

Introduction to Apache Spark

Architecture, RDD, DataFrames, SQL, Spark
Streaming & MLlib

Prasanth Kothuri, CERN

IT-Database Group

Hadoop and Spark Service

Course Outline

What is Apache Spark?

Spark Abstractions

Spark Architecture

Spark Data APIs

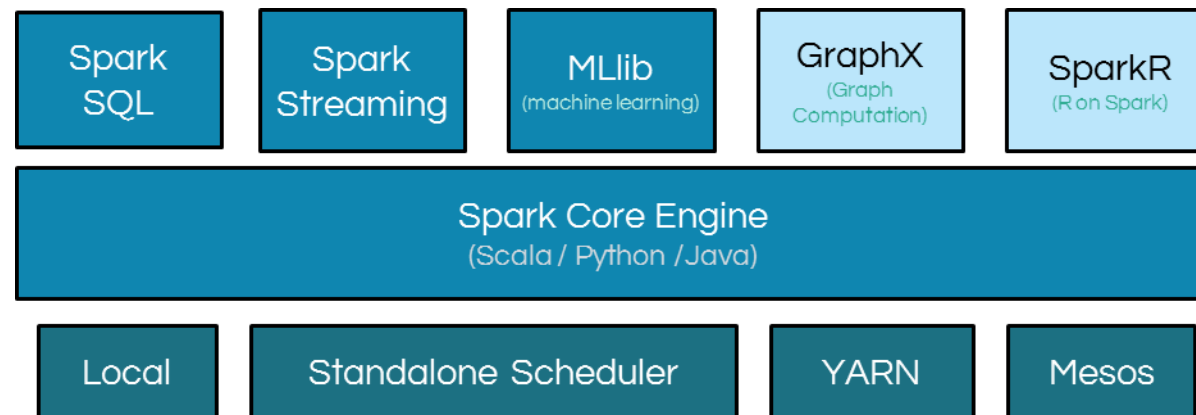
RDD, Dataframe and SQL

Spark Streaming

Spark Mlib

What is Apache Spark?

- Apache Spark is an open-source *parallel processing framework* with expressive development APIs (in multiple languages) that allows for *sophisticated analytics, real-time streaming* and *machine learning on large datasets*
- Spark ecosystem



Support for kubernetes since Spark 2.3

Apache Spark's Philosophy

- Unified: provide a unified platform for writing big data applications
 - Data loading, SQL queries, streaming computation and machine learning
- Computing engine: limits the scope to computations
 - Only handles loading data from storage systems and performing computations on it
 - Can be used with wide variety of persistent storage systems, including HDFS, Amazon S3, Azure Storage and EOS
- Libraries: In addition to standard libraries, spark supports wide array of external libraries (spark-packages.org)

Spark Abstractions

Two main abstractions of Spark

- RDD – Resilient Distributed Dataset [1]
- DAG – Direct Acyclic Graph

[1] DataFrames – higher level structured APIs

RDD

Resilient Distributed Datasets (RDDs) are the primary abstraction in Spark – a immutable distributed collection of records that can be operated on in parallel

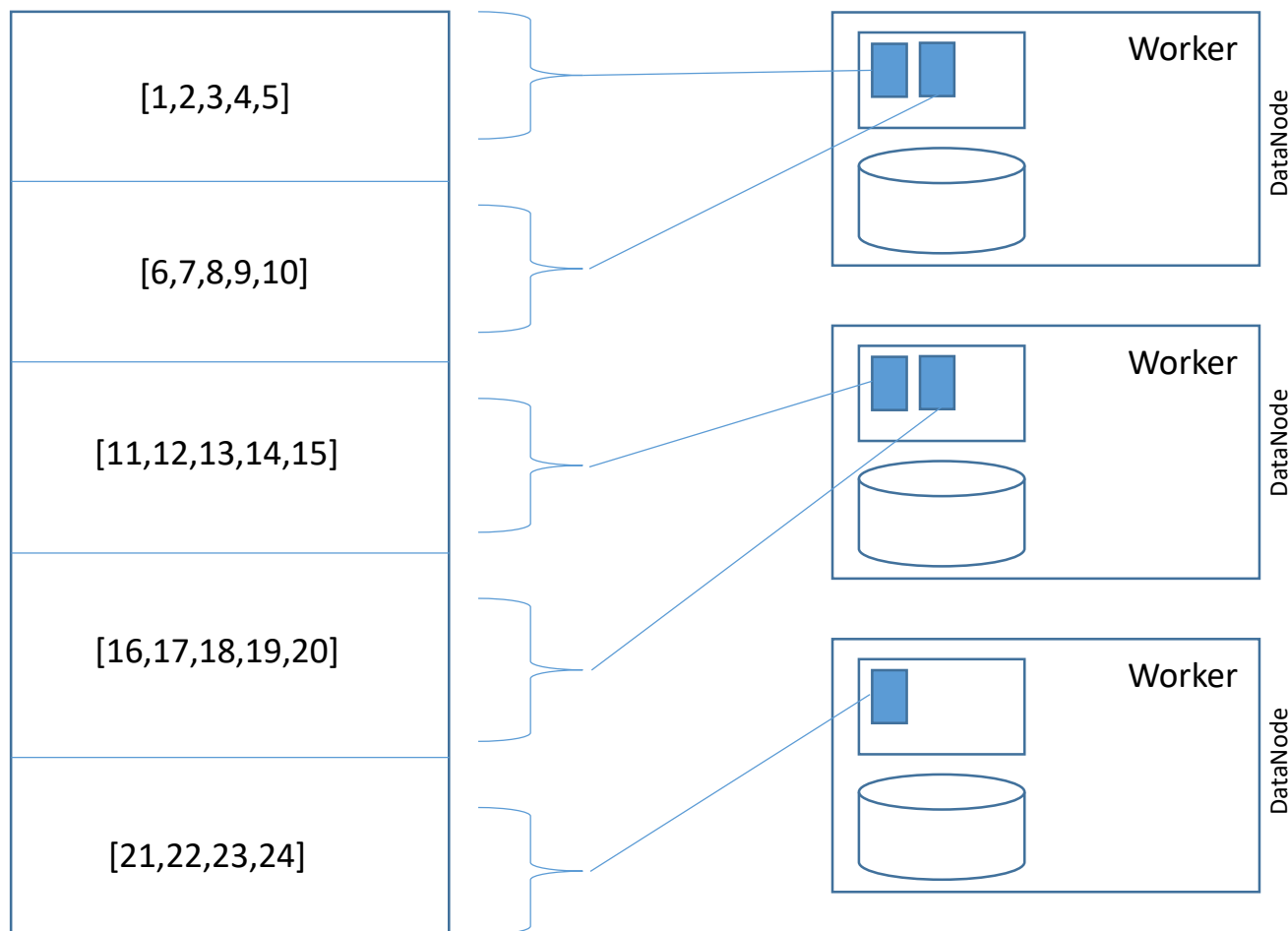
There are currently two ways to create them:

- **parallelized collections** – take an existing python/scala collection and run functions on it in parallel
- **External datasets** – run functions on each record of a file(s) in External/Hadoop distributed file system

RDD: Example

```
rdd = spark.parallelize(range(1,25),5)
```

rdd is split into
partitions



Partitioned: RDD is partitioned and distributed across worker nodes of the cluster

In-Memory : RDD is stored in memory as much (size) and long (time) as possible

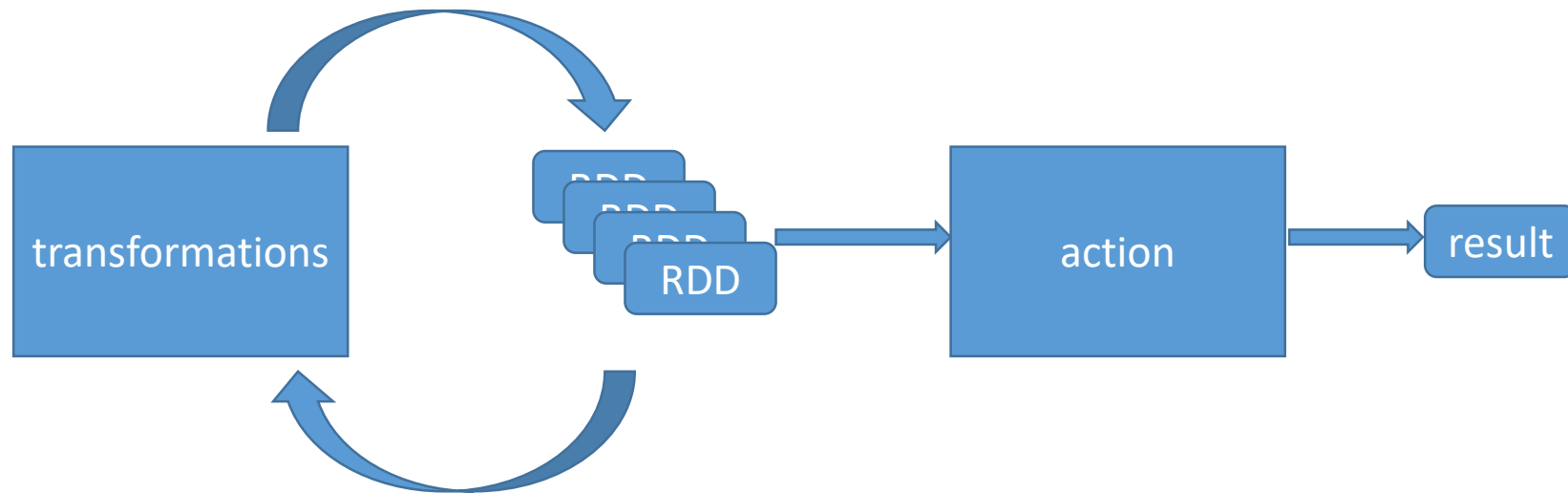
Immutable: does not change once created, can only be transformed into new RDDs

Typed: RDDs have types

Cacheable: hold all the data in a persistent storage like memory (preferable) or disk

RDD: Actions and Transformations

- Two types of operations on RDDs:
transformations and actions
- transformations – lazy operations that return another RDD
- actions – operations that trigger computation and return value



DAG

Direct Acyclic Graph – sequence of computations performed on data

- Node – RDD partition
- Edge – transformation on top of data
- Acyclic – graph cannot return to the older partition
- Direct – transformation is an action that transitions data partition state (from X to Y)

