

# AN INTRODUCTION TO SPARK AND TO ITS PROGRAMMING MODEL

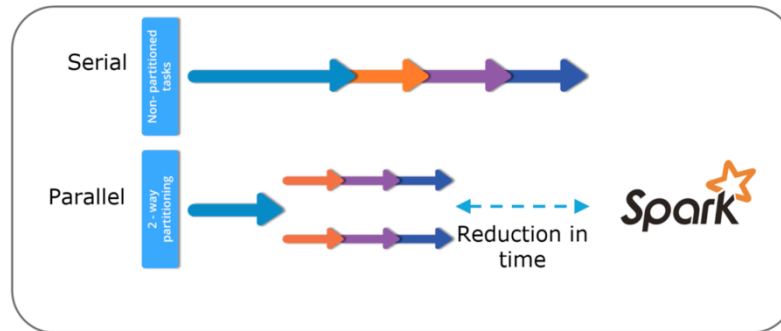
---

# Introduction to Apache Spark

- Fast, expressive cluster computing system compatible with Apache Hadoop
- It is much faster and much easier than Hadoop MapReduce to use due its rich APIs
- Large community
- Goes far beyond batch applications to support a variety of workloads:
  - including interactive queries, streaming, machine learning, and graph processing



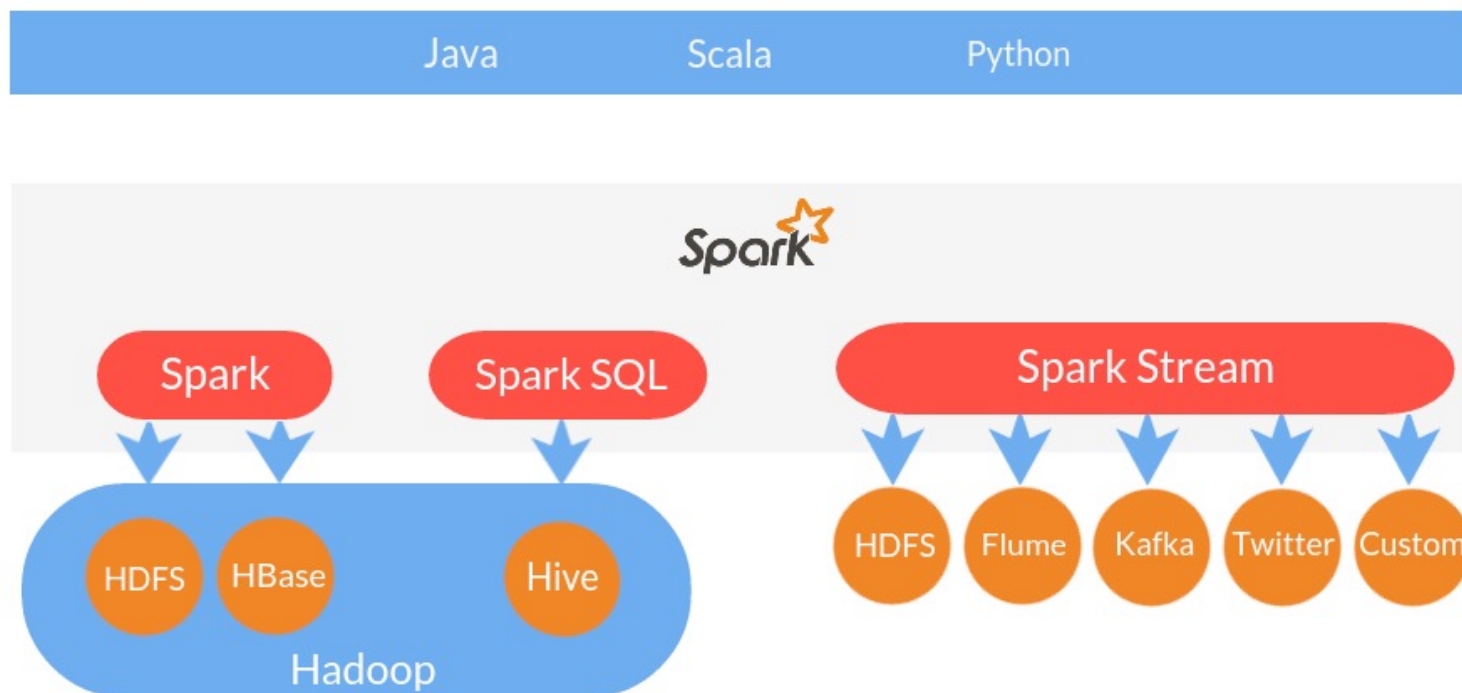
**Figure:** Real Time Processing In Spark



