# Intro to Apache Spark

http://databricks.com/

download slides:
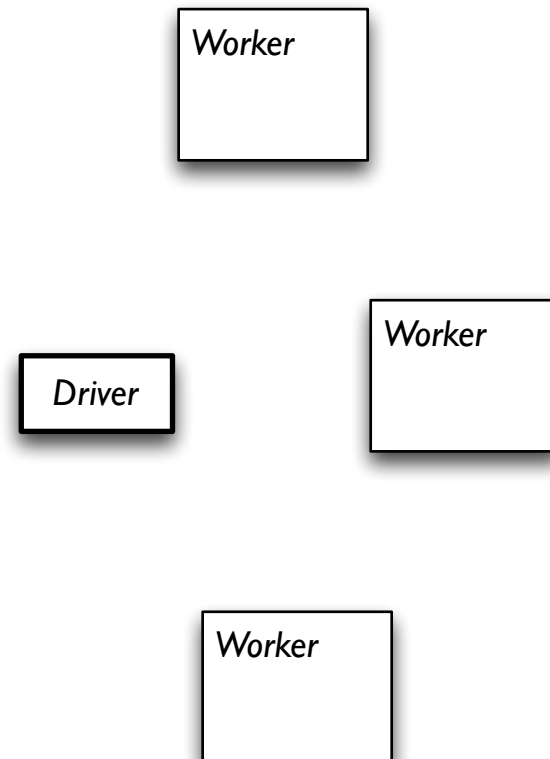training.databricks.com/workshop/itas_workshop.pdf

## Spark Deconstructed: *Log Mining Example*

```scala
// load error messages from a log into memory
// then interactively search for various patterns
// https://gist.github.com/ceteri/8ae5b9509a08c08a1132

// base RDD
val lines = sc.textFile("hdfs://...")

// transformed RDDs
val errors = lines.filter(_.startsWith("ERROR"))
val messages = errors.map(_.split("\t")).map(r => r(1))
messages.cache()

// action 1
messages.filter(_.contains("mysql")).count()

// action 2
messages.filter(_.contains("php")).count()
```

**Spark Deconstructed:** *Log Mining Example*

We start with Spark running on a cluster… submitting code to be evaluated on it:

Worker

Worker

Driver

Worker

# Spark Deconstructed: *Log Mining Example*

```scala
// base RDD
val lines = sc.textFile("hdfs://...")

// transformed RDDs
val errors = lines.filter(_.startsWith("ERROR"))
val messages = errors.map(_.split("\t")).map(r => r(1))
messages.cache()

// action 1
messages.filter(_.contains("mysql")).count()

// action 2
messages.filter(_.contains("php")).count()
```

discussing the other part

19

**Spark Deconstructed:** *Log Mining Example*

At this point, take a look at the transformed
RDD *operator graph*:

```
scala> messages.toDebugString
res5: String =
MappedRDD[4] at map at <console>:16 (3 partitions)
  MappedRDD[3] at map at <console>:16 (3 partitions)
    FilteredRDD[2] at filter at <console>:14 (3 partitions)
      MappedRDD[1] at textFile at <console>:12 (3 partitions)
        HadoopRDD[0] at textFile at <console>:12 (3 partitions)
```

# Spark Deconstructed: *Log Mining Example*

```scala
// base RDD
val lines = sc.textFile("hdfs://...")

// transformed RDDs
val errors = lines.filter(_.startsWith("ERROR"))
val messages = errors.map(_.split("\t")).map(r => r(1))
messages.cache()

// action 1
messages.filter(_.contains("mysql")).count()

// action 2
messages.filter(_.contains("php")).count()
```

discussing the other part

Worker

Worker

Driver

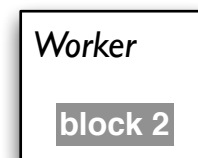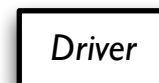Worker

