

# Developer Training for Spark and Hadoop

申建忠

Frank\_shen@uuu.com.tw

# Hadoop介紹

- 創作者 : Doug Cutting
- 開發語言 : Java
- 靈感啟發 : 2004 Google所發表的2篇論文
  - MapReduce -> MapReduce
  - GFS              -> HDFS
- 2006 Apache Top Project
  - MapReduce : Max utilization of CPU
  - HDFS : Max utilization of Disk

# Core Hadoop(MR v1)



Map Reduce  
(mapred API/Framework/Resource Management)

HDFS  
(Hadoop Distributed FileSystem)

# Core Hadoop(MR v2)



Map Reduce  
(mapreduce API/Framework)

YARN-Yet Another Resource Negotiator  
(YARN API/Resource Management)

HDFS

# HDFS+YARN+SPARK

SPARK



YARN-Yet Another Resource Negotiator  
(YARN API/Resource Management)

HDFS

# HDFS+SPARK

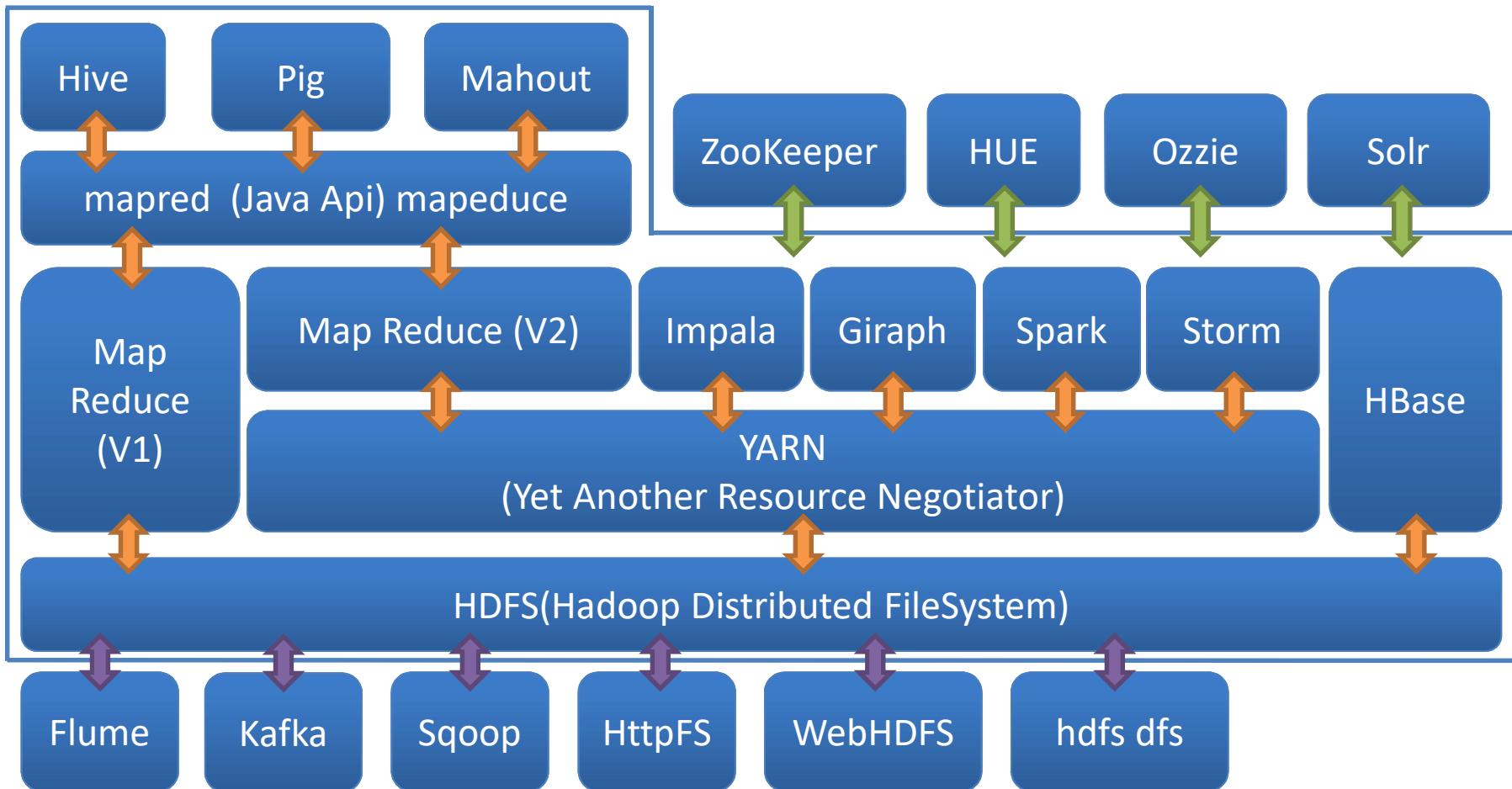
SPARK

SPARK Standalone/[YARN](#)/Apache Mesos

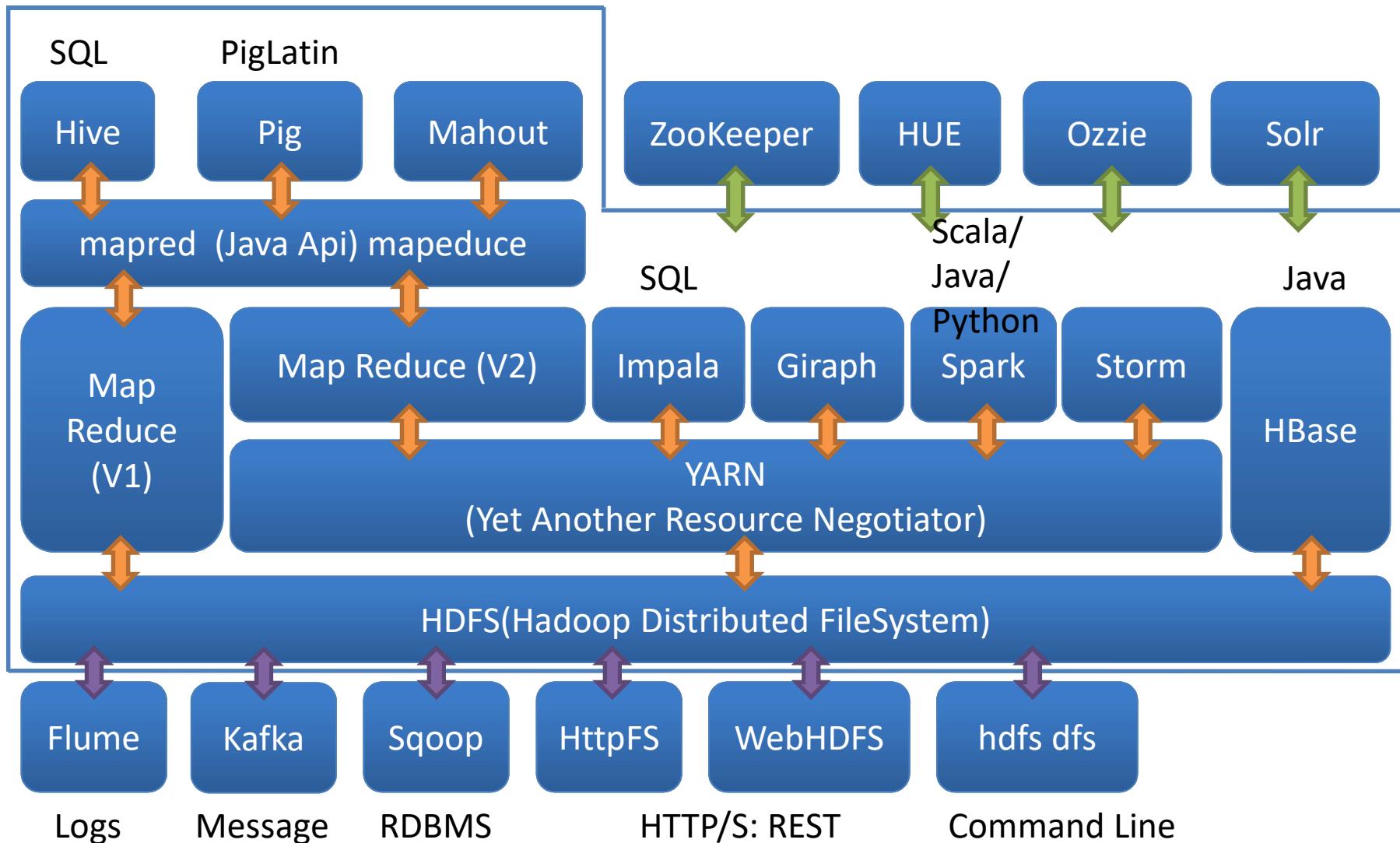
HDFS



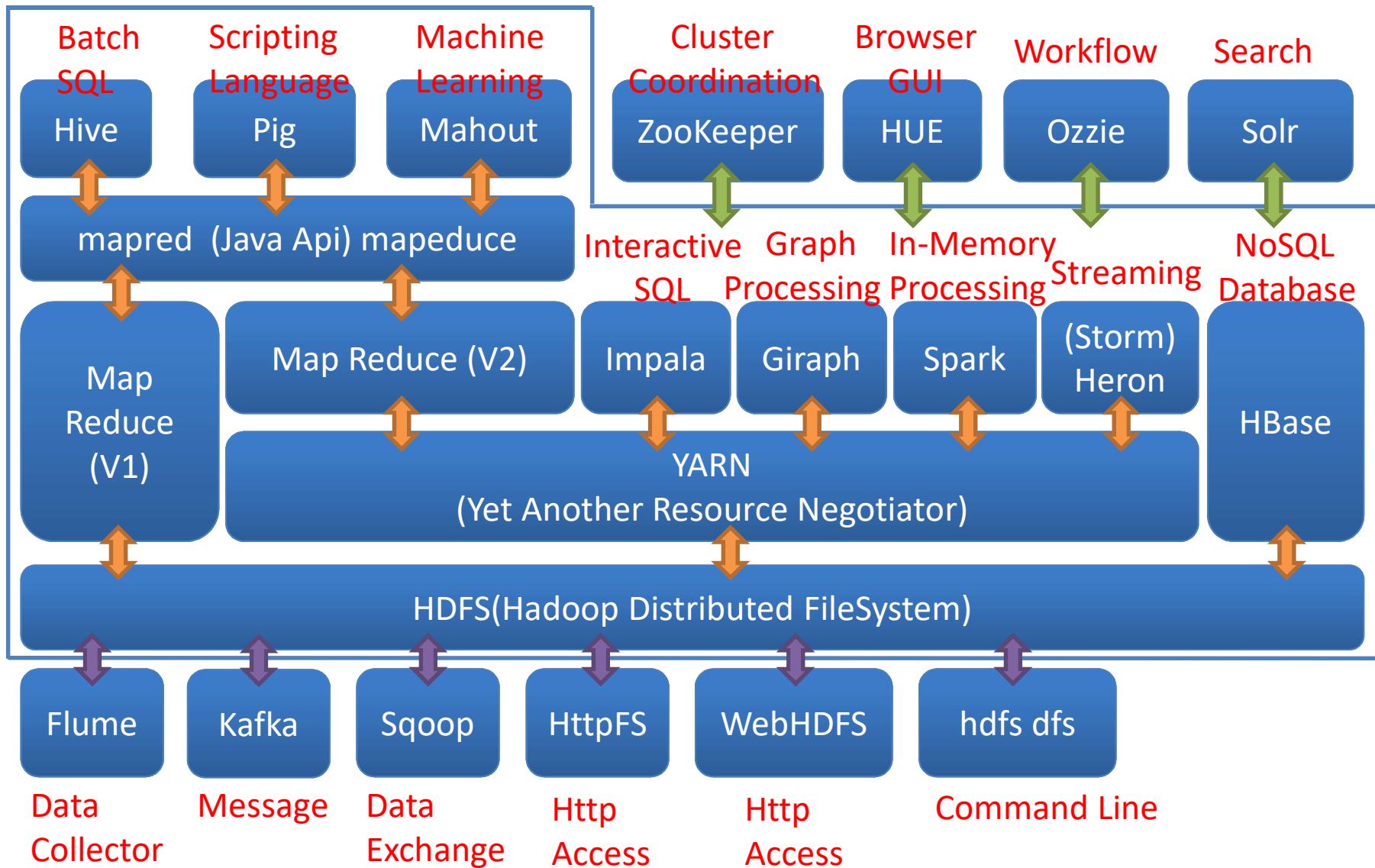
# Hadoop Ecosystem



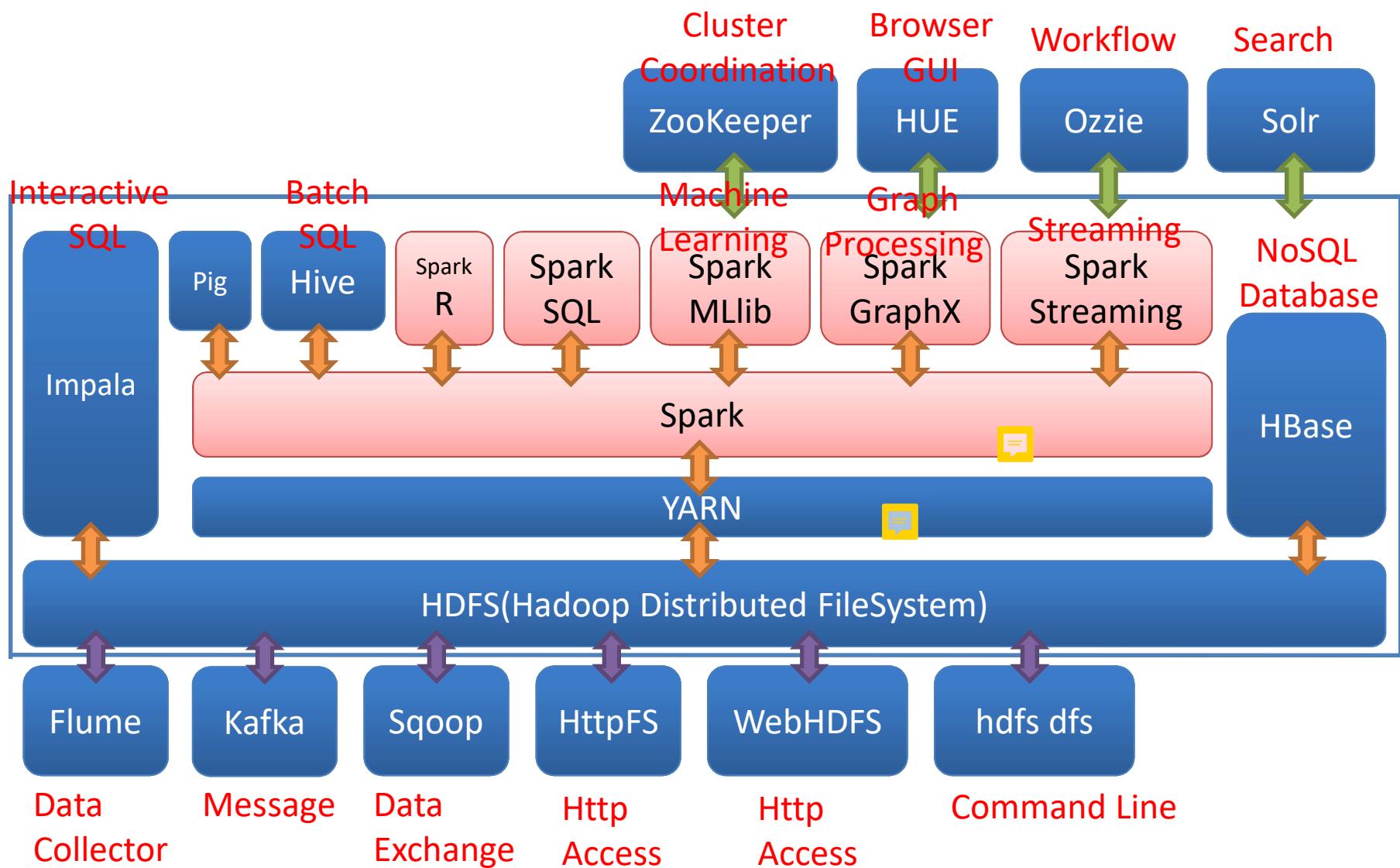
# Hadoop Ecosystem

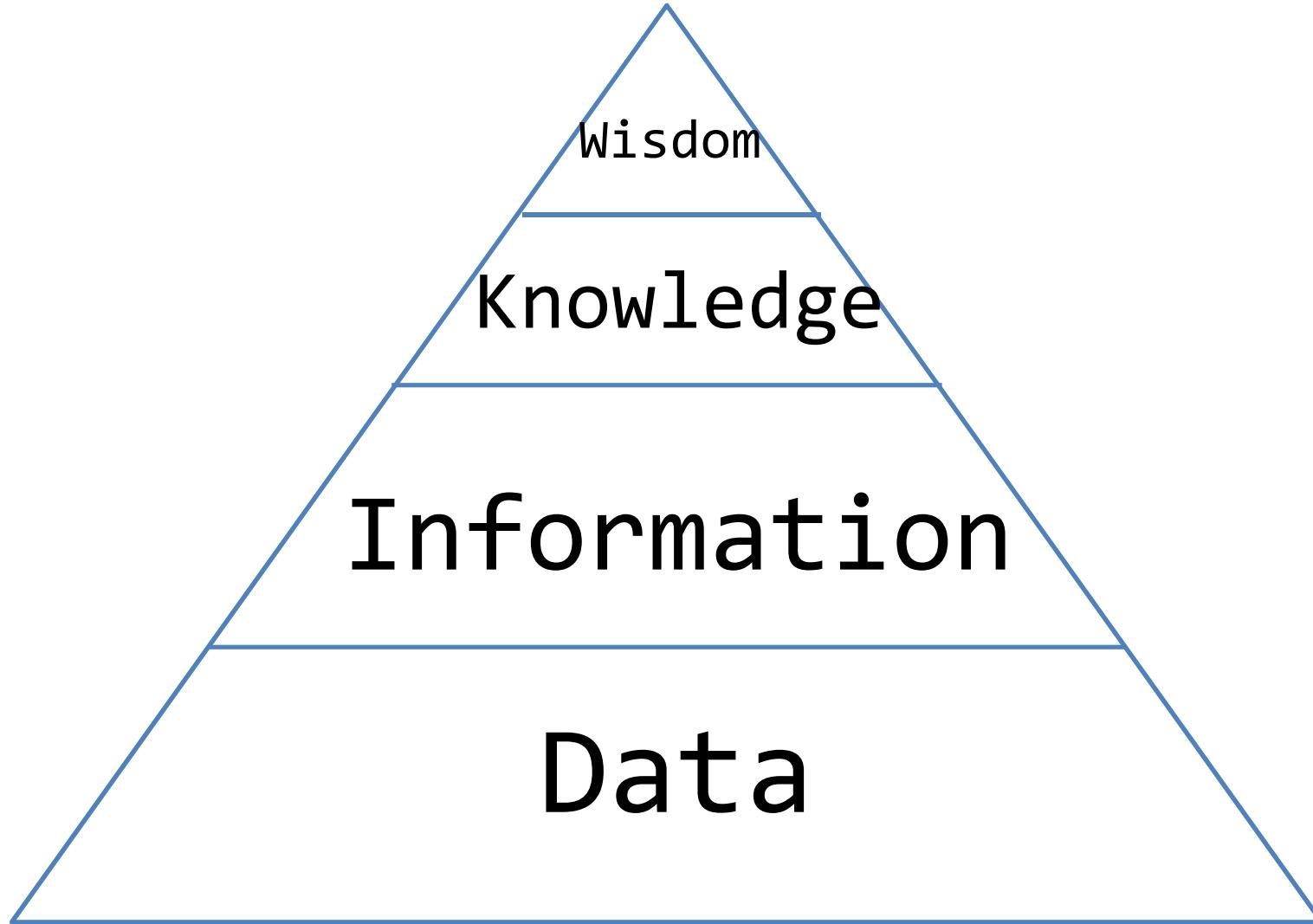


# Hadoop Ecosystem



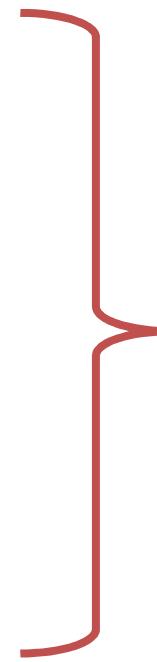
# Hadoop Ecosystem+SPARK





# Big Data

- Volume
  - Data amount
- Variety
  - Data type
- Velocity
  - Data processing Speed
  - Generated



Value



# How big is BIG

- 大數據（或Megadata），或稱巨量資料、海量資料、大資料，指的是所涉及的資料量規模巨大到無法透過人工，在合理時間內達到擷取、管理、處理、並整理成為人類所能解讀的資訊。
- 在總資料量相同的情況下，與個別分析獨立的小型資料集（data set）相比，將各個小型資料集合併後進行分析可得出許多額外的資訊和資料關聯性，可用來察覺商業趨勢、判定研究品質、避免疾病擴散、打擊犯罪或測定即時交通路況等；這樣的用途正是大型資料集盛行的原因。

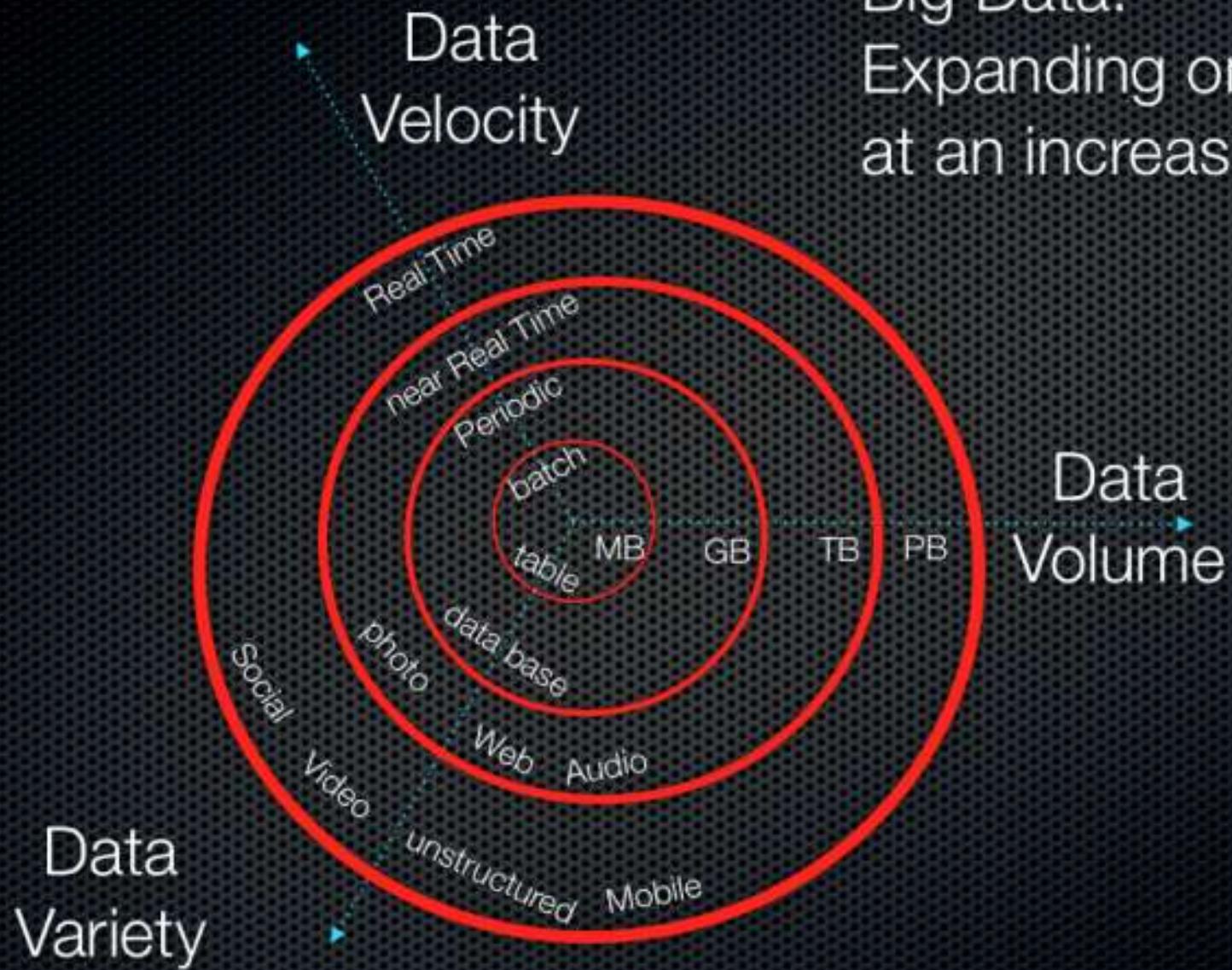
# How many is Variety

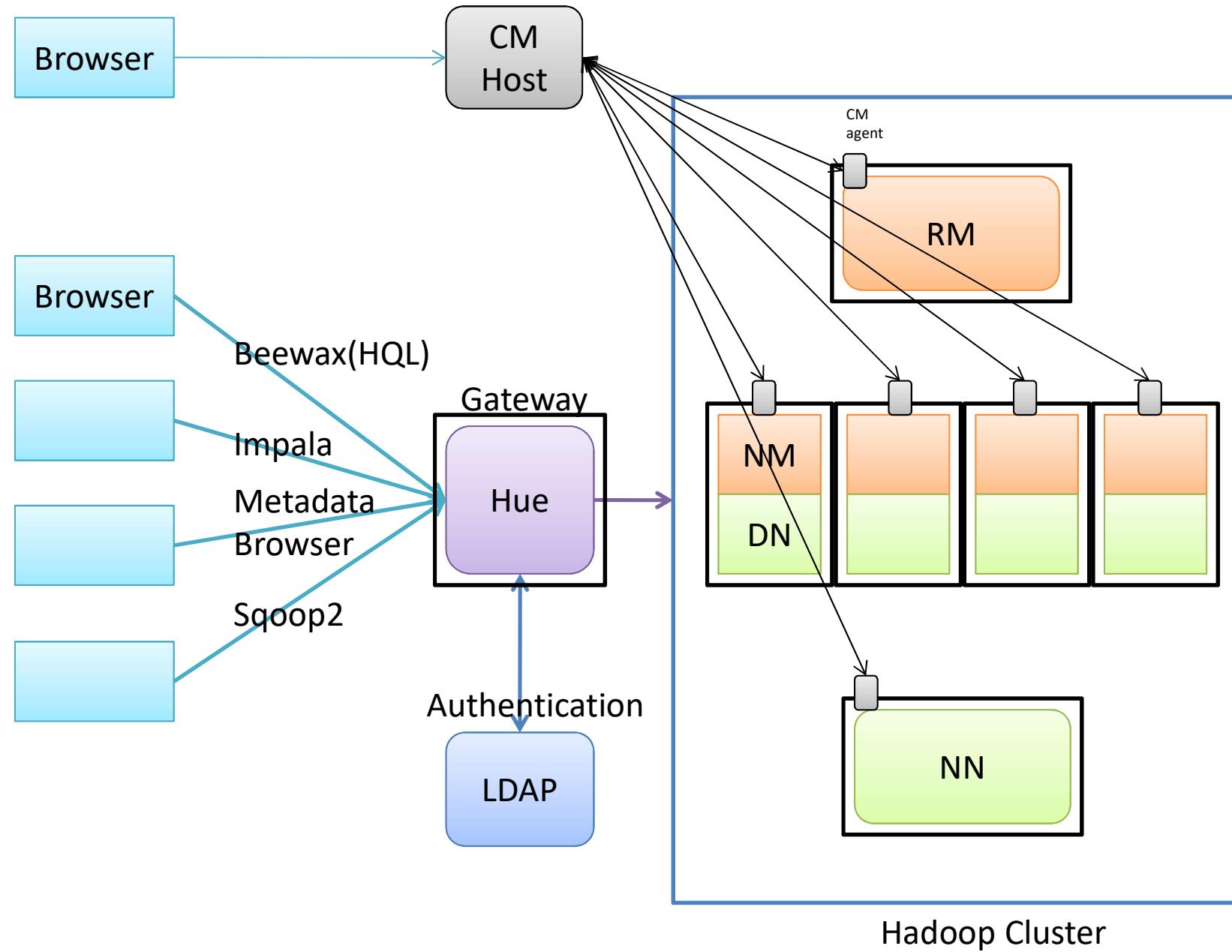
- 大數據的來源種類包羅萬象，十分多樣化，如果一定要把資料分類的話，最簡單的方法是分兩類，結構化與非結構化。早期的非結構化資料主要是文字，隨著網路的發展，又擴展到電子郵件、網頁、社交媒體、視訊、音樂、圖片等等，這些非結構化的資料造成儲存（storage）、探勘（mining）、分析（analyzing）上的困難。當然還有原本就已經儲存在傳統資料庫的資料。
- 資料集大小增長的部分原因來自於資訊持續從各種來源被廣泛收集，這些來源包括搭載感測裝置的行動裝置、高空感測科技（遙感）、軟體記錄、相機、麥克風、無線射頻辨識（RFID）和無線感測網路。
- 自1980年代起，現代科技可儲存資料的容量每40個月即增加一倍；截至2012年，全世界每天產生 $2.5 \text{ EB}$  ( $2.5 \times 10^{18}$ ) 的資料。

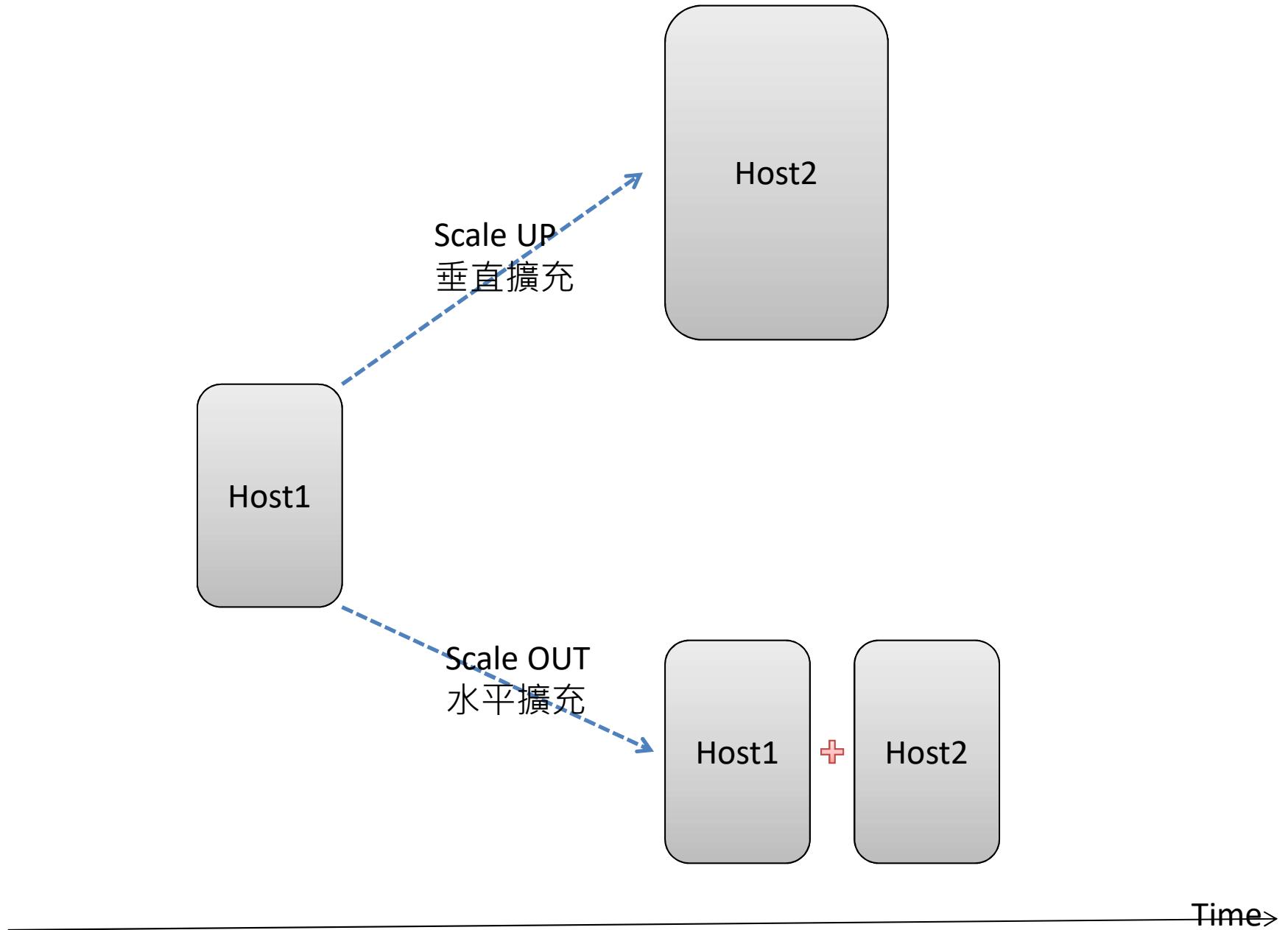
# How fast is Velocity

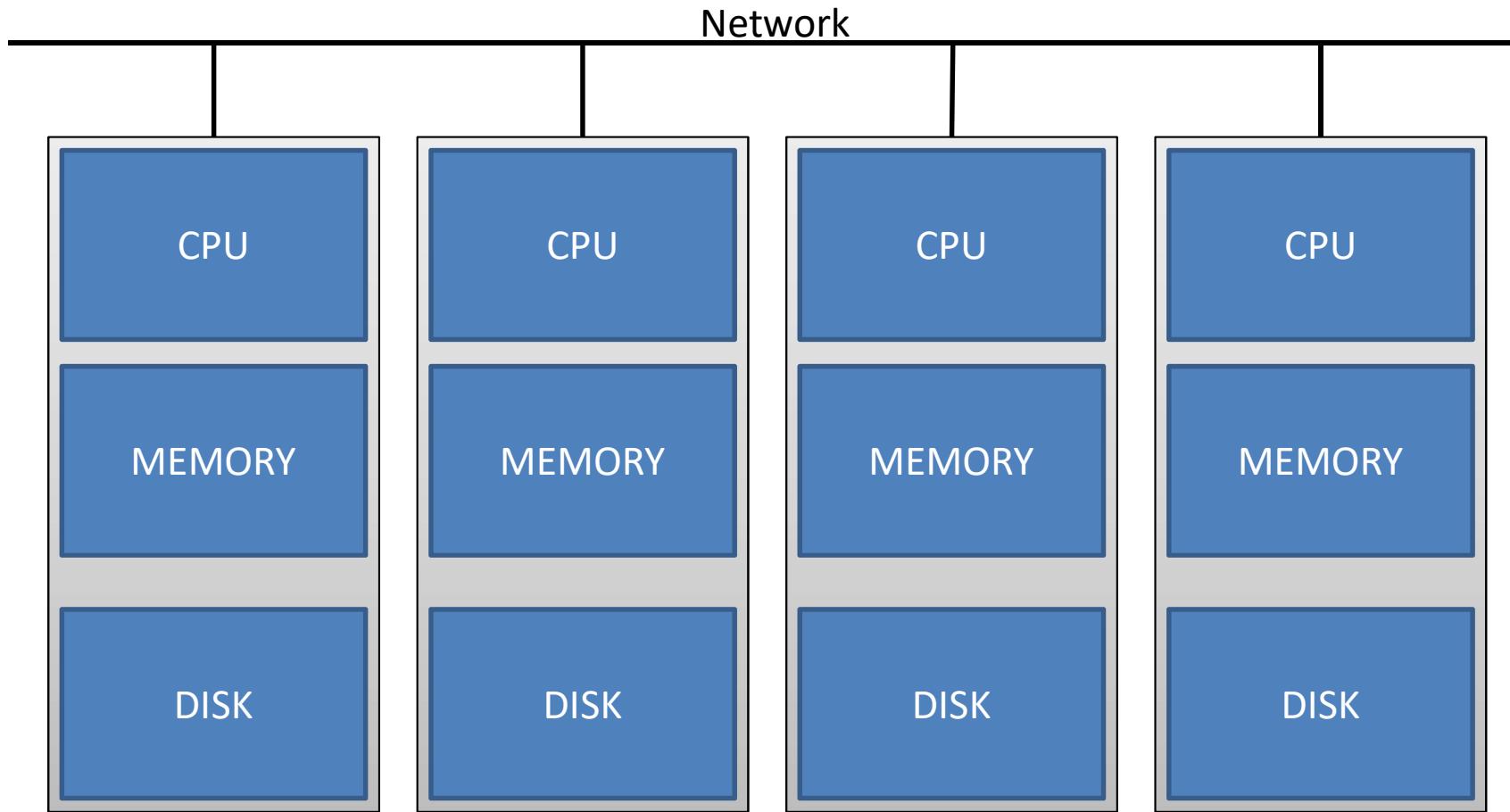
- 資料的傳輸流動（ data streaming ）是連續且快速的，隨著越來越多的機器、網路使用者，社群網站、搜尋結果每秒都在成長，每天都在輸出更多的內容。公司跟機構要處理龐大的資訊大潮向他們襲來，而回應、反應這些資料的速度也成為他們最大的挑戰，許多資料要能即時得到結果才能發揮最大的價值，因此也有人認為 Velocity 是「時效性」。

Big Data:  
Expanding on 3 fronts  
at an increasing rate.

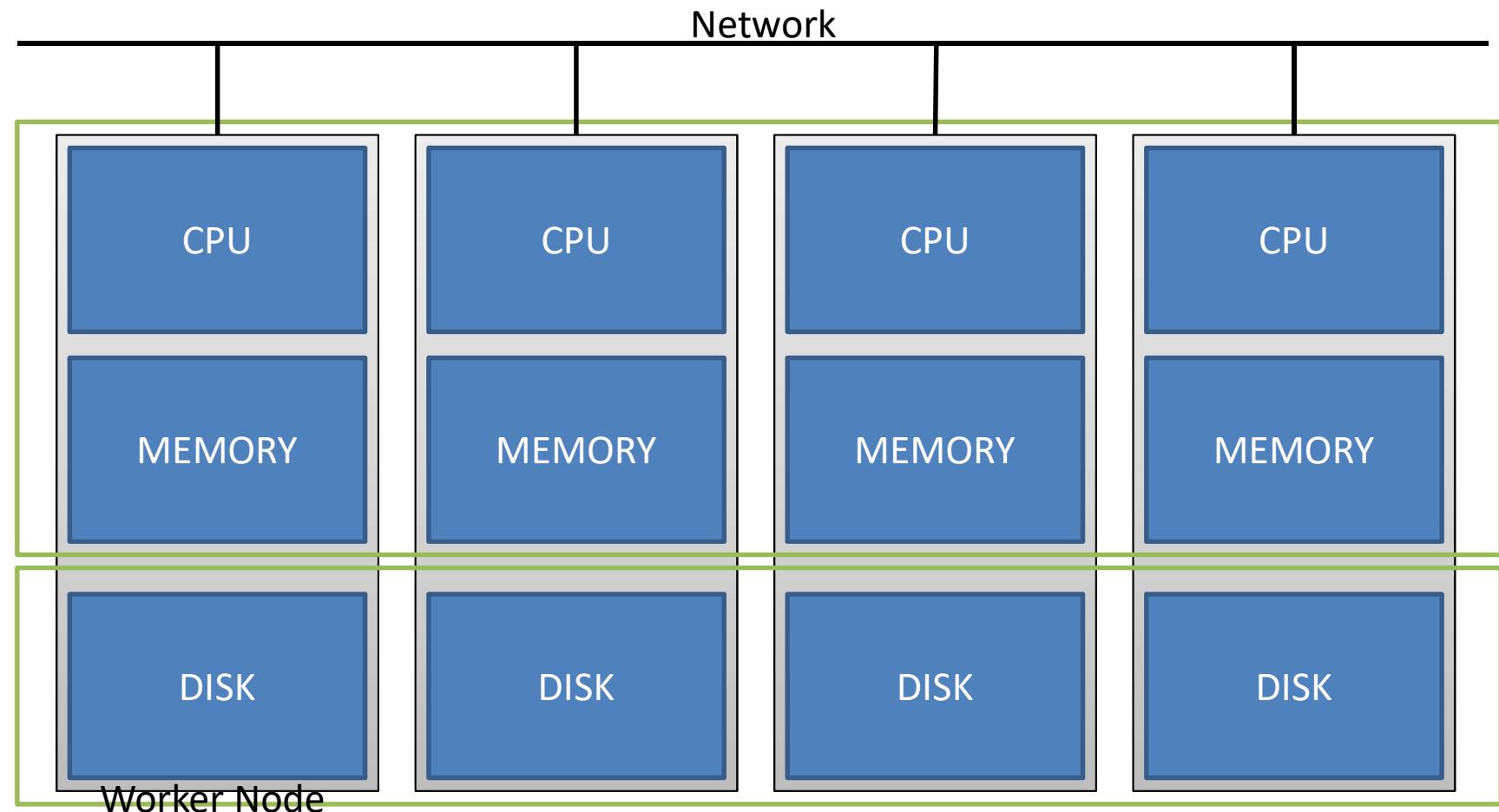






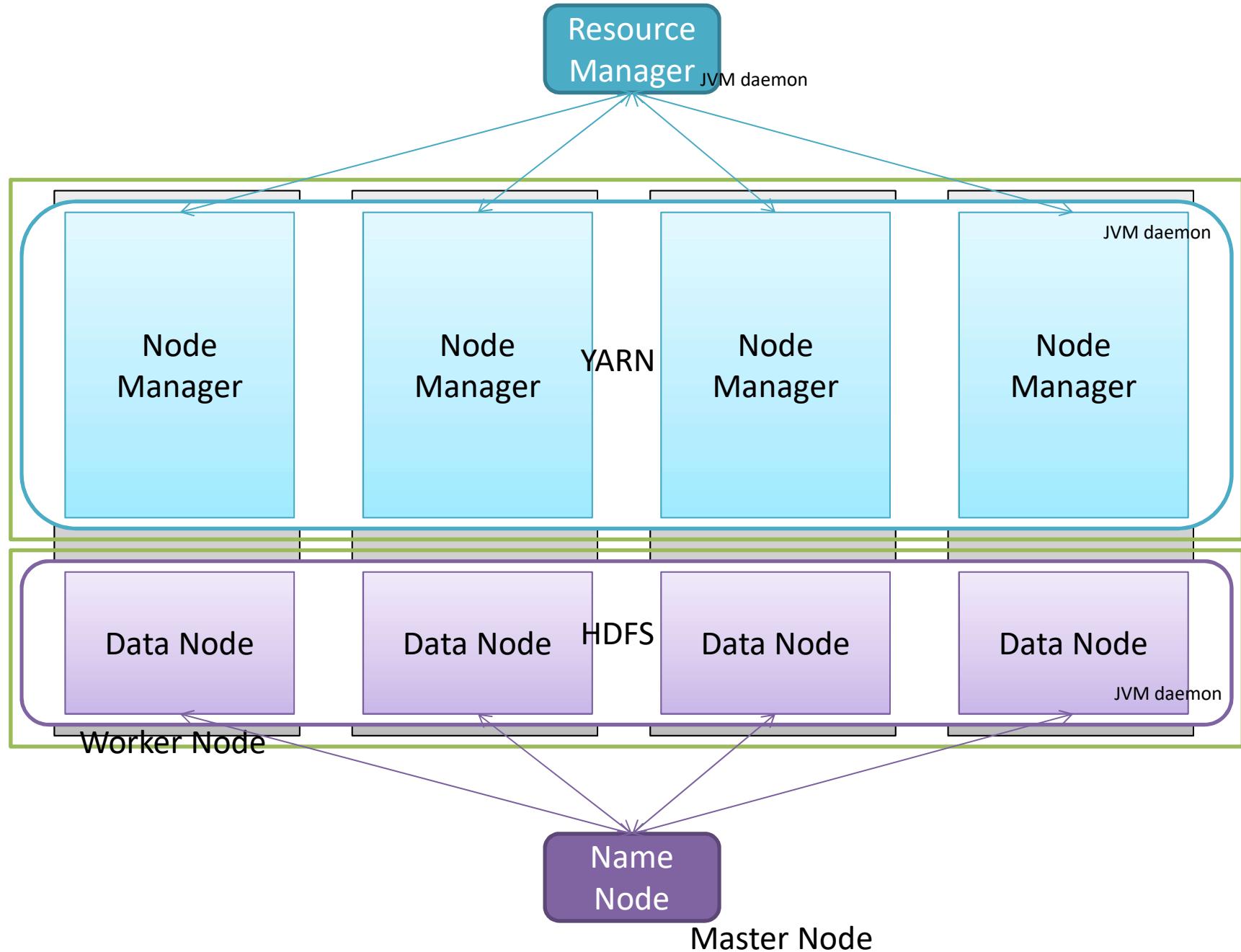


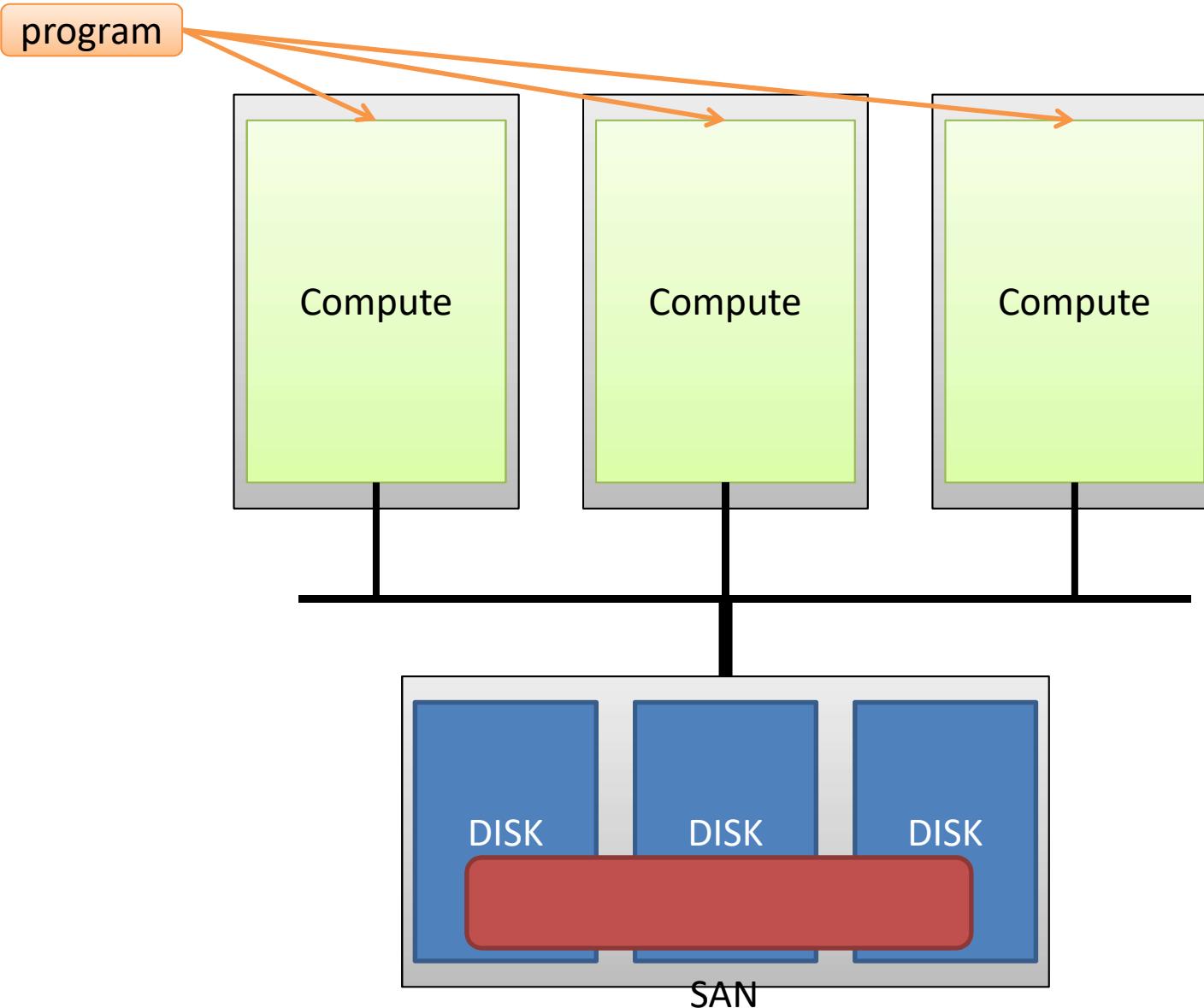
Share Nothing Cluster

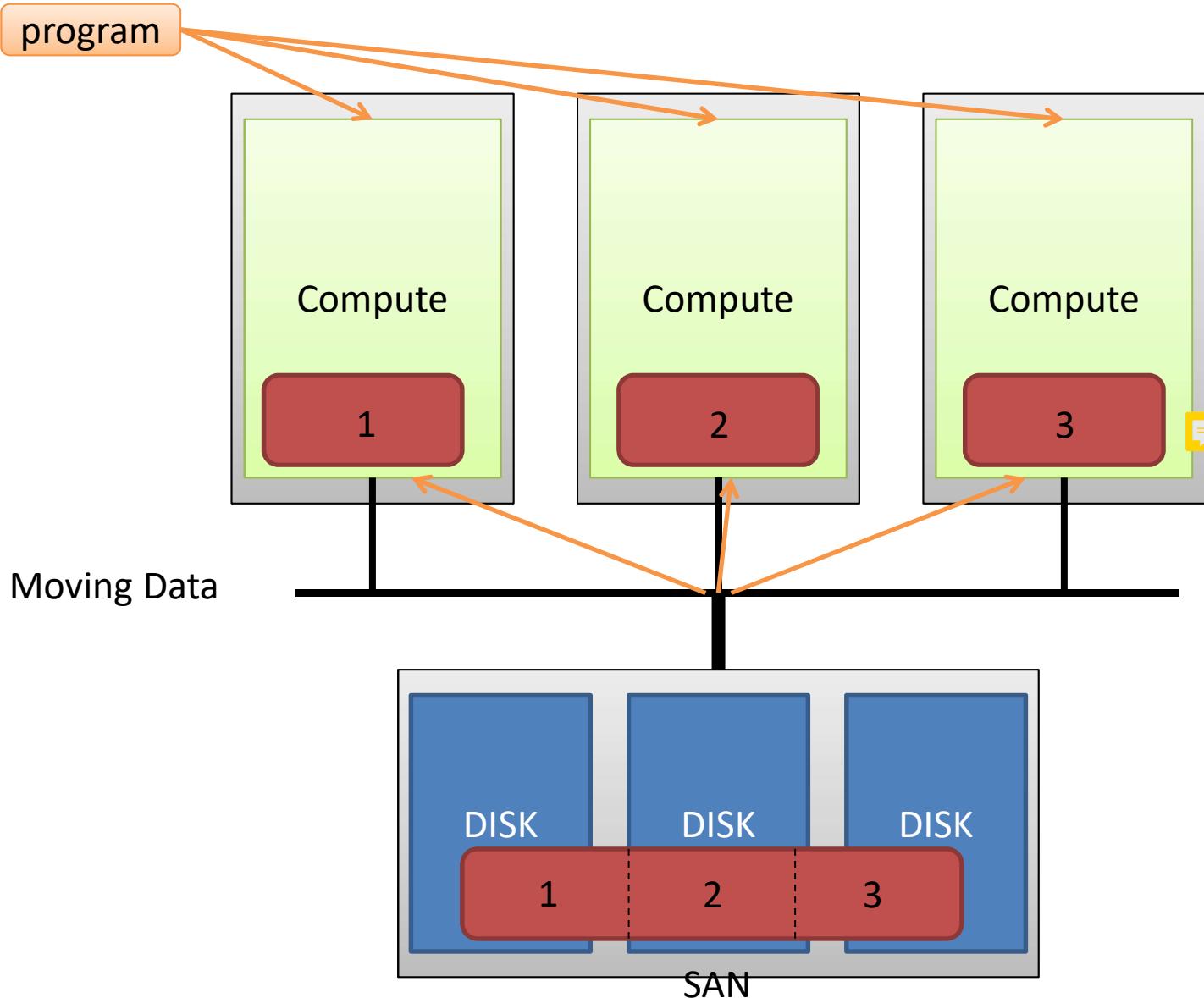


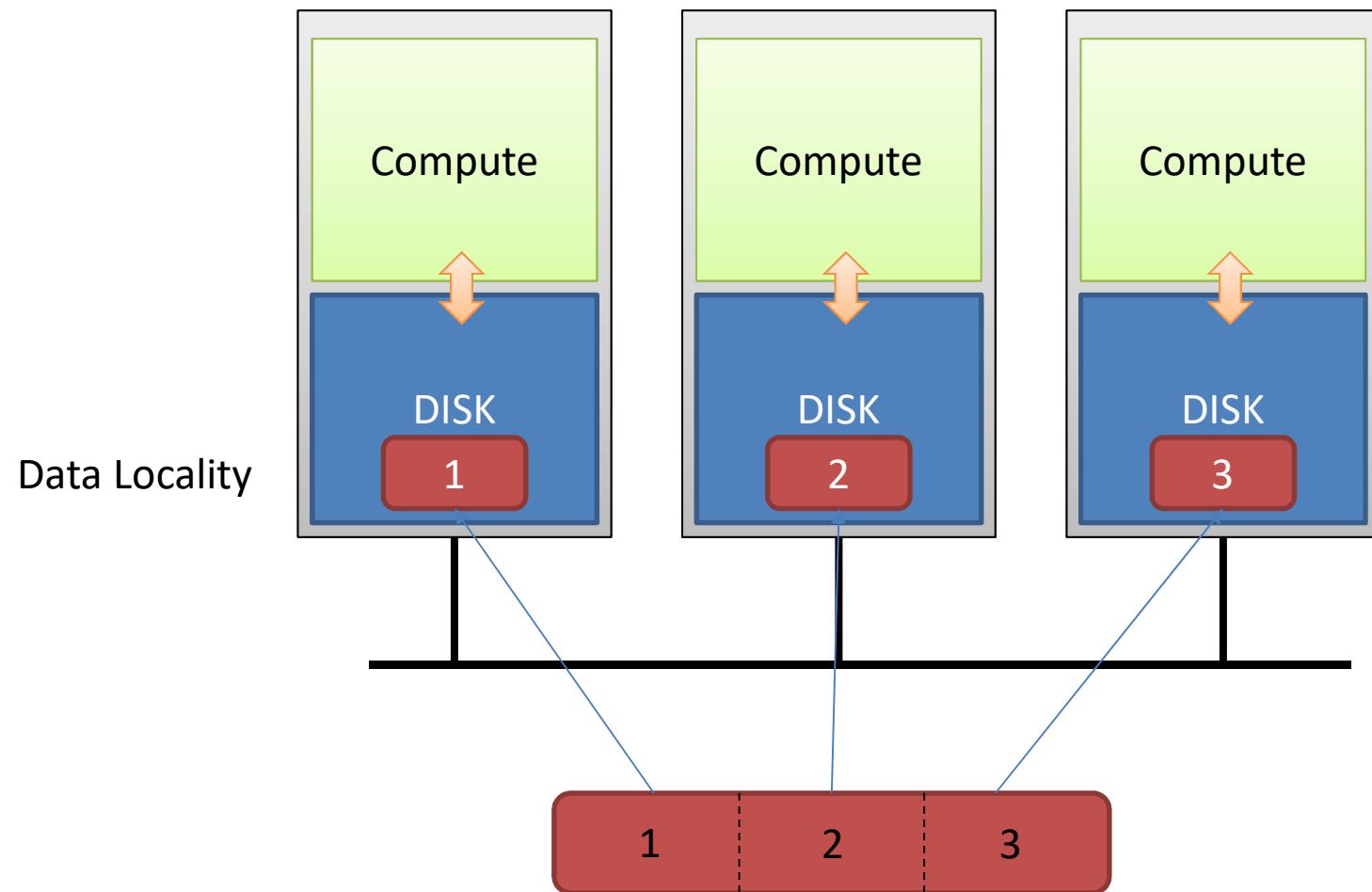
# Hadoop components

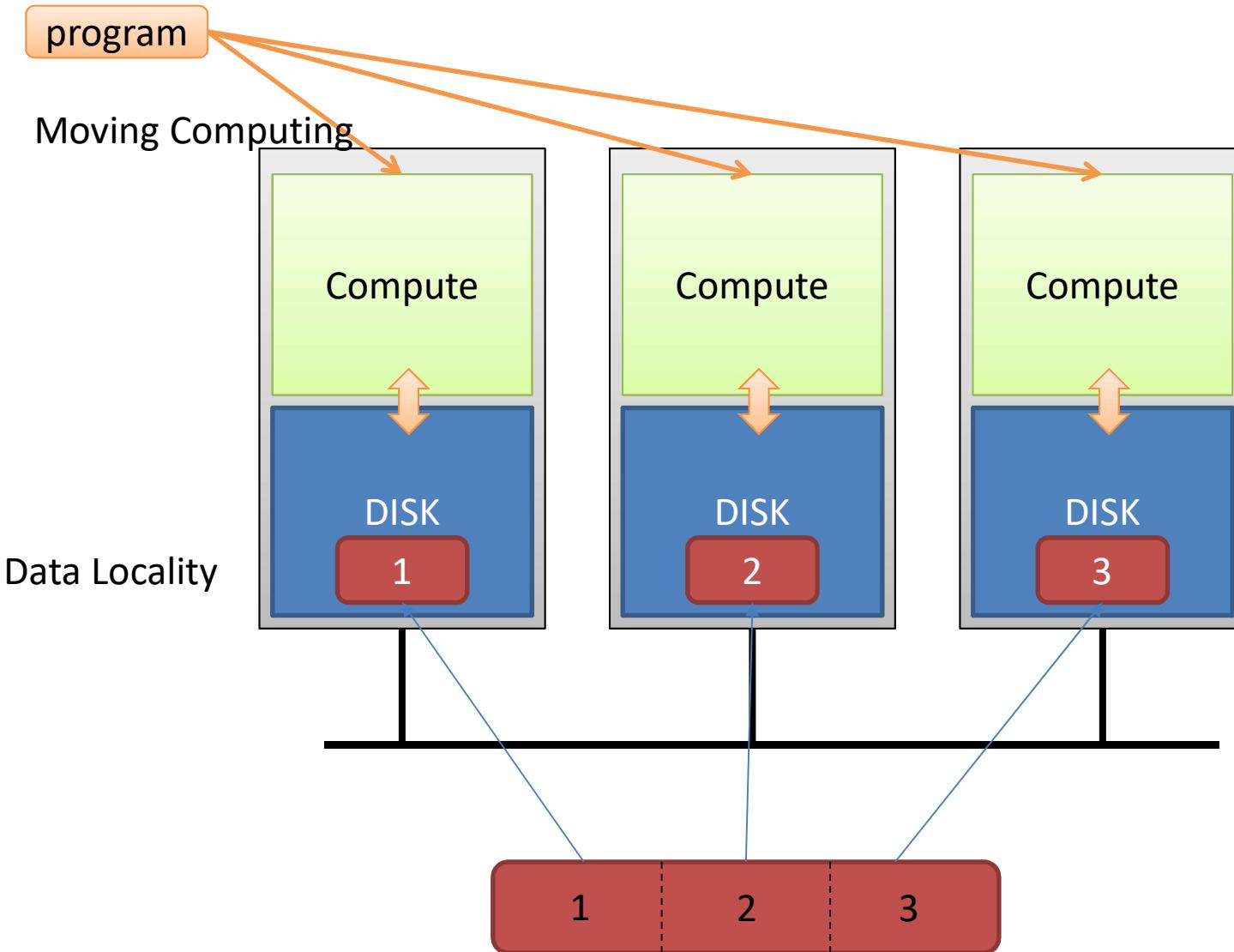
- Master Node 
  - HDFS : Name Node
  - YARN : Resource Manager/Job History
- Worker Node 
  - HDFS : Data Node
  - YRAN : Node Manager 





