# Spark

Stony Brook University
CSE545, Spring 2019

# Situations where MapReduce is not efficient

DFS ➡ Map ➡ LocalFS ➡ Network ➡ Reduce ➡ DFS ➡ Map ➡ ...

# Situations where MapReduce is not efficient

- Long pipelines sharing data

- Interactive applications

- Streaming applications

- Iterative algorithms (optimization problems)

DFS ➡ Map ➡ LocalFS ➡ Network ➡ Reduce ➡ DFS ➡ Map ➡ ...

(Anytime where MapReduce would need to write and read from disk a lot).

# Situations where MapReduce is not efficient

- Long pipelines sharing data

- Interactive applications

- Streaming applications

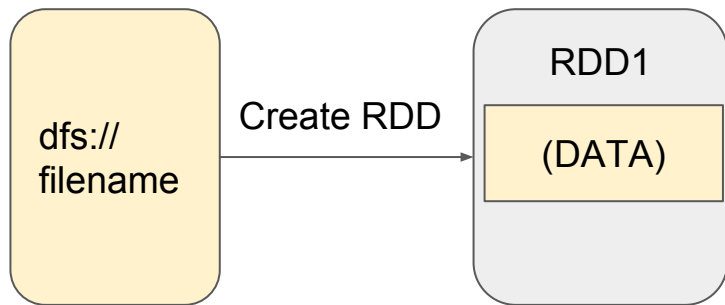- Iterative algorithms (optimization problems)

DFS ➡ Map ➡ LocalFS ➡ Network ➡ Reduce ➡ DFS ➡ Map ➡ ...

(Anytime where MapReduce would need to write and read from disk a lot).

# Spark's Big Idea

Resilient Distributed Datasets (RDDs) -- Read-only partitioned collection of records (like a DFS) but with a record of how the dataset was created as combination of *transformations* from other dataset(s).
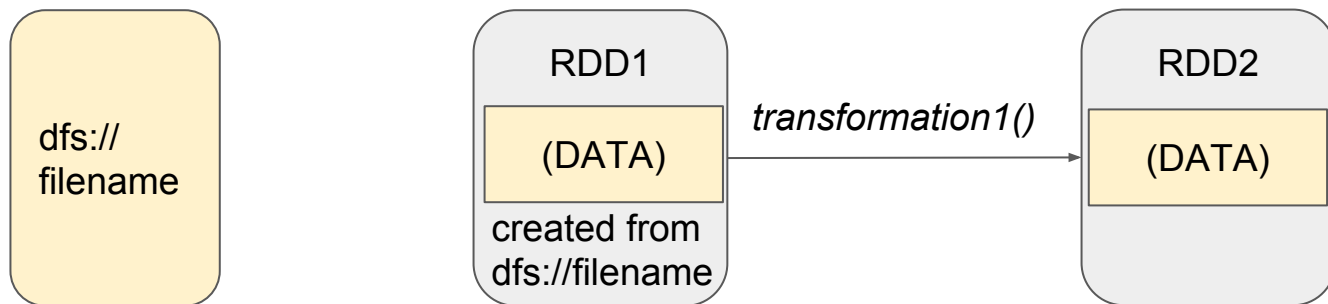
# Spark's Big Idea

Resilient Distributed Datasets (RDDs) -- Read-only partitioned collection of records (like a DFS) but with a record of how the dataset was created as combination of *transformations* from other dataset(s).

```
┌──────────┐                    ┌──────────────┐
│          │                    │    RDD1      │
│ dfs://   │   Create RDD       │ ┌──────────┐ │
│ filename │ ─────────────────▶ │ │ (DATA)   │ │
│          │                    │ └──────────┘ │
│          │                    │              │
└──────────┘                    └──────────────┘
```

# Spark's Big Idea

Resilient Distributed Datasets (RDDs) -- Read-only partitioned collection of records (like a DFS) but with a record of how the dataset was created as combination of *transformations* from other dataset(s).

# Spark's Big Idea

Resilient Distributed Datasets (RDDs) -- Read-only partitioned collection of records (like a DFS) but with a record of how the dataset was created as combination of *transformations* from other dataset(s).

dfs://
filename

**RDD1**

(can drop the data)

created from
dfs://filename

**RDD2**

(DATA)