# Spark Deconstructed: *Log Mining Example*

```scala
// base RDD
val lines = sc.textFile("hdfs://...")

// transformed RDDs
val errors = lines.filter(_.startsWith("ERROR"))
val messages = errors.map(_.split("\t")).map(r => r(1))
messages.cache()

// action 1
messages.filter(_.contains("mysql")).count()

// action 2
messages.filter(_.contains("php")).count()
```

discussing the other part

Worker

block 1

Worker

block 2

Driver

Worker

block 3
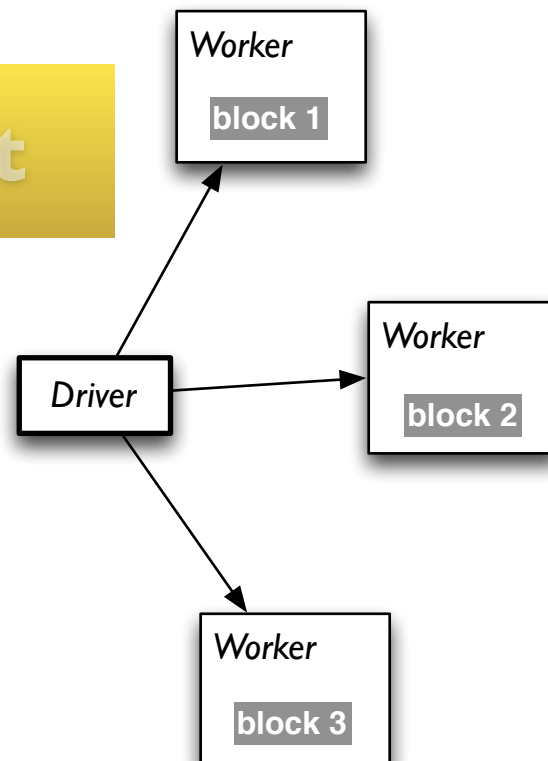
22

# Spark Deconstructed: *Log Mining Example*

```scala
// base RDD
val lines = sc.textFile("hdfs://...")

// transformed RDDs
val errors = lines.filter(_.startsWith("ERROR"))
val messages = errors.map(_.split("\t")).map(r => r(1))
messages.cache()

// action 1
messages.filter(_.contains("mysql")).count()

// action 2
messages.filter(_.contains("php")).count()
```

discussing the other part

Worker

block 1

Worker

block 2

Driver

Worker

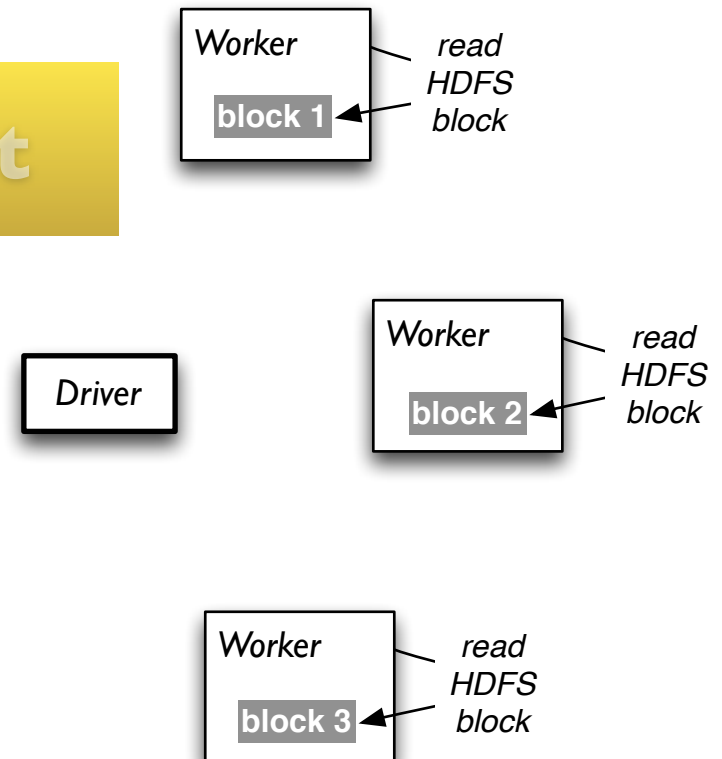block 3

# Spark Deconstructed: *Log Mining Example*

```scala
// base RDD
val lines = sc.textFile("hdfs://...")


// transformed RDDs
val errors = lines.filter(_.startsWith("ERROR"))
val messages = errors.map(_.split("\t")).map(r => r(1))
messages.cache()


// action 1
messages.filter(_.contains("mysql")).count()


// action 2
messages.filter(_.contains("php")).count()
```