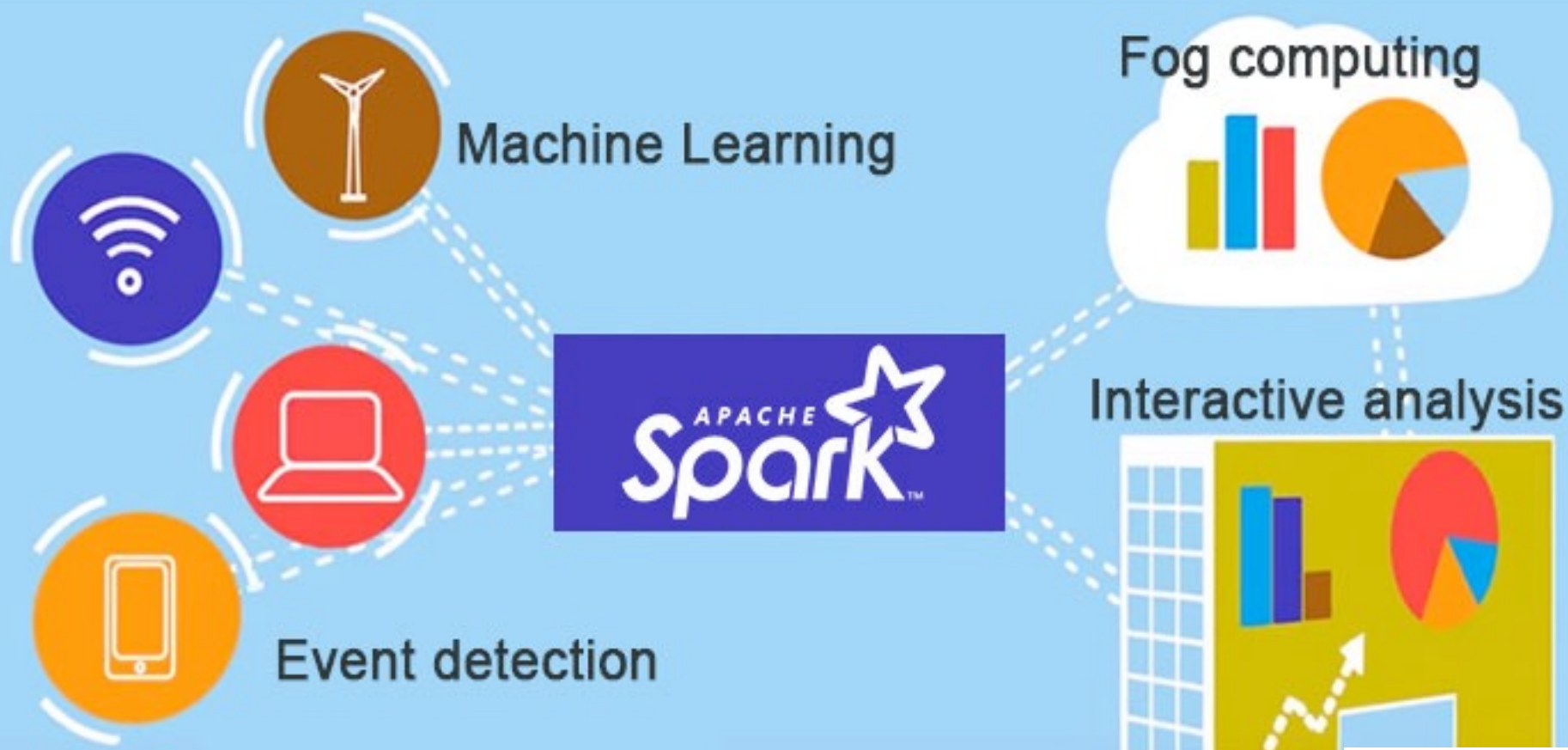
**Figure:** Data Parallelism In Spark

# Introduction to Apache Spark

- General-purpose cluster in-memory computing system
- Provides high-level APIs in Java, Scala, python



# Apache Spark Applications



## Uses Cases

The Netflix logo, featuring the word "NETFLIX" in a bold, white, sans-serif font with a black outline, set against a solid red background.

Uses Spark Streaming to provide the best-in-class movie streaming and recommendation tool to its users.



Uses Spark to collect TBs of raw and unstructured data every day from its users to convert it into structured data. This makes it ready for further complex analytics.



Feeds real-time data into Spark via Spark Streaming to get instant insights on how users are engaging with Pins globally. This makes Pinterest's recommendations (i.e. to show Pins) to be accurate.



# Uses Cases

## Spark Use Cases

edureka!