*transformation1*
from RDD1

*transformation2()*

**RDD3**
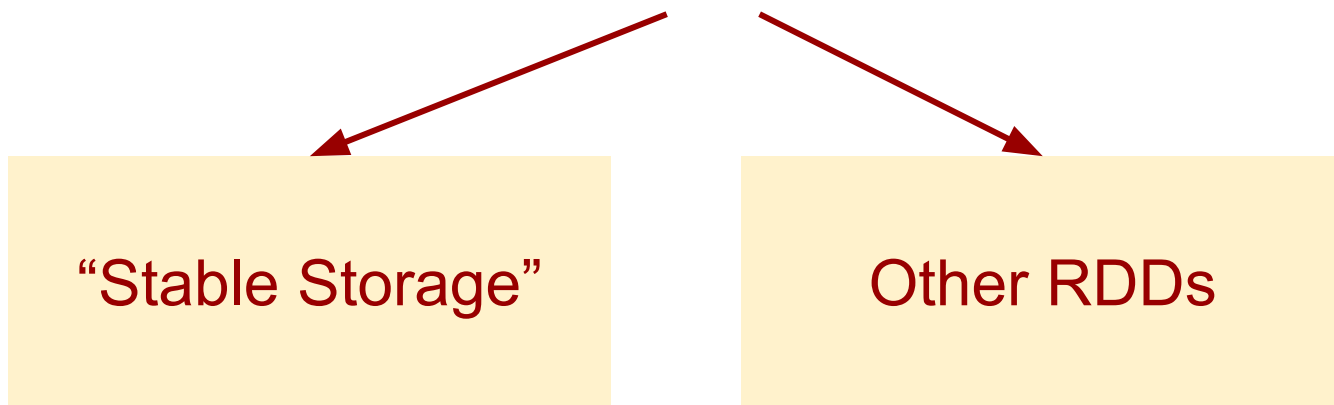
(DATA)

*transformation2*
from RDD2

# Spark's Big Idea

Resilient Distributed Datasets (RDDs) -- Read-only partitioned collection of records (like a DFS) but with a record of how the dataset was created as combination of *transformations* from other dataset(s).

- Enables rebuilding datasets on the fly.
- Intermediate datasets not stored on disk
  (and only in memory if needed and enough space)

⇒ Faster communication and I O

# The Big Idea

Resilient Distributed Datasets (RDDs) -- Read-only partitioned collection of records (like a DFS) but with a record of how the dataset was created as combination of *transformations* from other dataset(s).
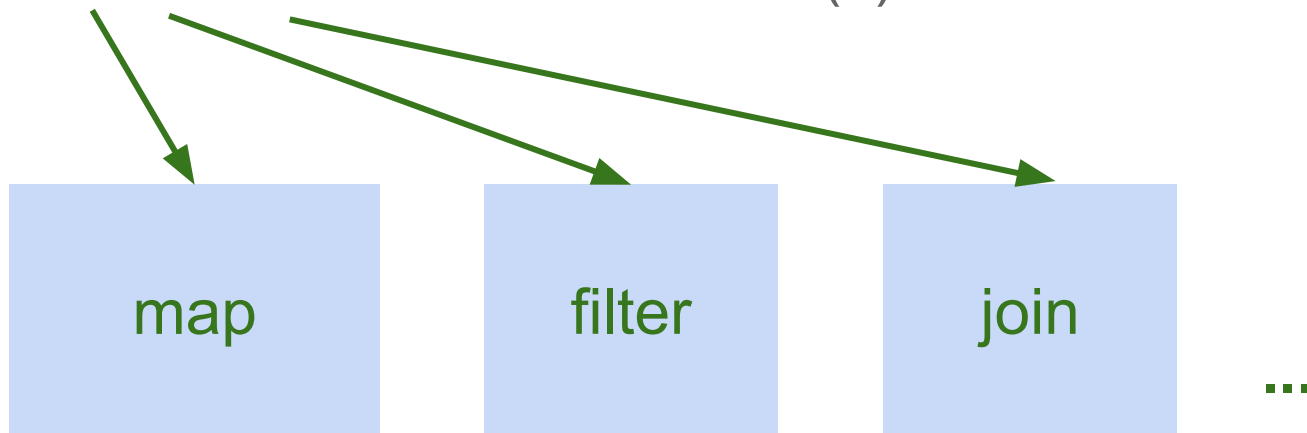
"Stable Storage"

Other RDDs

# The Big Idea

Resilient Distributed Datasets (RDDs) -- Read-only partitioned collection of records (like a DFS) but with a record of how the dataset was created as combination of *transformations* from other dataset(s).

| map | filter | join | ... |
|-----|--------|------|-----|

# Spark's Big Idea

Resilient Distributed Datasets (RDDs) -- Read-only partitioned collection of records (like a DFS) but with a record of how the dataset was created as combination of *transformations* from other dataset(s).
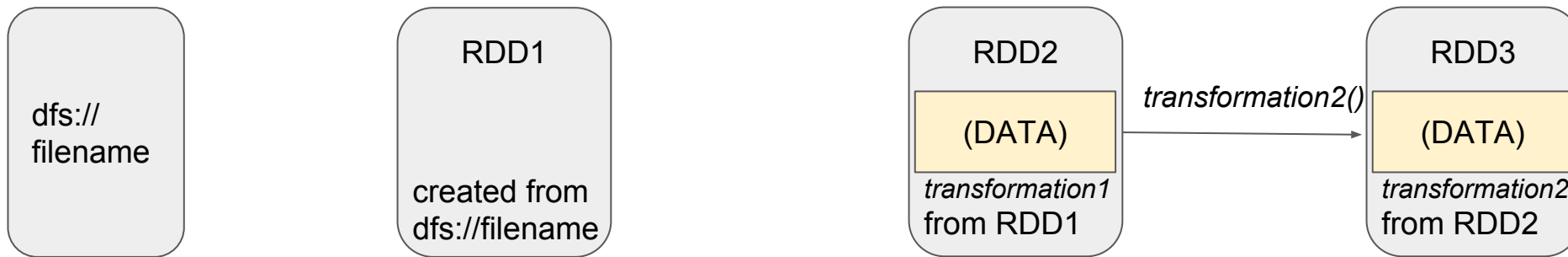
# Spark's Big Idea

Resilient Distributed Datasets (RDDs) -- Read-only partitioned collection of records (like a DFS) but with a record of how the dataset was created as combination of *transformations* from other dataset(s).

dfs:// filename

RDD1

created from dfs://filename

RDD2

*transformation1* from RDD1

*transformation2()*

RDD3

(DATA)

