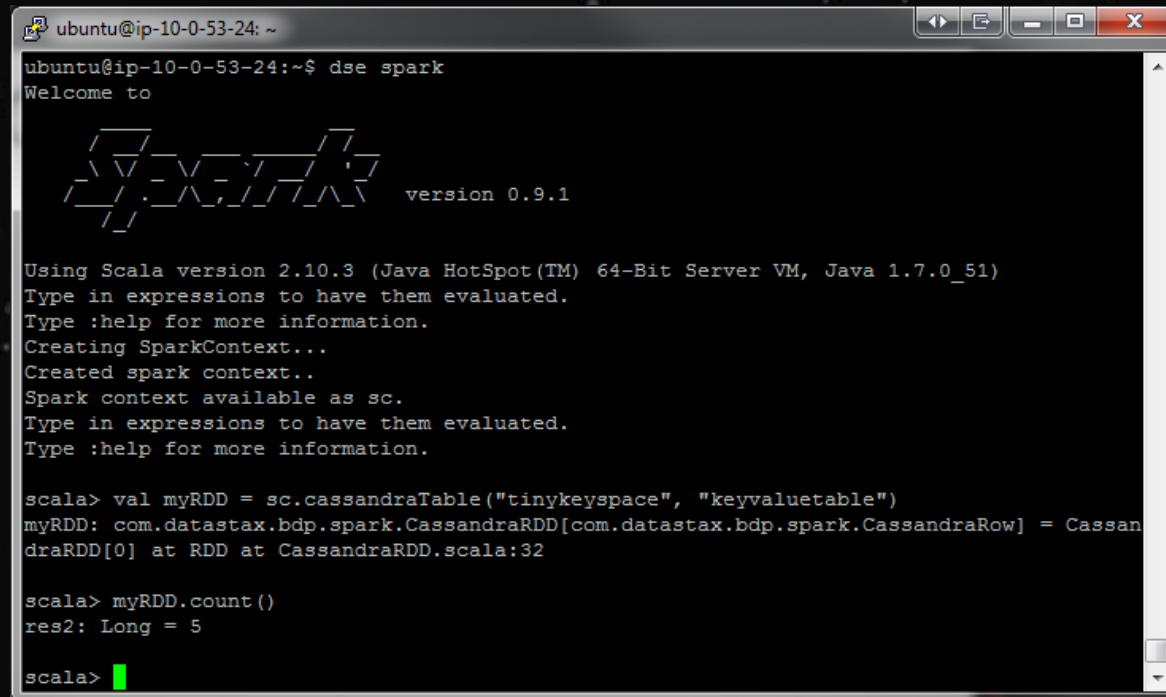


# RDD FUNDAMENTALS



# INTERACTIVE SHELL

A terminal window titled 'ubuntu@ip-10-0-53-24: ~' with standard window controls. The prompt is 'ubuntu@ip-10-0-53-24:~\$'. The user enters 'dse spark'. The output shows a 'Welcome to' message, a stylized 'DSE' logo, and 'version 0.9.1'. It then displays Scala version and JVM information, followed by instructions to type expressions or ':help'. The SparkContext is created and named 'sc'. The user enters a Scala command to create a Cassandra table, followed by a count command which returns 'res2: Long = 5'. The prompt 'scala>' is shown with a green cursor.

```
ubuntu@ip-10-0-53-24:~$ dse spark
Welcome to

  DSE  version 0.9.1

Using Scala version 2.10.3 (Java HotSpot(TM) 64-Bit Server VM, Java 1.7.0_51)
Type in expressions to have them evaluated.
Type :help for more information.
Creating SparkContext...
Created spark context..
Spark context available as sc.
Type in expressions to have them evaluated.
Type :help for more information.

scala> val myRDD = sc.cassandraTable("tinykeyspace", "keyvaluetable")
myRDD: com.datastax.bdp.spark.CassandraRDD[com.datastax.bdp.spark.CassandraRow] = Cassan
draRDD[0] at RDD at CassandraRDD.scala:32

scala> myRDD.count()
res2: Long = 5

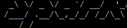
scala>
```

(Scala & Python only)