discussing the other part

```
Worker
  block 1          read
                   HDFS
                   block
```

```
Driver
```

```
Worker
  block 2          read
                   HDFS
                   block
```

```
Worker
  block 3          read
                   HDFS
                   block
```
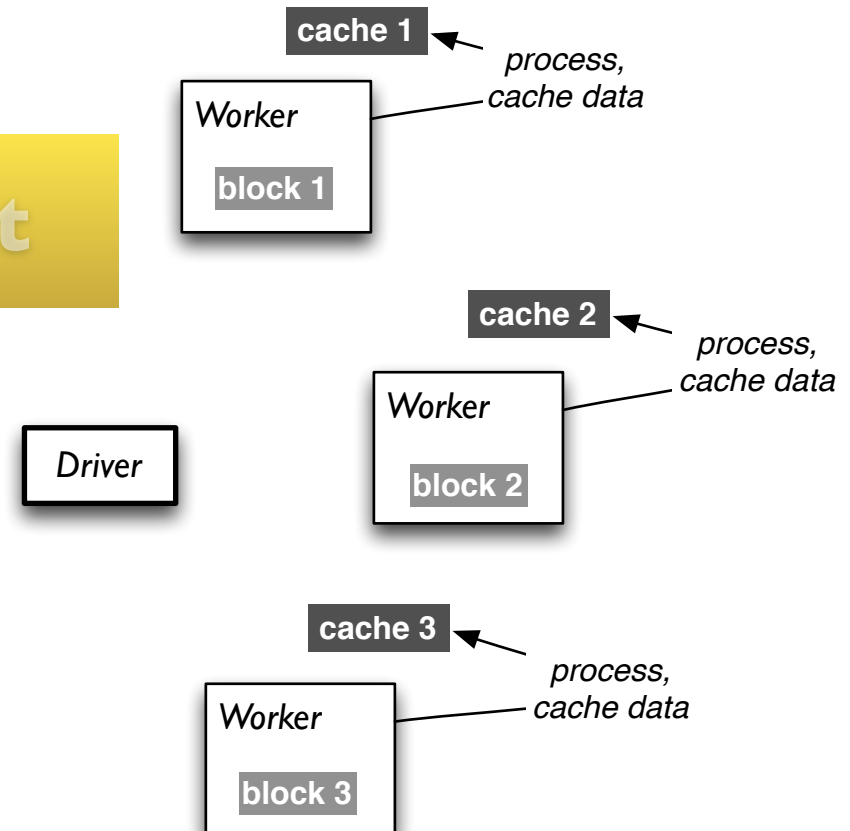
# Spark Deconstructed: *Log Mining Example*

```scala
// base RDD
val lines = sc.textFile("hdfs://...")

// transformed RDDs
val errors = lines.filter(_.startsWith("ERROR"))
val messages = errors.map(_.split("\t")).map(r => r(1))
messages.cache()

// action 1
messages.filter(_.contains("mysql")).count()

// action 2
messages.filter(_.contains("php")).count()
```

