

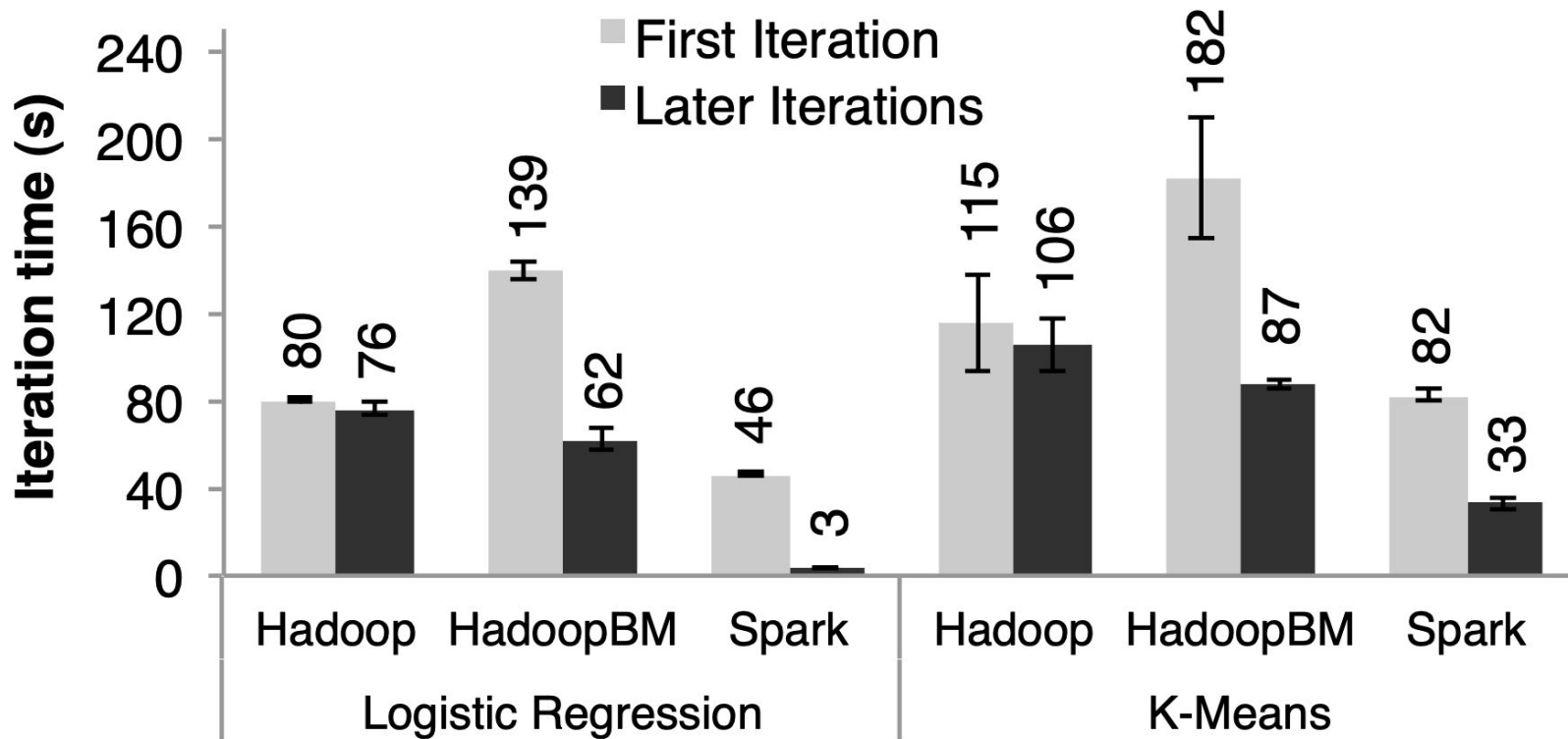
MapReduce offers linear scalability and fault tolerance for processing very large data sets

Spark maintains this revolutionary approach brought about by MapReduce

It also improves it in four different ways:

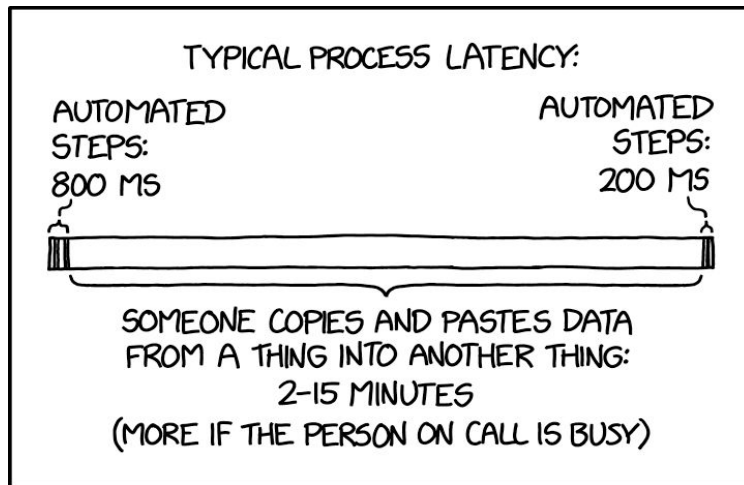
- Executes a series of operations specified in a DAG
 - Allows one stage of MR to send the results to the next (*Similar to Microsoft Dryad*)
- Provides a rich set of operations to express computation more naturally (*Similar to Pig*)
- Improves on in-memory computations through its **Resilient Distributed Data (RDD)**
 - Future steps dealing with the same data do not have to reload it from the disk
- Well-suited for highly iterative computing

Sample of Performance Results



Programming Productivity

Biggest bottleneck in data applications is not CPU, disk, or network but analyst productivity



Programming Productivity

Biggest bottleneck in data applications is not CPU, disk, or network but analyst productivity

If only we could collapse the entire pipeline from pre-processing of data to model evaluation into a single programming environment...

Spark transitions seamlessly between exploratory analytics and operational analytics

Word Count in Spark (Python API)

```
text_file = spark.textFile("hdfs://...")

text_file.flatMap(lambda line: line.split())
            .map(lambda word: (word, 1))
            .reduceByKey(lambda a, b: a+b)
```

Spark APIs

Spark runs on Hadoop, Mesos, standalone, or in the cloud. It can access diverse data sources including HDFS, Cassandra, HBase, and S3

It also provides APIs in many common languages:

- Scala API
- Java API
- Python API
- Dataframes API
- R API

Python Lambda Functions

Spark's Python API allows the use of lambda functions to transform data

```
x = lambda a : a + 10  
print(x(5))
```

```
x = lambda a, b : a * b  
print(x(5, 6))
```