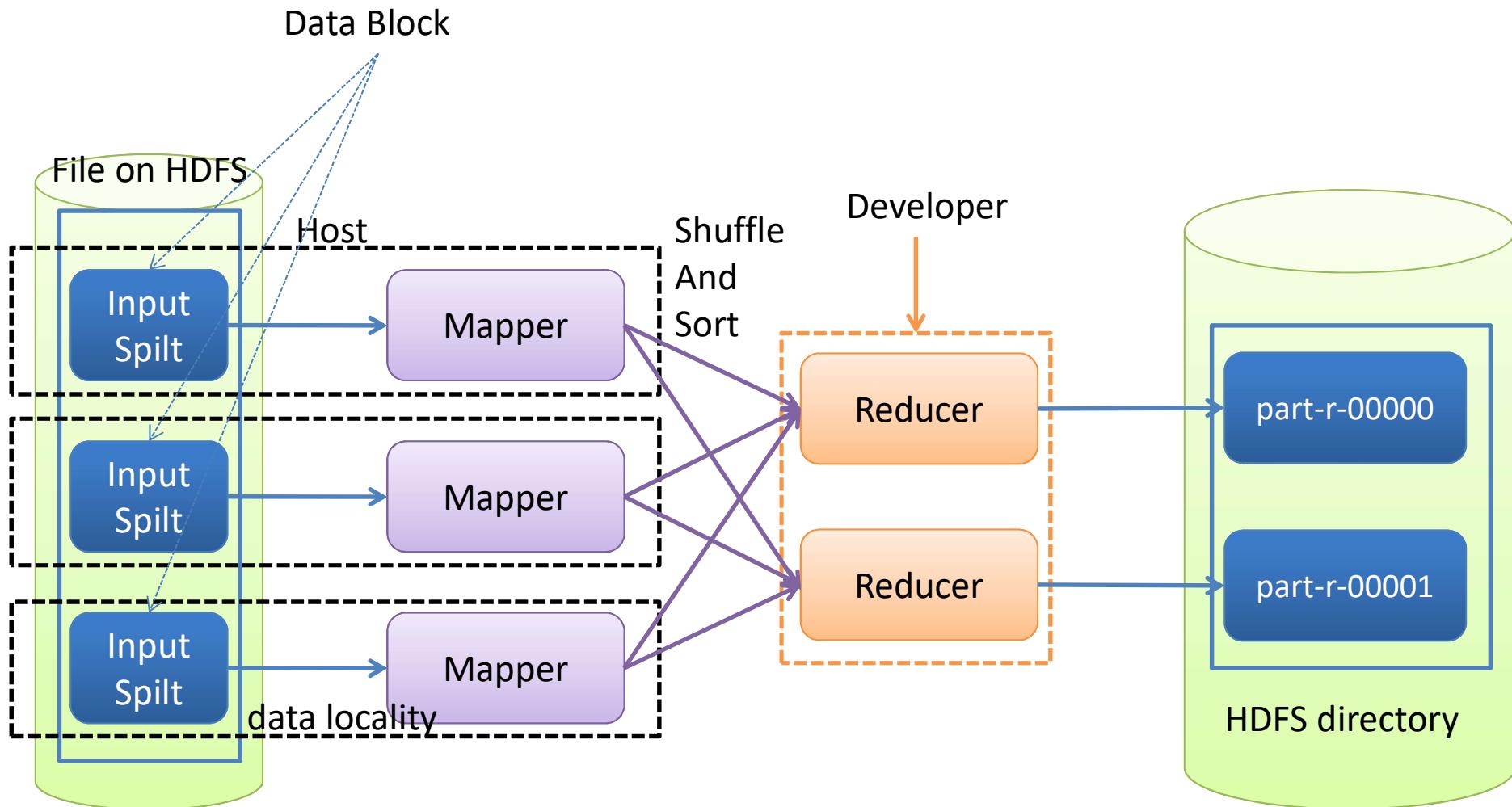


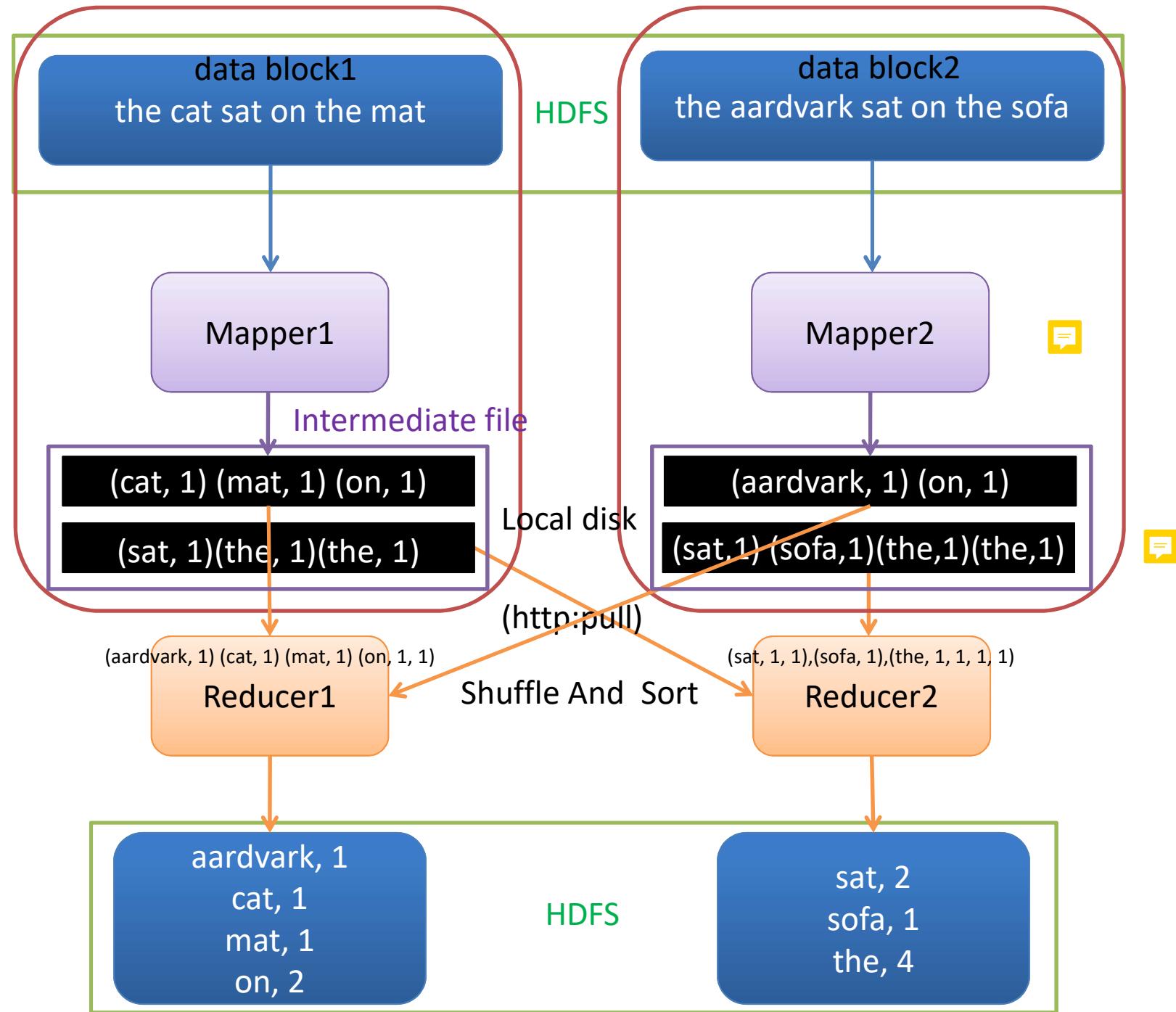




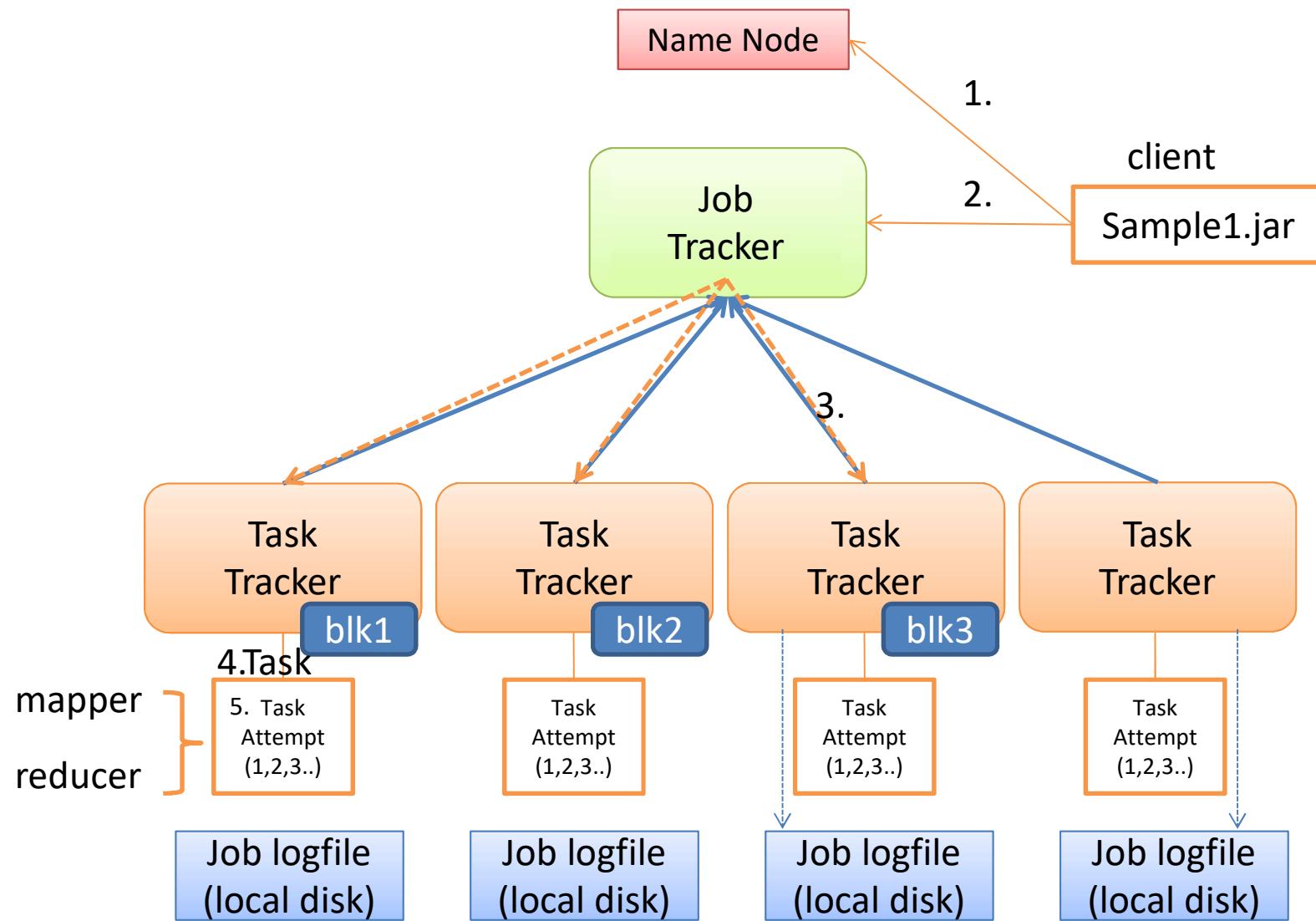
# HDFS

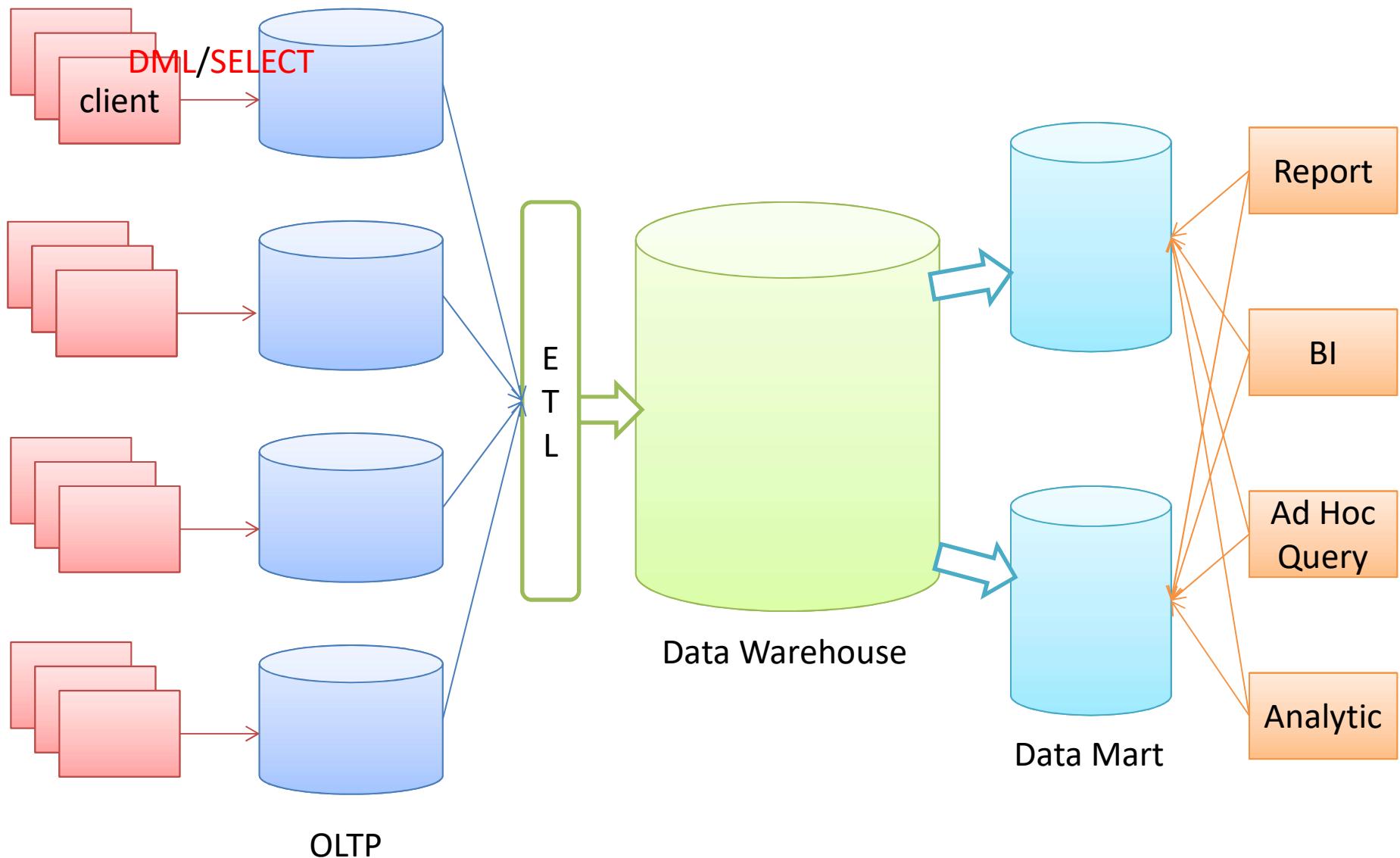
- \$ hdfs dfs -put f30M.txt /user/training/
  - hdfs file : f30M.txt (1 hdfs block : 30M)
- \$ hdfs dfs -put f100M.txt /user/training/
  - hdfs file : f100M.txt (1 hdfs block : 100M)
- \$ hdfs dfs -put f300M.txt /user/training/
  - hdfs file : f300M.txt (3 hdfs block : 128/128/44M)



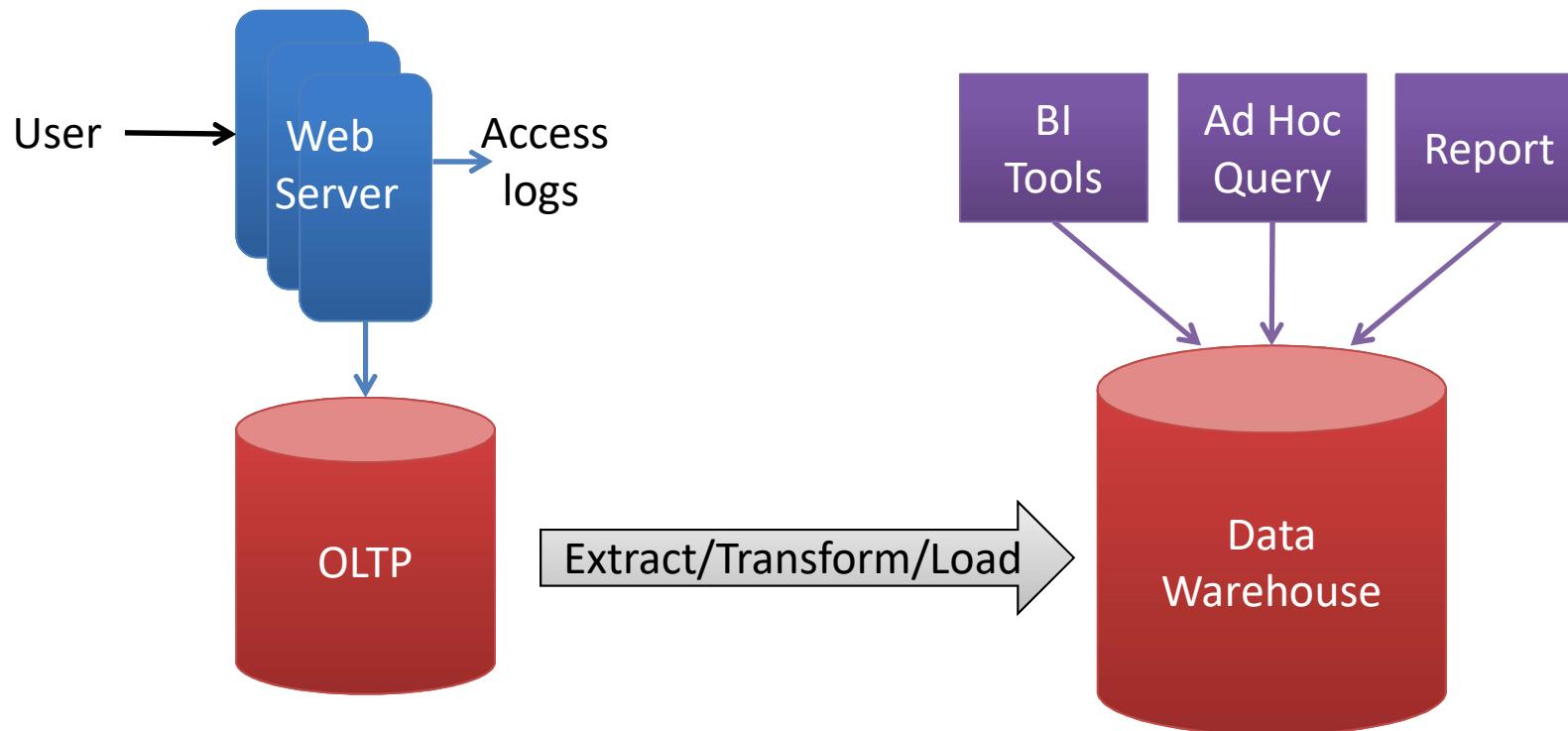


# Map Reduce v1

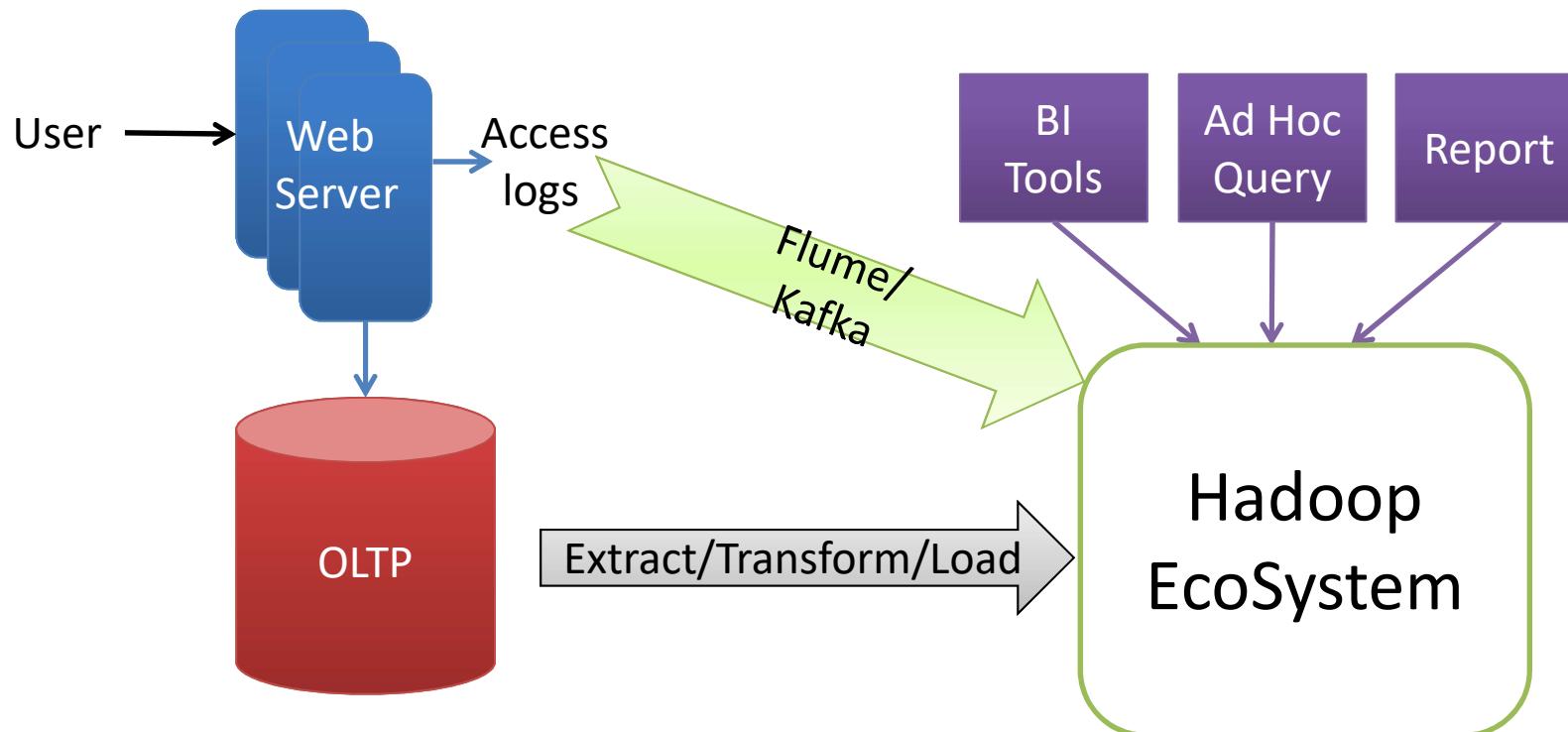




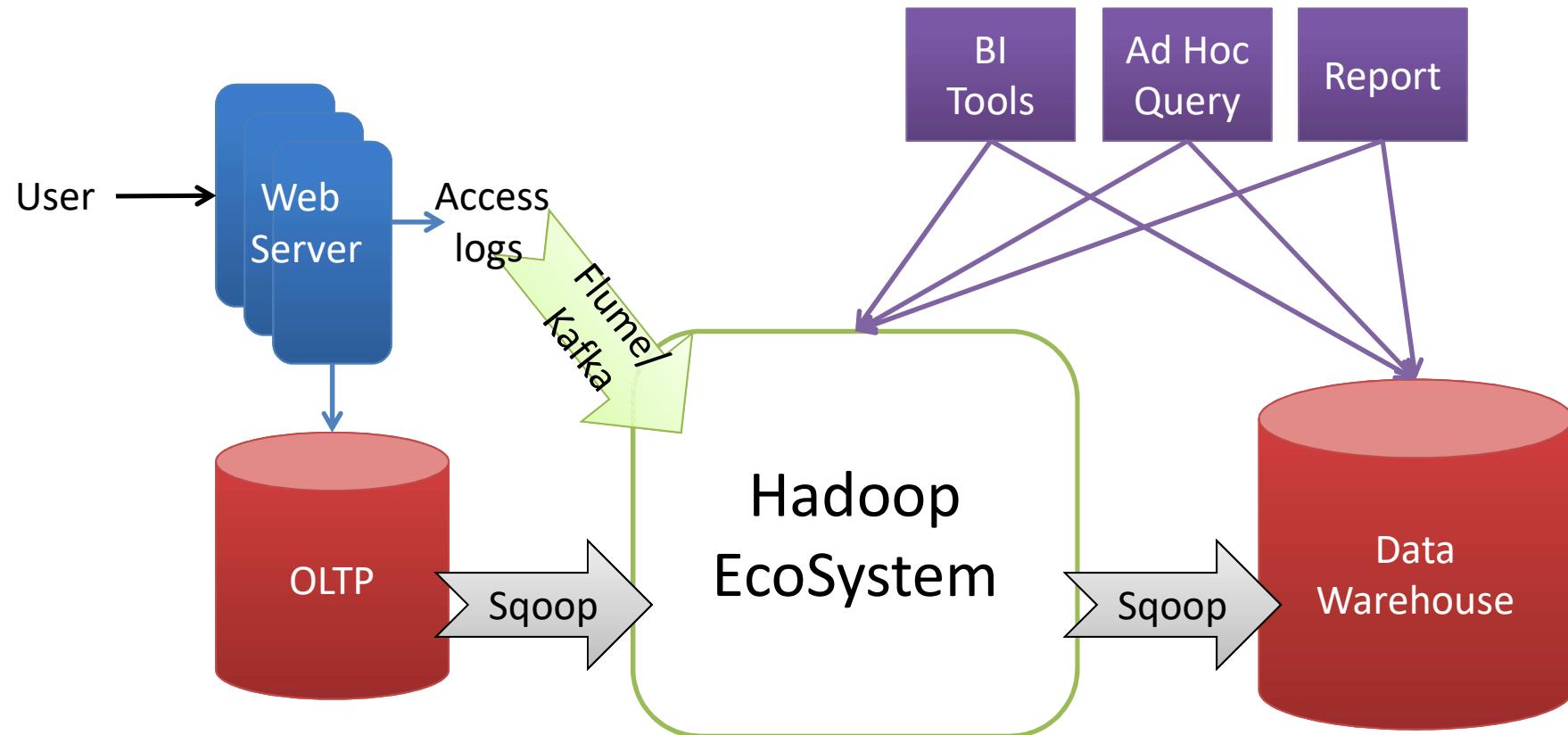
# Data Warehouse based on RDBMS



# Data Warehouse based on Hadoop



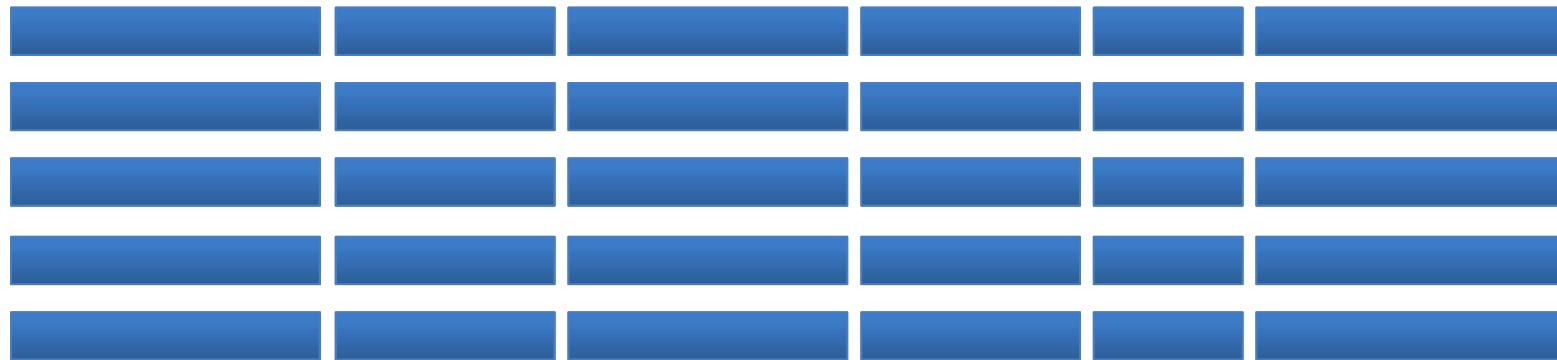
# Integrated RDBMS with Hadoop



# SQL on Hadoop

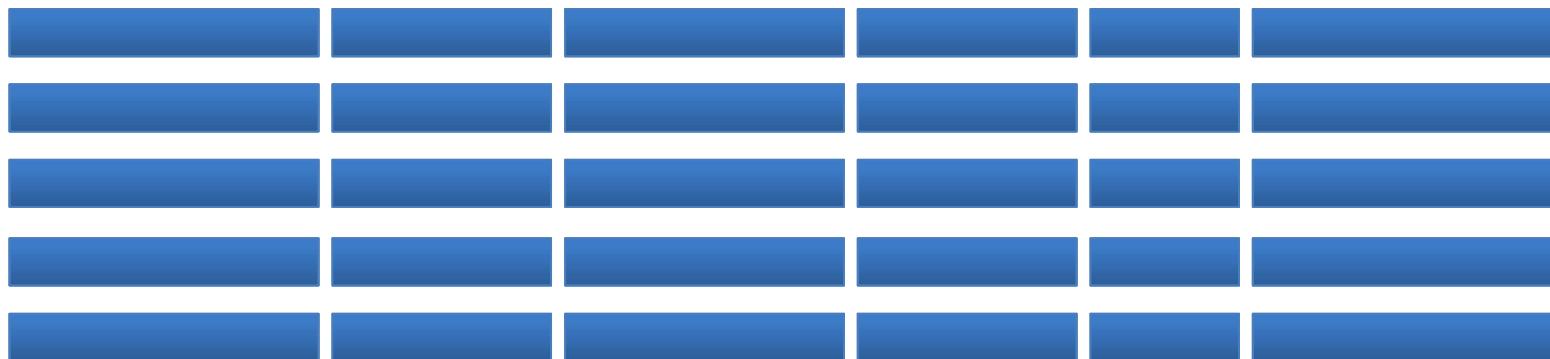
- Batch SQL
  - Hive(Facebook)
    - Based on MapReduce
    - Based on Tez(Hortonworks)
- Interactive SQL
  - Spark SQL(Hive based on Spark)
  - Impala(Cloudera)(Google Dremel)
  - Apache Drill(Google Dremel)
  - Presto(Facebook)
  - Pinot(LinkedIn:OLAP)
  - Kylin(Ebay China:OLAP)
- Operational SQL
  - Apache Phoenix(SQL on Hbase)

# Schema on Write

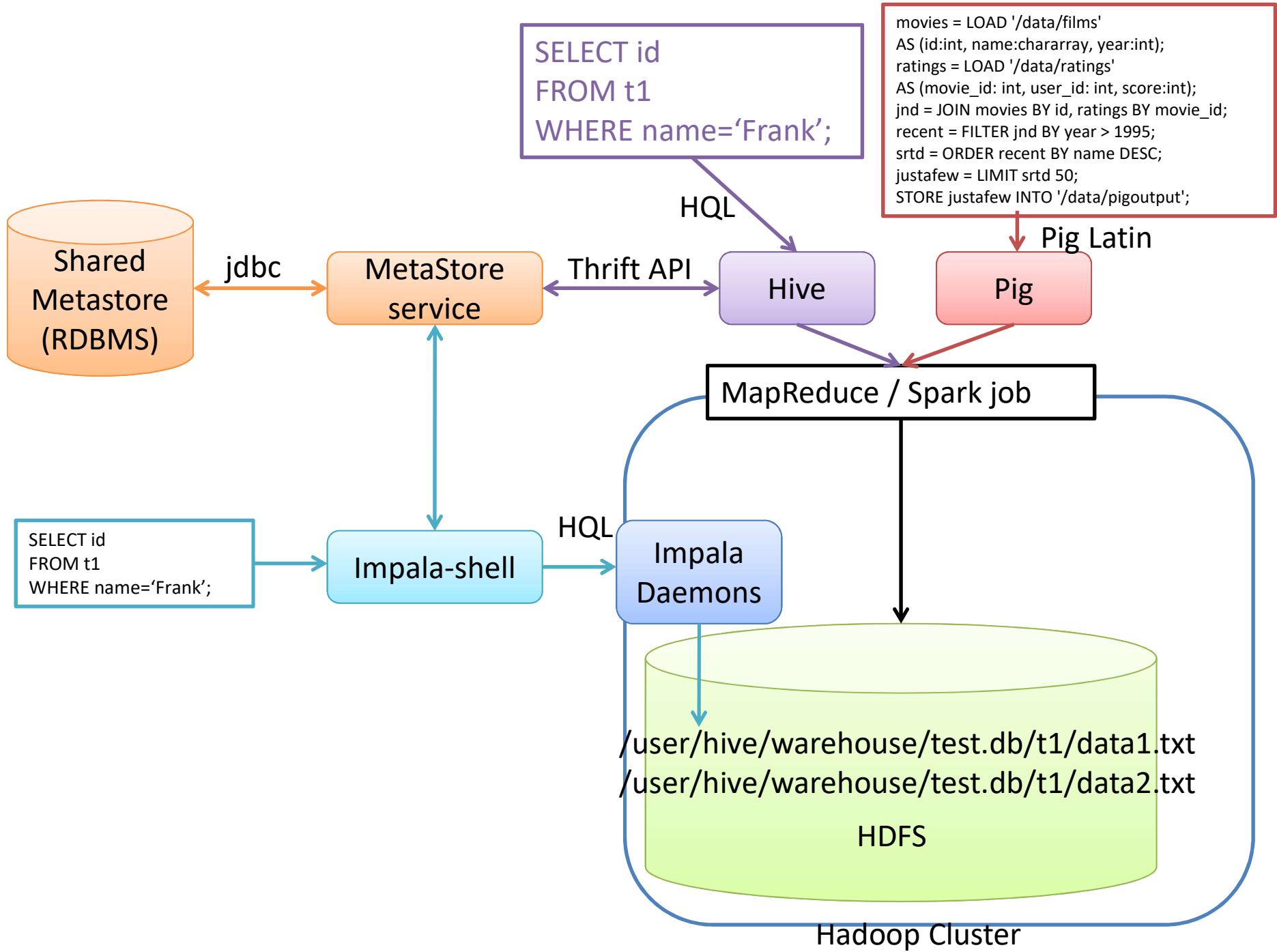


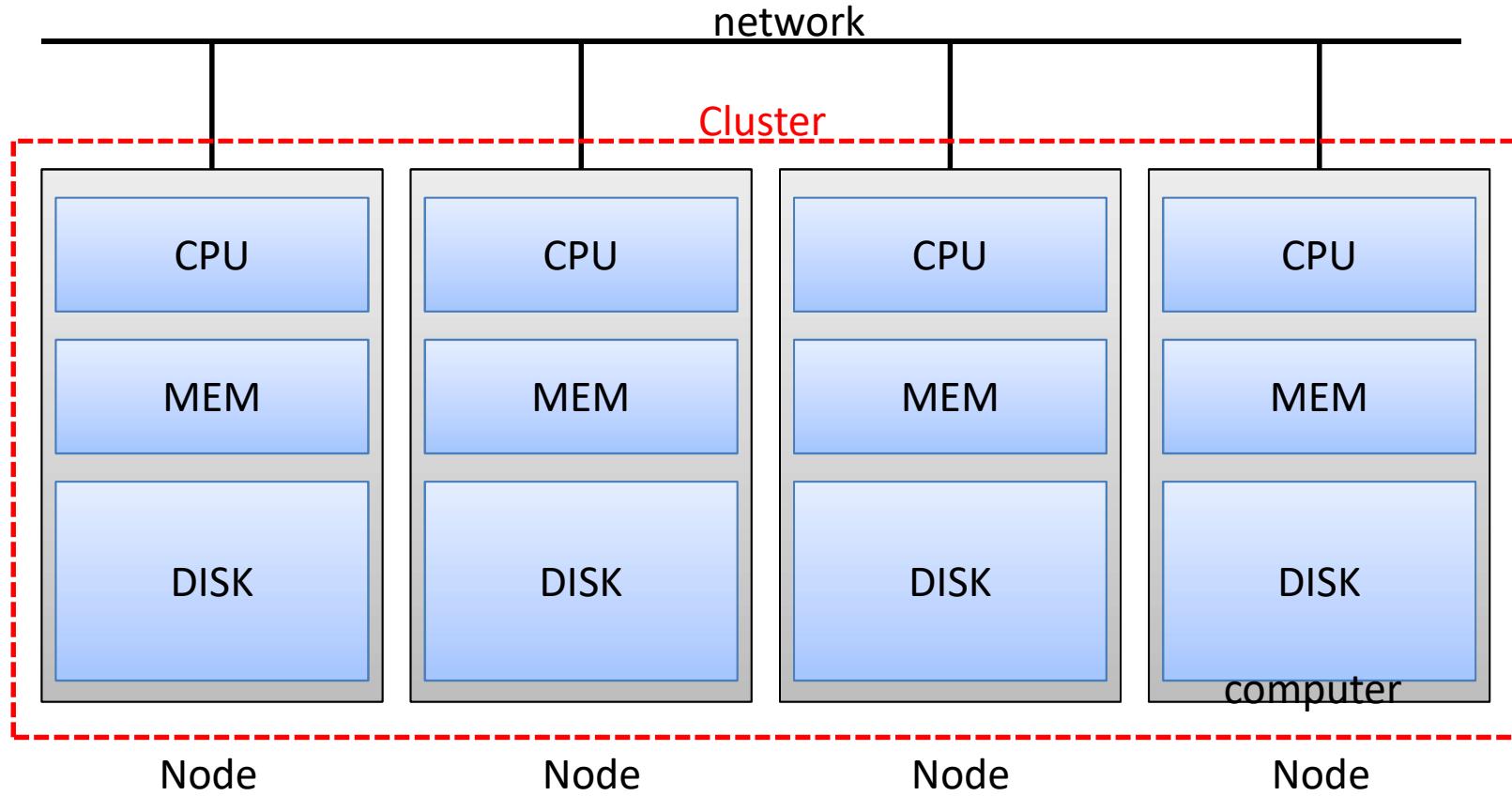
Empid	Ename	Salary	Hire_date

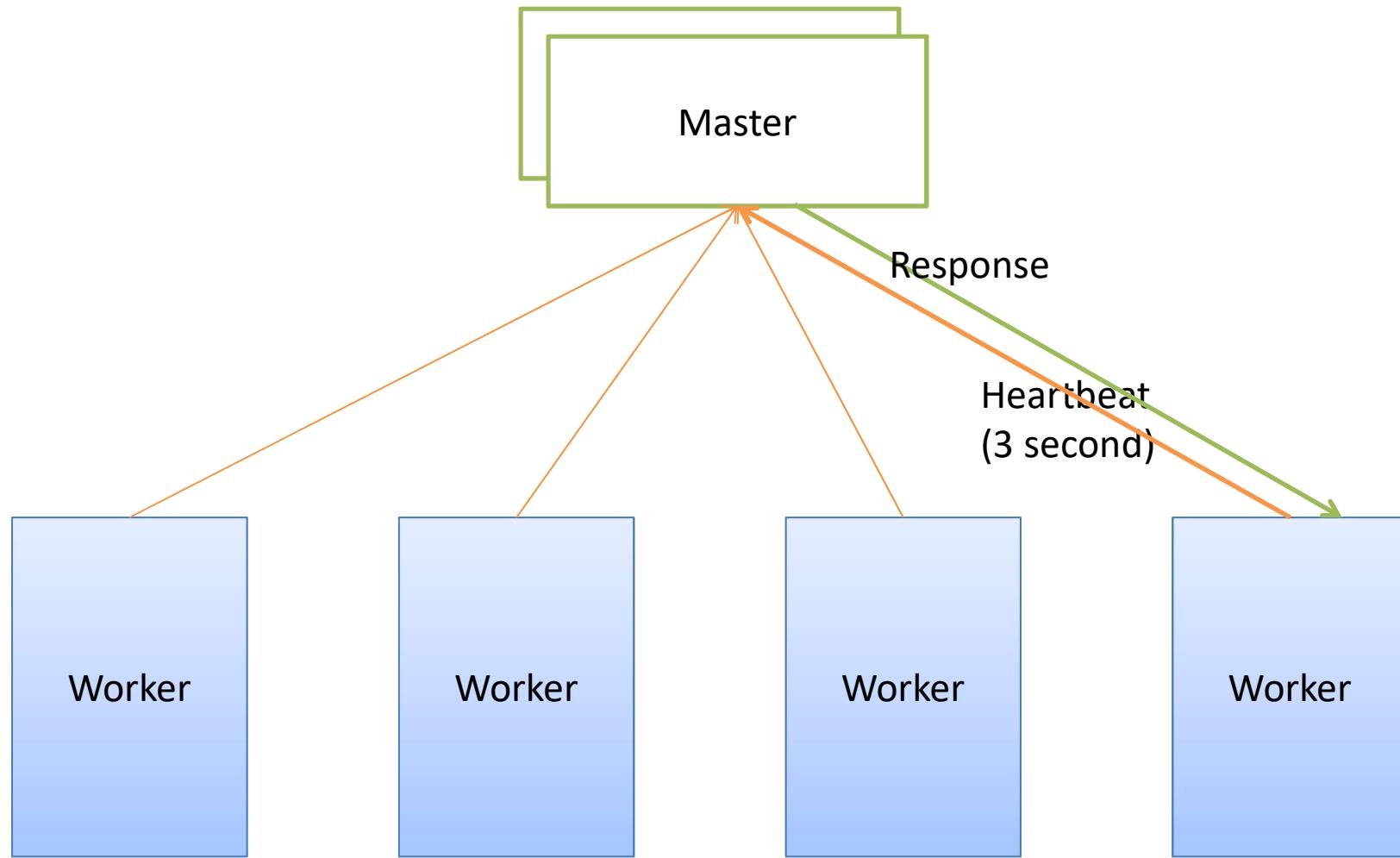
# Schema on Read

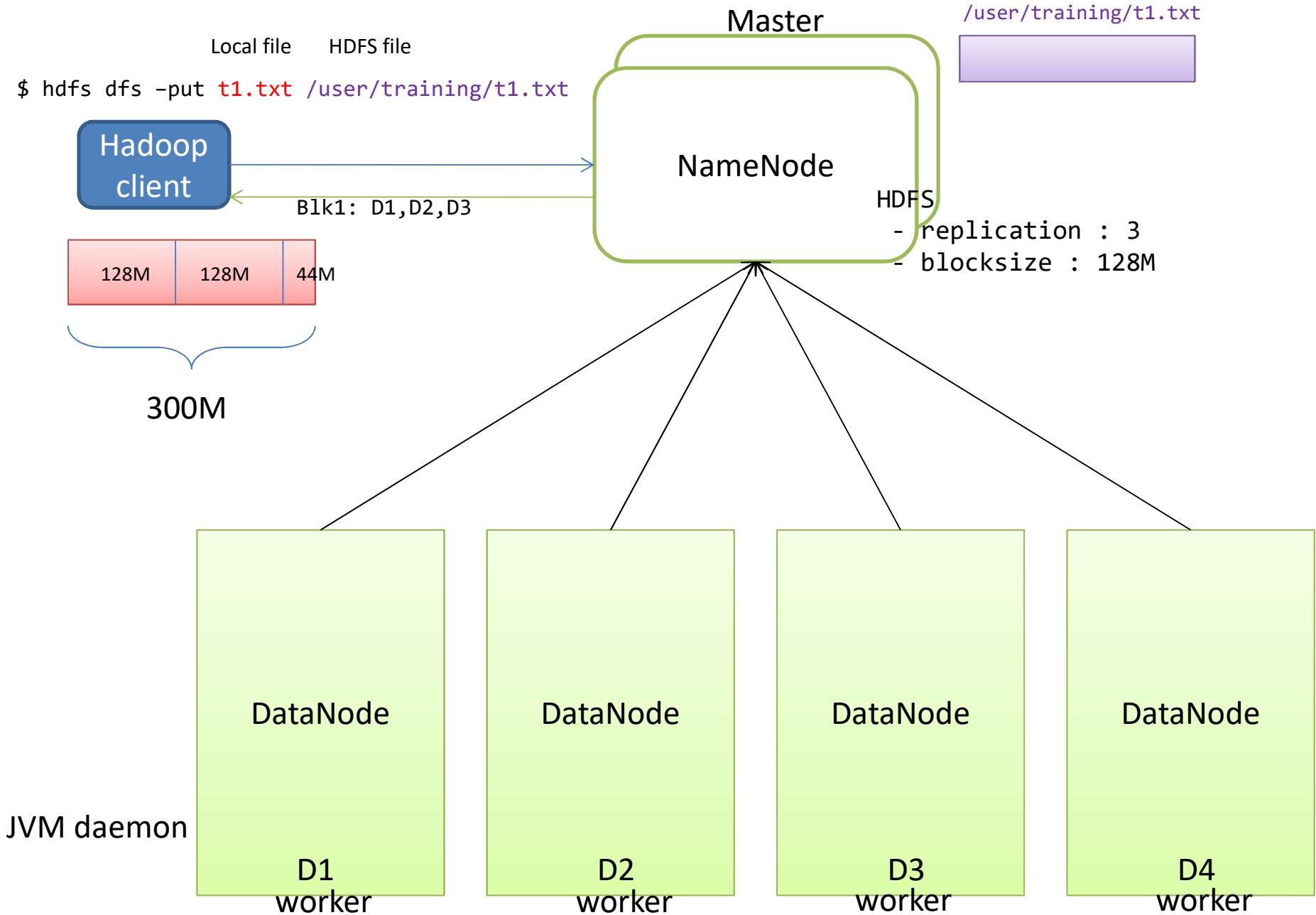


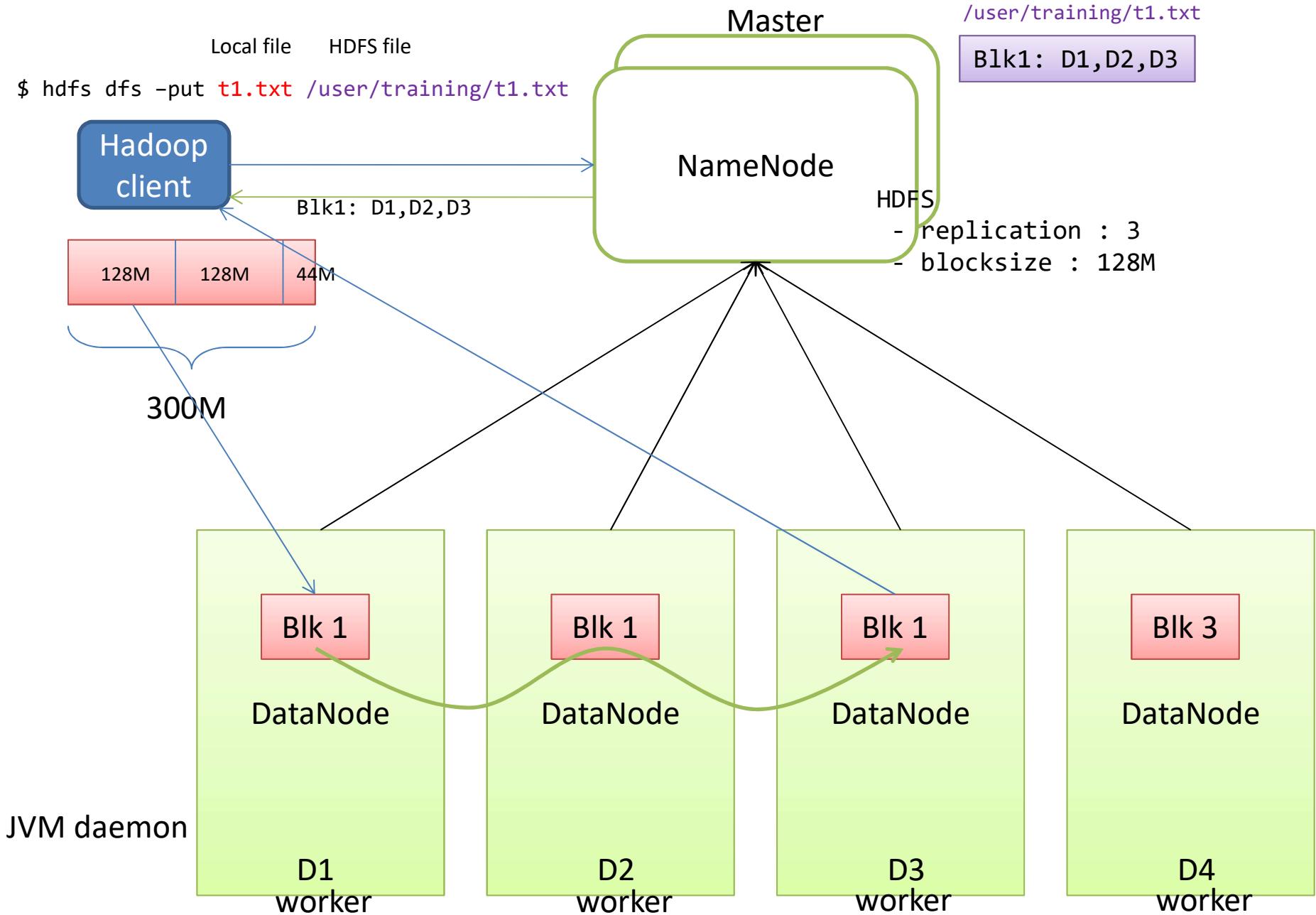
Empid	Ename	Job	Salary	Dept	Hire_date

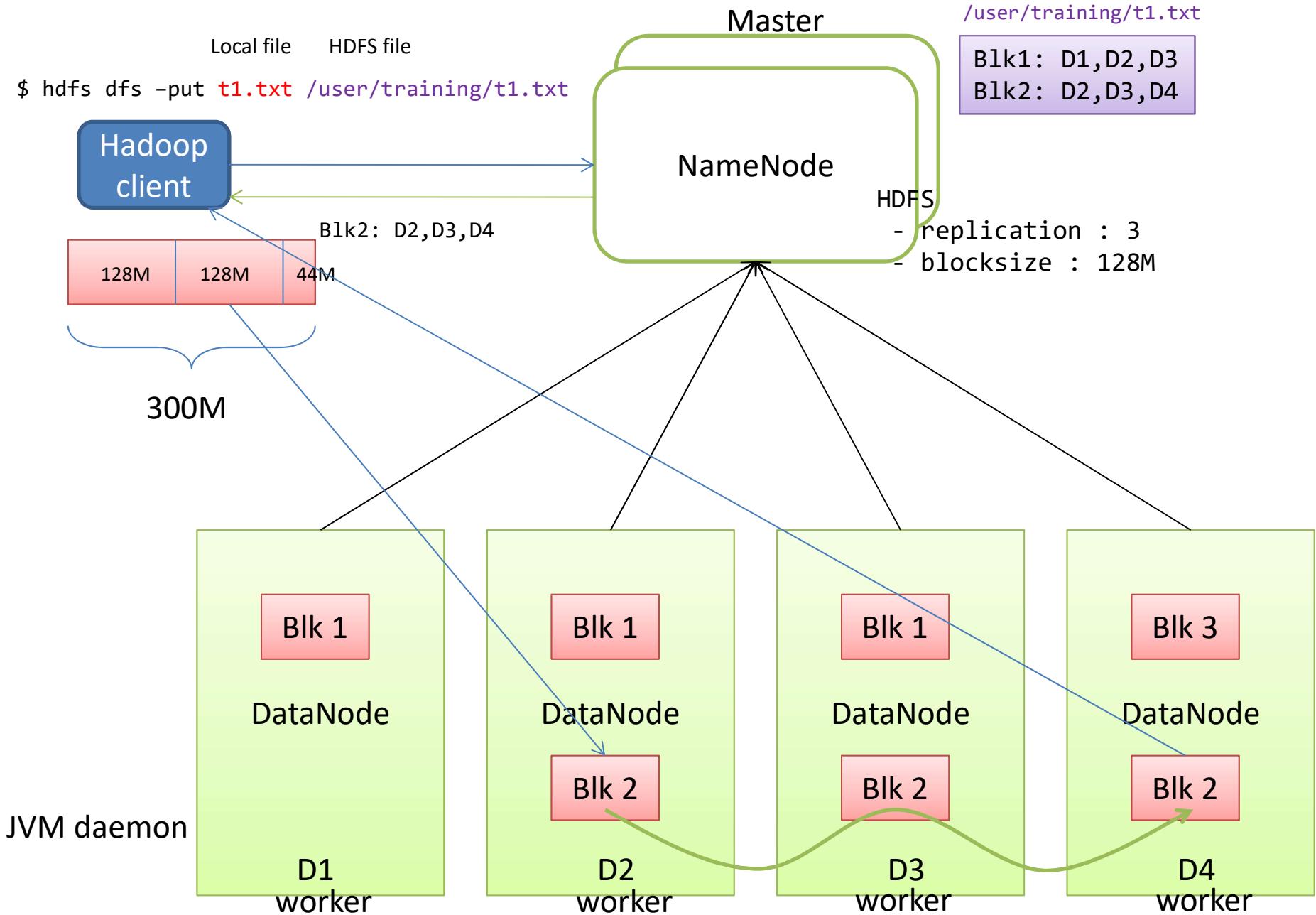


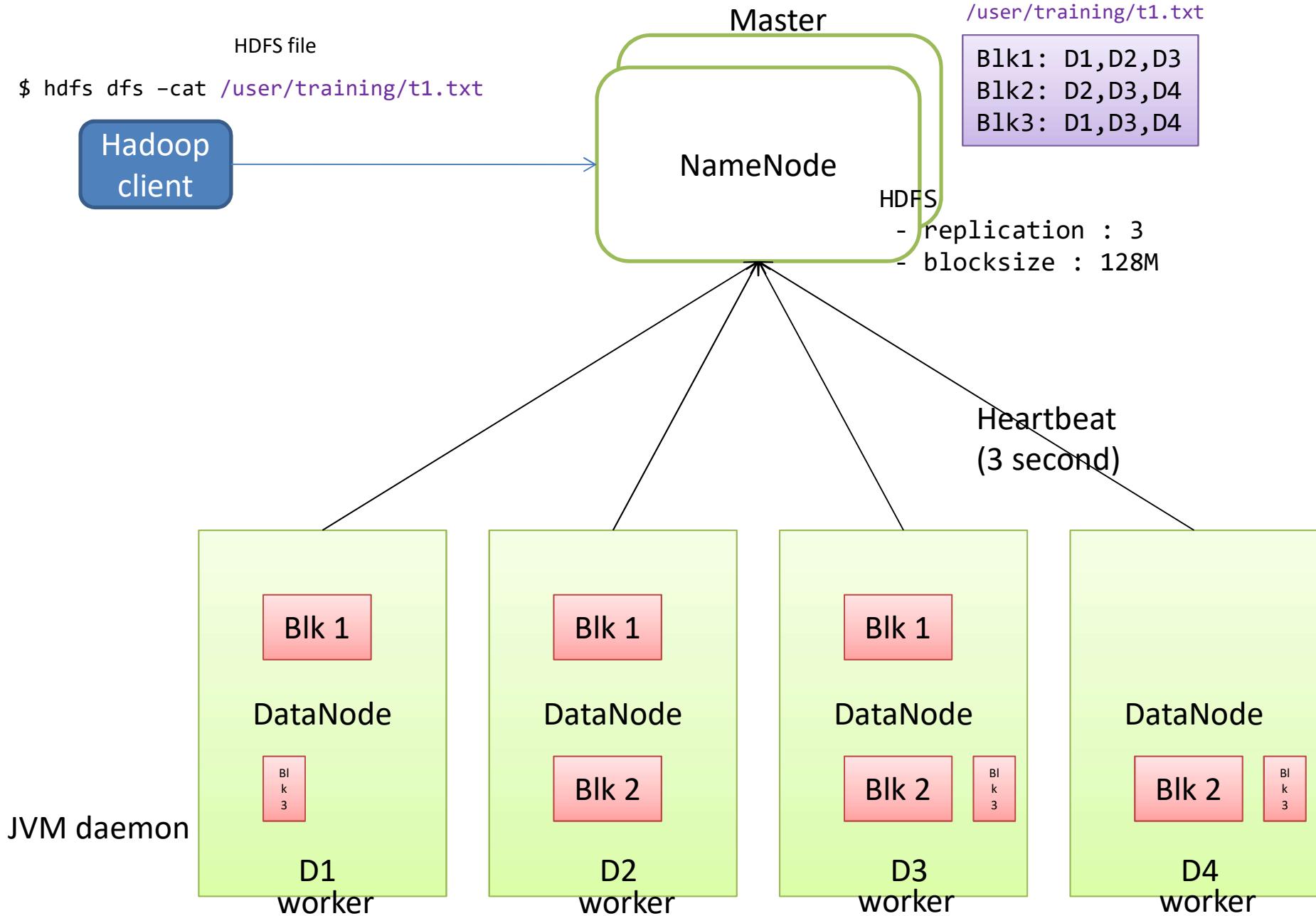


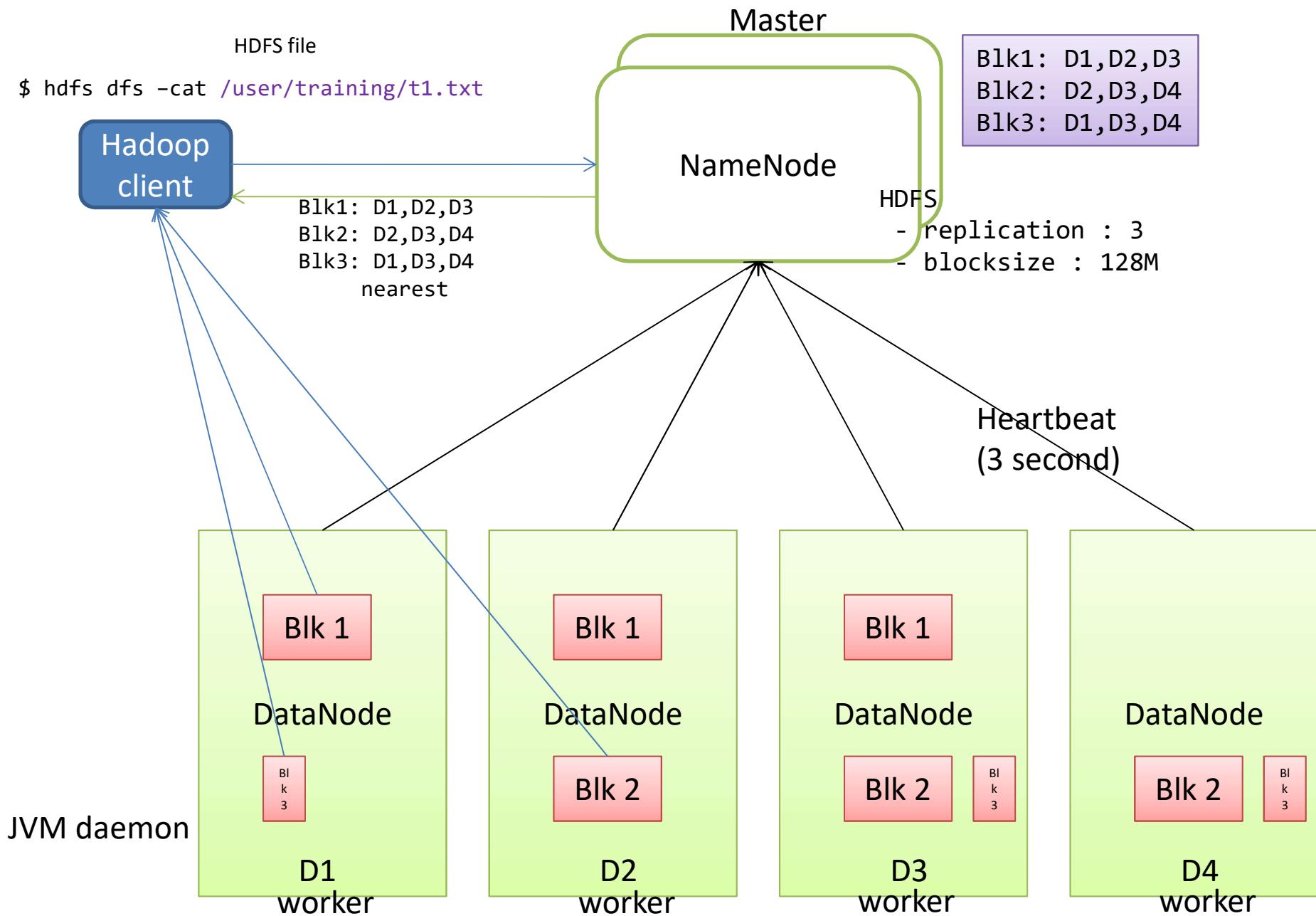


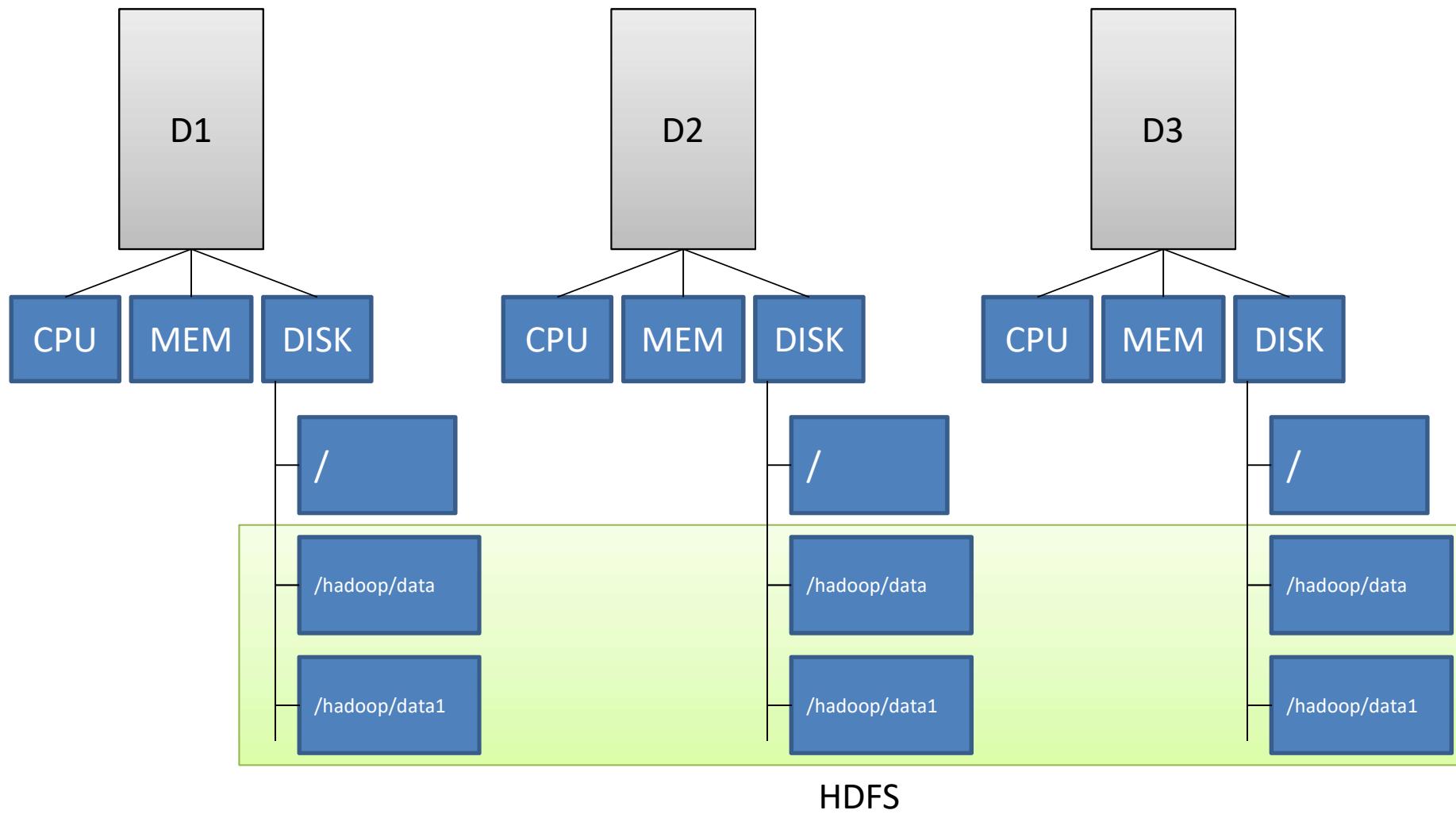


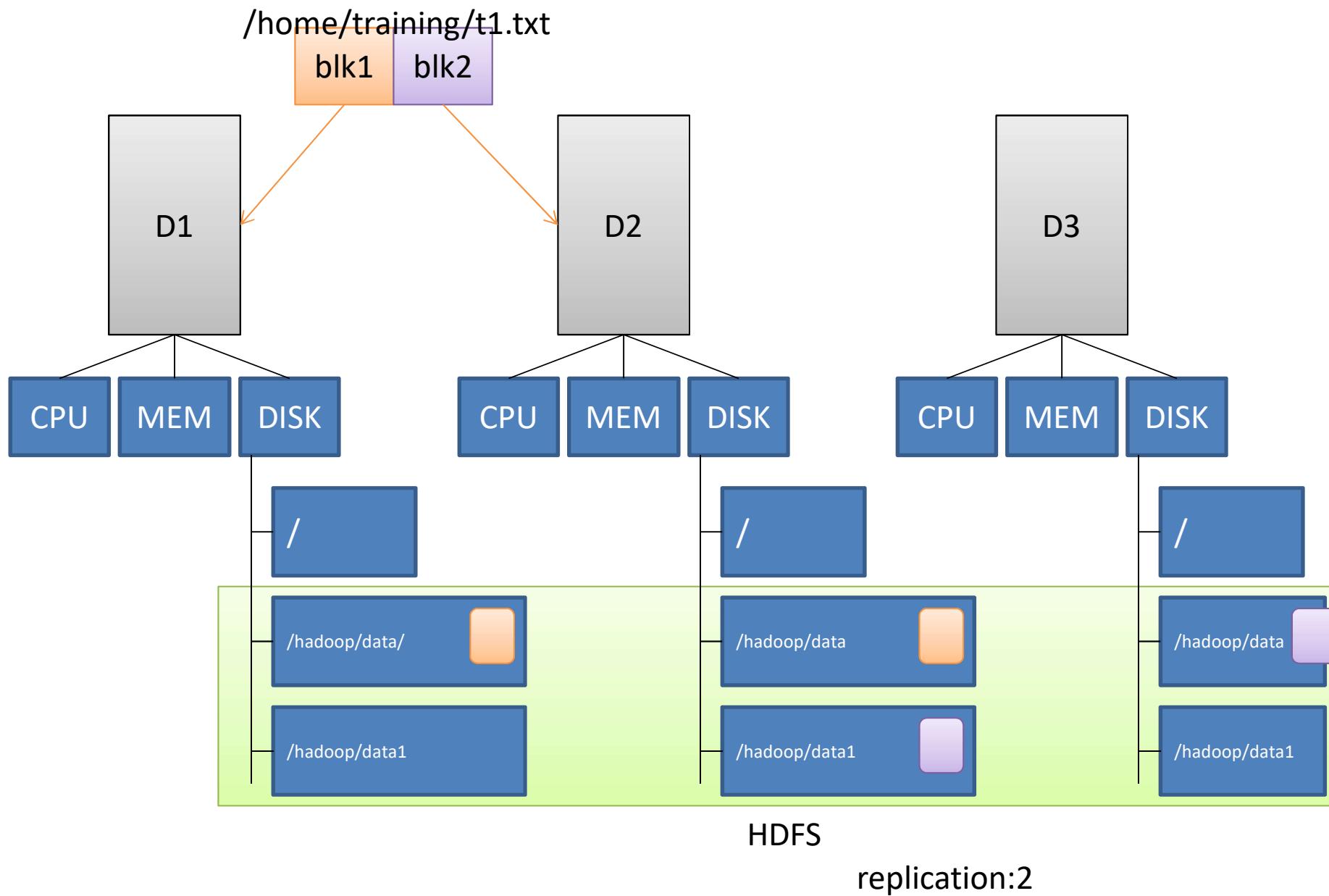












# DataNode

- Store block
- Block Report after DataNode startup
- Heart Beat(3 seconds) to NameNode