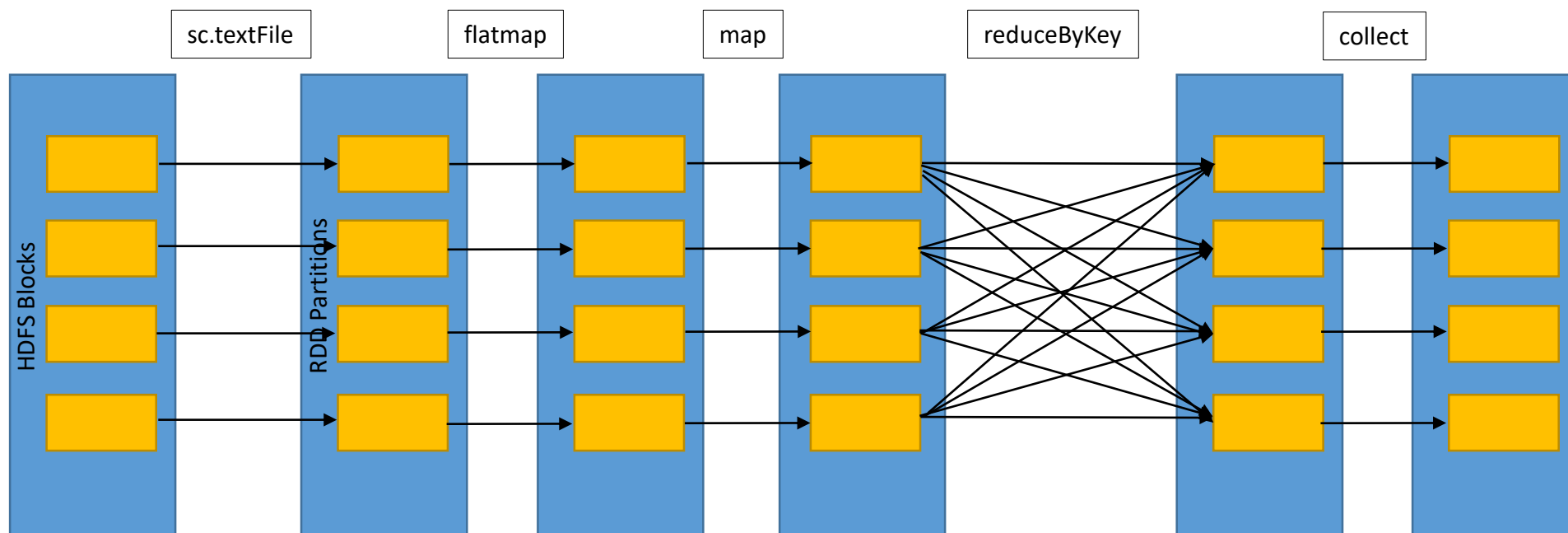
DAG

WordCount example

```
text_file =  
sc.textFile("hdfs:///user/pkothuri/sparkTraining/LICENSE.txt")  
counts = text_file.flatMap(lambda line: line.split(" ")) \  
                    .map(lambda word: (word, 1)) \  
                    .reduceByKey(lambda a, b: a + b)  
for x in counts.collect():  
    print x
```

DAG

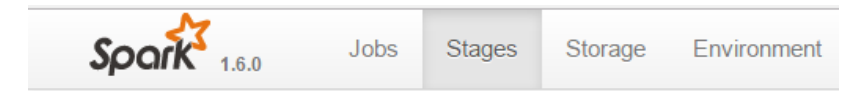
WordCount example



Spark WEB UI

Every SparkContext launches a web UI, by default on port 4040, that displays useful information about the application. This includes:

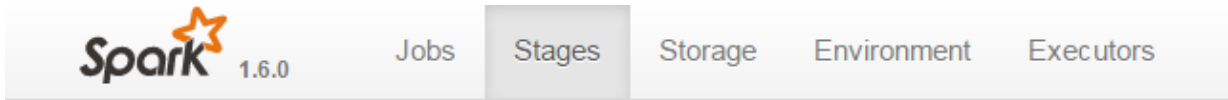
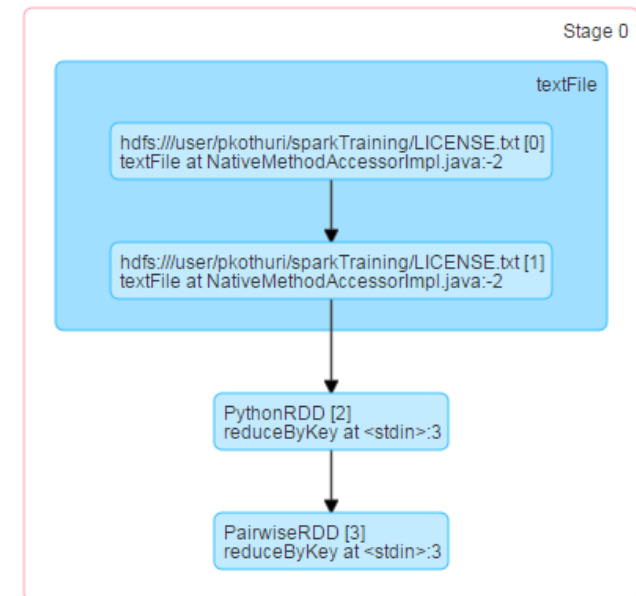
- A list of scheduler stages and tasks
- A summary of RDD sizes and memory usage
- Environmental information.
- Information about the running executors



Details for Stage 0 (Attempt 0)

Total Time Across All Tasks: 7 s
Locality Level Summary: Rack local: 2
Input Size / Records: 5.5 KB / 202
Shuffle Write: 9.8 KB / 24

[DAG Visualization](#)



Stages for All Jobs

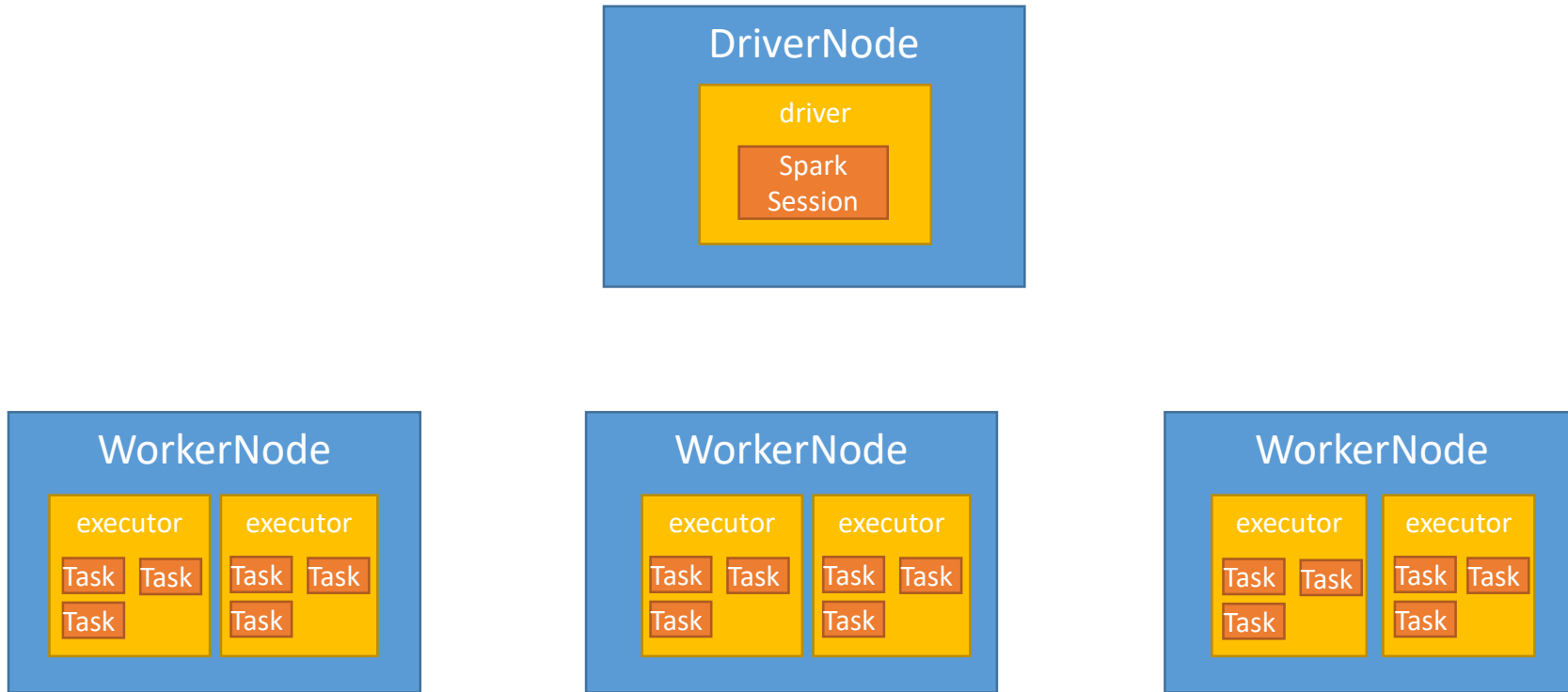
Completed Stages: 2

Completed Stages (2)

Stage Id	Description
1	collect at <stdin>:1
0	reduceByKey at <stdin>:3

Spark Architecture

Spark Cluster



Spark Architecture

Driver

- Entry point for Spark Shell (Scala, Python)
- SparkContext is created here and resides here
- Graph is built and submitted to DAGScheduler
- DAGScheduler divides it into stages and tasks
- Schedules tasks and controls their execution

Spark Architecture

Executor

- Reads data from HDFS (or external sources)
- Stores the data in cache in JVM heap or on HDDs
- Performs all data processing
- Writes data to HDFS (or external sources)

Spark Architecture

Application Decomposition

- Application
 - Single instance of SparkContext that stores some data processing logic and can schedule series of jobs, sequentially or in parallel
- Job
 - Complete set of transformations on RDD that finishes with action or data saving, triggered by the driver application

Spark Architecture

Application Decomposition

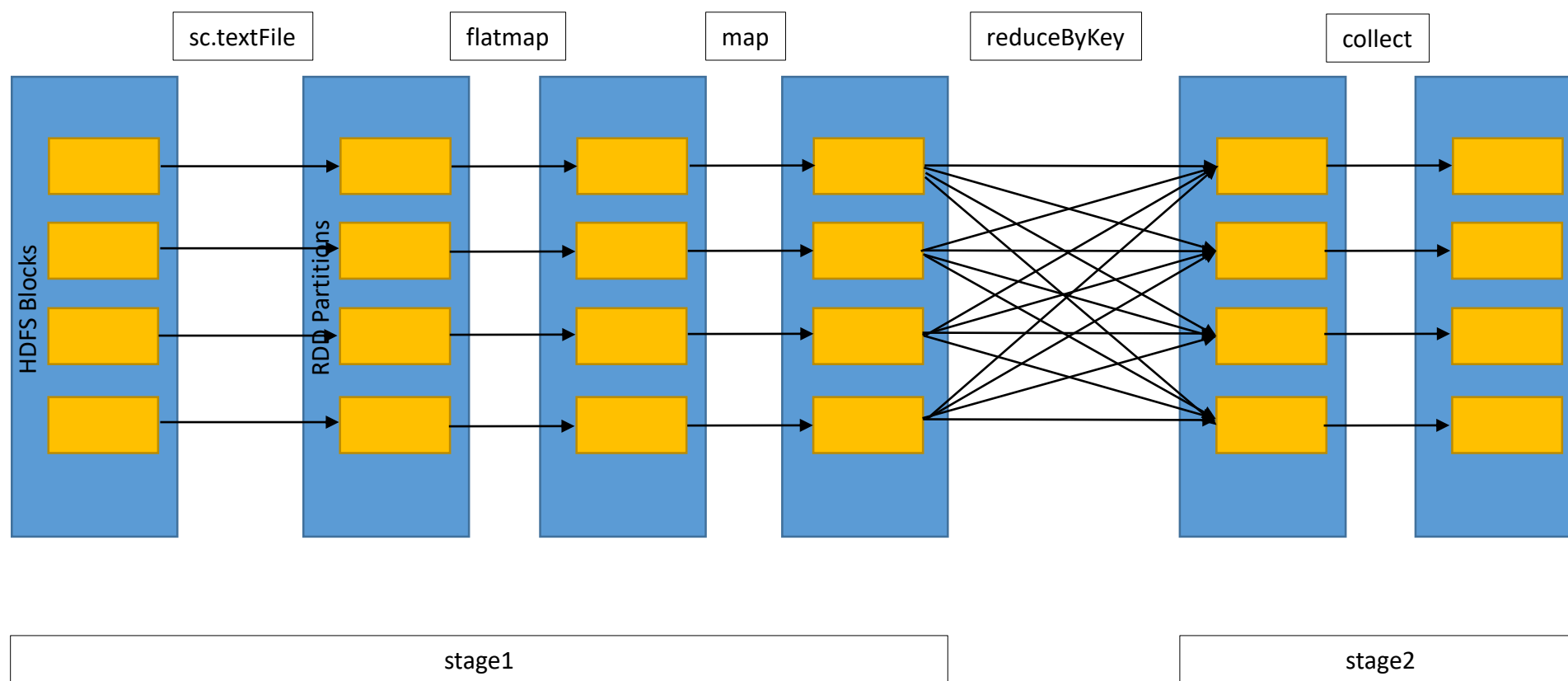
- Stage
 - Set of transformations that can be pipelined and executed by a single independent worker. Usually it is app the transformations between «read», «shuffle», «action» and «save»
- Task
 - Execution of the stage on a single data partition. Basic unit of scheduling

