



# 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

NSDI 2012





# Contents

---

1

**Motivation**

2

**Resilient Distributed Datasets (RDDs)**

3

**Spark programming**

4

**Evaluation**

5

**Conclusion**

# Motivation

---

**MapReduce** greatly simplifies “big data” analysis on large, unreliable clusters

Problems:

- Iterative algorithm & graph processing
- Interactive data mining

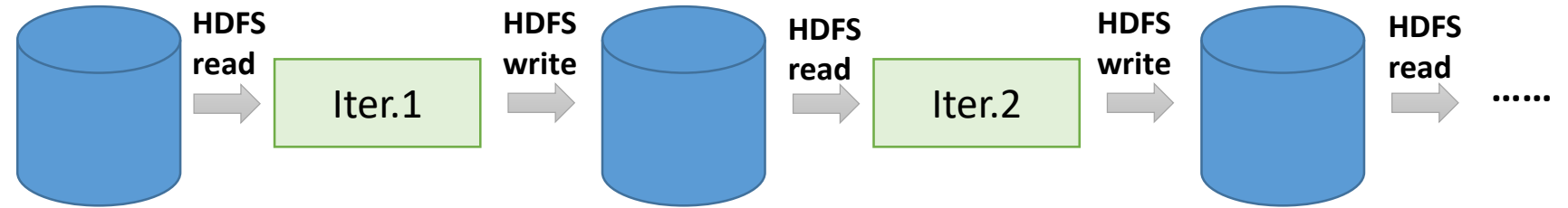


Cause:

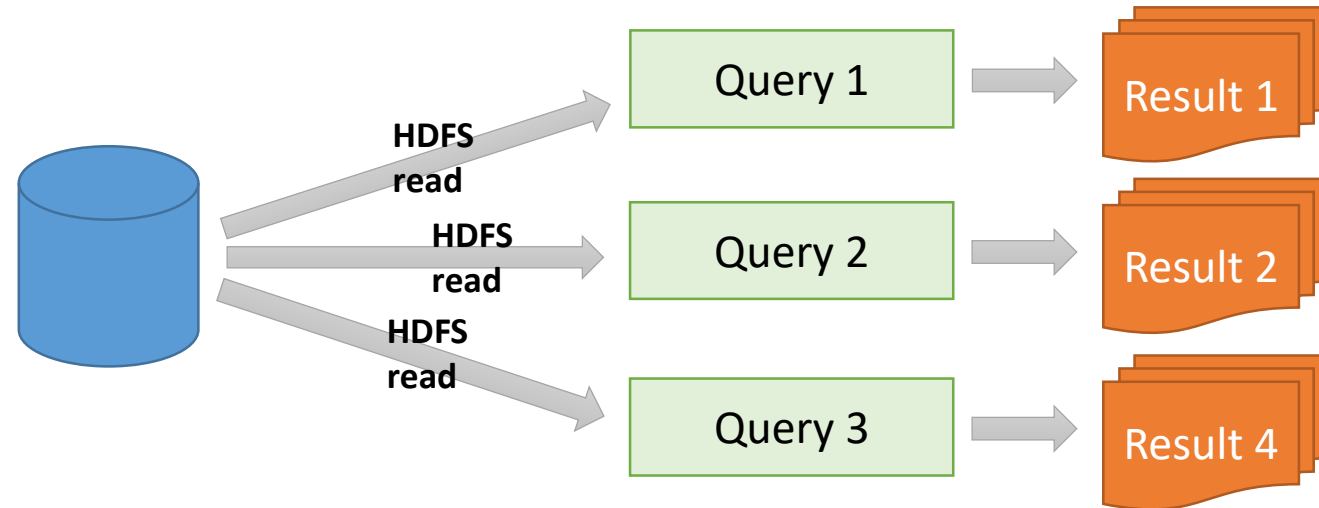
Inefficient **data sharing**

# Motivation

Iterative algorithm



Interactive data mining



Replication & Disk I/O

# Motivation



## GOAL DESIGN

- Efficient in-memory cluster computing (iterative algorithm & interactive data mining)
- Fault-tolerant
- General-purpose abstraction
- Explicitly persist intermediate results in memory
- Control partitioning to optimize data placement

**Replication & Disk I/O**

# Resilient Distributed Datasets (RDDs)

---

An RDD is a **read-only, partitioned** collection of records, which provides a **restricted form** of shared memory.