### Twitter Sentiment Analysis With Spark

Trending Topics can be used to create campaigns and attract larger audience

Sentiment helps in crisis management, service adjusting and target marketing



### NYSE: Real Time Analysis of Stock Market Data



### Banking: Credit Card Fraud Detection

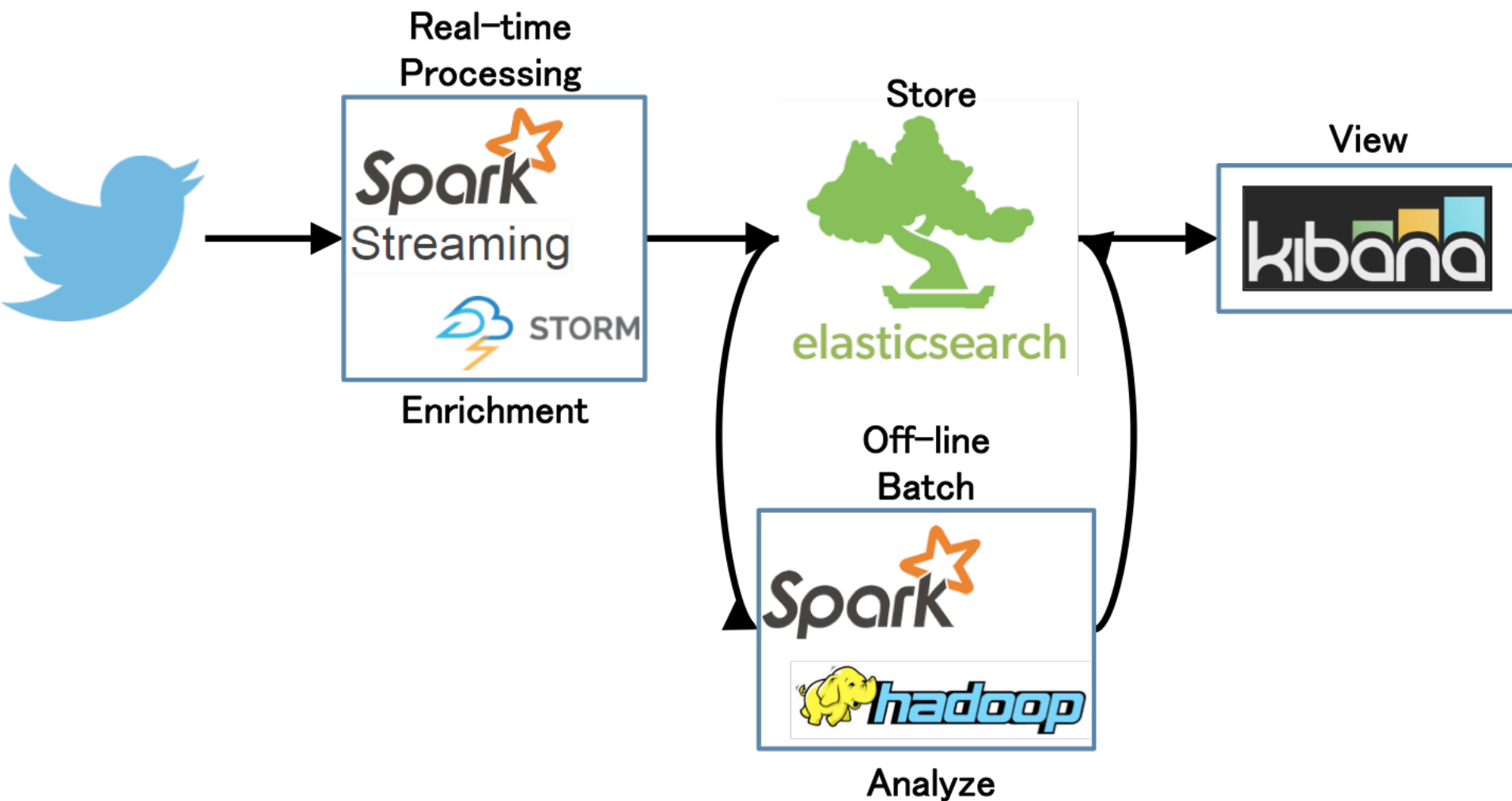


### Genomic Sequencing

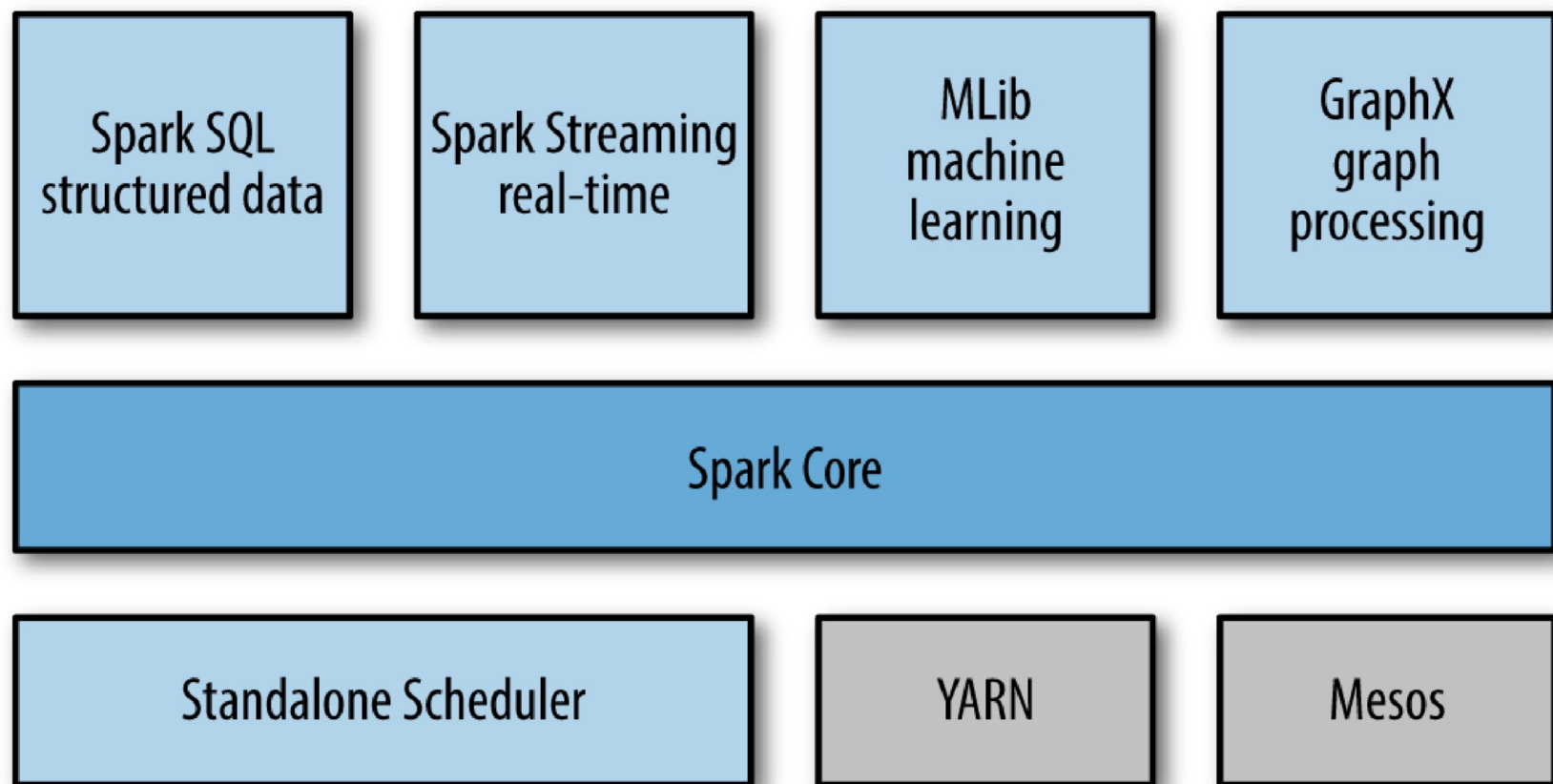


# Real Time Data Architecture for analyzing tweets

## - Twitter Sentiment Analysis



# Spark Ecosystem





## Spark Core

- Contains the basic functionality for
  - task scheduling,
  - memory management,
  - fault recovery,
  - interacting with storage systems,
  - and more.
- Defines the Resilient Distributed Data sets (RDDs)
  - main Spark programming abstraction.