DAG scheduler

- When an action is called on the RDD, Spark creates DAG and submits to the DAG scheduler
- The DAG scheduler divides operators into stages of tasks
- The stages are created based on the transformations, the narrow transformations are grouped (pipelined) into a single stage
- The DAG scheduler submits the stages to the task scheduler
 - The number of tasks depend on number of partitions
 - The number of tasks submitted depends on number of available executors

DAG

WordCount example



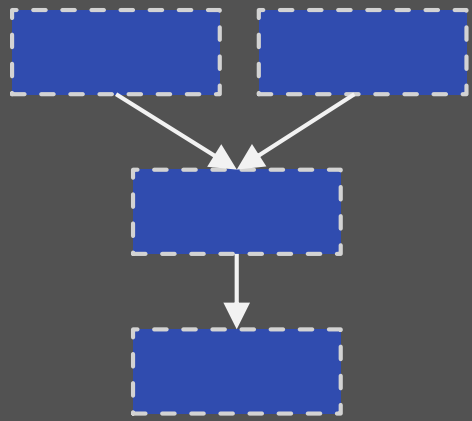
SCHEDULING PROCESS

RDD Objects

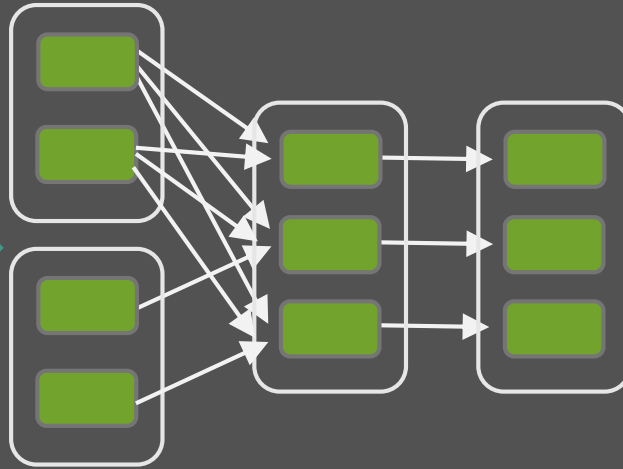
DAG Scheduler

Task Scheduler

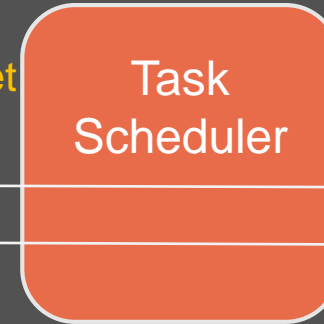
Executor



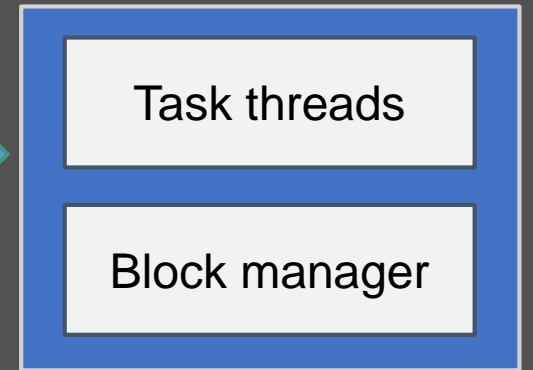
DAG




TaskSet



Task



 Rdd1.join(rdd2)
.groupBy(...)
.filter(...)

- Build operator DAG

- Split graph into *stages* of tasks
- Stage boundaries are at wide transformations (shuffle)
- Submit each stage as ready

- Launches individual tasks
- Retry failed or straggling tasks

- Execute tasks
- Store and serve blocks

Essential Spark RDD Operations

TRANSFORMATIONS

General

- map
- filter
- flatMap
- mapPartitions
- groupBy
- sortBy
- flatMapValues
- groupByKey
- reduceByKey
- foldByKey
- sortByKey
- combineByKey

Math / Stats

- sample
- sampleByKey
- randomSplit

Set Theory

- union
- intersection
- subtract
- distinct
- cartesian
- zip
- join
- fullOuterJoin
- leftOuterJoin
- rightOuterJoin

Data Structure / I/O

- keyBy
- zipWithIndex
- zipWithUniqueId
- zipPartitions
- coalesce
- repartition
- repartitionAndSortWithinPartitions
- pipe
- partitionBy

ACTIONS

- reduce
- collect
- aggregate
- fold
- first
- take
- foreach
- top
- treeAggregate
- treeReduce
- foreachPartition
- collectAsMap
- keys
- values

- count
- takeSample
- max
- min
- sum
- histogram
- mean
- variance
- stdev
- sampleVariance
- countApprox
- countApproxDistinct

- takeOrdered

- saveAsTextFile
- saveAsSequenceFile
- saveAsObjectFile
- saveAsHadoopDataset
- saveAsHadoopFile
- saveAsNewAPIHadoopDataset
- saveAsNewAPIHadoopFile

Transformations: **map** and **filter**

MAP

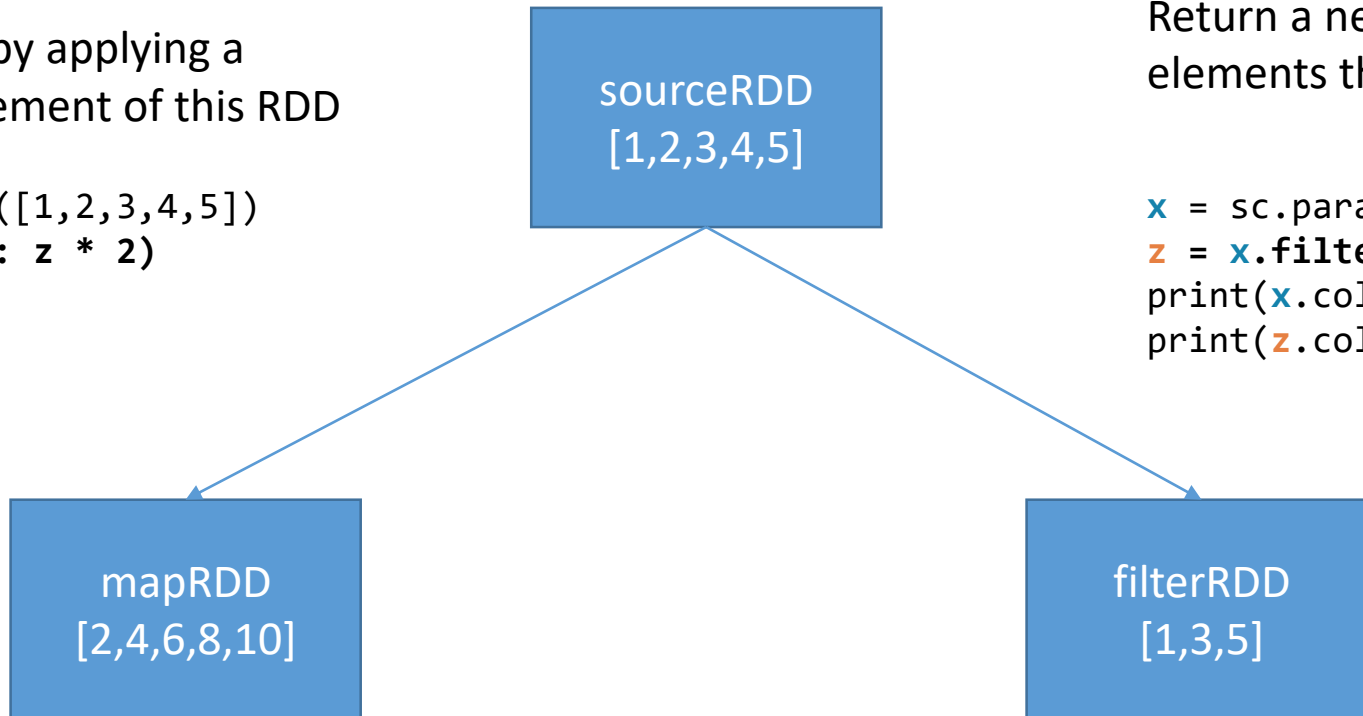
Return a new RDD by applying a function to each element of this RDD

```
x = sc.parallelize([1,2,3,4,5])
y = x.map(lambda z: z * 2)
print(x.collect())
print(y.collect())
```

FILTER

Return a new RDD that only has elements that pass the filter() function

```
x = sc.parallelize([1,2,3,4,5])
z = x.filter(lambda z: z % 2 != 0)
print(x.collect())
print(z.collect())
```

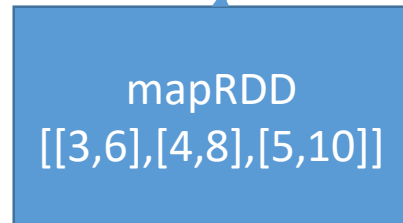


Transformations: **map** and **flatMap**

MAP

Return a new RDD by applying a function to each element of this RDD

```
x = sc.parallelize([3,4,5])
y = x.map(lambda z: [z,z*2])
print(x.collect())
print(y.collect())
```



FLATMAP

Return a new RDD by applying a function to each element of the RDD, and then flattening the results. Also, function in flatMap can return a list of elements (0 or more)

```
x = sc.parallelize([3,4,5])
z = x.flatMap(lambda z: [z,z*2])
print(x.collect())
print(z.collect())
```



Transformations: `reduceByKey` and `groupByKey`

REDUCEBYKEY

Return a new RDD by combining the values with the same key with a given function

`[('a',3), ('b',4), ('c',5), ('a',4), ('b',-6)]`

reduceRDD

`[('a', 7), ('c', 5), ('b', -2)]`

```
x = sc.parallelize([('a',3), ('b',4), ('c',5), ('a',4), ('b',-6)])
y = x.reduceByKey(add)
print(x.collect())
print(y.collect())
```