*transformation2* from RDD2

# Spark's Big Idea

Resilient Distributed Datasets (RDDs) -- Read-only partitioned collection of records (like a DFS) but with a record of how the dataset was created as combination of *transformations* from other dataset(s).
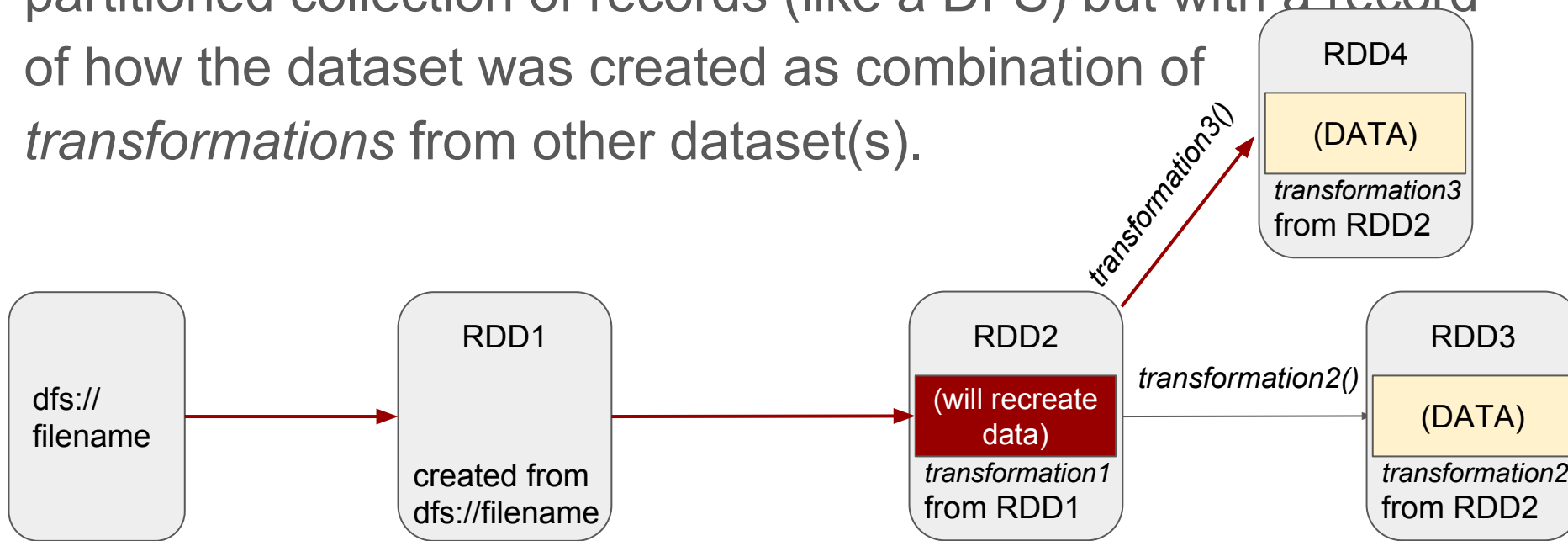
# Spark's Big Idea

Resilient Distributed Datasets (RDDs) -- Read-only partitioned collection of records (like a DFS) but with a record of how the dataset was created as combination of *transformations* from other dataset(s).

| RDD4 |
| --- |
| (DATA) |

*transformation3*
from RDD2

*transformation3()*

| dfs:// filename | → | RDD1<br><br>created from dfs://filename | → | RDD2<br>(will recreate data)<br>*transformation1* from RDD1 | *transformation2()* → | RDD3<br>(DATA)<br>*transformation2* from RDD2 |

# Original Transformations: RDD to RDD

| **Transformations** | | |
|---|---|---|
| $map(f : T \Rightarrow U)$ | : | $RDD[T] \Rightarrow RDD[U]$ |
| $filter(f : T \Rightarrow Bool)$ | : | $RDD[T] \Rightarrow RDD[T]$ |
| $flatMap(f : T \Rightarrow Seq[U])$ | : | $RDD[T] \Rightarrow RDD[U]$ |
| $sample(fraction : Float)$ | : | $RDD[T] \Rightarrow RDD[T]$  (Deterministic sampling) |
| $groupByKey()$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, Seq[V])]$ |
| $reduceByKey(f : (V, V) \Rightarrow V)$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, V)]$ |
| $union()$ | : | $(RDD[T], RDD[T]) \Rightarrow RDD[T]$ |
| $join()$ | : | $(RDD[(K, V)], RDD[(K, W)]) \Rightarrow RDD[(K, (V, W))]$ |
| $cogroup()$ | : | $(RDD[(K, V)], RDD[(K, W)]) \Rightarrow RDD[(K, (Seq[V], Seq[W]))]$ |
| $crossProduct()$ | : | $(RDD[T], RDD[U]) \Rightarrow RDD[(T, U)]$ |
| $mapValues(f : V \Rightarrow W)$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, W)]$  (Preserves partitioning) |
| $sort(c : Comparator[K])$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, V)]$ |
| $partitionBy(p : Partitioner[K])$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, V)]$ |

