# RDD FUNDAMENTALS

databricks

# INTERACTIVE SHELL



(Scala, Python and R only)

more partitions = more parallelism

RDD

| item-1  | item-6  | item-11 | item-16 | item-21 |
| item-2  | item-7  | item-12 | item-17 | item-22 |
| item-3  | item-8  | item-13 | item-18 | item-23 |
| item-4  | item-9  | item-14 | item-19 | item-24 |
| item-5  | item-10 | item-15 | item-20 | item-25 |

W

Ex
RDD
RDD

W

Ex
RDD
RDD

W

Ex
RDD

RDD w/ 4 partitions

| Error, ts, msg1<br>Warn, ts, msg2<br>Error, ts, msg1 | Info, ts, msg8<br>Warn, ts, msg2<br>Info, ts, msg8 | Error, ts, msg3<br>Info, ts, msg5<br>Info, ts, msg5 | Error, ts, msg4<br>Warn, ts, msg9<br>Error, ts, msg1 |
|---|---|---|---|

logLinesRDD

A base RDD can be created 2 ways:

- Parallelize a collection
- Read data from an external source (S3, C*, HDFS, etc)

# PARALLELIZE

```python
# Parallelize in Python
wordsRDD = sc.parallelize(["fish", "cats", "dogs"])
```

```scala
// Parallelize in Scala
val wordsRDD= sc.parallelize(List("fish", "cats", "dogs"))
```

```java
// Parallelize in Java
JavaRDD<String> wordsRDD = sc.parallelize(Arrays.asList("fish", "cats", "dogs"));
```

- Take an existing in-memory collection and pass it to SparkContext's parallelize method

- Not generally used outside of prototyping and testing since it requires entire dataset in memory on one machine

READ FROM TEXT FILE

```python
# Read a local txt file in Python
linesRDD = sc.textFile("/path/to/README.md")
```