discussing the other part

cache 1

process, cache data

Worker

block 1

cache 2

process, cache data

Worker

block 2

Driver

cache 3

process, cache data

Worker

block 3
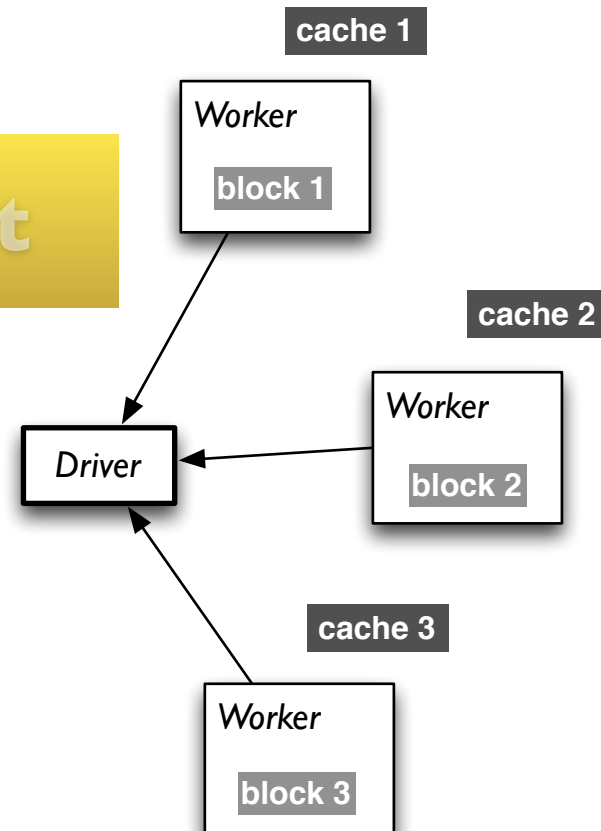
25

# Spark Deconstructed: *Log Mining Example*

```scala
// base RDD
val lines = sc.textFile("hdfs://...")

// transformed RDDs
val errors = lines.filter(_.startsWith("ERROR"))
val messages = errors.map(_.split("\t")).map(r => r(1))
messages.cache()

// action 1
messages.filter(_.contains("mysql")).count()

// action 2
messages.filter(_.contains("php")).count()
```

discussing the other part

cache 1

**Worker**

block 1

cache 2

**Worker**

block 2

Driver

cache 3

**Worker**

block 3
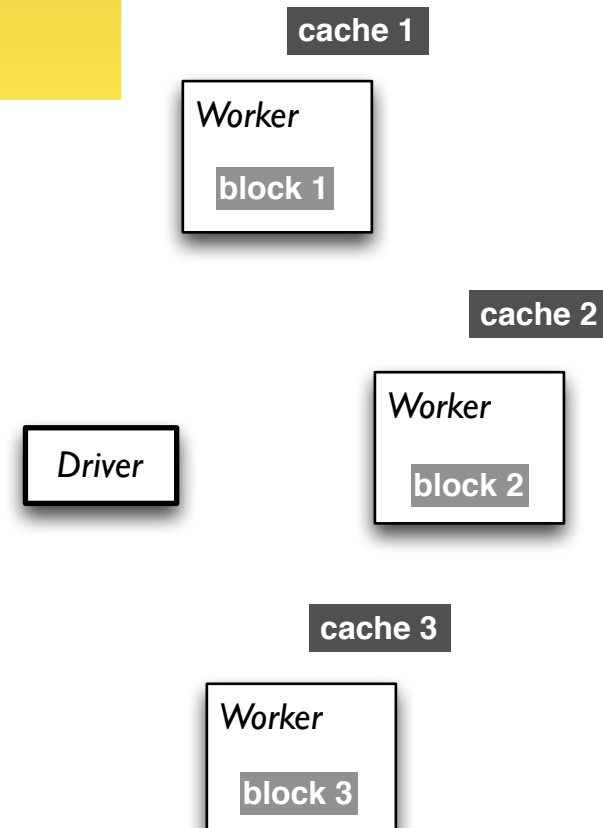
26

# Spark Deconstructed: *Log Mining Example*

```scala
// base RDD
val lines = sc.textFile("hdfs://...")

// transformed RDDs
val errors = lines.filter(_.startsWith("ERROR"))
val messages = errors.map(_.split("\t")).map(r => r(1))
messages.cache()

// action 1
messages.filter(_.contains("mysql")).count()
```

discussing the other part

```scala
// action 2
messages.filter(_.contains("php")).count()
```