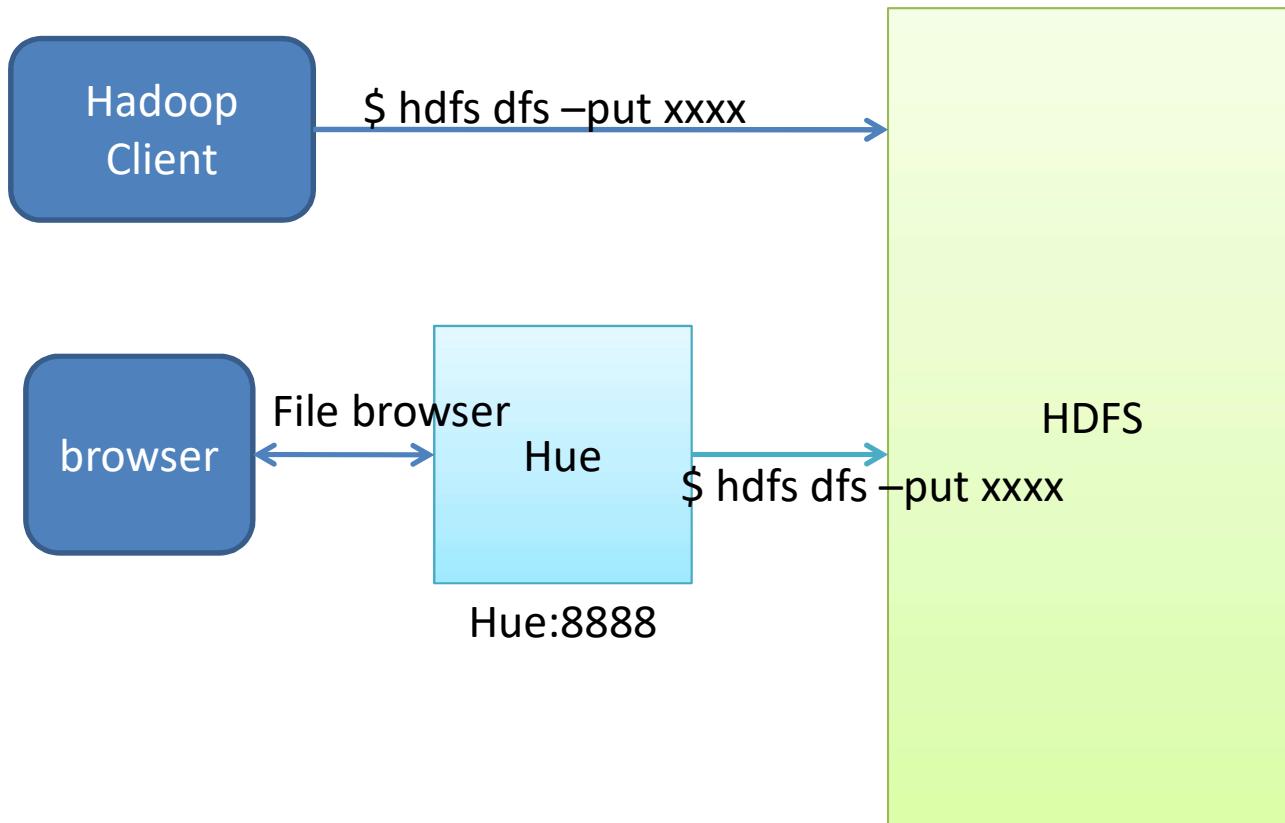
```
[cloudera@quickstart ~]$ cat /etc/hadoop/conf/hdfs-site.xml|grep fs.datanode.data.dir -B 1 -A 2
<property>
  <name>dfs.datanode.data.dir</name>
  <value>/var/lib/hadoop-hdfs/cache/${user.name}/dfs/data</value>
</property>
[cloudera@quickstart ~]$ sudo ls -l /var/lib/hadoop-hdfs/cache/hdfs/dfs/data/current/
total 8
drwx----- 4 hdfs hadoop 4096 Oct  1 06:23 BP-989008105-127.0.0.1-1433846136903
-rw-r--r-- 1 hdfs hadoop  229 Oct  1 06:23 VERSION
[cloudera@quickstart ~]$ sudo ls -l /var/lib/hadoop-hdfs/cache/hdfs/dfs/data/current/BP-989008105-127
.0.0.1-1433846136903/current/finalized/subdir0/subdir0 |more
total 220432
-rw-r--r-- 1 hdfs hadoop    11175 Jun  9 03:39 blk_1073741825
-rw-r--r-- 1 hdfs hadoop      95 Jun  9 03:39 blk_1073741825_1001.meta
-rw-r--r-- 1 hdfs hadoop    48741 Jun  9 03:39 blk_1073741826
-rw-r--r-- 1 hdfs hadoop     391 Jun  9 03:39 blk_1073741826_1002.meta
-rw-r--r-- 1 hdfs hadoop   16046 Jun  9 03:39 blk_1073741827
-rw-r--r-- 1 hdfs hadoop     135 Jun  9 03:39 blk_1073741827_1003.meta
-rw-r--r-- 1 hdfs hadoop  6377448 Jun  9 03:39 blk_1073741828
-rw-r--r-- 1 hdfs hadoop   49831 Jun  9 03:39 blk_1073741828_1004.meta
-rw-r--r-- 1 hdfs hadoop  12312 Jun  9 03:39 blk_1073741829
-rw-r--r-- 1 hdfs hadoop     107 Jun  9 03:39 blk_1073741829_1005.meta
```

# Block Replication policy

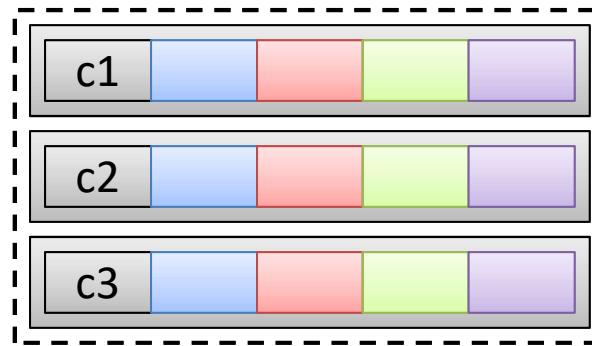
- First Replication
  - 放在上傳文件的DN。如果不是上傳者不是DN，則隨機選擇Disk不太滿且CPU不太忙的DN。
- Secondary Replication
  - 放在與第一個Replication不同的機架上的DN。
- Third Replication
  - 放在與第二個Replication相同機架的不同DN。

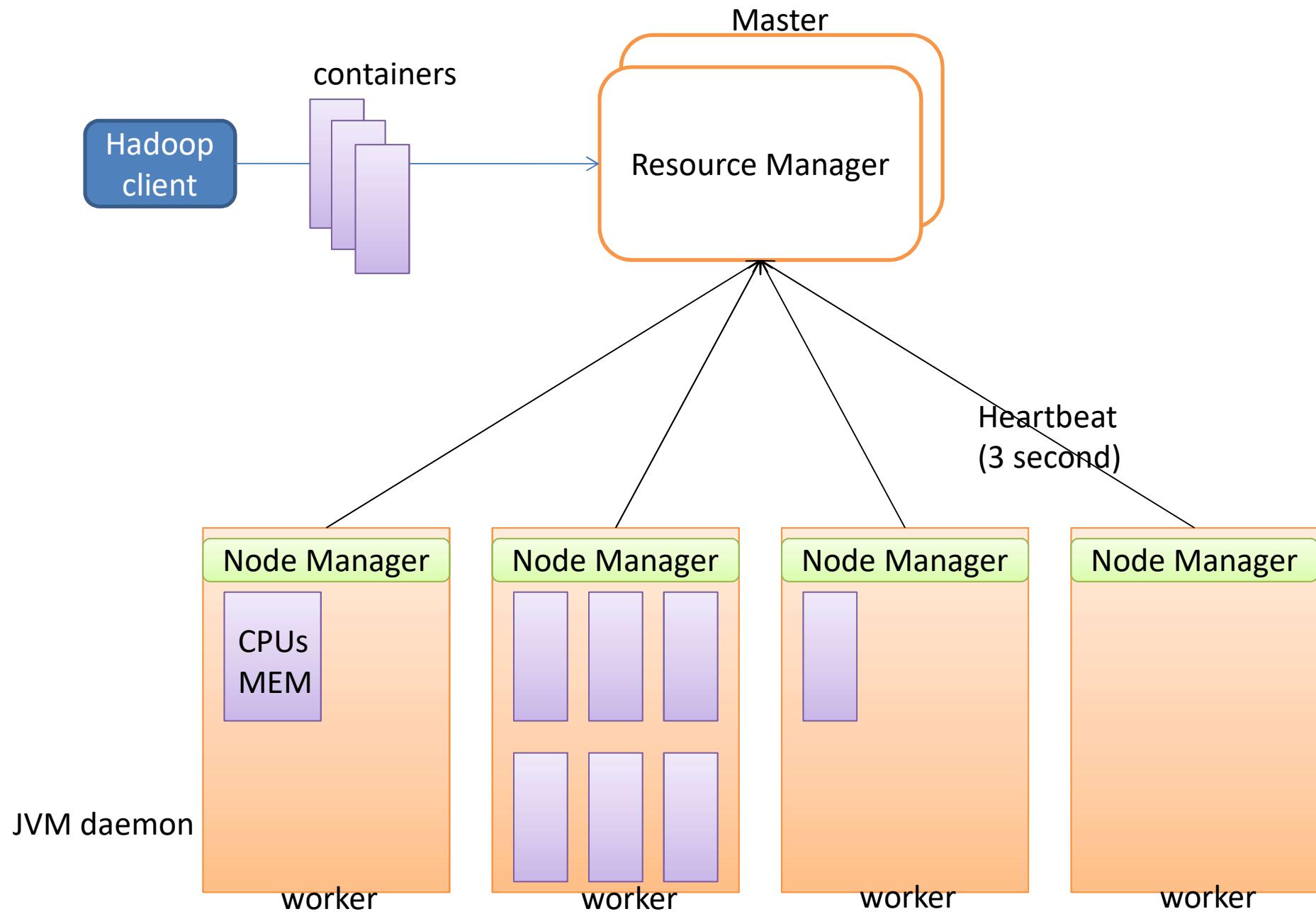
# Block size and File size

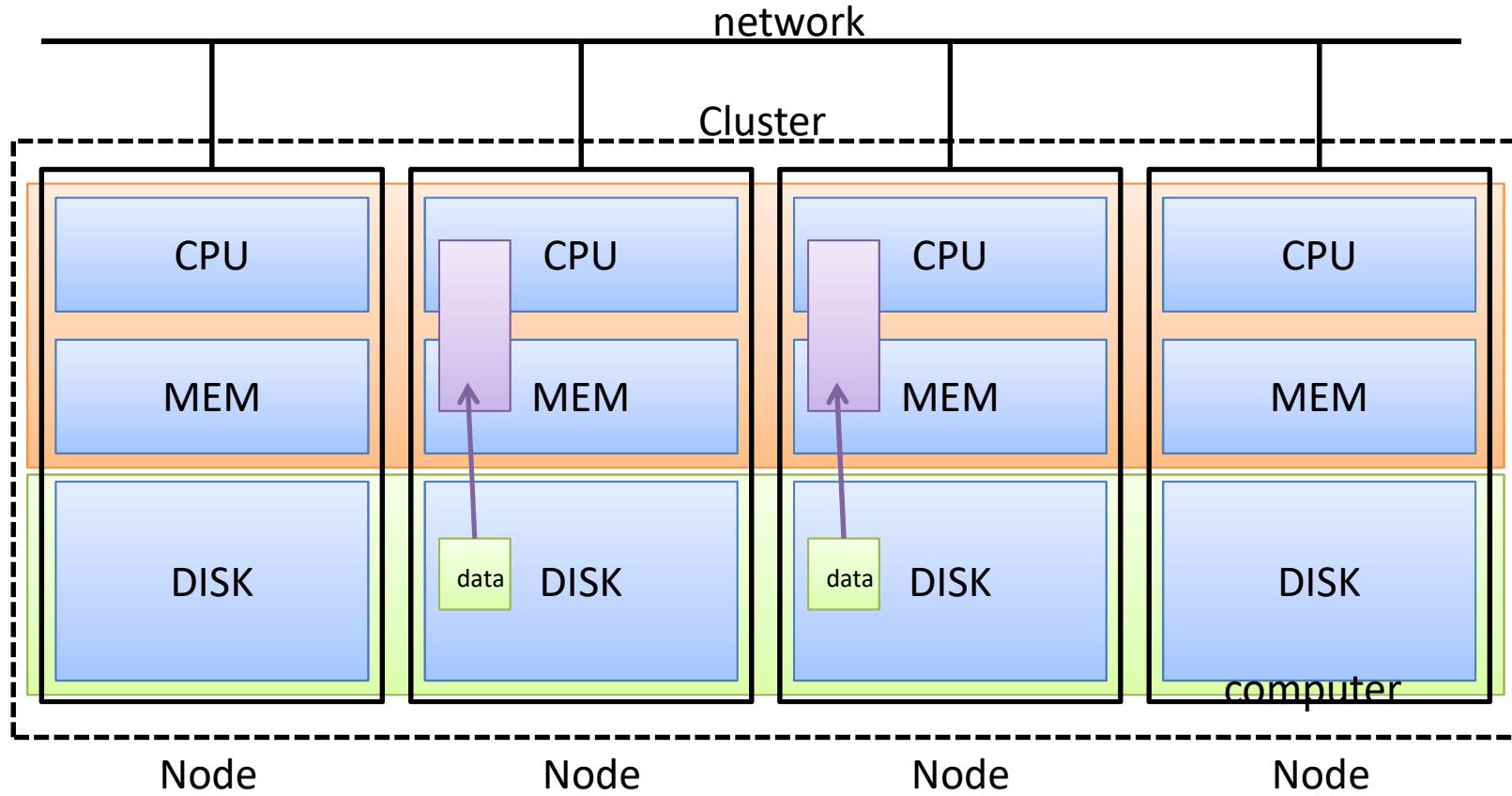
- 假設目前一個Block Size為128M，如果上傳文件小於128M，則該文件只由一個Block組成，但該Block以實際文件大小決定。
- 但每個file所佔的metadata資料空間還是一樣。所以為節省NameNode的記憶體，HDFS file應該為大量的大檔案，而不是極大量的小檔案所組成。



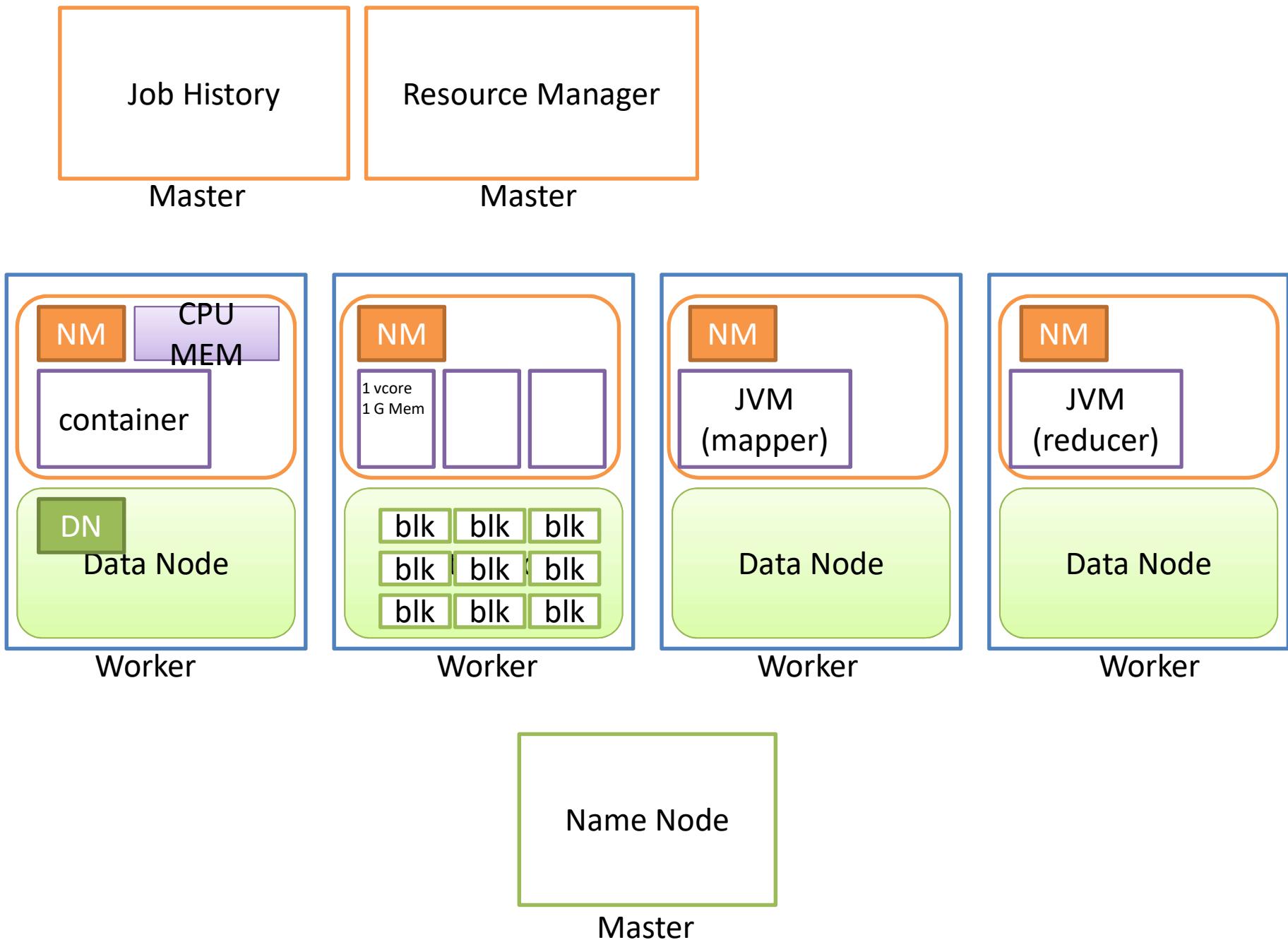
c1	c2	c3







Data Locality



## Spark Standalone

Spark Master

Spark History Server

Master

Master

Spark  
Worker

Spark  
Worker

Spark  
Worker

Spark  
Worker

Data Node

Data Node

Data Node

Data Node

Worker

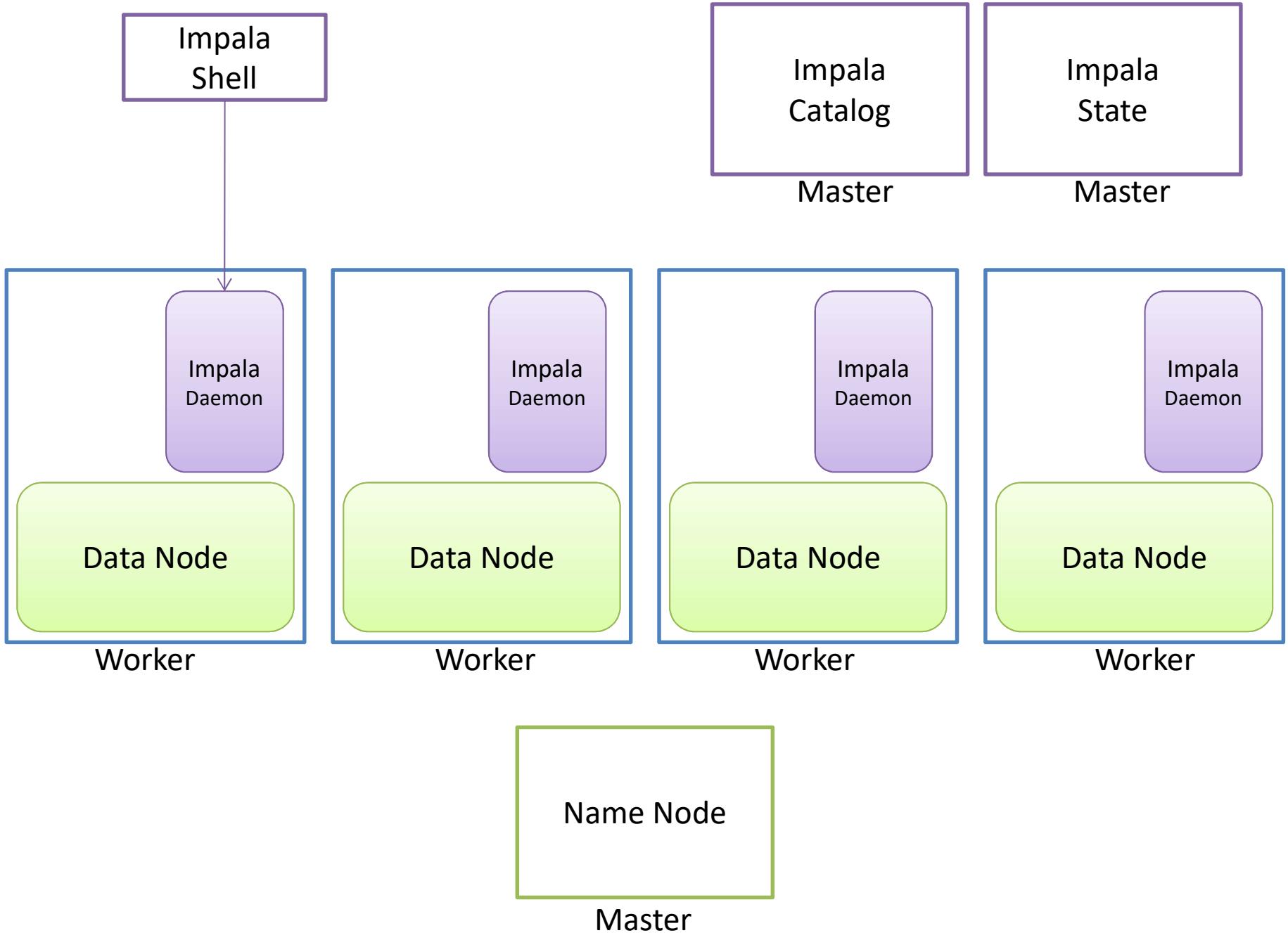
Worker

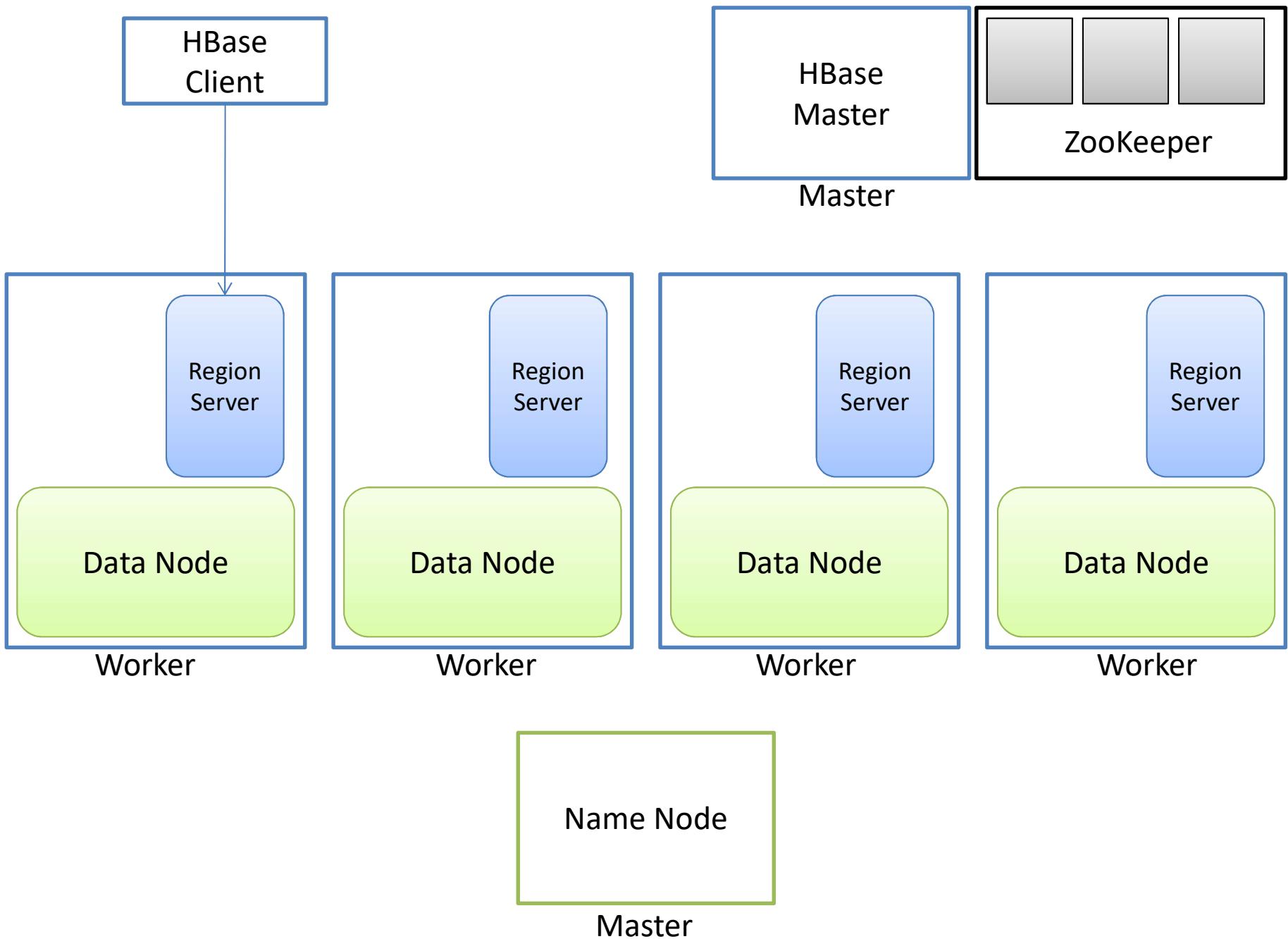
Worker

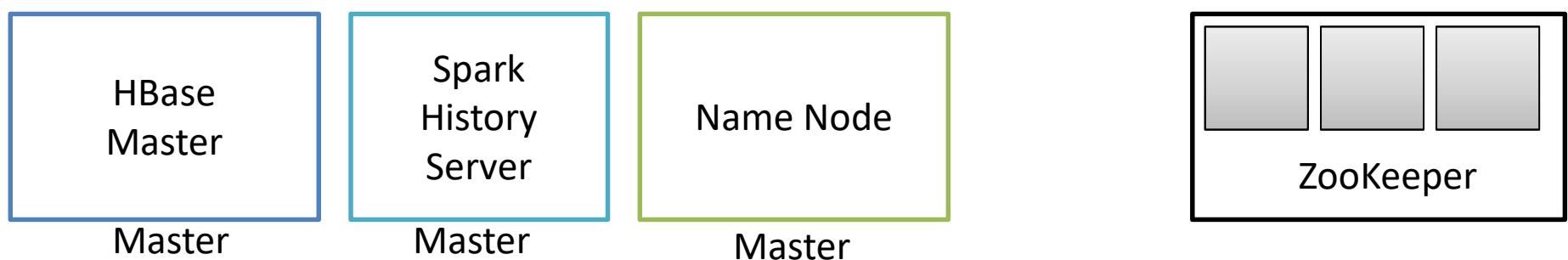
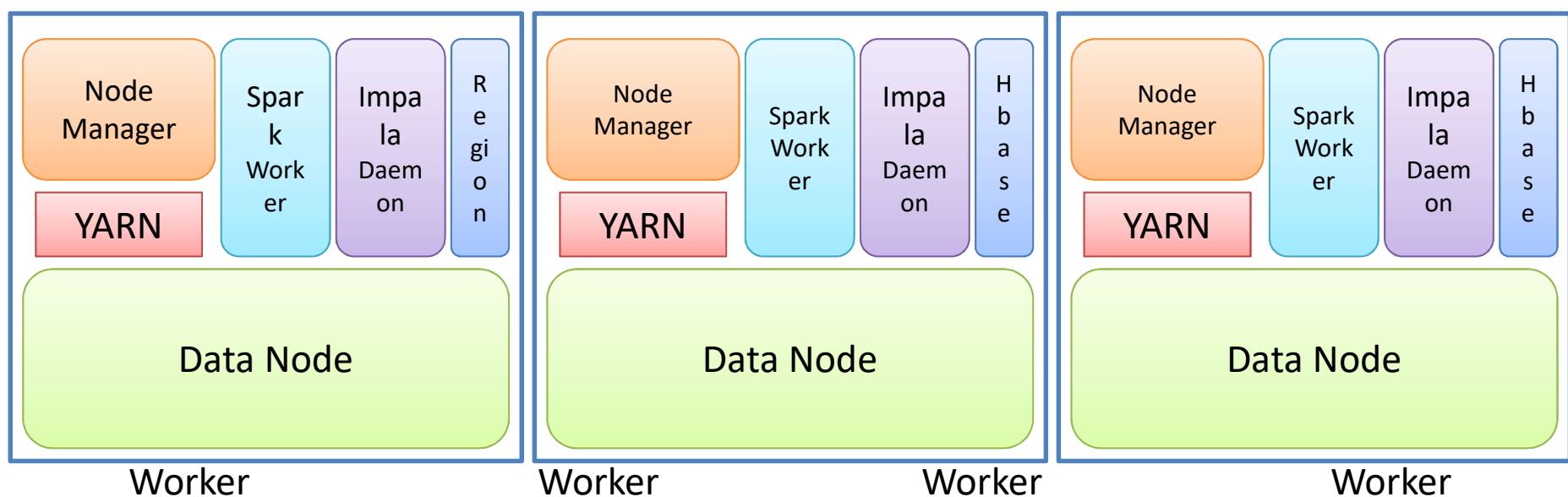
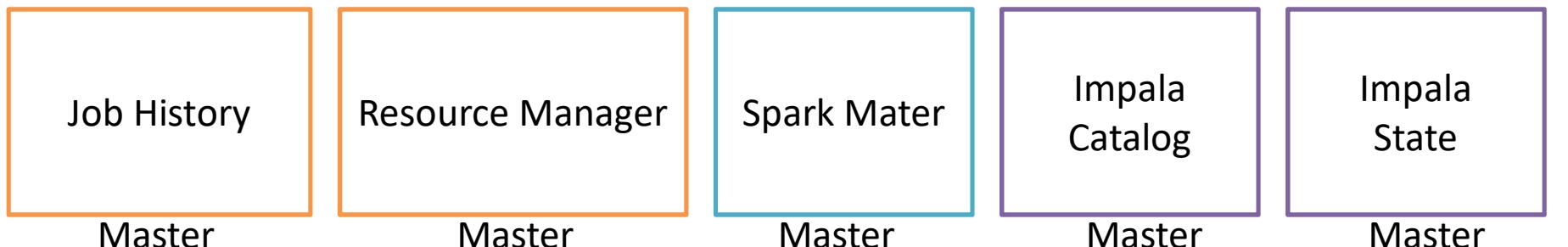
Worker

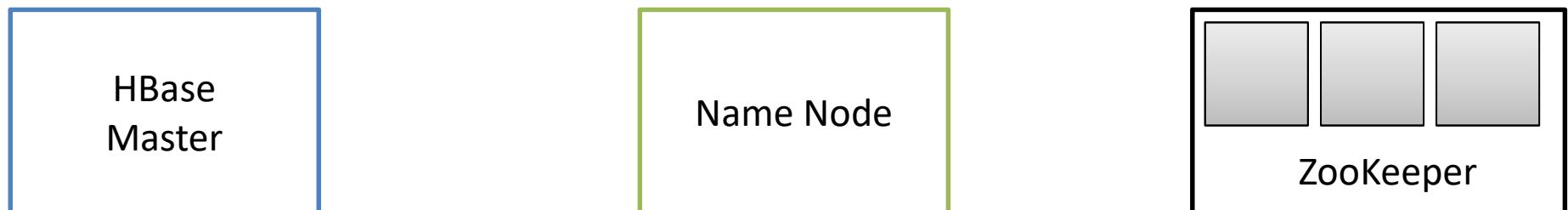
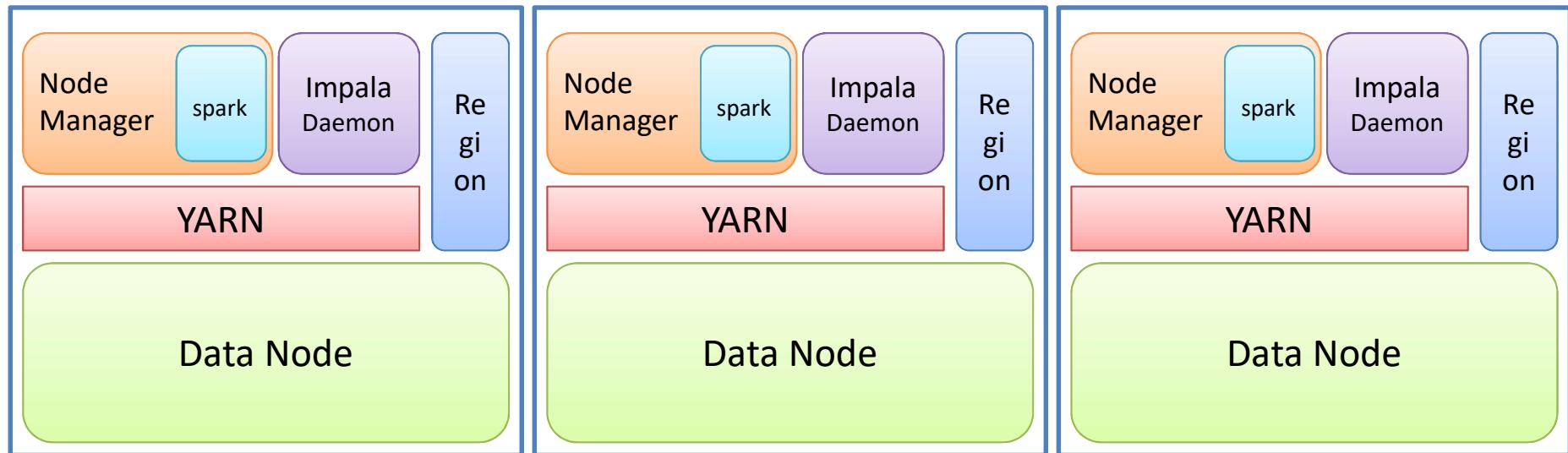
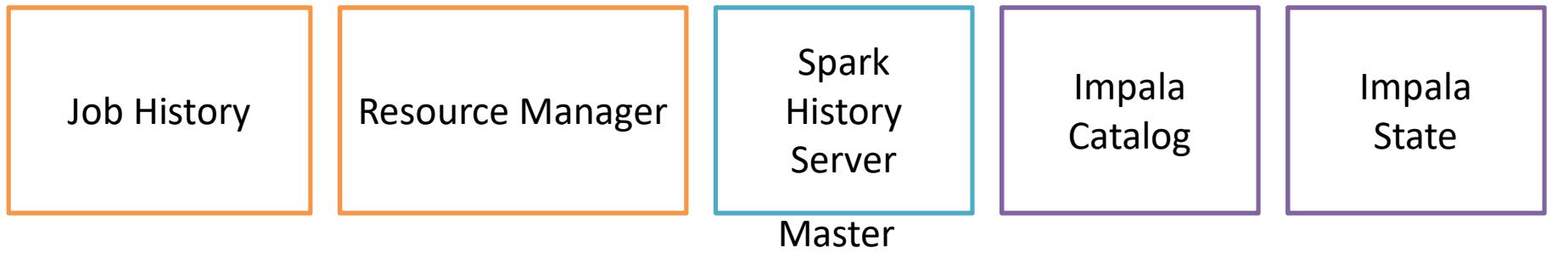
Name Node

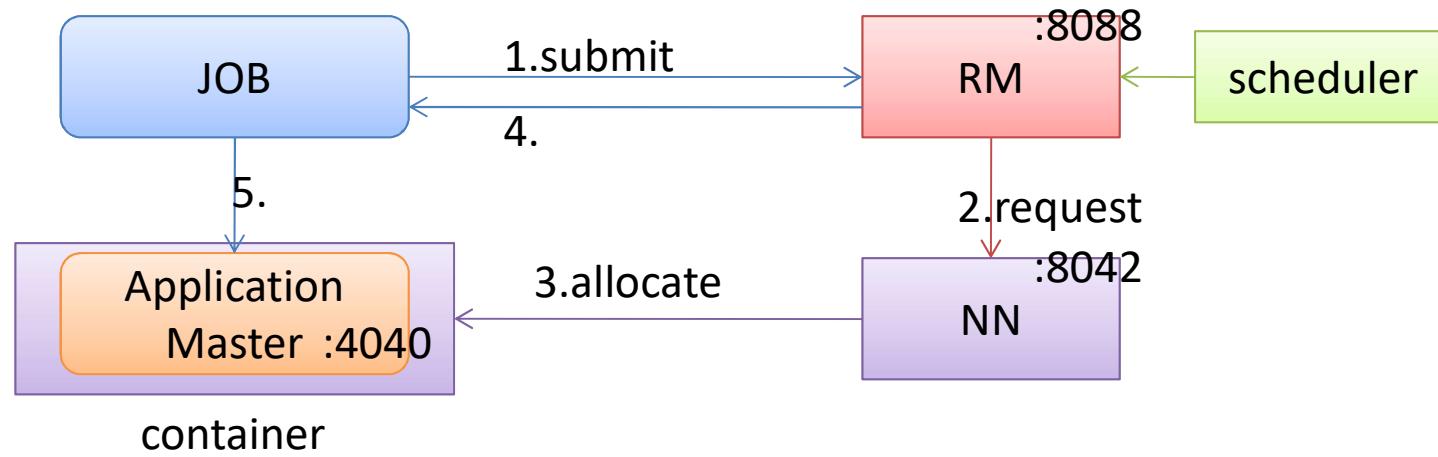
Master

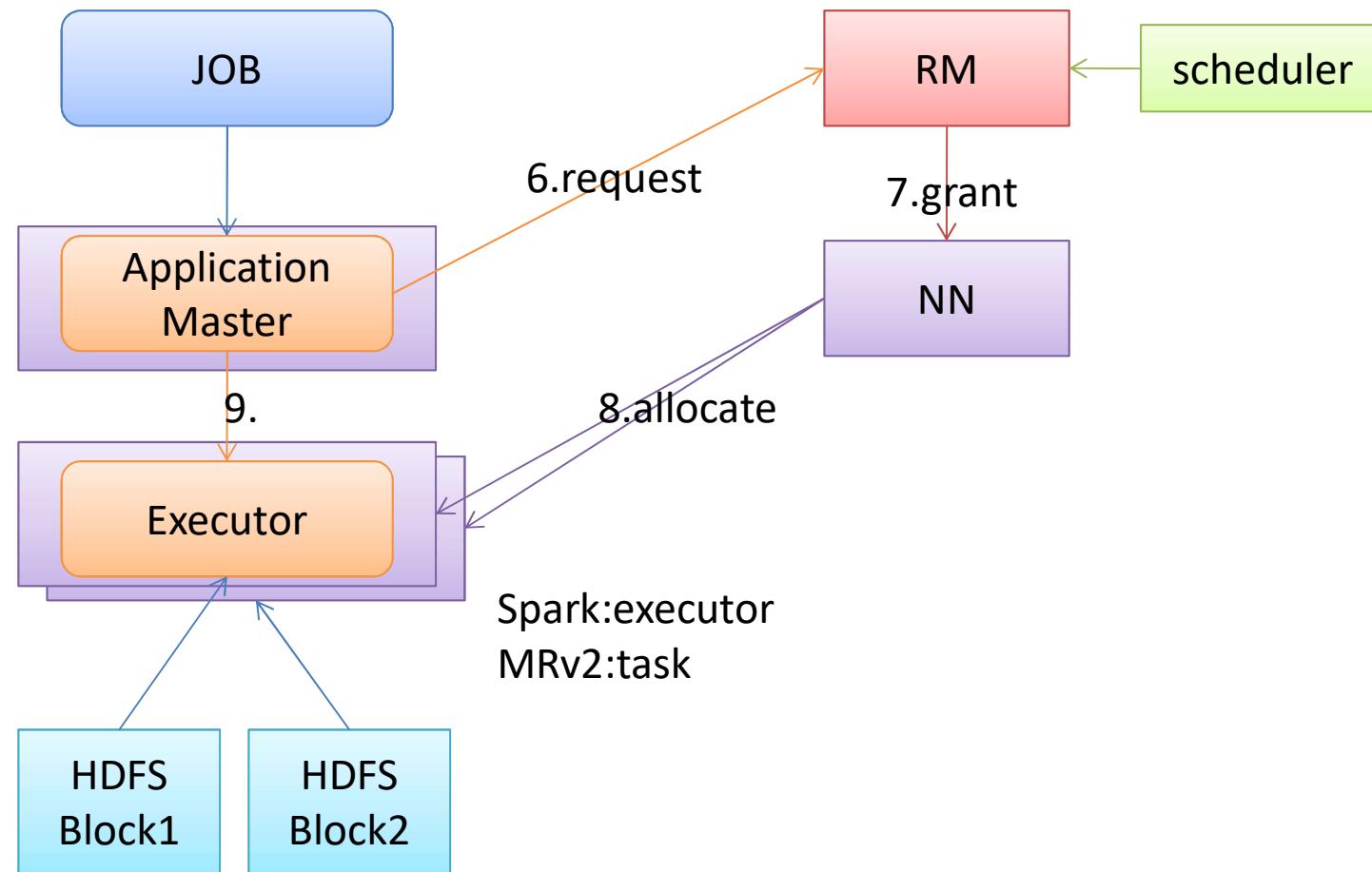


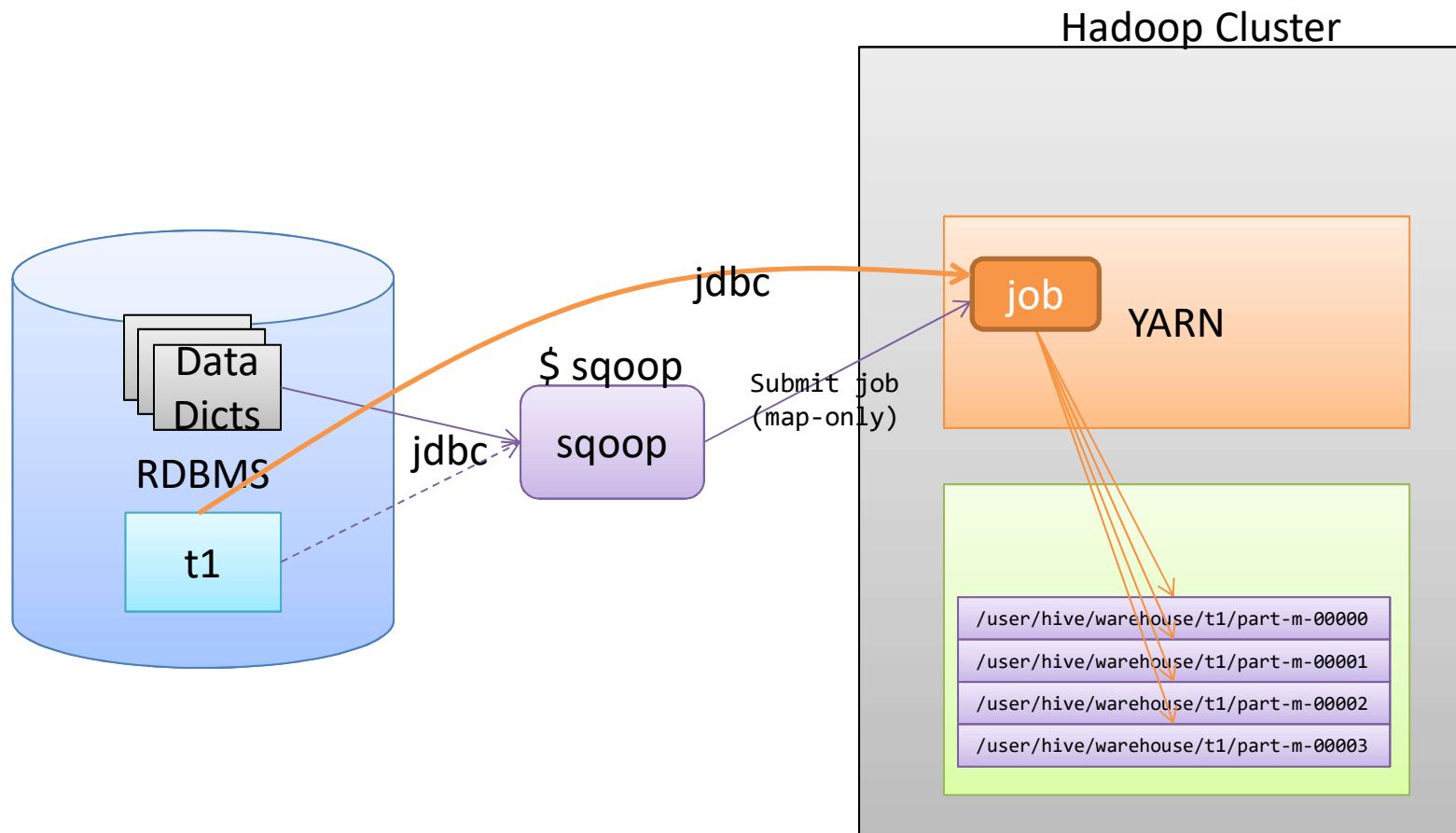


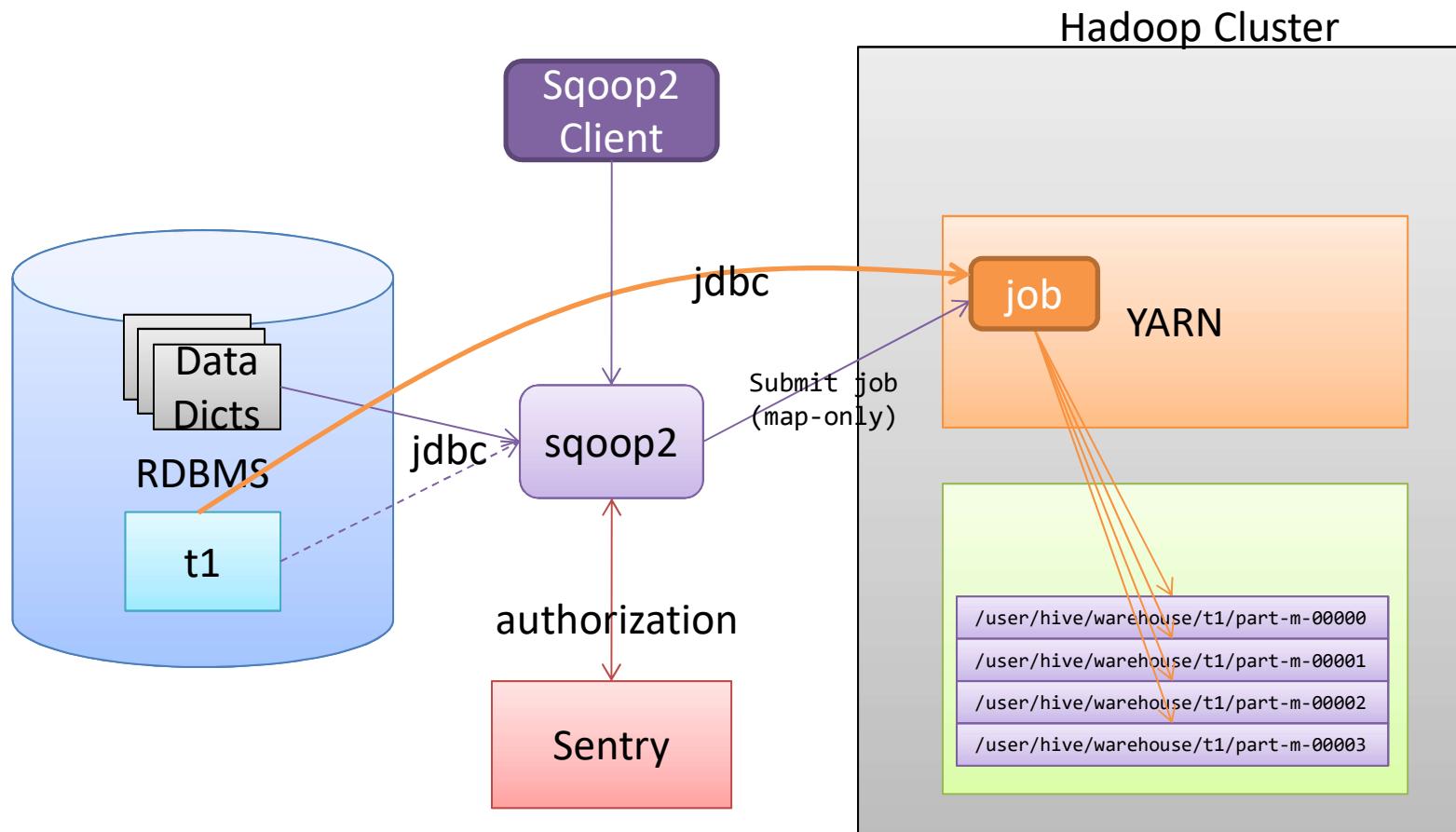












```
$ sqoop import \
--connect jdbc:mysql://Localhost/Loudacre \
--username dbuser --password pw \
--table accounts
--where "state='CA'" \
--columns "id,first_name,last_name,state" \
```

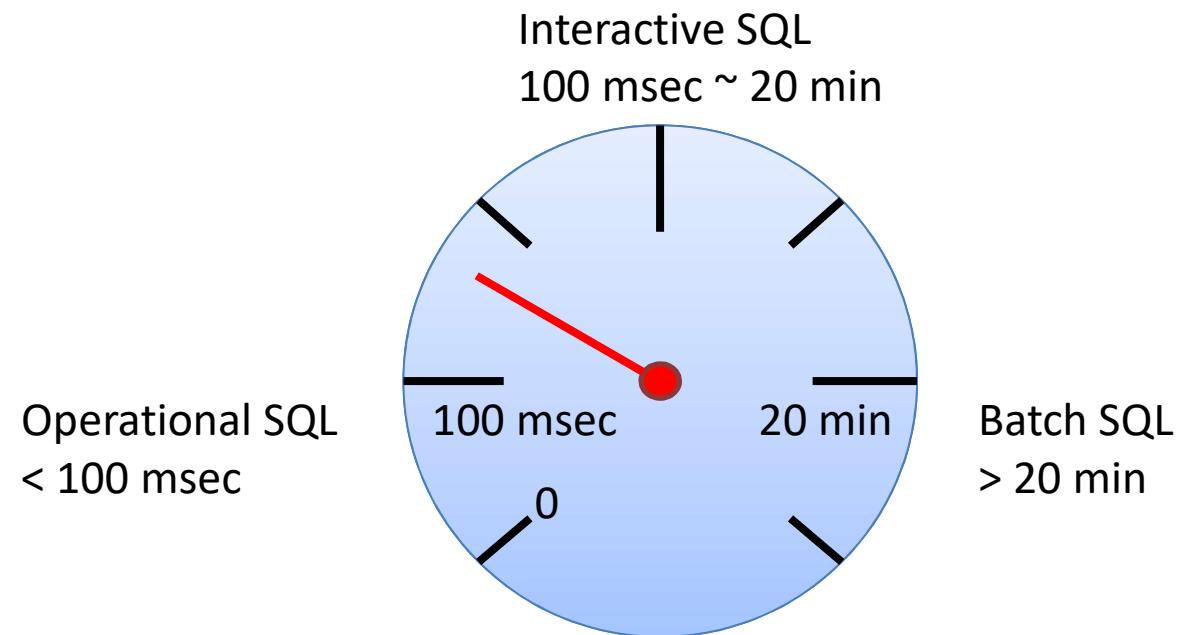
```
SELECT id,first_name,last_name,state
FROM accounts
WHERE state= 'CA '
```

```
$ sqoop import \
--connect jdbc:mysql://Localhost/Loudacre \
--username dbuser --password pw \
--target-dir /data/loudacre/payable
--split-by account.id
--query " SELECT
id,first_name,last_name,state
FROM accounts
WHERE \$CONDITION AND state= 'CA '"
```

```
$ sqoop import \
--connect jdbc:mysql://dbhost/Loudacre \
--username dbuser --password pw \
--target-dir /data/loudacre/payable \
--split-by accounts.id \
--query 'SELECT accounts.id, first_name, last_name,
bill_amount FROM accounts JOIN invoices ON
(accounts.id = invoices.cust_id) WHERE $CONDITIONS AND
bill_amount >= 40'
```

```
$ sqoop import \
--connect jdbc:mysql://dbhost/Loudacre \
--username dbuser --password pw \
--target-dir /data/loudacre/payable \
--split-by accounts.id \
--query "SELECT accounts.id, first_name, last_name,
bill_amount FROM accounts JOIN invoices ON
(accounts.id = invoices.cust_id) WHERE \$CONDITIONS AND
bill_amount >= 40"
```

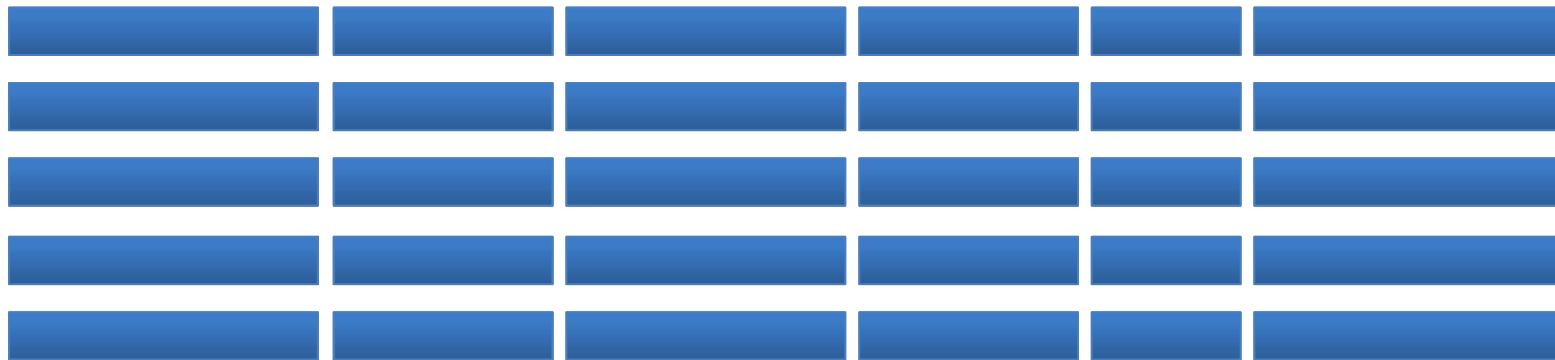
# SQL on Hadoop



# SQL on Hadoop

- Batch SQL
  - Hive
    - Based on MapReduce
    - Based on Tez(Hortonworks)
- Interactive SQL
  - Spark SQL(Hive on Spark)
  - Impala(Cloudera)(Google Dremel)
  - Apache Drill(Google Dremel)
  - Presto(Facebook) (Google Dremel)
- OLAP
  - Pinot(LinkedIn)
  - Kylin(Ebay)
- Operational SQL
  - Apache Phoenix(SQL on HBase)