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

# Original Transformations: RDD to RDD

| Transformations | | | |
|---|---|---|---|
| | $map(f : T \Rightarrow U)$ | : | $RDD[T] \Rightarrow RDD[U]$ |
| | $filter(f : T \Rightarrow Bool)$ | : | $RDD[T] \Rightarrow RDD[T]$ |
| | $flatMap(f : T \Rightarrow Seq[U])$ | : | $RDD[T] \Rightarrow RDD[U]$ |
| | $sample(fraction : Float)$ | : | $RDD[T] \Rightarrow RDD[T]$  (Deterministic sampling) |
| | $groupByKey()$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, Seq[V])]$ |
| | $reduceByKey(f : (V, V) \Rightarrow V)$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, V)]$ |
| | $union()$ | : | $(RDD[T], RDD[T]) \Rightarrow RDD[T]$ |
| | $join()$ | : | $(RDD[(K, V)], RDD[(K, W)]) \Rightarrow RDD[(K, (V, W))]$ |
| | $cogroup()$ | : | $(RDD[(K, V)], RDD[(K, W)]) \Rightarrow RDD[(K, (Seq[V], Seq[W]))]$ |
| | $crossProduct()$ | : | $(RDD[T], RDD[U]) \Rightarrow RDD[(T, U)]$ |
| | $mapValues(f : V \Rightarrow W)$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, W)]$  (Preserves partitioning) |
| | $sort(c : Comparator[K])$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, V)]$ |
| | $partitionBy(p : Partitioner[K])$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, V)]$ |

*Multiple Records*

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

Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, Ion Stoica. "Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing.". *NSDI 2012*. April 2012.

# Original Transformations: RDD to RDD

| Transformations | | | |
|---|---|---|---|
| | $map(f : T \Rightarrow U)$ | : | $RDD[T] \Rightarrow RDD[U]$ |
| | $filter(f : T \Rightarrow Bool)$ | : | $RDD[T] \Rightarrow RDD[T]$ |
| | $flatMap(f : T \Rightarrow Seq[U])$ | : | $RDD[T] \Rightarrow RDD[U]$ |
| | $sample(fraction : Float)$ | : | $RDD[T] \Rightarrow RDD[T]$  (Deterministic sampling) |
| | $groupByKey()$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, Seq[V])]$ |
| | $reduceByKey(f : (V, V) \Rightarrow V)$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, V)]$ |
| | $union()$ | : | $(RDD[T], RDD[T]) \Rightarrow RDD[T]$ |
| | $join()$ | : | $(RDD[(K, V)], RDD[(K, W)]) \Rightarrow RDD[(K, (V, W))]$ |
| | $cogroup()$ | : | $(RDD[(K, V)], RDD[(K, W)]) \Rightarrow RDD[(K, (Seq[V], Seq[W]))]$ |
| | $crossProduct()$ | : | $(RDD[T], RDD[U]) \Rightarrow RDD[(T, U)]$ |
| | $mapValues(f : V \Rightarrow W)$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, W)]$  (Preserves partitioning) |
| | $sort(c : Comparator[K])$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, V)]$ |
| | $partitionBy(p : Partitioner[K])$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, V)]$ |

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

Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, Ion Stoica. "Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing.". *NSDI 2012*. April 2012.

# Original *Transformations*: RDD to RDD

| Transformations | | | |
|---|---|---|---|
| | $map(f : T \Rightarrow U)$ | : | $RDD[T] \Rightarrow RDD[U]$ |
| | $filter(f : T \Rightarrow Bool)$ | : | $RDD[T] \Rightarrow RDD[T]$ |
| | $flatMap(f : T \Rightarrow Seq[U])$ | : | $RDD[T] \Rightarrow RDD[U]$ |
| | $sample(fraction : Float)$ | : | $RDD[T] \Rightarrow RDD[T]$ (Deterministic sampling) |
| | $groupByKey()$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, Seq[V])]$ |
| | $reduceByKey(f : (V, V) \Rightarrow V)$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, V)]$ |
| | $union()$ | : | $(RDD[T], RDD[T]) \Rightarrow RDD[T]$ |
| | $join()$ | : | $(RDD[(K, V)], RDD[(K, W)]) \Rightarrow RDD[(K, (V, W))]$ |
| | $cogroup()$ | : | $(RDD[(K, V)], RDD[(K, W)]) \Rightarrow RDD[(K, (Seq[V], Seq[W]))]$ |
| | $crossProduct()$ | : | $(RDD[T], RDD[U]) \Rightarrow RDD[(T, U)]$ |
| | $mapValues(f : V \Rightarrow W)$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, W)]$ (Preserves partitioning) |
| | $sort(c : Comparator[K])$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, V)]$ |
| | $partitionBy(p : Partitioner[K])$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, V)]$ |

# Original *Actions*: RDD to Value, Object, or Storage

| Actions | | | |
|---|---|---|---|
| | $count()$ | : | $RDD[T] \Rightarrow Long$ |
| | $collect()$ | : | $RDD[T] \Rightarrow Seq[T]$ |
| | $reduce(f : (T, T) \Rightarrow T)$ | : | $RDD[T] \Rightarrow T$ |
| | $lookup(k : K)$ | : | $RDD[(K, V)] \Rightarrow Seq[V]$ (On hash/range partitioned RDDs) |
| | $save(path : String)$ | : | Outputs RDD to a storage system, *e.g.*, HDFS |

Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, Ion Stoica. "Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing.". *NSDI 2012*. April 2012.

# Current Transformations and Actions

http://spark.apache.org/docs/latest/rdd-programming-guide.html#transformations

common transformations: *filter, map, flatMap, reduceByKey, groupByKey*

http://spark.apache.org/docs/latest/rdd-programming-guide.html#actions

common actions: *collect, count, take*

# An Example

Count errors in a log file:

*TYPE    MESSAGE    TIME*

lines

filter.(_.startsWith("ERROR"))

errors

count()

Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, Ion Stoica. "Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing.". *NSDI 2012*. April 2012.