GROUPBYKEY

Return a new RDD by grouping the values with the same key

sourceRDD

groupRDD

`[('a', [3, 4]), ('c', [5]), ('b', [-6, 4])]`

```
x = sc.parallelize([('a',3), ('b',4), ('c',5), ('a',4), ('b',-6)])
I = x.groupByKey()
y = I.map(lambda z: (z[0], sum(z[1])))
print(x.collect())
print(y.collect())
```



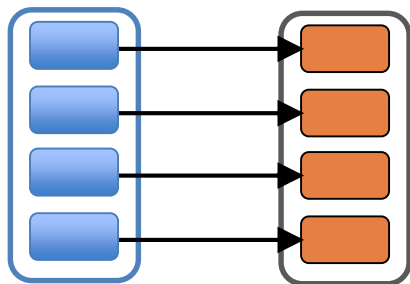

VS



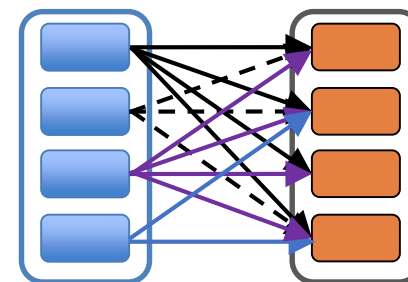
narrow

wide

each partition of the parent RDD is used by at most one partition of the child RDD

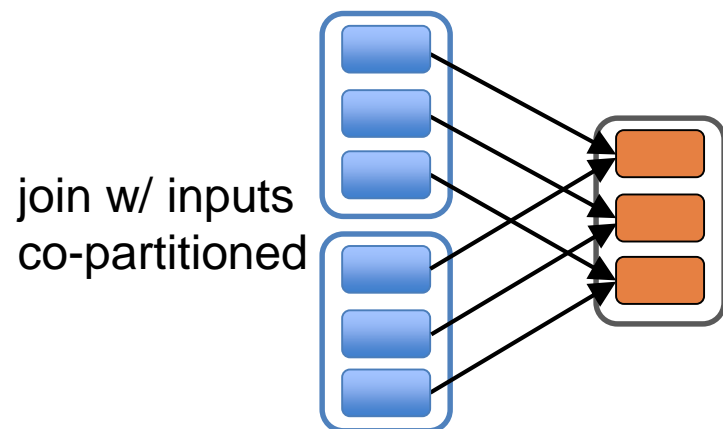
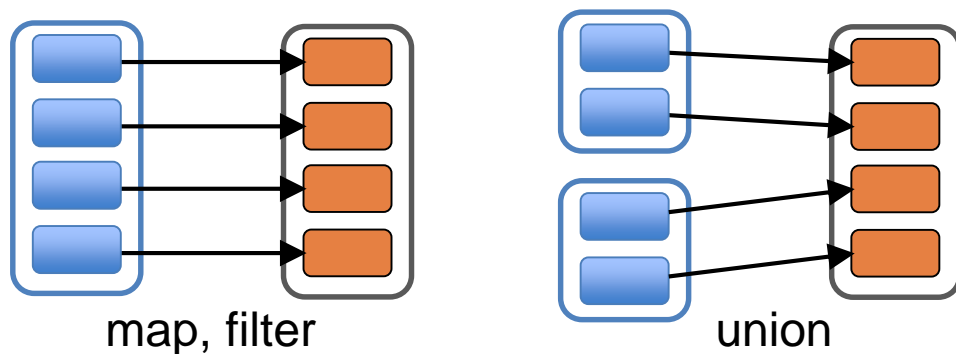


multiple child RDD partitions may depend on a single parent RDD partition



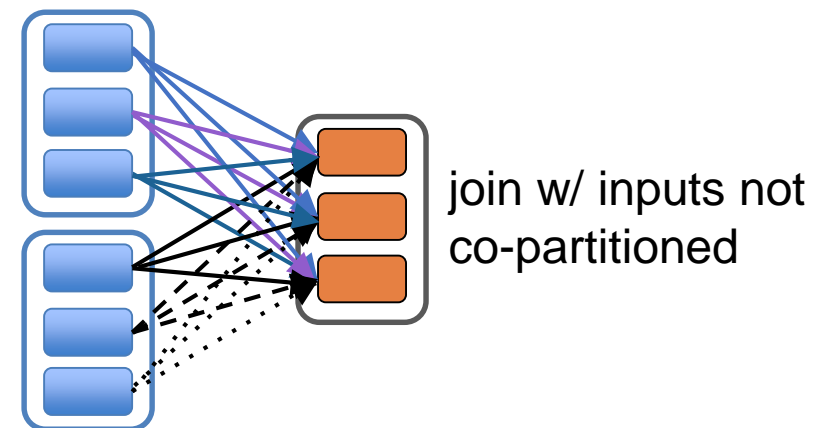
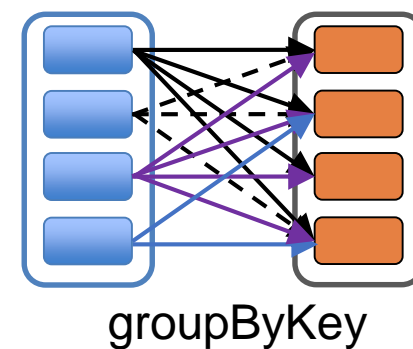
narrow

each partition of the parent RDD is used by at most one partition of the child RDD



wide

multiple child RDD partitions may depend on a single parent RDD partition



Narrow vs wide transformations

Transformations with (usually) Narrow dependencies:

map
mapValues
flatMap
filter
mapPartitions
mapPartitionsWithIndex

Transformations with (usually) Wide dependencies:
(might cause a shuffle)

cogroup
groupWith
join
leftOuterJoin
rightOuterJoin
groupByKey
reduceByKey
combineByKey
distinct
intersection
repartition
coalesce

This list usually holds, but as seen above, in case of join, depending on the use case, the dependency of an operation may be different from the above lists

Actions

Actions trigger the computations and resulting result must fit in the driver JVM

<code>collect():</code>	returns all the elements of the RDD as an array to the driver, should be done after a filter or other operation
<code>count()</code>	returns the number of elements in the RDD
<code>countByKey()</code>	for pair RDDs returns (K, Int) pairs with the count of each key
<code>first()</code>	returns the first element of the RDD
<code>take(n)</code>	returns an array with first n elements
<code>saveAsTextFile()</code>	writes the elements of the RDD as a text file to HDFS or local filesystem
<code>getNumPartitions()</code>	returns the number of partitions of the RDD

RDD: Hands On

Login to lxplus

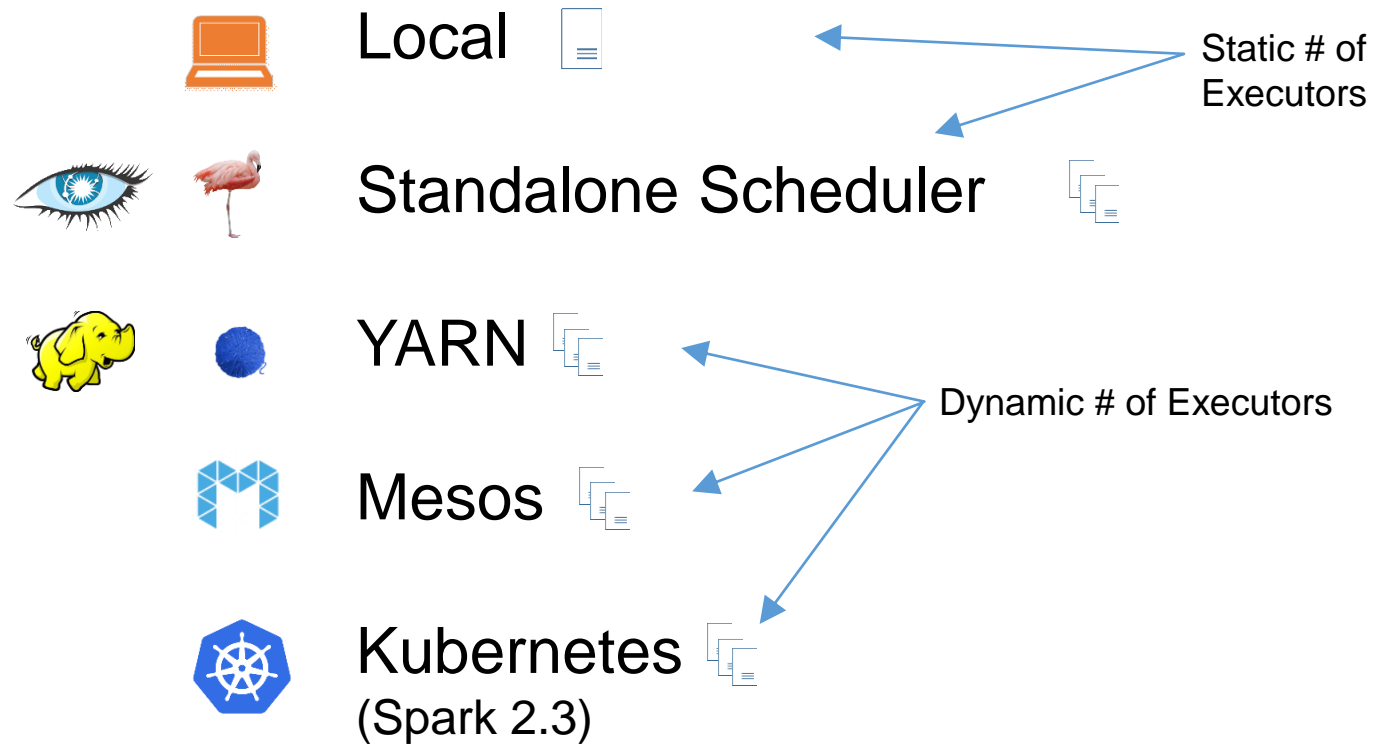
```
cd EOS_HOME (cd /eos/user/p/pkothuri)  
git clone https://github.com/prasanthkothuri/sparkTraining.git
```

Login to SWAN

<https://swan.cern.ch/>

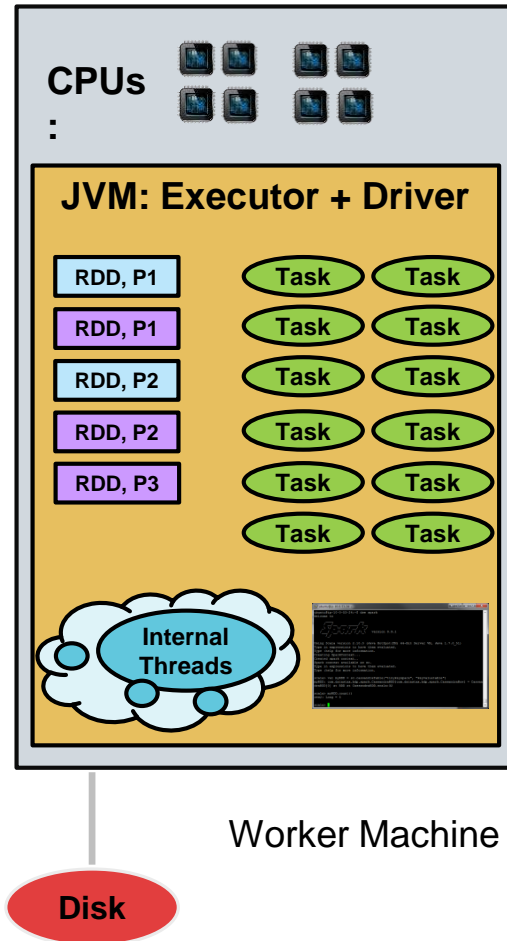
open sparkTraining->notebooks->Tutorial_RDD_Short-Final.ipynb

Resource Managers





Local Mode



Worker Machine



LOCAL MODE

3 options:

- local
- local[N]
- local[*]



```
> ./bin/spark-shell --master local[12]
```

```
> ./bin/spark-submit --name "MyFirstApp"  
--master local[12]  
myApp.jar
```

```
val conf = new SparkConf()  
    .setMaster("local[12]")  
    .setAppName("MyFirstApp")  
    .set("spark.executor.memory", "3g")  
val sc = new SparkContext(conf)
```



YARN Mode

What is YARN?