# Driver Program

```

$ ./bin/cp --cp 10-0-12-60...
kudu@spark10-0-12-60:~$ cd /
kudu@spark10-0-12-60:~$ java -jar spark
Welcome to

 version 1.1.0

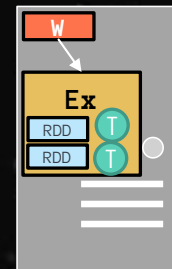
Using Scala version 2.10.4 (Java HotSpot(TM) 64-Bit Server VM, Java 1.7.0_71)
in expressions to have them evaluated.
Type :help for more information.
Creating SparkContext...
Created spark context...
Spark context available as sc.
Type in expressions to have them evaluated.
Type :help for more information.

scala> val keyValuesRDD = sc.cassandraTable("tinykeyspace", "keyvaluetable")
scala> sc.dataTable(spark.connector.rdd.CassandraRDD[com.datastax.spark.
connector.CassandraRow] = CassandraRDD[0] at RDD at CassandraRDD.scala:49

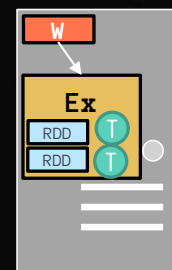
scala> keyValuesRDD.count()
res2: Long = 4

scala>

```



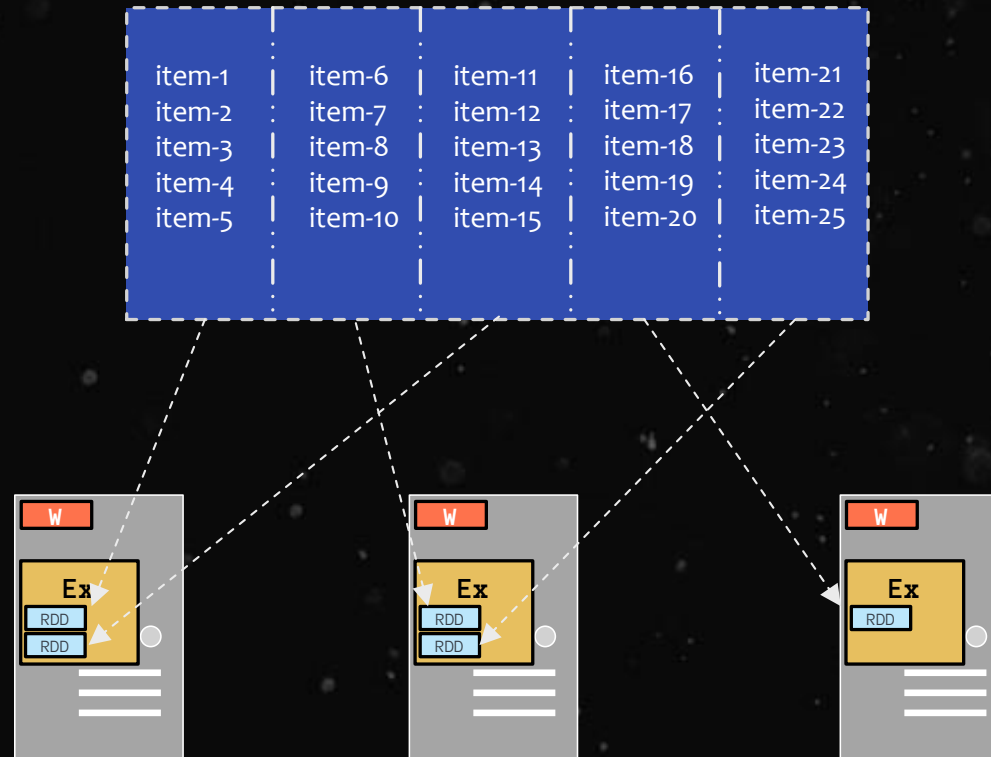
## Worker Machine



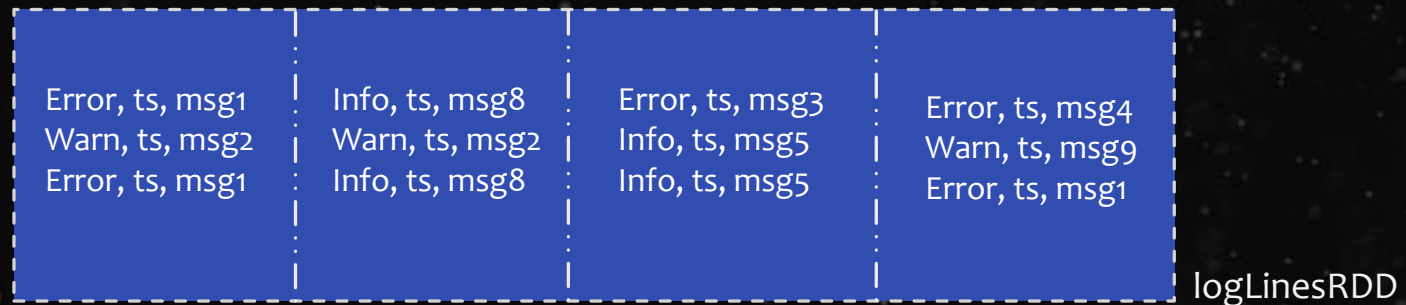
## Worker Machine

*more partitions = more parallelism*

## RDD



## RDD w/ 4 partitions



An RDD can be created 2 ways:

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

# PARALLELIZE



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

---



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

---



```
// 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



```
# 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.