# An Example

Count errors in a log file:

*TYPE    MESSAGE    TIME*

Pseudocode:

```
lines = sc.textFile("dfs:...")
errors =
    lines.filter(_.startswith("ERROR"))
errors.count
```
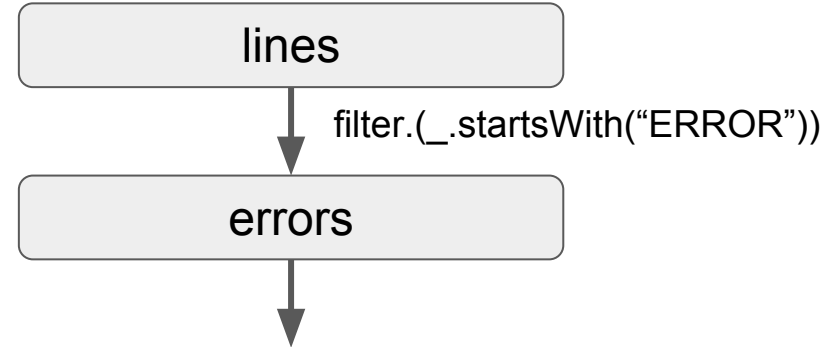
```
lines
```

filter.(_.startsWith("ERROR"))

```
errors
```

count()

Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, Ion Stoica. "Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing.". *NSDI 2012*. April 2012.

# An Example

Collect times of hdfs-related errors

*TYPE    MESSAGE    TIME*

```
Pseudocode:

lines = sc.textFile("dfs:...")
errors =
    lines.filter(_.startswith("ERROR"))
errors.persist
errors.count
…
```

```
lines
```
filter.(_.startsWith("ERROR"))
```
errors
```

Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, Ion Stoica. "Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing.". *NSDI 2012*. April 2012.

# An Example

Collect times of hdfs-related errors

*TYPE    MESSAGE    TIME*

```
Pseudocode:

lines = sc.textFile("dfs:...")
errors =
    lines.filter(_.startswith("ERROR"))
errors.persist
errors.count
…
```

**Persistance**

Can specify that an RDD "persists" in memory so other queries can use it.
Can specify a priority for persistance; lower priority => moves to disk, if needed, earlier

Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, Ion Stoica. "Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing.". *NSDI 2012*. April 2012.

# An Example

Collect times of hdfs-related errors

*TYPE      MESSAGE      TIME*

```
Pseudocode:

lines = sc.textFile("dfs:...")
errors =
    lines.filter(_.startswith("ERROR"))
errors.persist
errors.count
…
```

**Persistance**

Can specify that an RDD "persists" in memory so other queries can use it.
Can specify a priority for persistance; lower priority => moves to disk, if needed, earlier

[parameters for persist](parameters for persist)

Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, Ion Stoica. "Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing.". *NSDI 2012*. April 2012.
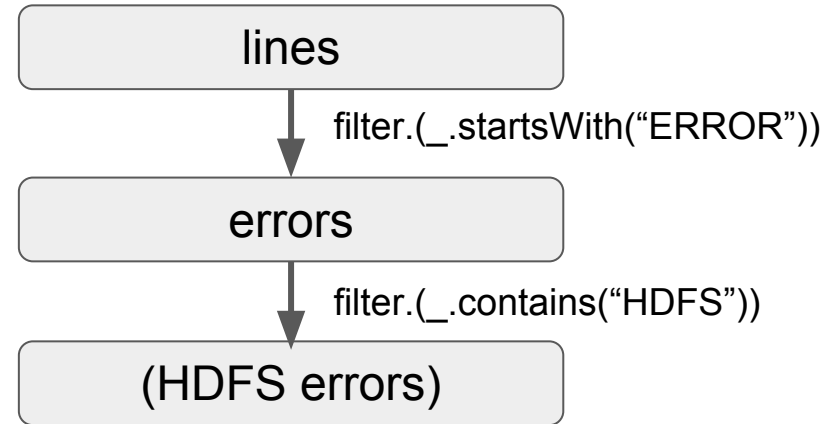
# An Example

Collect times of hdfs-related errors

*TYPE     MESSAGE    TIME*

```
Pseudocode:

lines = sc.textFile("dfs:...")
errors =
    lines.filter(_.startswith("ERROR"))
errors.persist
errors.count
errors.filter(_.contains("HDFS"))
    ...
```

Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, Ion Stoica. "Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing.". *NSDI 2012*. April 2012.
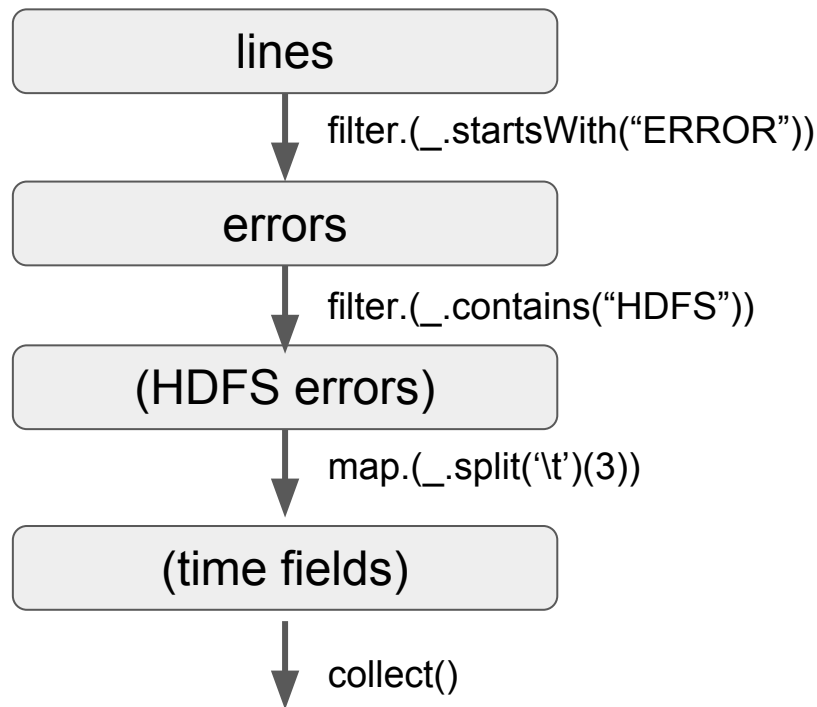
# An Example