Resource Management and negotiation so that multiple apps can live and operate together in a Hadoop cluster

Manage CPU and memory and allocate it to different apps

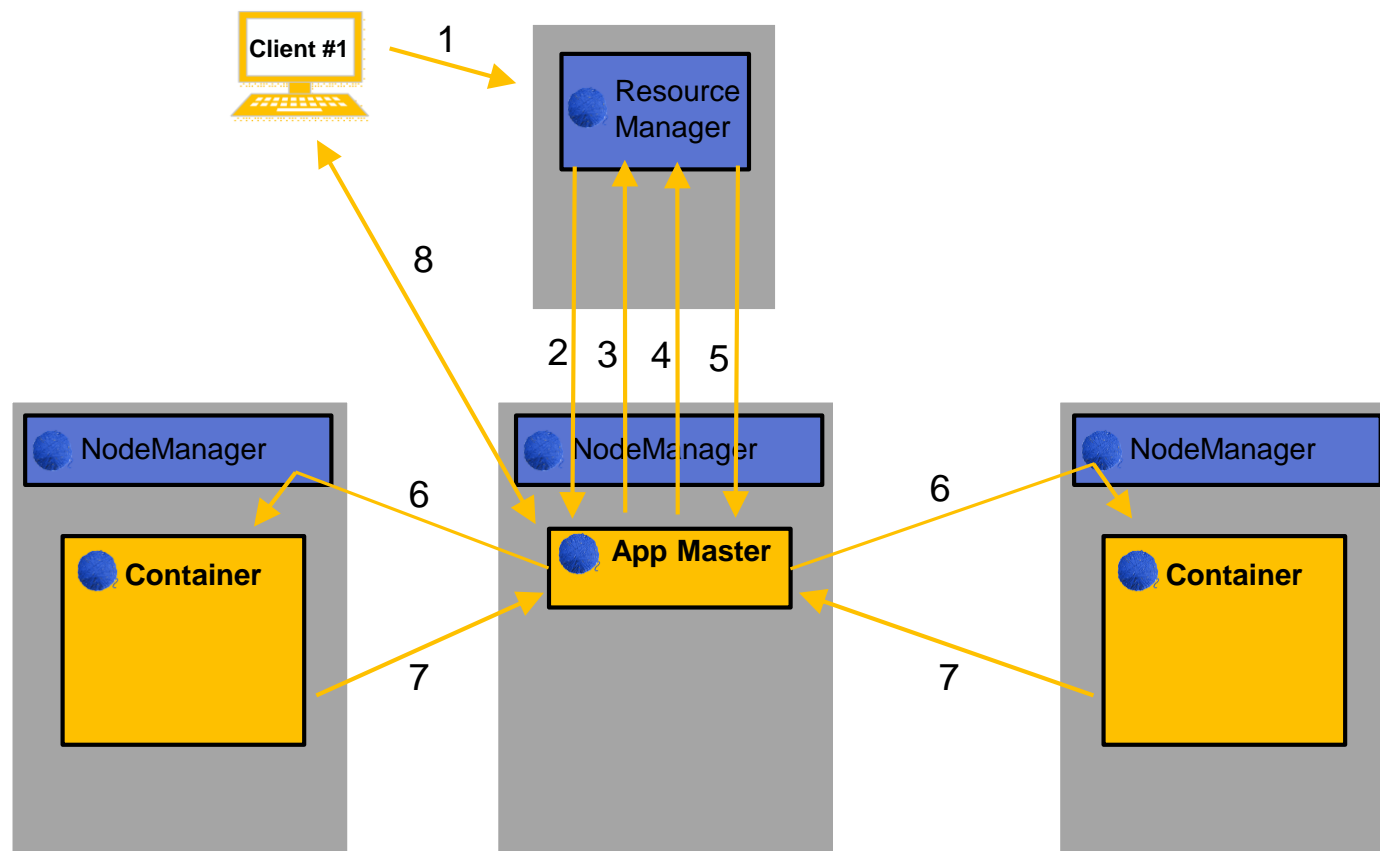
Multiple users

Multiple apps

Lots of advanced features



YARN Mode





YARN Mode - Settings

- `--num-executors`: controls how many executors will be allocated
- `--executor-memory`: RAM for each executor
- `--executor-cores`: CPU cores for each executor

Ways To Run Spark

Interactive Analysis and Data Exploration

spark-shell (spark scala REPL)

pyspark (spark python REPL)

Jupyter Notebooks

Batch applications

spark-submit

Lifecycle of a Spark program

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

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

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

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

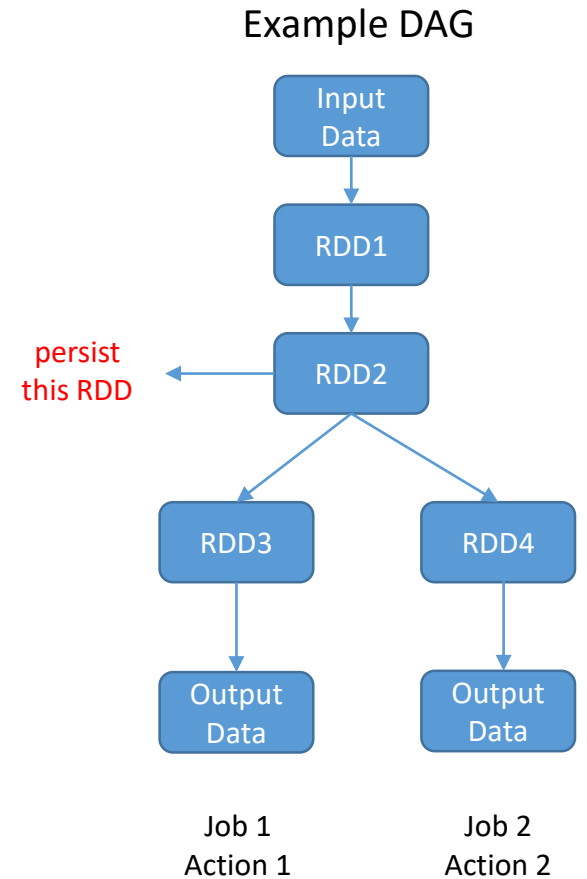
Persistence and Cache

`persist()` is an action which triggers computation and persists a dataset in memory across operations

You can persist with different storage level options (e.g `MEMORY_ONLY` or `MEMORY_AND_DISK`)

`persist()` or `cache()` is key for iterative algorithms

`cache()` is same as `persist()` with `MEMORY_ONLY` storage option



DataFrames

- An extension to the existing RDD API
- DataFrame is an RDD with schema
- DataFrames have numerous optimizations that make them much faster than RDD (predicate pushdown, bloom-filter)
- Write less code – solve problems concisely using dataframe functions
- Inspired by data frames in Python (pandas) and R

High-Level Operations

Solve common problems concisely using DataFrame functions:

- selecting columns and filtering
- joining different data sources
- aggregations (count, sum, average, etc.)
- descriptive statistics
- plotting results (e.g., with Pandas)

Write Less Code: DataFrame vs RDD

Using RDDs

```
w_rdd.map(lambda record: (record.NAMELAST, 1) ) \
    .reduceByKey(add) \
    .map(lambda (x,y):(y,x)) \
    .sortByKey(False) \
    .collect()
```

Using DataFrames

```
w_df.groupBy("NAMELAST") \
    .count() \
    .orderBy("count") \
    .show()
```


Unified Interface for I/O

```
df = spark.read \
    .format("json") \
    .option("samplingRatio", "0.1") \
    .load("/user/pkothuri/sparkTraining/meetup.json")
```

```
df.write \
    .format("parquet") \
    .mode("append") \
    .partitionBy("year") \
    .saveAsTable("my-data")
```

read and **write**
functions create
new builders for
doing I/O

Construct a DataFrame

```
# Create a DataFrame from json file
```

```
df = spark.read.json("/tmp/read_my_shiny.json")
```

```
# Create a DataFrame by loading a parquet file
```

```
df = spark.read.parquet("/tmp/path_to_the_parquet_file")
```

```
df = sqlContext.read.format('jdbc') \  
  .options(driver='oracle.jdbc.driver.OracleDriver',url='jdbc:oracle:thin:username/password@host  
:port:service_name',dbtable='table_name') \  
  .load()
```

```
# Convert a DataFrame
```

```
df = rdd.toDF()
```

Schema Inference

DataFrames have schemas and can infer the schema from the type of the data being read

A Parquet file has a schema (column names and types) that DataFrames can use.

```
>>> df.printSchema()
root
|-- LocID: string (nullable = true)
|-- Location: string (nullable = true)
|-- VarID: string (nullable = true)
|-- Variant: string (nullable = true)
|-- Time: string (nullable = true)
|-- MidPeriod: string (nullable = true)
|-- SexID: string (nullable = true)
|-- Sex: string (nullable = true)
|-- AgeGrp: string (nullable = true)
|-- AgeGrpStart: string (nullable = true)
|-- AgeGrpSpan: string (nullable = true)
|-- Value: string (nullable = true)
```

What if the data doesn't have schema (e.g csv)

create an RDD of particular type using python namedtuple, dict and
convert RDD to DataFrame

You can also specify the column names in .toDF() function

DataFrame: Transformations and Actions

DataFrames are *lazy*. *Transformations* contribute to the query plan, but they don't execute anything.

Actions cause the execution of the query.

Execution of the query means:

- Spark initiates a distributed read of the data source
- The data flows through the transformations
- The result of the action is pulled back into the driver JVM.