Spark Ecosystem

Spark
SQL

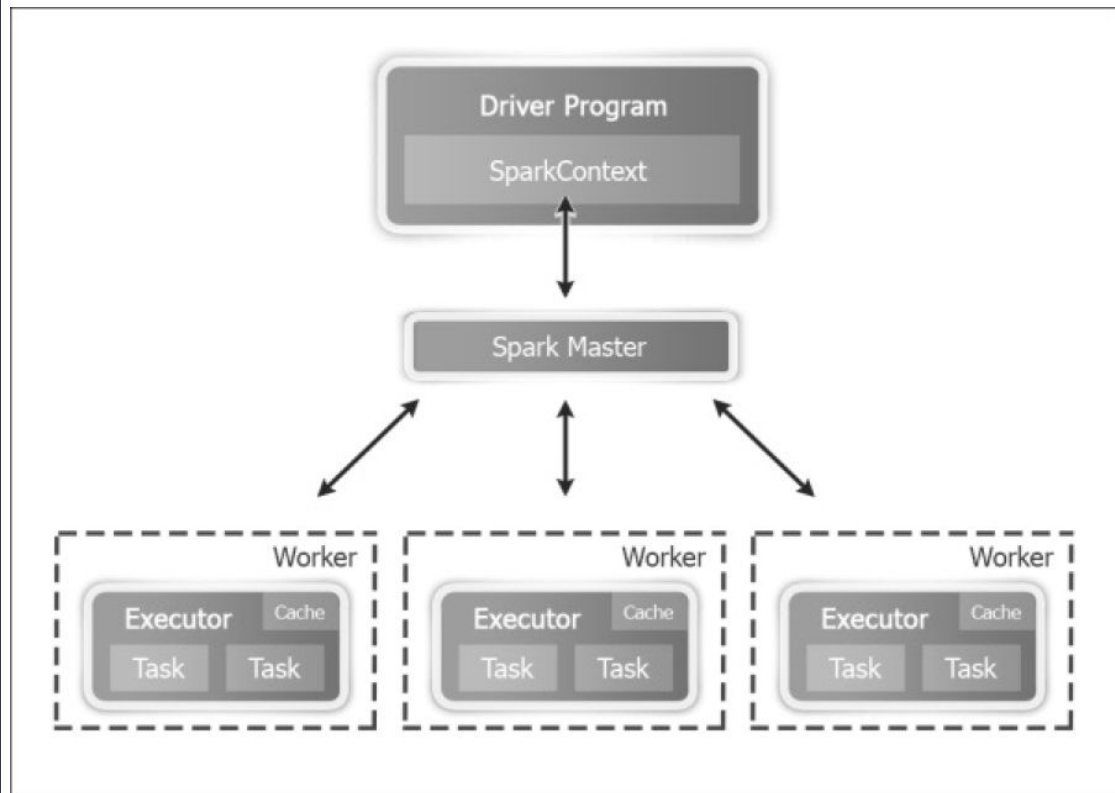
Spark
Streaming

MLlib
(machine
learning)

GraphX
(graph)

Apache Spark

Spark Architecture



Programming Model

- **Spark Context:** sc
- **RDD:** Resilient Distributed Datasets
 - Transformations and actions on RDDs
- **Diverse set of data sources:** HDFS, relational databases
- **Diverse APIs:** JavaAPI, Python API, Scala API, Dataframe API, R API

Spark Context

Spark Context (sc) is an object

- Main entry point for Spark applications
- Just like any object it has methods associated with it
- We can see what those methods are (in Scala `sc.[\t]`)
- Some of the methods:
 - `getConf`
 - `runJob`
 - `addFile`
 - `cancelAllJobs`
 - `makeRDD`

Resilient Distributed Datasets (RDDs)

A core Spark concept is **Resilient Distributed Datasets (RDD)** which is a fault tolerant collection of elements that can be operated in parallel

An RDD is a convenient way to describe the computations that we want to perform in small independent steps and in parallel

There are two ways to create RDDs:

- Parallelizing an existing collection in the driver program; performing a transformation on one or more existing RDDs, like filtering records, aggregating records by a common key or by joining multiple RDDs together.
- Using SparkContext to create an RDD from an external dataset in an external storage system such as a shared filesystem, HDFS, HBase or any data source offering a Hadoop input format

Examples

A few transformations to build a dataset and store into a file:

```
text-file = sc.textFile("hdfs://...")
counts = text_file.flatMap(lambda line: line.split(" "))
    .map(lambda word: (word,1))
    .reduceByKey(lambda a,b: a+b)
counts.saveAsTextFile("hdfs://..")
```

Resilient Distributed Datasets (RDDs)

The building block of the Spark API

(<http://spark.apache.org/docs/latest/programming-guide.html#resilient-distributed-datasets-rdds>)

In RDD API there are two types of operations:

1. *Transformations* that define a new data set based on previous ones
2. *Actions* which kick off a job to execute on a cluster

RDD

Transformations and Actions

Transformations

- map (func)
- flatMap(func)
- filter(func)
- groupByKey()
- reduceByKey(func)
- mapValues(func)
- sample(...)
- union(other)
- distinct()
- sortByKey()
- ...

Actions

- reduce(func)
- collect()
- count()
- first()
- take(n)
- saveAsTextFile(path)
- countByKey()
- foreach(func)
- ...

