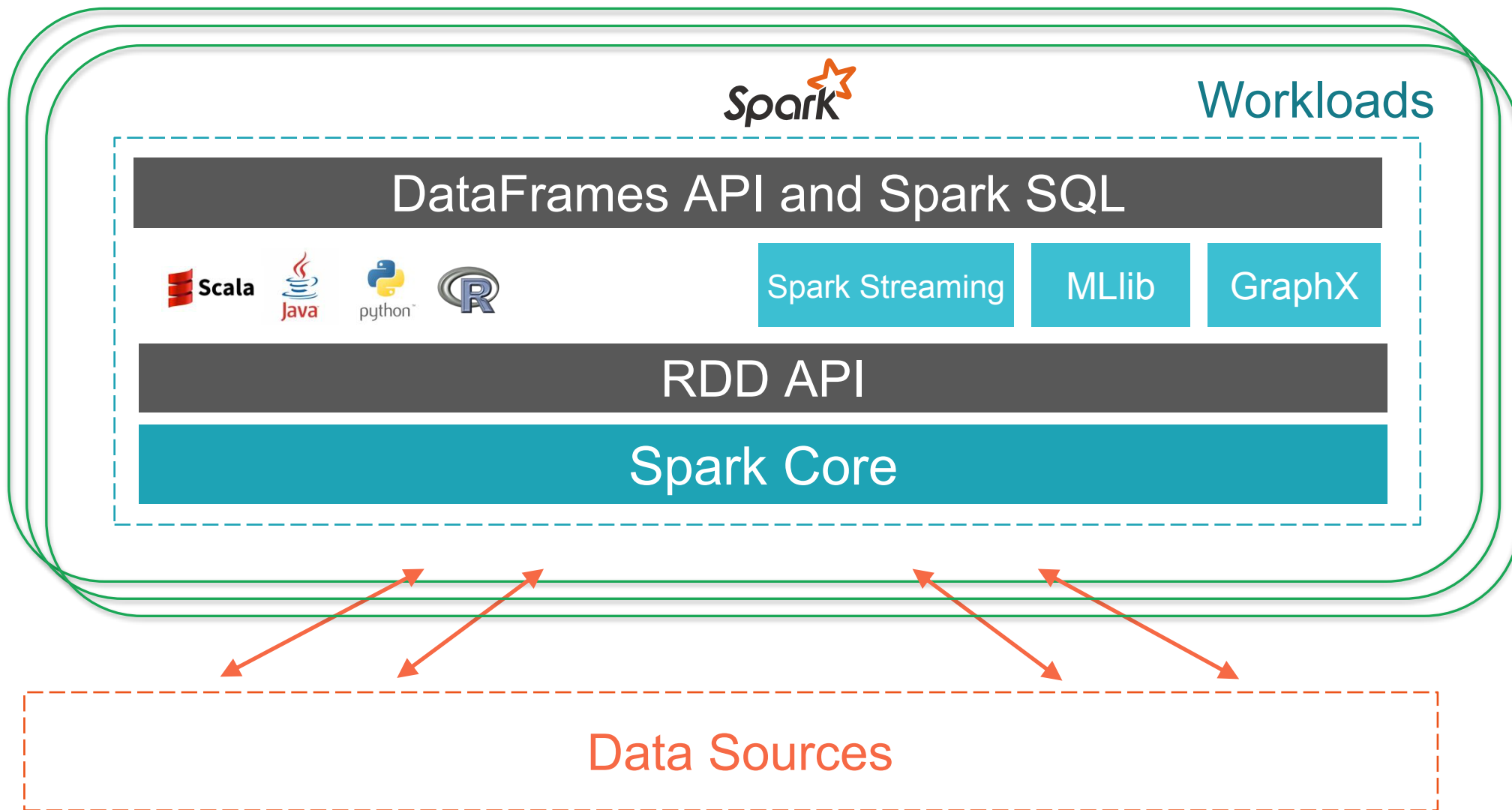


RDD Fundamentals



Driver Program

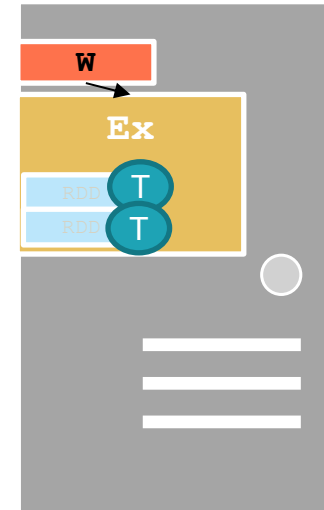
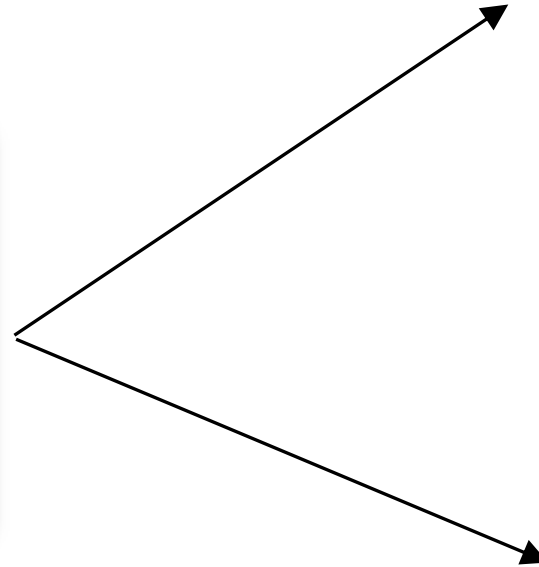
```
$ pyspark
Python 2.7.9 (default, Jan 7 2015, 11:49:12)
Type "copyright", "credits" or "license()" for more information.

Python 2.7.9 -- An enhanced Interactive Python.
?              -> Introduction and overview of Python's features.
help()         -> Quick reference.
help()         -> Python's own help system.
object?        -> Details about 'object', use 'object??' for extra details.
Welcome to

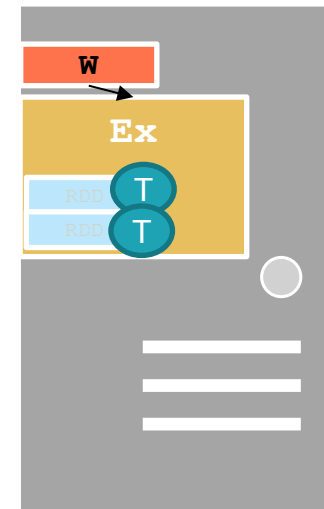
Spark version 1.4.0

Using Python version 2.7.9 (default, Jan 7 2015 11:49:12)
SparkContext available as sc, HiveContext available as sqlContext.

In [1]:
```