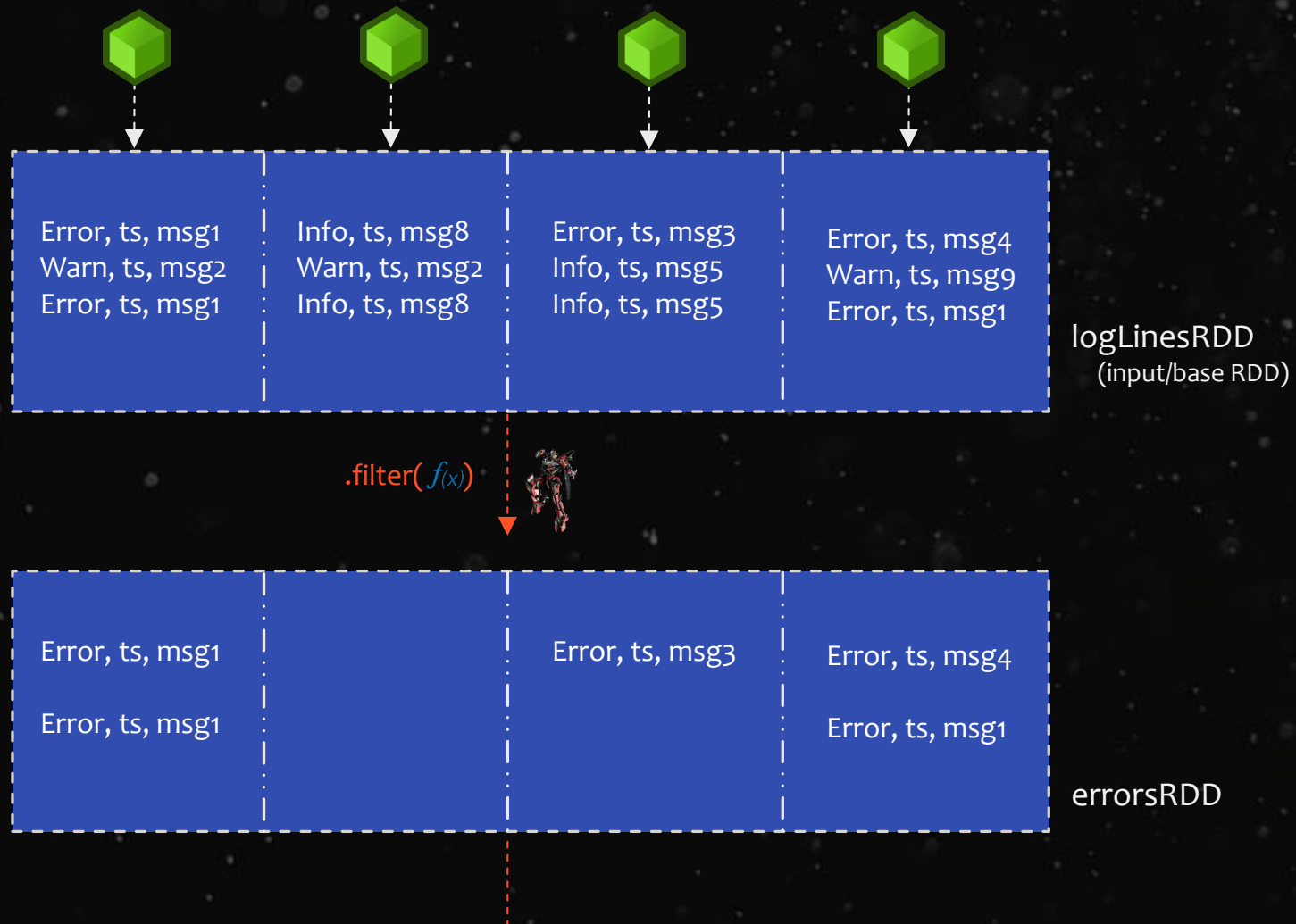


```
// Read a local txt file in Scala  
val linesRDD = sc.textFile("/path/to/README.md")
```