RDD Transformations and Actions

Transformations	<ul style="list-style-type: none"><i>map</i>($f : T \Rightarrow U$) : $RDD[T] \Rightarrow RDD[U]$<i>filter</i>($f : T \Rightarrow \text{Bool}$) : $RDD[T] \Rightarrow RDD[T]$<i>flatMap</i>($f : T \Rightarrow \text{Seq}[U]$) : $RDD[T] \Rightarrow RDD[U]$<i>sample</i>(<i>fraction</i> : Float) : $RDD[T] \Rightarrow RDD[T]$ (Deterministic sampling)<i>groupByKey</i>() : $RDD[(K, V)] \Rightarrow RDD[(K, \text{Seq}[V])]$<i>reduceByKey</i>($f : (V, V) \Rightarrow V$) : $RDD[(K, V)] \Rightarrow RDD[(K, V)]$<i>union</i>() : $(RDD[T], RDD[T]) \Rightarrow RDD[T]$<i>join</i>() : $(RDD[(K, V)], RDD[(K, W)]) \Rightarrow RDD[(K, (V, W))]$<i>cogroup</i>() : $(RDD[(K, V)], RDD[(K, W)]) \Rightarrow RDD[(K, (\text{Seq}[V], \text{Seq}[W]))]$<i>crossProduct</i>() : $(RDD[T], RDD[U]) \Rightarrow RDD[(T, U)]$<i>mapValues</i>($f : V \Rightarrow W$) : $RDD[(K, V)] \Rightarrow RDD[(K, W)]$ (Preserves partitioning)<i>sort</i>($c : \text{Comparator}[K]$) : $RDD[(K, V)] \Rightarrow RDD[(K, V)]$<i>partitionBy</i>($p : \text{Partitioner}[K]$) : $RDD[(K, V)] \Rightarrow RDD[(K, V)]$
Actions	<ul style="list-style-type: none"><i>count</i>() : $RDD[T] \Rightarrow \text{Long}$<i>collect</i>() : $RDD[T] \Rightarrow \text{Seq}[T]$<i>reduce</i>($f : (T, T) \Rightarrow T$) : $RDD[T] \Rightarrow T$<i>lookup</i>($k : K$) : $RDD[(K, V)] \Rightarrow \text{Seq}[V]$ (On hash/range partitioned RDDs)<i>save</i>(<i>path</i> : String) : Outputs RDD to a storage system, <i>e.g.</i>, HDFS

Table 2: Transformations and actions available on RDDs in Spark. $\text{Seq}[T]$ denotes a sequence of elements of type T .

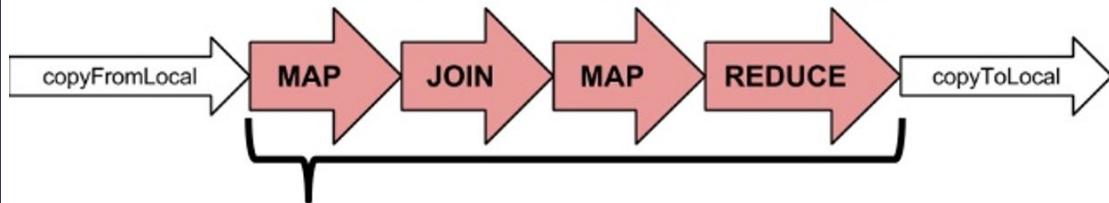
Resilient Distributed Datasets (RDDs)

A **distributed memory abstraction** that enables **in-memory computations** on large clusters in a **fault-tolerant** manner

- Motivation: iterative algorithms, interactive data mining tools
 - In both cases above keeping data in memory will help enormously for performance improvement
- RDDs are parallel data structures allowing coarse grained transformations
- It ***provides fault tolerance by storing the lineage as opposed to the actual data*** as done in Hadoop

Transformations VS Actions

RDD Transformation vs. Action



- **Transformations are lazy:** nothing actually happens when this code is evaluated
- **RDDs are computed only when an *action* is called on them, e.g.,**
 - Calculate statistics over the elements of an RDD (count, mean)
 - Save the RDD to a file (saveAsTextFile)
 - Reduce elements of an RDD into a single object or value (reduce)
- **Allows you to define partitioning/caching behavior *after* defining the RDD but *before* calculating its contents**