Worker Machine

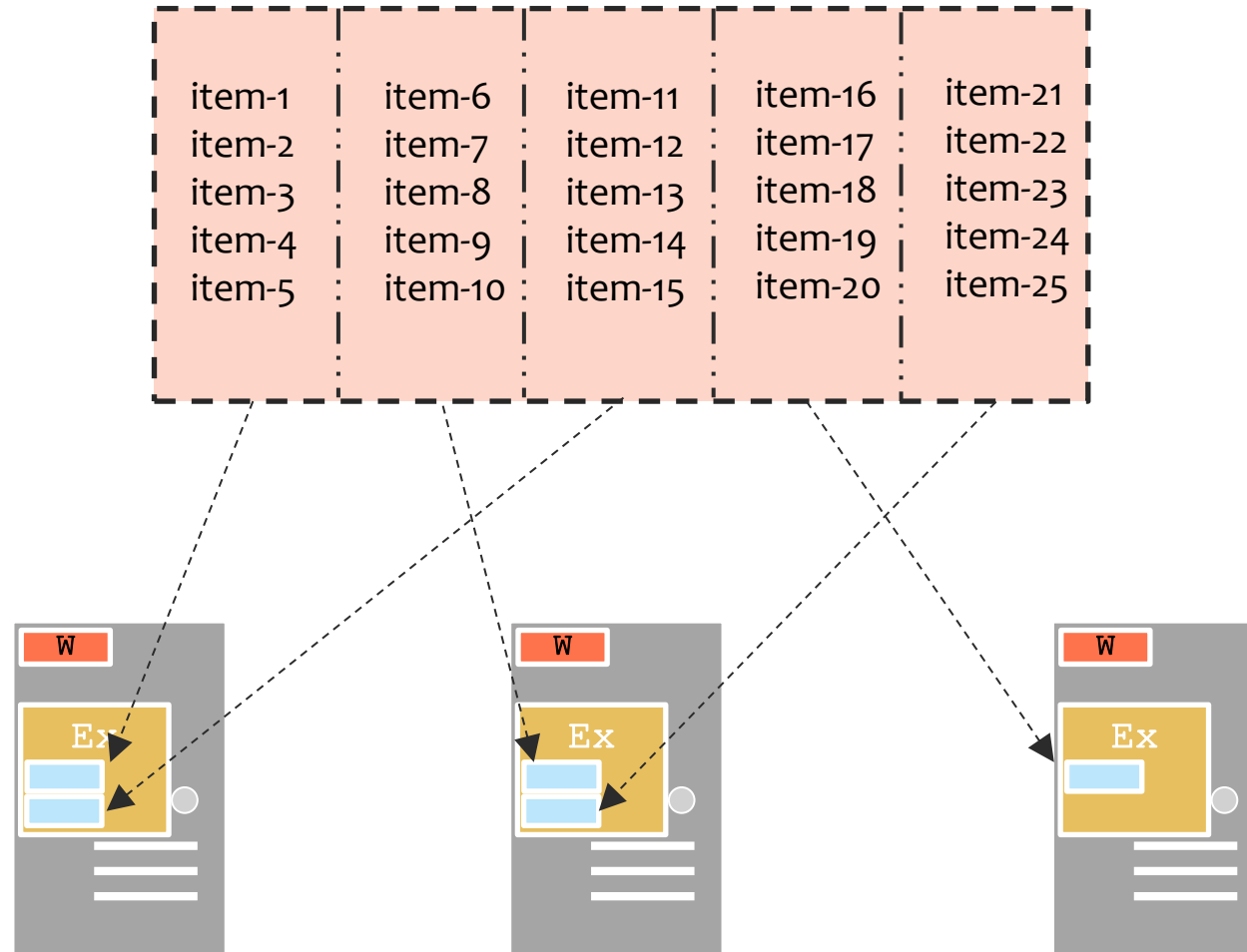


Worker Machine

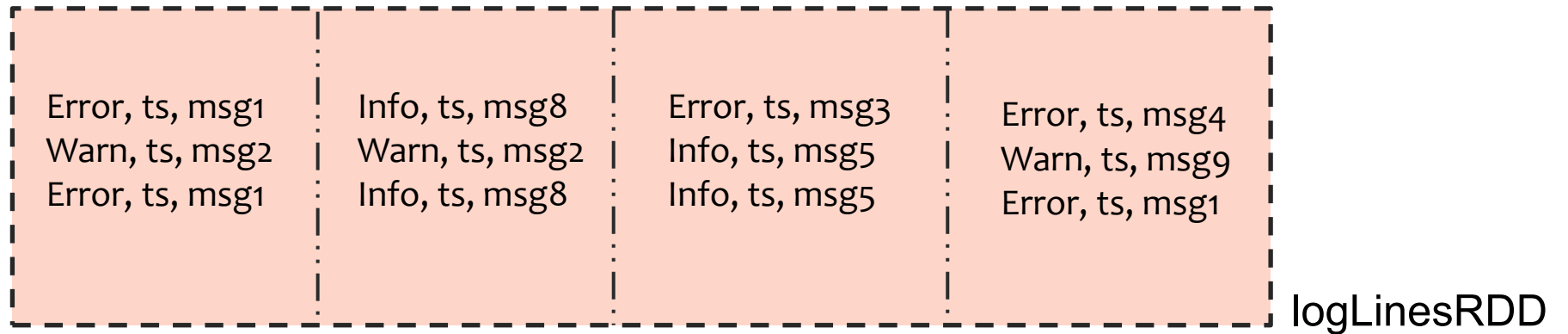
Resilient Distributed Datasets (RDDs)

- Write programs in terms of operations on distributed data
- Partitioned collections of objects spread across a cluster
- Diverse set of parallel transformations and actions
- Fault tolerant

RDD



RDD w/ 4 partitions



A base RDD can be created 2 ways:

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

Create a Base RDD



Parallelize in Python

```
wordsRDD = sc.parallelize(["fish", "cats", "dogs"])
```

Parallelize

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



Read a local txt file in Python

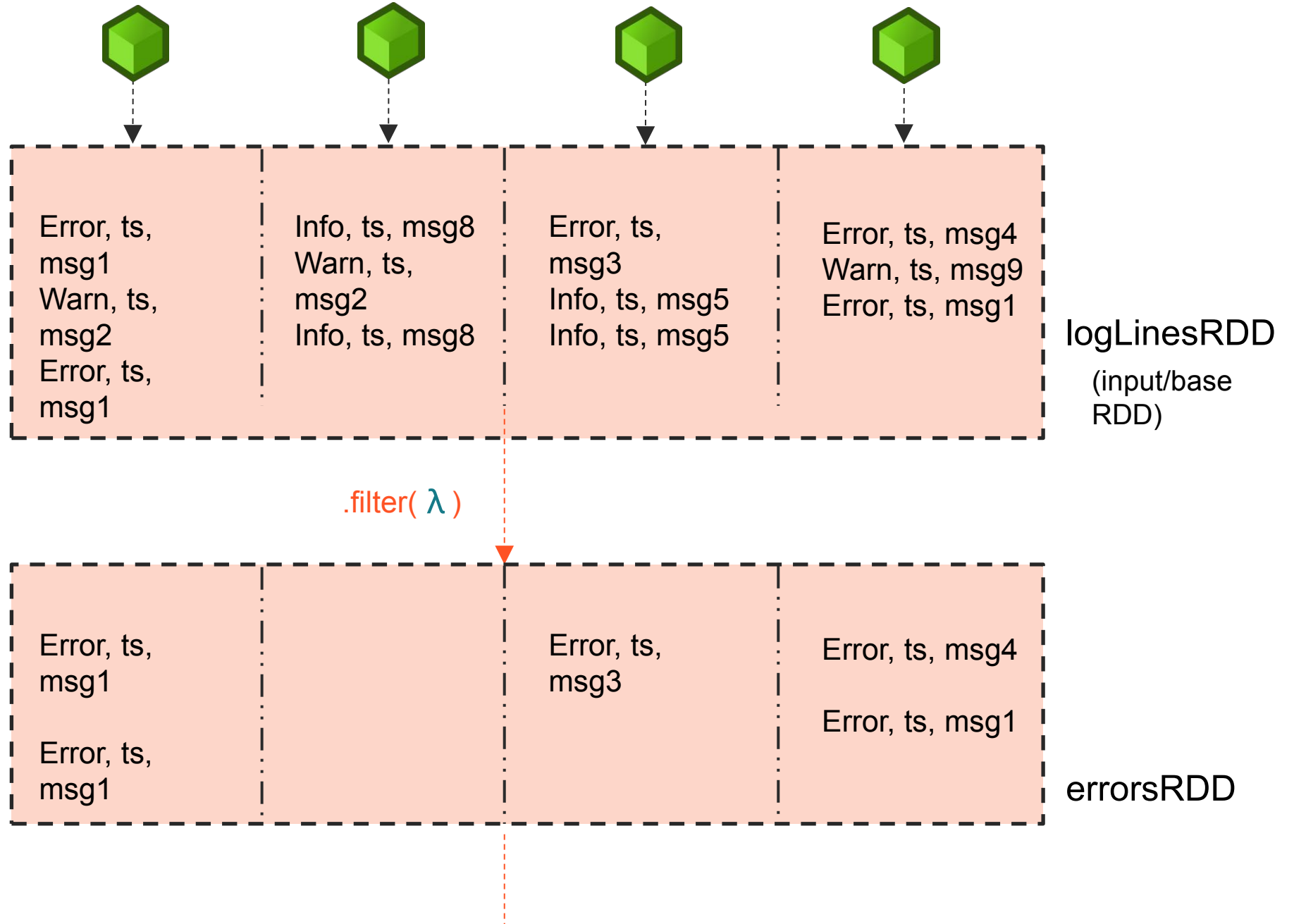
```
linesRDD = sc.textFile("/path/to/README.md")
```