cache 1

Worker

block 1

cache 2

Worker

block 2

Driver

cache 3

Worker

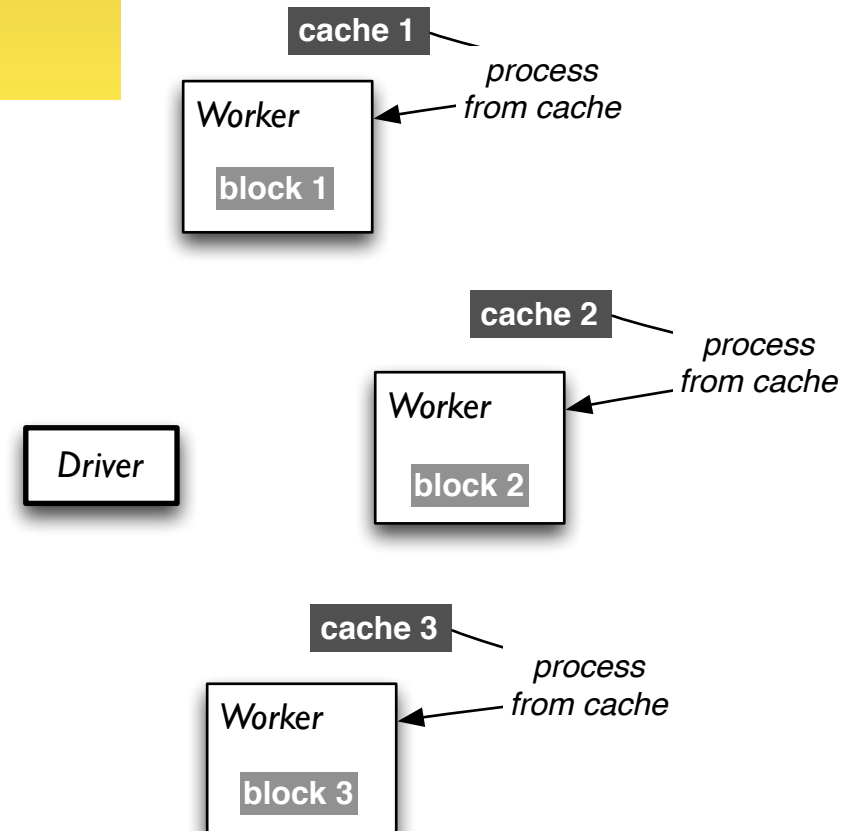block 3

# Spark Deconstructed: *Log Mining Example*

```scala
// base RDD
val lines = sc.textFile("hdfs://...")

// transformed RDDs
val errors = lines.filter(_.startsWith("ERROR"))
val messages = errors.map(_.split("\t")).map(r => r(1))
messages.cache()

// action 1
messages.filter(_.contains("mysql")).count()
```

*discussing the other part*

```scala
// action 2
messages.filter(_.contains("php")).count()
```

**cache 1**

process from cache

**Worker**

**block 1**

**cache 2**

process from cache

**Worker**

**block 2**

**Driver**

**cache 3**

process from cache
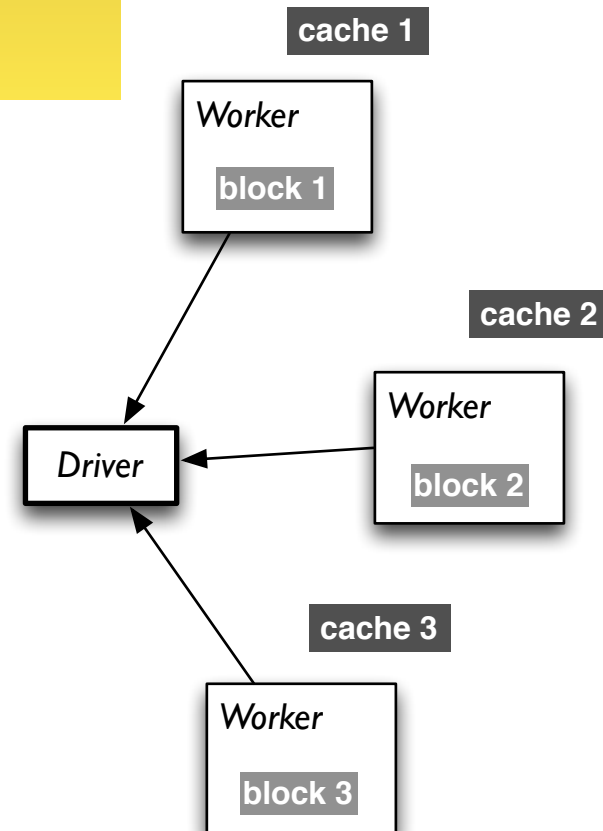
**Worker**

**block 3**

# Spark Deconstructed: *Log Mining Example*

```scala
// base RDD
val lines = sc.textFile("hdfs://...")

// transformed RDDs
val errors = lines.filter(_.startsWith("ERROR"))
val messages = errors.map(_.split("\t")).map(r => r(1))
messages.cache()

// action 1
messages.filter(_.contains("mysql")).count()
```

discussing the other part

```scala
// action 2
messages.filter(_.contains("php")).count()
```