# Schema on Write

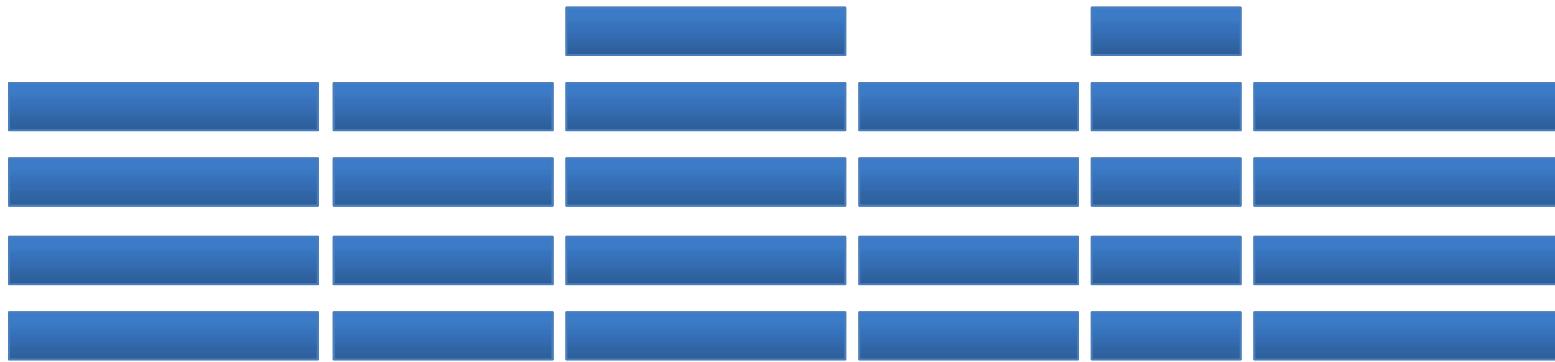


2.Load data into table

Empid	Ename	Salary	Hire_date

1.Create table

# Schema on Write



2.Load data into table

Empid	Ename	Salary	Hire_date

1.Create table

# Schema on Write

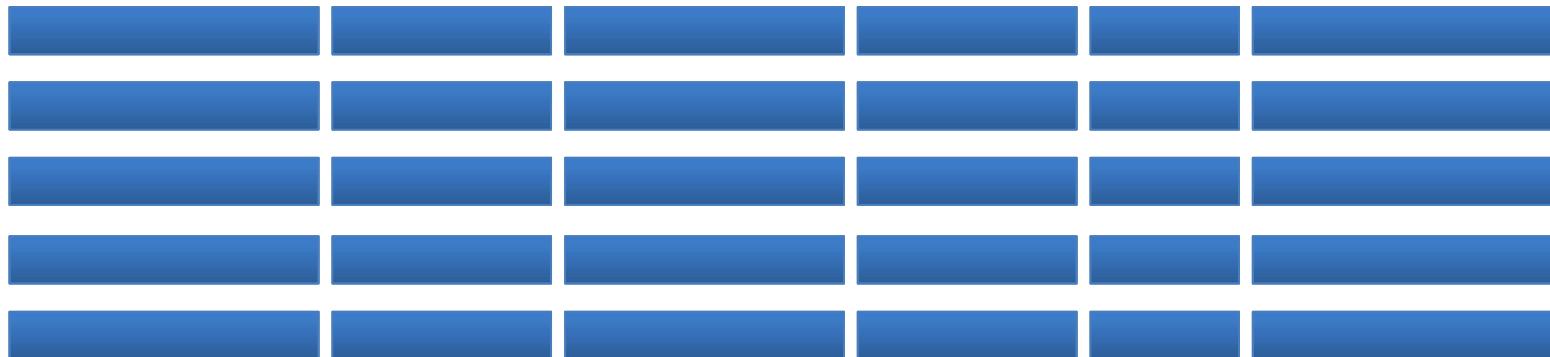


2.Load data into table

Empid	Ename	Salary	Hire_date

1.Create table

# Schema on Read

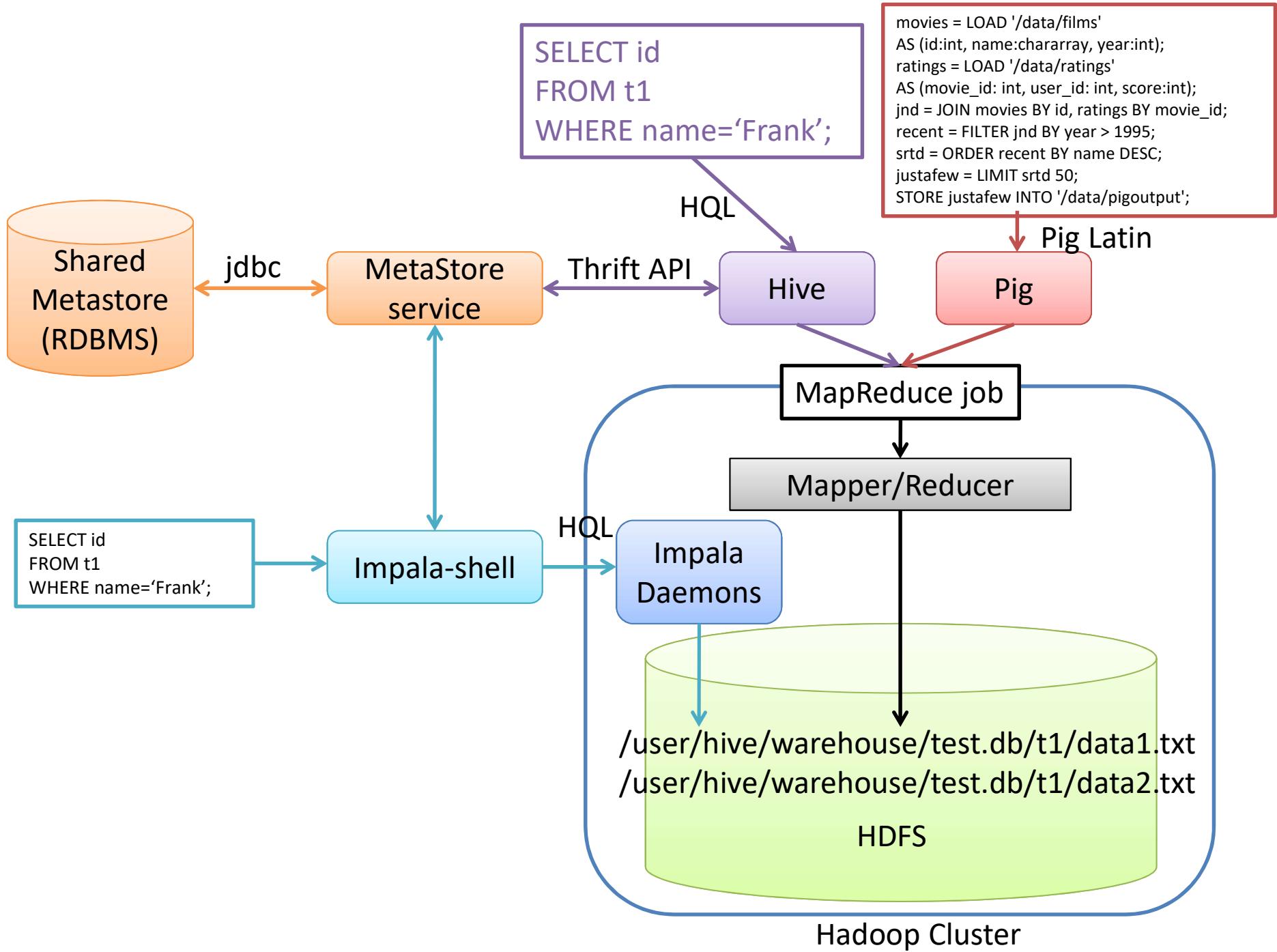


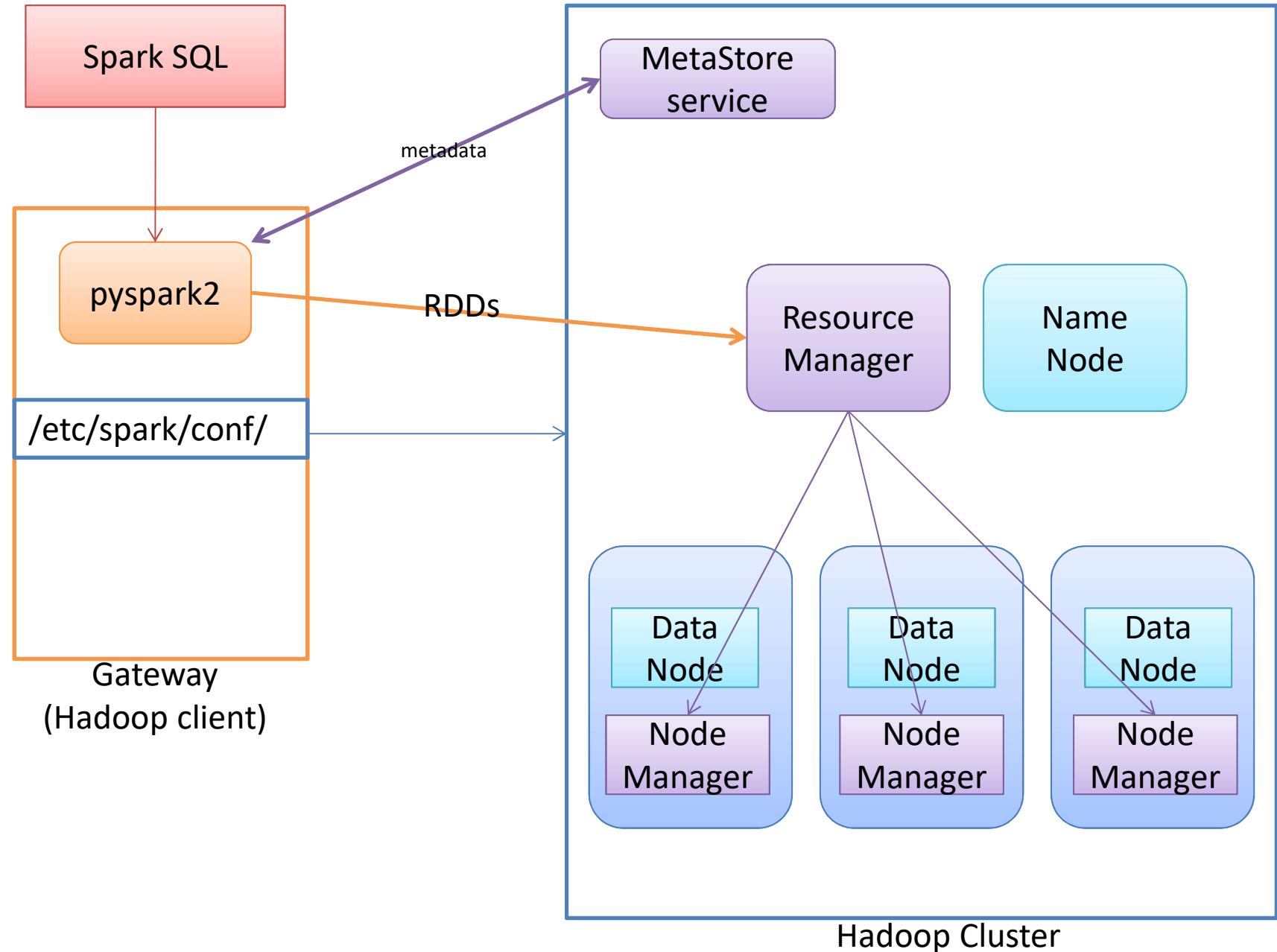
1.Load data into HDFS

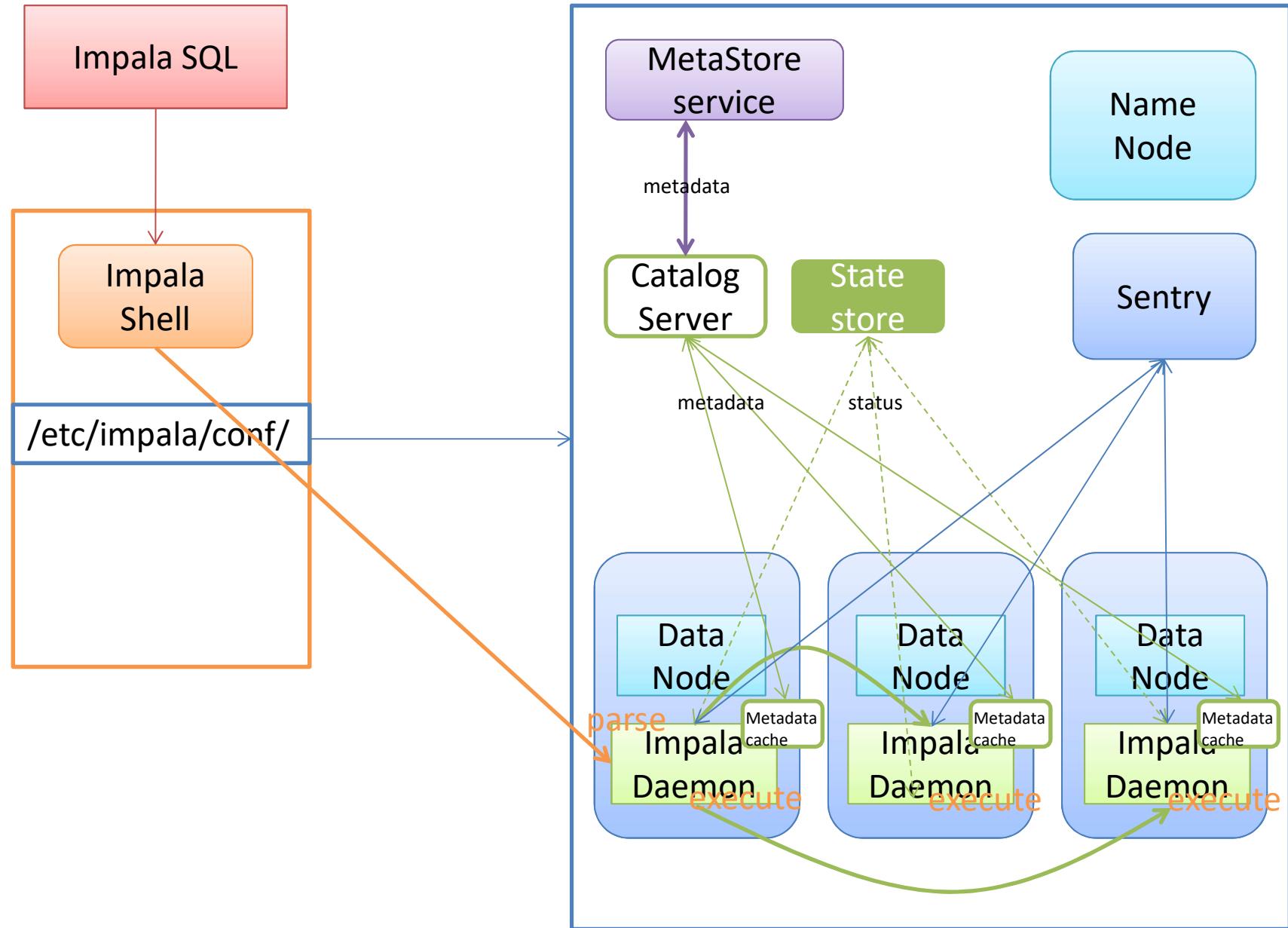
# Schema on Read

Empid	Ename	Job	Salary	Dept	Hire_date

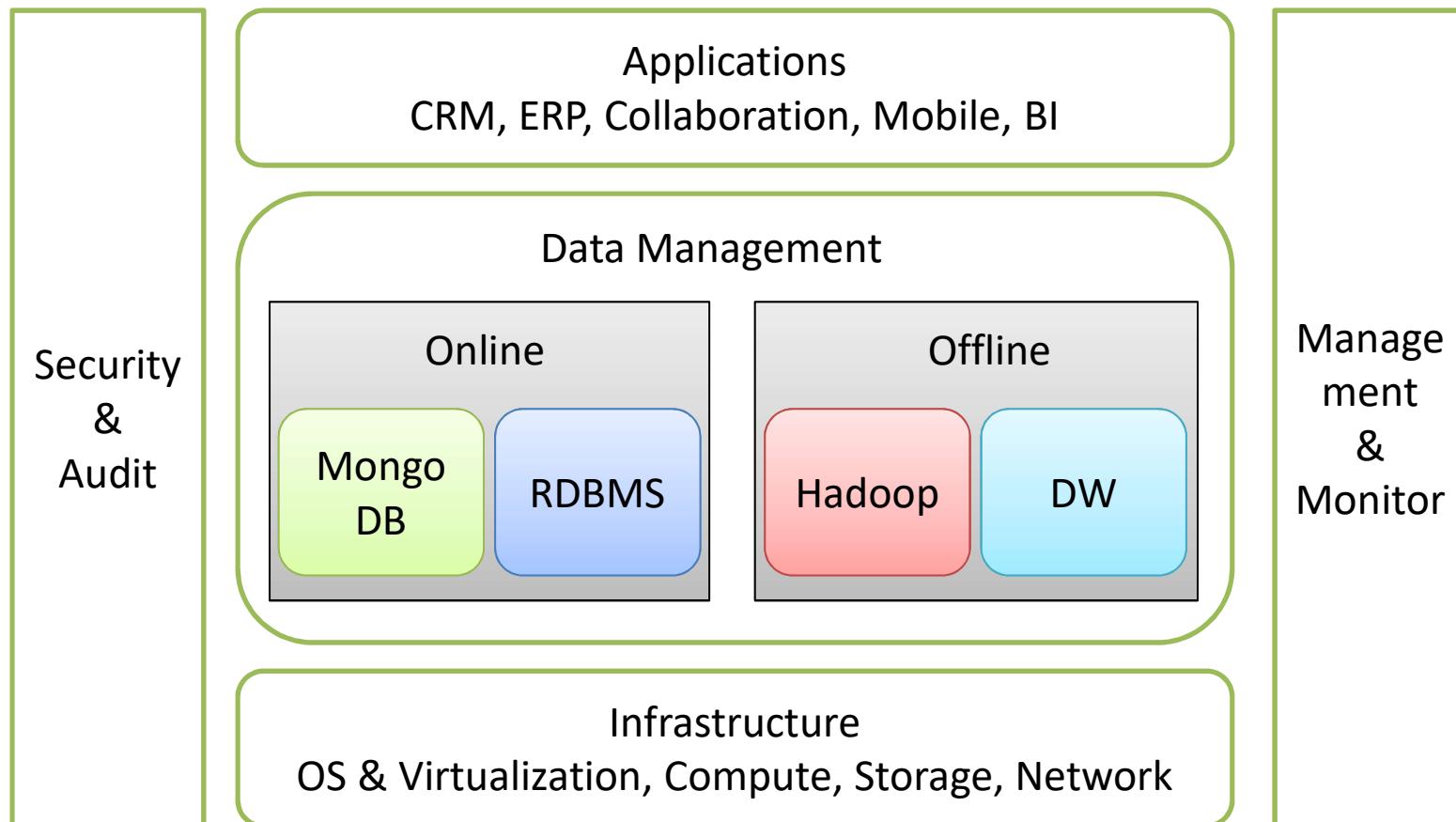
2. Create Hive table mapping HDFS file







# Enterprise Big Data Stack



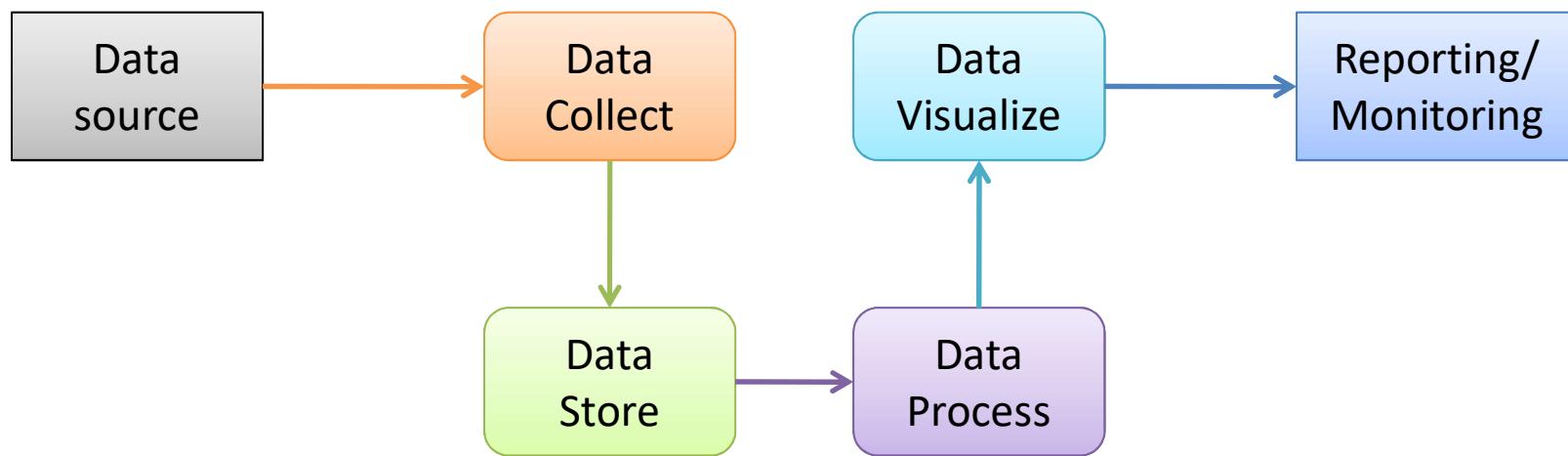
# Document是類JSON格式

```
{  
    "name"      : "Frank Shen",  
    "age"       : 45,  
    "address"   : [  
        {"city"     : "Taipei City",  
         "country" : "Taiwan"},  
        {"city"     : "Taoyuan City",  
         "country" : "Taiwan"}  
    ]  
}
```

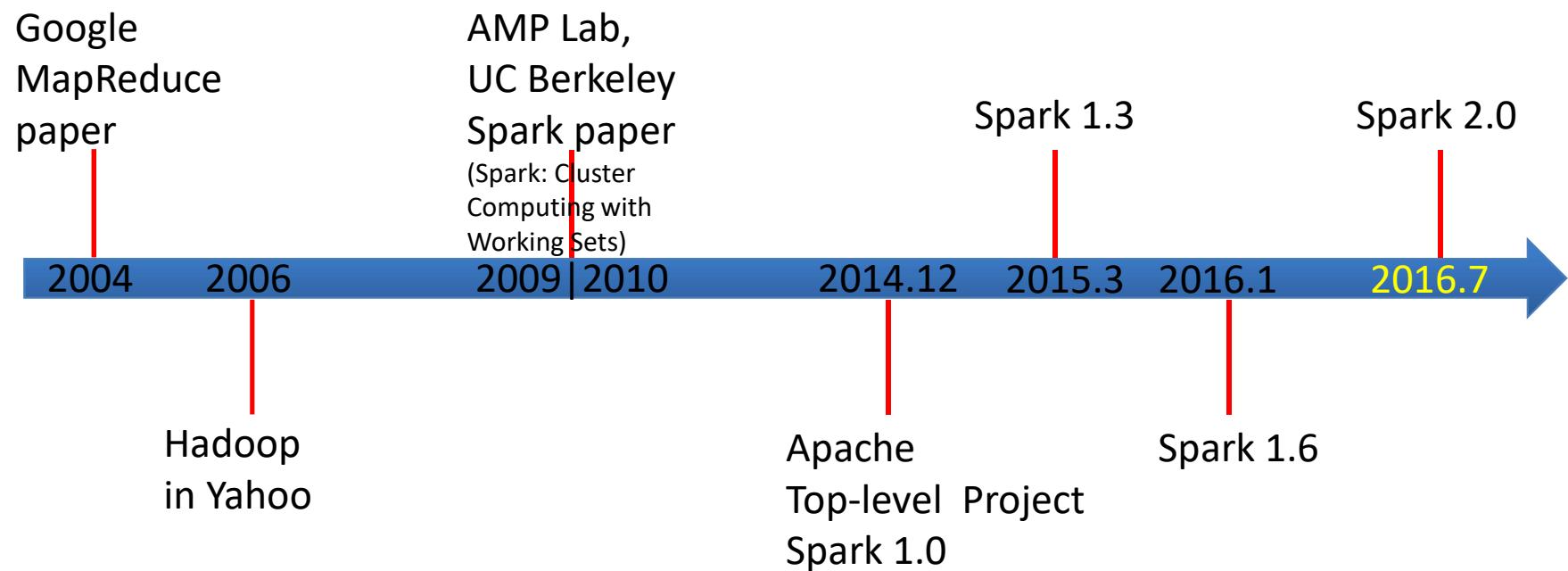
```
{  
    "name"      : "Wilson Chang",  
    "age"       : 38  
}
```

# Document 結構





# Brief History of Spark



# What and Why Spark

- Fast and expressive cluster computing system interoperable with Apache Hadoop
- Improves efficiency through:
  - In-memory computing primitives
  - General computation graphs
- Improves usability through:
  - Rich APIs in Scala, Java, Python
  - Interactive shell 

Up to 100x faster  
(2 ~ 10x on Disks)

Often 5x less on code

```

public class WordCount {

    public static class TokenizerMapper
        extends Mapper<Object, Text, Text, IntWritable>{

        private final static IntWritable one = new IntWritable(1);
        private Text word = new Text();

        public void map(Object key, Text value, Context context
                        ) throws IOException, InterruptedException {
            StringTokenizer itr = new StringTokenizer(value.toString());
            while (itr.hasMoreTokens()) {
                word.set(itr.nextToken());
                context.write(word, one);
            }
        }

        public static class IntSumReducer
            extends Reducer<Text,IntWritable,Text,IntWritable> {
            private IntWritable result = new IntWritable();

            public void reduce(Text key, Iterable<IntWritable> values,
                               Context context
                               ) throws IOException, InterruptedException
            {
                int sum = 0;
                for (IntWritable val : values) {
                    sum += val.get();
                }
                result.set(sum);
                context.write(key, result);
            }

            public static void main(String[] args) throws Exception {
                Configuration conf = new Configuration();
                Job job = Job.getInstance(conf, "word count");
                job.setJarByClass(WordCount.class);
                job.setMapperClass(TokenizerMapper.class);
                job.setCombinerClass(IntSumReducer.class);
                job.setReducerClass(IntSumReducer.class);
                job.setOutputKeyClass(Text.class);
                job.setOutputValueClass(IntWritable.class);
                FileInputFormat.addInputPath(job, new Path(args[0]));
                FileOutputFormat.setOutputPath(job, new Path(args[1]));
                System.exit(job.waitForCompletion(true) ? 0 : 1);
            }
        }
    }
}

```

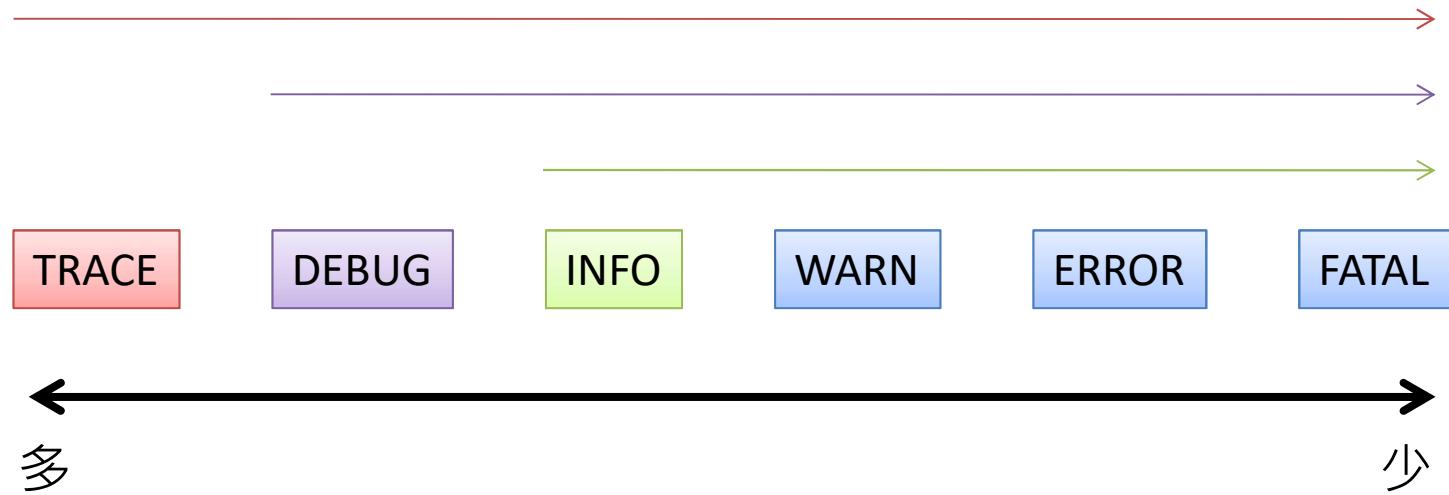
Java

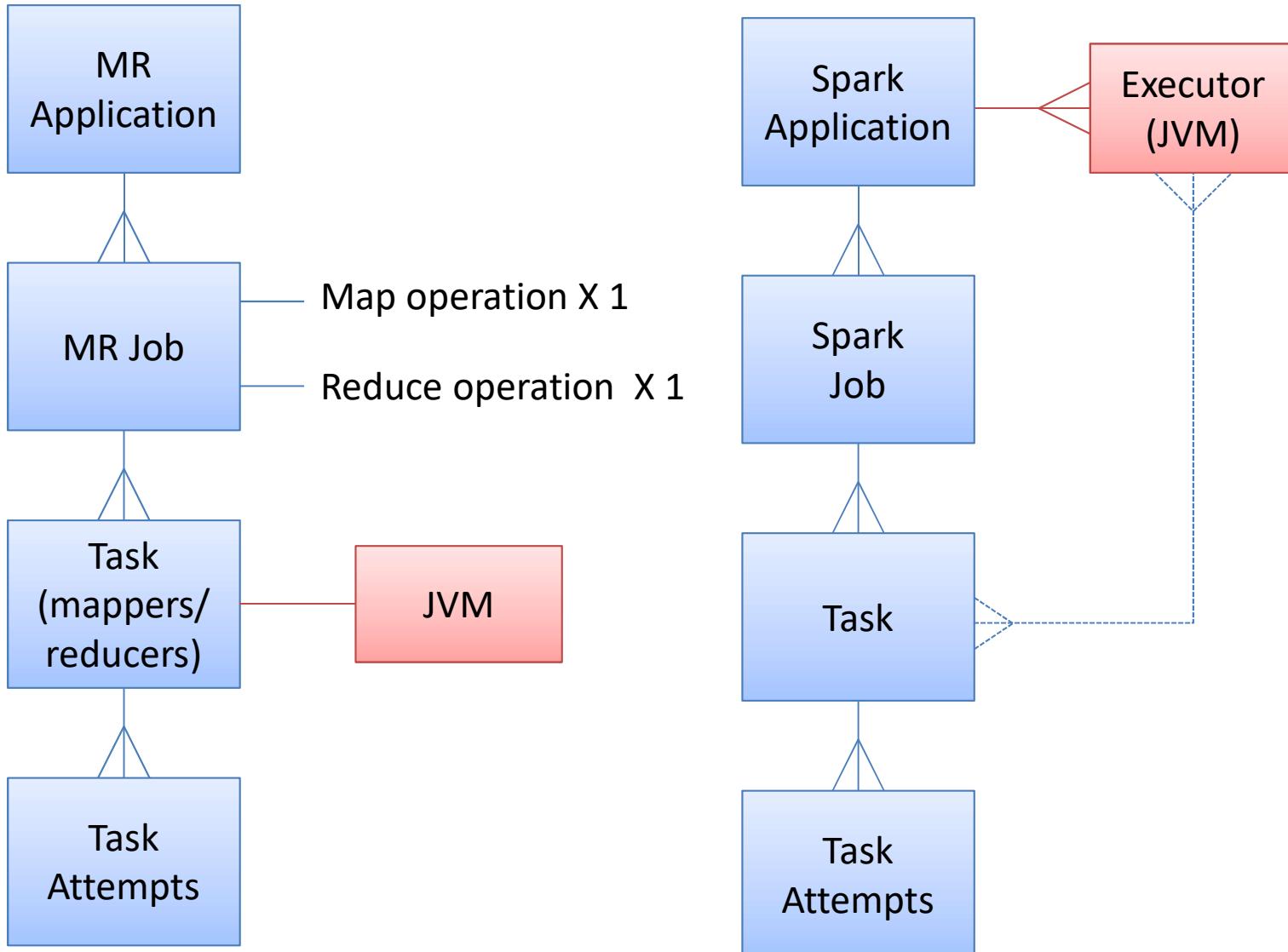
Spark

```

text_file = sc.textFile("hdfs://...")
counts = text_file.flatMap(lambda line: line.split(" "))
                  .map(lambda word: (word, 1))
                  .reduceByKey(lambda a, b: a + b)
counts.saveAsTextFile("hdfs://...")

```





# RDD to DataFrame

---

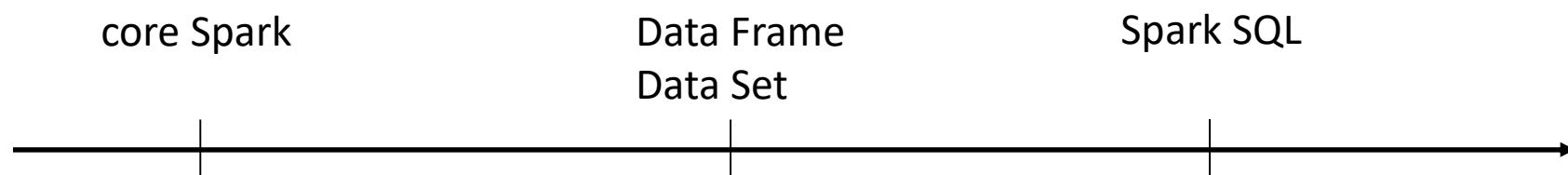
```
pdata.map(lambda x: (x.dept, [x.age, 1])) \  
    .reduceByKey(lambda x, y: [x[0] + y[0], x[1] + y[1]]) \  
    .map(lambda x: [x[0], x[1][0] / x[1][1]]) \  
    .collect()
```

```
data.groupBy("dept").avg("age")
```

# SQL integrate with Spark

```
sqlcontext = HiveContext(sc)
results = sqlcontext.sql("SELECT country, count(*) FROM people
GROUP BY country")
names = results.map(lambda p: p.name)
```

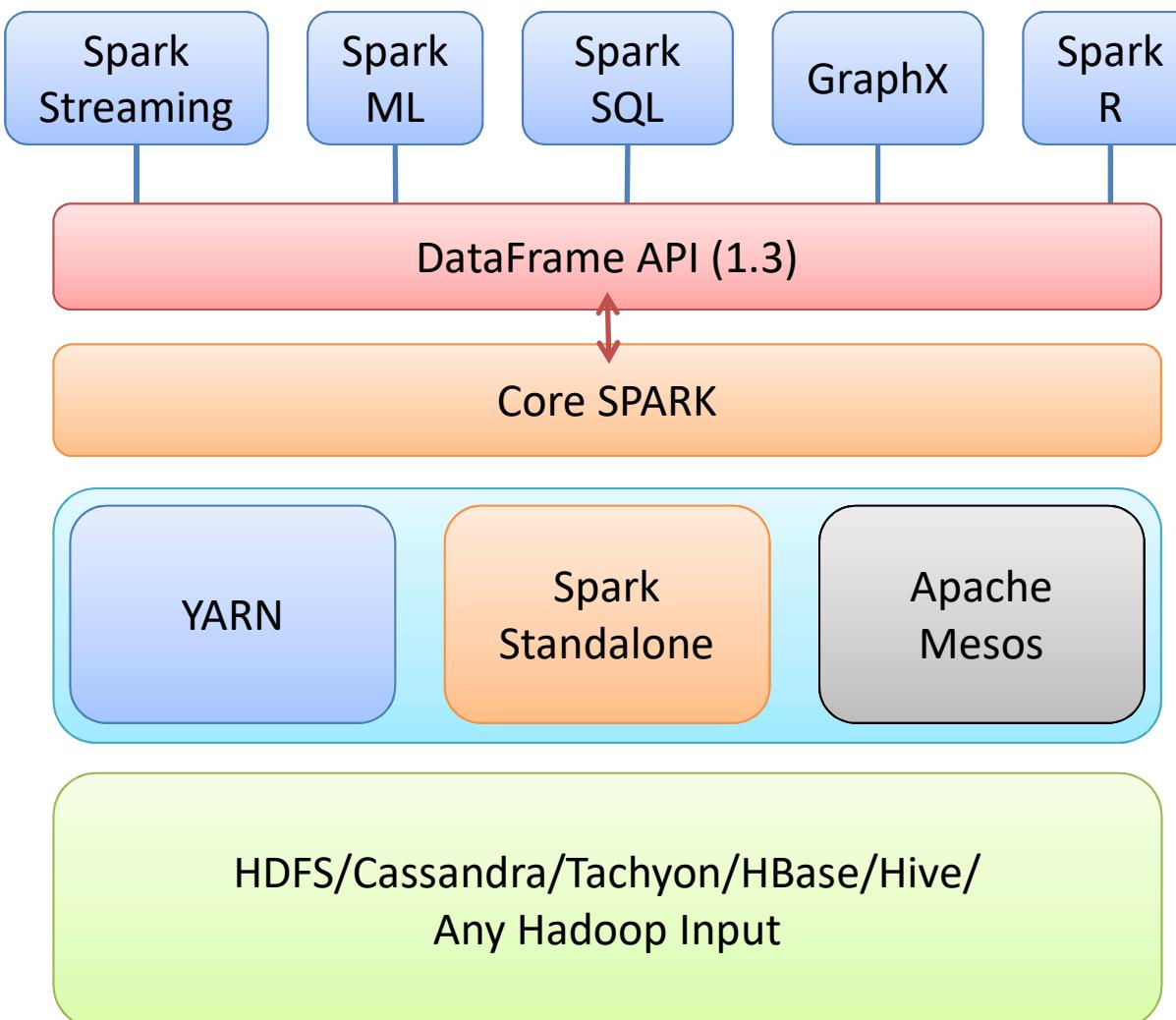
# Spark Coding 演進



# What is DataFrame

- Distributed collection of data grouped into named columns
  - Inspired by data frames in R and Python Pandas
  - Conceptually equivalent to a table in RDBMS
- Supported Data Formats and Sources
  - Structured files by Parquet 、 JSON 、 S3 、 HDFS 、 Cassandra
  - Tables in Hive 、 RDBMS
- Available in Python 、 Scala 、 Java and R

# Dataframe is future of Spark(1.3)

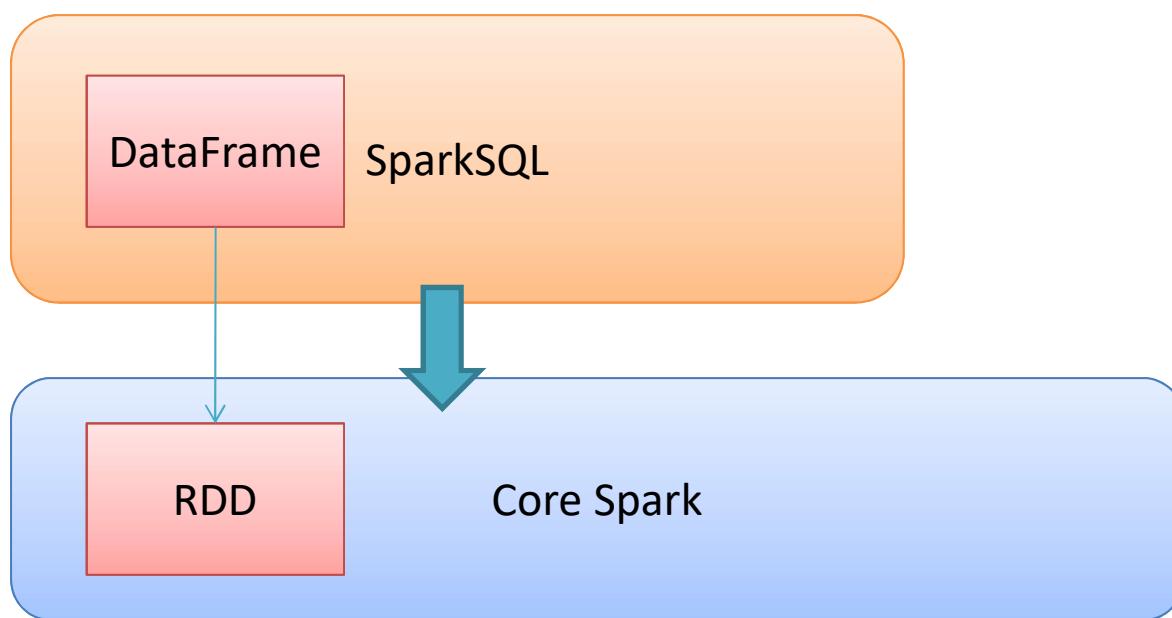


# RDD to DataFrame

---

```
pdata.map(lambda x: (x.dept, [x.age, 1])) \  
    .reduceByKey(lambda x, y: [x[0] + y[0], x[1] + y[1]]) \  
    .map(lambda x: [x[0], x[1][0] / x[1][1]]) \  
    .collect()
```

```
data.groupBy("dept").avg("age")
```



~~SELECT~~

~~FROM~~

~~WHERE~~

~~GROUP BY~~

~~HAVING~~

~~ORDER BY~~

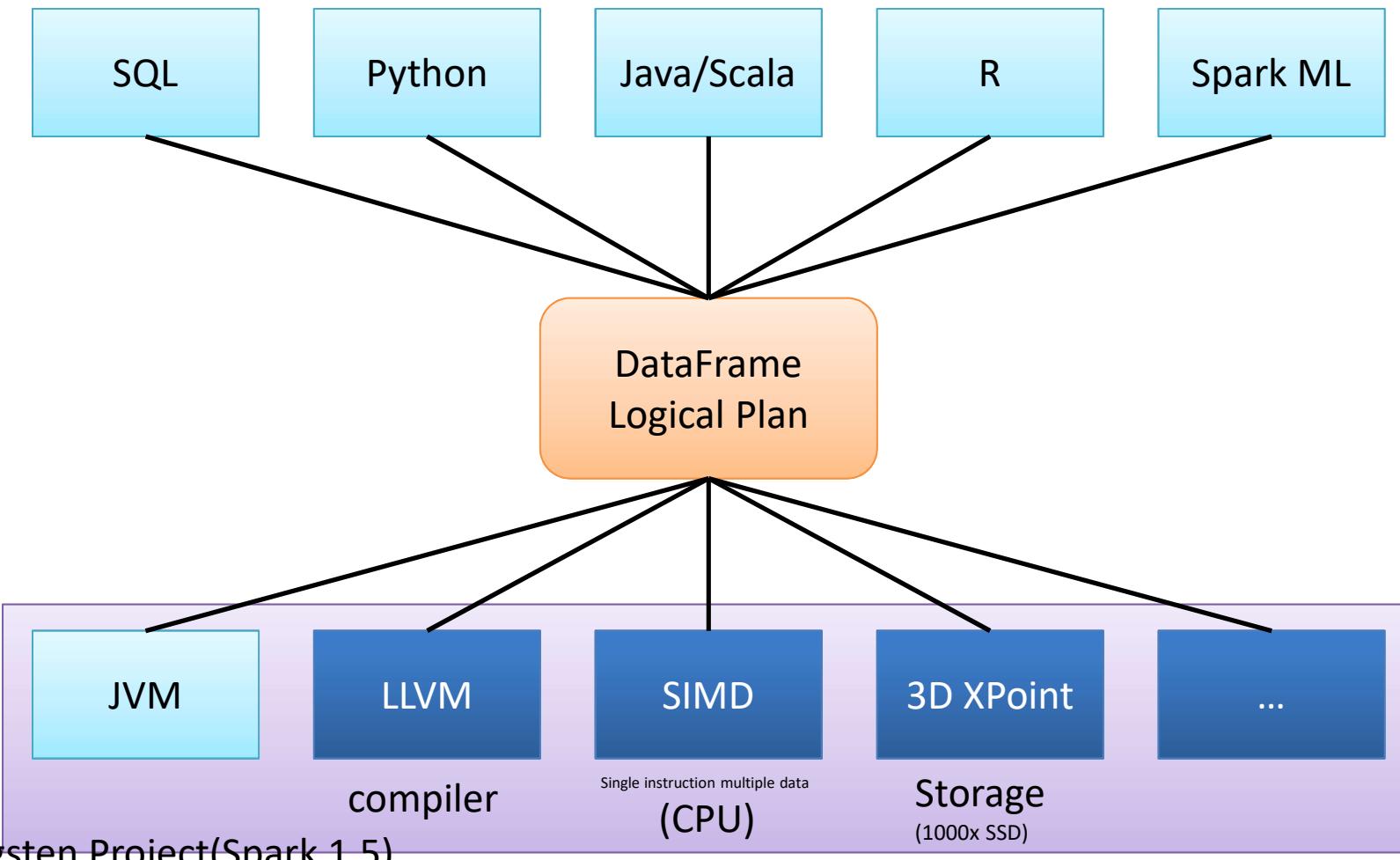
~~OFFSET/FETCH~~

# SQL integrate with Spark

```
sqlcontext = HiveContext(sc)
results = sqlcontext.sql("SELECT * FROM customers")
names = results.map(lambda c: c.name)
```



# Unified API. One Engine. Automatically Optimized



# Hardware Trends

2010

2015



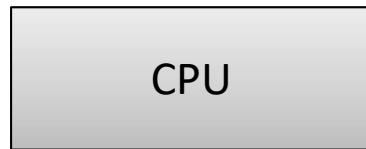
50+ MB/s  
(HDD)

500+ MB/s  
(SSD)



1Gbps

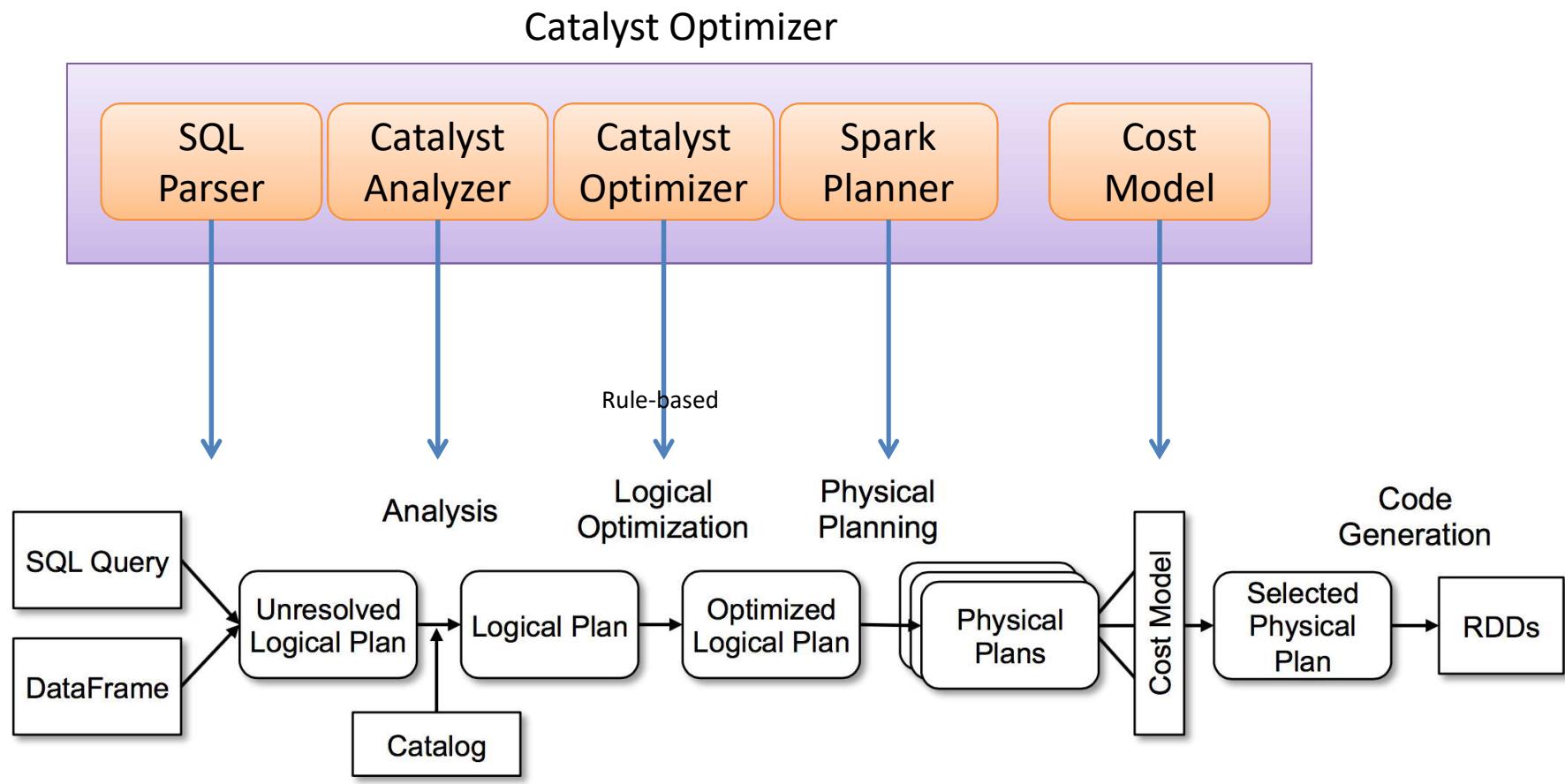
10Gbps



~3GHz

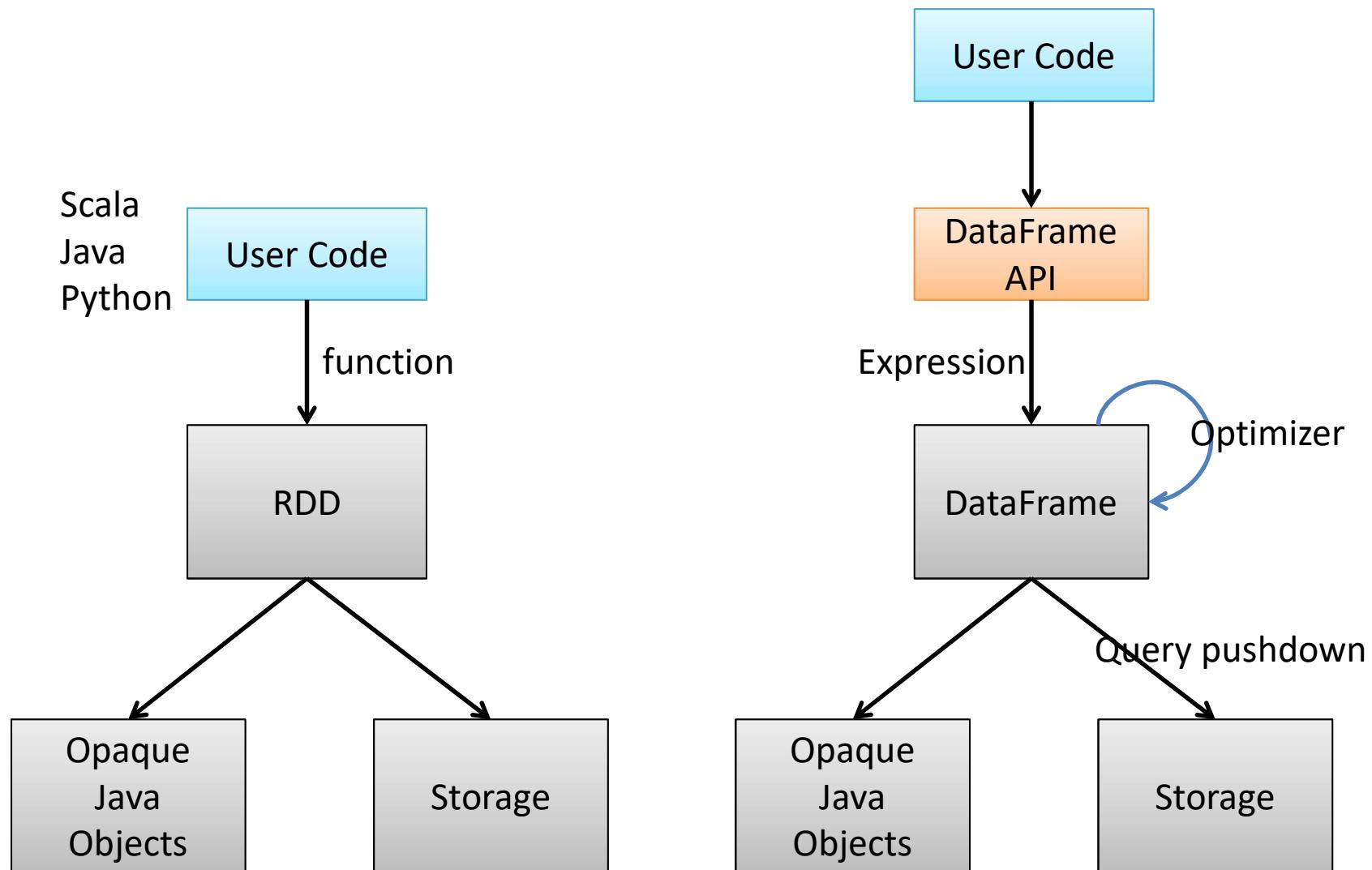
~3GHz





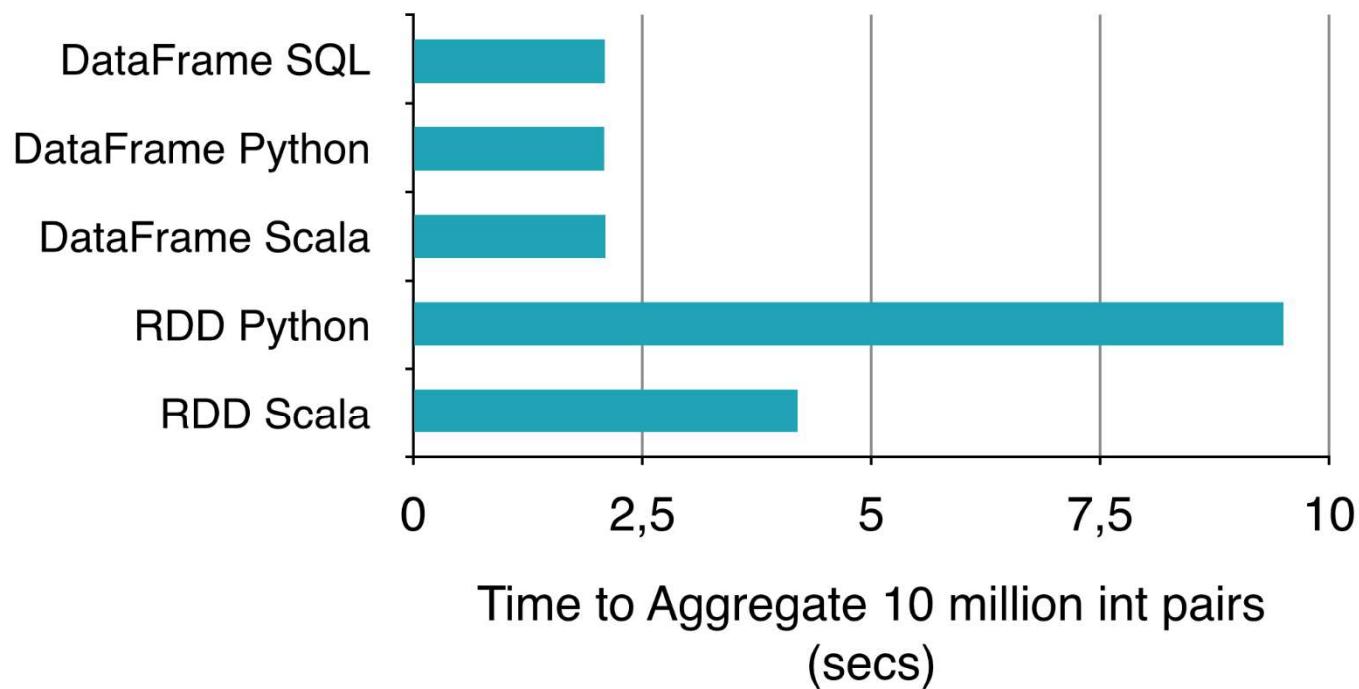
**Figure 3: Phases of query planning in Spark SQL. Rounded rectangles represent Catalyst trees.**

資料來源:Spark SQL: Relational Data Processing in Spark  
sigmod 2015



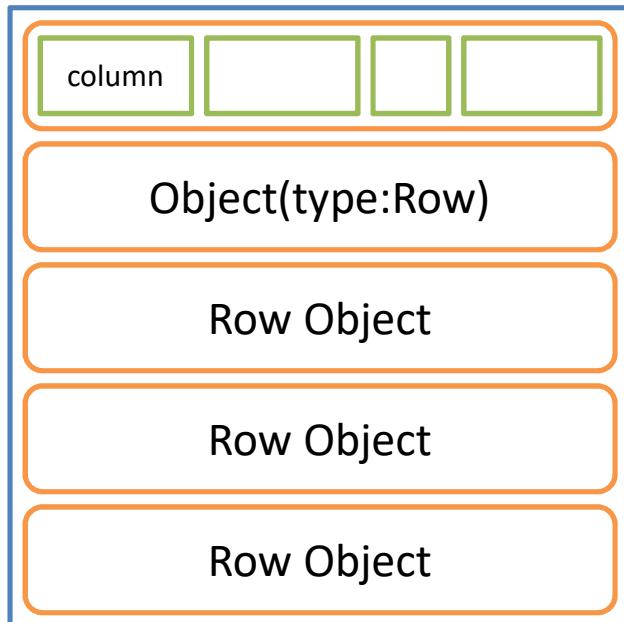
# Benefit of Logical Plan

Physical Execution: Unified Across Languages



**Schema** { Column Name  
Column Type

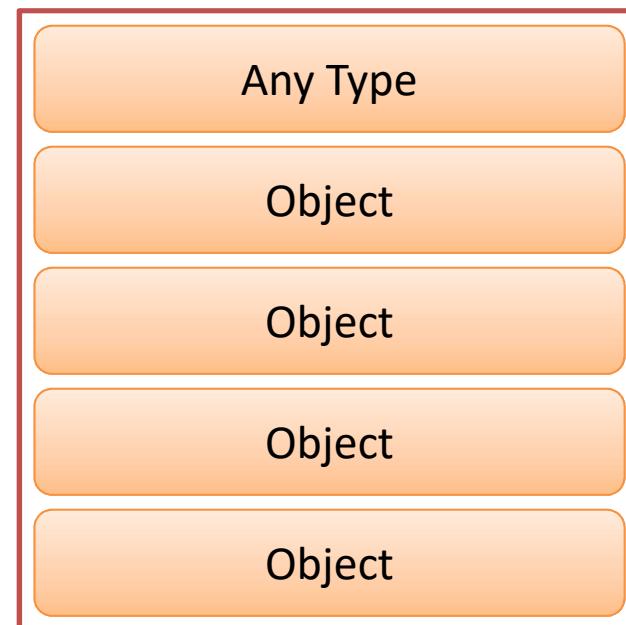
Python/Java/Scala



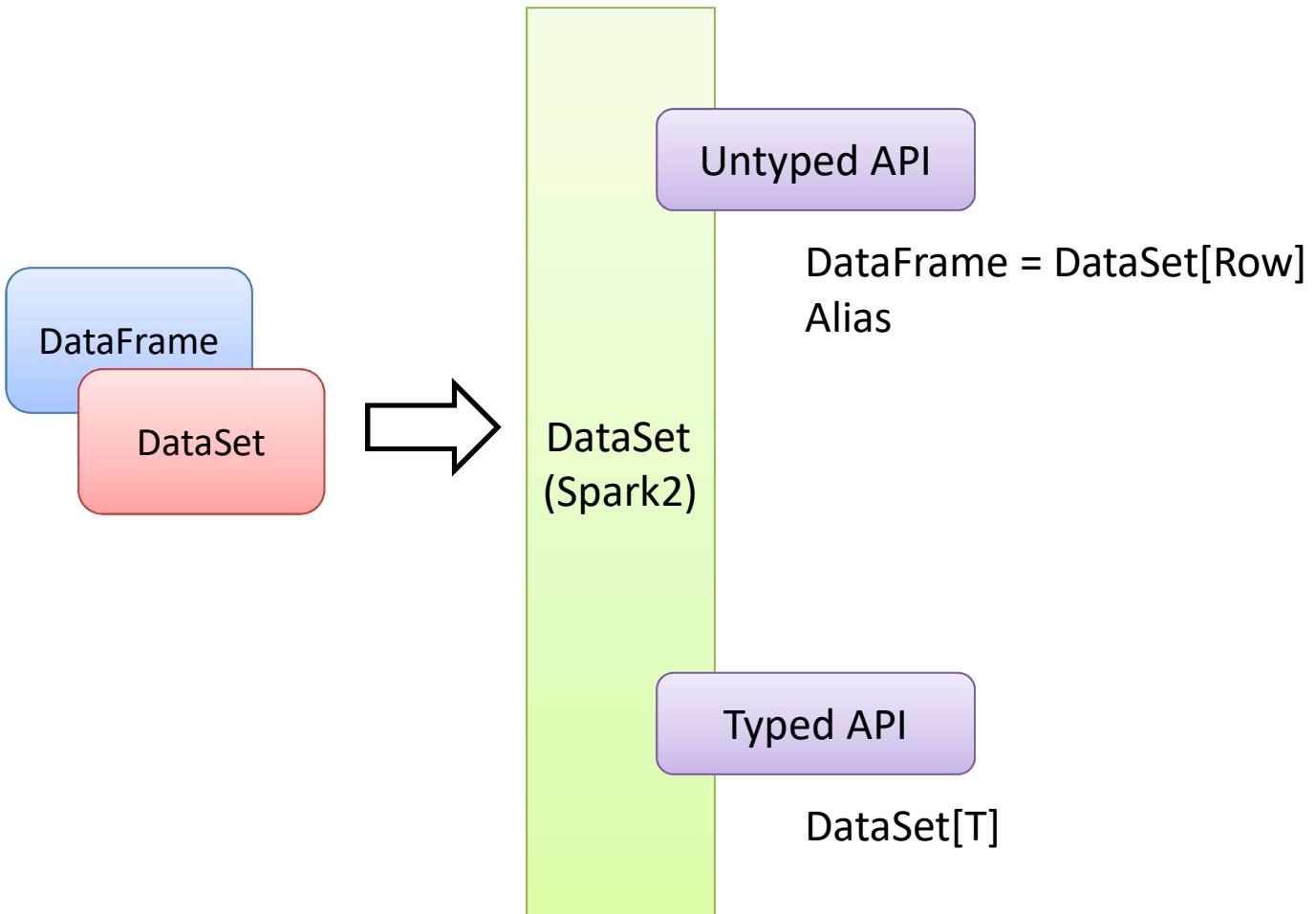
DataFrame

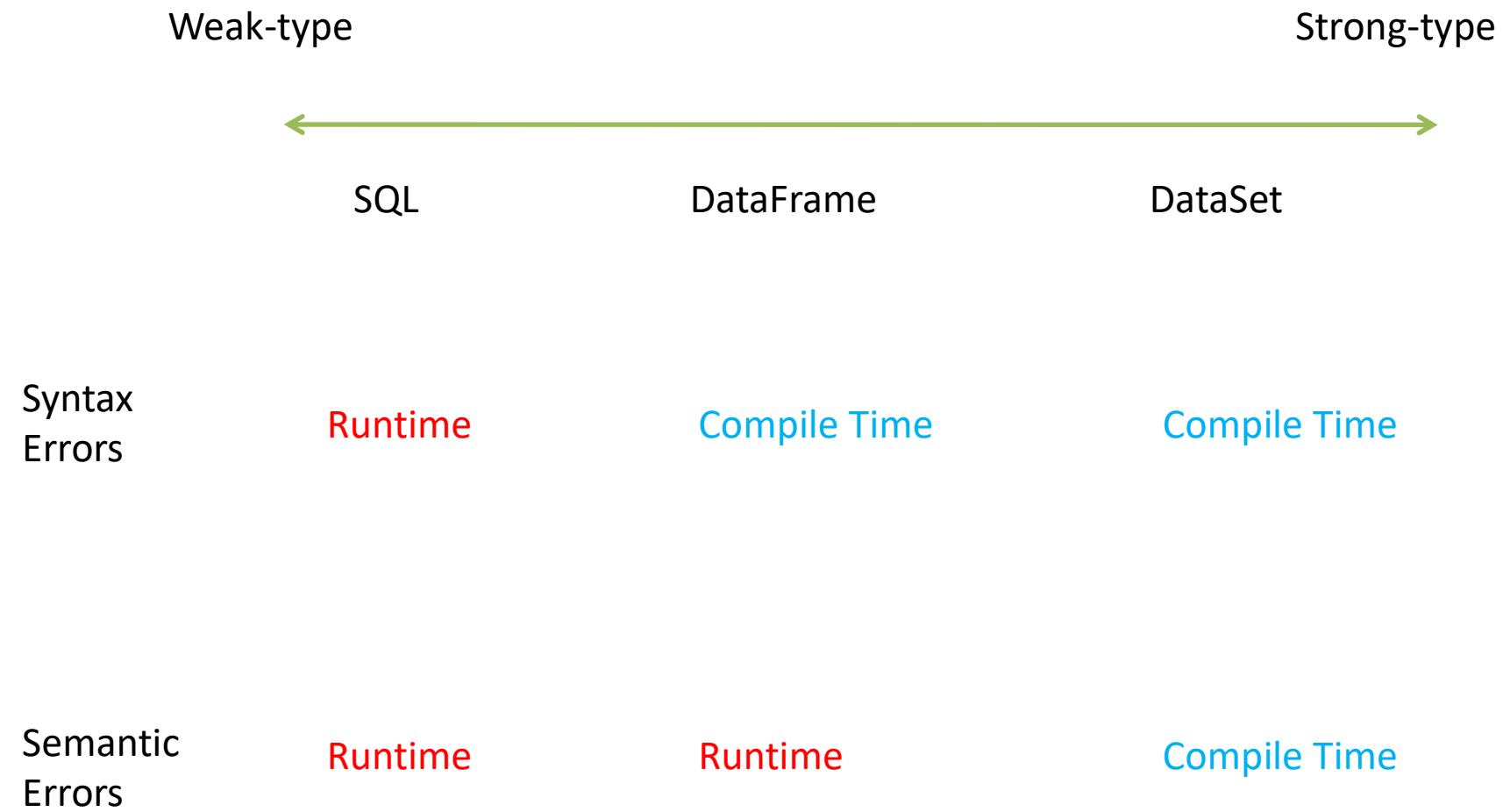
DataSet[Row]

Java/Scala



DataSet

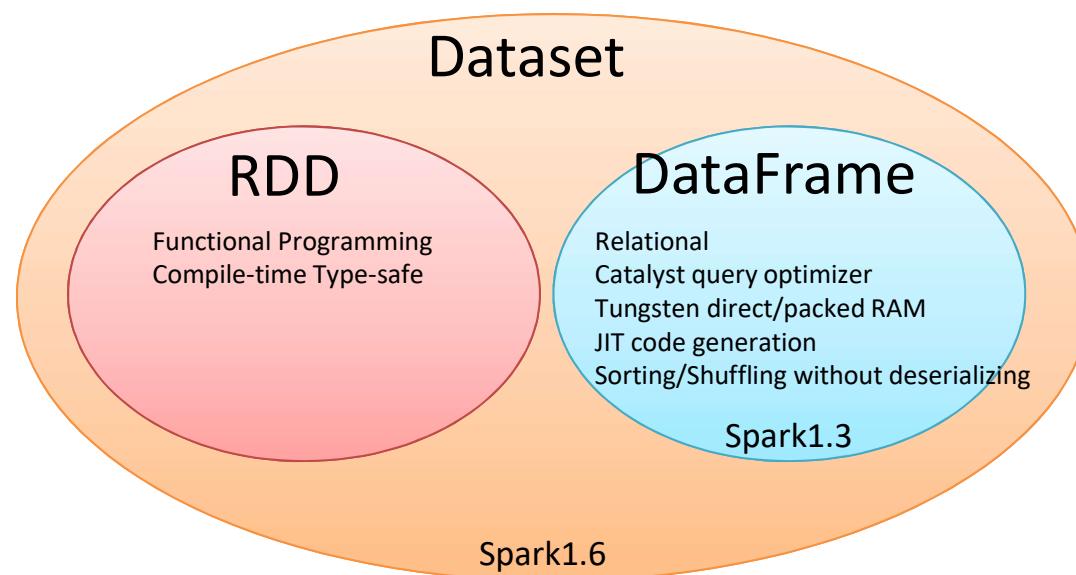




# DataSet

- Spark 1.6增加DataSet概念
  - 類似RDD
  - 既有Spark SQL的好處，又允許使用者自行定義物件與使用lambda函數
  - 提供編譯階段的型別檢查，與物件導向式API。
  - 可以自由地在DataSet、DataFrame、RDD間轉換。
- Spark 2.0開始，統一DataSet與DataFrame API。
  - DataFrame被定義為Dataset[ROW]

# DataSet是什麼？



# RDD、DataFrame、DataSet之間的關係

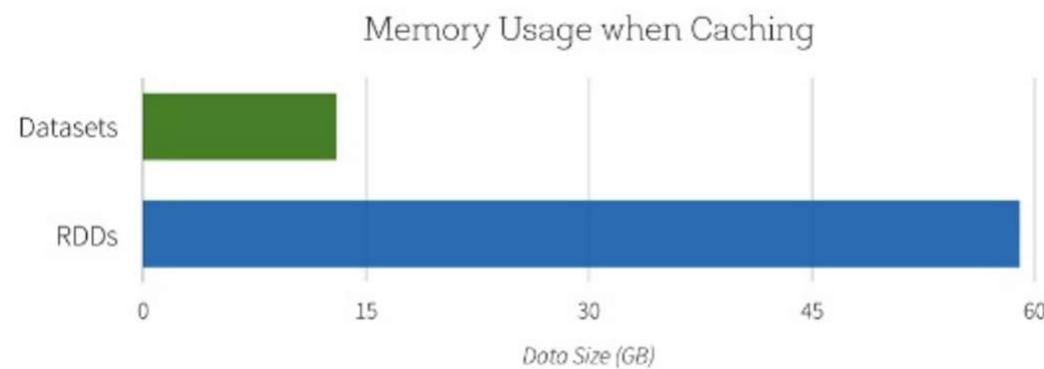
- DataSet是分散式的資料集合。
  - 提供RDD的優點(strong type、能夠使用lambda函數)與Spark SQL執行引擎的優點。
  - DataSet可以由JVM物件建構並使用轉換功能(map、flatMap、filter等等)。
  - DataSet API可以在Scala與Java API使用，但Python不支援靜態型態，所以不能使用。
- DataFrame是由DataSet組成的Row物件。
  - 概念等同關聯式資料庫的表格或R/Python的DataFrame。

# Typed API與Untyped API

語言	主要的抽象
Scala	DataSet[T] & DataFrame(alias for Dataset[ROW])
Java	DataSet[T]
Python	DataFrame
R	DataFrame

因為Python與R並沒有編譯階段的型態安全，所以只能使用DataFrame

## Space Efficiency



資料來源 : A Tale of Three Apache Spark APIs:RDDs vs DataFrames and DataSets

# 何時該用RDD

- 希望能夠對資料集進行最基本的轉換、處理和控制
- 資料結構為非結構化，例如：media stream或text stream
- 希望使用函數程式設計(functional programming)而不是特定領域表示(domain specific expression)來處理資料
- 不想要使用綱要、行格式(columnar)、以及使用名字或欄位存取資料
- 不想得到利用DataFrame/DataSet操作結構化、半結構化資料所得到的程式優化與效能提升