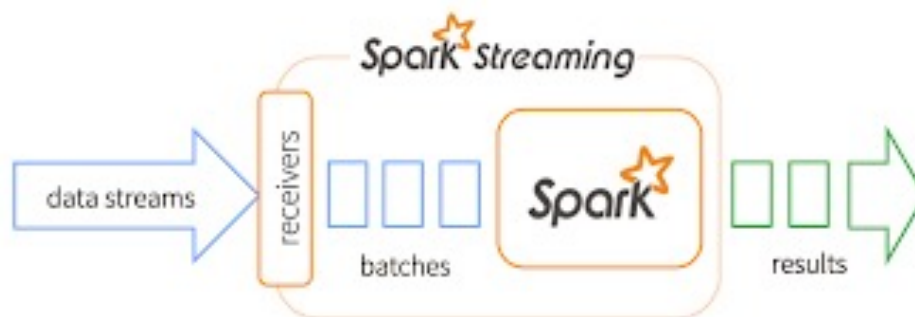
## Spark SQL

- For working with structured data.
- View datasets as relational tables
- Define a schema of columns for a dataset
- Perform SQL queries
- Supports many sources of data
  - Hive tables, Parquet and JSON
- DataFrame



# Spark Streaming

- Data analysis of streaming data
  - e.g. log files generated by production web servers
- Aimed at high-throughput and fault-tolerant stream processing
- Dstream → Stream of datasets that contain data from certal interval

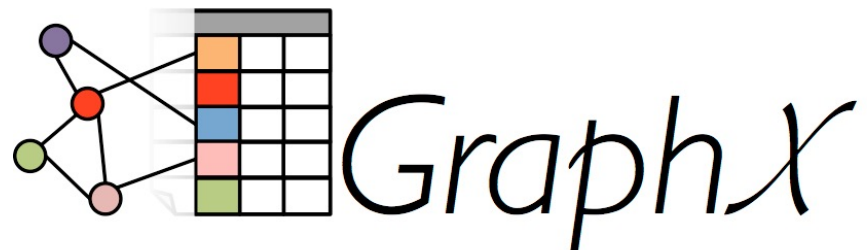


## Spark MLlib

- MLlib is a library that contains common Machine Learning (ML) functionality:
  - Basic statistics
  - Classification (Naïve Bayes, decision tress, LR)
  - Clustering (k-means, Gaussian mixture, ...)
  - And many others!
- All the methods are designed to scale out across a cluster.

# Spark GraphX

- Graph Processing Library
- Defines a graph abstraction
  - Directed multi-graph
  - Properties attached to each edge and vertex
  - RDDs for edges and vertices
- Provides various operators for manipulating graphs (e.g. subgraph and mapVertices)



## Programming with Python Spark (pySpark)

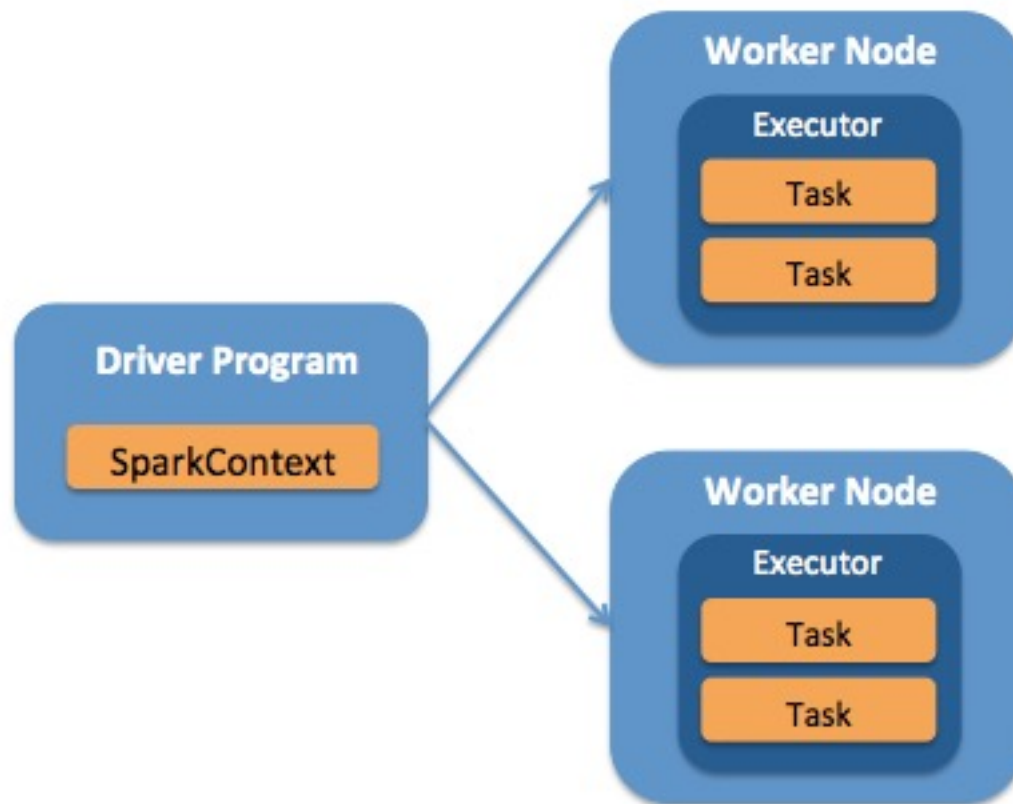
- We will use Python's interface to Spark called **pySpark**
- **A driver program** accesses the Spark environment through a **SparkContext** object
- Their **key concept** in Spark are datasets called **RDDs** (**R**esilient **D**istributed **D**ataset )
- Basic idea: We load our data into RDDs and perform some **operations**