cache 1

Worker

block 1

cache 2

Worker

block 2

Driver

cache 3

Worker

block 3

29

**Spark Deconstructed:**

# Looking at the RDD transformations and actions from another perspective…
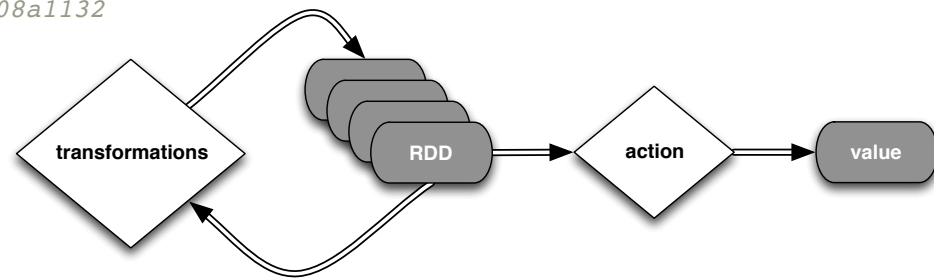
```
// load error messages from a log into memory
// then interactively search for various patterns
// https://gist.github.com/ceteri/8ae5b9509a08c08a1132

// base RDD
val lines = sc.textFile("hdfs://...")

// transformed RDDs
val errors = lines.filter(_.startsWith("ERROR"))
val messages = errors.map(_.split("\t")).map(r => r(1))
messages.cache()

// action 1