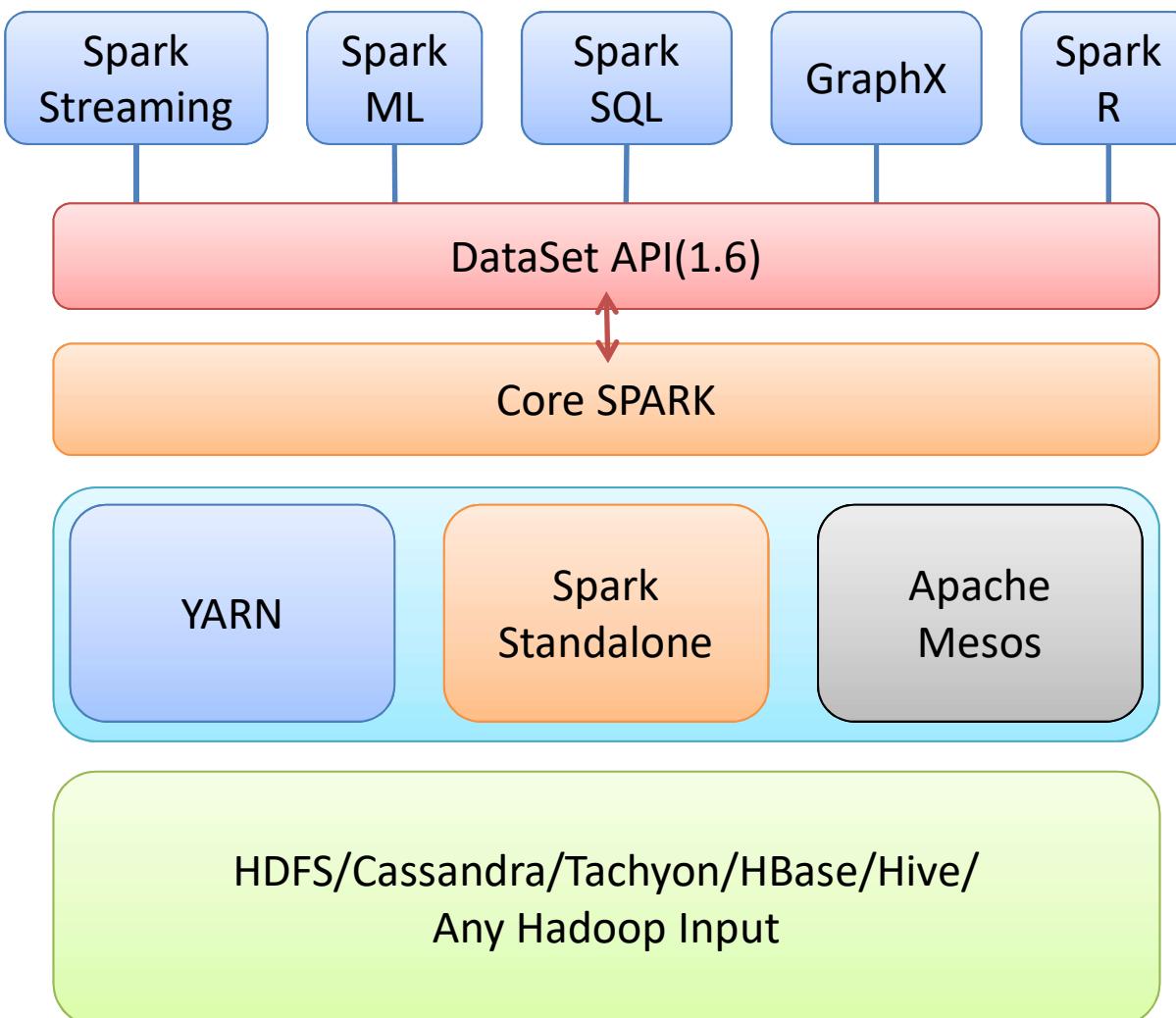
資料來源：A Tale of Three Apache Spark APIs:RDDs vs DataFrames and DataSets

# 何時該用 DataFrame/DataSet

- 需要豐富的語義、高階抽象與特定領域專用的API
- 需要對半結構化資料進行高階處理如：filter、map、aggregation、average、sum、SQL查詢、columnar存取或使用lambda函數
- 需要編譯階段的類型安全、typed JVM物件、使用Catalyst優化並利用Tungsten所產生的高效程式碼，則建議使用Dataset
- 需要在不同的Spark程式庫之間使用一致和簡化的API
- 使用R撰寫Spark程式碼，則建議使用DataFrame
- 使用Python撰寫Spark程式碼，一般建議使用DataFrame。但若要進行更精細的操作，則可以用RDD。

資料來源：A Tale of Three Apache Spark APIs:RDDs vs DataFrames and DataSets

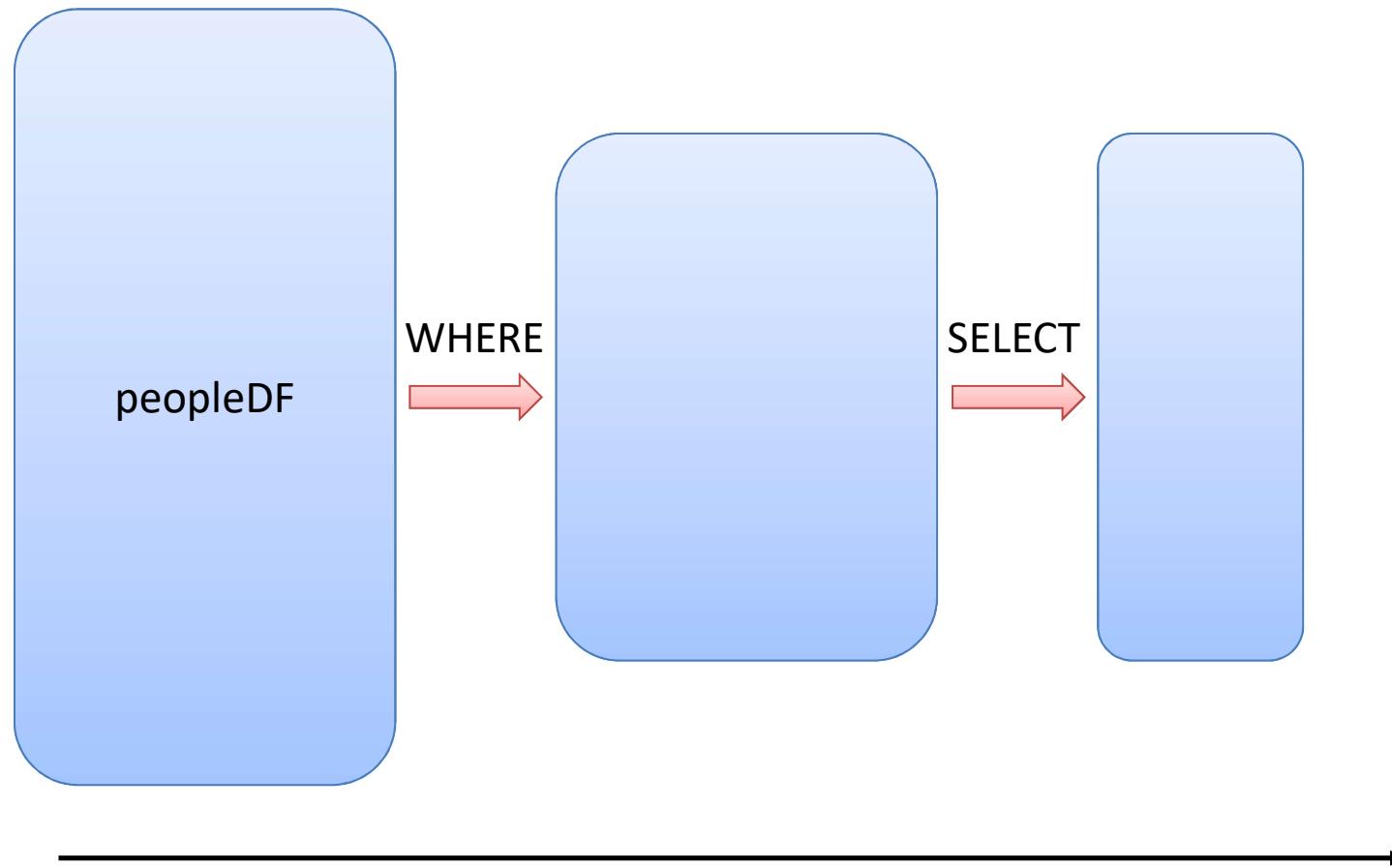
# DataSet is future of Spark 1.6



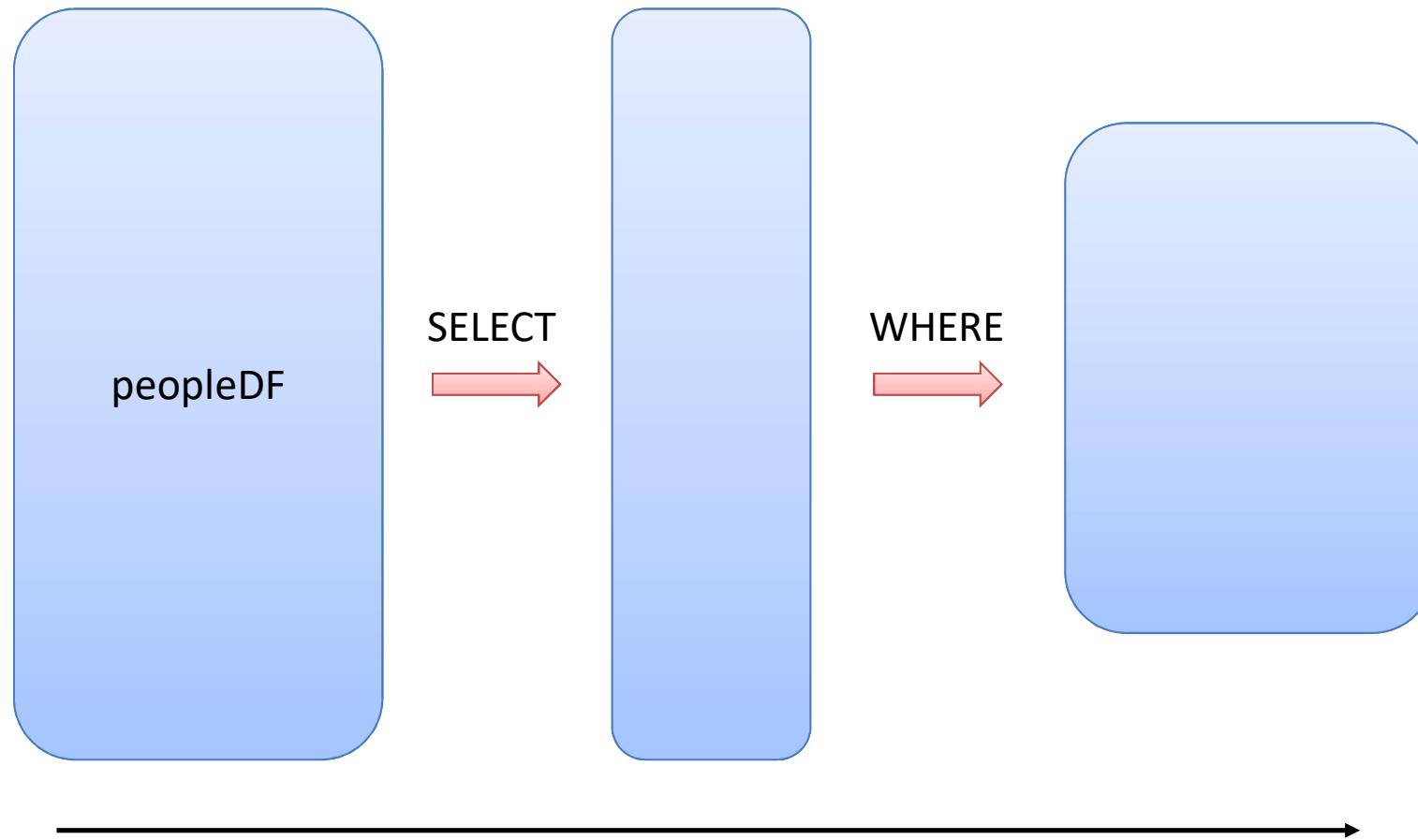
# SPJ

- SELECTION
  - WHERE: 過濾 rows
- PROJECTION
  - SELECT: 摳取 columns
- JOIN
  - 結合兩個 DataFrame

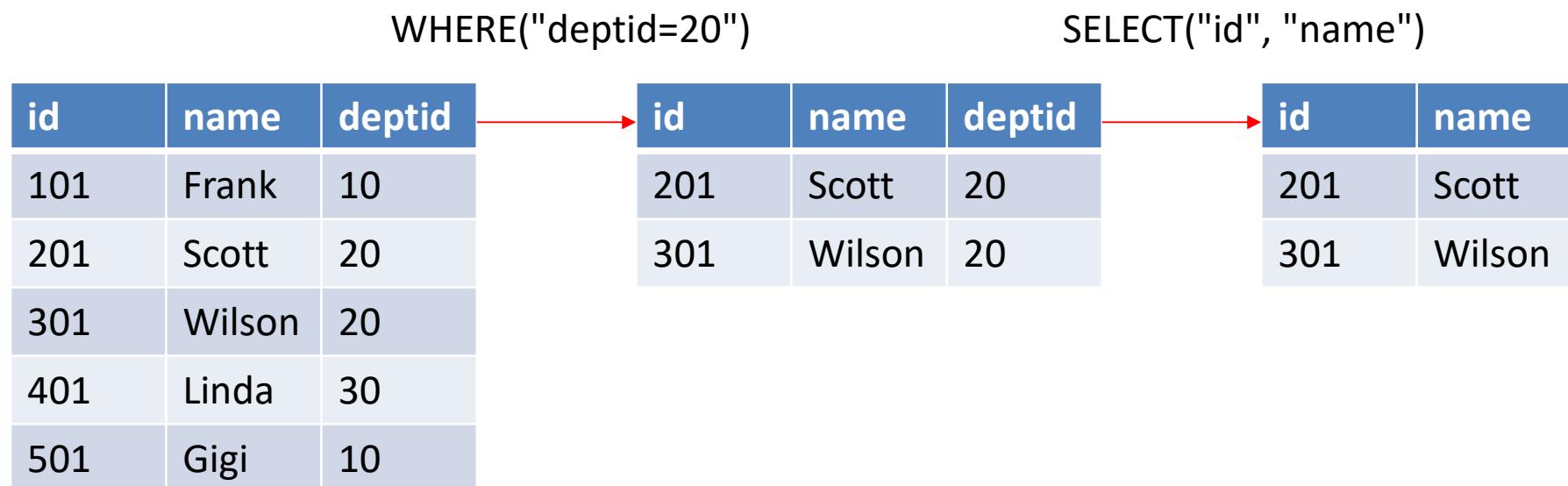
# Transformation order



# Transformation order



```
df.where("deptid=20").select("id", "name")
```



```
df1.join(df2,join_conds,join_type)
```

- INNER JOIN
- OUTER JOIN
  - LEFT OUTER JOIN
  - RIGHT OUTER JOIN
  - FULL OUTER JOIN
- CROSS JOIN

# JOIN\_TYPES

```
'inner', 'outer', 'full', 'fullouter', 'full_outer', 'leftouter', 'left',
'left_outer', 'rightouter', 'right', 'right_outer', 'leftsemi',
'left_semi', 'leftanti', 'left_anti', 'cross'
```

# empDF,deptDF

```
>>> empDF.show()
+---+-----+-----+-----+-----+-----+
|EMPNO| ENAME |      JOB | MGR | HIREDATE | SAL |COMM| DEPTNO |
+---+-----+-----+-----+-----+-----+
| 7839 | KING | PRESIDENT | null | 17-NOV-81 | 5000 |null|     10 |
| 7698 | BLAKE |    MANAGER | 7839 |01-MAY-81 | 2850 |null|     30 |
| 7782 | CLARK |    MANAGER | 7839 |09-JUN-81 | 2450 |null|     10 |
| 7566 | JONES |    MANAGER | 7839 |02-APR-81 | 2975 |null|     20 |
| 7788 | SCOTT | ANALYST | 7566 |19-APR-87 | 3000 |null|     20 |
| 7902 | FORD | ANALYST | 7566 |03-DEC-81 | 3000 |null|     20 |
| 7369 | SMITH |    CLERK | 7902 |17-DEC-80 |  800 |null|     20 |
| 7499 | ALLEN | SALESMAN | 7698 |20-FEB-81 | 1600 | 300|     30 |
| 7521 | WARD | SALESMAN | 7698 |22-FEB-81 | 1250 | 500|     30 |
| 7654 |MARTIN | SALESMAN | 7698 |28-SEP-81 | 1250 |1400|     30 |
| 7844 |TURNER | SALESMAN | 7698 |08-SEP-81 | 1500 |   0|     30 |
| 7876 | ADAMS |    CLERK | 7788 |23-MAY-87 | 1100 |null|     20 |
| 7900 | JAMES |    CLERK | 7698 |03-DEC-81 |  950 |null|     30 |
| 7934 |MILLER |    CLERK | 7782 |23-JAN-82 | 1300 |null|     10 |
| 7901 | FRANK | TRAINER | 7839 |23-JAN-85 | 3300 |null|     50 |
+---+-----+-----+-----+-----+-----+
```

```
>>> deptDF.show()
+---+-----+
|deptno|      dname |      loc |
+---+-----+
| 10 |ACCOUNTING | NEW YORK |
| 20 |    RESEARCH |    DALLAS |
| 30 |       SALES | CHICAGO |
| 40 |OPERATIONS |    BOSTON |
+---+-----+
```

# INNER JOIN(equality join)

```
>>> empDF.join(deptDF,empDF.deptno==deptDF.deptno,'inner').show()
+-----+-----+-----+-----+-----+-----+-----+-----+
|empno|  name|    job|  mgr| hiredate|   sal|comm|deptno|deptno|      dname|      loc|
+-----+-----+-----+-----+-----+-----+-----+-----+
| 7839| KING |PRESIDENT| null|17-NOV-81|5000| null| 10| 10|ACCOUNTING| NEW YORK|
| 7698| BLAKE| MANAGER| 7839|01-MAY-81|2850| null| 30| 30|SALES| CHICAGO|
| 7782| CLARK| MANAGER| 7839|09-JUN-81|2450| null| 10| 10|ACCOUNTING| NEW YORK|
| 7566| JONES| MANAGER| 7839|02-APR-81|2975| null| 20| 20|RESEARCH| DALLAS|
| 7788| SCOTT| ANALYST| 7566|19-APR-87|3000| null| 20| 20|RESEARCH| DALLAS|
| 7902| FORD | ANALYST| 7566|03-DEC-81|3000| null| 20| 20|RESEARCH| DALLAS|
| 7369| SMITH| CLERK| 7902|17-DEC-80| 800| null| 20| 20|RESEARCH| DALLAS|
| 7499| ALLEN| SALESMAN| 7698|20-FEB-81|1600| 300| 30| 30|SALES| CHICAGO|
| 7521| WARD | SALESMAN| 7698|22-FEB-81|1250| 500| 30| 30|SALES| CHICAGO|
| 7654|MARTIN| SALESMAN| 7698|28-SEP-81|1250|1400| 30| 30|SALES| CHICAGO|
| 7844|TURNER| SALESMAN| 7698|08-SEP-81|1500| 0| 30| 30|SALES| CHICAGO|
| 7876| ADAMS | CLERK| 7788|23-MAY-87|1100| null| 20| 20|RESEARCH| DALLAS|
| 7900| JAMES | CLERK| 7698|03-DEC-81| 950| null| 30| 30|SALES| CHICAGO|
| 7934|MILLER| CLERK| 7782|23-JAN-82|1300| null| 10| 10|ACCOUNTING| NEW YORK|
+-----+-----+-----+-----+-----+-----+-----+-----+
```

# INNER JOIN(non-equality join)

```
>>> empDF.join(deptDF,empDF.deptno>deptDF.deptno,'inner').show()
+-----+-----+-----+-----+-----+-----+-----+-----+
|empno|  name|    job|  mgr| hiredate|   sal|comm|deptno|deptno|      dname|      loc|
+-----+-----+-----+-----+-----+-----+-----+-----+
| 7698| BLAKE| MANAGER|7839|01-MAY-81|2850|null|     30|      10|ACCOUNTING|NEW YORK|
| 7698| BLAKE| MANAGER|7839|01-MAY-81|2850|null|     30|      20|  RESEARCH| DALLAS|
| 7566| JONES| MANAGER|7839|02-APR-81|2975|null|     20|      10|ACCOUNTING|NEW YORK|
| 7788| SCOTT| ANALYST|7566|19-APR-87|3000|null|     20|      10|ACCOUNTING|NEW YORK|
| 7902| FORD| ANALYST|7566|03-DEC-81|3000|null|     20|      10|ACCOUNTING|NEW YORK|
| 7369| SMITH| CLERK|7902|17-DEC-80|  800|null|     20|      10|ACCOUNTING|NEW YORK|
| 7499| ALLEN| SALESMAN|7698|20-FEB-81|1600|  300|     30|      10|ACCOUNTING|NEW YORK|
| 7499| ALLEN| SALESMAN|7698|20-FEB-81|1600|  300|     30|      20|  RESEARCH| DALLAS|
| 7521| WARD| SALESMAN|7698|22-FEB-81|1250|   500|     30|      10|ACCOUNTING|NEW YORK|
| 7521| WARD| SALESMAN|7698|22-FEB-81|1250|   500|     30|      20|  RESEARCH| DALLAS|
| 7654|MARTIN| SALESMAN|7698|28-SEP-81|1250| 1400|     30|      10|ACCOUNTING|NEW YORK|
| 7654|MARTIN| SALESMAN|7698|28-SEP-81|1250| 1400|     30|      20|  RESEARCH| DALLAS|
| 7844|TURNER| SALESMAN|7698|08-SEP-81|1500|     0|     30|      10|ACCOUNTING|NEW YORK|
| 7844|TURNER| SALESMAN|7698|08-SEP-81|1500|     0|     30|      20|  RESEARCH| DALLAS|
| 7876| ADAMS| CLERK|7788|23-MAY-87|1100|null|     20|      10|ACCOUNTING|NEW YORK|
| 7900| JAMES| CLERK|7698|03-DEC-81|  950|null|     30|      10|ACCOUNTING|NEW YORK|
| 7900| JAMES| CLERK|7698|03-DEC-81|  950|null|     30|      20|  RESEARCH| DALLAS|
+-----+-----+-----+-----+-----+-----+-----+-----+
```

# LEFT OUTER JOIN

```
>>> empDF.join(deptDF,empDF.deptno==deptDF.deptno,'left_outer').show()
+-----+-----+-----+-----+-----+-----+-----+-----+
|empno|  name|    job|  mgr| hiredate|   sal|comm|deptno|deptno|      dname|      loc|
+-----+-----+-----+-----+-----+-----+-----+-----+
| 7839| KING |PRESIDENT| null|17-NOV-81| 5000| null|     10|     10|ACCOUNTING| NEW YORK|
| 7698| BLAKE |MANAGER | 7839|01-MAY-81| 2850| null|     30|     30|SALES | CHICAGO|
| 7782| CLARK |MANAGER | 7839|09-JUN-81| 2450| null|     10|     10|ACCOUNTING| NEW YORK|
| 7566| JONES |MANAGER | 7839|02-APR-81| 2975| null|     20|     20|RESEARCH'| DALLAS|
| 7788| SCOTT |ANALYST | 7566|19-APR-87| 3000| null|     20|     20|RESEARCH'| DALLAS|
| 7902| FORD  |ANALYST | 7566|03-DEC-81| 3000| null|     20|     20|RESEARCH'| DALLAS|
| 7369| SMITH |CLERK   | 7902|17-DEC-80|  800| null|     20|     20|RESEARCH'| DALLAS|
| 7499| ALLEN |SALESMAN| 7698|20-FEB-81| 1600|  300|     30|     30|SALES | CHICAGO|
| 7521| WARD  |SALESMAN| 7698|22-FEB-81| 1250|  500|     30|     30|SALES | CHICAGO|
| 7654|MARTIN |SALESMAN| 7698|28-SEP-81| 1250| 1400|     30|     30|SALES | CHICAGO|
| 7844|TURNER |SALESMAN| 7698|08-SEP-81| 1500|    0|     30|     30|SALES | CHICAGO|
| 7876| ADAMS |CLERK   | 7788|23-MAY-87| 1100| null|     20|     20|RESEARCH'| DALLAS|
| 7900| JAMES  |CLERK   | 7698|03-DEC-81|  950| null|     30|     30|SALES | CHICAGO|
| 7934|MILLER |CLERK   | 7782|23-JAN-82| 1300| null|     10|     10|ACCOUNTING| NEW YORK|
| 7901| FRANK  |TRAINER | 7839|23-JAN-85| 3300| null|     50|    null|    null|    null|
+-----+-----+-----+-----+-----+-----+-----+-----+
```

# RIGHT OUTER JOIN

```
>>> empDF.join(deptDF,empDF.deptno==deptDF.deptno,'right_outer').show()
+-----+-----+-----+-----+-----+-----+-----+-----+
|empno| name | job | mgr | hiredate| sal|comm|deptno|deptno|      dname|      loc|
+-----+-----+-----+-----+-----+-----+-----+-----+
| 7934|MILLER| CLERK|7782|23-JAN-82|1300|null|    10|      10|ACCOUNTING| NEW YORK|
| 7782| CLARK| MANAGER|7839|09-JUN-81|2450|null|    10|      10|ACCOUNTING| NEW YORK|
| 7839| KING| PRESIDENT|null|17-NOV-81|5000|null|    10|      10|ACCOUNTING| NEW YORK|
| 7876| ADAMS| CLERK|7788|23-MAY-87|1100|null|    20|      20| RESEARCH| DALLAS|
| 7369| SMITH| CLERK|7902|17-DEC-80|  800|null|    20|      20| RESEARCH| DALLAS|
| 7902| FORD| ANALYST|7566|03-DEC-81|3000|null|    20|      20| RESEARCH| DALLAS|
| 7788| SCOTT| ANALYST|7566|19-APR-87|3000|null|    20|      20| RESEARCH| DALLAS|
| 7566| JONES| MANAGER|7839|02-APR-81|2975|null|    20|      20| RESEARCH| DALLAS|
| 7900| JAMES| CLERK|7698|03-DEC-81|  950|null|    30|      30| SALES| CHICAGO|
| 7844|TURNER| SALESMAN|7698|08-SEP-81|1500|   0|    30|      30| SALES| CHICAGO|
| 7654|MARTIN| SALESMAN|7698|28-SEP-81|1250|1400|    30|      30| SALES| CHICAGO|
| 7521| WARD| SALESMAN|7698|22-FEB-81|1250|  500|    30|      30| SALES| CHICAGO|
| 7499| ALLEN| SALESMAN|7698|20-FEB-81|1600|  300|    30|      30| SALES| CHICAGO|
| 7698| BLAKE| MANAGER|7839|01-MAY-81|2850|null|    30|      30| SALES| CHICAGO|
| null|  null|     null|null|     null|null|null|      40|OPERATIONS| BOSTON|
+-----+-----+-----+-----+-----+-----+-----+-----+
```

# FULL OUTER JOIN

```
>>> empDF.join(deptDF,empDF.deptno==deptDF.deptno,'full_outer').show()
+-----+-----+-----+-----+-----+-----+-----+-----+
|empno|  name|    job|  mgr| hiredate|   sal|comm|deptno|deptno|      dname|      loc|
+-----+-----+-----+-----+-----+-----+-----+-----+
|  7698| BLAKE|  MANAGER|7839|01-MAY-81|2850|null| 30| 30|      SALES| CHICAGO|
|  7499| ALLEN| SALESMAN|7698|20-FEB-81|1600| 300| 30| 30|      SALES| CHICAGO|
|  7521|  WARD| SALESMAN|7698|22-FEB-81|1250| 500| 30| 30|      SALES| CHICAGO|
|  7654|MARTIN| SALESMAN|7698|28-SEP-81|1250|1400| 30| 30|      SALES| CHICAGO|
|  7844|TURNER| SALESMAN|7698|08-SEP-81|1500|    0| 30| 30|      SALES| CHICAGO|
|  7900| JAMES|    CLERK|7698|03-DEC-81|  950|null| 30| 30|      SALES| CHICAGO|
| null|  null|    null|null|    null|null|  null| 40|OPERATIONS| BOSTON|
|  7566|  JONES|  MANAGER|7839|02-APR-81|2975|null| 20| 20| RESEARCH| DALLAS|
|  7788| SCOTT| ANALYST|7566|19-APR-87|3000|null| 20| 20| RESEARCH| DALLAS|
|  7902|   FORD| ANALYST|7566|03-DEC-81|3000|null| 20| 20| RESEARCH| DALLAS|
|  7369| SMITH|    CLERK|7902|17-DEC-80|  800|null| 20| 20| RESEARCH| DALLAS|
|  7876| ADAMS|    CLERK|7788|23-MAY-87|1100|null| 20| 20| RESEARCH| DALLAS|
|  7839|   KING| PRESIDENT|null|17-NOV-81|5000|null| 10| 10|ACCOUNTING| NEW YORK|
|  7782| CLARK|  MANAGER|7839|09-JUN-81|2450|null| 10| 10|ACCOUNTING| NEW YORK|
|  7934|MILLER|    CLERK|7782|23-JAN-82|1300|null| 10| 10|ACCOUNTING| NEW YORK|
|  7901| FRANK|    TRAINER|7839|23-JAN-85|3300|null| 50| null|    null|    null|
+-----+-----+-----+-----+-----+-----+-----+-----+
```

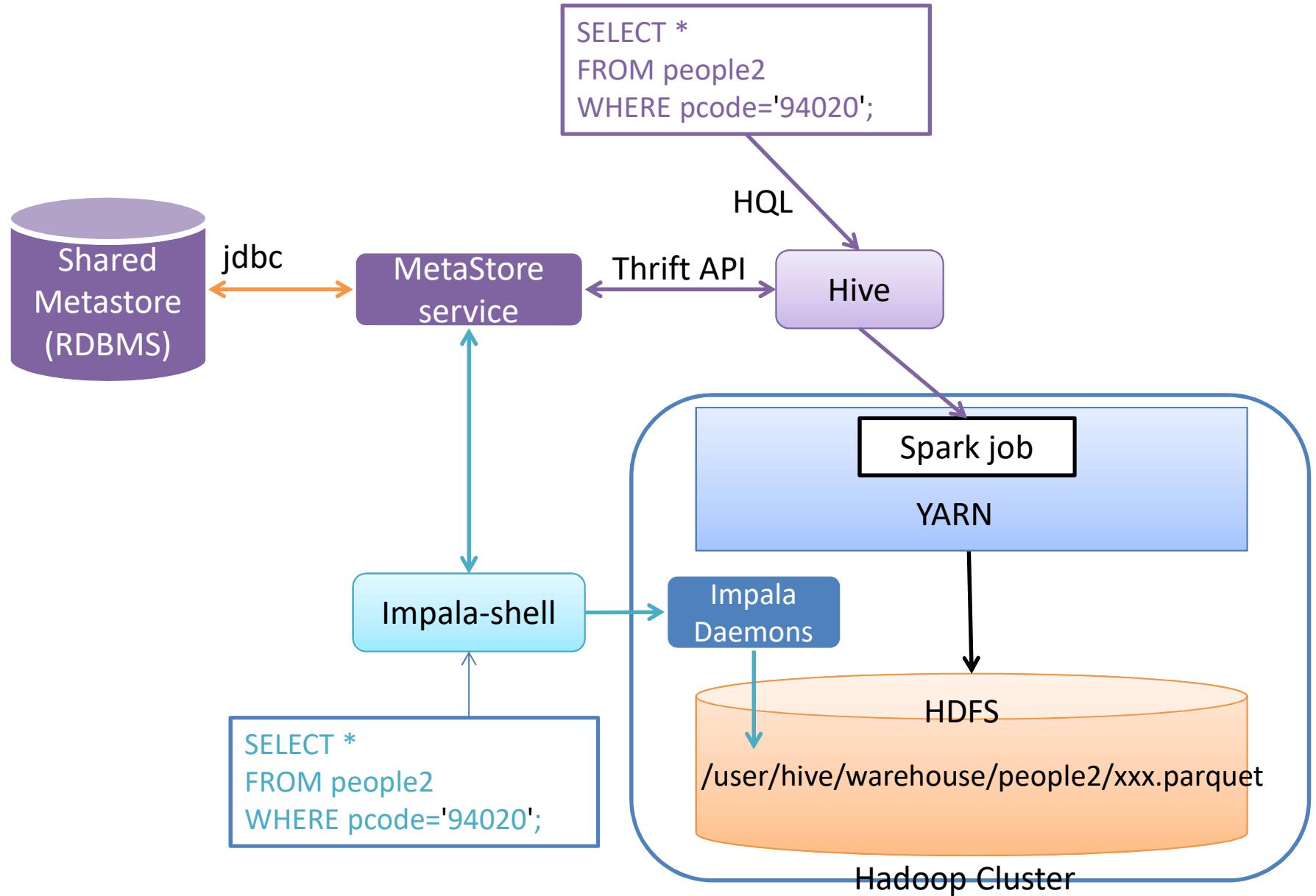
# CROSS JOIN

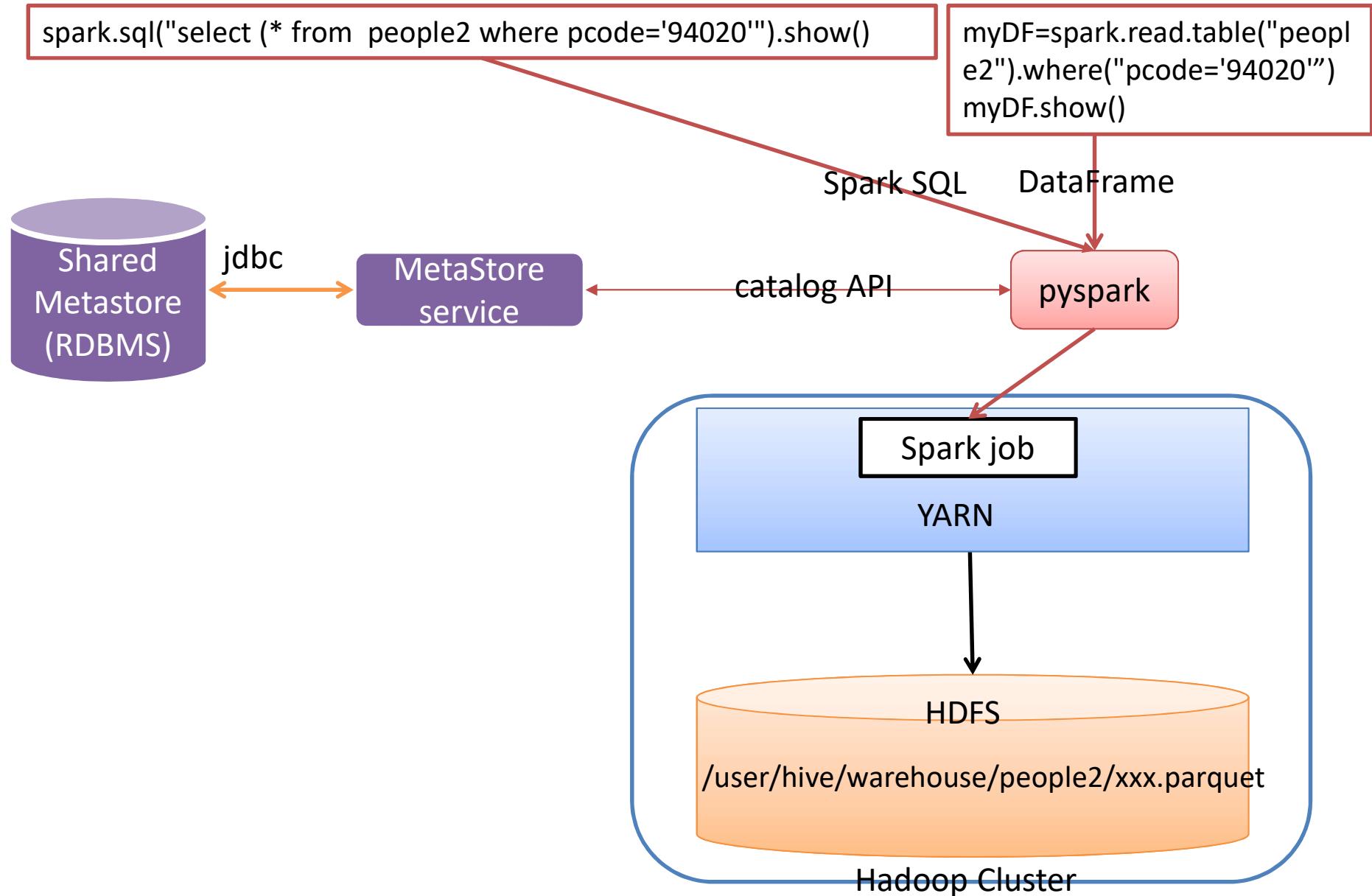
```
>>> empDF.crossJoin(deptDF).show()
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|empno| name | job | mgr | hiredate | sal | comm | deptno | deptno |      dname |      loc |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| 7839 | KING | PRESIDENT | null | 17-NOV-81 | 5000 | null | 10 | 10 | ACCOUNTING | NEW YORK |
| 7839 | KING | PRESIDENT | null | 17-NOV-81 | 5000 | null | 10 | 20 | RESEARCH | DALLAS |
| 7839 | KING | PRESIDENT | null | 17-NOV-81 | 5000 | null | 10 | 30 | SALES | CHICAGO |
| 7839 | KING | PRESIDENT | null | 17-NOV-81 | 5000 | null | 10 | 40 | OPERATIONS | BOSTON |
| 7698 | BLAKE | MANAGER | 7839 | 01-MAY-81 | 2850 | null | 30 | 10 | ACCOUNTING | NEW YORK |
| 7698 | BLAKE | MANAGER | 7839 | 01-MAY-81 | 2850 | null | 30 | 20 | RESEARCH | DALLAS |
| 7698 | BLAKE | MANAGER | 7839 | 01-MAY-81 | 2850 | null | 30 | 30 | SALES | CHICAGO |
| 7698 | BLAKE | MANAGER | 7839 | 01-MAY-81 | 2850 | null | 30 | 40 | OPERATIONS | BOSTON |
| 7782 | CLARK | MANAGER | 7839 | 09-JUN-81 | 2450 | null | 10 | 10 | ACCOUNTING | NEW YORK |
| 7782 | CLARK | MANAGER | 7839 | 09-JUN-81 | 2450 | null | 10 | 20 | RESEARCH | DALLAS |
| 7782 | CLARK | MANAGER | 7839 | 09-JUN-81 | 2450 | null | 10 | 30 | SALES | CHICAGO |
| 7782 | CLARK | MANAGER | 7839 | 09-JUN-81 | 2450 | null | 10 | 40 | OPERATIONS | BOSTON |
| 7566 | JONES | MANAGER | 7839 | 02-APR-81 | 2975 | null | 20 | 10 | ACCOUNTING | NEW YORK |
| 7566 | JONES | MANAGER | 7839 | 02-APR-81 | 2975 | null | 20 | 20 | RESEARCH | DALLAS |
| 7566 | JONES | MANAGER | 7839 | 02-APR-81 | 2975 | null | 20 | 30 | SALES | CHICAGO |
| 7566 | JONES | MANAGER | 7839 | 02-APR-81 | 2975 | null | 20 | 40 | OPERATIONS | BOSTON |
| 7788 | SCOTT | ANALYST | 7566 | 19-APR-87 | 3000 | null | 20 | 10 | ACCOUNTING | NEW YORK |
| 7788 | SCOTT | ANALYST | 7566 | 19-APR-87 | 3000 | null | 20 | 20 | RESEARCH | DALLAS |
| 7788 | SCOTT | ANALYST | 7566 | 19-APR-87 | 3000 | null | 20 | 30 | SALES | CHICAGO |
| 7788 | SCOTT | ANALYST | 7566 | 19-APR-87 | 3000 | null | 20 | 40 | OPERATIONS | BOSTON |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 20 rows
```

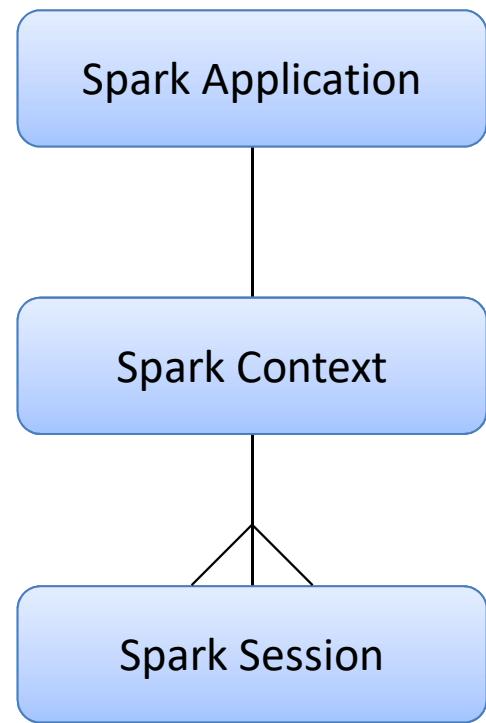
# LEFT-SEMI JOIN

```
>>> empDF.join(deptDF,empDF.deptno==deptDF.deptno,'left_semi').show()
+-----+-----+-----+-----+-----+-----+
|empno|    name|      job|  mgr| hiredate|   sal|comm|deptno|
+-----+-----+-----+-----+-----+-----+
| 7839|    KING|PRESIDENT| null|17-NOV-81|5000|null|     10|
| 7698|  BLAKE|  MANAGER|7839|01-MAY-81|2850|null|     30|
| 7782| CLARK|  MANAGER|7839|09-JUN-81|2450|null|     10|
| 7566|  JONES|  MANAGER|7839|02-APR-81|2975|null|     20|
| 7788| SCOTT| ANALYST|7566|19-APR-87|3000|null|     20|
| 7902|   FORD| ANALYST|7566|03-DEC-81|3000|null|     20|
| 7369| SMITH|    CLERK|7902|17-DEC-80|  800|null|     20|
| 7499|  ALLEN|SALESMAN|7698|20-FEB-81|1600| 300|     30|
| 7521|   WARD|SALESMAN|7698|22-FEB-81|1250|  500|     30|
| 7654|MARTIN|SALESMAN|7698|28-SEP-81|1250|1400|     30|
| 7844|TURNER|SALESMAN|7698|08-SEP-81|1500|    0|     30|
| 7876| ADAMS|    CLERK|7788|23-MAY-87|1100|null|     20|
| 7900| JAMES|    CLERK|7698|03-DEC-81|  950|null|     30|
| 7934|MILLER|    CLERK|7782|23-JAN-82|1300|null|     10|
+-----+-----+-----+-----+-----+-----+
```

FROM  
WHERE  
GROUP BY  
HAVING  
SELECT  
ORDER BY  
LIMIT







# Application Deployment mode

- --master yarn
  - --deploy-mode cluster
  - --deploy-mode client
- --master meson
- --master spark
- --master local

# RDD(Resilient Distributed Dataset)

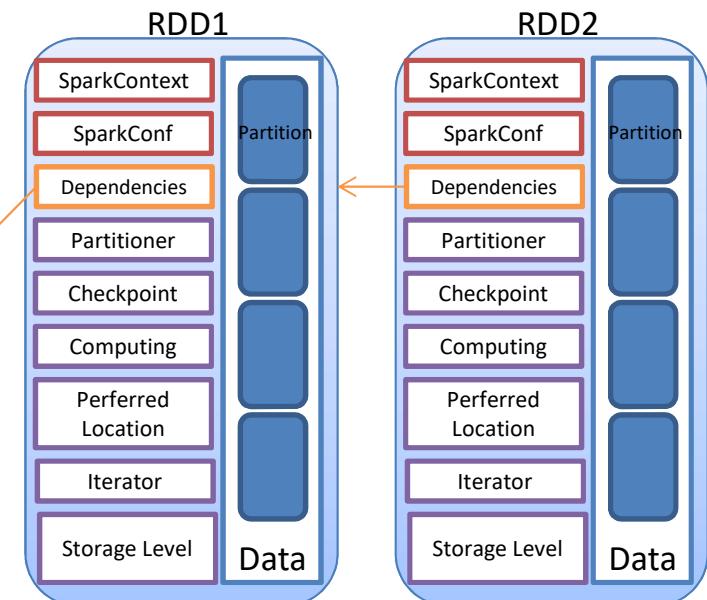
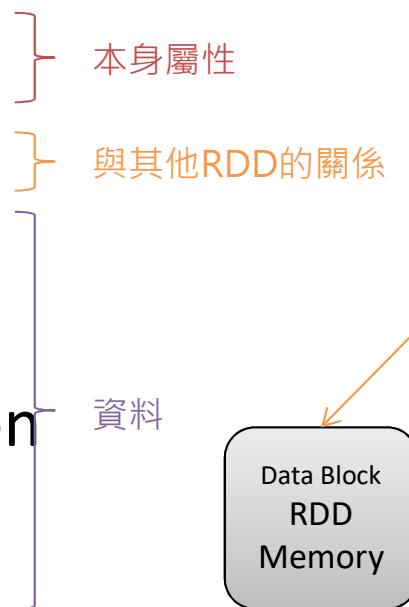
- 定義
  - 由可容錯元件所組成的集合，這個集合可以被平行操作。
    - RDDs are fault-tolerate, parallel data structures.
  - 是一個唯讀、被分區的紀錄集合。
    - An RDD is a read-only, partitioned collection of records.

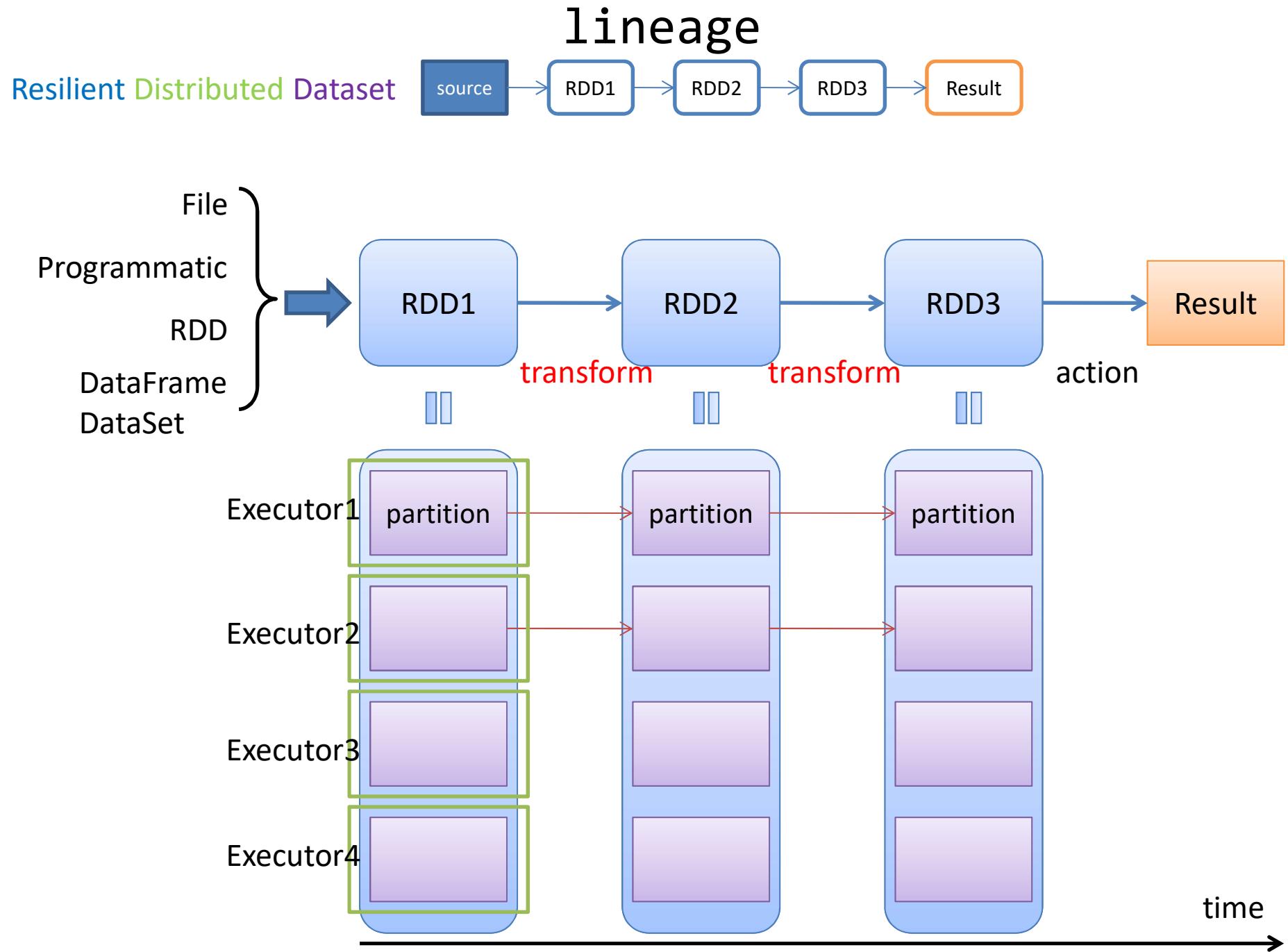
"Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing" 2012

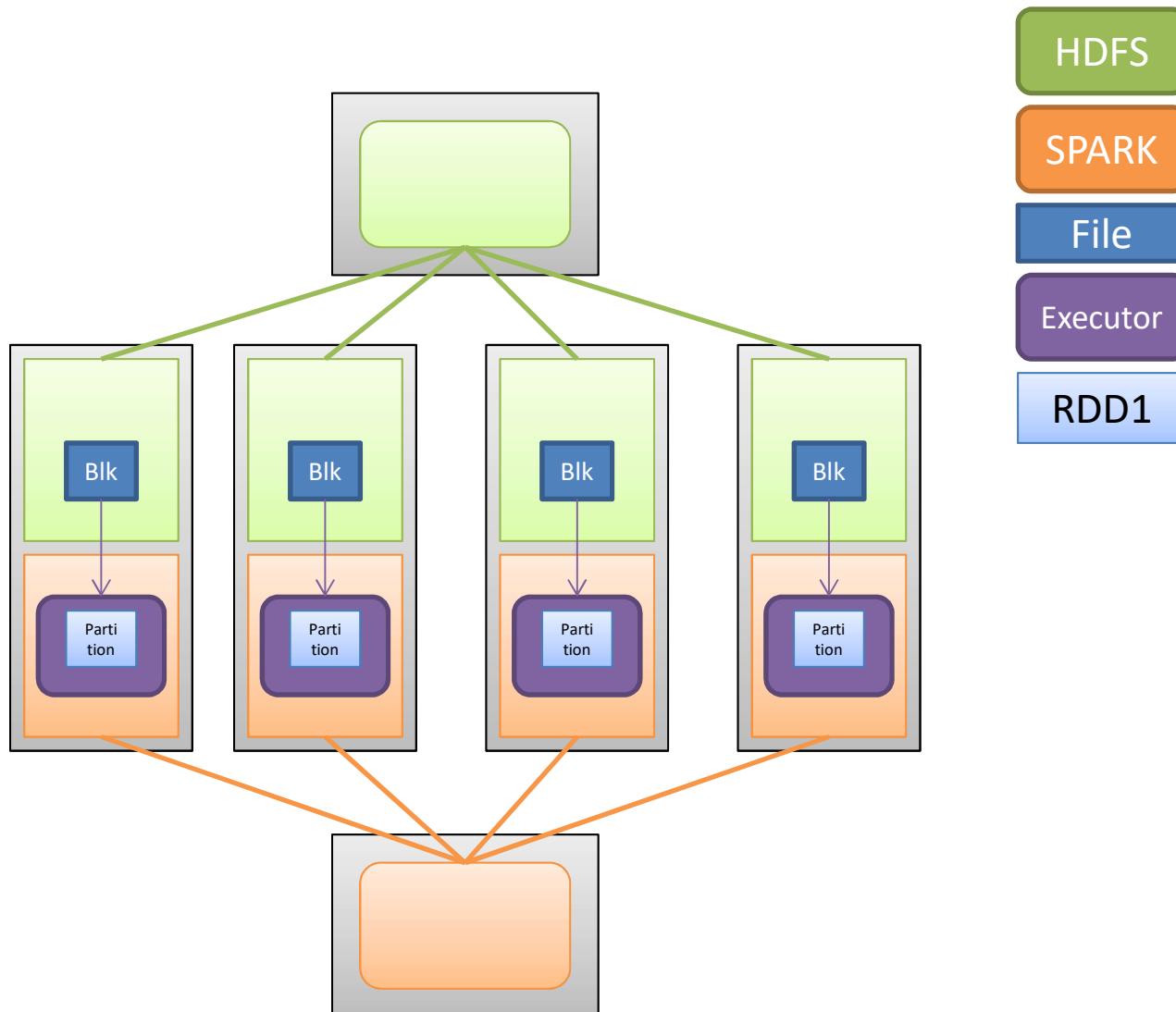
- Resilient
  - 有彈性的，可彈回。
  - 表示發生錯誤後，透過lineage可以重新產生RDD。
- Distributed
  - RDD由partition組成，而partition被分散到多個節點上執行。
- Dataset
  - 不是真正的“集合”，而是Spark處理資料的邏輯。
  - 資料放在partition裡。

# RDD的內部組成元件

- SparkContext
- SparkConf
- Dependencies
- Partitioner
- Checkpoint
- Computing
- Preferred Location
- Iterator
- Storage Level
- Data





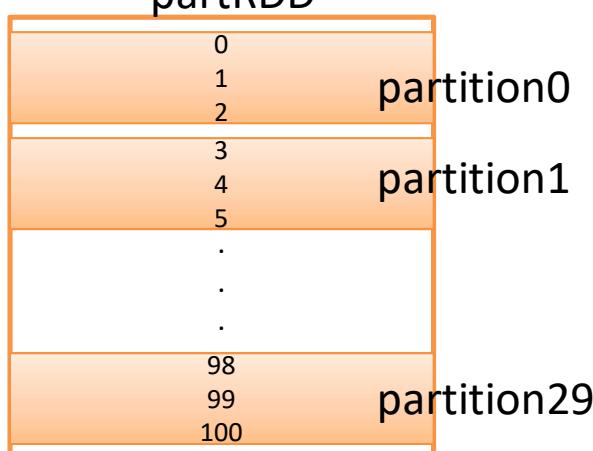


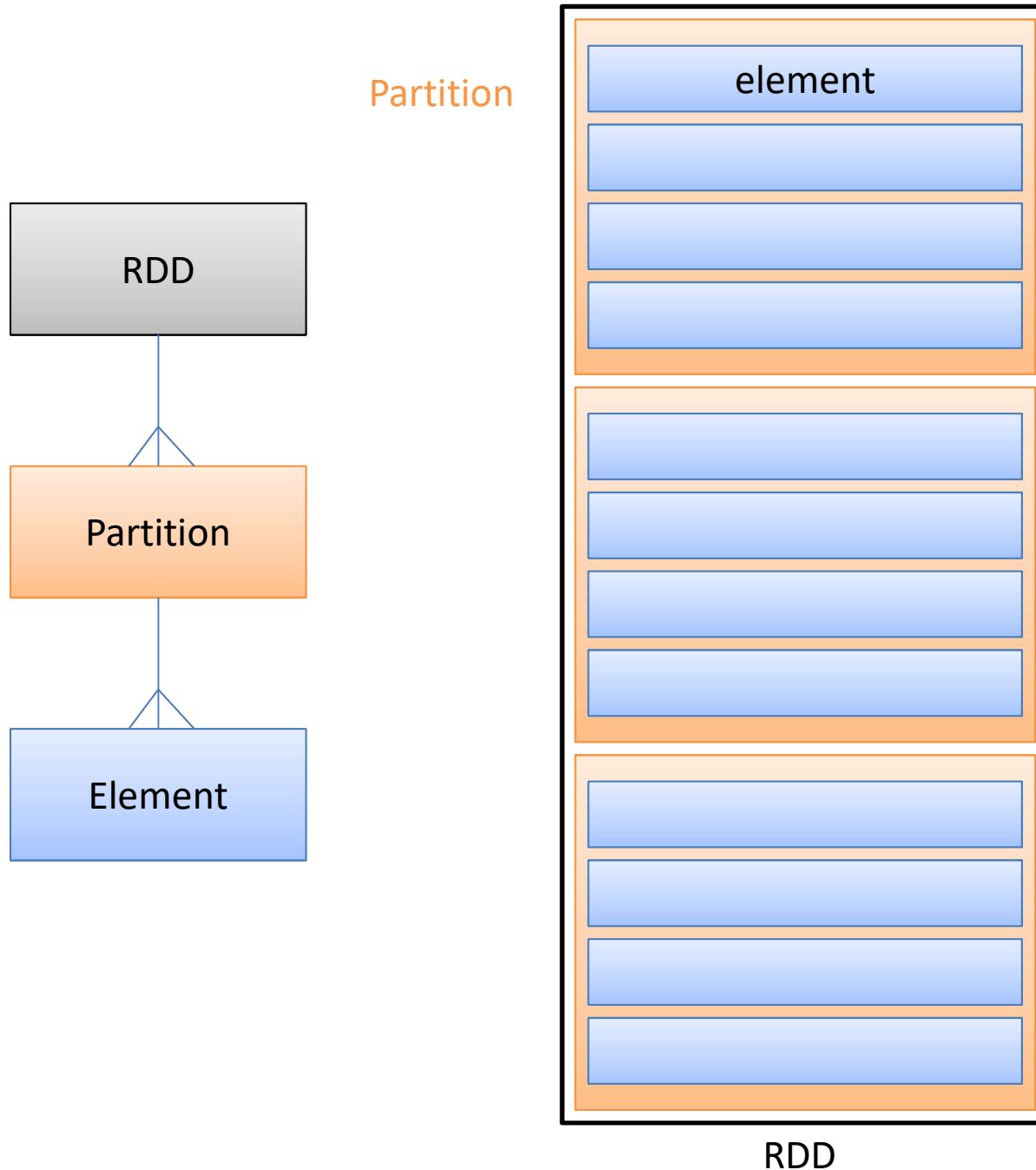
# Partition

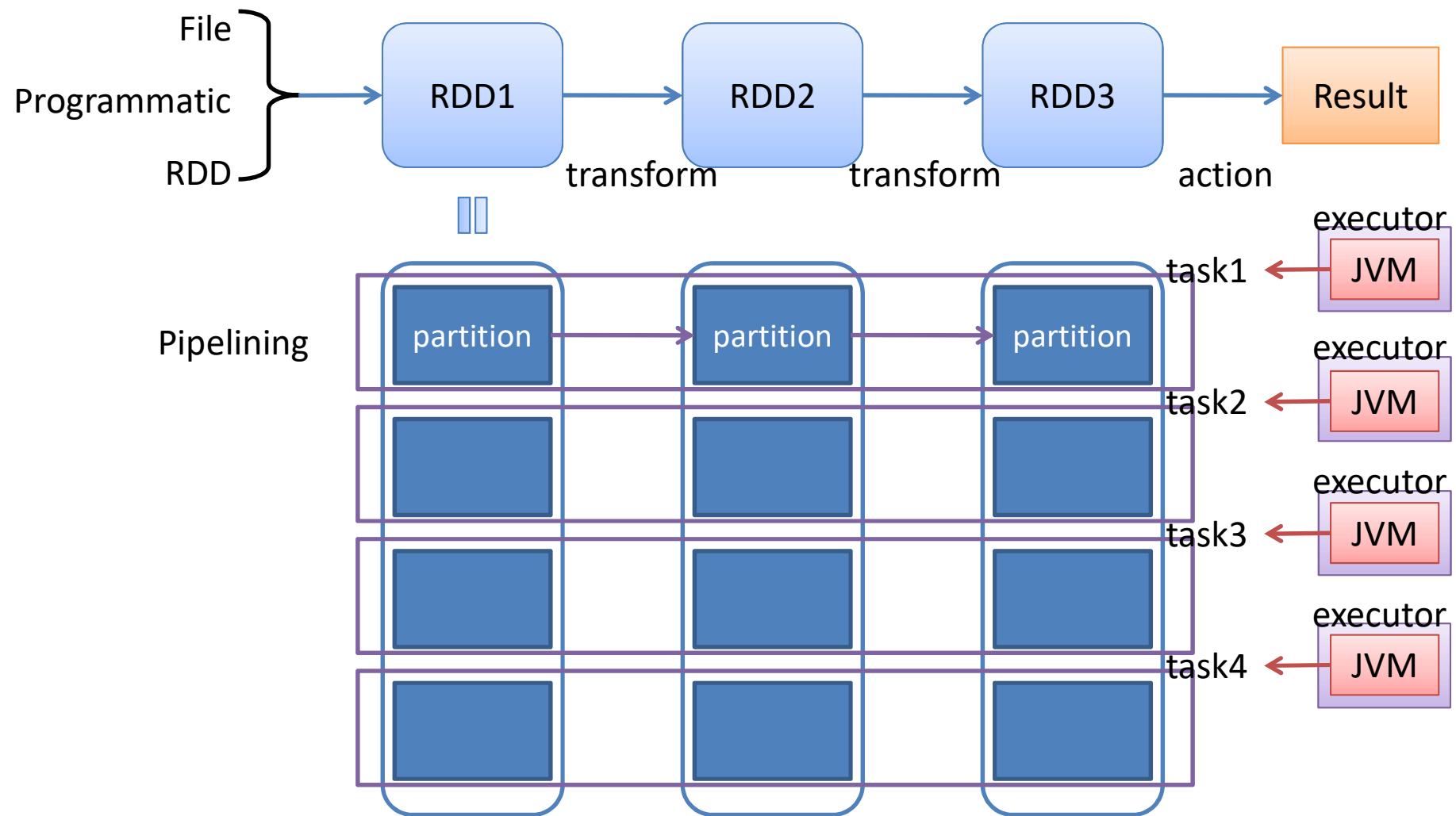
- 產生RDD時，明確指定  
`sc.textFile( "inputPath" ,minPartitions)`
- 若沒有明確指定partition數量，則由RDD生成方式決定預設partition數量
  - 由collection方式產生：
    - 為此application所分配到的CPU數量(Executors)
  - 由HDFS檔案產生：
    - 檔案的HDFS Block數量，在Cluster架構預設值為2
  - 由RDD產生
    - 繼承父系RDD的數量

# How Many Partitions

```
>>> range(101)
[0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17,
18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33,
34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49,
50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65,
66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81,
82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97,
98, 99, 100]
>>> partRDD = sc.parallelize(range(101), 30)
>>> partRDD.getNumPartitions()
30
```







```
In [1]: mydata = sc.parallelize(range(10))
```

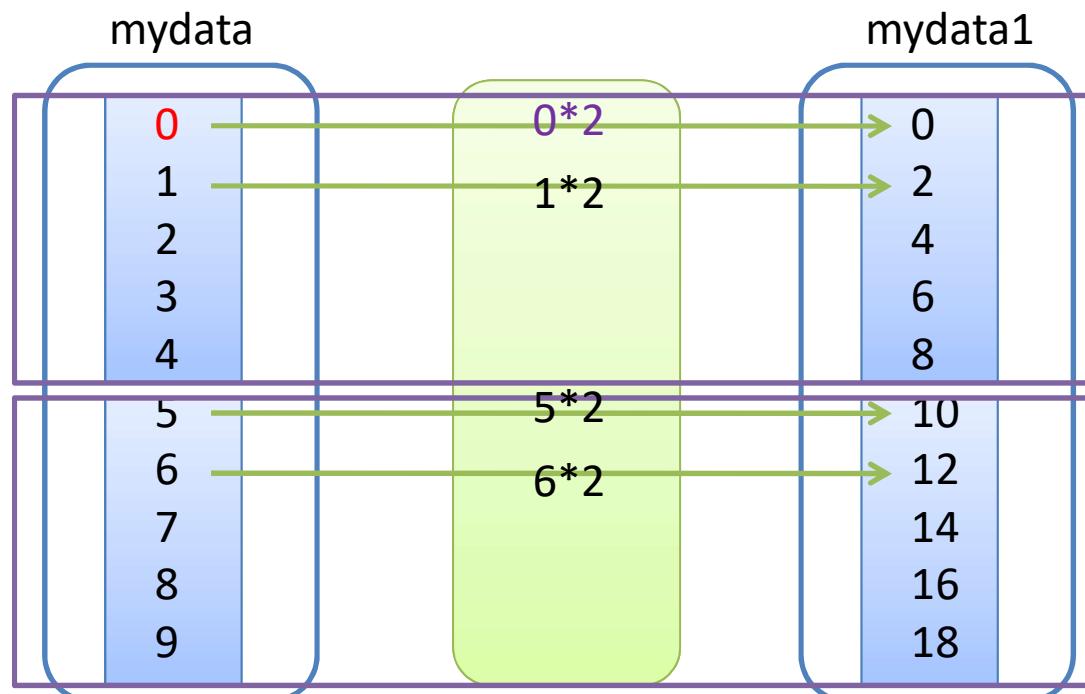
```
In [2]: mydata.getNumPartitions()
```

```
Out[2]: 2
```