## Programming environment - Spark concepts

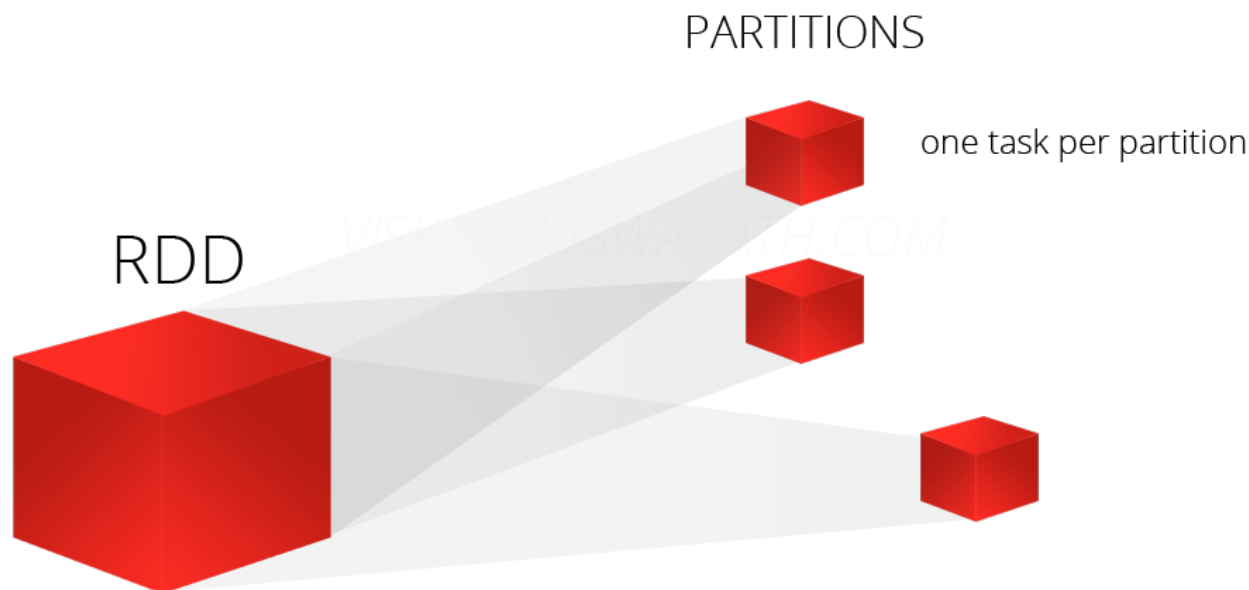
- Driver programs access Spark through a **SparkContext** object which represents a connection to the computing cluster.
- In a *shell* the **SparkContext** is created for you and available as the variable **sc**.
- You can use it to build **Resilient Distributed Data (RDD)** objects.
- Driver programs manage a number of *worker nodes* called *executors*.

## Programming environment- Spark concepts



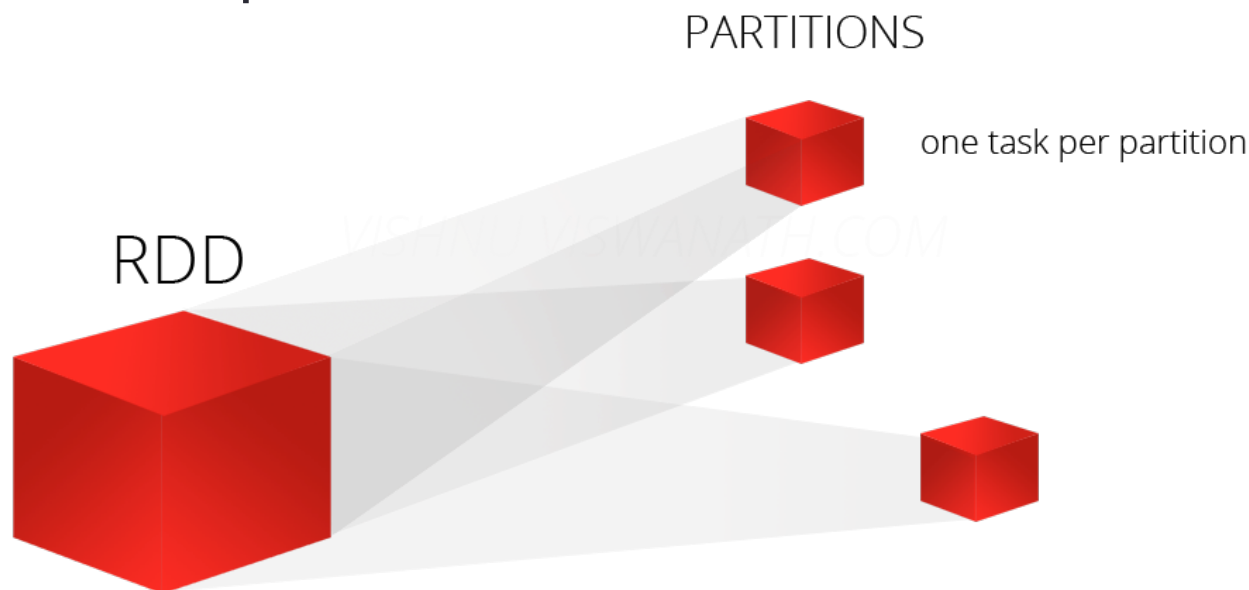
## RDD abstraction

- Represent data or transformations on data
- It is distributed collection of items - partitions
- Read-only → they are immutable
- Enables operations to be performed in parallel



## RDD abstraction

- Fault tolerant:
  - Lineage of data is preserved, so data can be re-created on a new node at any time
- Caching dataset in memory
  - different storage levels available
  - fallback to disk possible