Transformation examples

- filter
- select
- drop
- intersect
- join

Action examples

- count
- collect
- show
- head
- take

Full list - <https://spark.apache.org/docs/2.2.0/api/python/pyspark.sql.html>

Transformations: select() and selectExpr()

select()

similar to SQL SELECT, allows you to limit the results to specific columns

```
// single column
df.select("NAMELAST").show()

// multiple columns
df.select("NAMELAST", "NAMEFIRST").show()

// column alias with expr
df.select(expr("NAMELAST as lastname")).show()

// further manipulation with expr
df.select(expr("NAMELAST as lastname").alias("NAMELAST")).show()
```

selectExpr()

simple way to build complex expressions, can add any non-aggregating functions

```
// column alias
df.selectExpr("NAMELAST as lastname", "NAMELAST").show()

// add new column
df.selectExpr(
    "*", # all original columns
    "(Total_People < 50) as Total_less_than_50")\
    .show()

// aggregations over entire dataframe
df.selectExpr("avg(count)",
    "count(distinct(DEST_COUNTRY_NAME))").show()
```

Transformations: lit(), withColumn(), withColumnRenamed() and drop()

lit()

allows you to add constant value column

withColumnRenamed()

allows you rename column in the dataframe

```
// add constant value column
df.select(expr("*"), lit(1).alias("One")).show()

// rename a column
df.withColumnRenamed("OLD_COL_NAME", "newName")
```

withColumn()

allows you to add a column to the dataframe

drop()

remove column(s) from the dataframe

```
// add a column to the dataframe
df.withColumn("withinCountry",
  expr("ORIGIN_COUNTRY_NAME == DEST_COUNTRY_NAME"))\
  .show(2)

// drop a column
df.drop("COL1", "COL2")
```

Transformations: filter, sort and cast

filter(), where()

we create an expression that evaluates to true or false and filter out the rows with an expression that is equal to false

```
// with filter()

df.filter(col("age") < 18).show()

// with where()

df.where("age < 18")\
  .where("lastName = 'Steve'")\
  .show()
```

sort(), orderBy()

sort the values in a dataframe, default ascending order

show()

displays the first n elements in the DataFrame (n defaults to 20)

```
// with filter()
df.sort("age").show()

// with orderBy()
df.orderBy("age", "lastName").show(5)

// show(), limit()
df.show(10, False)
df.limit(5).show()
```

Transformations: aggregations

Following are the most common aggregation functions

count()

countDistinct()

approx_count_distinct()

first()

last()

min()

max()

sum()

sumDistinct()

avg()

variance(), var_pop()

stddev(), stddev_pop()

skewness()

kurtosis()

Transformations: grouping data and window functions

groupBy()

perform calculations (e.g
aggregations) on groups of data

window functions

Compute aggregations on a specific
“window” of data

```
// groupby lastName

>>>w_df.groupBy(w_df.NAMELAST) \
      .count() \
      .orderBy("count",ascending=False) \
      .show(10)

// groupby with maps

df.groupBy("InvoiceNo").agg(expr("avg(Quantity)"),expr("stddev_pop(Quantity)"))\
  .show()
```

Transformations: Joins

r_DF – is a DataFrame holding movie ratings [userId,movieId,rating,timestamp]

m_DF – is a DataFrame holding movie information [movieId,title,genres]

These DataFrames can be joined as below to obtain number of reviews per genre

```
>>>r_DF.join(m_DF,r_DF.movieId == m_DF.movieId) \  
      .groupBy(m_DF.genres) \  
      .count() \  
      .orderBy("count",ascending =False) \  
      .show(5)
```

```
+-----+-----+  
|          genres|count|  
+-----+-----+  
|          Drama| 5832| | |
|          Comedy| 5648|  
|    Comedy|Romance| 3194|  
|    Drama|Romance| 2649|  
|Comedy|Drama|Romance| 2486|  
+-----+-----+  
only showing top 5 rows
```

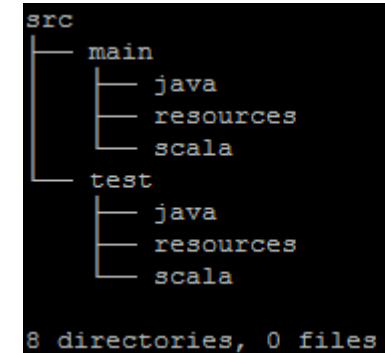
Demo: Simple Application in Spark and Scala

1) Install Spark and sbt (not required in our case as we ssh to a node where these are available)

```
ssh hadalytic
```

2) Create the following directory structure

```
mkdir ~/sparkTest
cd ~/sparkTest
mkdir -p src/main/scala
mkdir -p src/main/java
mkdir -p src/main/resources
mkdir -p src/test/scala
mkdir -p src/test/java
mkdir -p src/test/resources
```

A terminal window with a black background and white text showing the creation of a directory structure. The root directory is 'src', which contains two subdirectories: 'main' and 'test'. Each of these subdirectories contains three subdirectories: 'java', 'resources', and 'scala'. At the bottom of the terminal, it says '8 directories, 0 files'.

```
src
├── main
│   ├── java
│   ├── resources
│   └── scala
└── test
    ├── java
    ├── resources
    └── scala

8 directories, 0 files
```

3) Create .sbt file

```
vi build.sbt

name := "dfAgg"
version := "1.0"
scalaVersion := "2.11.8"
val sparkVersion = "2.3.0"

libraryDependencies ++= Seq(
  "org.apache.spark" %% "spark-core" % sparkVersion,
  "org.apache.spark" %% "spark-sql" % sparkVersion,
  "org.apache.hadoop" % "hadoop-client" % "2.7.0"
)
```

Demo: Simple Application in Spark and Scala

4) Create sample spark scala application as below

```
mkdir -p src/main/scala/ch/cern/sparkTraining/examples/  
vi src/main/scala/ch/cern/sparkTraining/examples/dfExample.scala  
  
/* dataframe aggregation Example */  
  
package ch.cern.sparkTraining.examples  
  
import org.apache.spark.sql.SparkSession  
  
object dfAggExample {  
  def main(args: Array[String]) {  
  
    if (args.length < 1) {  
      System.err.println("Usage: dfAggExample <file>")  
      System.exit(1)  
    }  
  
    val spark = SparkSession  
      .builder  
      .appName("dfAgg")  
      .getOrCreate()  
  
    val df = spark.read.parquet(args(0))  
    val cnt = df.count()  
    println("Number of records in the file is: %s".format(cnt))  
  }  
}
```

Demo: Simple Application in Spark and Scala

5) Build the application / package

```
sbt package
```

6) Run the application

```
export SPARK_HOME=/usr/hdp/spark-2.3.0
export PATH=$SPARK_HOME/bin:$PATH

spark-submit \
  --class "ch.cern.sparkTraining.examples.dfAggExample" \
  --master yarn \
  --deploy-mode client \
  target/scala-2.11/dfagg_2.11-1.0.jar /Training/Spark/WH_VR.parquet
```

Dataframes: Hands On

Login to Ixplus

```
cd EOS_HOME (cd /eos/user/p/pkothuri)  
git clone https://github.com/prasanthkothuri/sparkTraining.git
```

Login to SWAN

<https://swan.cern.ch/>

open sparkTraining->notebooks->Tutorial_DataFrame1_Final.ipynb

open sparkTraining->notebooks->Tutorial_DataFrame2_Final.ipynb

Spark SQL and DataFrames

DataFrames and Spark SQL are essentially tied to each other

- The DataFrames API provides a programmatic interface – really, a domain-specific language (DSL) – for interacting with data
- Spark SQL provides a SQL-like interface
- Whatever you can do in DataFrames, you can do in Spark SQL and vice versa

Spark SQL contd.

Spark SQL allows you to manipulate distributed data with SQL queries. Currently, two SQL dialects are supported.

- If you're using a Spark SQLContext, the only supported dialect is "sql", a rich subset of SQL 92.
- If you're using a HiveContext, the default dialect is "hiveql", corresponding to Hive's SQL dialect. "sql" is also available, but "hiveql" is a richer dialect.

Spark SQL contd.

- You issue SQL queries through a `SparkSession` or `SQLContext` or `HiveContext`, using the `sql()` method.
- The `sql()` method returns a `DataFrame`.
- You can mix `DataFrame` methods and SQL queries in the same code.
- To use SQL, you must either:
 - query a persisted Impala or Hive table, or
 - make a table alias for a `DataFrame`, using `registerTempTable()` or `createOrReplaceTempView()`

Spark SQL - Example

To issue SQL against an existing DataFrame, create a temporary table, which essentially gives the DataFrame a *name* that's usable within a query.

```
>>> df = spark.read.parquet("/tmp/WH_VRecord/part-r-00000-5396c70a-ff5b-4dda-9306-3a5e8bd9167a.gz.parquet")

>>> df.registerTempTable("VistorRecords")

>>> sql_df = spark.sql("SELECT NAMELAST,NAMEFIRST,APPT_START_DATE,APPT_END_DATE FROM VistorRecords")

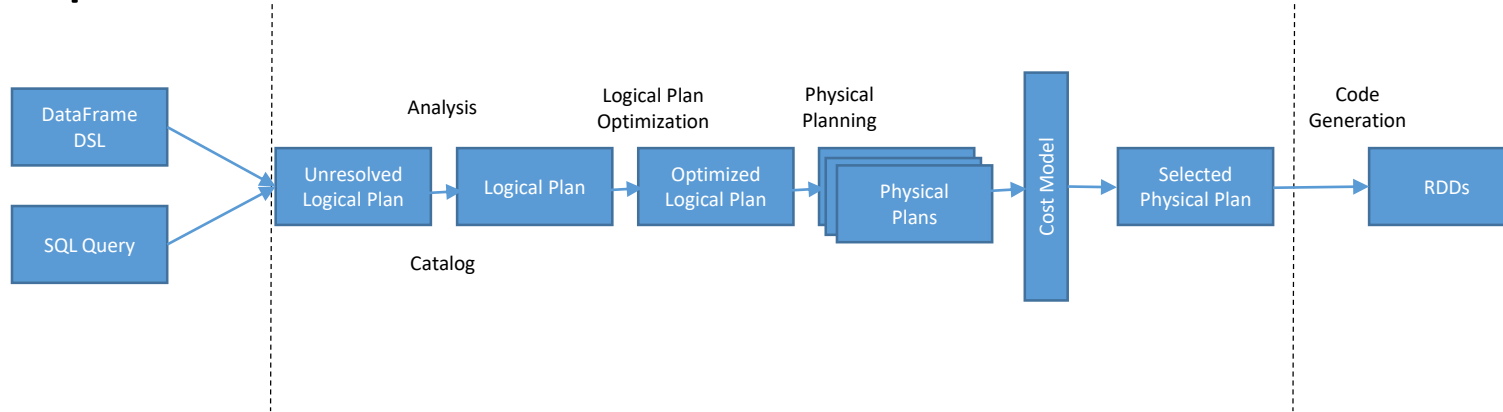
>>> sql_df.show(5)
```

NAMELAST	NAMEFIRST	APPT_START_DATE	APPT_END_DATE
Adamopoulos	Stella	5/1/15	5/1/15
Brosman	Muriel	5/1/15	5/1/15
Brumfield	Avery	5/1/15	5/1/15
Chipman	Catherine	5/1/15	5/1/15
Chubb	Steven	5/1/15	5/1/15

```
only showing top 5 rows
```

Catalyst : Spark's Optimizer

- Spark SQL uses catalyst to optimize all the queries written both in spark sql and dataframe dsl



Analysis

- phase where attribute references or relations are resolved
- e.g: column validity, column type
- catalog object tracks the tables in all data sources

Catalyst : Spark's Optimizer

Logical Optimizations

- Standard rule-based optimizations
- e.g: predicate pushdown, project pruning, null propagation etc

Physical Planning

- generated one or more physical plans
- cost model is used to select a plan

Code Generation

- generate java bytecode to speed up execution

Further Reading

- <https://databricks.com/blog/2015/04/13/deep-dive-into-spark-sqls-catalyst-optimizer.html>

DataFrames and RDDs

- DataFrames are built on top of the Spark RDD API.
 - This means you can use normal RDD operations on DataFrames.
- However, stick with the DataFrame API, wherever possible.
 - Using RDD operations will often give you back an RDD, not a DataFrame.
 - The DataFrame API is more efficient and underlying operations are optimized with Catalyst.