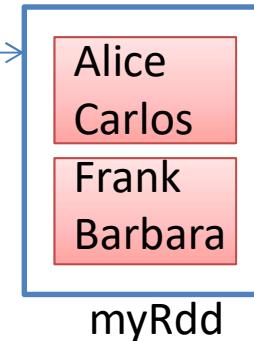
```
In [3]: mydata1 = mydata.map(lambda x: x*x)
```

```
In [4]: mydata1.collect()
```

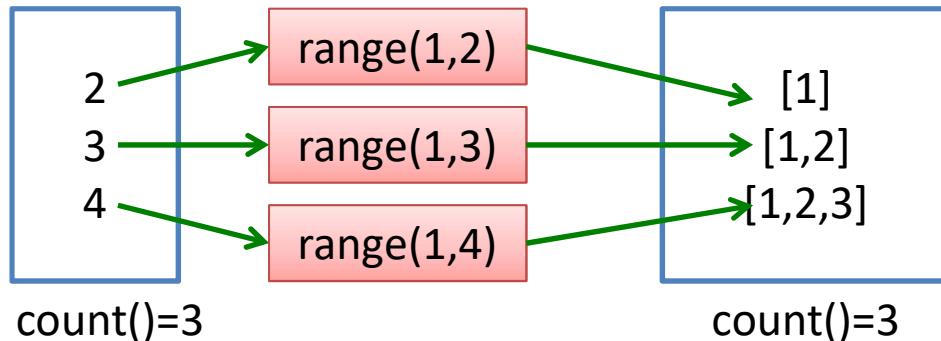
```
Out[4]: [0, 2, 4, 6, 8, 10, 12, 14, 16, 18]
```



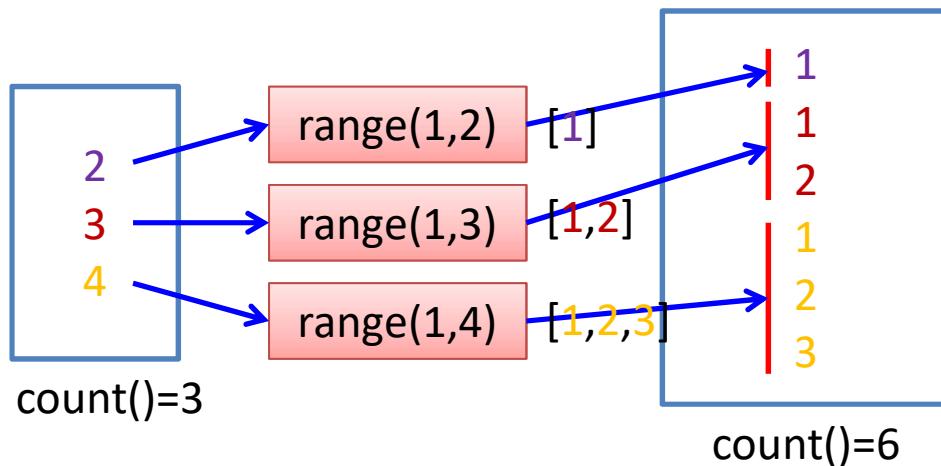
```
> myData = ["Alice","Carlos","Frank","Barbara"]  
  
> myRdd = sc.parallelize(myData) →  
  
> myRdd  
ParallelCollectionRDD[15] at parallelize at  
PythonRDD.scala:376  
  
> myRdd.take(2)  
['Alice', 'Carlos']  
  
> myRdd.count()  
4  
  
> myRdd.getNumPartitions()  
2  
  
> myRdd.collect()  
['Alice', 'Carlos', 'Frank', 'Barbara']  
  
> myRdd.sortBy(lambda x: x[0]).collect()  
['Alice', 'Barbara', 'Carlos', 'Frank']
```



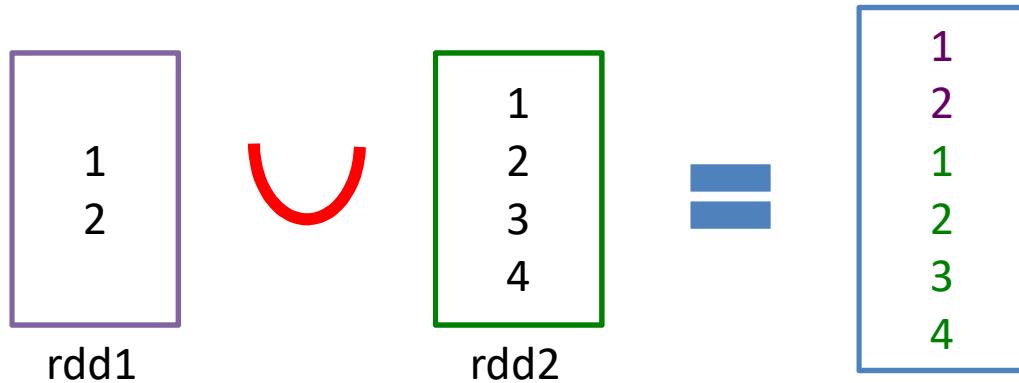
```
In [1]: sc.parallelize([2,3,4]).map(lambda x:range(1,x)).collect()  
Out[1]: [[1], [1, 2], [1, 2, 3]]
```



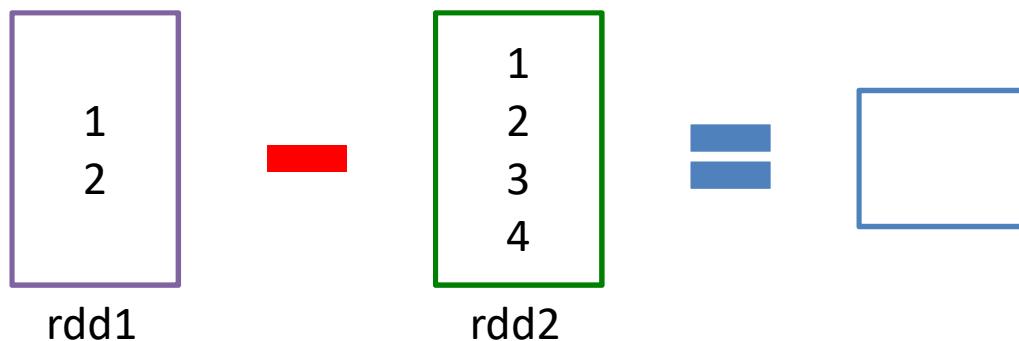
```
In [2]: sc.parallelize([2,3,4]).flatmap(lambda x:range(1,x)).collect()  
Out[2]: [1, 1, 2, 1, 2, 3]
```



```
In [1]: rdd1 = sc.parallelize(range(1,3))
In [2]: rdd2 = sc.parallelize(range(1,5))
In [3]: rdd1.union(rdd2).collect()
Out[3]: [1, 2, 1, 2, 3, 4]
```

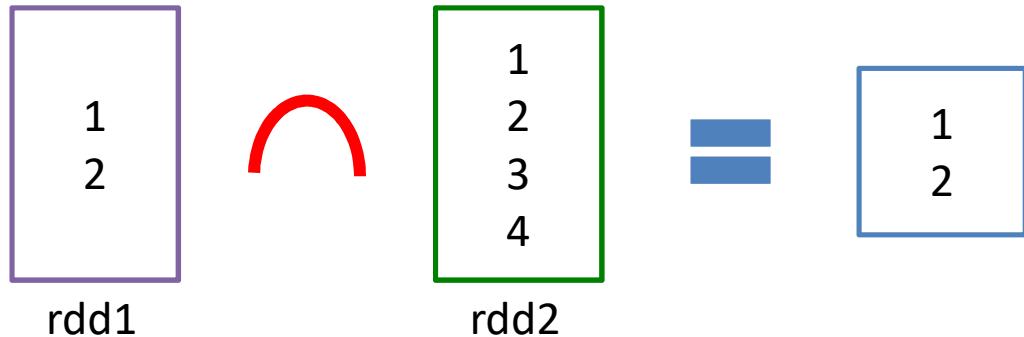


```
In [4]: rdd1.subtract(rdd2).collect()
Out[4]: []
In [5]: rdd2.subtract(rdd1).collect()
Out[5]: [3,4]
```



```
In [6]: rdd1.intersection(rdd2).collect()
```

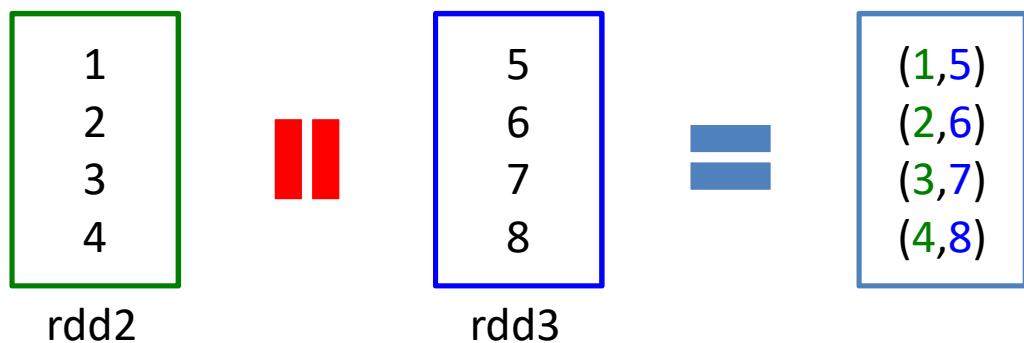
```
Out[6]: [1, 2, 1, 2, 3, 4]
```



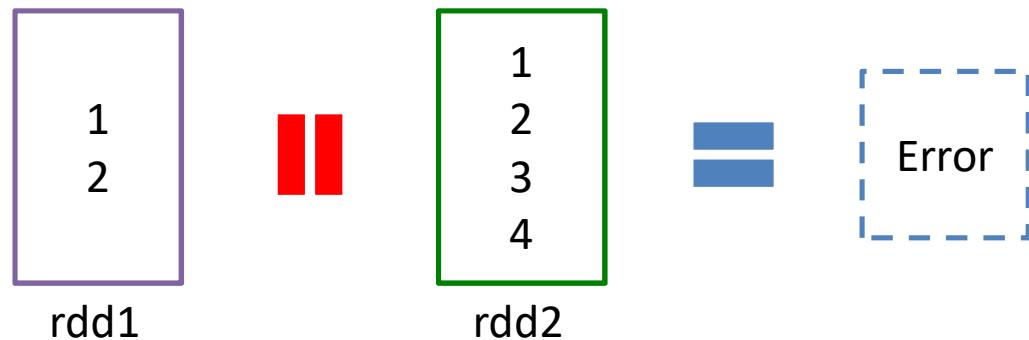
```
In [7]: rdd3=sc.parallelize(range(5,9))
```

```
In [8]: rdd2.zip(rdd3).collect()
```

```
Out[8]: [(1, 5), (2, 6), (3, 7), (4, 8)]
```



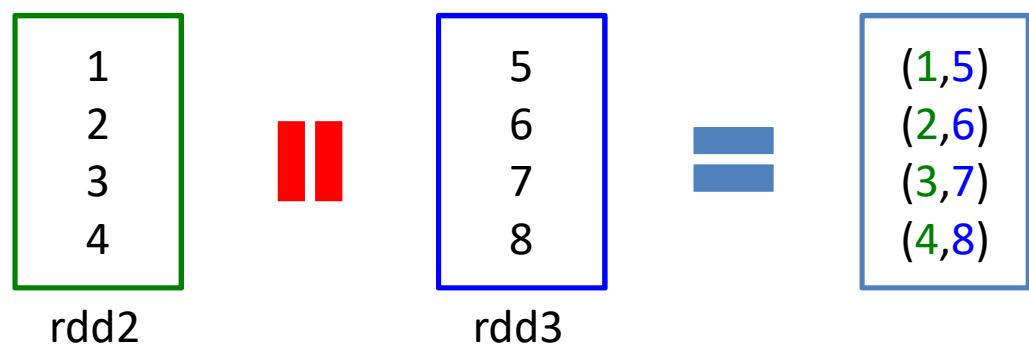
```
In [9]: rdd1.zip(rdd2).collect()
```



```
In [10]: rdd3=sc.parallelize(range(5,9))
```

```
In [11]: rdd2.zip(rdd3).collect()
```

```
Out[11]: [(1, 5), (2, 6), (3, 7), (4, 8)]
```



```
In [12]: rdd4 = rdd1.union(rdd2)
```

```
In [12]: rdd4.collect()
```

```
Out[12]: [1, 2, 1, 2, 3, 4]
```

1
2
1
2
3
4

rdd4

```
In [13]: rdd4.distinct.collect()
```

```
Out[13]: [1,2,3,4]
```

```
In [14]: rdd1.cartesian(rdd2).collect()
```

```
Out[14]: [(1, 1), (1, 2), (1, 3), (1, 4), (2, 1), (2, 2), (2, 3), (2, 4)]
```

1
2



rdd1

1
2
3
4

rdd2



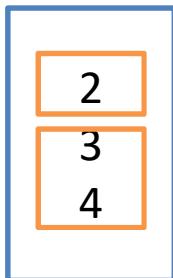
(1,1)
(1,2)
(1,3)
(1,4)
(2,1)
(2,2)
(2,3)
(2,4)

```
In [15]: sc.parallelize([2,3,4]).getNumPartitions()
```

```
Out[15]: 2
```

```
In [16]: sc.parallelize([2,3,4]).glom().collect()
```

```
Out[16]: [[2], [3, 4]]
```



```
In [17]: def f(iterator):
....:     yield sum(iterator)
....:
```

```
In [18]: sc.parallelize([2,3,4]).mapPartitions(f).collect()
```

```
Out[18]: [2, 7]
```



```
$ hdfs dfs -cat purplecow.txt  
I've never seen a purple cow.  
I never hope to see one;  
But I can tell you, anyhow,  
I'd rather see than be one.
```

```
$ pyspark2
```

```
In [1]: mydata = sc.textFile('purplecow.txt')
```

```
I've never seen a purple cow.  
I never hope to see one;  
But I can tell you, anyhow,  
I'd rather see than be one.
```

mydata

```
In [2]: mydata.collect()
```

```
Out[2]:
```

```
[u"I've never seen a purple cow.",  
 u'I never hope to see one;',  
 u'But I can tell you, anyhow,',  
 u"I'd rather see than be one."]
```

```
the cat sat on the mat  
the aardvark sat on the sofa
```

mydata

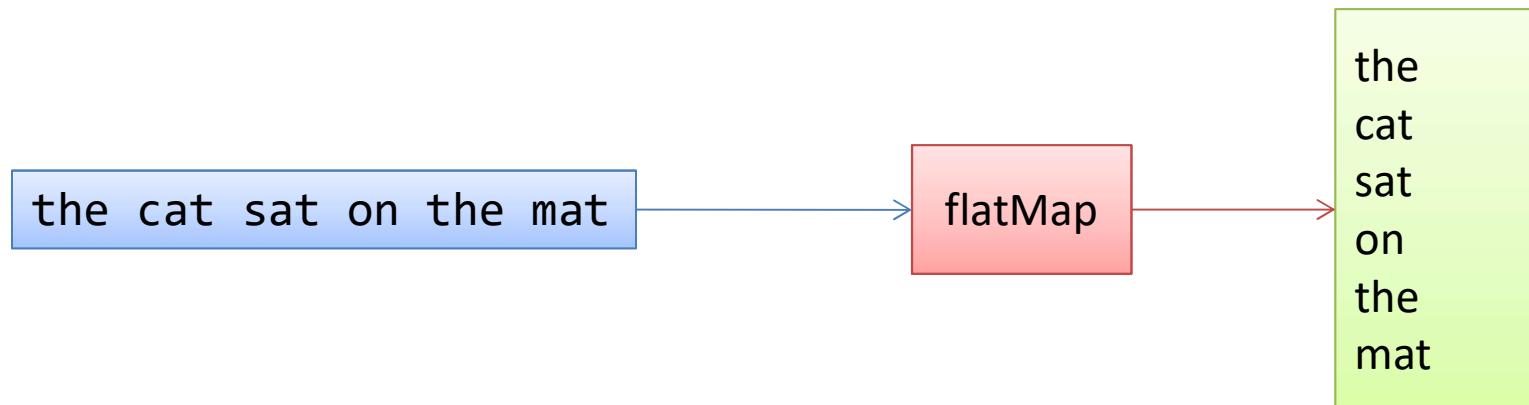
```
In [4]: mydata2 = mydata.map(lambda line: line.split(' '))
```

```
graph LR; A[the cat sat on the mat] --> B[map]; B --> C[[the,Cat,Sat,on ,the,mat]]
```

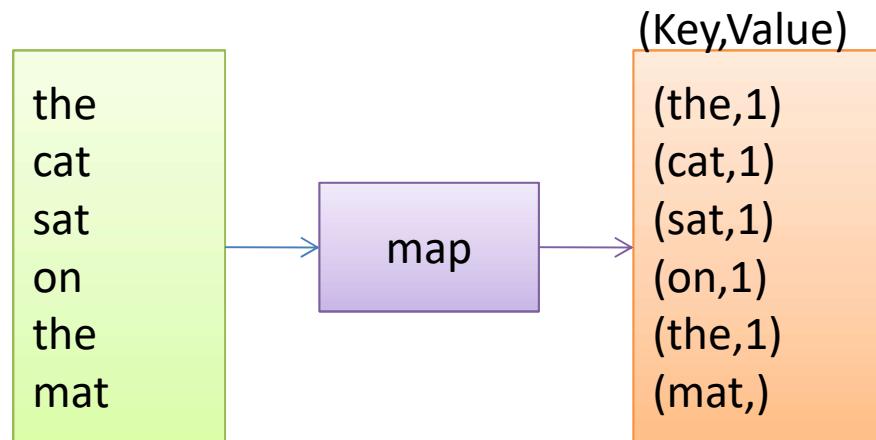
```
the cat sat on the mat  
the aardvark sat on the sofa
```

mydata

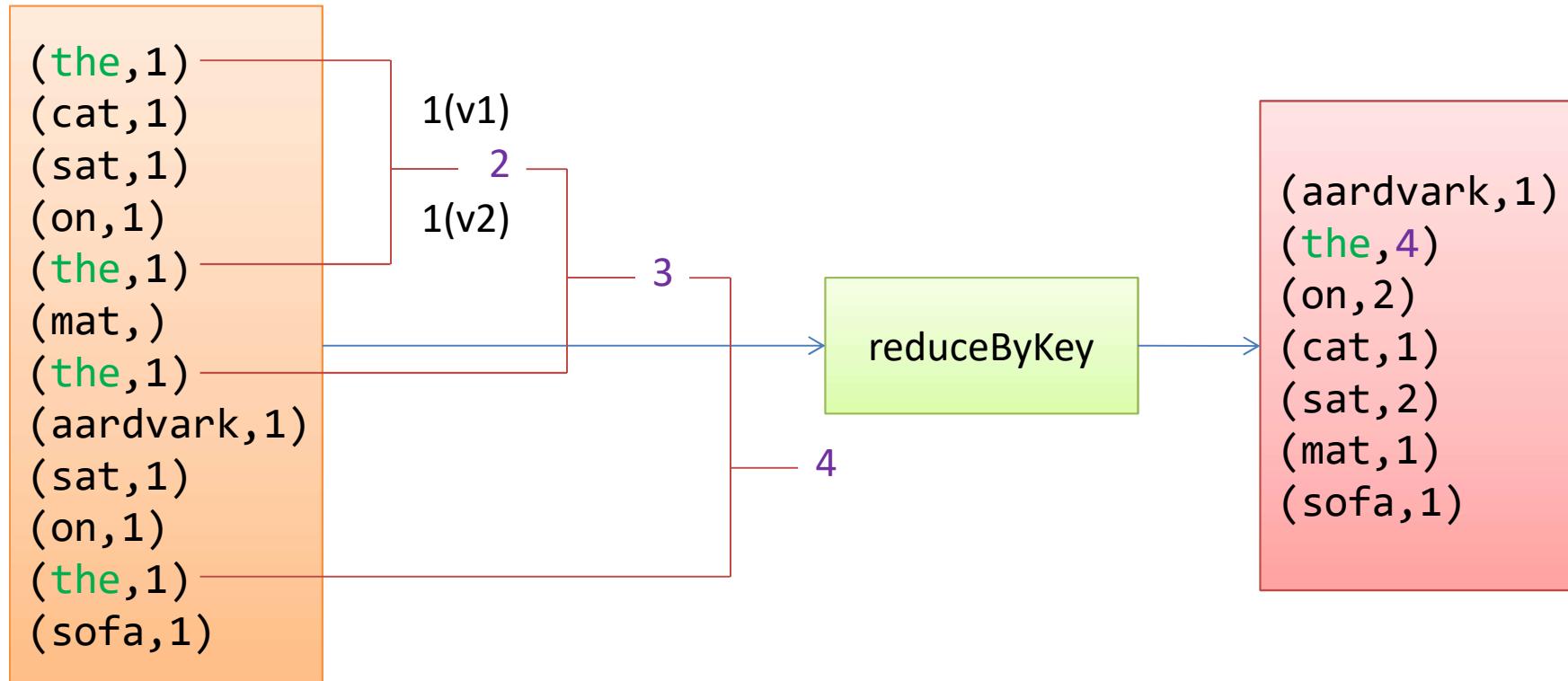
```
In [4]: mydata2 = mydata.flatMap(lambda line: line.split(' '))
```



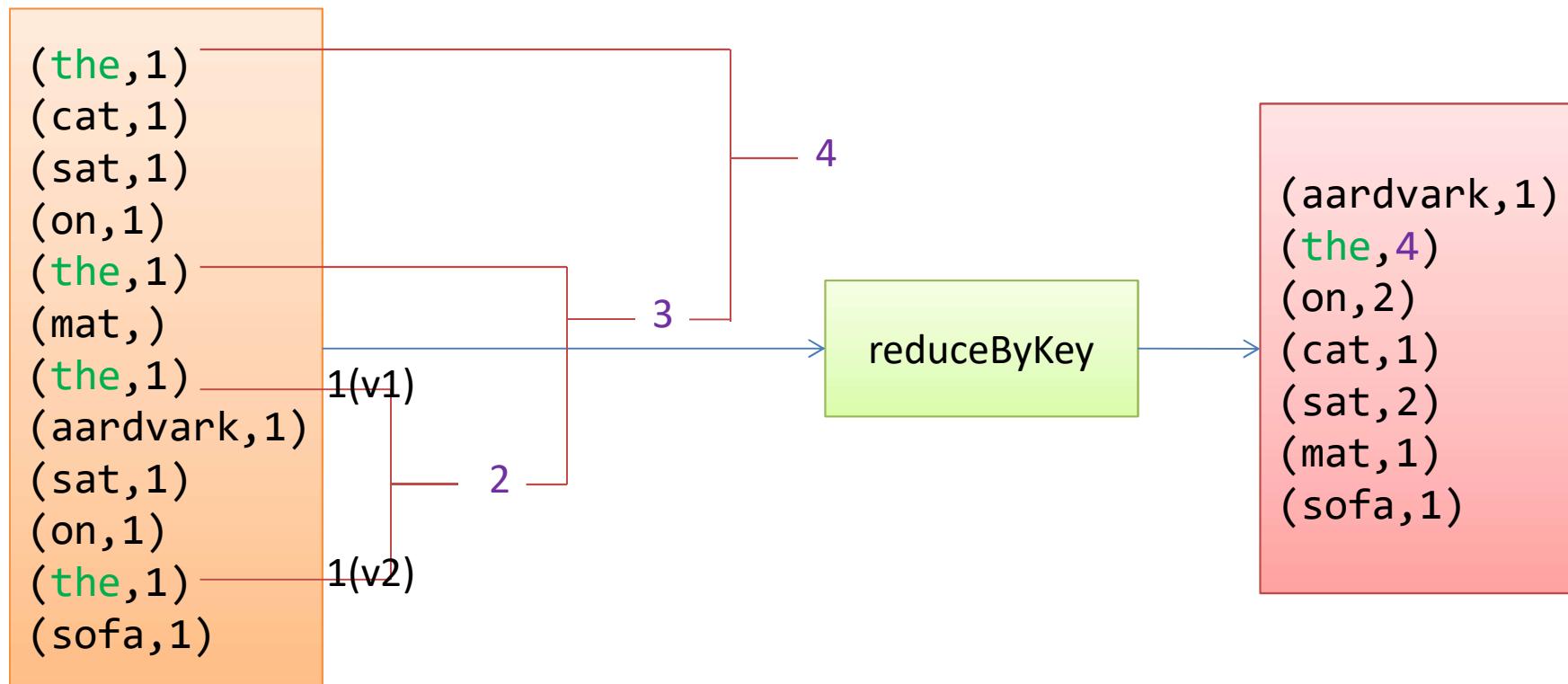
```
In [7]: mydata3 = mydata2.map(lambda word: (word,1))
```



```
In [10]: mydata4 = mydata3.reduceByKey(lambda v1,v2: v1+v2)
```



```
In [10]: mydata4 = mydata3.reduceByKey(lambda v1,v2: v1+v2)
```



```
user001\tFred Flintstone\n
user090\tBugs Bunny\n
user111\tHarry Potter\n
```

```
> users = sc.textFile(file) \
    .map(lambda line: line.split('\t')) \
    .map(lambda fields: (fields[0], fields[1]))
```

```
user001\tFred Flintstone
[user001,Fred Flintstone]
(user001,Fred Flintstone)
```

rdd1m

(100,Frank)
(200,Scott)
(300,Linda)

rdd2m

(100,1000)
(300,5000)
(400,20000)

rdd1m.**join**(rdd2m)

(300,(Linda,5000))
(100,(Frank,10000))

rdd1m.**rightOuterJoin**(rdd2m)

(300,(Linda,5000))
(100,(Frank,10000))
(400,(None,20000))

rdd1m.**leftOuterJoin**(rdd2m)

(300,(Linda,5000))
(100,(Frank,10000))
(200,(Scott,None))

Rdd1m.**fullOuterJoin**(rdd2m)

(300,(Linda,5000))
(100,(Frank,10000))
(200,(Scott,None))
(400,(None,20000))

0

1

2

3

4

5

```
56.38.234.188 - 99788 "GET /KBDOC-00157.html HTTP/1.0"
```

```
> sc.textFile(logfile).keyBy(lambda line: line.split(' ')[2])
```

```
(99788,56.38.234.188 - 99788 "GET /KBDOC-00157.html HTTP/1.0")
```

0

1

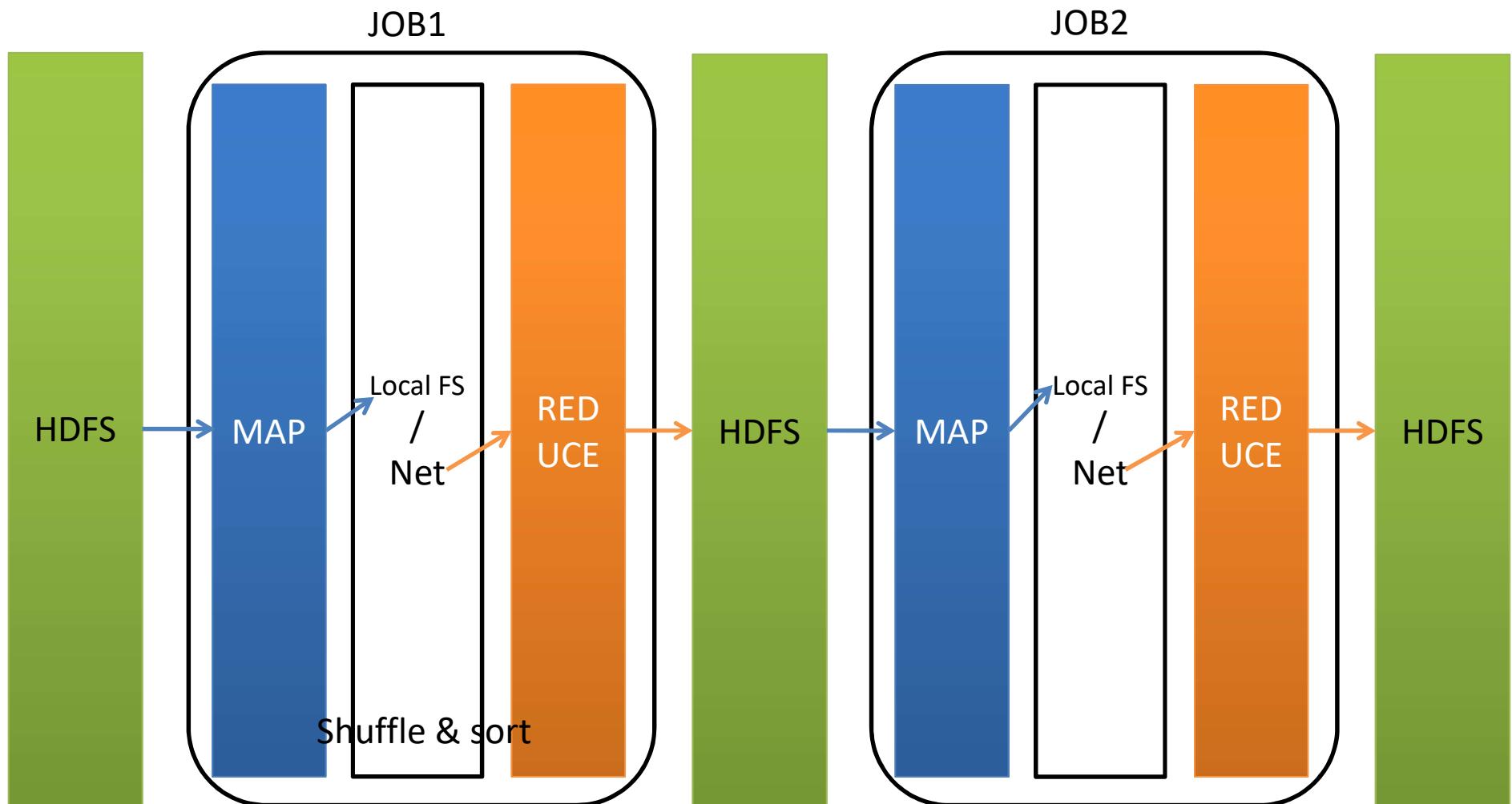
2

```
00210 43.005895 -71.013202
```

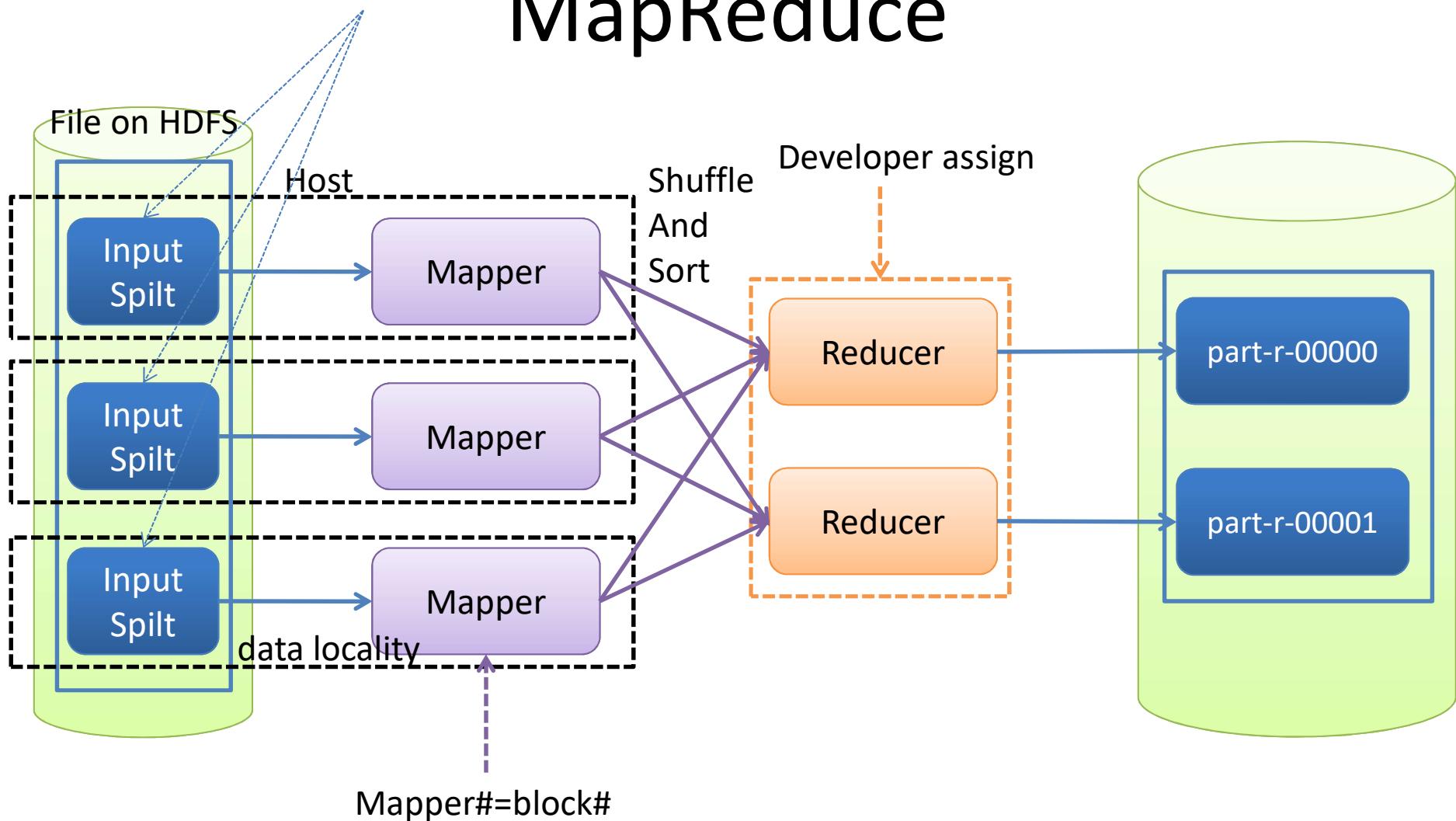
```
> sc.textFile(file).map(lambda line: line.split('')).\n    map(lambda fields: (fields[0],(fields[1],fields[2])))
```

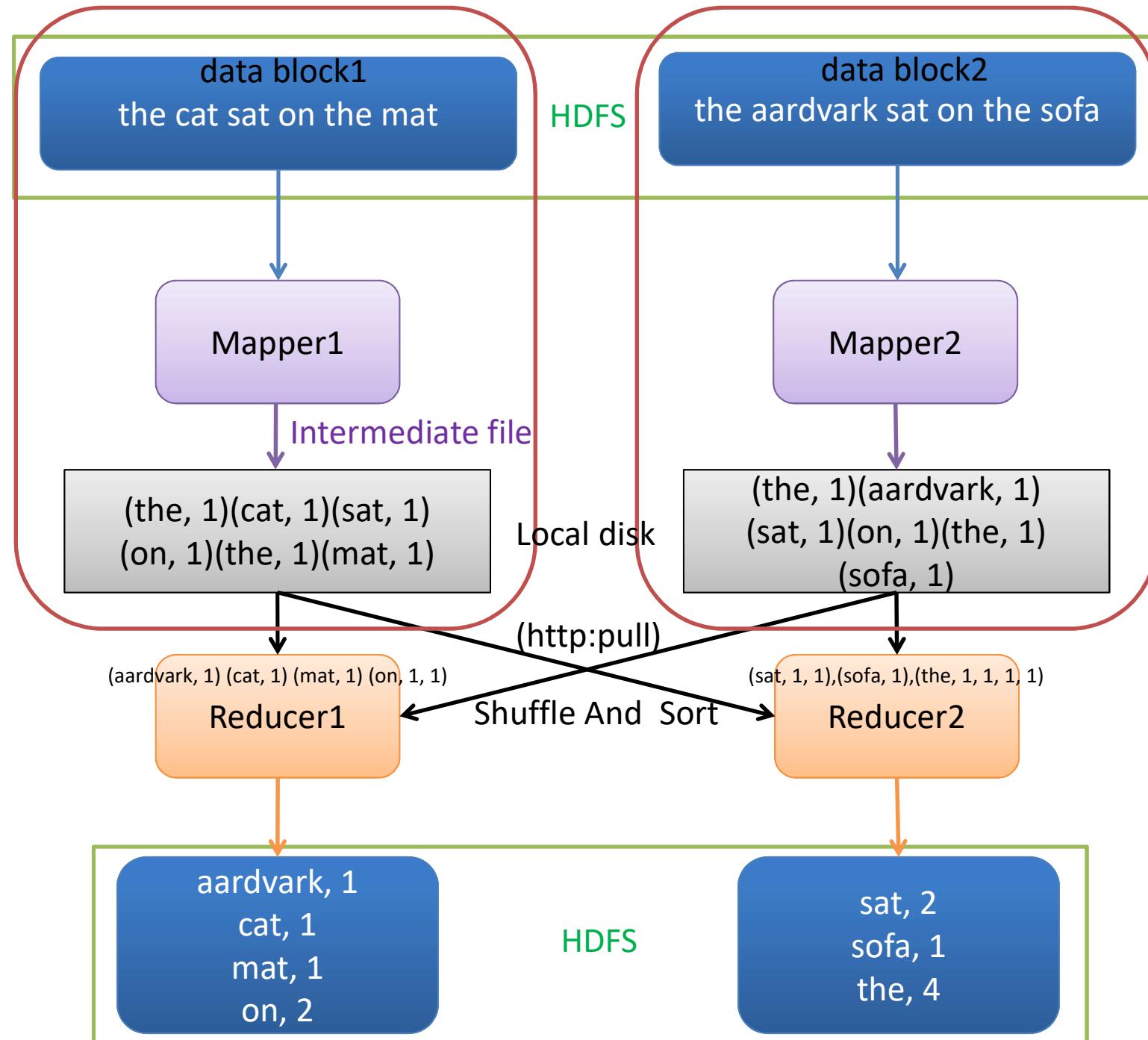
```
(00210,(43.005895,-71.013202))
```

# MapReduce

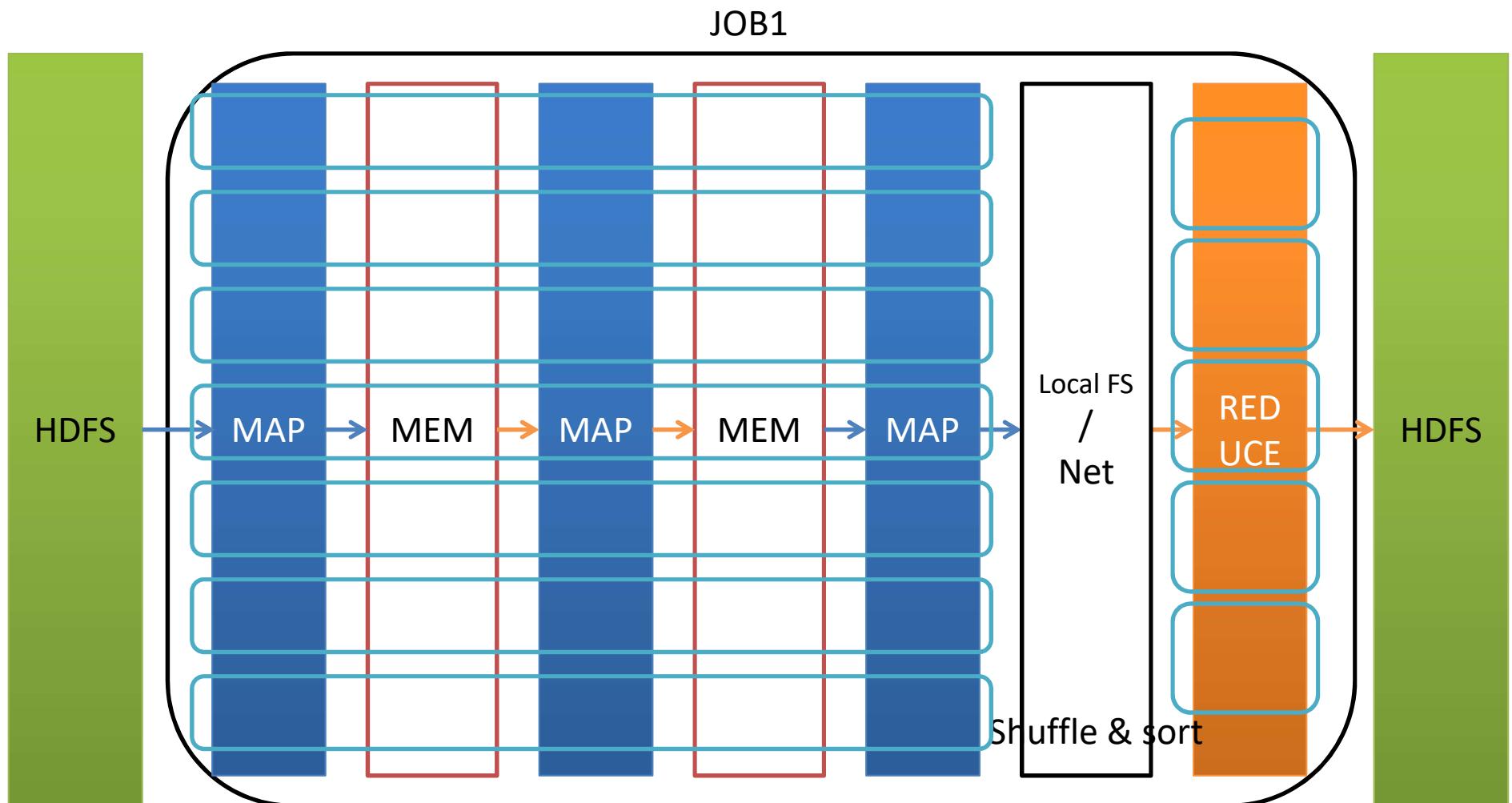


# MapReduce



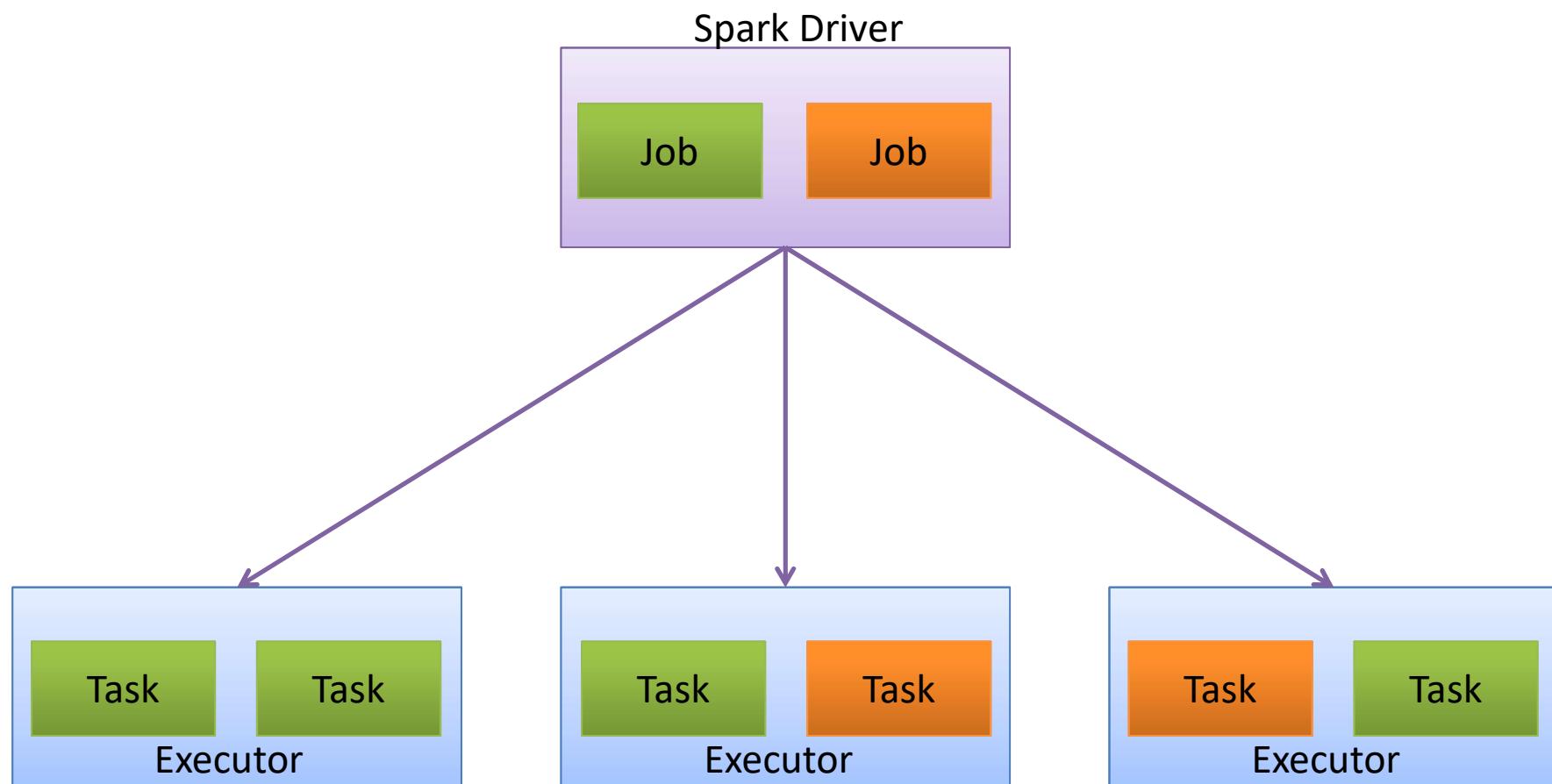


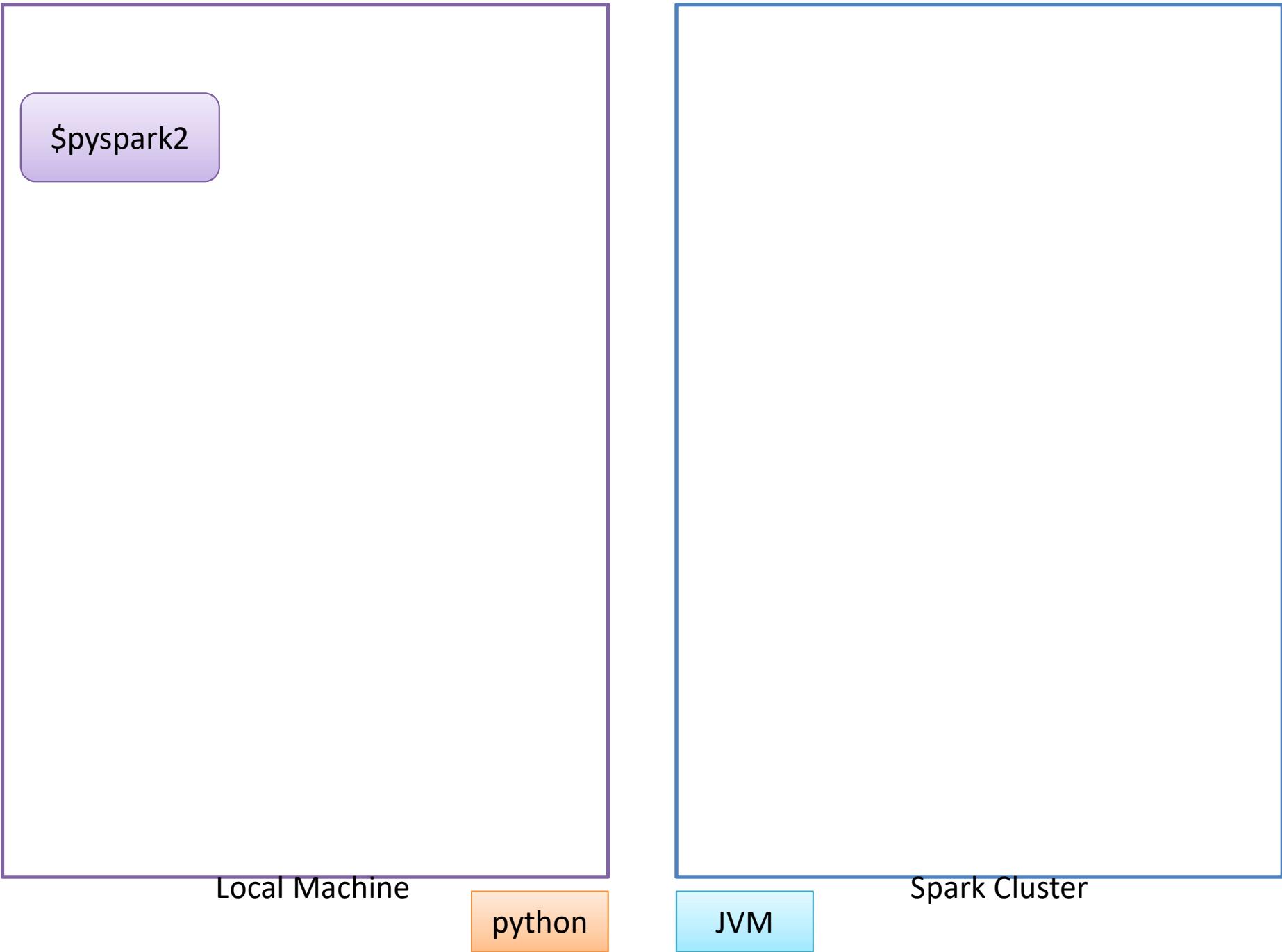
# Spark

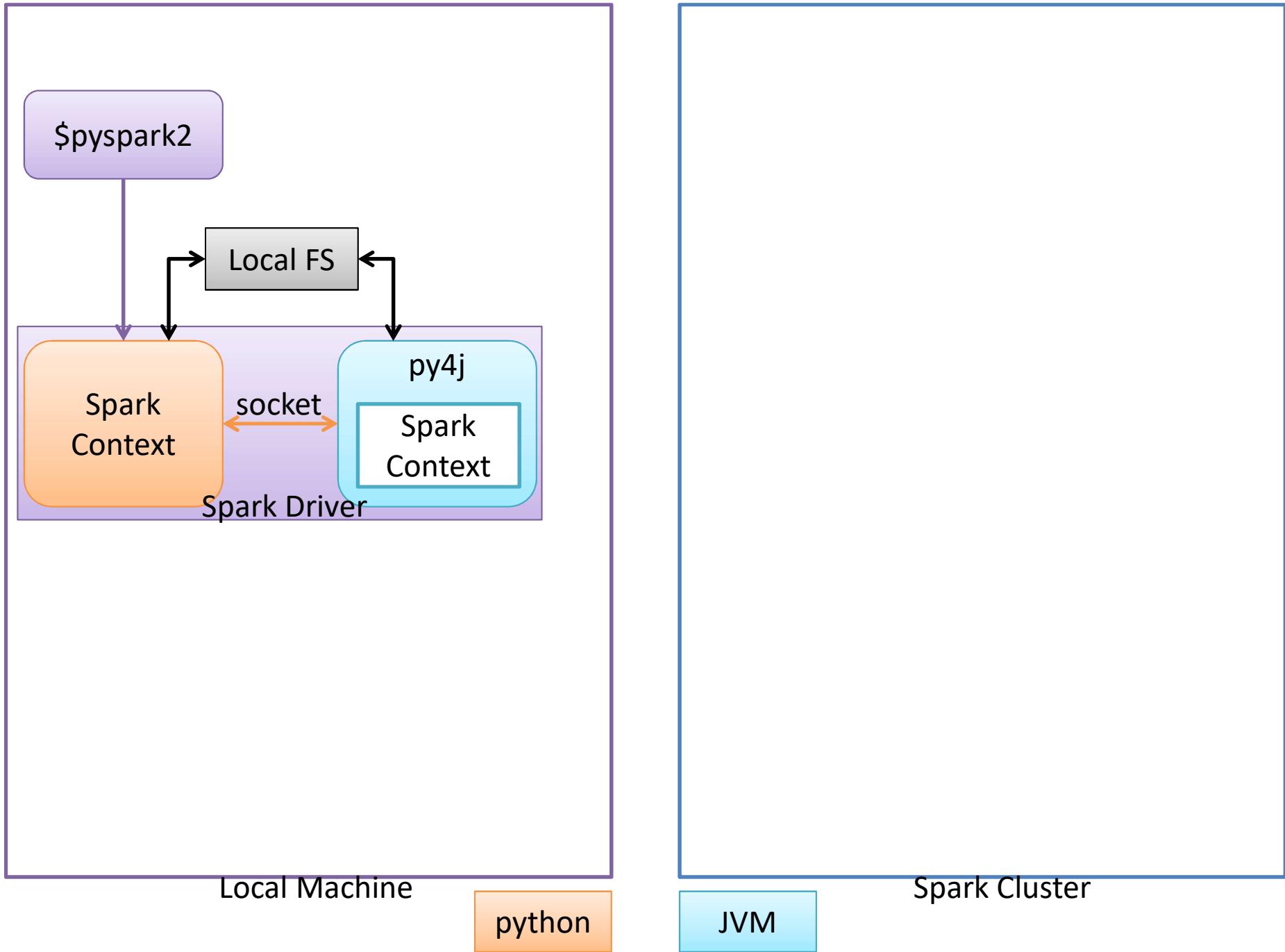


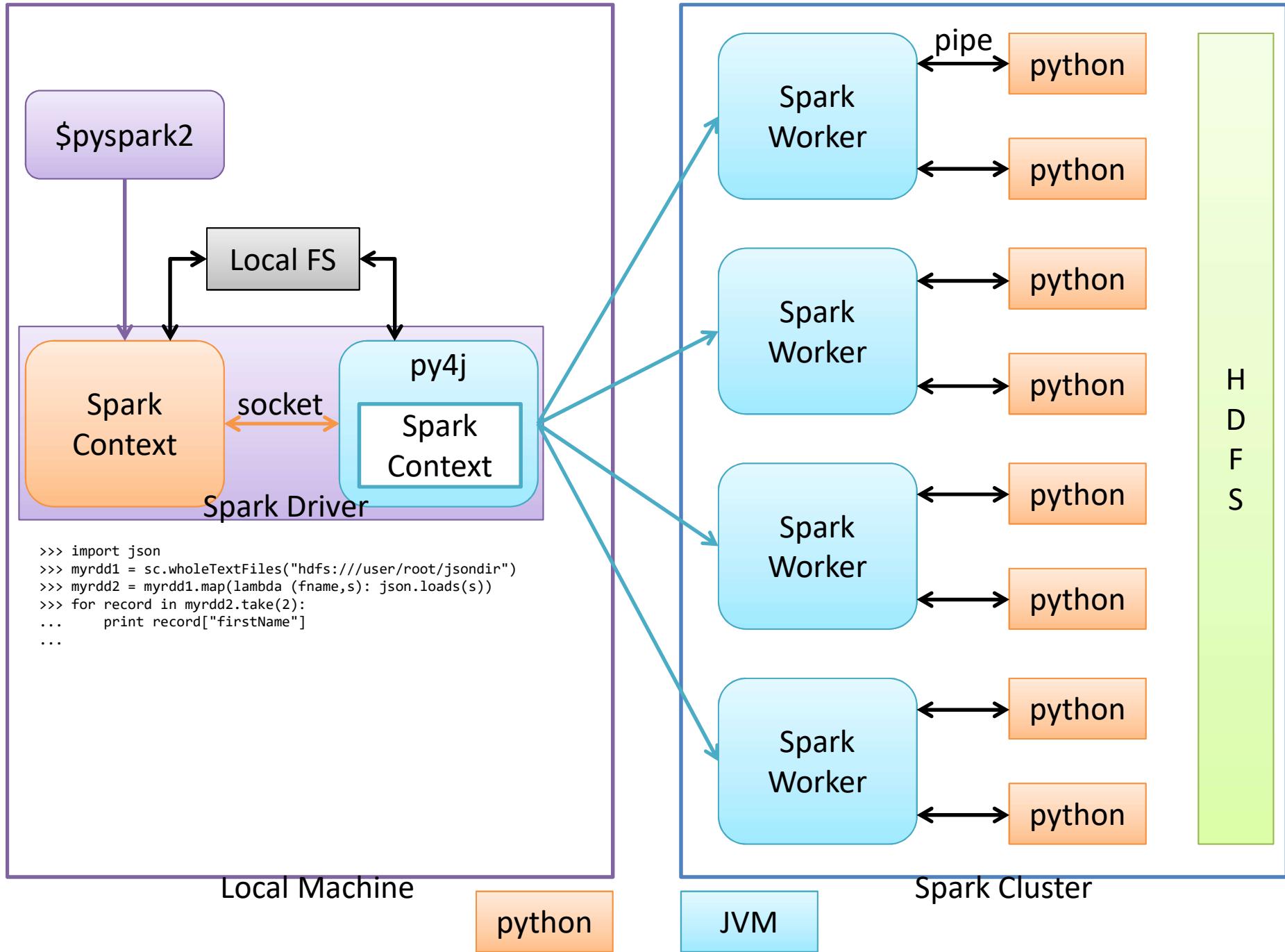
# Spark Application

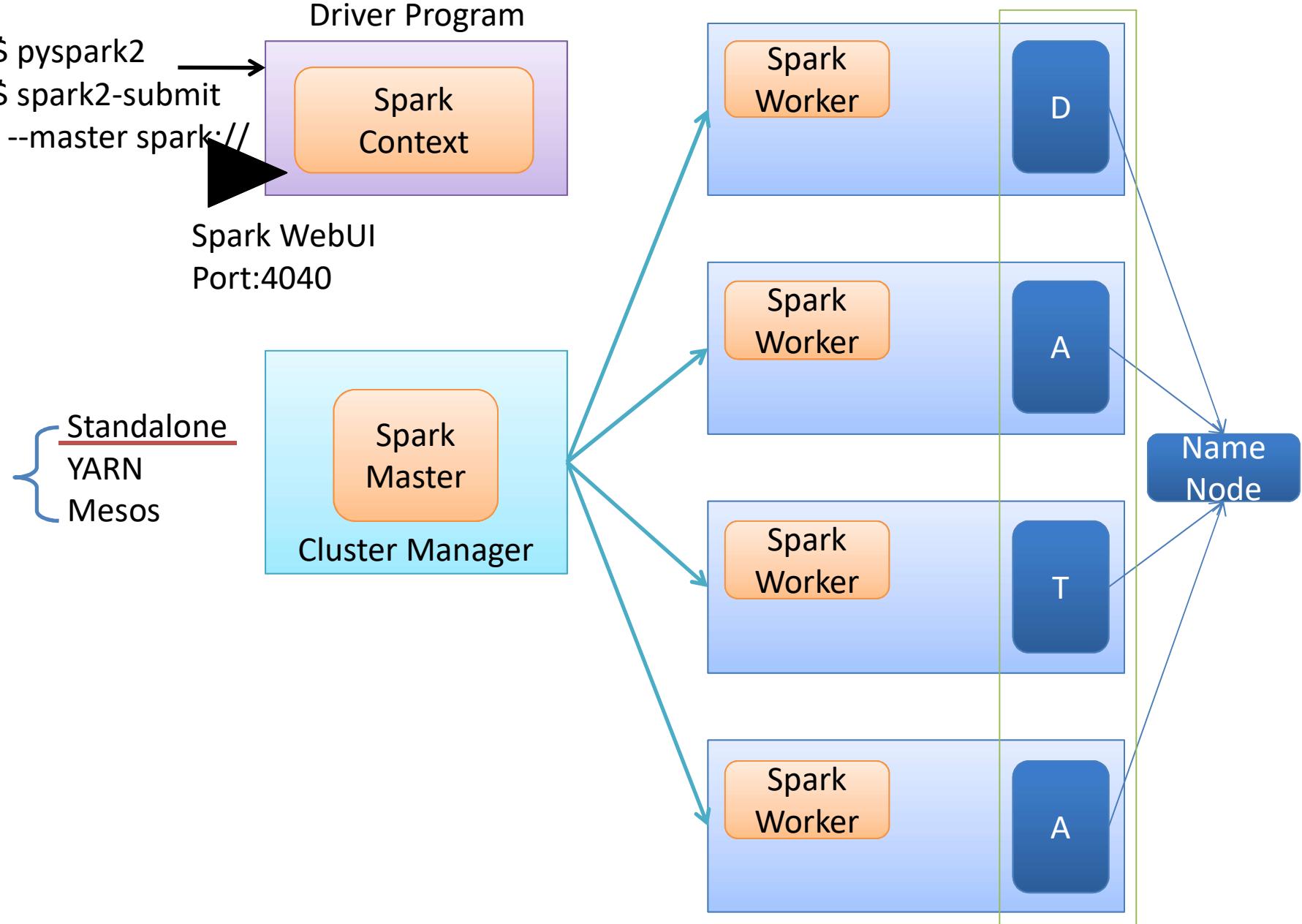
- 只能有一個sparkContext
- 可以一個或多個sparkSession