- There are other methods to read data from HDFS, C*, S3, HBase, etc.
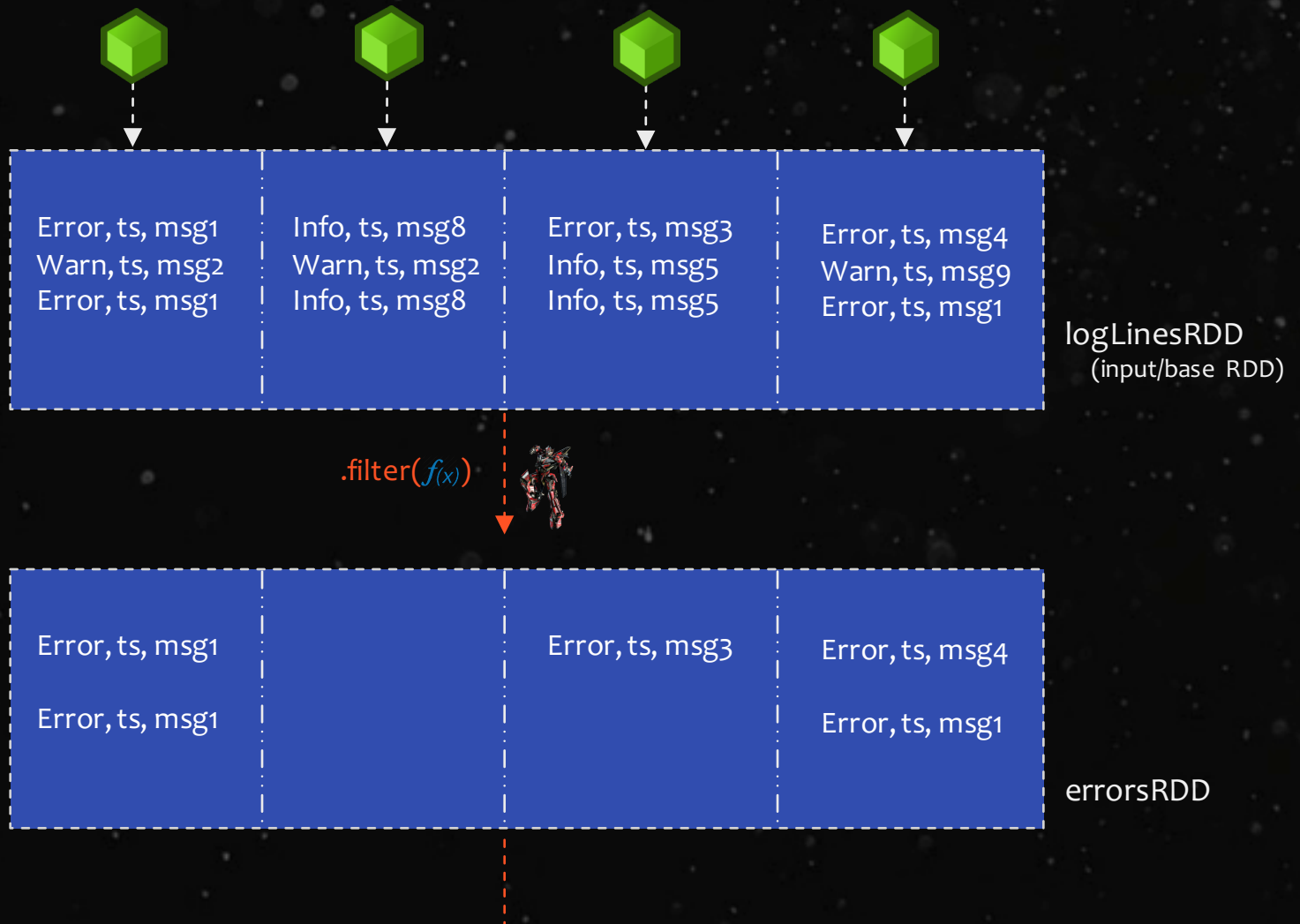
```scala
// Read a local txt file in Scala
val linesRDD = sc.textFile("/path/to/README.md")
```
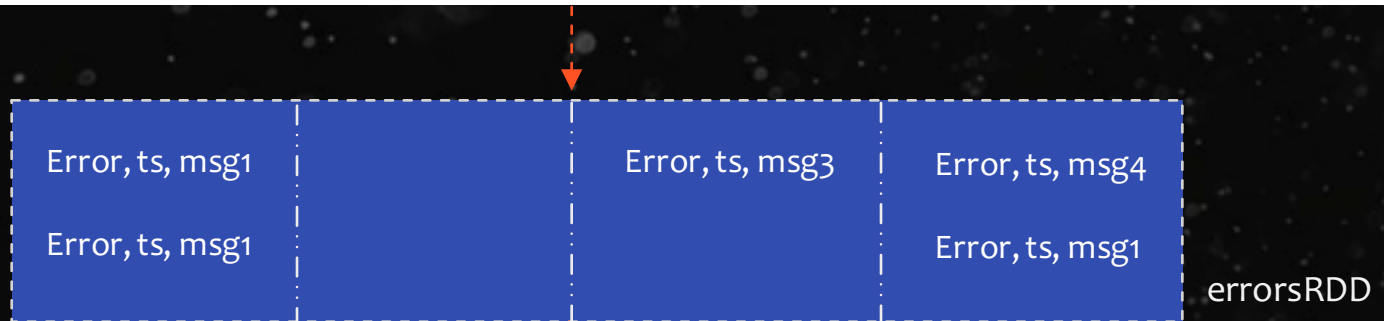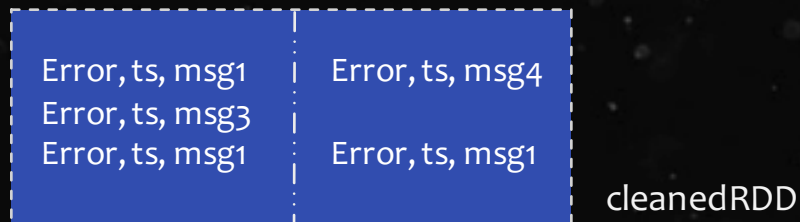
```java
// Read a local txt file in Java
JavaRDD<String> lines = sc.textFile("/path/to/README.md");
```

**hadoop** *HDFS*

logLinesRDD
(input/base RDD)

| Error, ts, msg1 | Info, ts, msg8 | Error, ts, msg3 | Error, ts, msg4 |
| Warn, ts, msg2 | Warn, ts, msg2 | Info, ts, msg5 | Warn, ts, msg9 |
| Error, ts, msg1 | Info, ts, msg8 | Info, ts, msg5 | Error, ts, msg1 |

.filter($f_{(x)}$)

| Error, ts, msg1 | | Error, ts, msg3 | Error, ts, msg4 |
| Error, ts, msg1 | | | Error, ts, msg1 |

errorsRDD

Error, ts, msg1

Error, ts, msg1

Error, ts, msg3

Error, ts, msg4

Error, ts, msg1

errorsRDD

.coalesce( 2 )

Error, ts, msg1
Error, ts, msg3
Error, ts, msg1

Error, ts, msg4

Error, ts, msg1

cleanedRDD

.collect( )

Driver

# Execute DAG!

.collect( )