Collect times of hdfs-related errors

*TYPE      MESSAGE      TIME*

```
Pseudocode:

lines = sc.textFile("dfs:...")
errors =
    lines.filter(_.startswith("ERROR"))
errors.persist
errors.count
errors.filter(_.contains("HDFS"))
    .map(_split('\t')(3))
    .collect()
```

```
lines
        │
        │ filter.(_.startsWith("ERROR"))
        ▼
      errors
        │
        │ filter.(_.contains("HDFS"))
        ▼
   (HDFS errors)
        │
        │ map.(_.split('\t')(3))
        ▼
   (time fields)
        │
        │ collect()
        ▼
```

Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, Ion Stoica. "Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing.". *NSDI 2012*. April 2012.
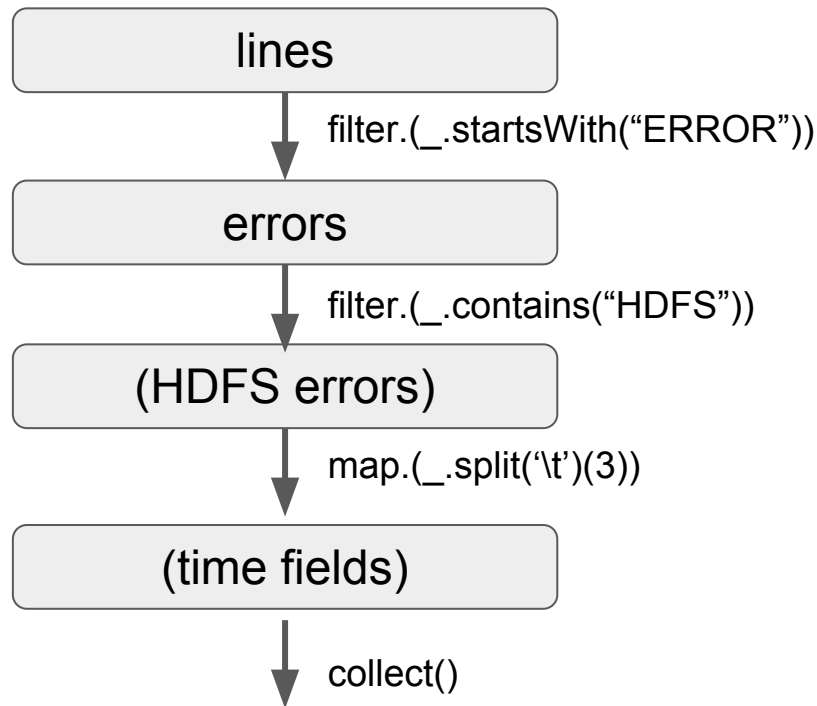
# An Example