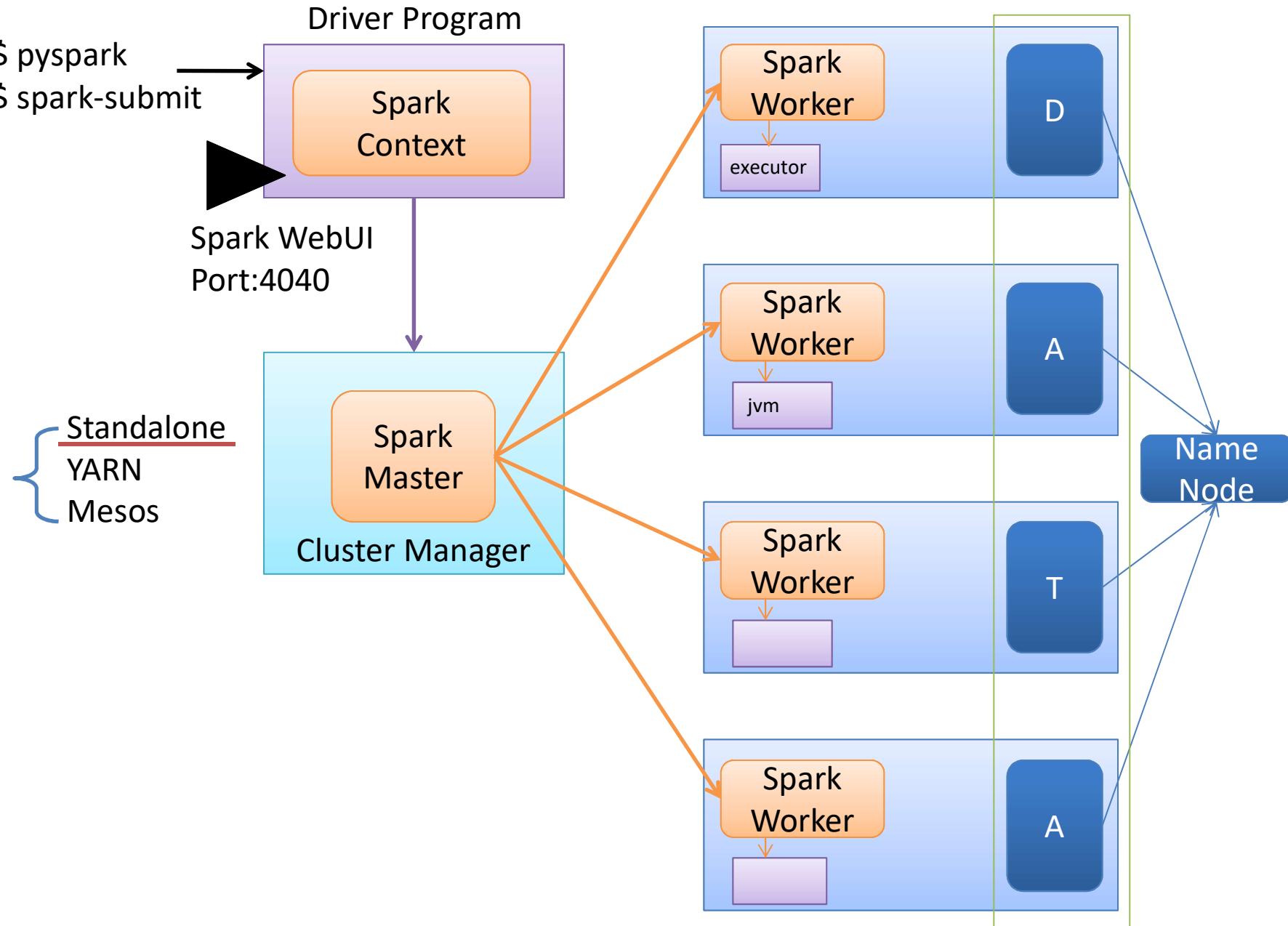


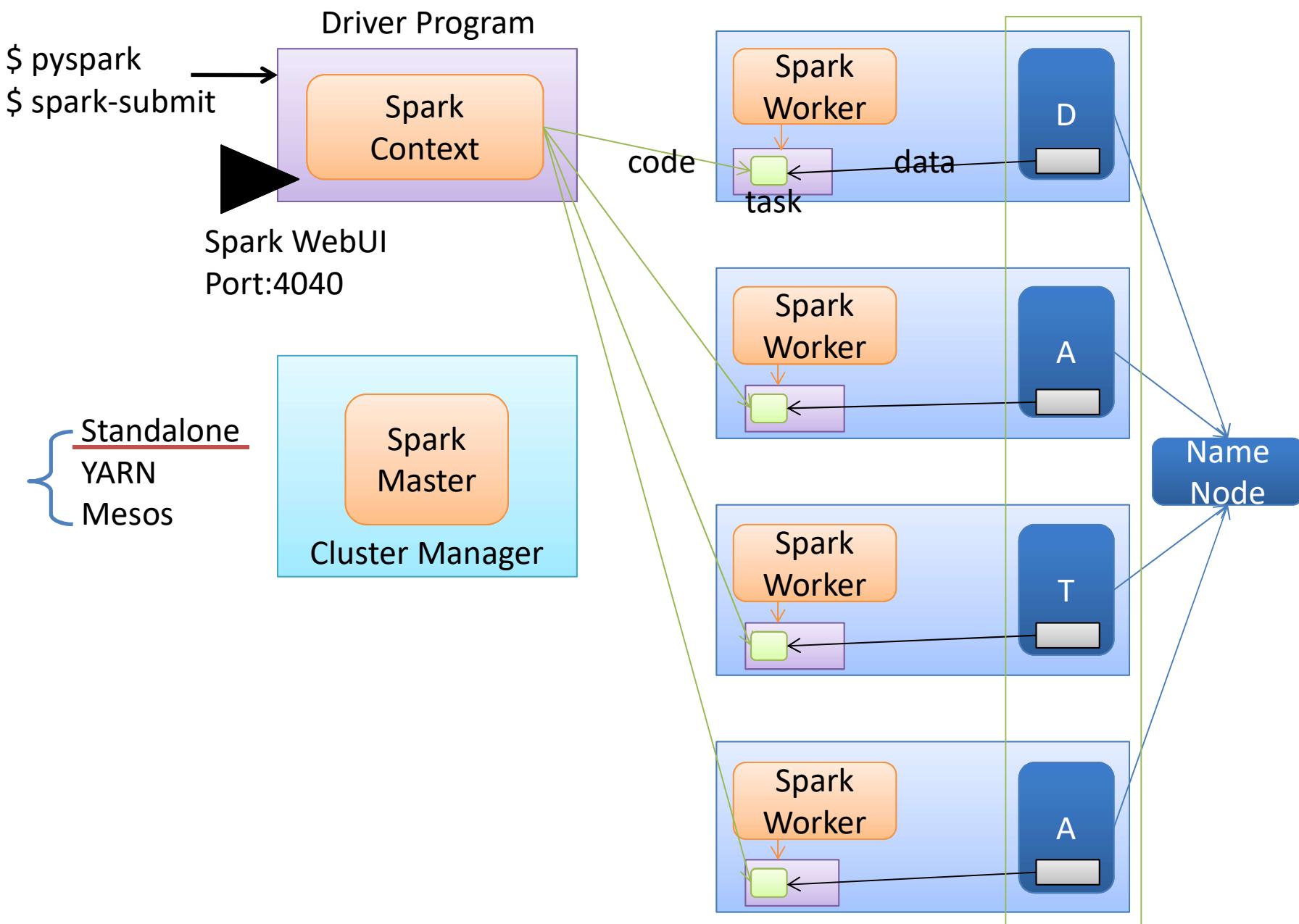


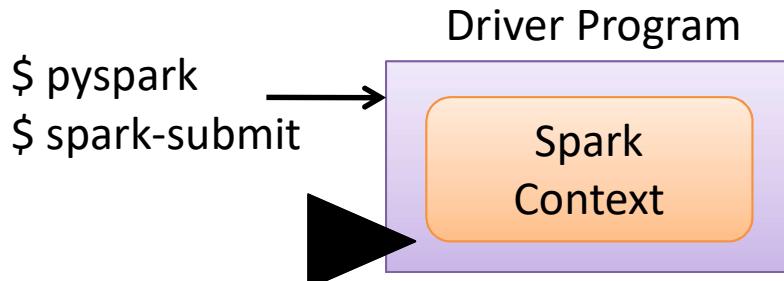




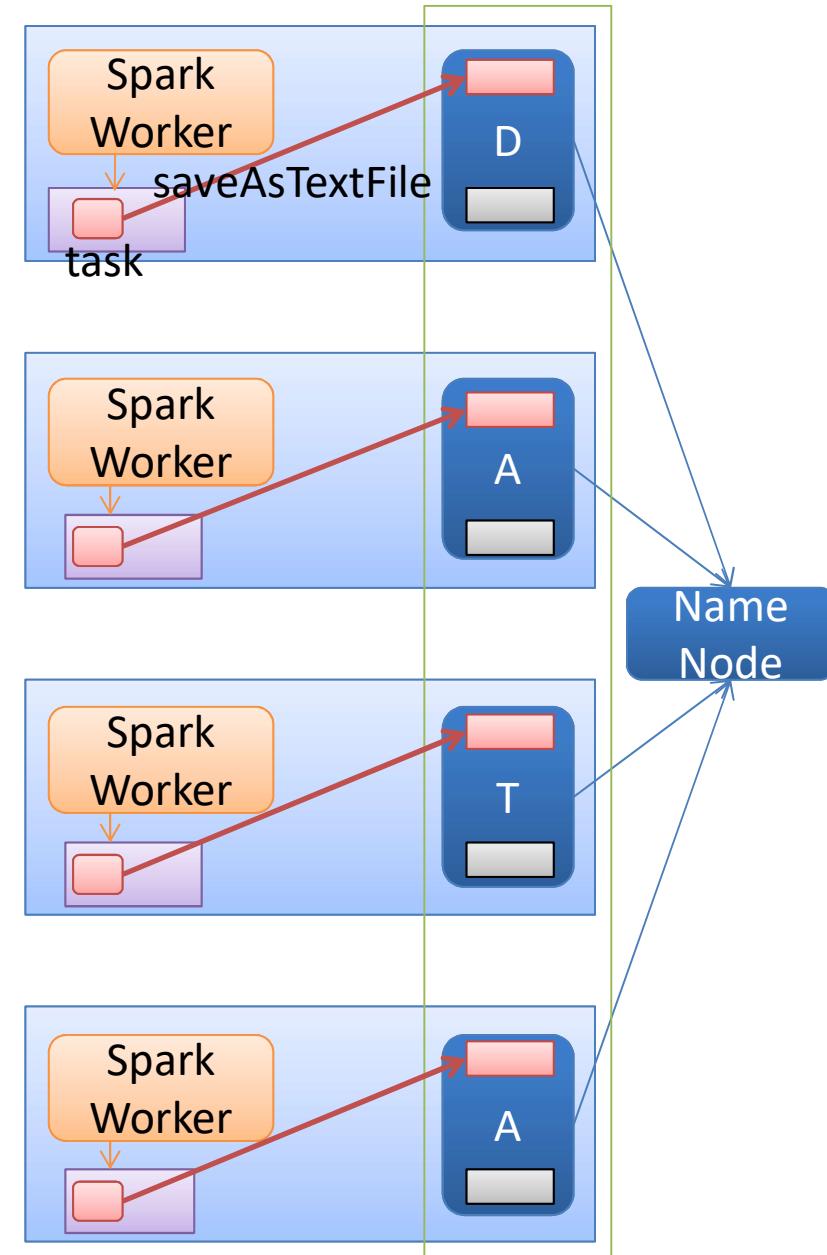


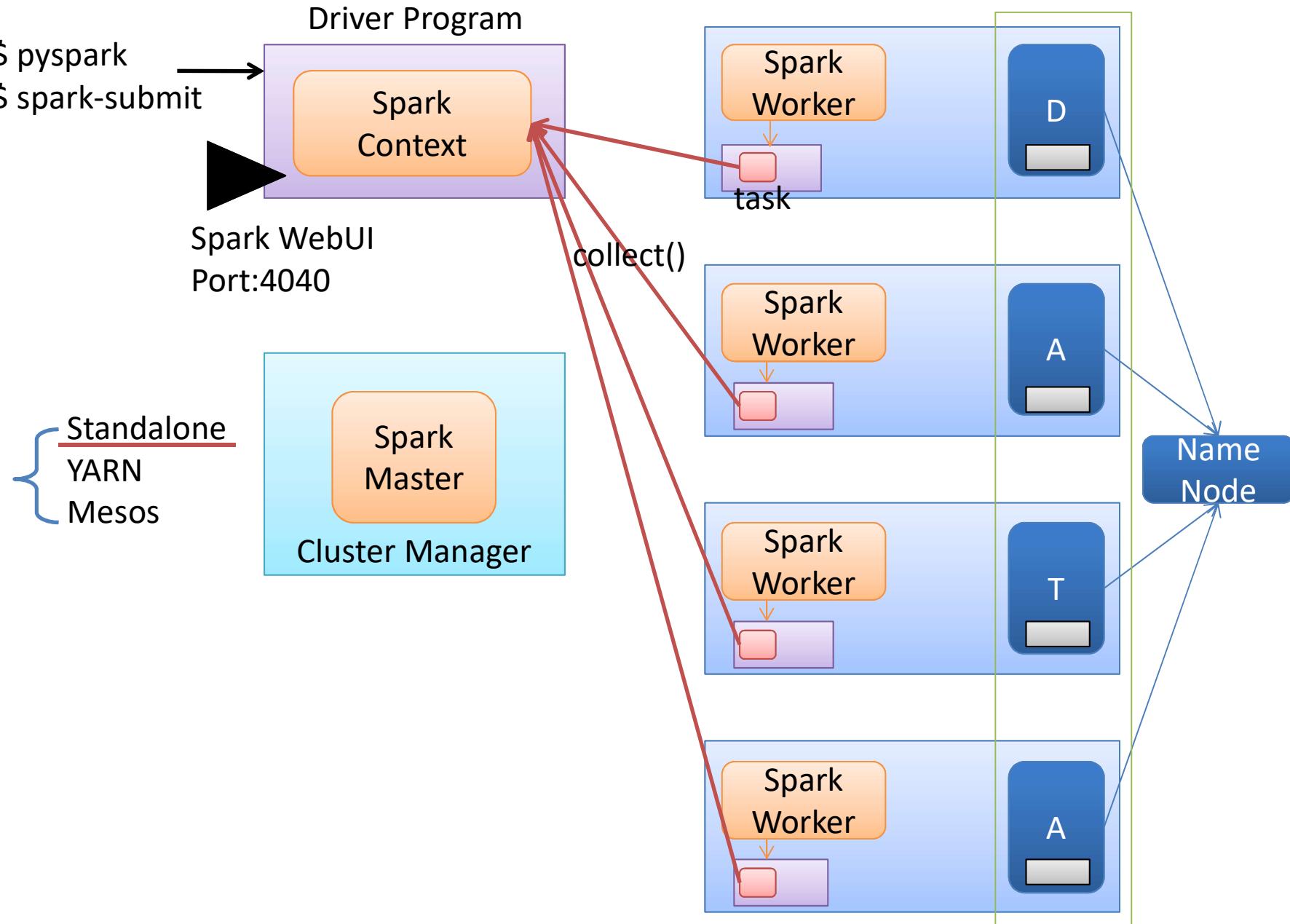


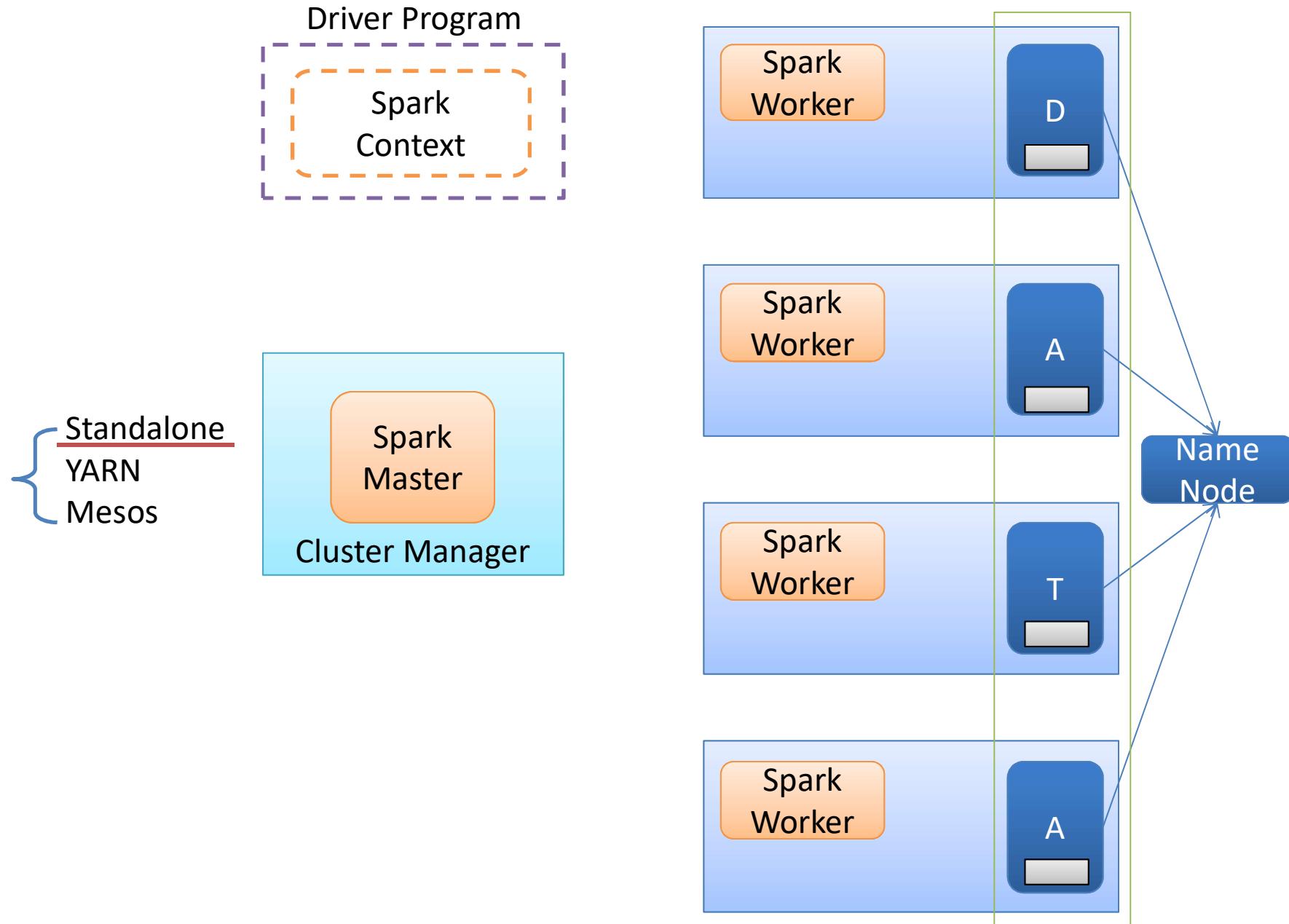


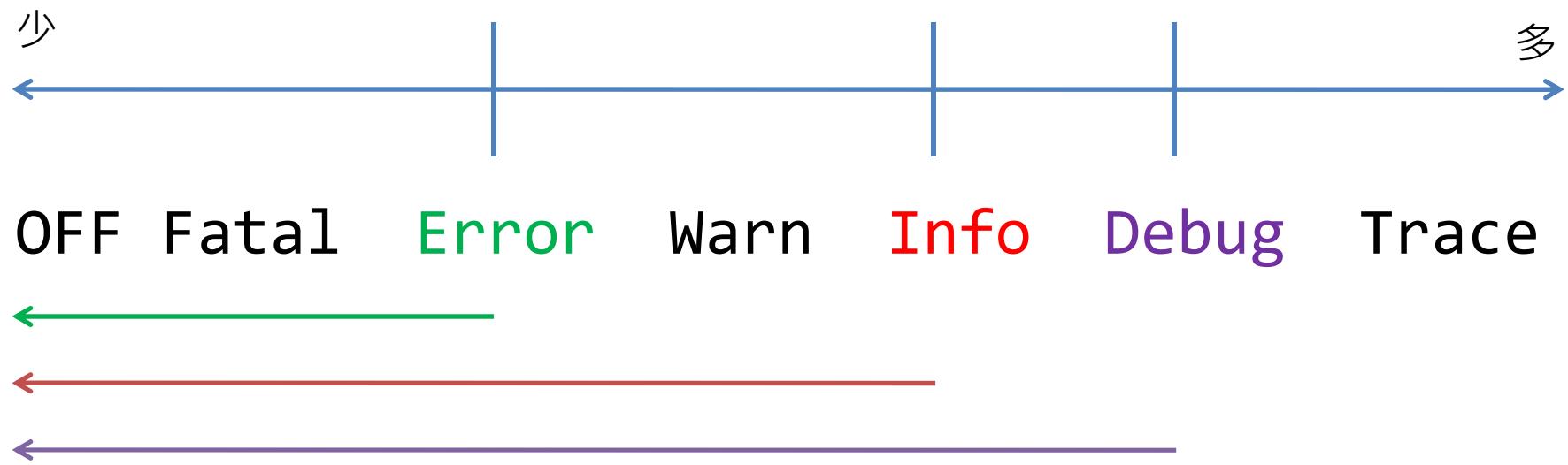


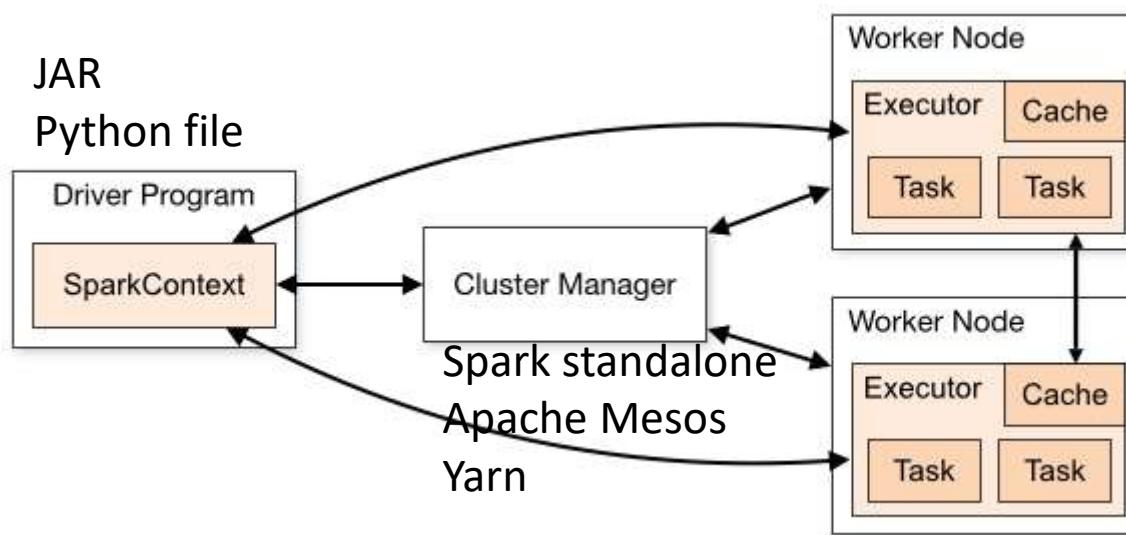
Spark WebUI  
Port:4040

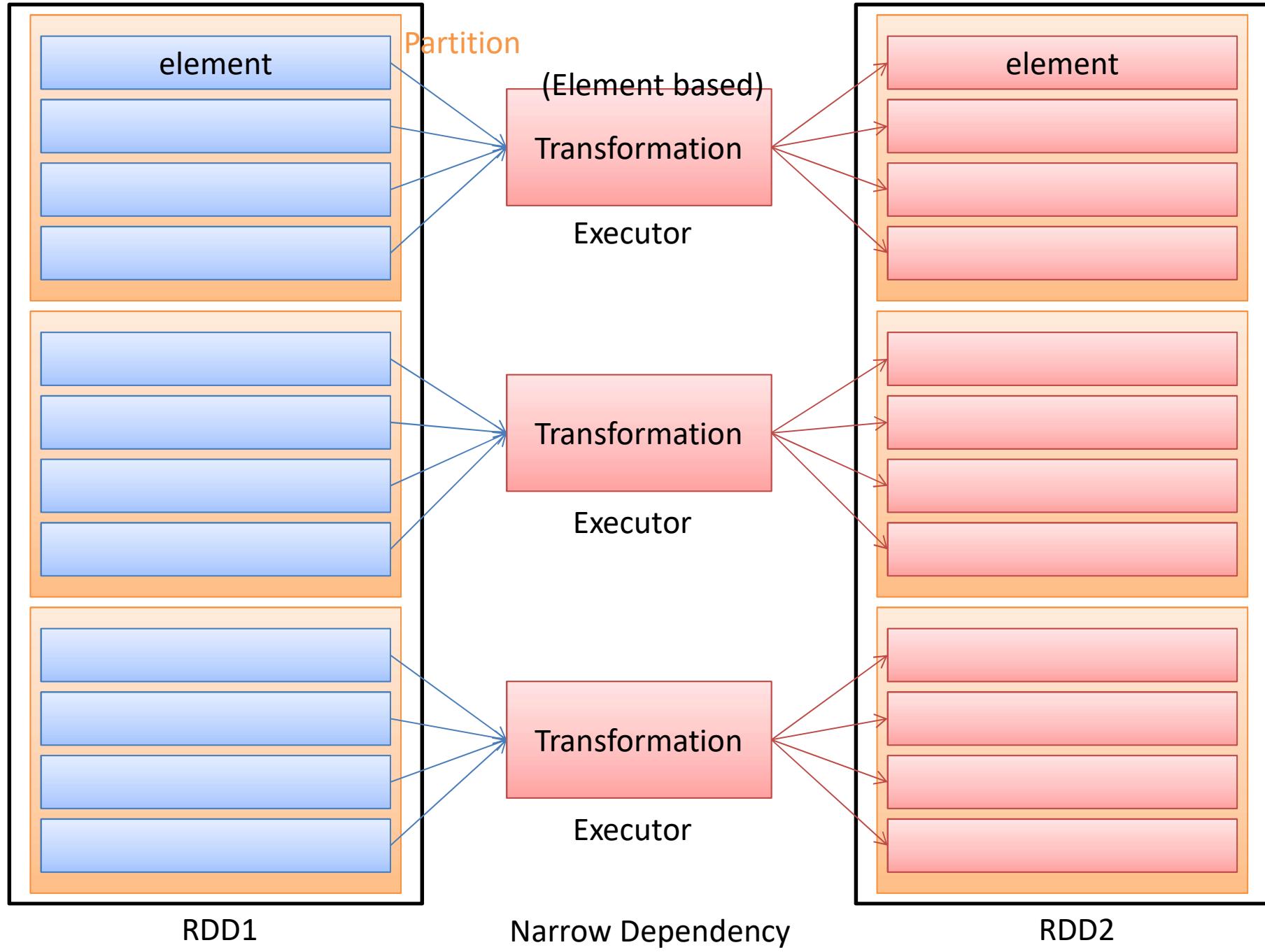


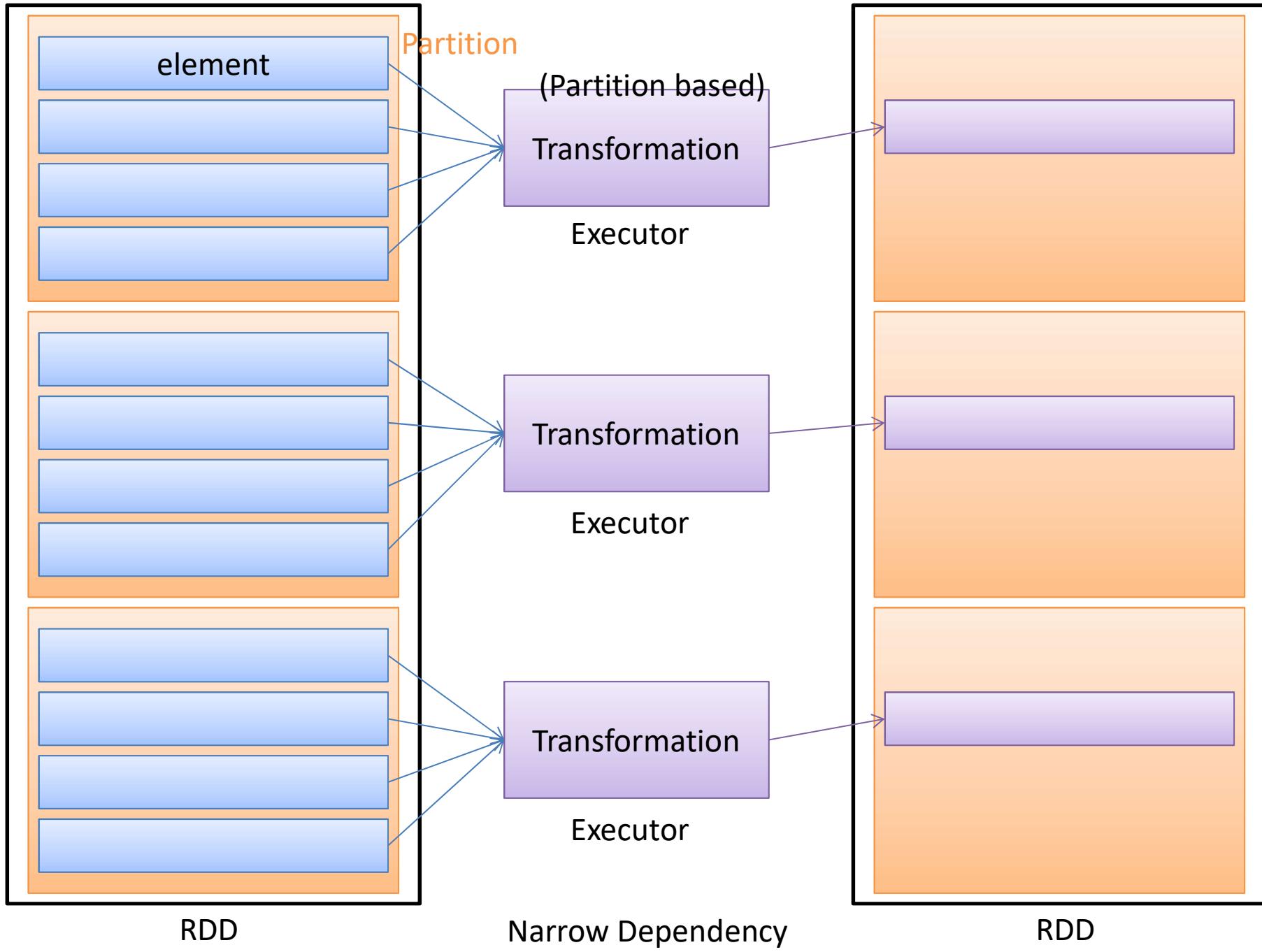


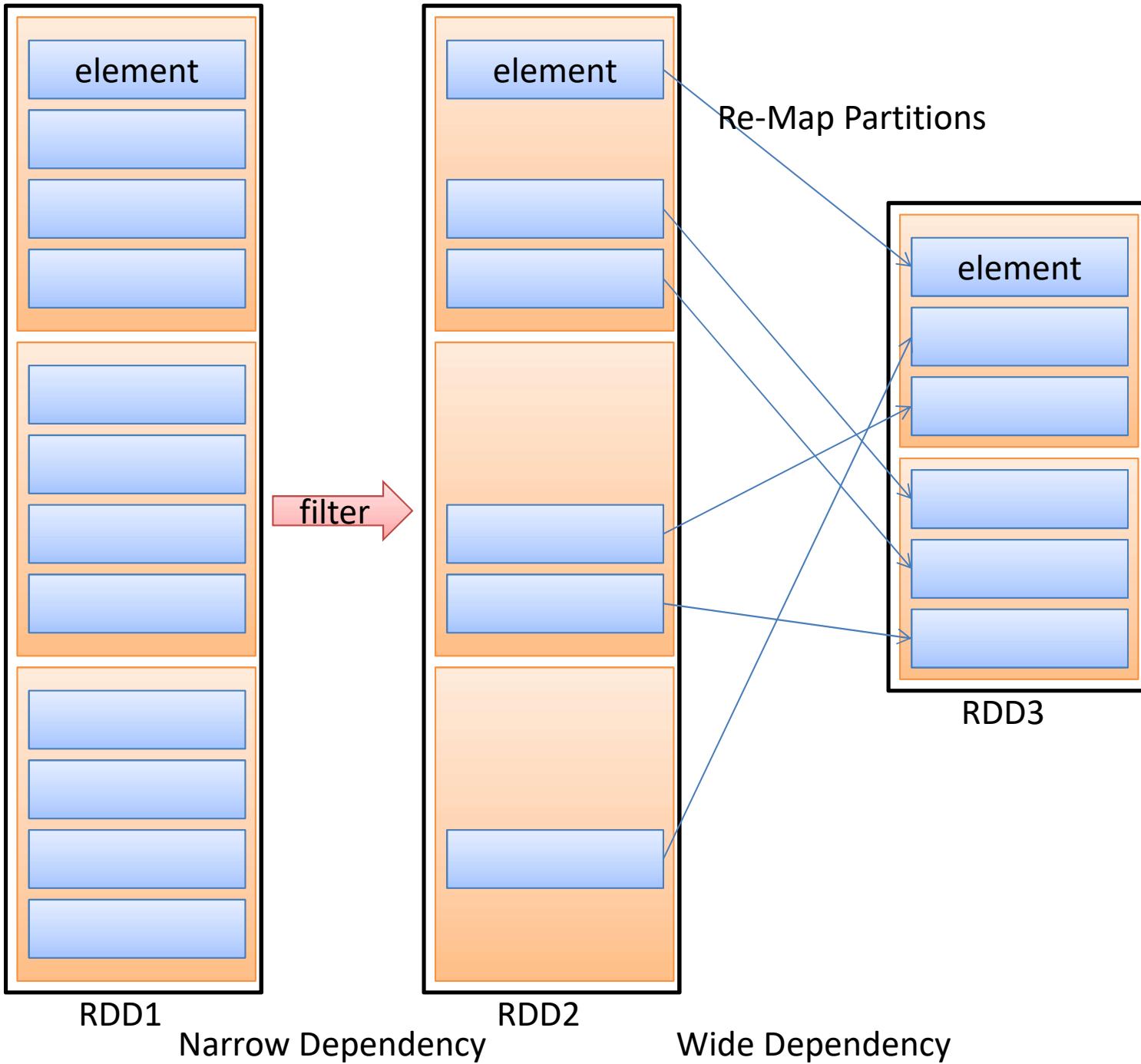


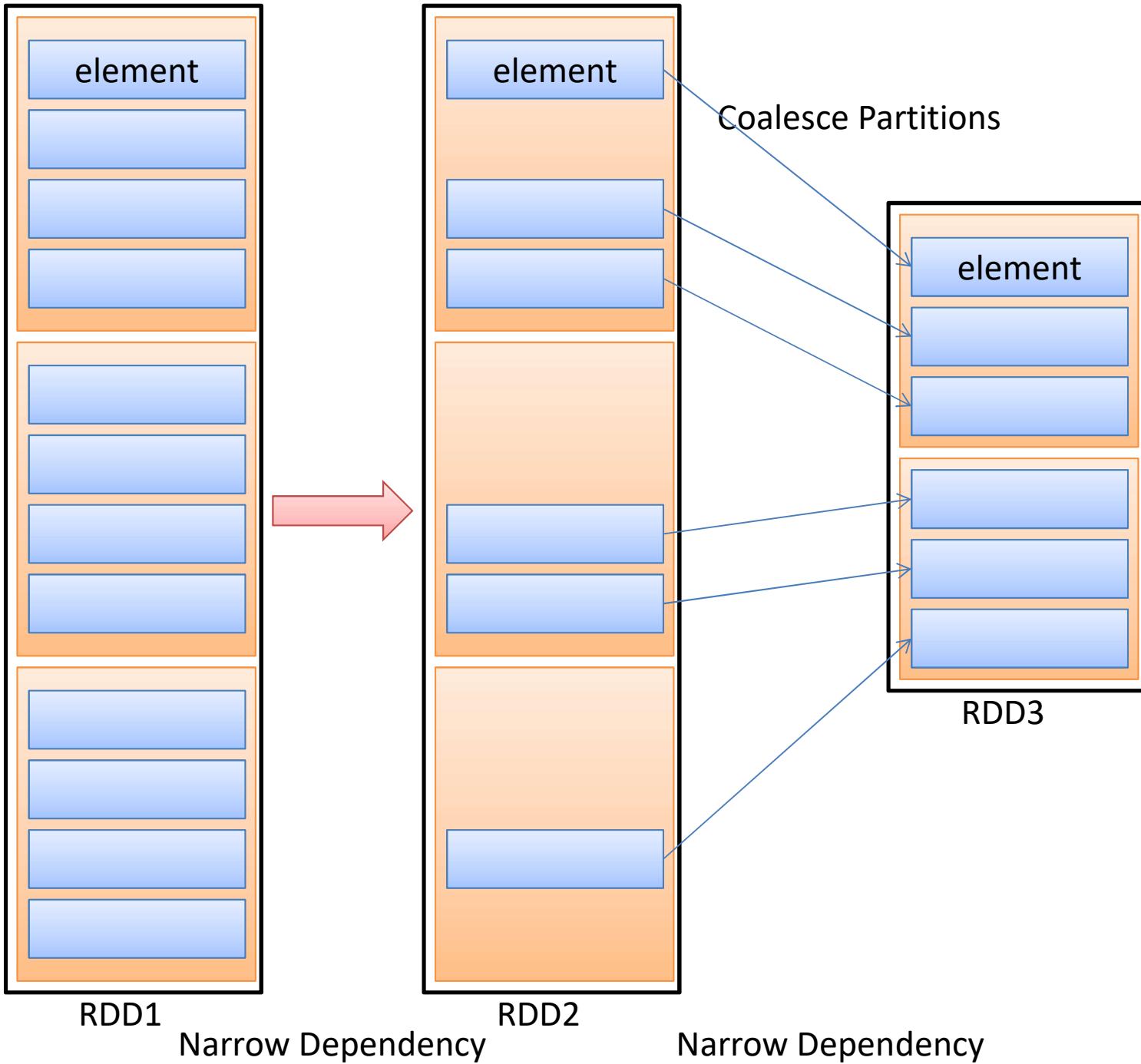


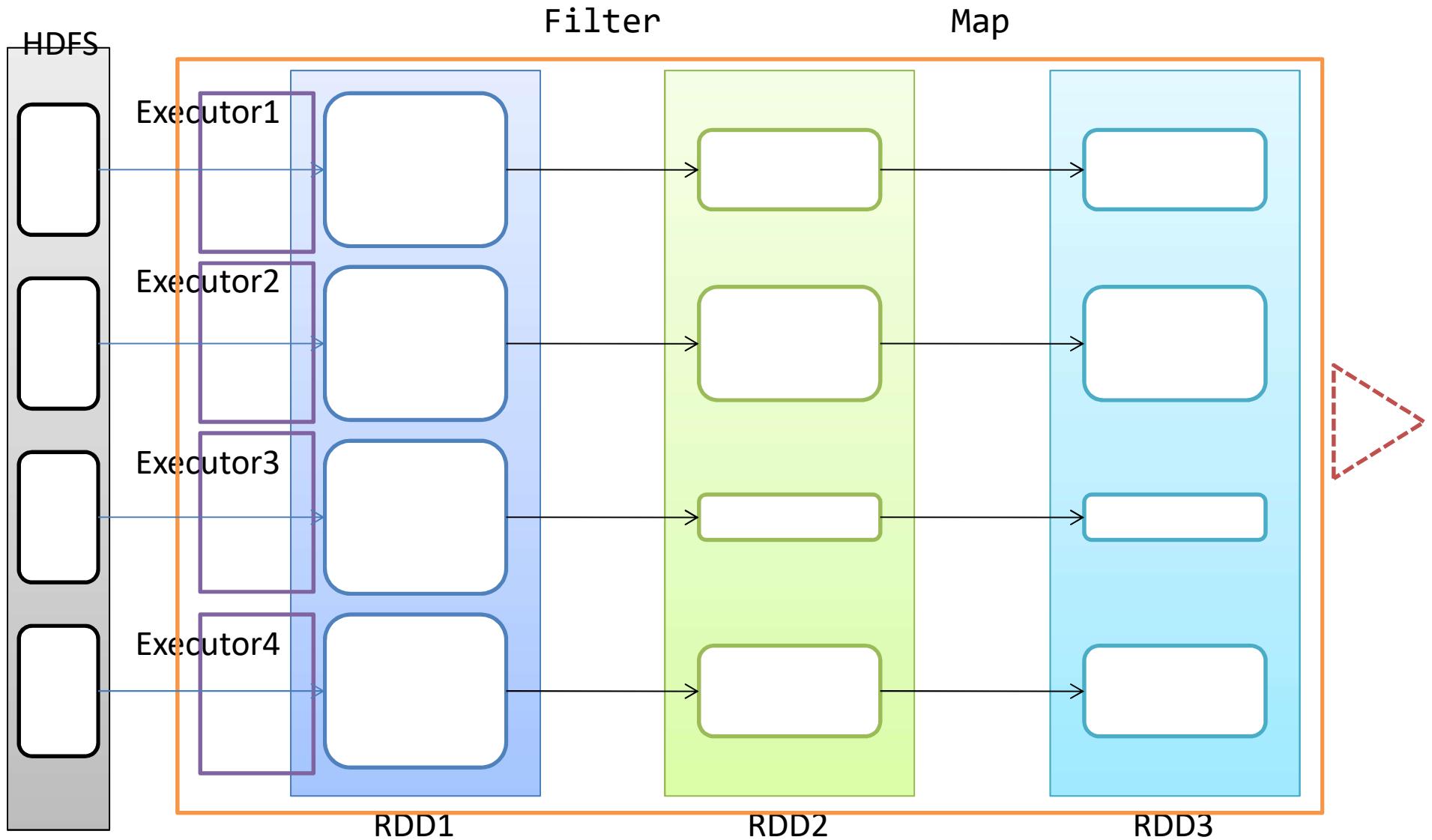


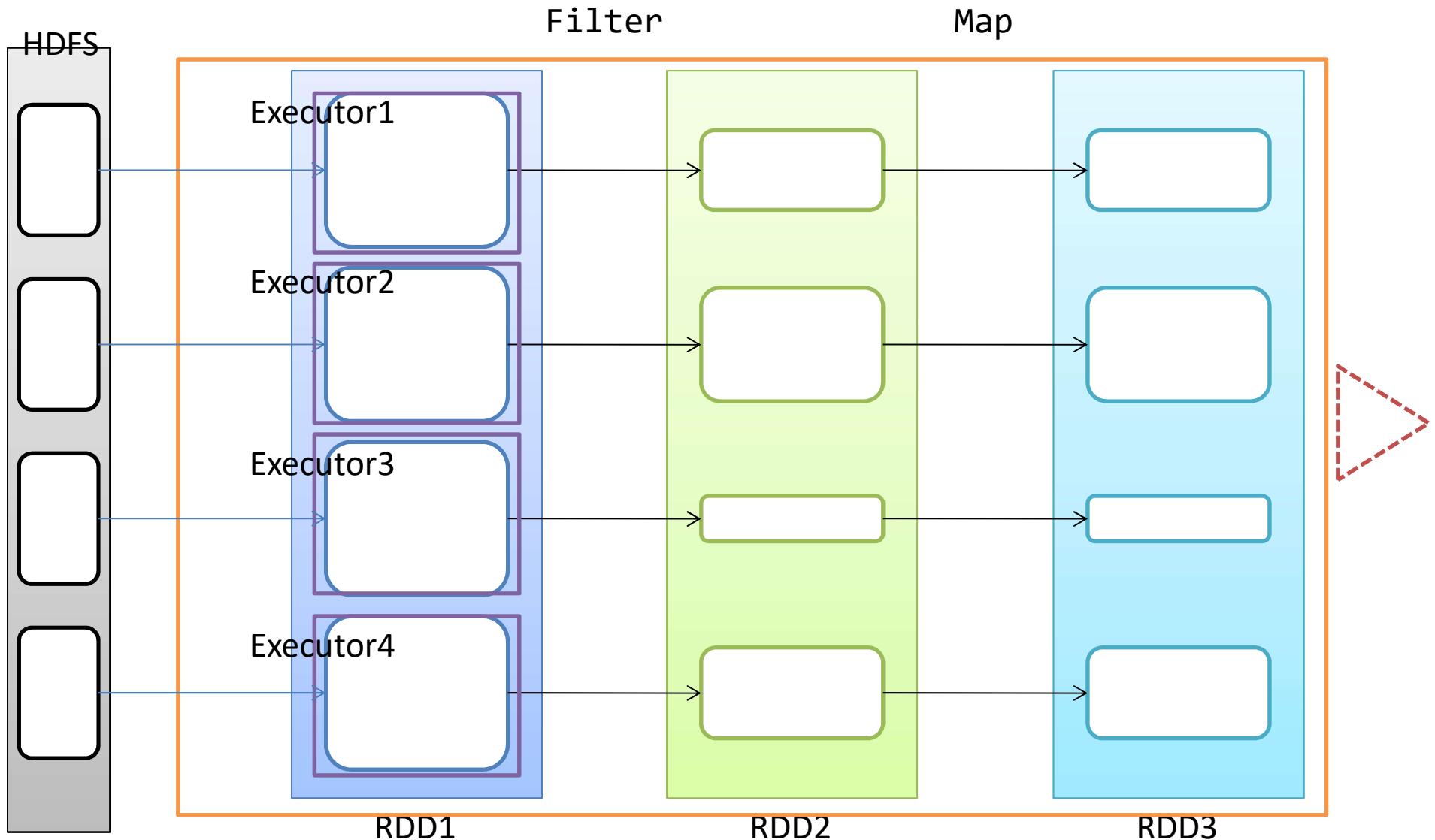


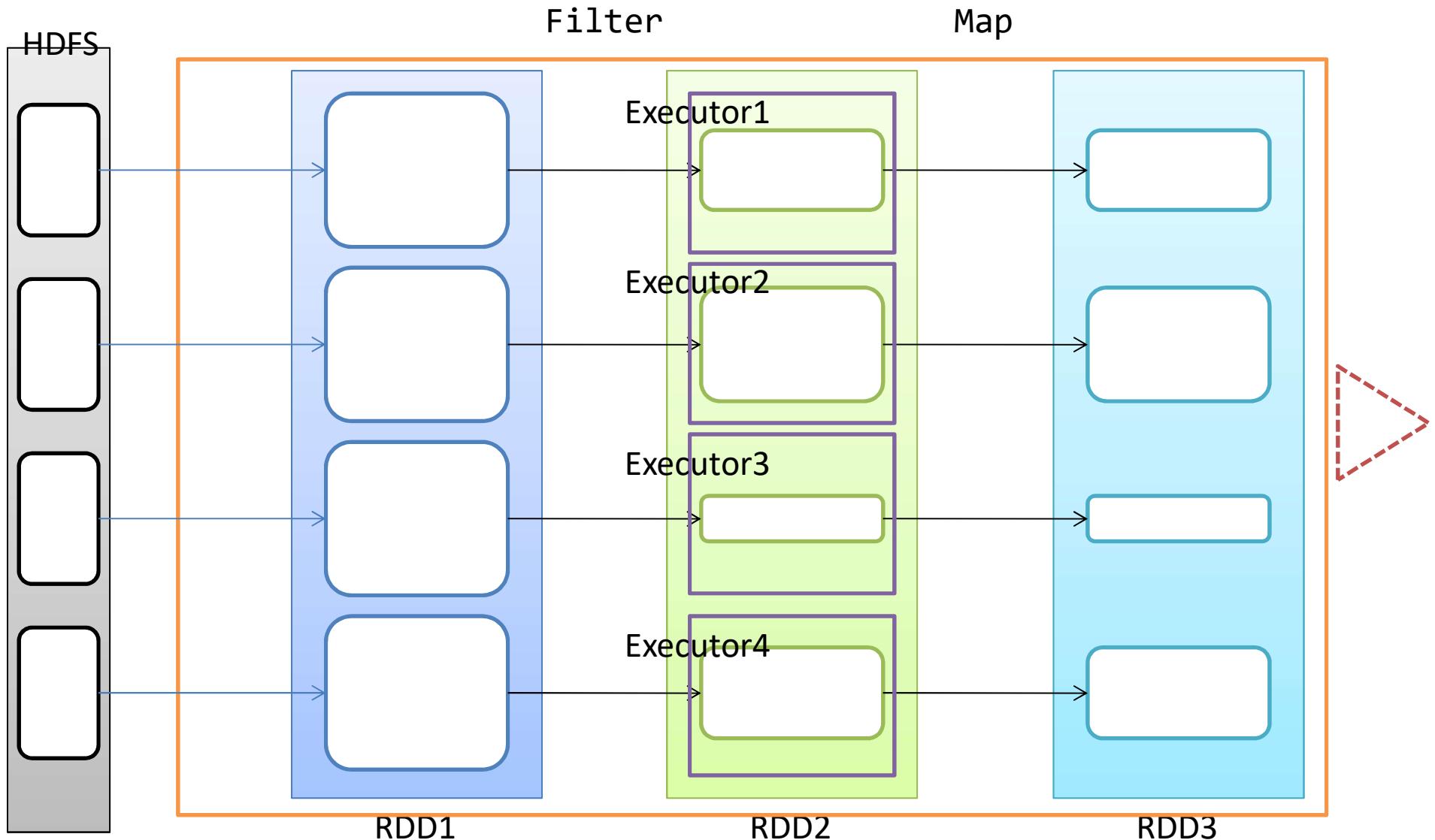


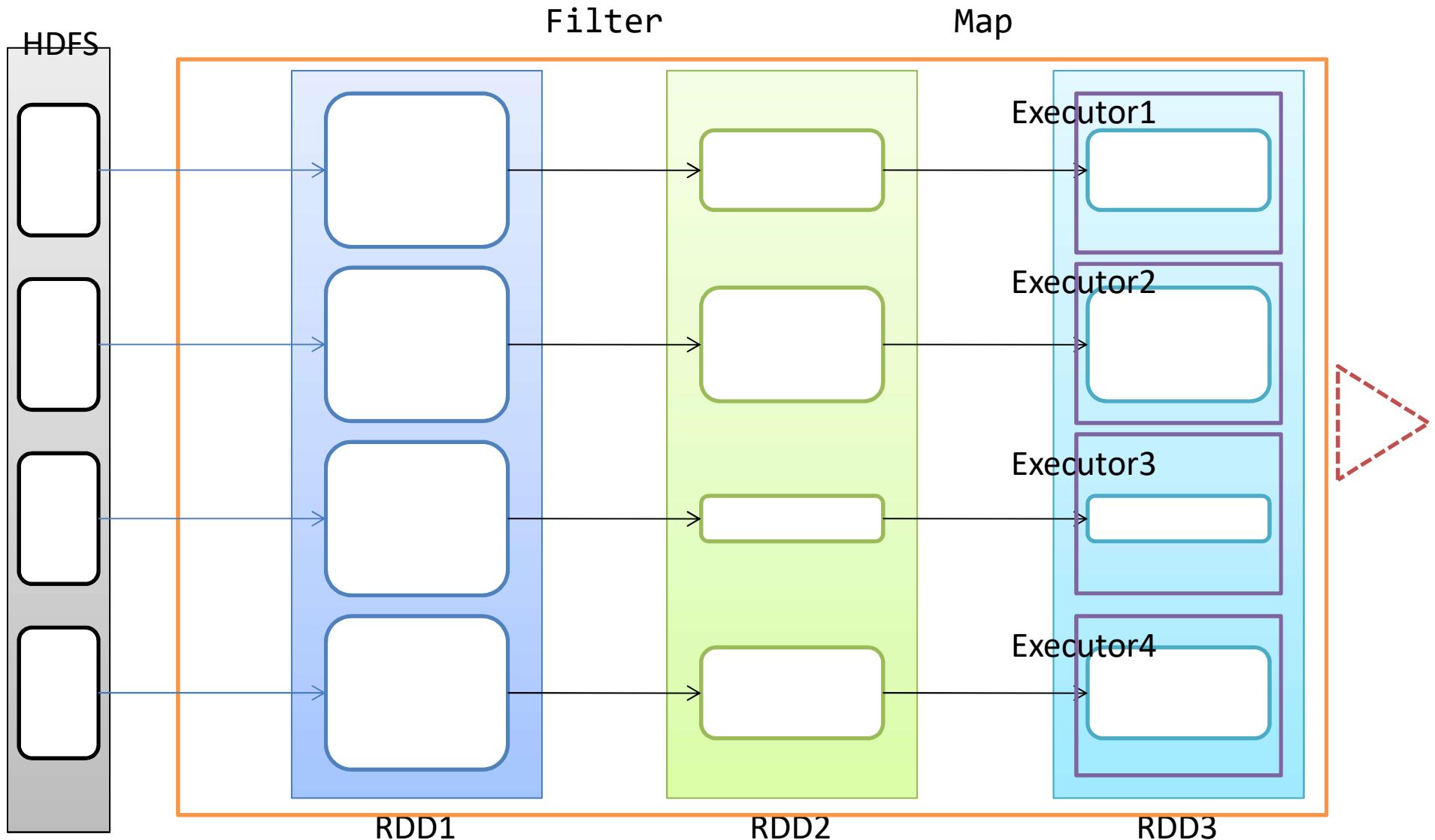


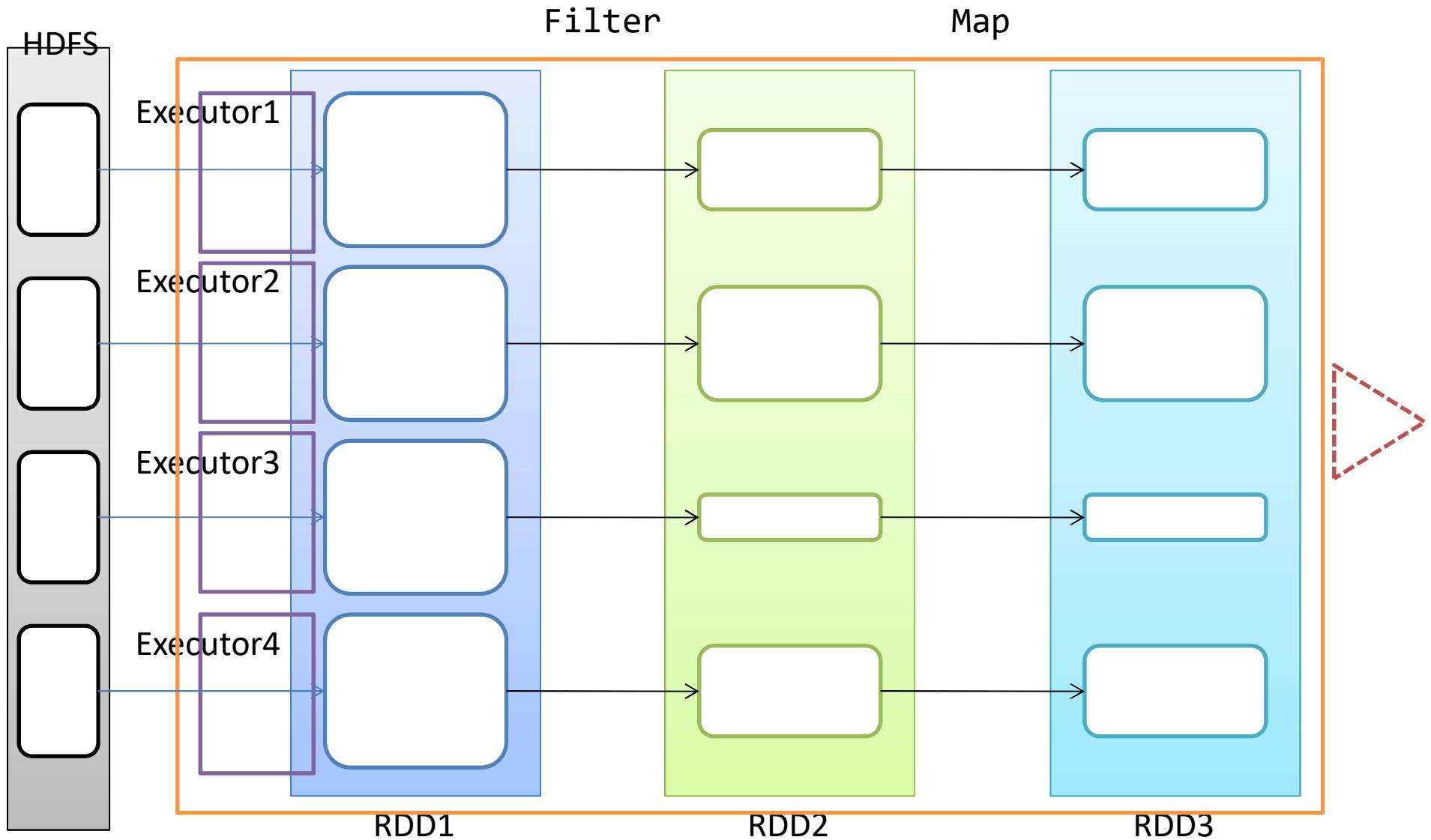


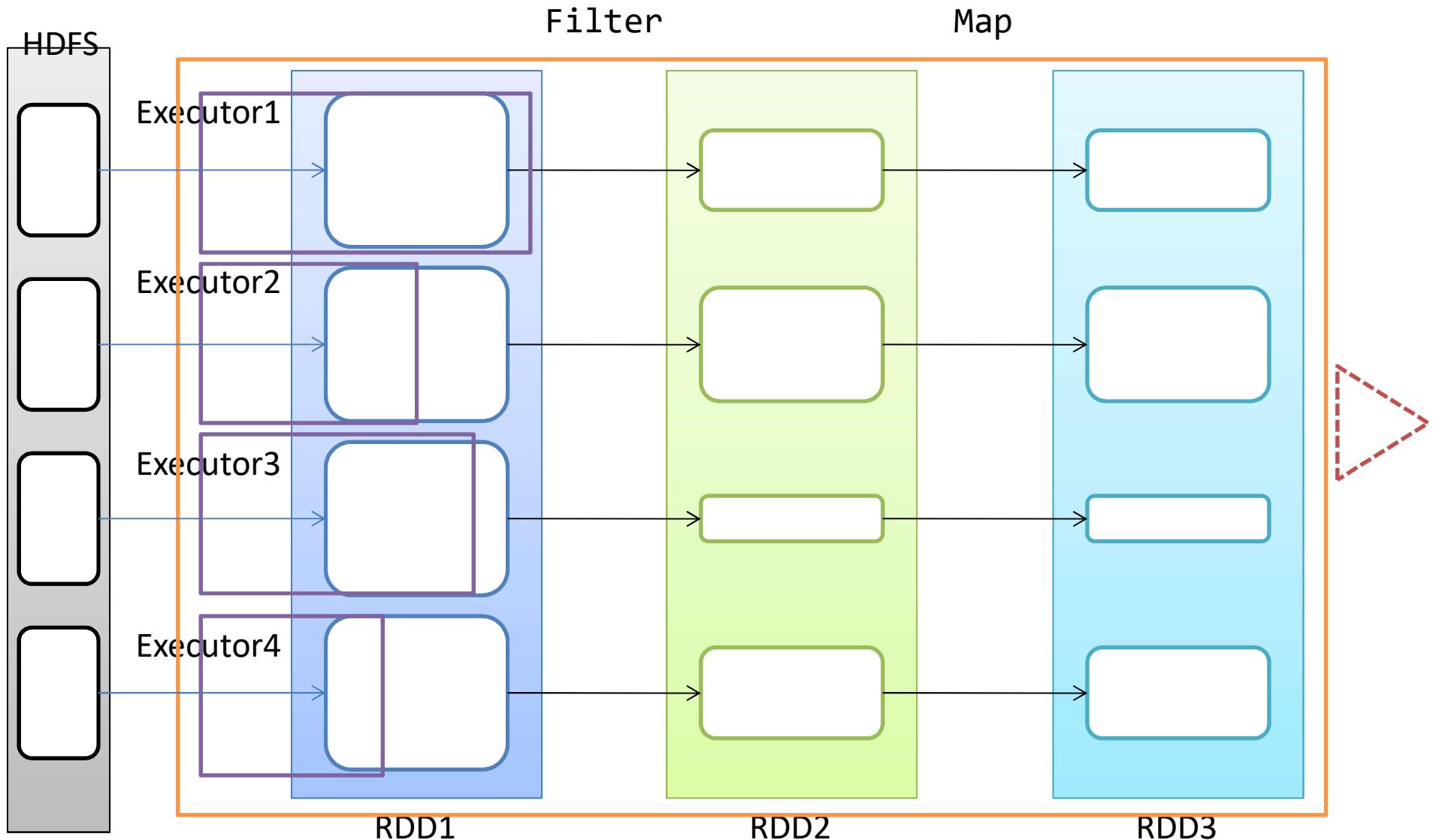


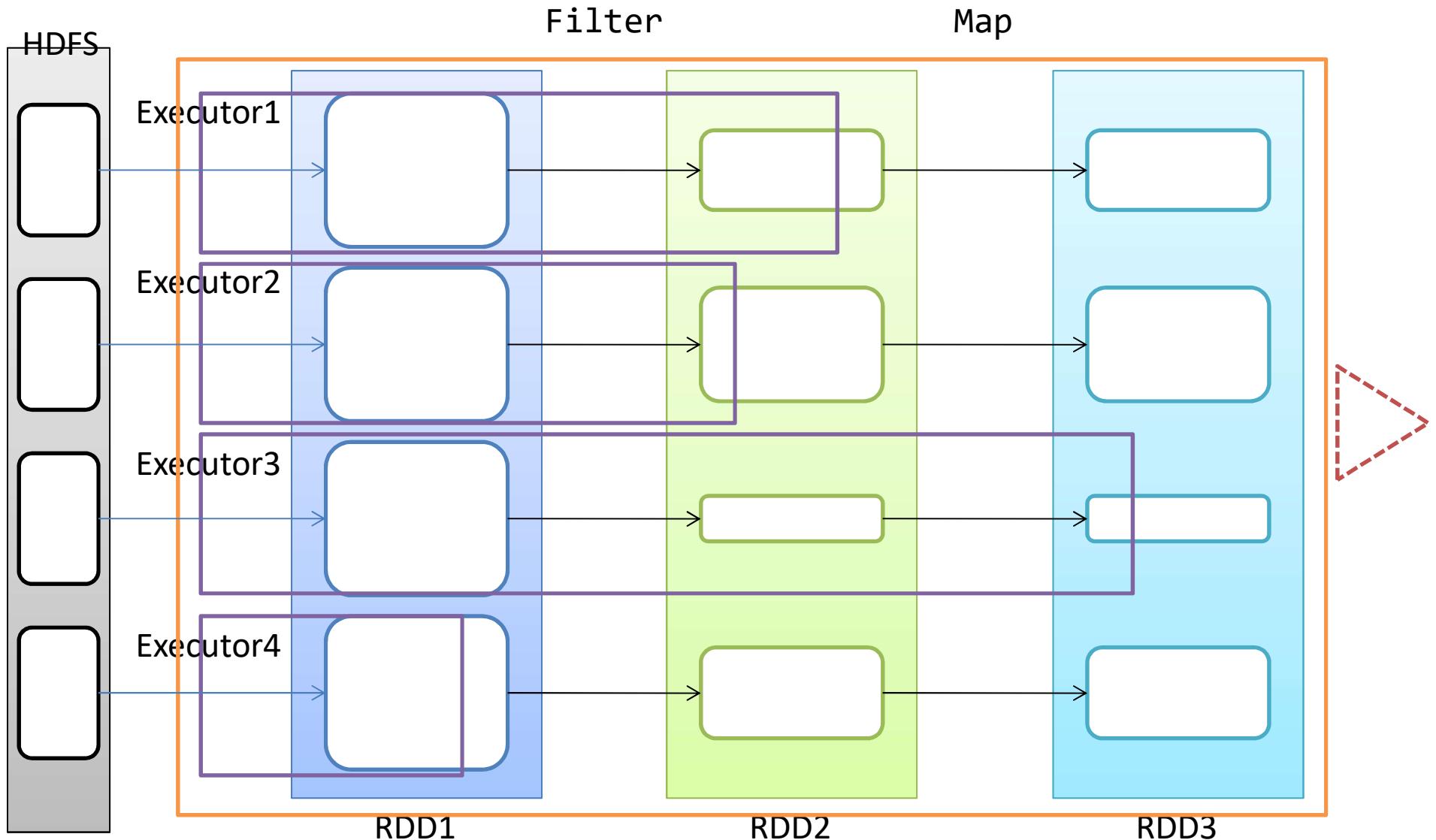


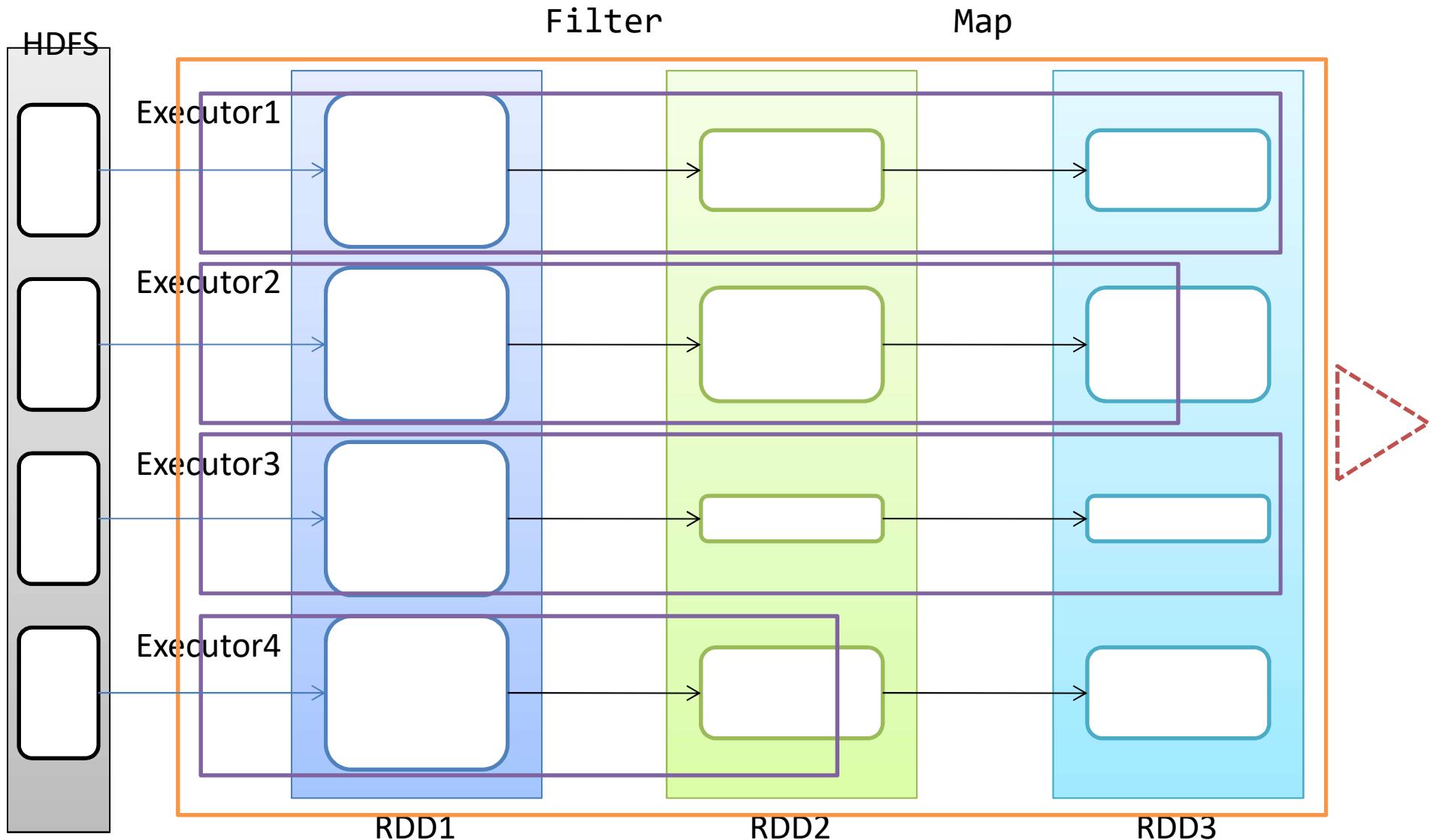


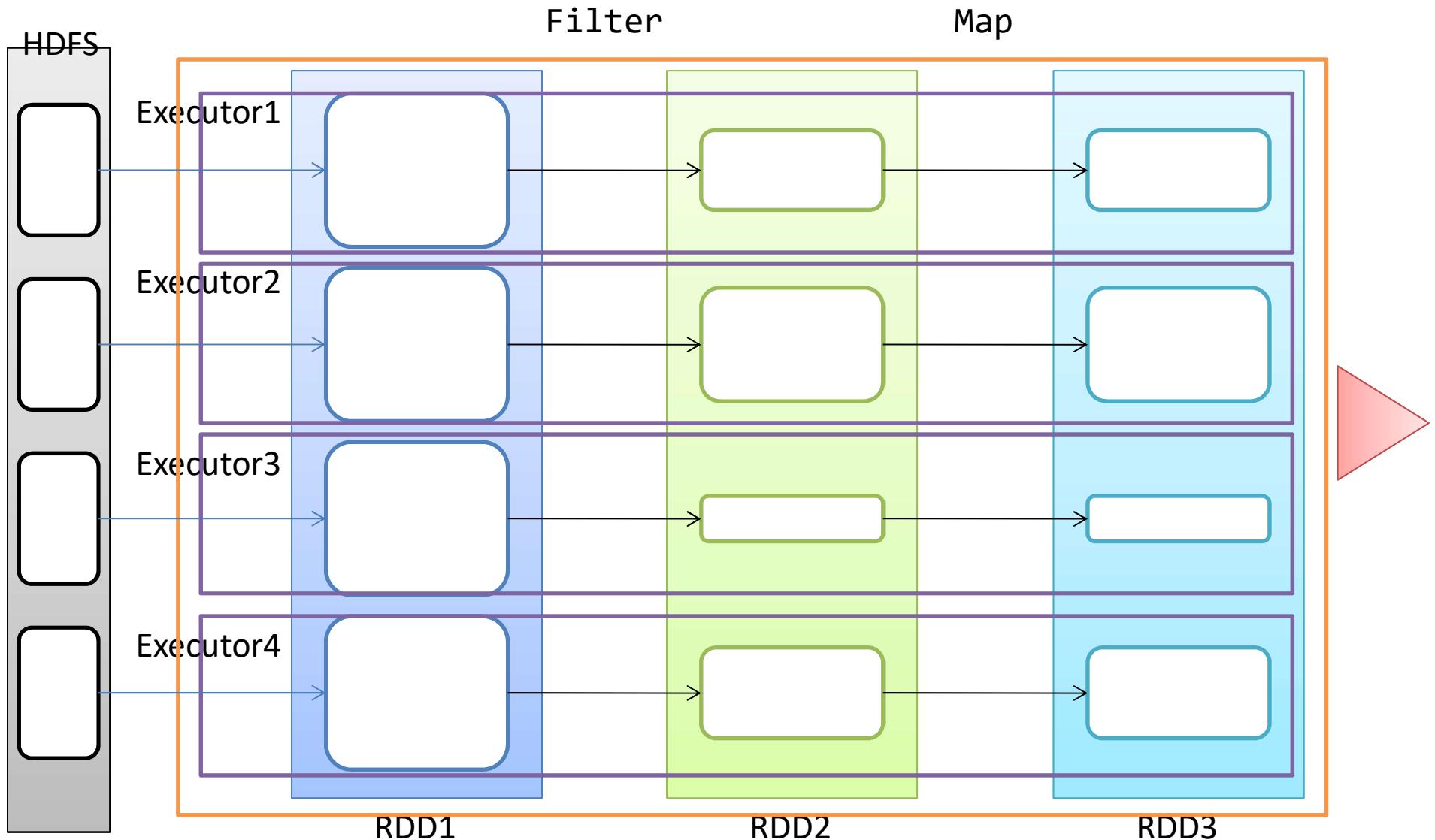


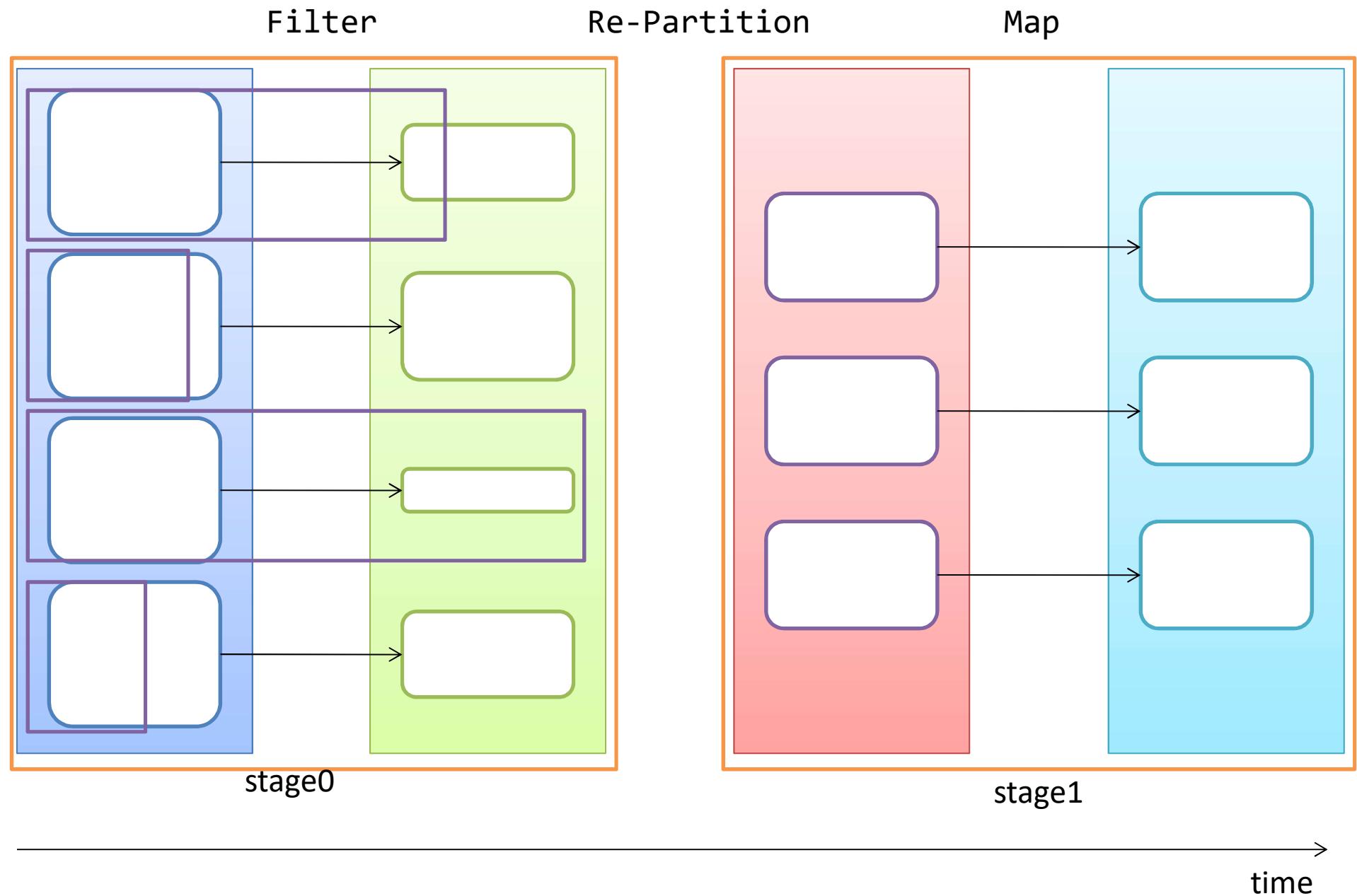


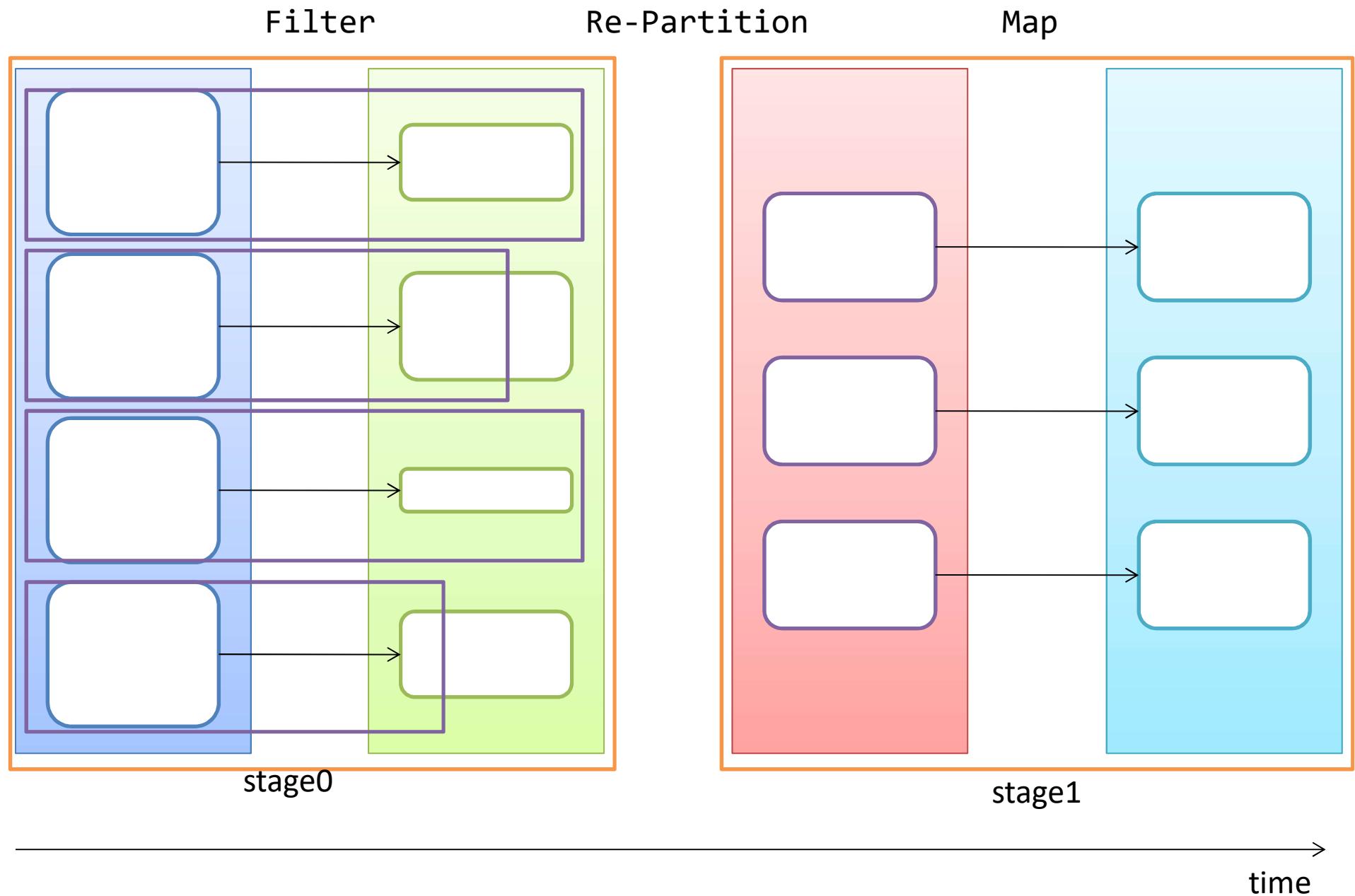


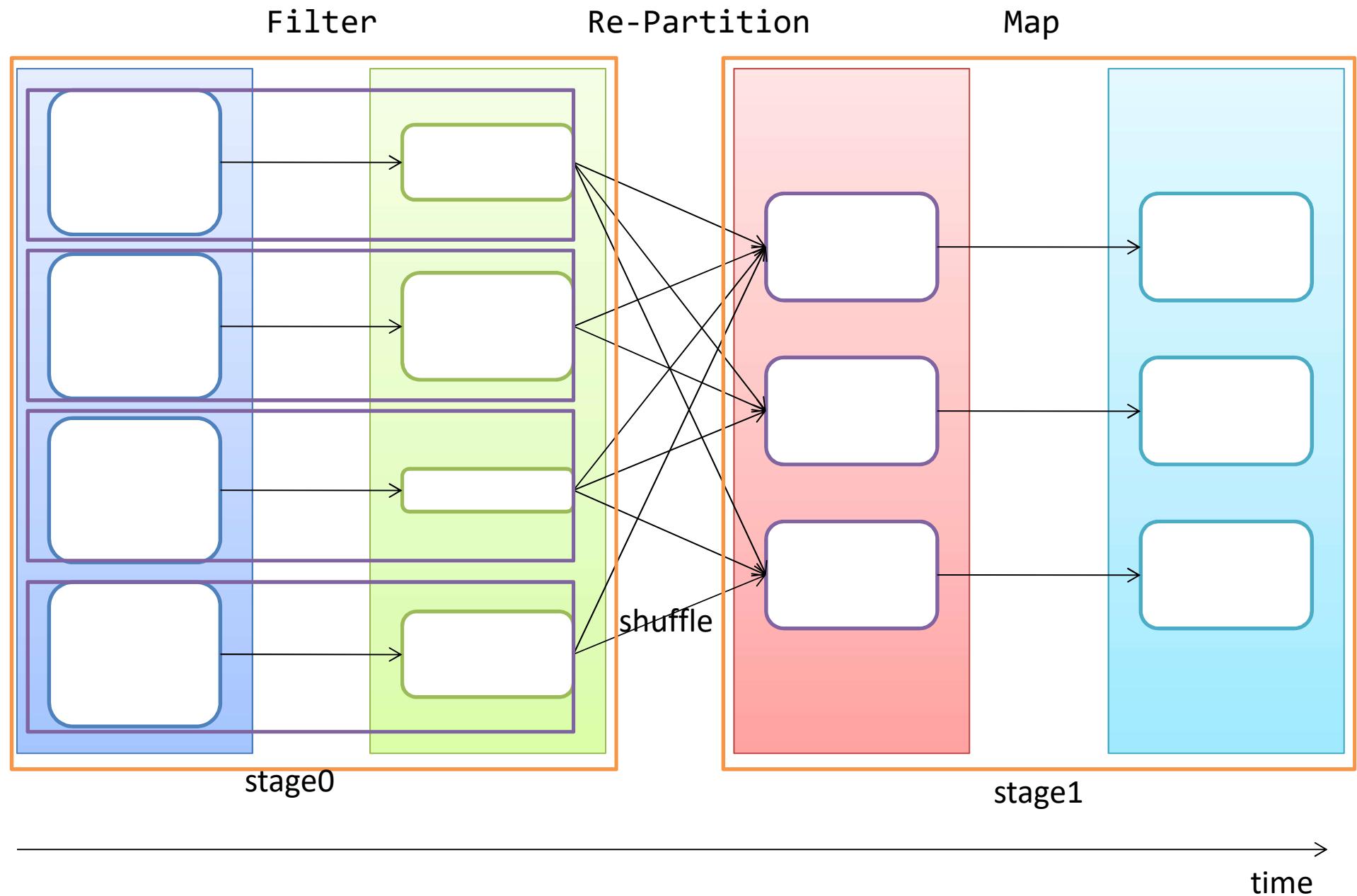


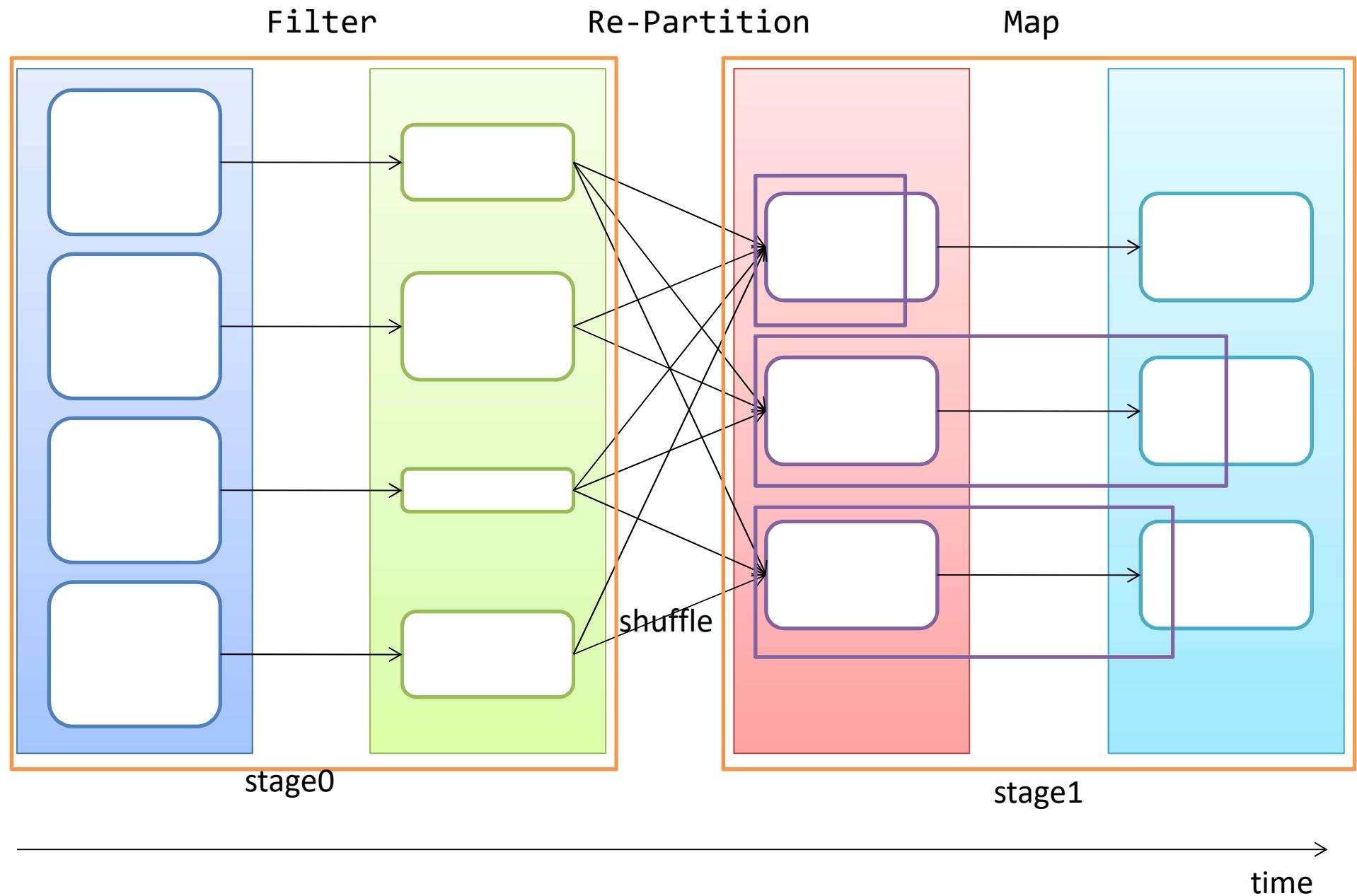






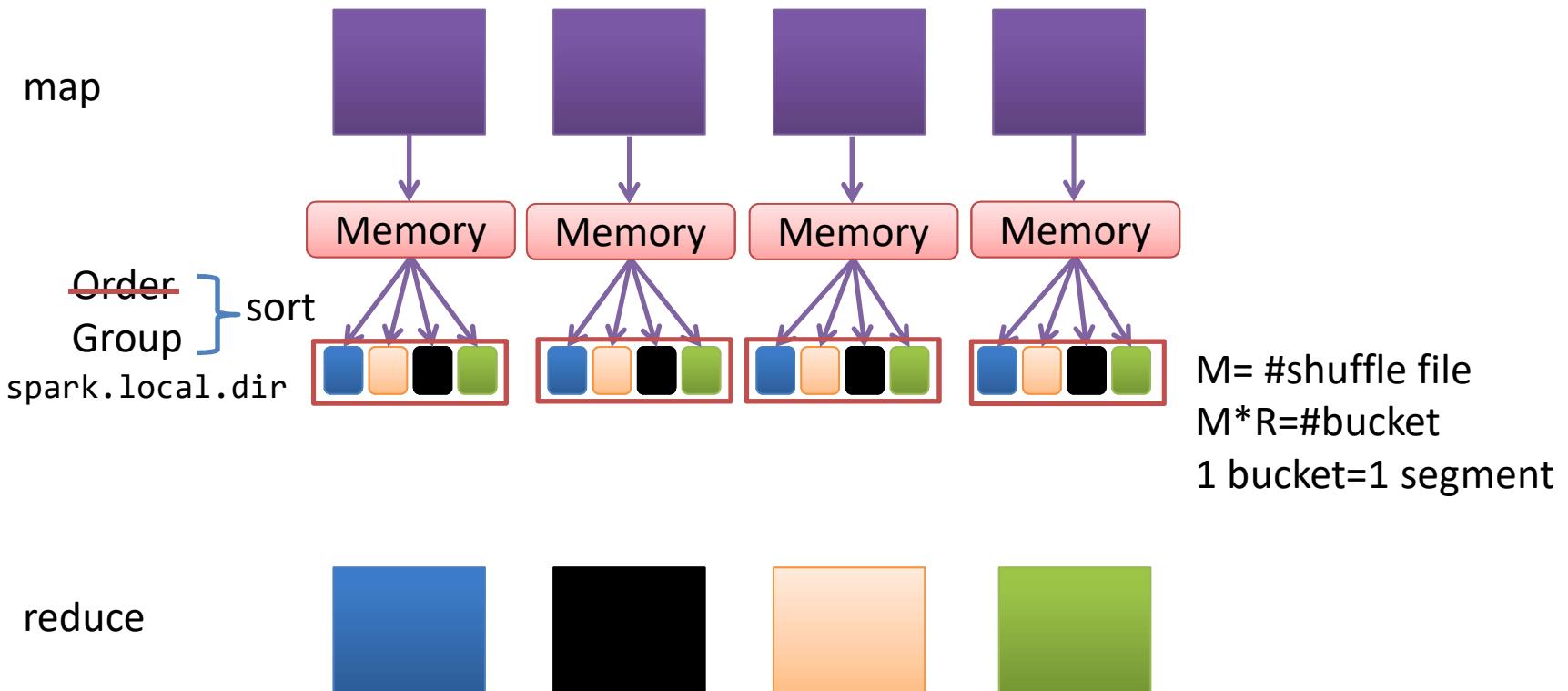




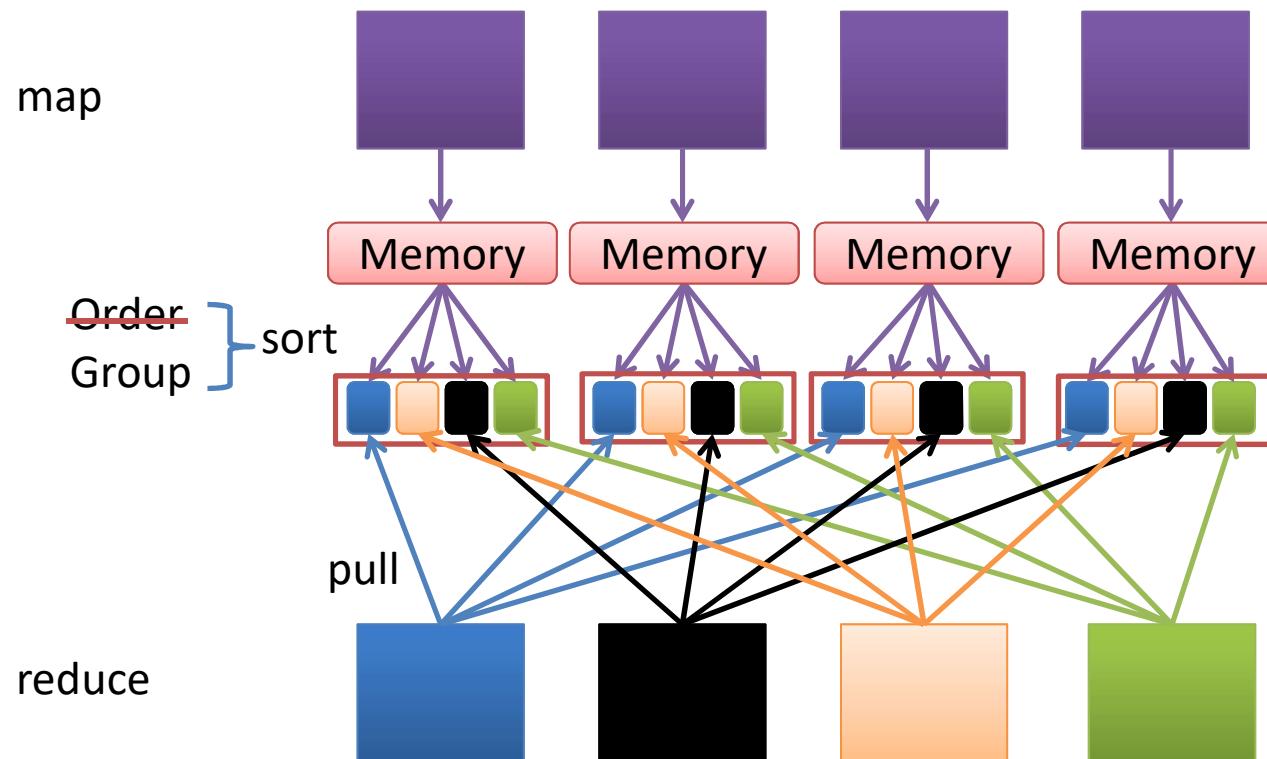


## Spark 1.1 Sort-based shuffle

Core=4 => 4 Tasks running



## Spark 1.1 Sort-based shuffle



- Application
- Driver Program
- Cluster Manager
- Worker Node
- Executor
- Job
- Stage
- Task
- Partition

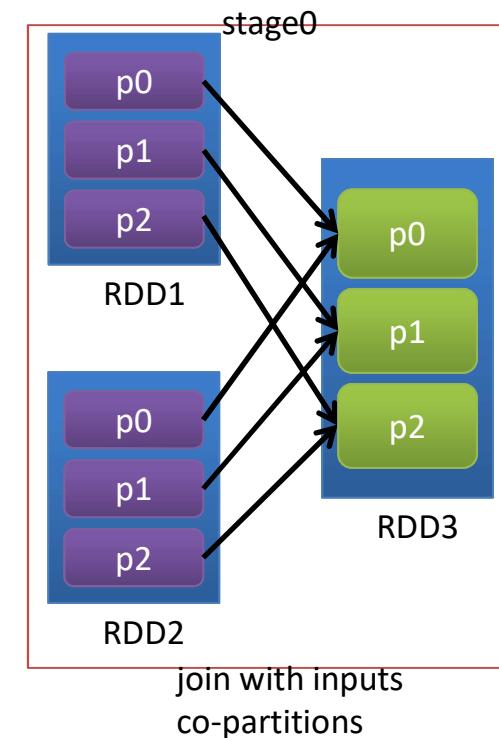
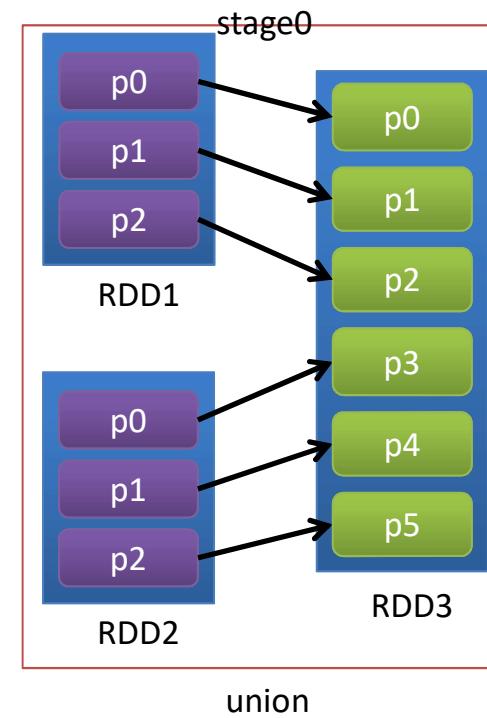
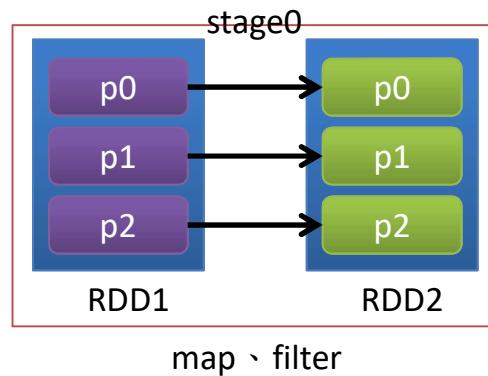