Restricted conditions:

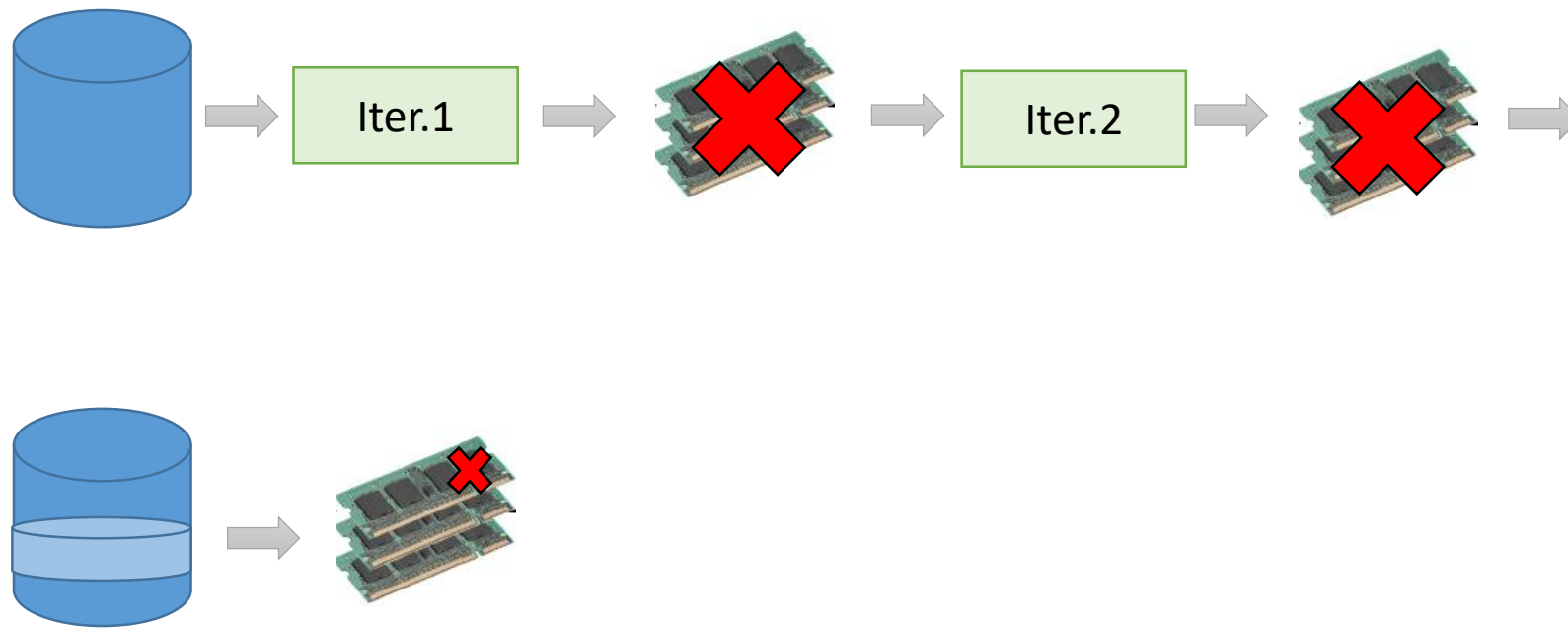
- Read-only, partitioned
- Built through **coarse-grained transformations** (e.g., map, filter and join)

Powerful properties:

- Generality of RDDs — Abstraction
- Running backup copies of tasks from the slow nodes
- Automatic schedule based on data locality
- Keeping intermediate data in memory
- Efficient fault recovery using lineage