Driver

logLinesRDD

logLinesRDD

.filter($f(x)$)

errorsRDD

.coalesce( 2 )

cleanedRDD

.collect( )

Error, ts, msg1    Error, ts, msg4
Error, ts, msg3
Error, ts, msg1    Error, ts, msg1

Driver

logLinesRDD

errorsRDD

cleanedRDD

Driver

logLinesRDD

errorsRDD

.saveAsTextFile( )

Error, ts, msg1
Error, ts, msg3
Error, ts, msg1

Error, ts, msg4

Error, ts, msg1

cleanedRDD

.filter($f_{(x)}$)

.count( )

5

Error, ts, msg1

Error, ts, msg1

Error, ts, msg1

Error, ts, msg1

errorMsg1RDD

.collect( )

logLinesRDD

errorsRDD

.cache( )

.saveAsTextFile( )

Error, ts, msg1
Error, ts, msg3
Error, ts, msg1

Error, ts, msg4

Error, ts, msg1

cleanedRDD

.filter($f(x)$)

.count( )

5

Error, ts, msg1

Error, ts, msg1

Error, ts, msg1

errorMsg1RDD

.collect( )

# PARTITION -> TASK -> PARTITION

.filter($f(x)$)

logLinesRDD
(HadoopRDD)

Task-1

Task-2

Task-3

Task-4

errorsRDD
(filteredRDD)
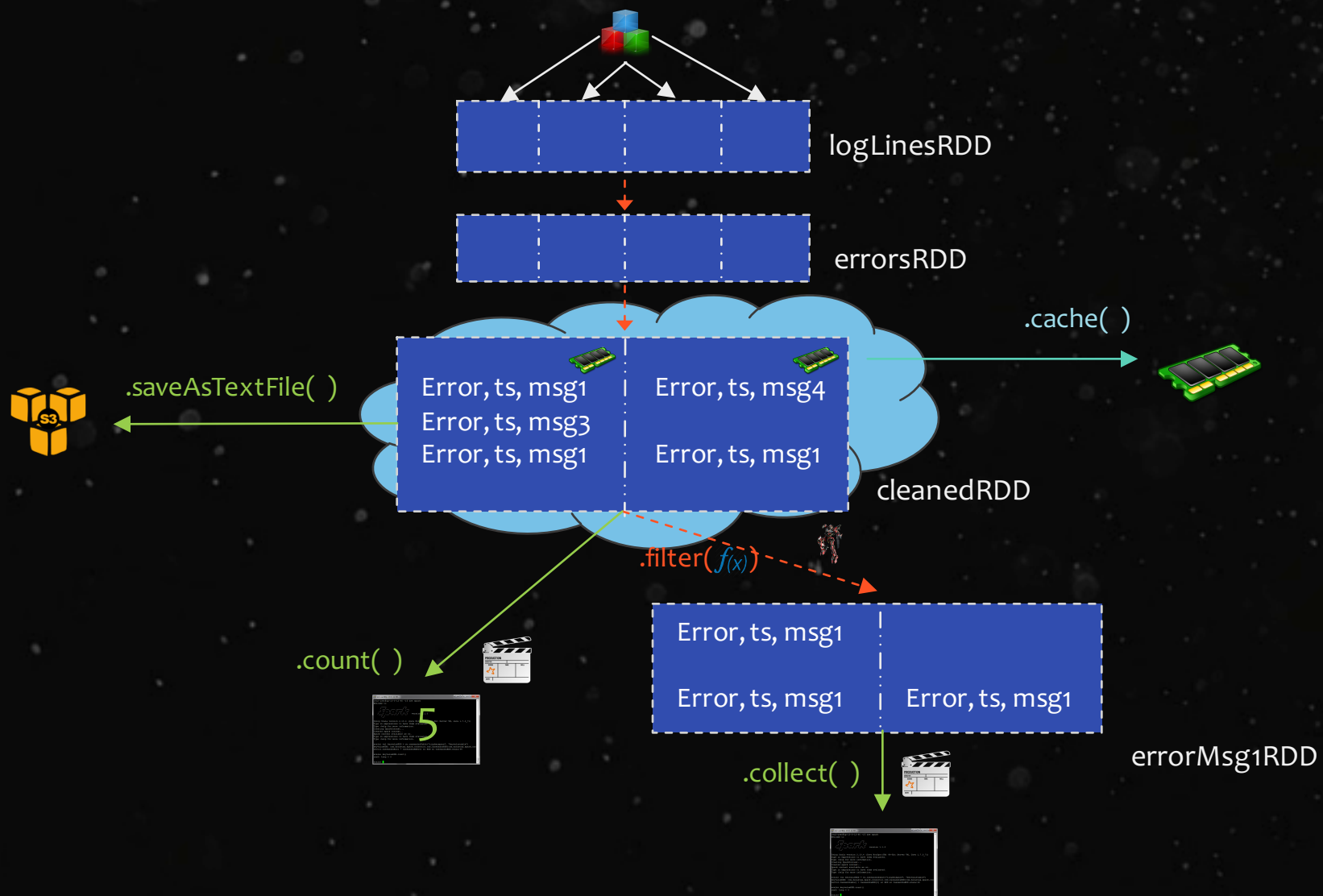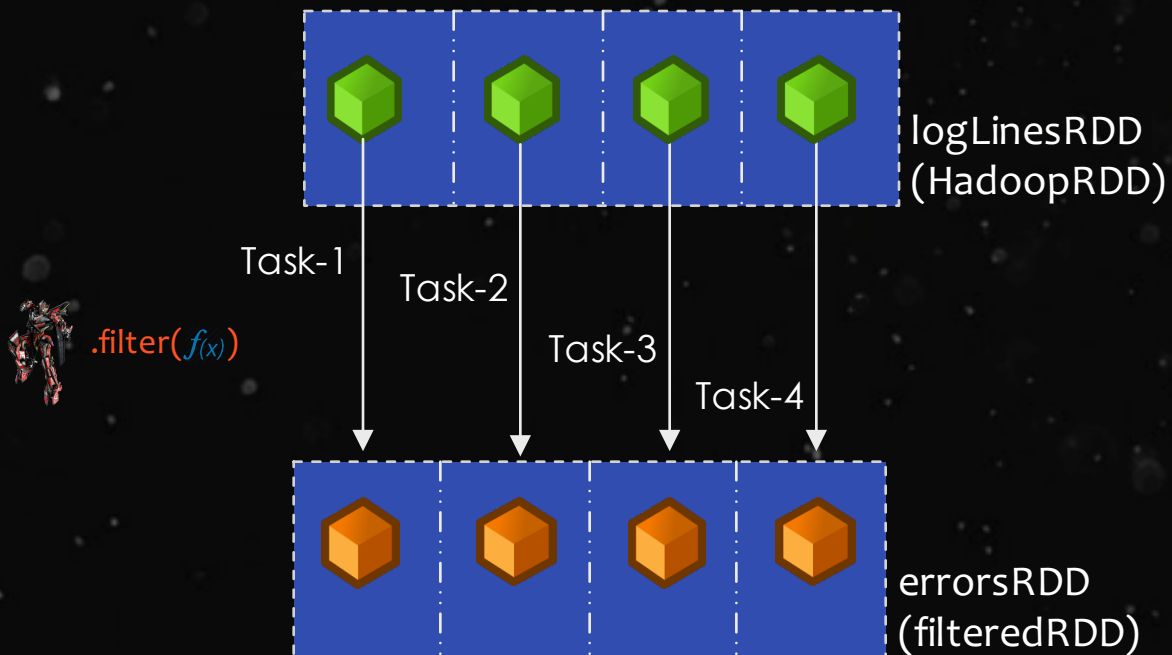
# LIFECYCLE OF A SPARK PROGRAM