Collect times of hdfs-related errors
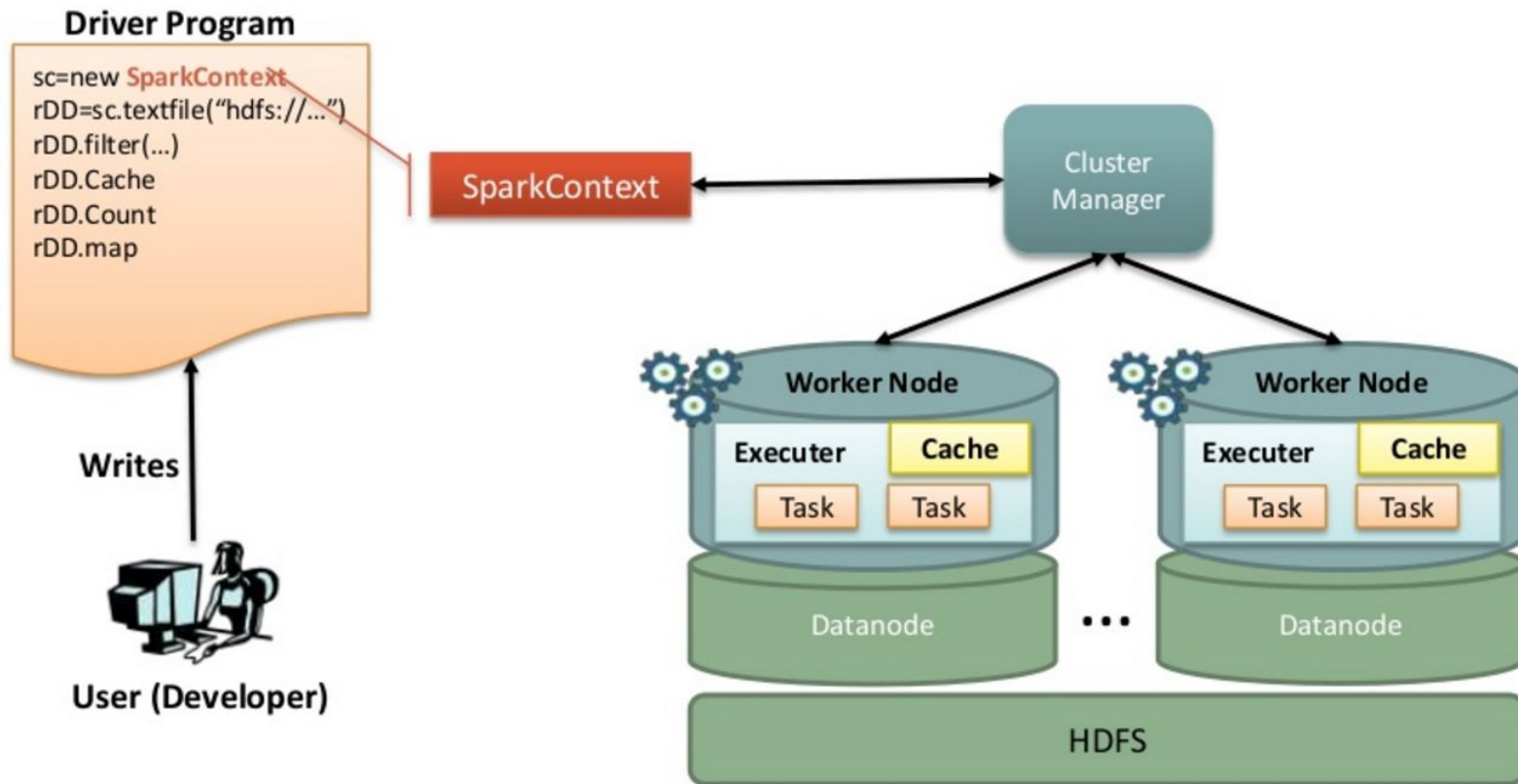
*TYPE     MESSAGE     TIME*

```
Pseudocode:

lines = sc.textFile("dfs:...")
errors =
     lines.filter(_.startswith("ERROR"))
errors.persist
errors.count
errors.filter(_.contains("HDFS"))
     .map(_split('\t')(3))
     .collect()
```

**Functional Programming**

```mermaid
lines
  | filter.(_.startsWith("ERROR"))
errors
  | filter.(_.contains("HDFS"))
(HDFS errors)
  | map.(_.split('\t')(3))
(time fields)
  | collect()
```

lines

filter.(_.startsWith("ERROR"))

errors

filter.(_.contains("HDFS"))

(HDFS errors)

map.(_.split('\t')(3))

(time fields)

collect()

Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, Ion Stoica. "Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing.". *NSDI 2012*. April 2012.

# The Spark Programming Model



Gupta, Manish. Lightening Fast Big Data Analytics using Apache Spark. *UniCom 2014.*