# Resilient Distributed Datasets (RDDs)

## Fault recovery



# Spark programming

Language-integrated API in **Scala** language

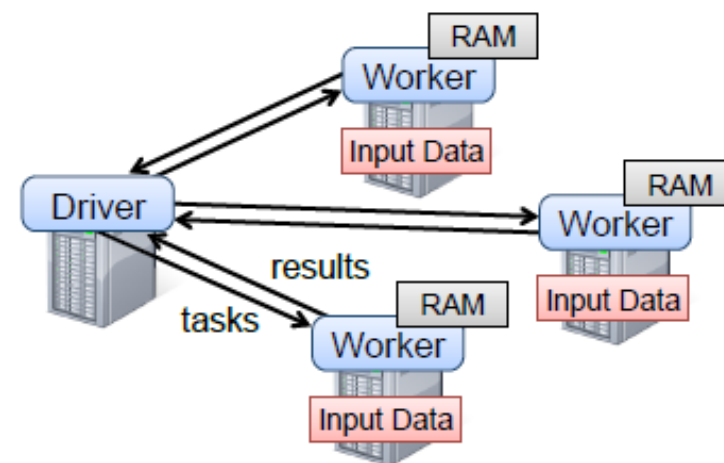
Runtime: a driver connects to a cluster of workers

Two type of operations:

- Transformations — create RDDs
- Actions — compute and output results

Control partitioning

Control persistence (storage in RAM, on disk, etc)





# Spark programming

Transformations	<i>map</i> ( <i>f</i> : <i>T</i> ⇒ <i>U</i> )	: RDD[ <i>T</i> ] ⇒ RDD[ <i>U</i> ]
	<i>filter</i> ( <i>f</i> : <i>T</i> ⇒ Bool)	: RDD[ <i>T</i> ] ⇒ RDD[ <i>T</i> ]
	<i>flatMap</i> ( <i>f</i> : <i>T</i> ⇒ Seq[ <i>U</i> ])	: RDD[ <i>T</i> ] ⇒ RDD[ <i>U</i> ]
	<i>sample</i> ( <i>fraction</i> : Float)	: RDD[ <i>T</i> ] ⇒ RDD[ <i>T</i> ] (Deterministic sampling)
	<i>groupByKey</i> ()	: RDD[( <i>K</i> , <i>V</i> )] ⇒ RDD[( <i>K</i> , Seq[ <i>V</i> ])]
	<i>reduceByKey</i> ( <i>f</i> : ( <i>V</i> , <i>V</i> ) ⇒ <i>V</i> )	: RDD[( <i>K</i> , <i>V</i> )] ⇒ RDD[( <i>K</i> , <i>V</i> )]
	<i>union</i> ()	: (RDD[ <i>T</i> ], RDD[ <i>T</i> ]) ⇒ RDD[ <i>T</i> ]
	<i>join</i> ()	: (RDD[( <i>K</i> , <i>V</i> )], RDD[( <i>K</i> , <i>W</i> )]) ⇒ RDD[( <i>K</i> , ( <i>V</i> , <i>W</i> ))]
	<i>cogroup</i> ()	: (RDD[( <i>K</i> , <i>V</i> )], RDD[( <i>K</i> , <i>W</i> )]) ⇒ RDD[( <i>K</i> , (Seq[ <i>V</i> ], Seq[ <i>W</i> ]))]
	<i>crossProduct</i> ()	: (RDD[ <i>T</i> ], RDD[ <i>U</i> ]) ⇒ RDD[( <i>T</i> , <i>U</i> )]
	<i>mapValues</i> ( <i>f</i> : <i>V</i> ⇒ <i>W</i> )	: RDD[( <i>K</i> , <i>V</i> )] ⇒ RDD[( <i>K</i> , <i>W</i> )] (Preserves partitioning)
	<i>sort</i> ( <i>c</i> : Comparator[ <i>K</i> ])	: RDD[( <i>K</i> , <i>V</i> )] ⇒ RDD[( <i>K</i> , <i>V</i> )]
	<i>partitionBy</i> ( <i>p</i> : Partitioner[ <i>K</i> ])	: RDD[( <i>K</i> , <i>V</i> )] ⇒ RDD[( <i>K</i> , <i>V</i> )]
Actions	<i>count</i> ()	: RDD[ <i>T</i> ] ⇒ Long
	<i>collect</i> ()	: RDD[ <i>T</i> ] ⇒ Seq[ <i>T</i> ]
	<i>reduce</i> ( <i>f</i> : ( <i>T</i> , <i>T</i> ) ⇒ <i>T</i> )	: RDD[ <i>T</i> ] ⇒ <i>T</i>
	<i>lookup</i> ( <i>k</i> : <i>K</i> )	: RDD[( <i>K</i> , <i>V</i> )] ⇒ Seq[ <i>V</i> ] (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. Seq[*T*] denotes a sequence of elements of type *T*.