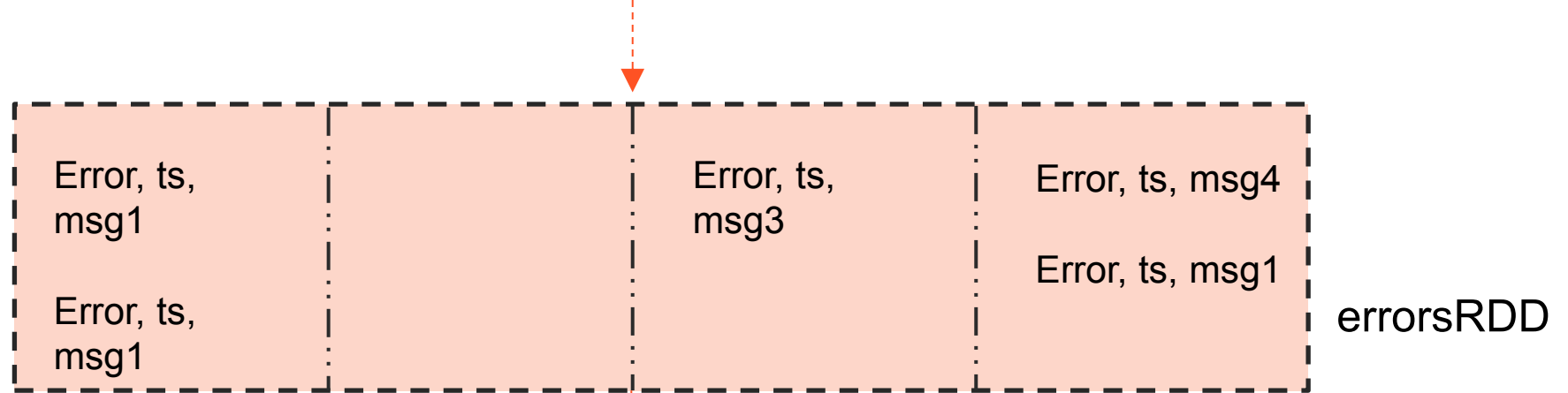
Read from Text File

There are other methods to read data from HDFS, C*, S3, HBase, etc.

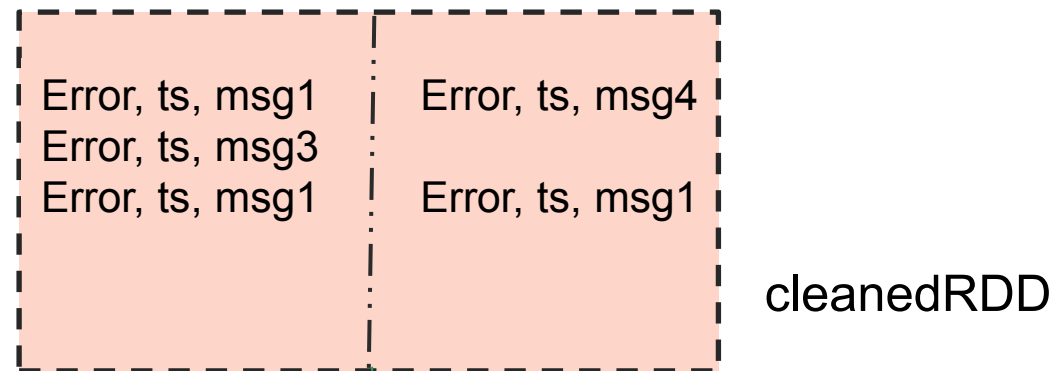
Operations on Distributed Data

- Two types of operations: *transformations* and *actions*
- Transformations are lazy (*not computed immediately*)
- Transformations are executed when an action is run
- Persist (cache) distributed data in memory or disk