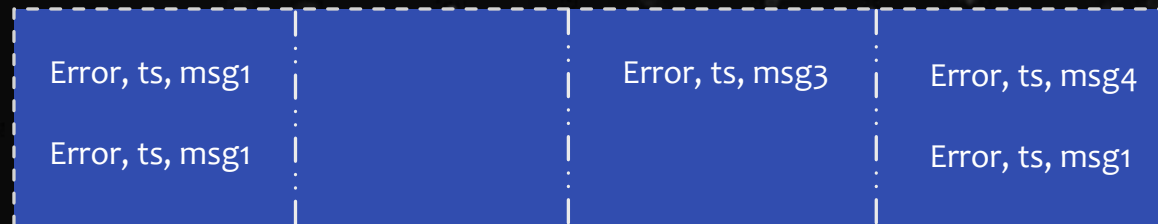
---



```
// Read a local txt file in Java  
JavaRDD<String> lines = sc.textFile("/path/to/README.md");
```

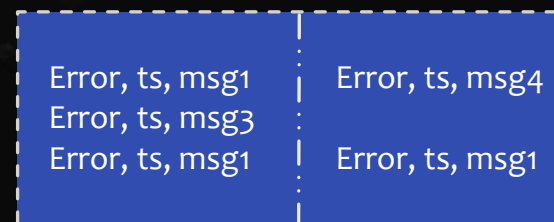






errorsRDD

.coalesce( 2 )



cleanedRDD

.collect( )

```
ec2-user@ip-10-0-12-40: ~$ ./bin/spark
Welcome to
      ____              __
     /  _ \__  _/   ___/  /
    /  /_  /_  /_  /_  /_  /_
   /_/_/_/_/_/_/_/_/_/_/_/_/_/_/_/_
Spark version 1.1.0

Using Scala version 2.10.4 (Java HotSpot(TM) 64-Bit Server VM, Java 1.7.0_71)
Type in expressions to have them evaluated.
Type help for more information.
Creating SparkContext...
Detected spark context...
Spark context available as sc.
Type in expressions to have them evaluated.
Type help for more information.

scala> val keyValueRDD = sc.cassandraTable("tinykeyspace", "keyvaluetable")
keyValueRDD: com.datastax.spark.connector.rdd.CassandraRDD[com.datastax.spark.connector.CassandraRow] = CassandraRDD() as RDD at CassandraRDD.scala:49

scala> keyValueRDD.count()
res2: Long = 4

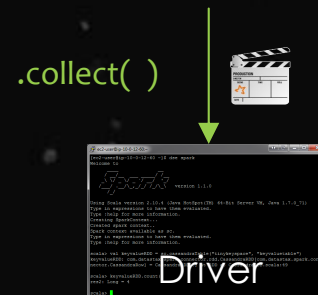
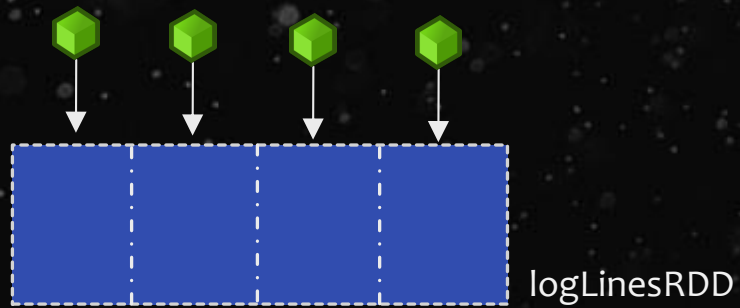
scala>
```