# Spark programming

## Example 1: Log Mining

lines = spark.textFile("hdfs://...")

Base RDD

errors = lines.filter(\_.startsWith("ERROR"))

Transformed RDD

messages = errors.map(\_.split('\t')(2))

messages.persist()

Declare RDD persist in memory

messages.count()

messages.filter(\_.contains("MySQL")).count()

Action

# Spark programming

## Example 2: PageRank

```
val links = spark.textFile(...).map(...).persist() // RDD of (URL, neighbors) pairs
```

```
var ranks = // RDD of (URL, rank) pairs
```

```
for (i <- 1 to ITERATIONS) {
```

```
  // Build an RDD of (targetURL, float) pairs
```

```
  // with the contributions sent by each page
```

```
  val contribs = links.join(ranks).flatMap {
```

```
    (url, (links, rank)) =>
```

```
      links.map(dest => (dest, rank/links.size))
```

```
  }
```

```
  // Sum contributions by URL and get new ranks
```

```
  ranks = contribs.reduceByKey((x,y) => x+y)
```

```
    .mapValues(sum =>  $\frac{a}{N} + (1-a)*sum$ )
```

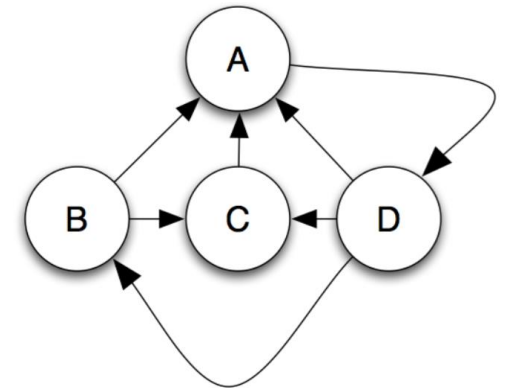
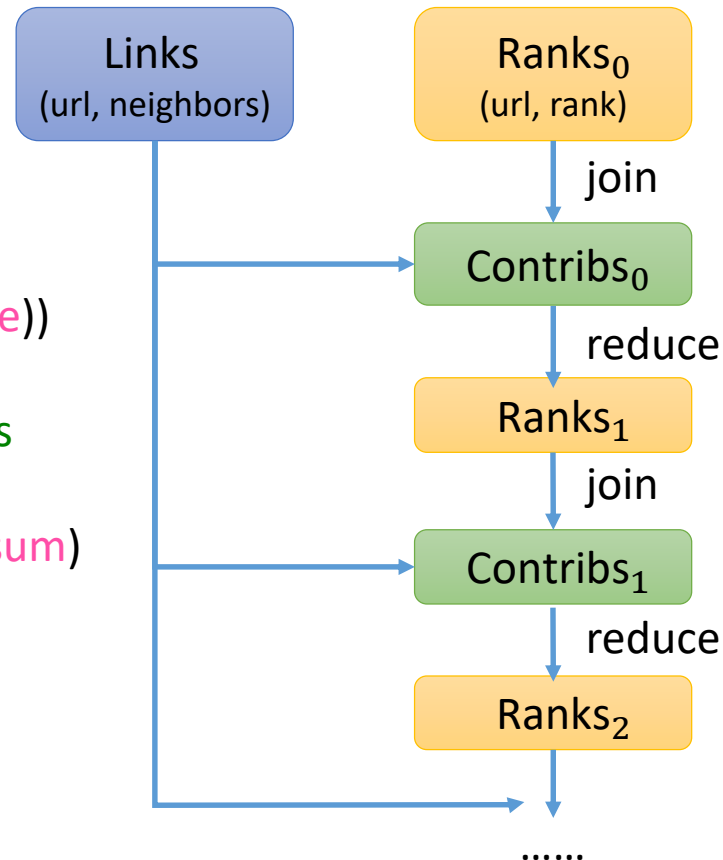
```
}
```

### Optimizing RDD placement

➤ hash-partition the links and ranks by URL

➤ custom Partitioner using partitionBy

e.g., links = spark.textFile(...).map(...) .partitionBy(myPartFunc).persist()