SPARK SQL: Hands On

Login to lxplus

```
cd EOS_HOME (cd /eos/user/p/pkothuri)
```

```
git clone https://github.com/prasanthkothuri/sparkTraining.git
```

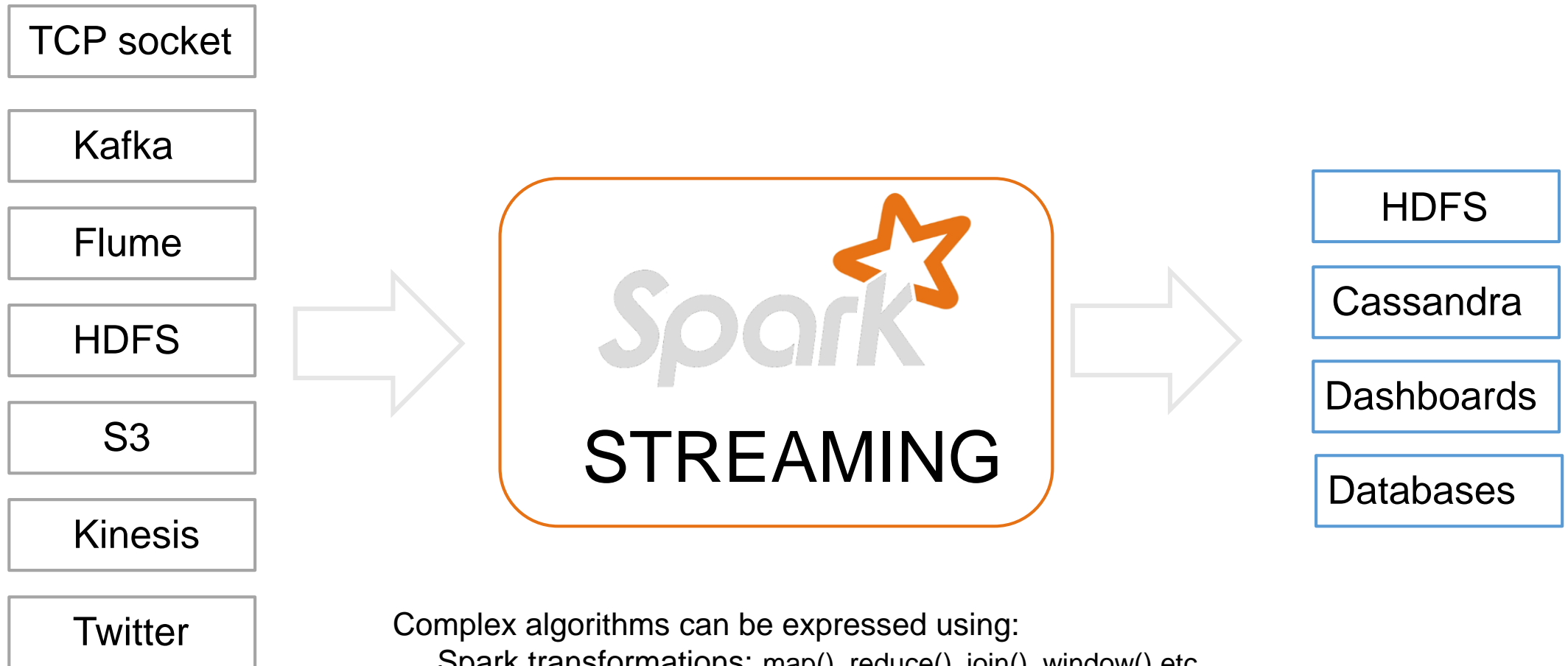
Login to SWAN

<https://swan.cern.ch/>

open sparkTraining->notebooks->Tutorial_SparkSQL_Final.ipynb

Spark Streaming – What is it?

Extension of the core Spark API that enables scalable, high-throughput and fault-tolerant processing of live data streams (unbounded data)



Complex algorithms can be expressed using:

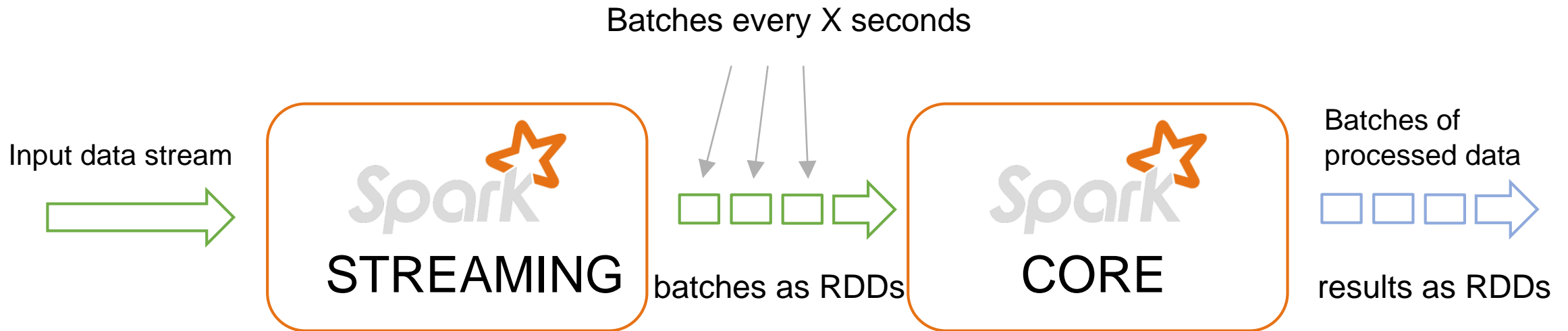
Spark transformations: `map()`, `reduce()`, `join()`, `window()` etc

MLlib + GraphX

SQL

Spark Streaming – How does it work?

Spark Streaming receives live input data streams and divides the data into batches, which are then processed by the Spark engine to generate the final stream of results in batches.



- breakup up data streams into batches of few secs
- each batch of data is treated as RDD and process them using RDD operations
- Processed results are pushed out in batches

Spark Streaming – Programming model

Discretized stream (DStream)