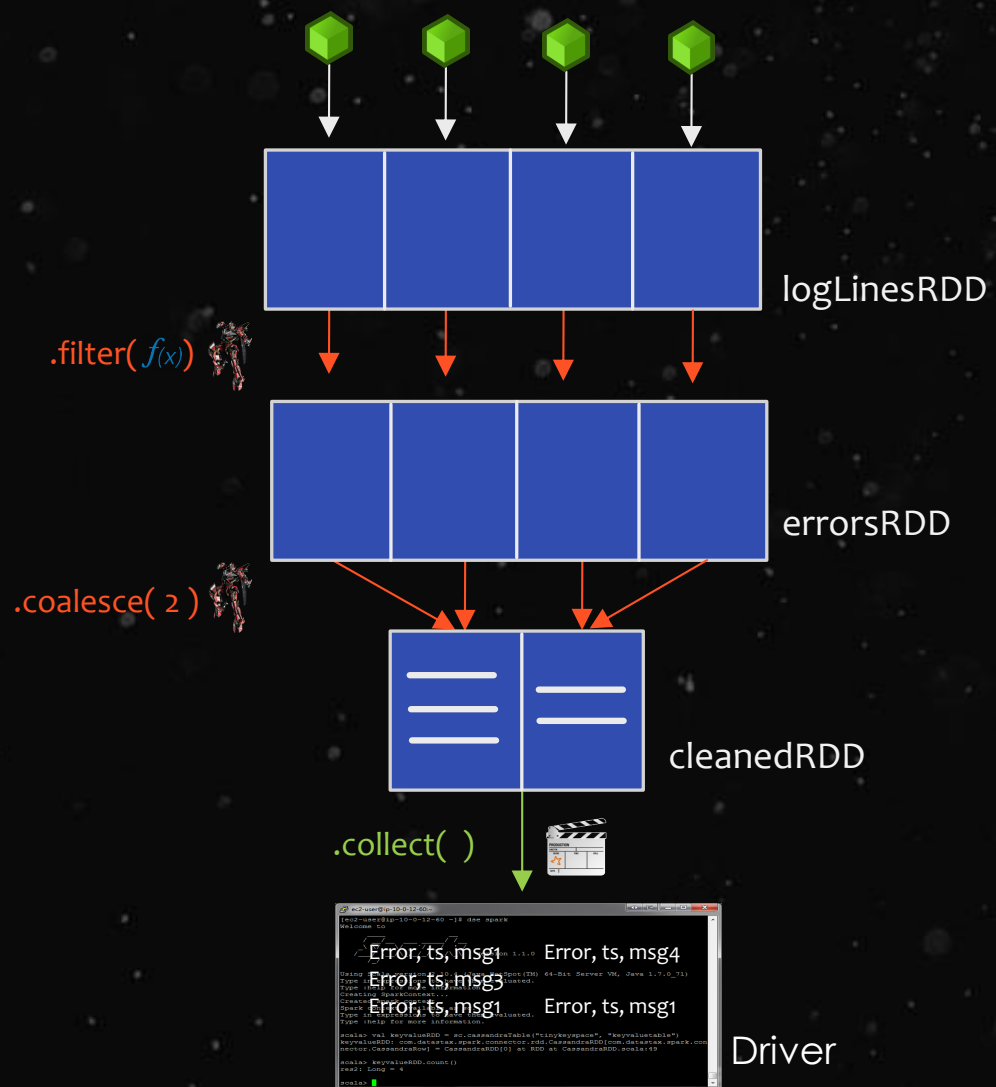
Driver

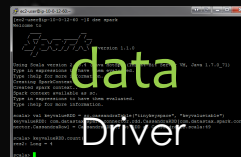
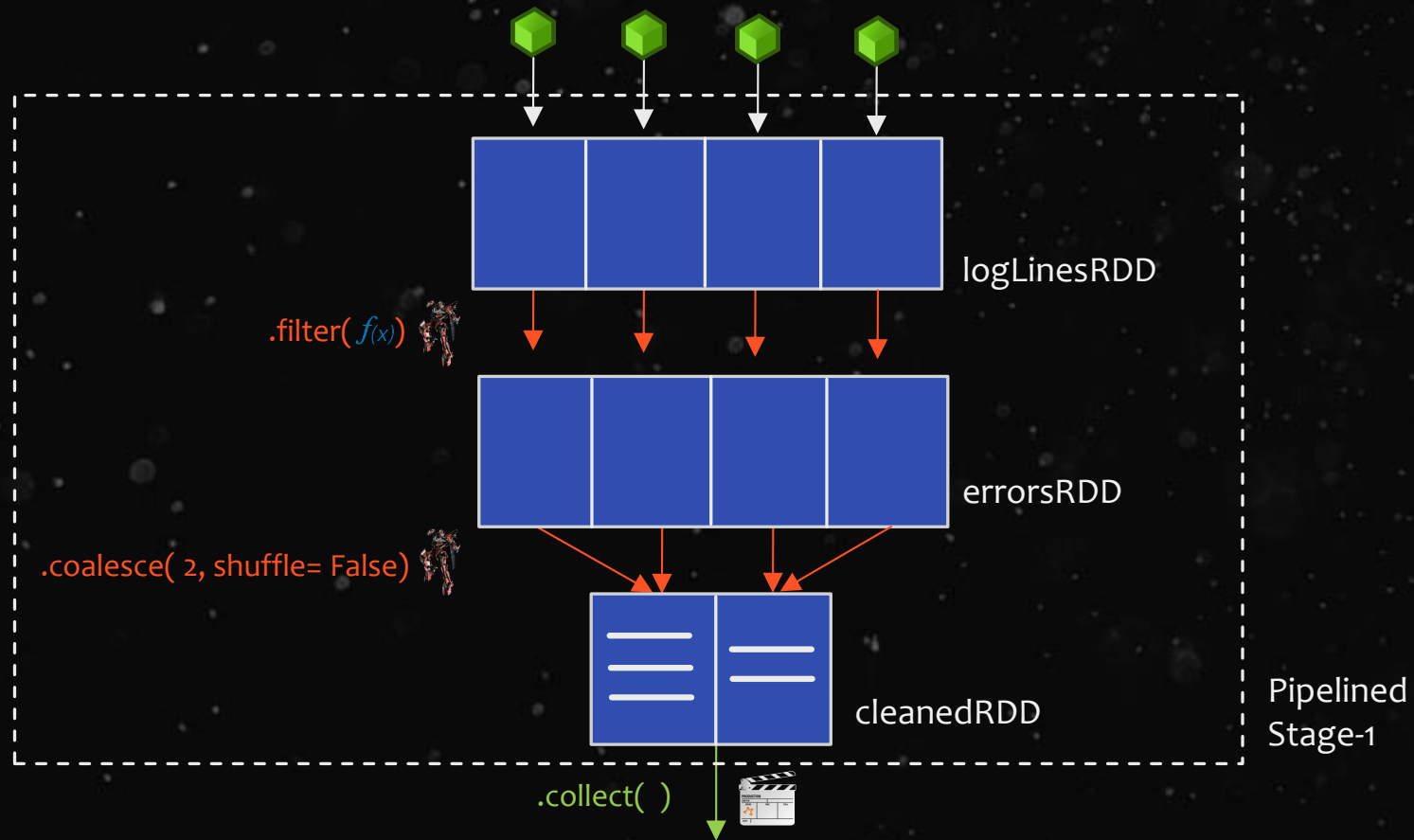