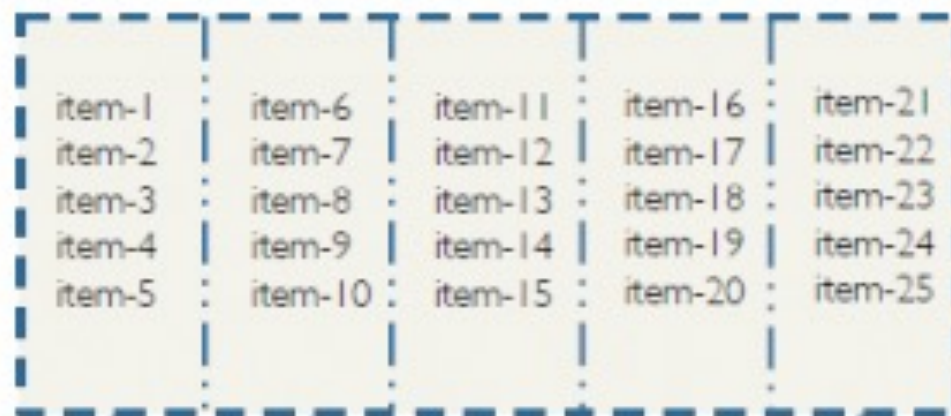


## Programming with RDDs

- All work is expressed as either:
  - creating new RDDs
  - transforming existing RDDs
  - calling operations on RDDs to compute a result.
- Distributes the data contained in RDDs across the nodes (executors) in the cluster and parallelizes the operations.
- Each RDD is split into multiple **partitions**, which can be computed on different nodes of the cluster.

# Partition

RDD split into 5 partitions



3 workers





# An RDD can be created 2 ways

- Parallelize a collection

```
# Parallelize in Python
wordsRDD = sc.parallelize(["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 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.

# First Program!

```
sc = SparkContext(master="local[*]")
```

**Context**

```
lines = sc.textFile("README.md", 4)
```

**RDD-1**

```
lines.count()
```

```
pythonLines = lines.filter(lambda line : "Python" in line)
```

**RDD-2**

```
pythonLines.first()
```

## RDD operations

- Once created, RDDs offer two types of operations:
  - **transformations**
    - transformations include *map*, *filter*, *join*
    - lazy operation to build RDDs from other RDDs
  - **actions**
    - actions include *count*, *collect*, *save*
    - return a result or write it to storage

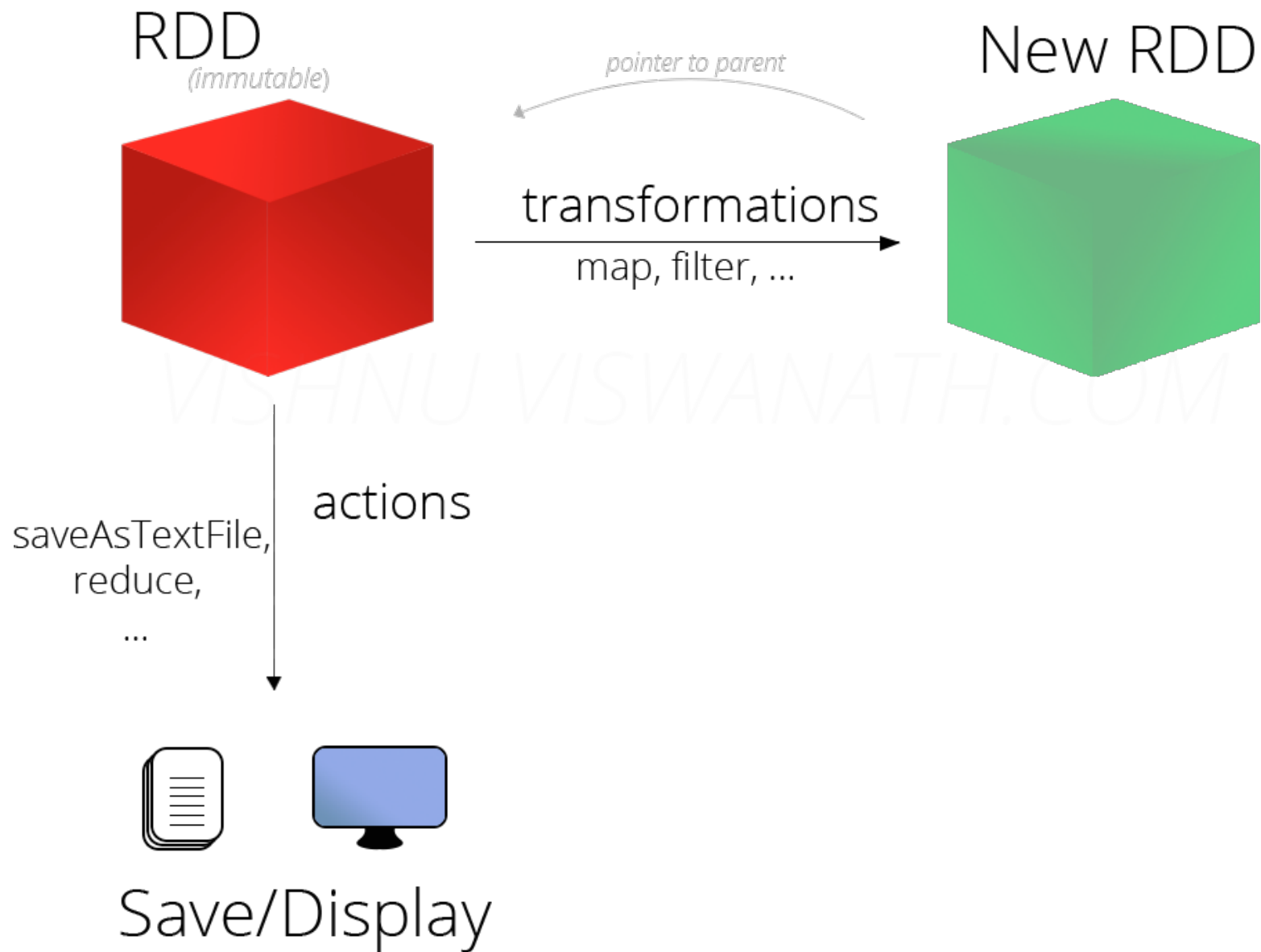
## Transformation vs Actions

### Transformations