`.coalesce(2)`

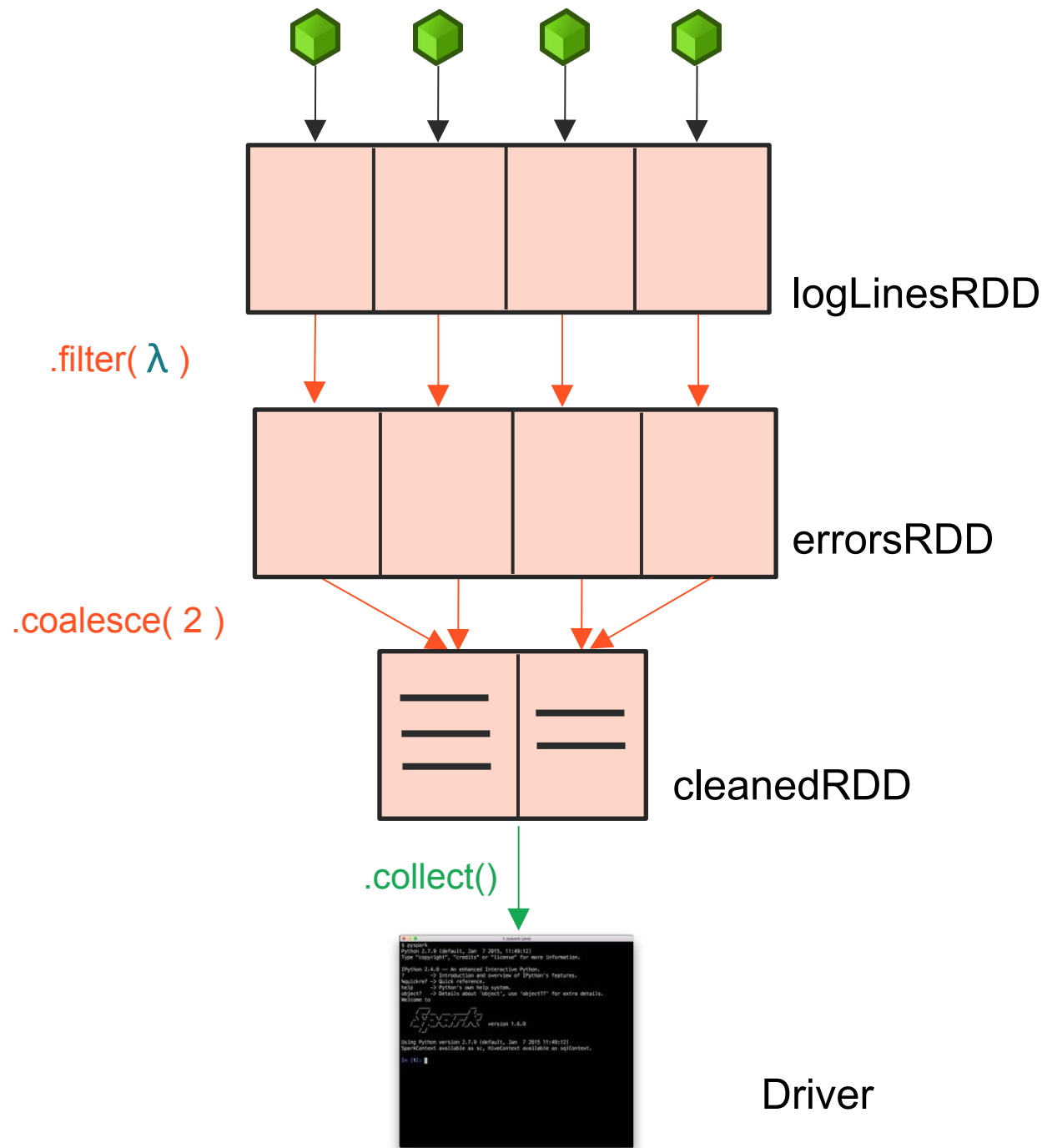


`.collect()`

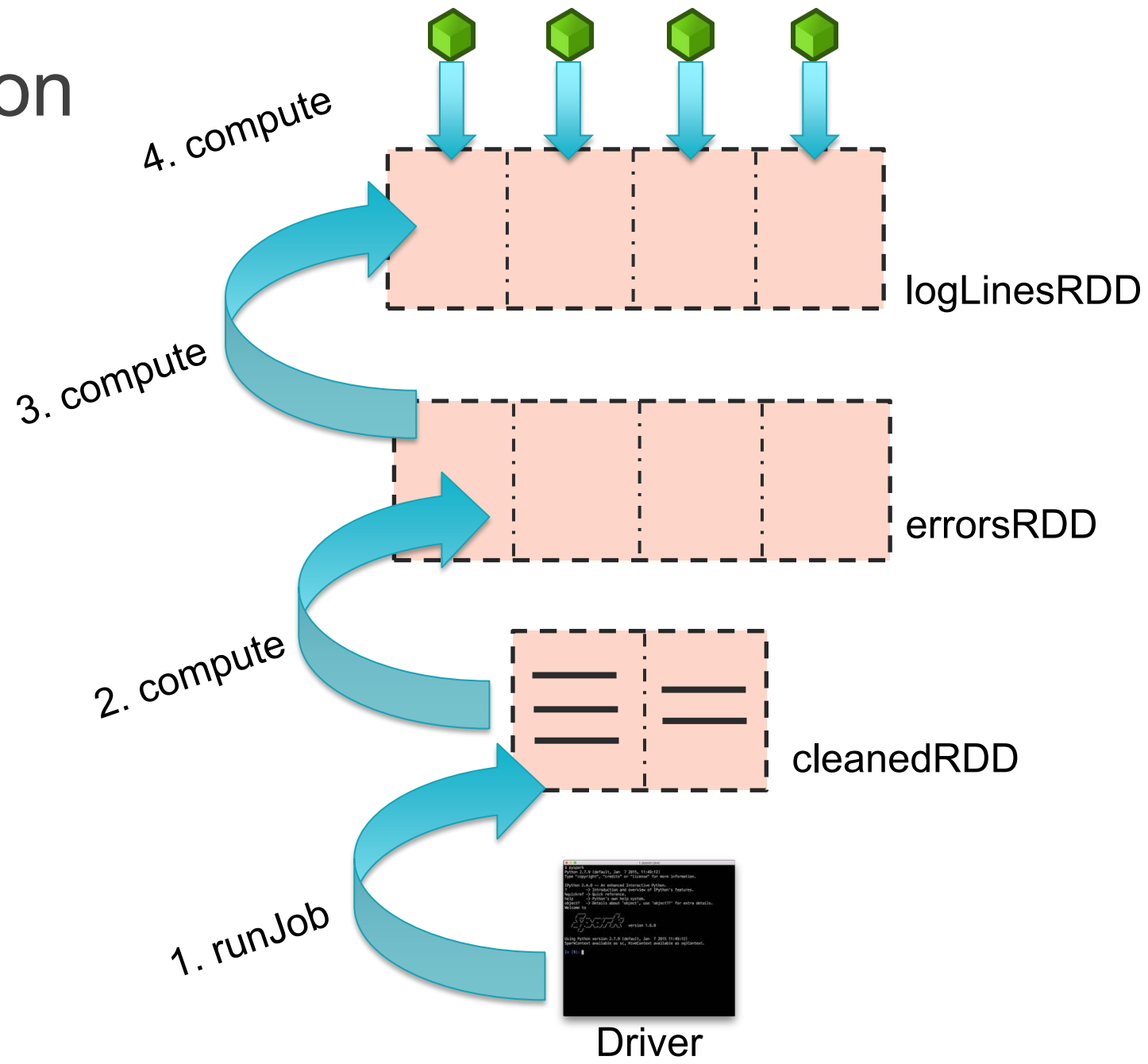
```
Python 2.7.9 (default, Jan 7 2015, 11:01:11)
Type "copyright", "credits" or "license()" for more information.
> Python 2.7.9 is an extended interactive Python.
> Help -> Introduction and overview of Python's features.
> Help -> Built-in modules.
> Help -> Python's own help system.
> Help -> Details about "object", with "object()" for extra details.
Welcome to
Spark version 1.6.0
Using Python version 2.7.9 (default, Jan 7 2015 11:01:11)
User-defined variables as in: http://www.databricks.com
In [1]:
```