Execute DAG!











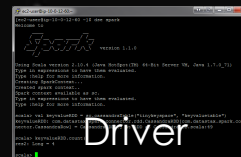
logLinesRDD



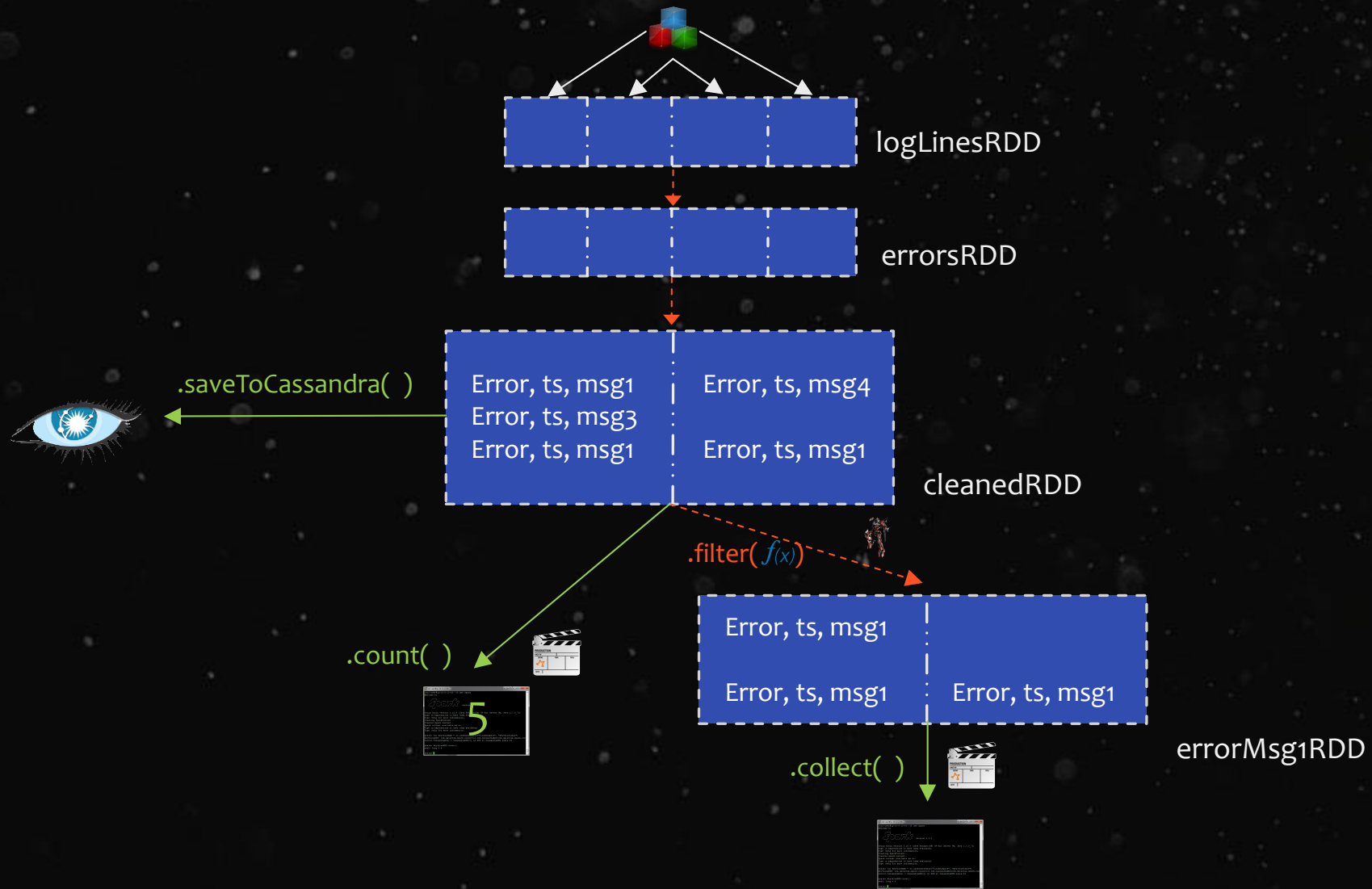