- Primary abstraction in spark streaming

- Continuous stream of RDD

- Every RDD contains data from a certain interval

DStream Operations

Transformations

- Stateless Transformations

- Stateful Transformations

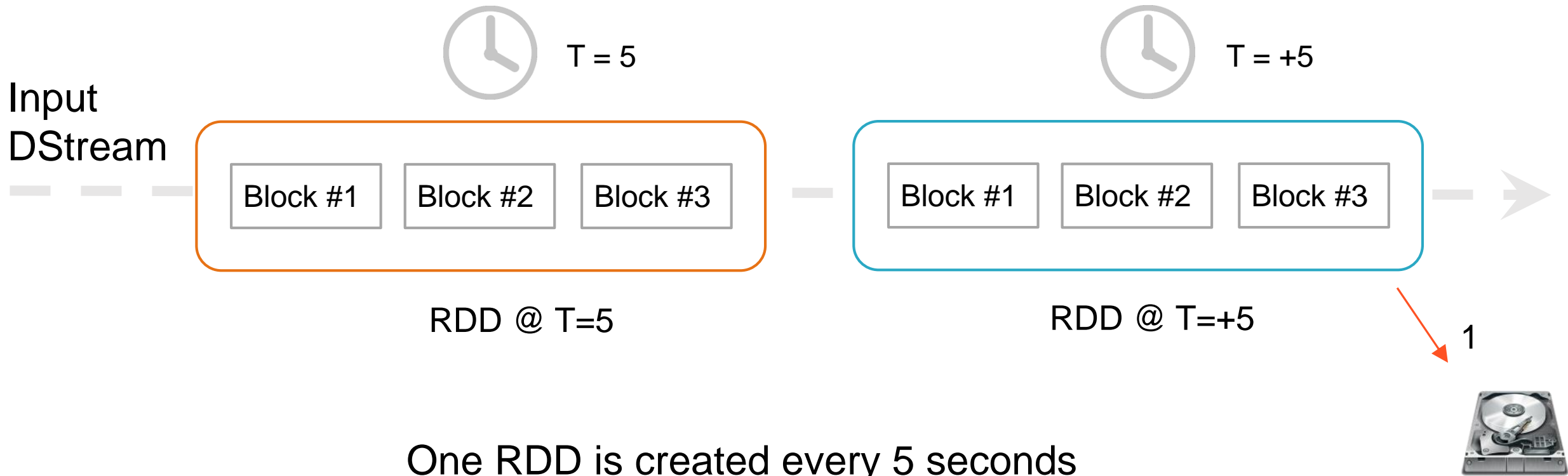
Output Operations

- Print to console

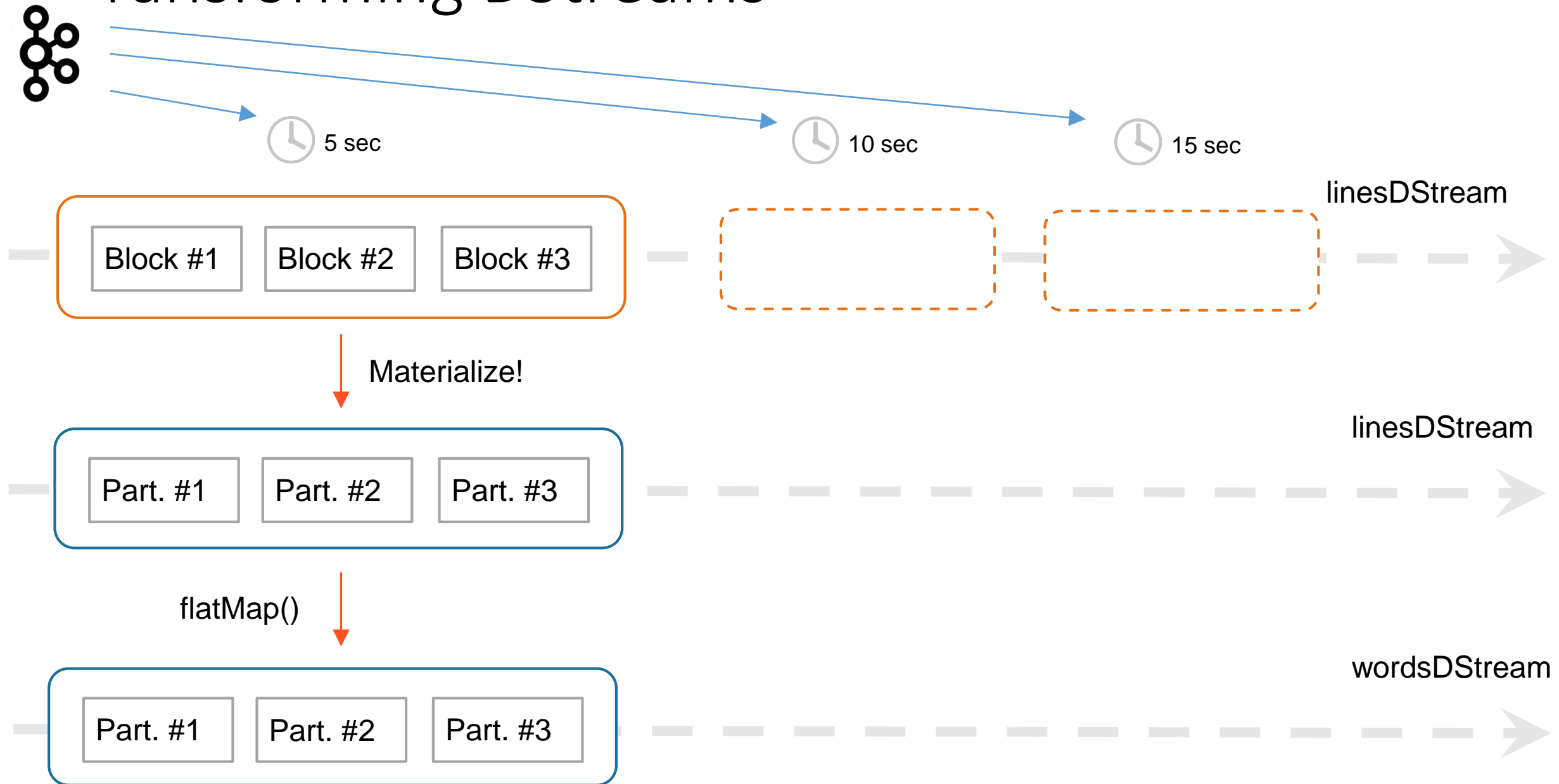
- Save to persistence storage

DSTREAM – Discretized Stream

Batch interval = 5 seconds



Transforming DStreams



Transformations on DStreams

Stateless Transformations

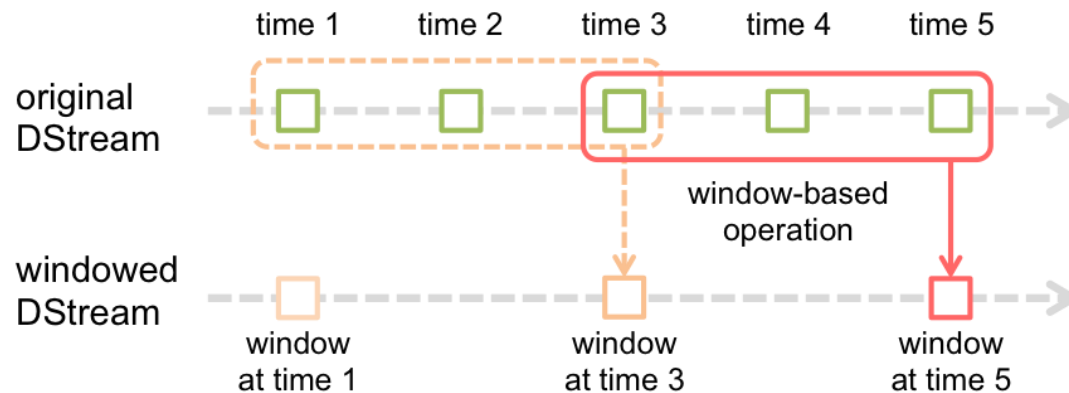
map($f(x)$)
flatMap($f(x)$)
filter($f(x)$)
count()
repartition(n)
reduce($f(x)$)
countByValue($f(x)$)
reduceByKey($f(x)$)
transform($f(x)$)

Stateful Transformations

updateStateByKey($f(x)$)
countByWindow(wL,sl)
reduceByWindow($f(x)$)
reduceByKeyAndWindow(wL,sl)
countByValueAndKey(wL,sl)

Windowed Operations

windowed operations allow you to apply transformations over a sliding window of data



The RDDs that fall within the window are combined and operated upon to produce RDDs for the window stream

window operation needs 2 additional parameters

window length – the duration of the window

sliding interval – the interval at which the window operation is performed

* both these parameters must be multiple of batch interval

Output Operations on DStreams

Output operations specify what needs to be done with the final transformed data in a stream

if no output operation is applied on a DStream and any of its descendants, then those DStreams will not be evaluated (lazy evaluation)

Output operations

- `print()`

- `saveAsTextFiles()`

- `saveAsHadoopFiles()`

- `foreachRDD ($f(x)$)`

Structured Streaming (since 2.2)

Stream processing engine built on top of Spark SQL

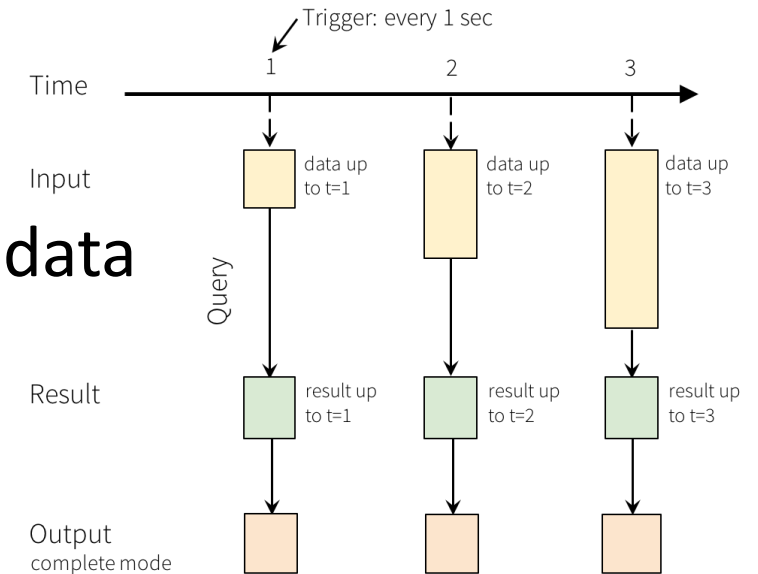
Input: data from source as append only table

Trigger: how frequently to check input for new data

Query: operations on input; usual map/filter/reduce and window operations

Result: final table updated every trigger interval

Output: what part of result to write to data sink after every trigger



Programming Model for Structured Streaming

SPARK Streaming: Hands On

Login to lxplus

```
cd EOS_HOME (cd /eos/user/p/pkothuri)  
git clone https://github.com/prasanthkothuri/sparkTraining.git
```

Login to SWAN

<https://swan.cern.ch/>

open sparkTraining->notebooks->Tutorial-SparkStreaming-Final.ipynb

Spark MLlib

Why Apache Spark for Machine Learning?

- Bigger than memory datasets
 - Able to train a model on large scale dataset
- General Purpose
 - Apart from the libraries for commonly used algorithms, libs for advanced data preparation, feature engineering etc
- Compatibility
 - Support for multiple languages and integrate well with python libs like pandas, scikit-learn etc

Dataframe and RDD based API

DataFrame-based API (spark.ml)

- Easier to construct a machine learning pipeline
- Flexible and versatile API compared to spark.mllib
- Will reach feature parity with spark.mllib in the next releases

RDD-based (spark.mllib)

- Original Machine Learning API
- No new features, only bugfixes
- Will be removed in Spark 3.0

Spark MLlib – Main Concepts

DataFrame

- Same as the Dataframe from Spark SQL / Dataframe API
- Used to hold ML dataset

Transformer

- Transforms one DataFrame into another DataFrame
- Examples
 - Feature transformer appends new column (features) to DataFrame
 - Learning model transforms a DataFrame with features into a DataFrame with predictions

Estimator

- Abstracts the learning algorithm; accepts a DataFrame and produces a Model

Pipeline

- A sequence of transformers and estimators to process and learn from data

Parameter

- Specifying parameters for transformers and estimators

Spark MLlib utilities and algorithms

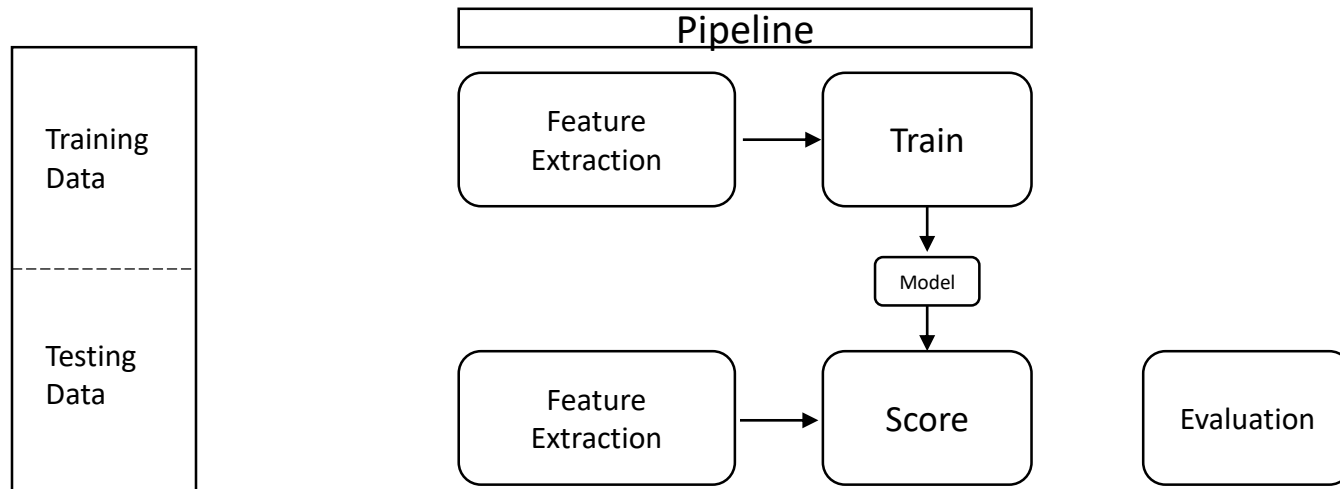
Distributed pre-processing workflow utilities

- Feature Engineering
 - Extraction:- Extracting features from raw data
 - Word2vec
 - countvectorizer
 - Transformations:- modifying, converting or scaling features
 - Tokenizer
 - StringIndexer
 - Standardization
 - Normalization
 - VectorAssembler
 - Selectors:- selecting a subset from a larger set of features

High performant ML algorithms (classification, regression, clustering etc)

- Full list - <http://spark.apache.org/mllib/>

Model Lifecycle



SPARK MLlib: Hands On

Login to lxplus

```
cd /eos/user/p/pkothuri
```

```
git clone https://github.com/prasanthkothuri/sparkTraining.git
```

Login to SWAN

<https://swan.cern.ch/>

open sparkTraining->notebooks->Tutorial-ML-Final.ipynb

Key Learnings

- We have covered the spark concepts; abstraction and architecture
- Introduced to Spark data APIs – RDD, DataFrame and Spark SQL
- Demonstration of using Spark data APIs for exploratory data analysis and data analytics
- Introduced to Spark Mllib for scalable machine learning
- Understood how spark can aid in distributed computing of VERY large datasets
- Several ways to interact with spark – spark-shell, pyspark, spark-submit and jupyter notebooks

Further study

- Slides and Notebooks for Hands On
 - <https://github.com/prasanthkothuri/sparkTraining>
- IT department provides Hadoop, Spark and Kafka services
 - <https://cern.service-now.com/service-portal/function.do?name=Hadoop-Components>
 - <https://cern.service-now.com/service-portal/function.do?name=Kafka>
- Access to Spark service @CERN
 - Build hadoop client machine
 - Ixplus – <https://cern.service-now.com/service-portal/article.do?n=KB0004426>
 - Jupyter notebooks (pySpark) – <https://swan.cern.ch>

Further study

- Coursera course – Big Data Analysis with Scala and Spark
 - <https://www.coursera.org/learn/scala-spark-big-data>
- Hadoop: The Definitive Guide, 4th Edition; By: Tom White
- Expert Hadoop Administration; By Sam R. Alapati
- Kafka: The Definitive Guide; By Todd Palino, Neha Narkhede
- Safari books online accessible for free from CERN IP
 - <http://proquest.tech.safaribooksonline.de/>