# Spark programming

## Implementation

- Representing RDDs

Interface exposes five pieces of information:

- a set of **partitions**
- a set of **dependencies** on parent RDDs
- a **function for computing** the dataset based on its parents
- the **metadata** about its partitioning scheme
- the **data placement**

- Job Scheduling

Delay scheduling

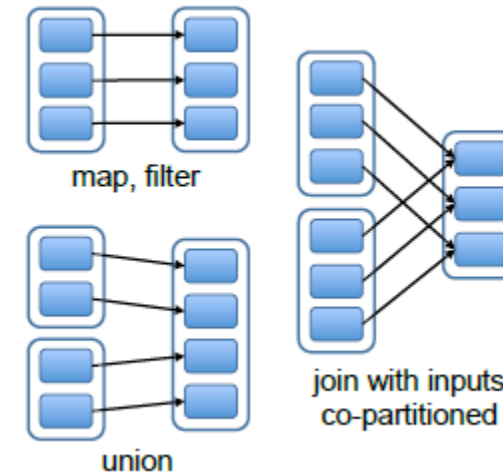
Schedule based on data locality

- Memory Management

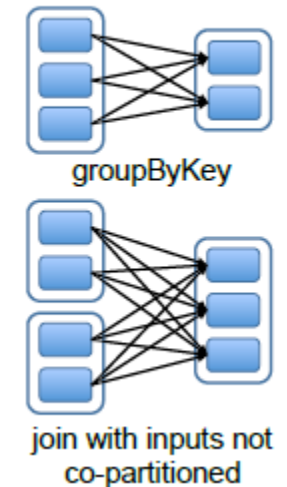
- in-memory storage as **deserialized Java objects**
- in-memory storage as **serialized data**
- on-**disk storage**

- Support for Checkpointing

Narrow Dependencies:



Wide Dependencies:



limited memory available — **LRU**

Persist — **REPLICATE flag**

# Evaluation

## Iterative Machine Learning Applications

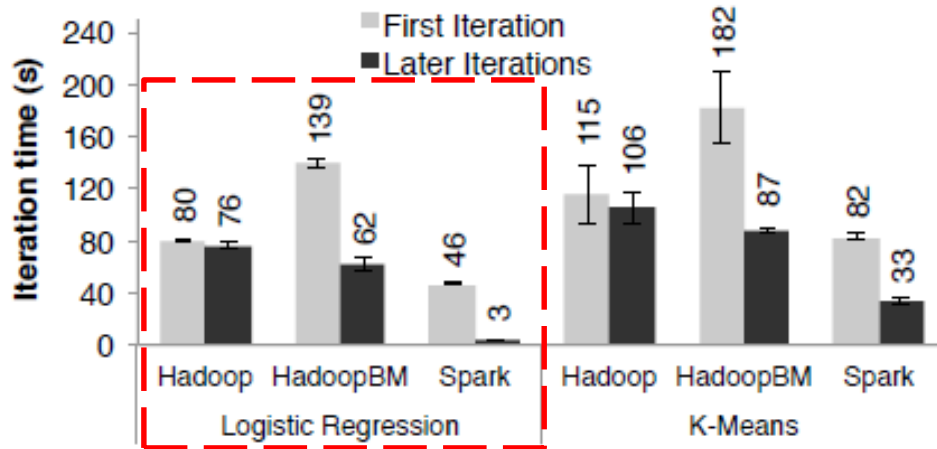


Figure 7: Duration of the first and later iterations in Hadoop, HadoopBinMem and Spark for logistic regression and k-means using 100 GB of data on a 100-node cluster.