errorsRDD



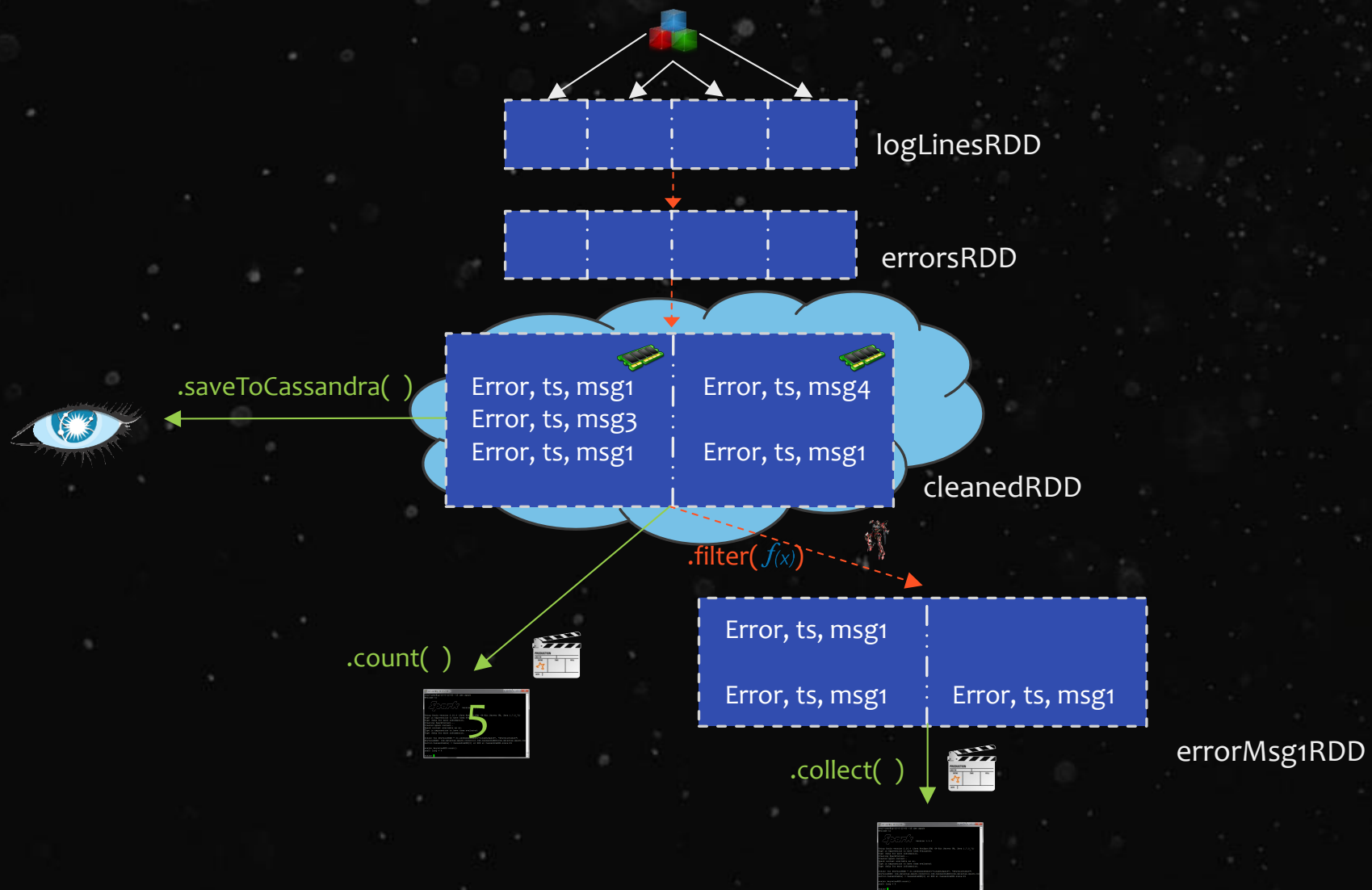
cleanedRDD



[illegible]