1) Create some input RDDs from external data or parallelize a
   collection in your driver program.

2) Lazily transform them to define new RDDs using
   transformations like `filter()` or `map()`

3) Ask Spark to `cache()` any intermediate RDDs that will need to
   be reused.

4) Launch actions such as `count()` and `collect()` to kick off a
   parallel computation, which is then optimized and executed
   by Spark.

# TRANSFORMATIONS (lazy)

| map() | intersection() | cartesion() |
| flatMap() | distinct() | pipe() |
| filter() | groupByKey() | coalesce() |
| mapPartitions() | reduceByKey() | repartition() |
| mapPartitionsWithIndex() | sortByKey() | partitionBy() |
| sample() | join() | ... |
| union() | cogroup() | ... |

# ACTIONS

reduce()

collect()

count()

first()

take()

takeSample()

saveToCassandra()

takeOrdered()

saveAsTextFile()

saveAsSequenceFile()

saveAsObjectFile()

countByKey()

foreach()

...

# TYPES OF RDDS

- HadoopRDD

- FilteredRDD

- MappedRDD

- PairRDD

- ShuffledRDD

- UnionRDD

- PythonRDD

- DoubleRDD

- JdbcRDD

- JsonRDD

- VertexRDD

- EdgeRDD

- CassandraRDD *(DataStax)*

- GeoRDD *(ESRI)*

- EsSpark *(ElasticSearch)*