RDD Lineage

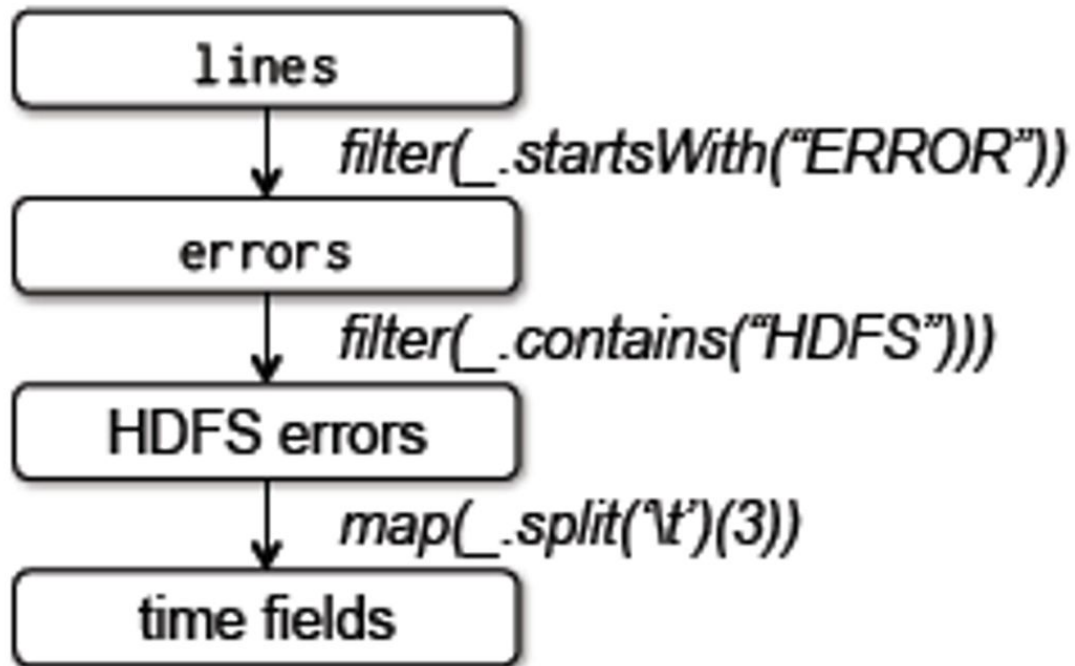
An RDD can depend on zero or more other RDDs

- ie when $x = y.map(\dots)$, x will depend on y
- These dependency relationships can be thought of as a graph.

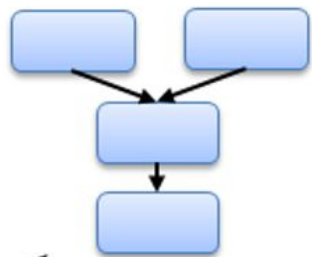
You can call this graph a lineage graph, as it represents the derivation of each RDD

- It is also necessarily a **DAG**, since a loop is impossible to be present in it.
- Narrow dependencies, where a shuffle is not required (think map and filter) can be collapsed into a single **stage**.
 - A stage is a unit of execution, generated by the scheduler from RDD dependency graph
 - Stages also depend on each other and the scheduler builds and uses this dependency graph (which is also necessarily a DAG) to schedule the stages

RDD Lineage



RDD Objects

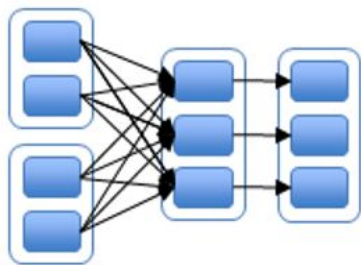


```
rdd1.join(rdd2)  
.groupBy(...)  
.filter(...)
```

build operator DAG

DAG

DAGScheduler



split graph into
stages of tasks

submit each
stage as ready

agnostic to
operators!

TaskScheduler

Cluster
manager

TaskSet

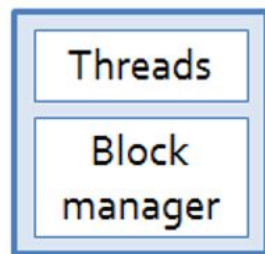
launch tasks via
cluster manager

retry failed or
straggling tasks

doesn't know
about stages

Task

Worker



execute tasks

store and serve
blocks

stage
failed

Representing RDDs

Each RDD is represented through a common interface that exposes 5 pieces of information:

1. A set of partitions, atomic pieces of datasets
2. Set of dependencies on the parent RDDs
3. Function for computing the RDD from the parents
4. Metadata about partitioning scheme
5. Data placement

See table 3 in the RDD paper.

Dependencies

Narrow dependencies: each parent RDD partition used by at most one child; ie map()

- allow pipelined execution: example map() and filter() in iterative fashion
- recovery after node failure is more efficient

Wide dependencies: multiple child partitions may depend on a parent RDD; ie join()

- Single failed node in a wide dependency lineage graph may cause loss of partition in many ancestral dependencies