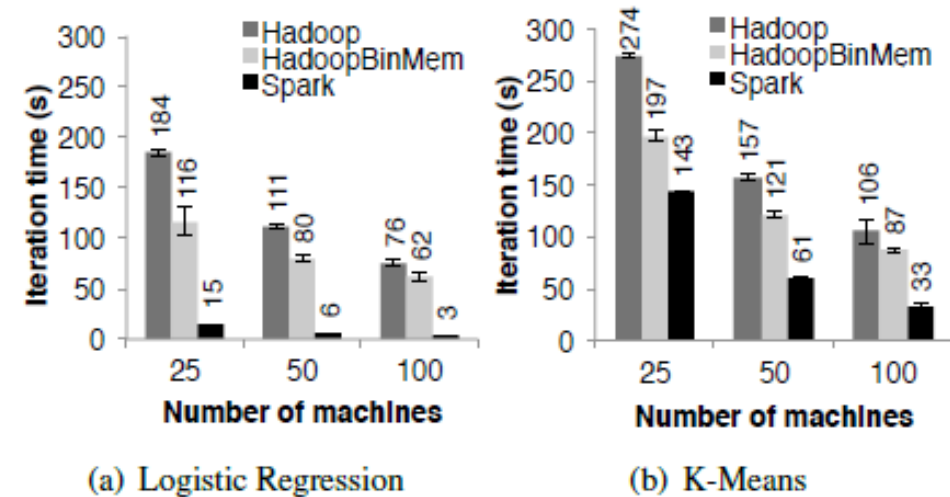
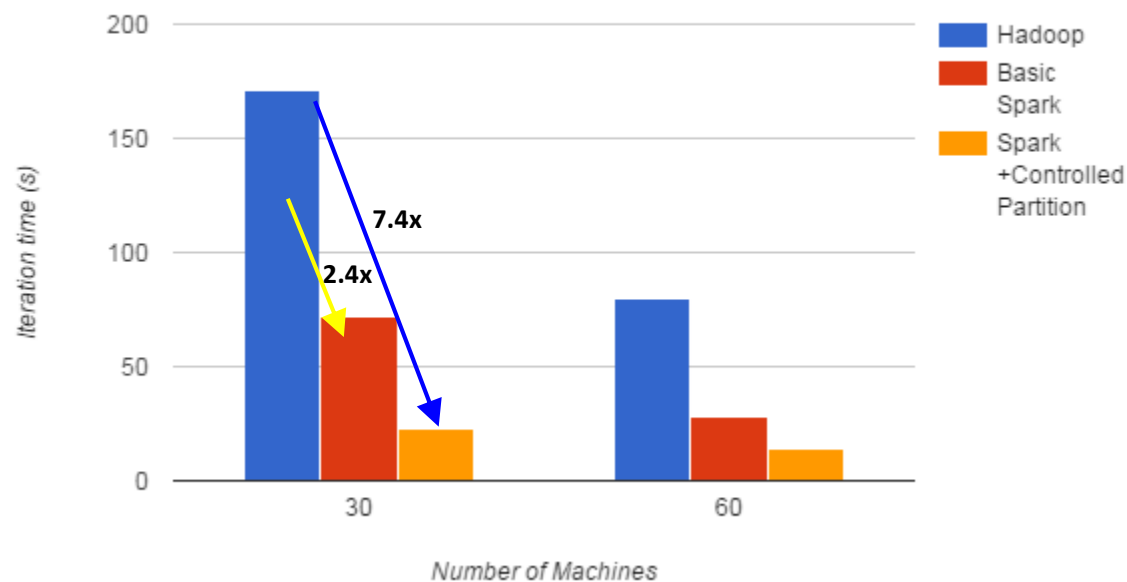


Figure 8: Running times for iterations after the first in Hadoop, HadoopBinMem, and Spark. The jobs all processed 100 GB.

- ✓ Hadoop: The Hadoop 0.20.2 stable release.
- ✓ HadoopBinMem: A Hadoop deployment that converts the input data into a **low-overhead binary format** in the first iteration and stores it in an **in-memory HDFS** instance.
- ✓ Spark: Our implementation of RDDs.