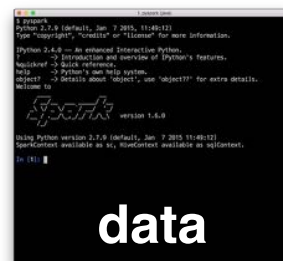
Driver

DAG



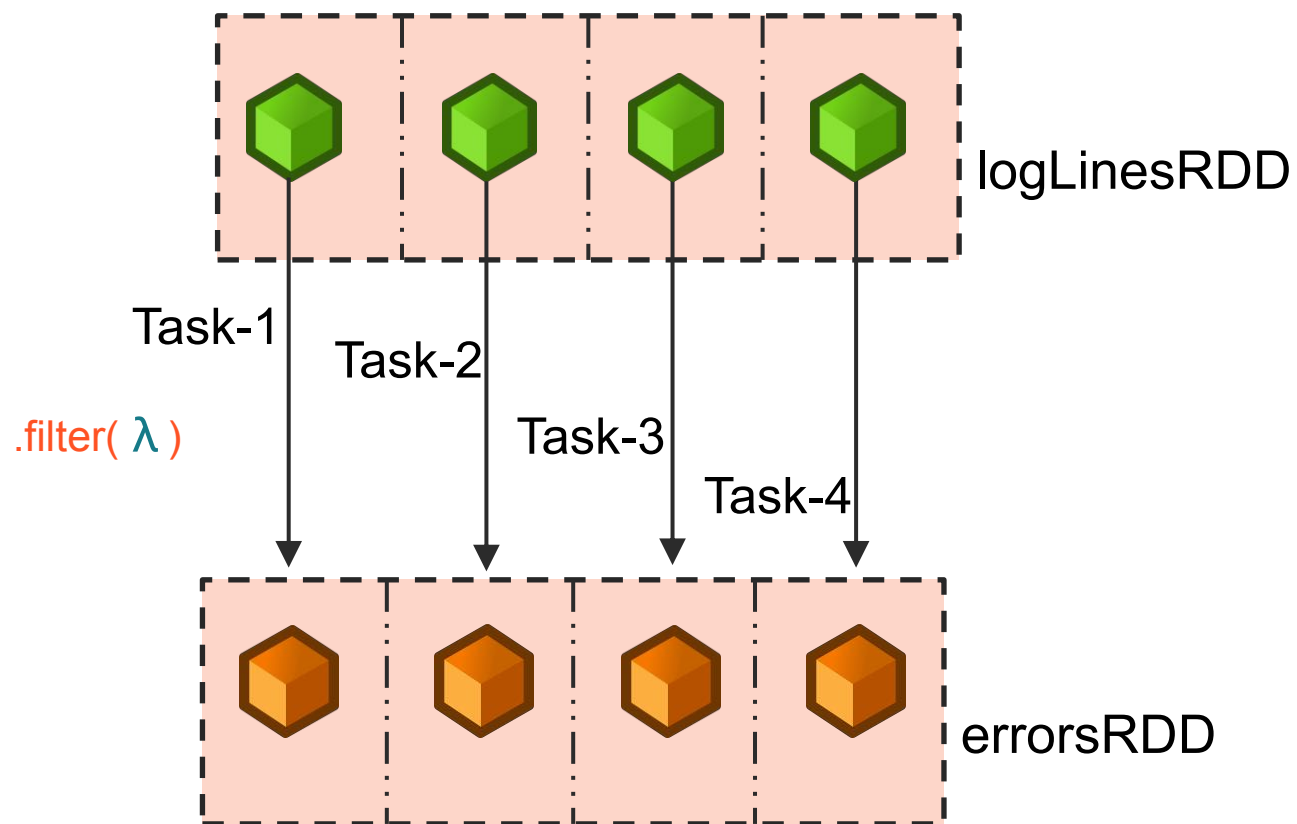
Execution





Driver

Partition >>> Task >>> Partition



Lifecycle of an RDD-based Spark Program

- 1) Create base RDD
- 2) Chain together transformations
- 3) Cache intermediate RDDs
- 4) Perform actions

Transformations

`map()`

`intersection()`

`zipWithIndex()`

`flatMap()`

`distinct()`

`pipe()`

`filter()`

`groupByKey()`

`coalesce()`

`mapPartitions()`

`reduceByKey()`

`...`

Actions

`reduce()`

`collect()`

`count()`

`first()`

`take()`

`takeOrdered()`

`saveAsTextFile()`

`...`

RDDs vs DataFrames

- RDDs provide a low-level interface into Apache Spark
- DataFrames have a schema
- DataFrames are cached using Tungsten format
- DataFrames are optimized via Catalyst
- DataFrames are built on top of the RDD and core APIs