```
map(func)
flatMap(func)
filter(func)
groupByKey()
reduceByKey(func)
mapValues(func)
sample(...)
union(other)
distinct()
sortByKey()
...
```

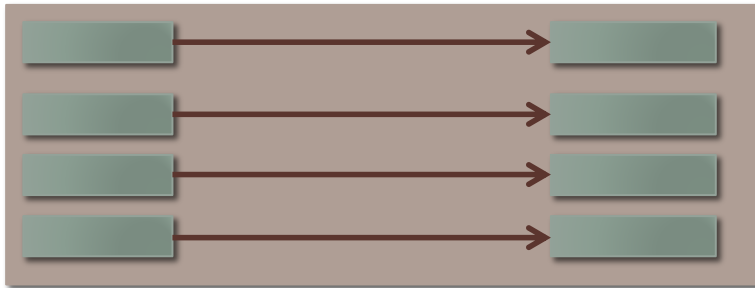
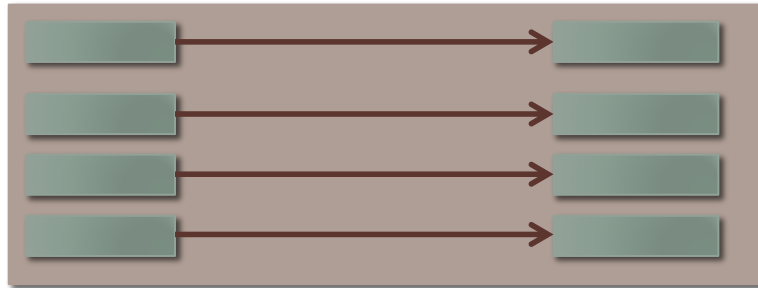
### Actions

```
reduce(func)
collect()
count()
first()
take(n)
saveAsTextFile(path)
countByKey()
foreach(func)
...
```



# Narrow Vs. Wide transformation

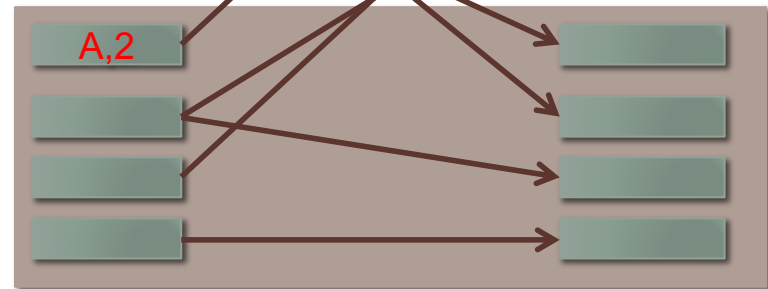
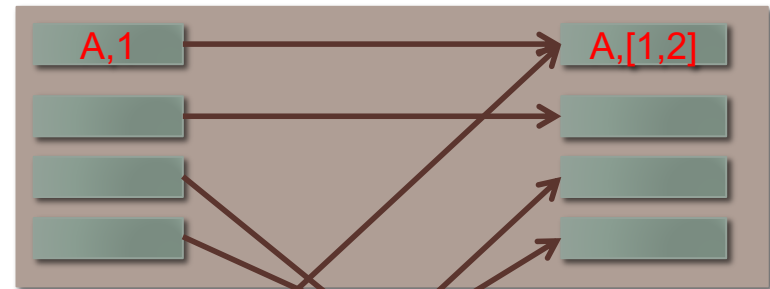
Narrow



Map

Vs.

Wide



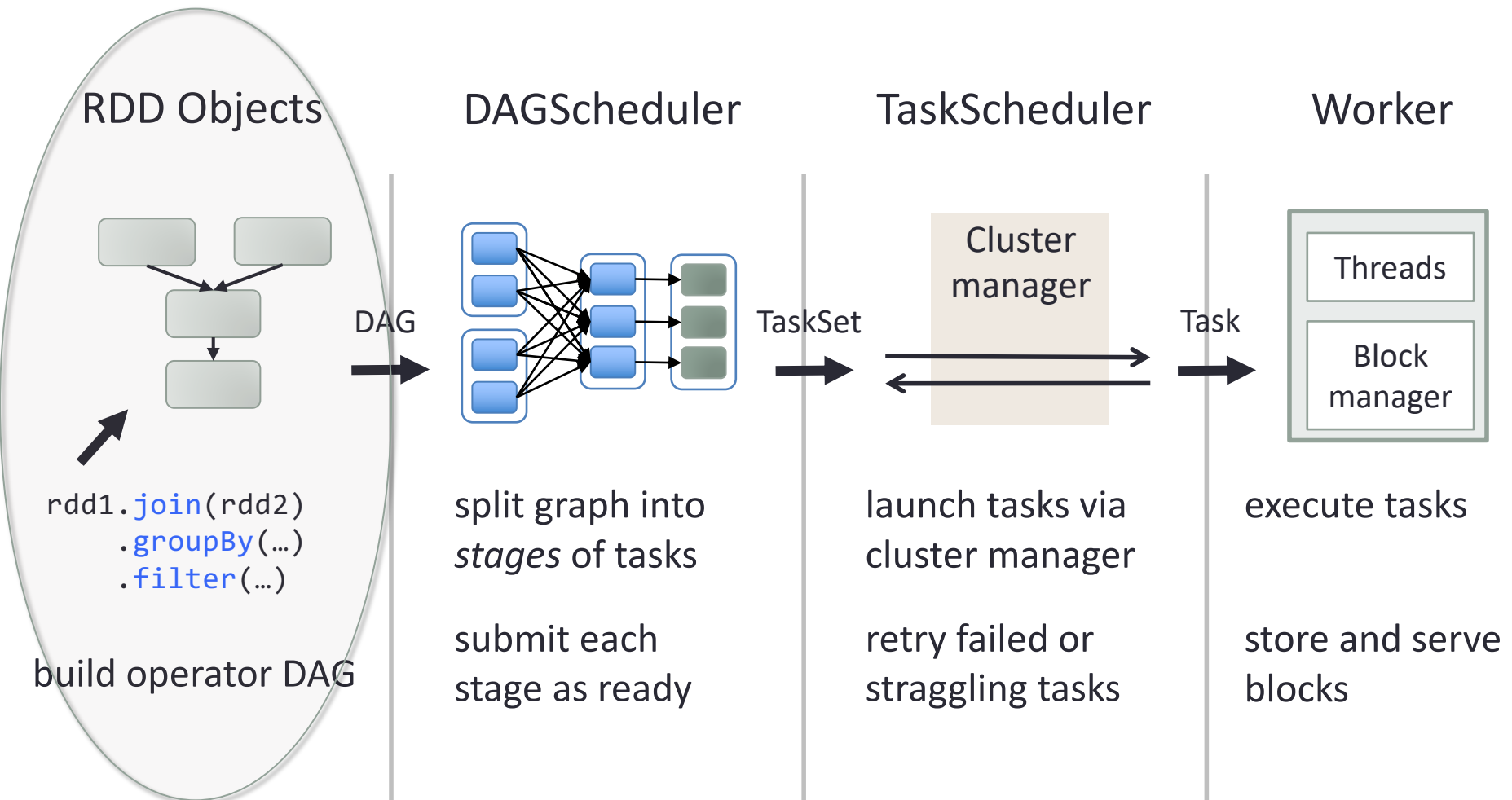
groupByKey



# Life cycle of 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()`
  - 1) Ask Spark to `cache()` any intermediate RDDs that will need to be reused.
  - 2) Launch actions such as `count()` and `collect()` to kick off a parallel computation, which is then optimized and executed by Spark.

# Job scheduling



## Example: Mining Console Logs

- Load error messages from a log into memory, then interactively search for patterns