messages.filter(_.contains("mysql")).count()

// action 2
messages.filter(_.contains("php")).count()
```
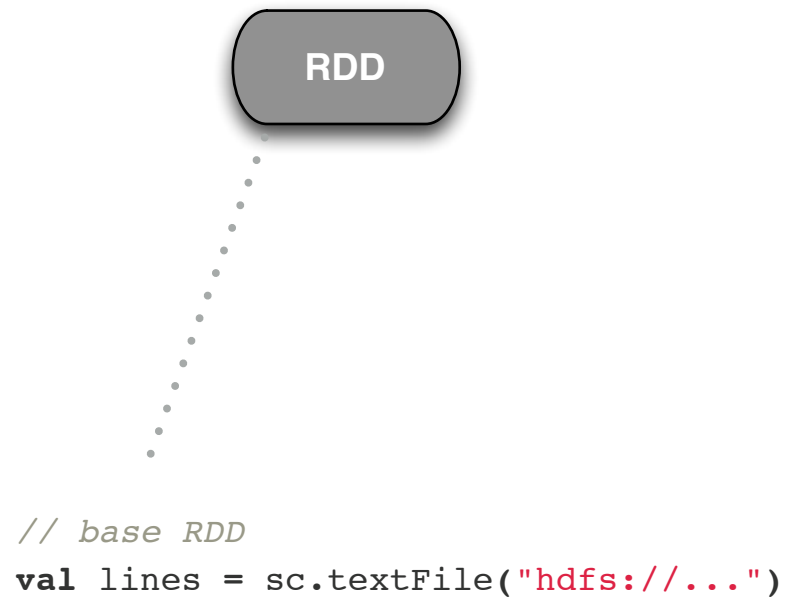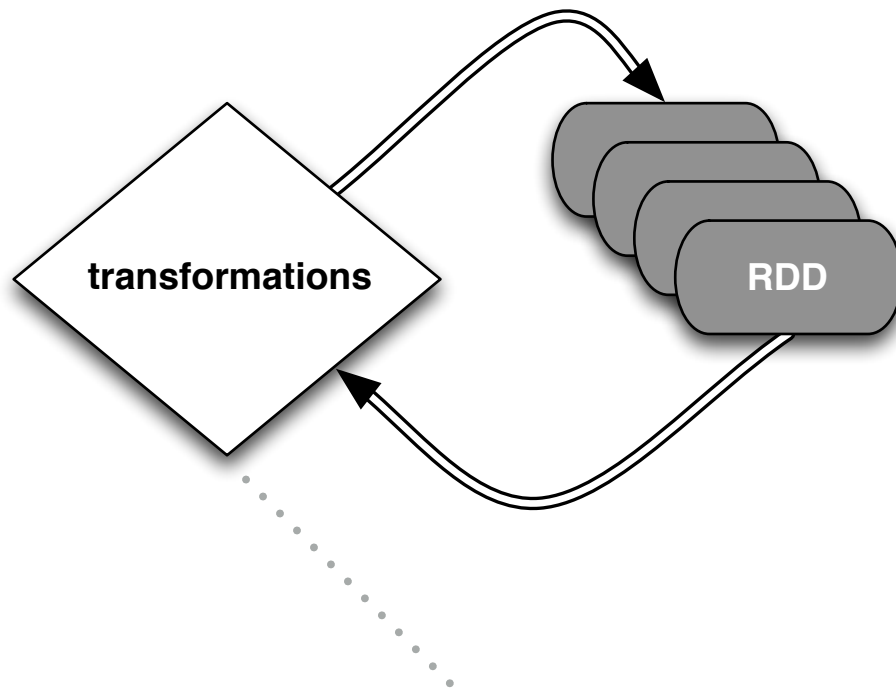
# Spark Deconstructed:
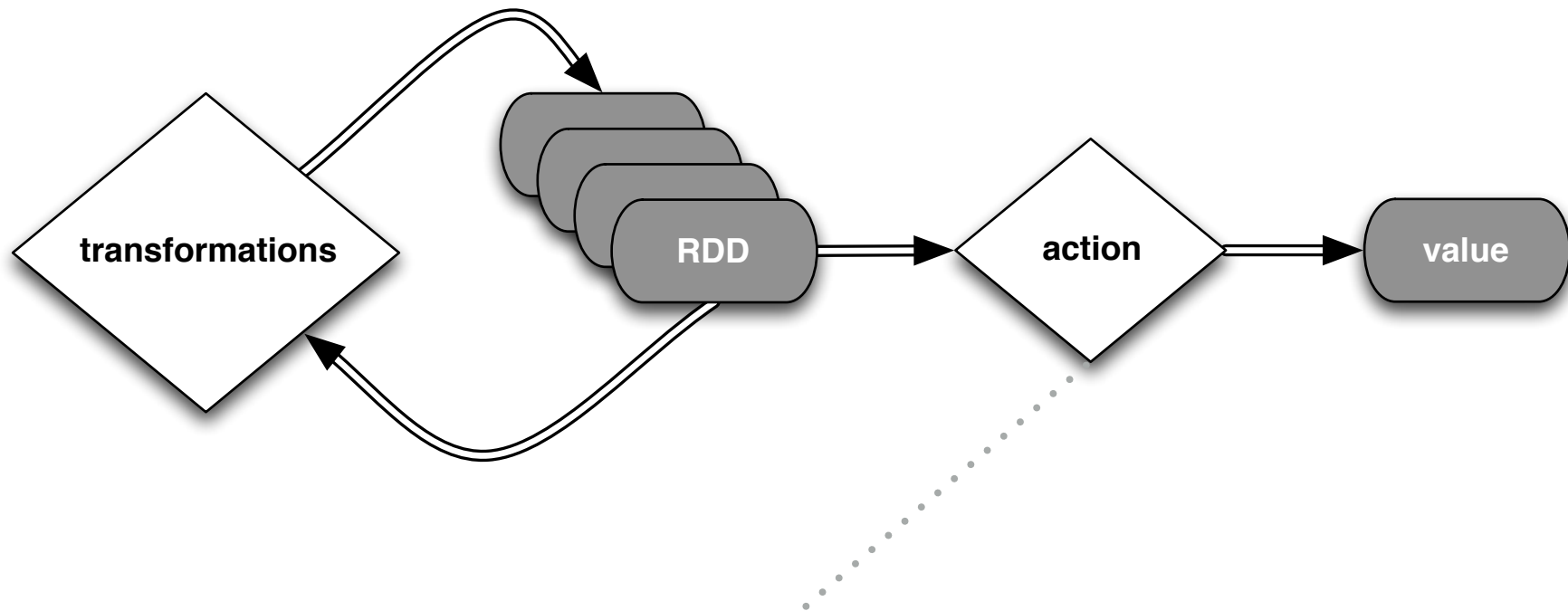
**RDD**

```
// base RDD
val lines = sc.textFile("hdfs://...")
```

# Spark Deconstructed:



```scala
// transformed RDDs
val errors = lines.filter(_.startsWith("ERROR"))
val messages = errors.map(_.split("\t")).map(r => r(1))
messages.cache()
```

# Spark Deconstructed:



```
// action 1
messages.filter(_.contains("mysql")).count()
```