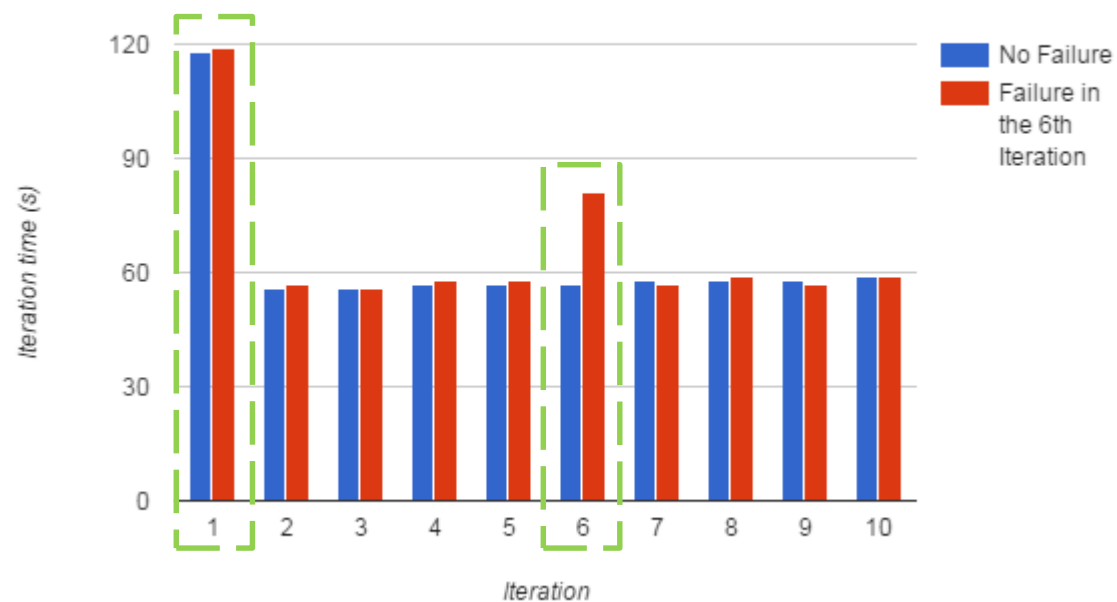
# Evaluation

## PageRank



Performance of PageRank on Hadoop and Spark

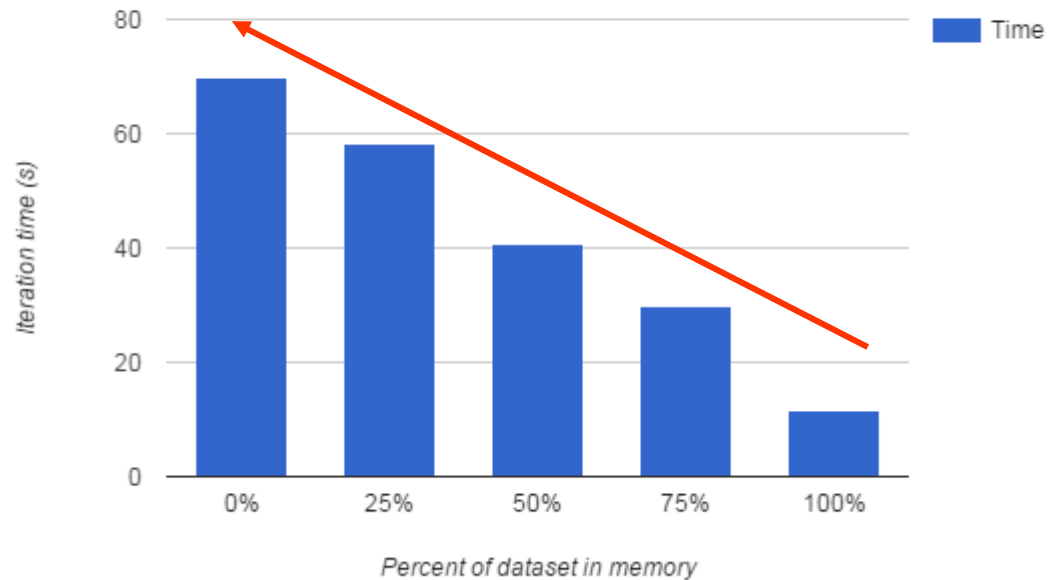
## Fault Recovery



Iteration times for k-means in presence of a failure. One machine was killed at the start of the 6th iteration, resulting in partial reconstruction of an RDD using lineage

# Evaluation

## Behavior with Insufficient Memory



**Performance of logistic regression using 100GB data on 25 machines with varying amounts of data in memory**



# Conclusion

---

- Resilient distributed datasets (RDDs) is an efficient, general-purpose and fault-tolerant abstraction for sharing data in cluster applications.
- RDDs offer an API based on coarse-grained transformations that lets them recover data efficiently using lineage.
- Implement RDDs in Spark that outperforms Hadoop by up to 20x in iterative applications.