```
lines = spark.textFile("hdfs://...")
errors = lines.filter(lambda s: s.startswith("ERROR"))
messages = errors.map(lambda s: s.split('\t')[2])
messages.cache()
```

```
messages.filter(lambda s: "foo" in s).count()
messages.filter(lambda s: "bar" in s).count()
...

```

Action

Base  
RDD

Transformed  
RDD

Cache 1

Worker

Block 1

tasks

results

Driver

Cache 2

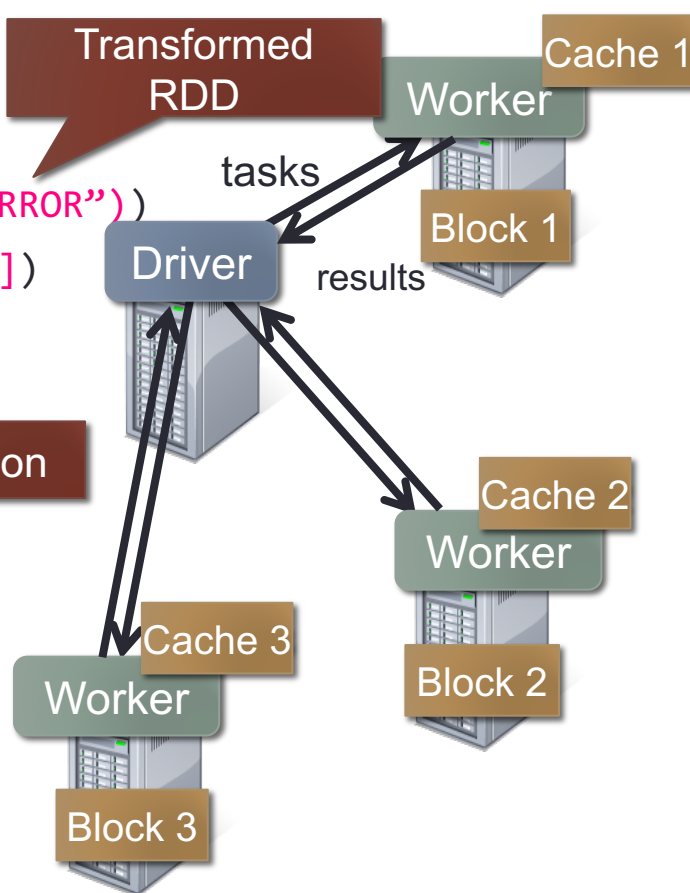
Worker

Block 2

Cache 3

Worker

Block 3



## Some Apache Spark tutorials

- <https://www.cloudera.com/documentation/enterprise/5-7-x/PDF/cloudera-spark.pdf>
- [https://stanford.edu/~rezab/sparkclass/slides/itas\\_workshop.pdf](https://stanford.edu/~rezab/sparkclass/slides/itas_workshop.pdf)
- <https://www.coursera.org/learn/big-data-essentials>
- <https://www.cloudera.com/documentation/enterprise/5-6-x/PDF/cloudera-spark.pdf>

# Spark: when not to use

- Even though Spark is versatile, that doesn't mean Spark's in-memory capabilities are the best fit for all use cases:
  - For many simple use cases Apache MapReduce and Hive might be a more appropriate choice
  - Spark was not designed as a multi-user environment
  - Spark users are required to know that memory they have is sufficient for a dataset
  - Adding more users adds complications, since the users will have to coordinate memory usage to run code