# Apache Spark 101

Nirmal Ranganathan

# Spark

General purpose engine for

large-scale data processing

# Brief History

- 2009 - AMPLab at the University of California, Berkeley

- 2013 - Apache incubation

- 2014 - Top Level Project

- 600+ contributors

Why so popular?
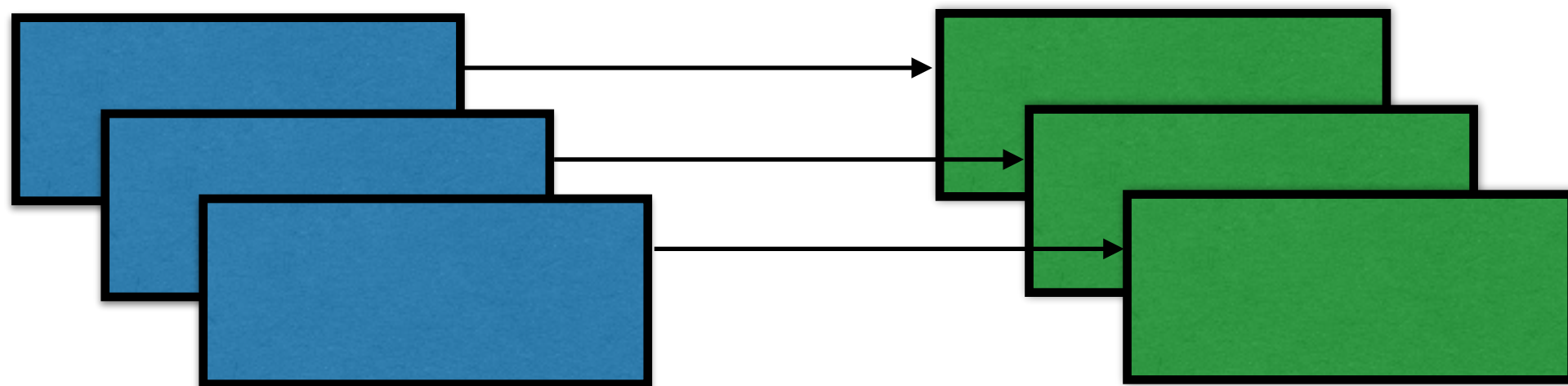
# RDDs

- **Transformations**

  - map

  - filter

  - groupByKey

  - ReduceByKey