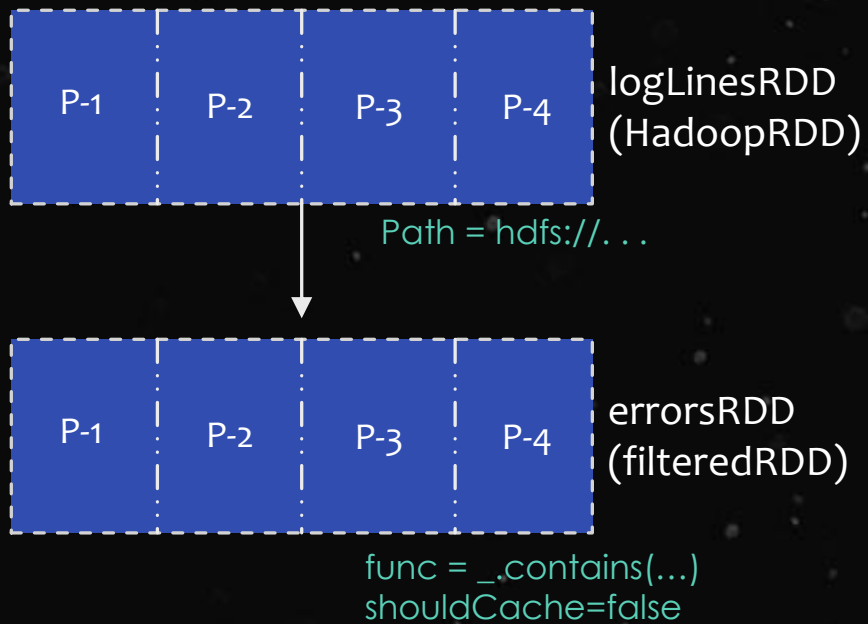




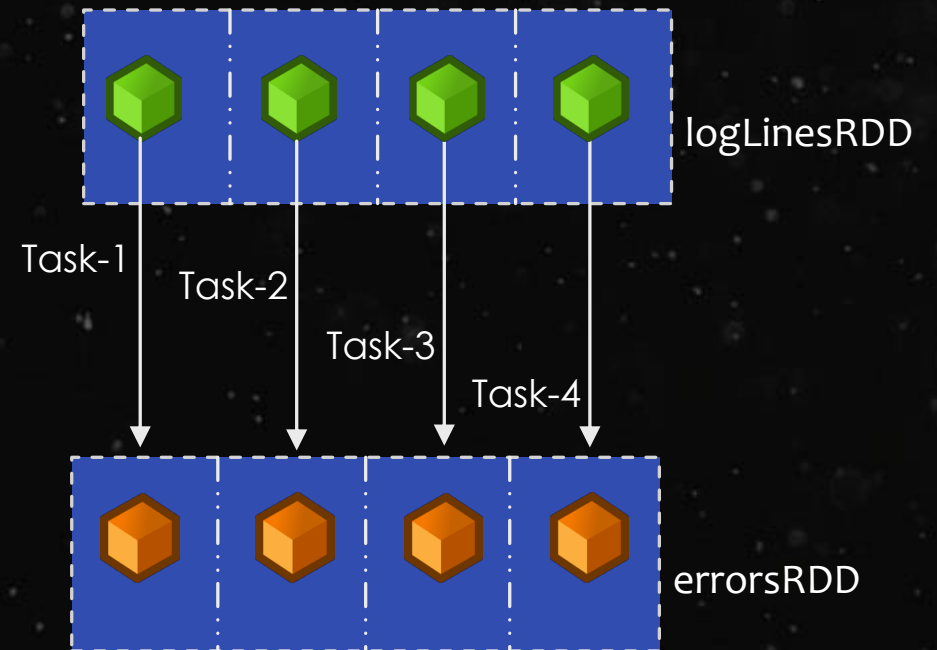


# RDD GRAPH

Dataset-level view:



Partition-level view:



# 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)

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

- Most transformations are element-wise (they work on one element at a time), but this is not true for all transformations

# 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
- SchemaRDD
- VertexRDD
- EdgeRDD
- **CassandraRDD** (*DataStax*)
- **GeoRDD** (*ESRI*)
- **EsSpark** (*ElasticSearch*)