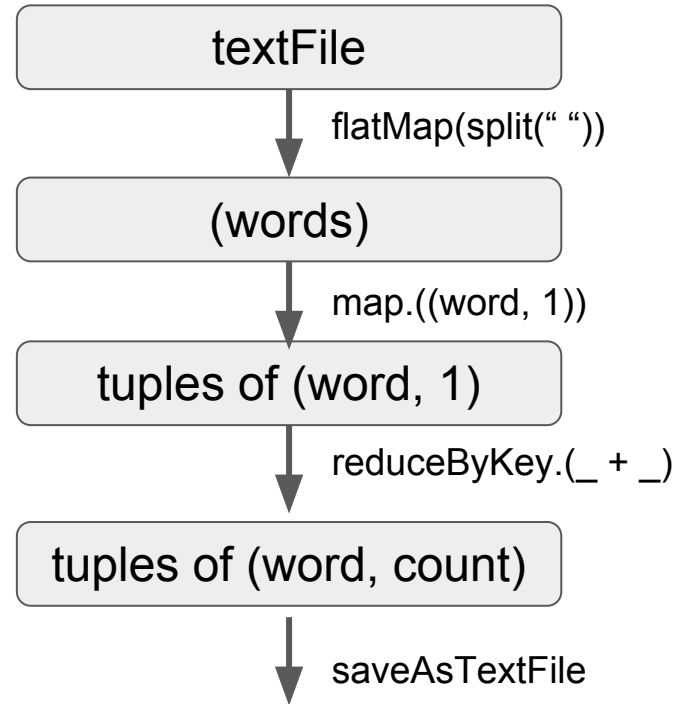# An Example

## Word Count

textFile

# An Example

## Word Count

```scala
Scala:

val textFile =
    sc.textFile("hdfs://...")
val counts = textFile
    .flatMap(line => line.split(" "))
    .map(word => (word, 1))
    .reduceByKey(_ + _)
counts.saveAsTextFile("hdfs://...")
```
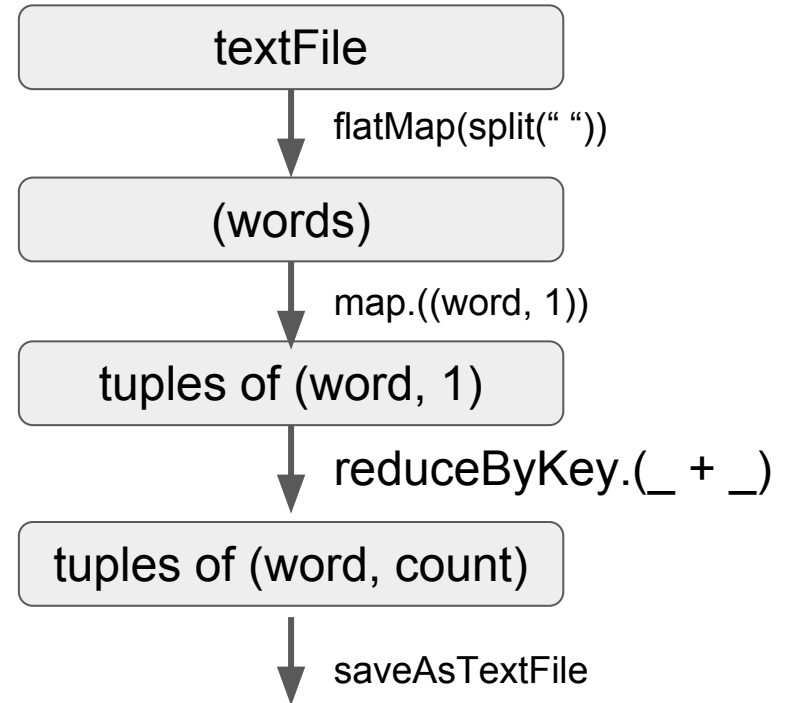
```
┌─────────────────────┐
│      textFile       │
└─────────────────────┘
           │  flatMap(split(" "))
           ▼
┌─────────────────────┐
│      (words)        │
└─────────────────────┘
           │  map.((word, 1))
           ▼
┌─────────────────────┐
│  tuples of (word, 1)│
└─────────────────────┘
           │  reduceByKey.(_ + _)
           ▼
┌─────────────────────┐
│tuples of (word, count)│
└─────────────────────┘
           │  saveAsTextFile
           ▼
```

# An Example

## Word Count

```
Python:

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

textFile

↓ flatMap(split(" "))

(words)

↓ map.((word, 1))

tuples of (word, 1)

↓ reduceByKey.(_ + _)

tuples of (word, count)

↓ saveAsTextFile

Apache Spark Examples
http://spark.apache.org/examples.html

# PySpark Demo



https://data.worldbank.org/data-catalog/poverty-and-equity-database

# Lazy Evaluation

Spark waits to **load data** and **execute transformations** until necessary -- *lazy*
Spark tries to complete **actions** as immediately as possible -- *eager*

Why?

- Only executes what is necessary to achieve action.

- Can optimize the complete *chain of operations* to reduce communication

# Lazy Evaluation

Spark waits to *load data* and *execute transformations* until necessary -- ***lazy***
Spark tries to complete actions as quickly as possible -- ***eager***

## Why?

- Only executes what is necessary to achieve action.

- Can optimize the complete *chain of operations* to reduce communication

e.g.

```
rdd.map(lambda r: r[1]*r[3]).take(5)  #only executes map for five records

rdd.filter(lambda r: "ERROR" in r[0]).map(lambda r: r[1]*r[3])
                        #only passes through the data once
```

# Broadcast Variables

Read-only objects can be shared across all nodes.

Broadcast variable is a wrapper: access object with .value

```
Python:

filterWords = ['one', 'two', 'three', 'four', …]
fwBC = sc.broadcast(set(filterWords))
```

# Broadcast Variables

Read-only objects can be shared across all nodes.
    Broadcast variable is a wrapper: access object with .value

```python
Python:

filterWords = ['one', 'two', 'three', 'four', …]
fwBC = sc.broadcast(set(filterWords))

textFile = sc.textFile("hdfs:...")
counts = textFile
    .map(lambda line: line.split(" "))
    .filter(lambda words: len(set(words) and word in fwBC.value) > 0)
    .flatMap(lambda word: (word, 1))
    .reduceByKey(lambda a, b: a + b)
counts.saveAsTextFile("hdfs:...")
```

# Accumulators

Write-only objects that keep a running aggregation

Default Accumulator assumes sum function

```
initialValue = 0
sumAcc = sc.accumulator(initialValue)
rdd.foreach(lambda i: sumAcc.add(i))
print(sumAcc.value)
```

# Accumulators

Write-only objects that keep a running aggregation

Default Accumulator assumes sum function

Custom Accumulator: Inherit (AccumulatorParam) as class and override methods

```python
initialValue = 0
sumAcc = sc.accumulator(initialValue)
rdd.foreeach(lambda i: sumAcc.add(i))
print(minAcc.value)


class MinAccum(AccumulatorParam):
    def zero(self, zeroValue = np.inf):#overwrite this
        return zeroValue
    def addInPlace(self, v1, v2):#overwrite this
        return min(v1, v2)
minAcc = sc.accumulator(np.inf, minAccum())
rdd.foreeach(lambda i: minAcc.add(i))
print(minAcc.value)
```

# Spark Overview

- RDD provides full recovery by backing up transformations from stable storage rather than backing up the data itself.

- RDDs, which are immutable, can be stored in memory and thus are often much faster.

- Functional programming is used to define transformation and actions on RDDs.

# Spark Overview

- RDD provides full recovery by backing up transformations from stable storage rather than backing up the data itself.

- RDDs, which are immutable, can be stored in memory and thus are often much faster.

- Functional programming is used to define transformation and actions on RDDs.

- Still need Hadoop (or some DFS) to hold original or resulting data efficiently and reliably.

- Lazy evaluation enables optimizing chain of operations.

- Memory across Spark cluster should be large enough to hold entire dataset to fully leverage speed.

  - MapReduce may still be more cost-effective for very large data that does not fit in memory.