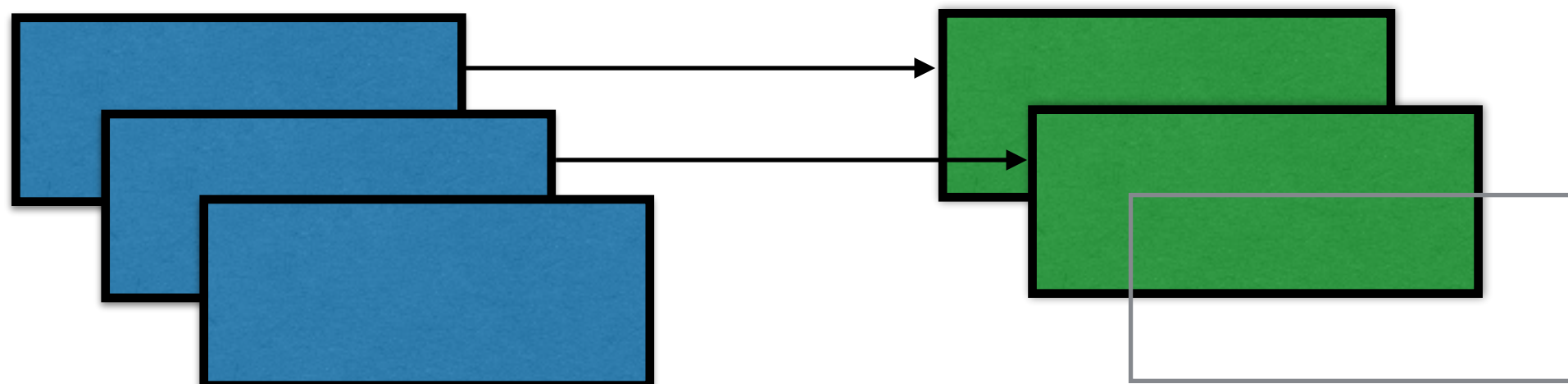- **Actions**

  - collect

  - reduce

  - count

  - save

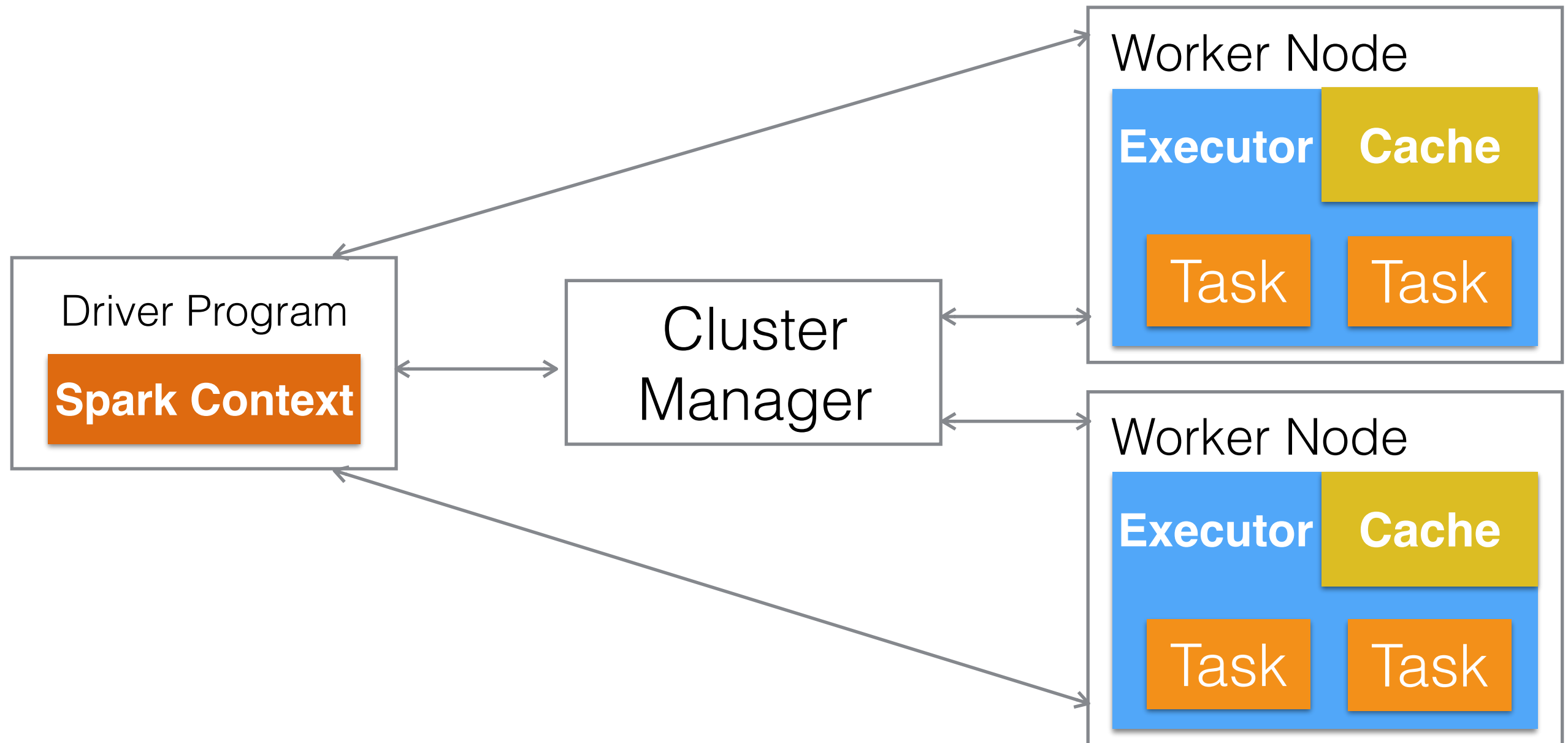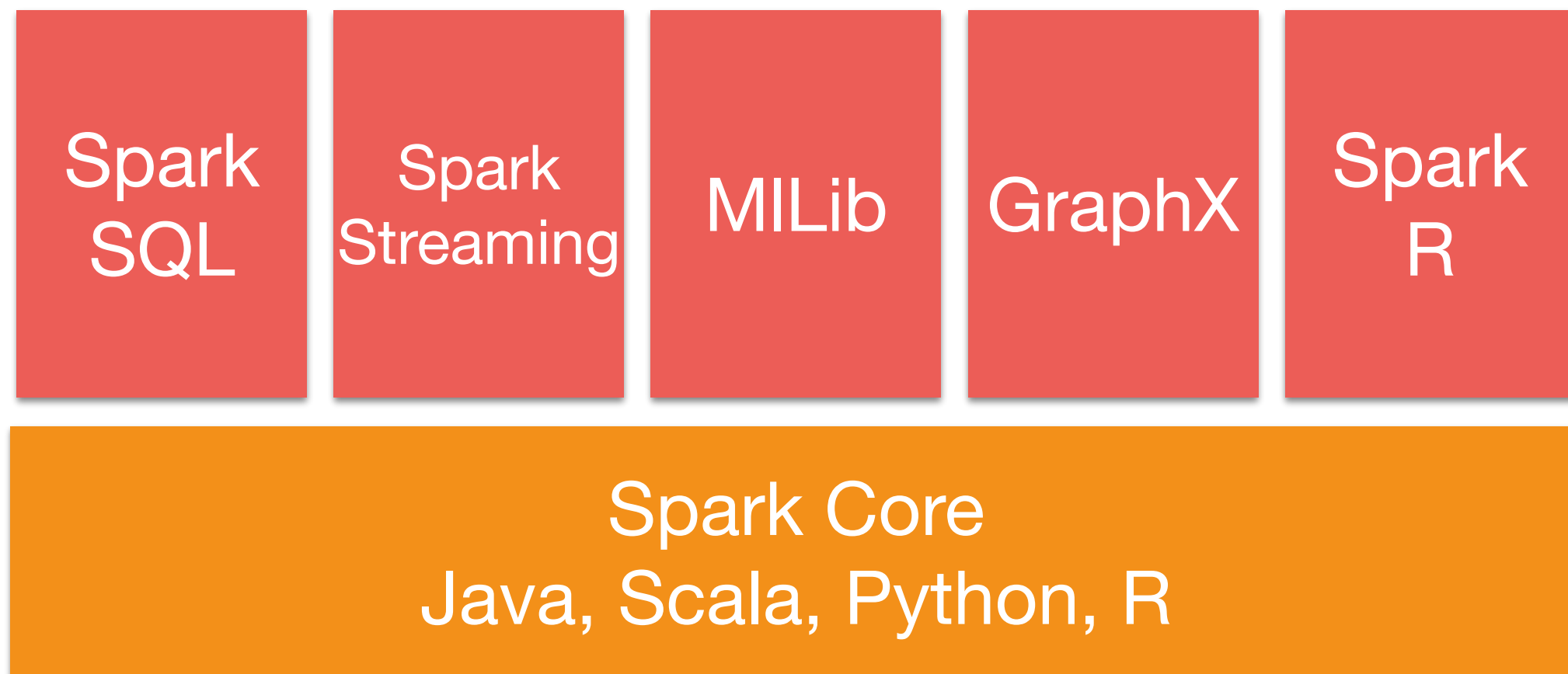# Transformation

map(func)

# Transformation

**filter(func)**

# Components

# Components

# Word Count

```python
text_file = sc.textFile("source.txt")

counts = text_file.flatMap(lambda line: line.split(" ")) \
                  .map(lambda word: (word, 1)) \
                  .reduceByKey(lambda a, b: a + b)

counts.saveAsTextFile("hdfs://host:port/results")
```

# Word Count

```scala
val text_file = sc.textFile("source.txt")

val counts = text_file.flatMap(line => line.split(" "))
                .map(word => (word, 1))
                .reduceByKey(_ + _)

counts.saveAsTextFile("hdfs://host:port/results")
```