# Job

- 在一個Spark Application中，以Action操作為邊界劃分為多個Job。
- 一個job中，可包含多個Transformation操作與一個Action操作。
- 一個job可由多個stage組成

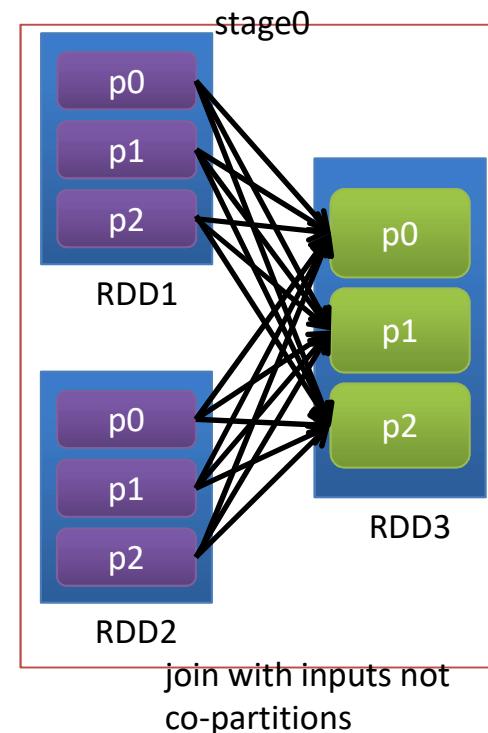
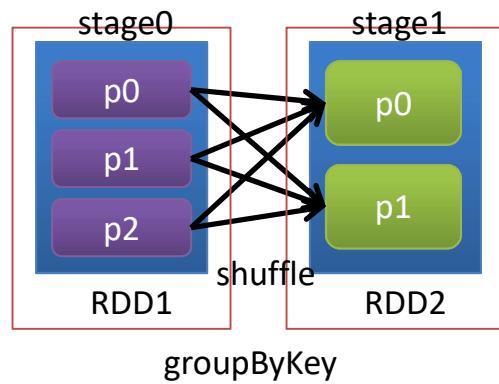
# Dependency(相依)

- Narrow Dependency(窄相依)
  - 父RDD的一個partition最多被一個子RDD partition所依賴
- Wide Dependency(寬相依)
  - 父RDD的一個partition被多個子RDD partition所依賴，此時會依據元素的Key值進行混洗。
  - mapPartitions、repartitionAndSortWithinPartitions、sortBy、repartition、~~coalesce~~、ByKey operations(groupByKey、reduceByKey)、join operations(cogroup、join)
- 寬相依是Job切分stage的參考點
- 窄相依之間可以進行pipeline操作

# Narrow Dependency(窄相依)



# Wide Dependency(寛相依)

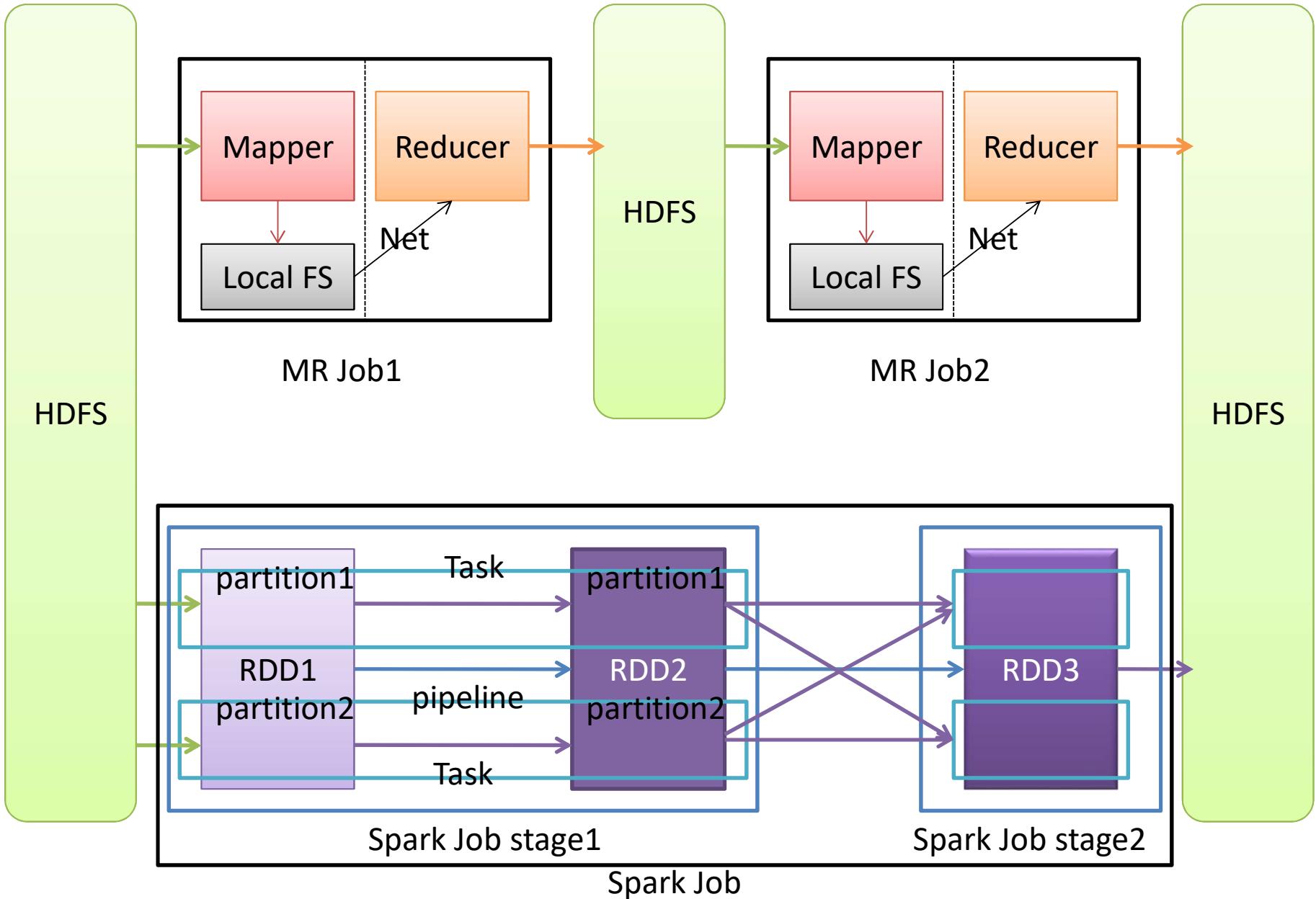


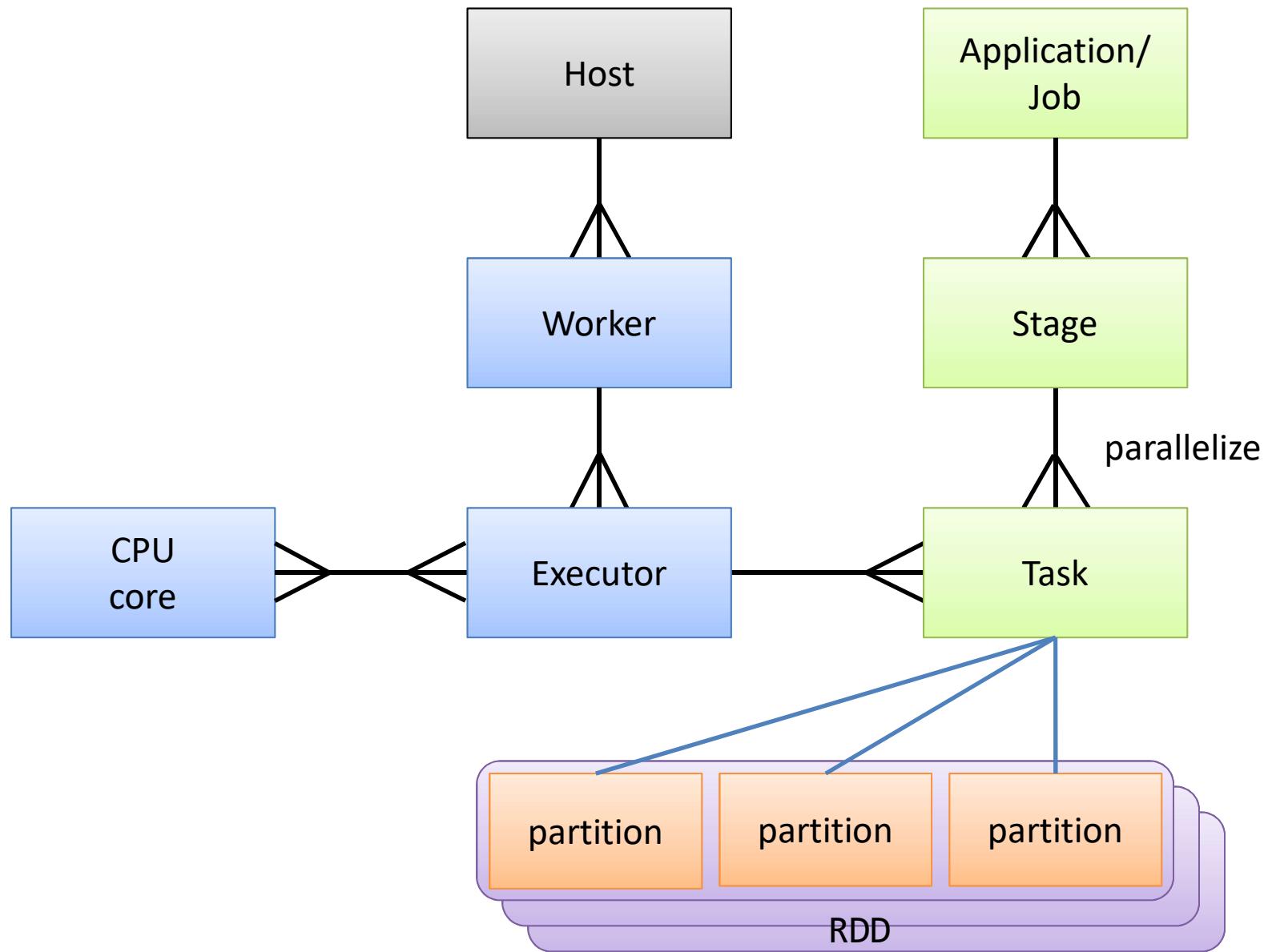
# Stage

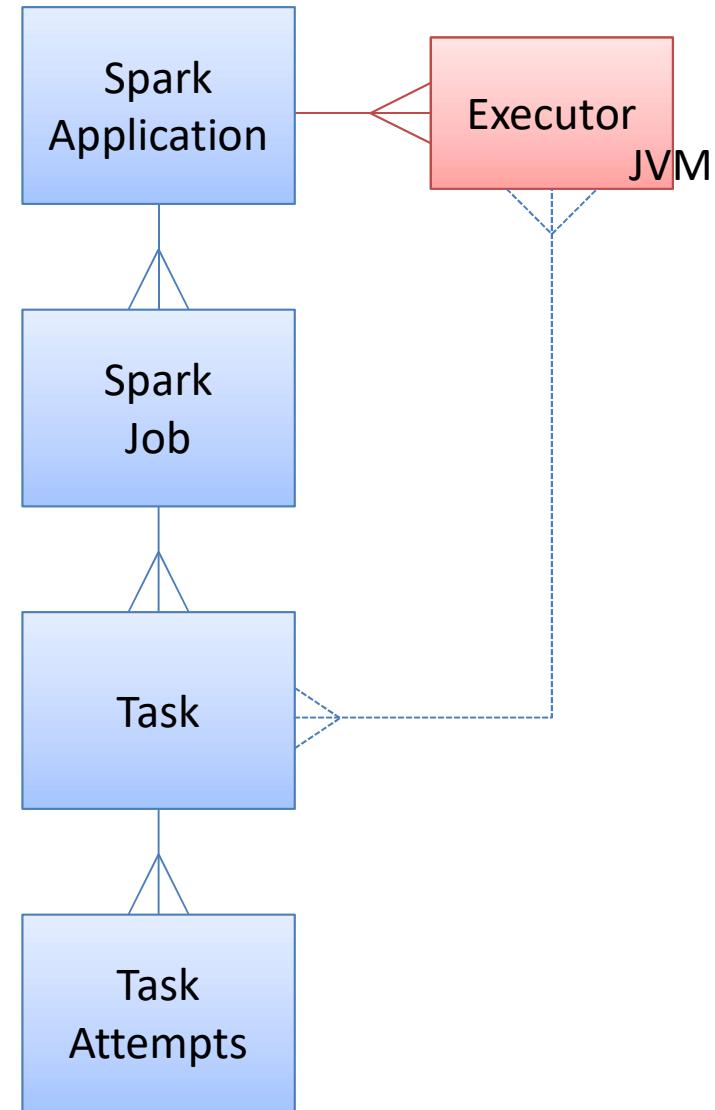
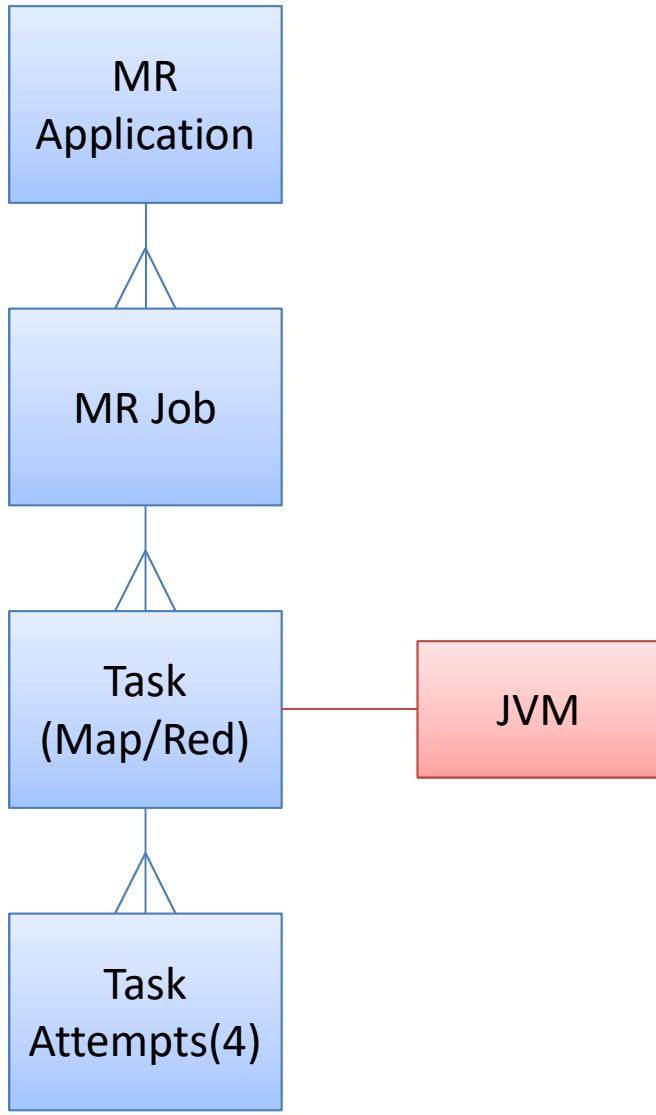
- 一個Spark job中以需要shuffle操作的轉換為邊界劃分出不同階段(stage)
- 每個階段包含一組管線化的窄相依操作
- 不同階段之間不能平行地同時進行，因為位在後面的階段需要位在前面的階段的執行結果。

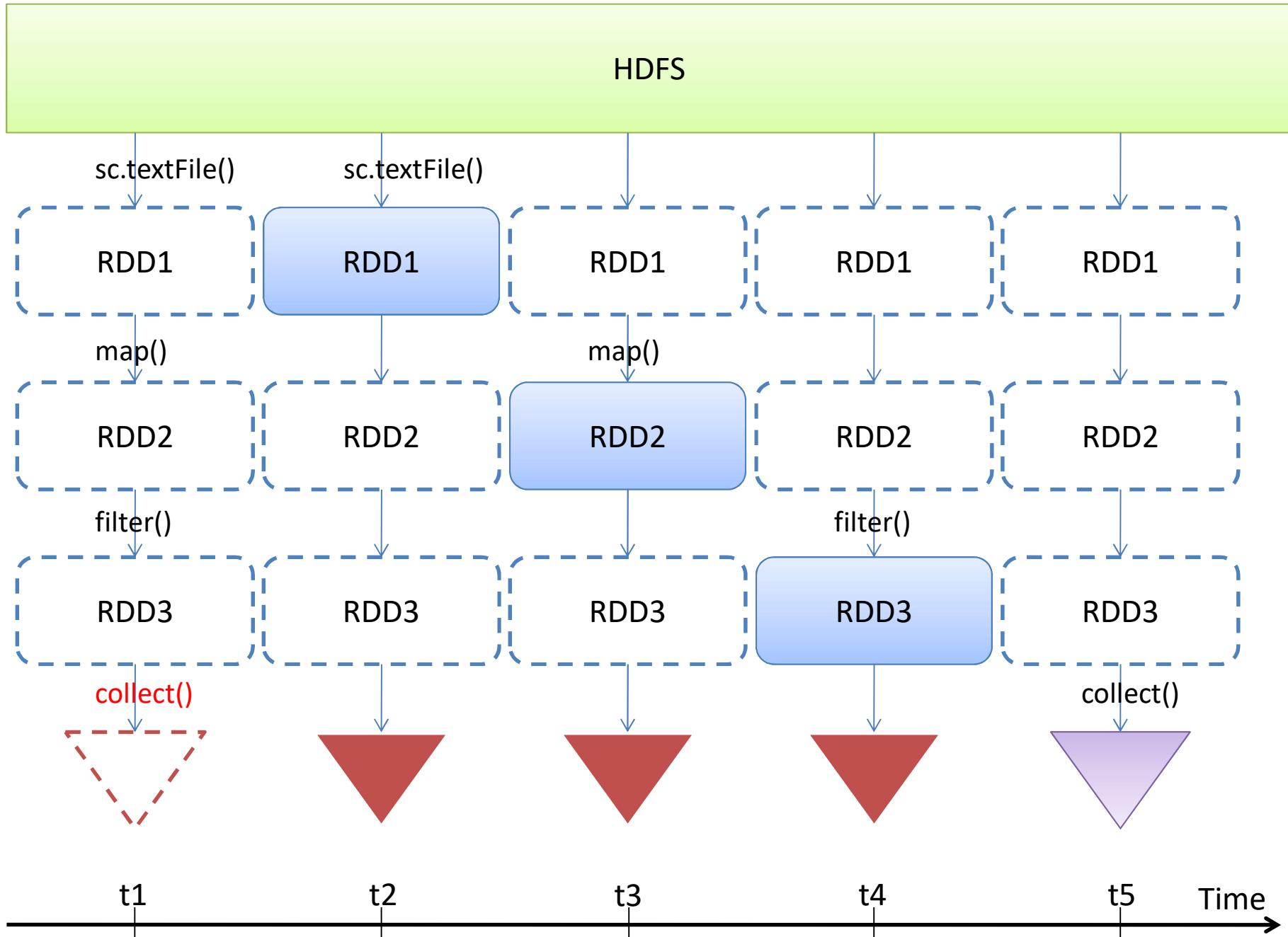
# Task

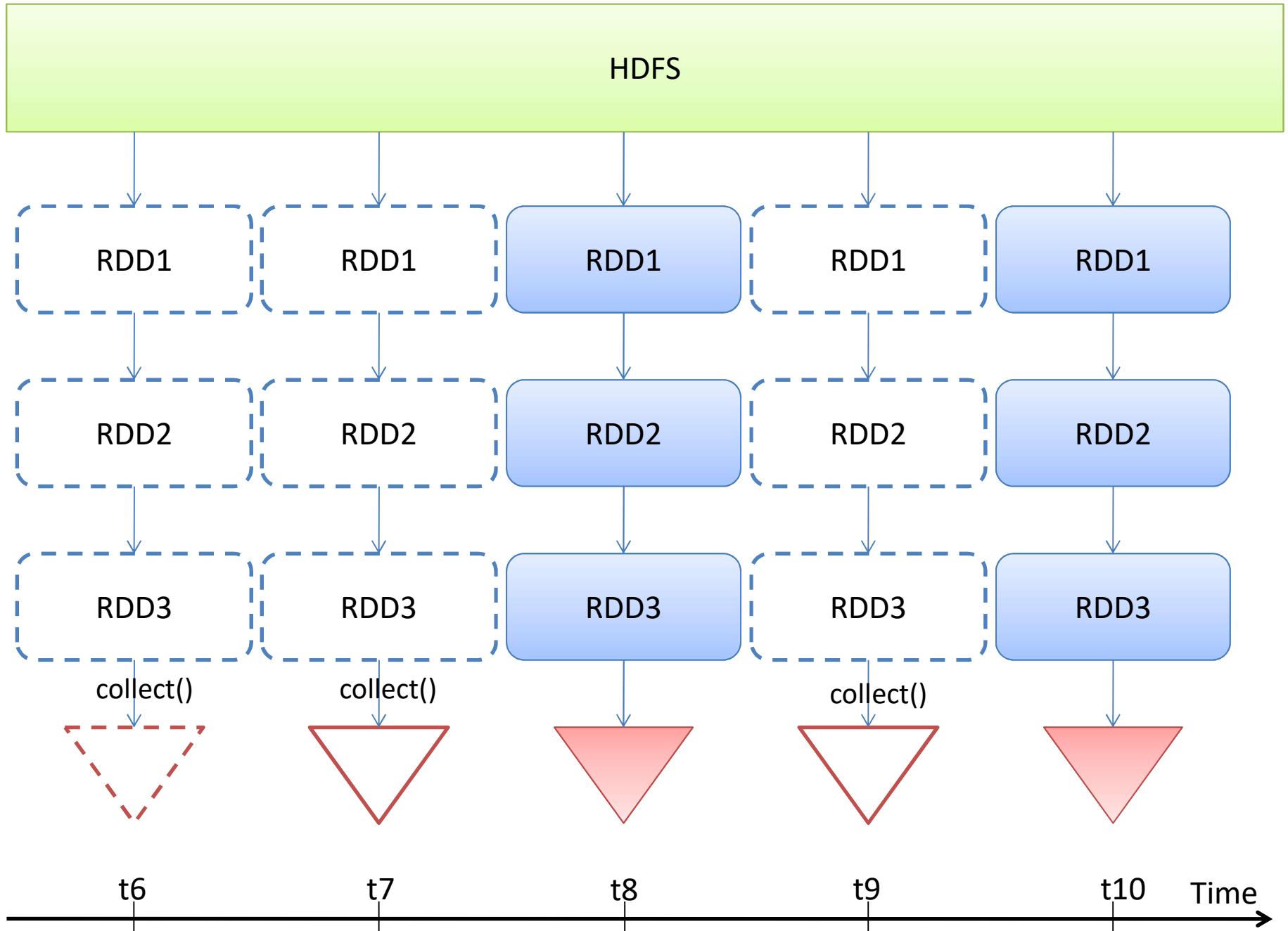
- 在一個stage的RDDs可以序列地管道化進行轉換操作。
- 每個RDD可以切分成多個Partition，每個Partition對應一個Task。所以一個stage可有多個Task組成(taskset)。這些task的執行邏輯完全相同，只是對應不同的partition。
- 一個stage的總task數，由此stage的最後一個RDD的partition數決定。
- 一個Executor可以平行地執行多少個task，由executor被配置的cpu數/每個task所需cpu數決定  
(`spark.executor.cores/spark.task.cpus`)
- Task分為ShuffleMapTask/ResultTask，位在最後一個stage的task為ResultTask，其他階段的task皆屬於ShuffleMapTask。

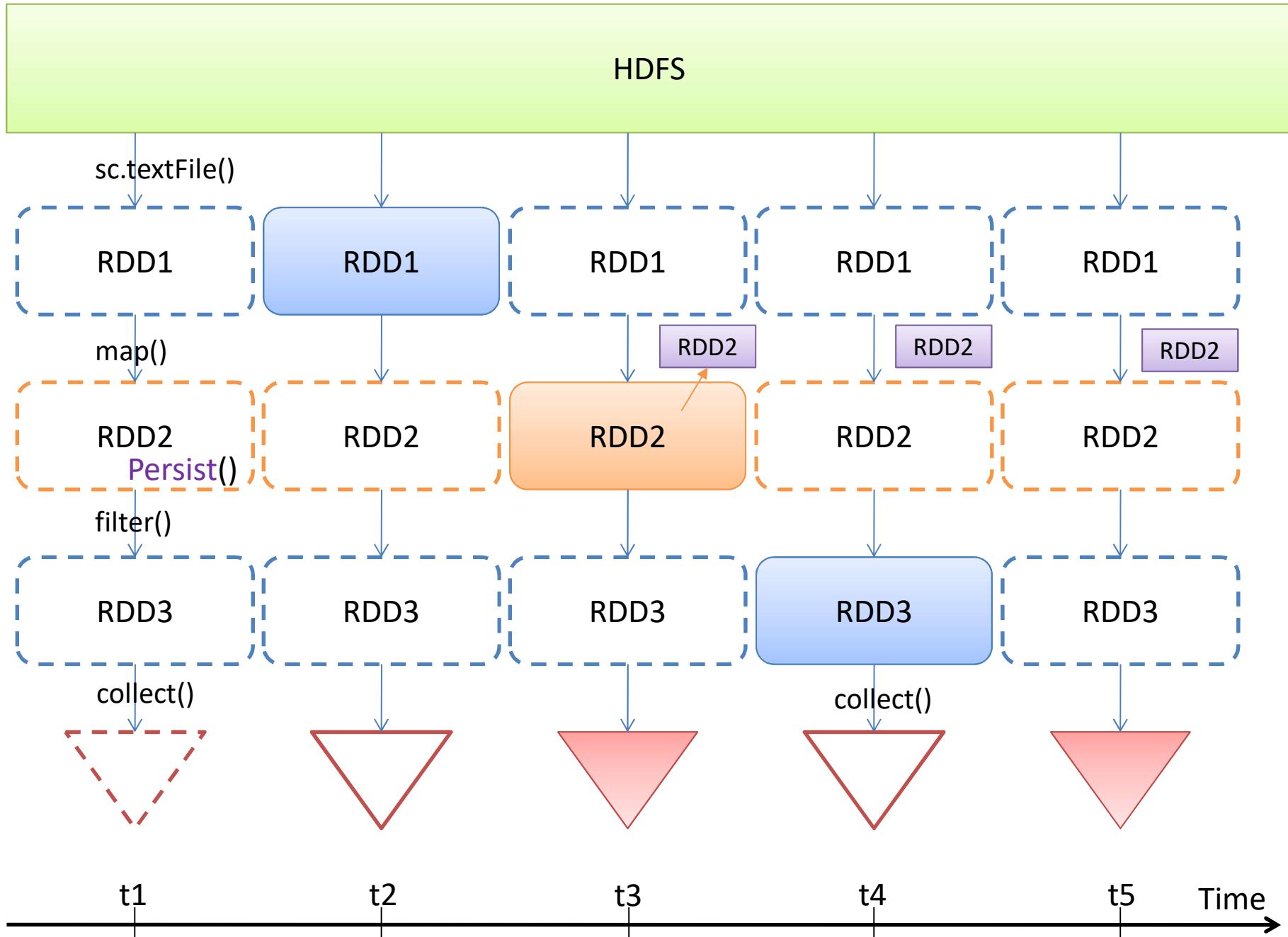


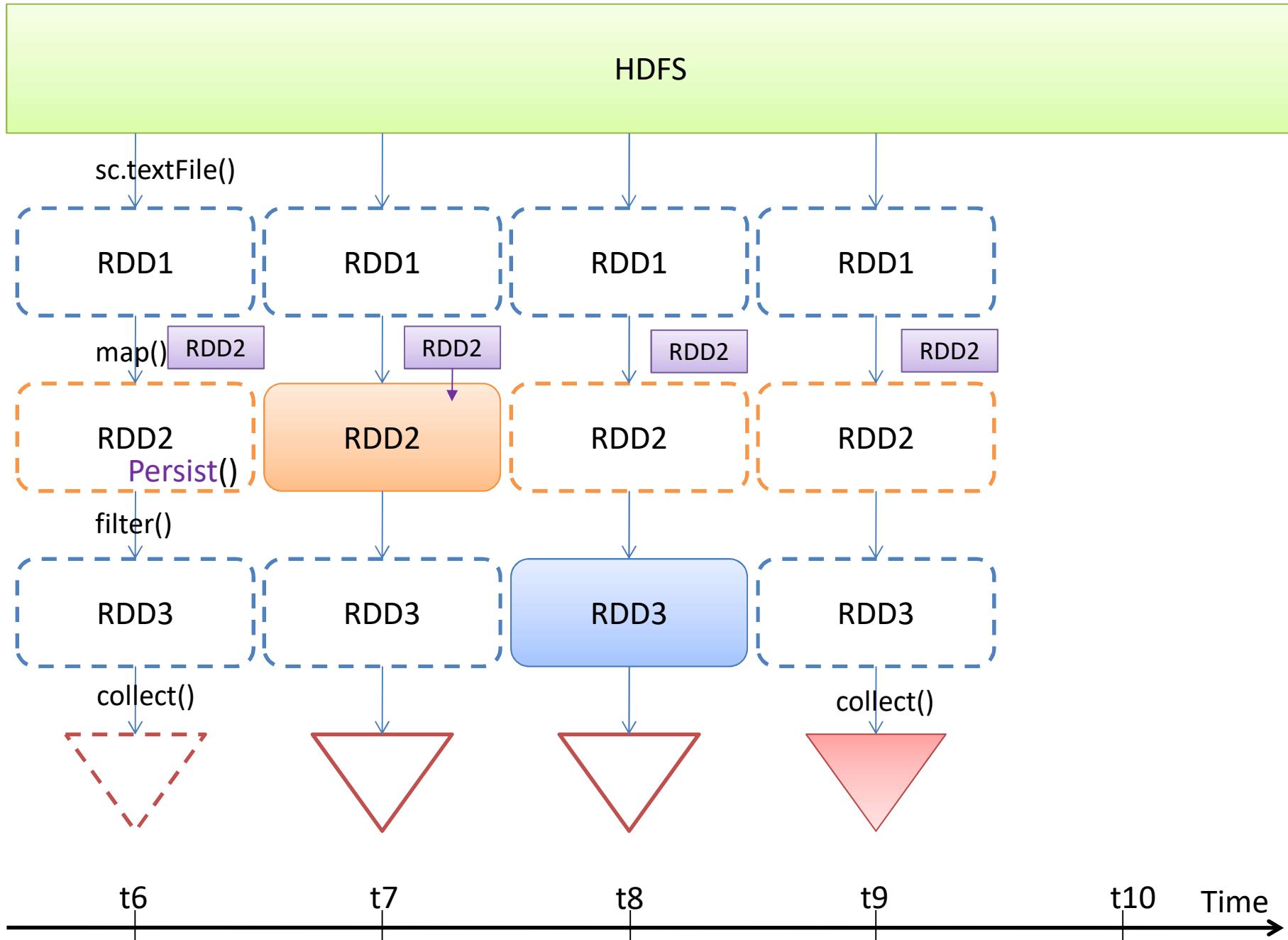


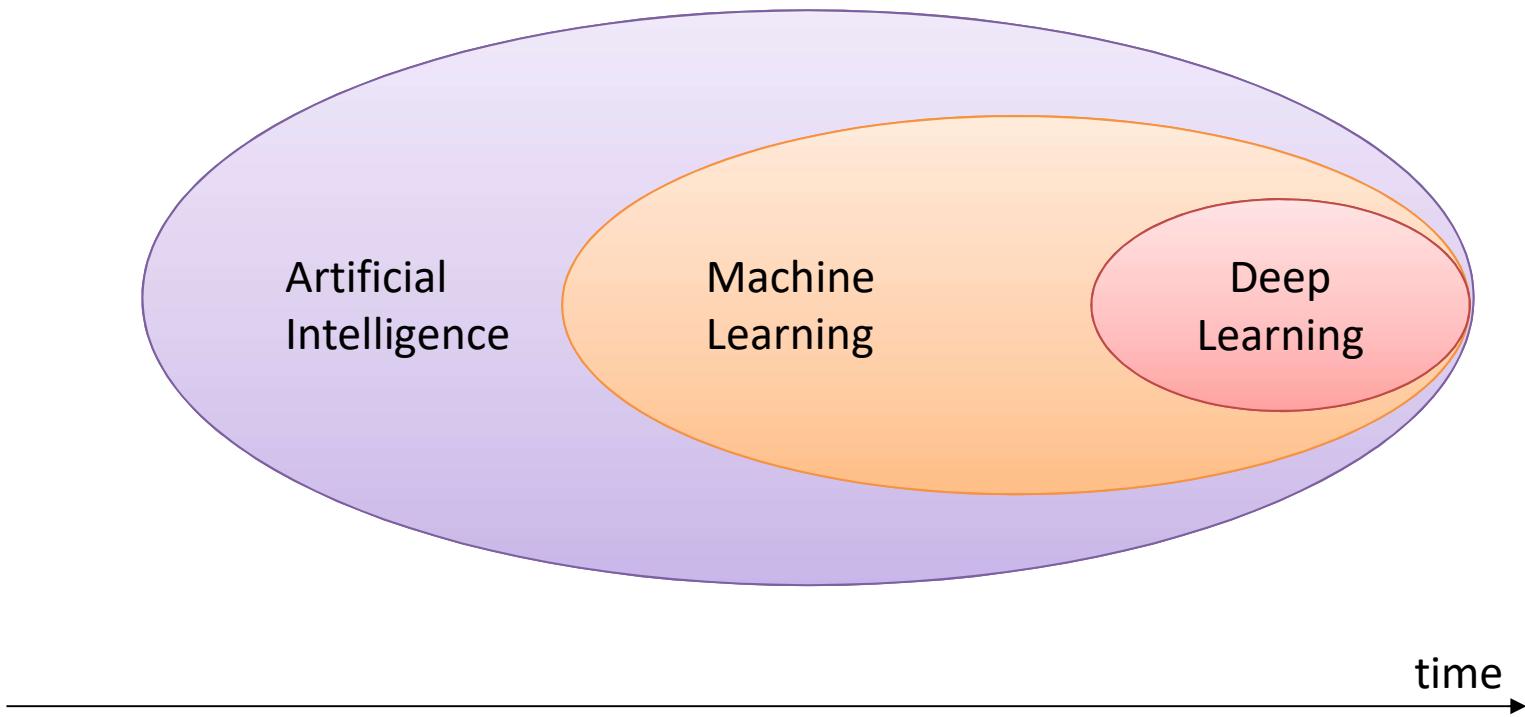


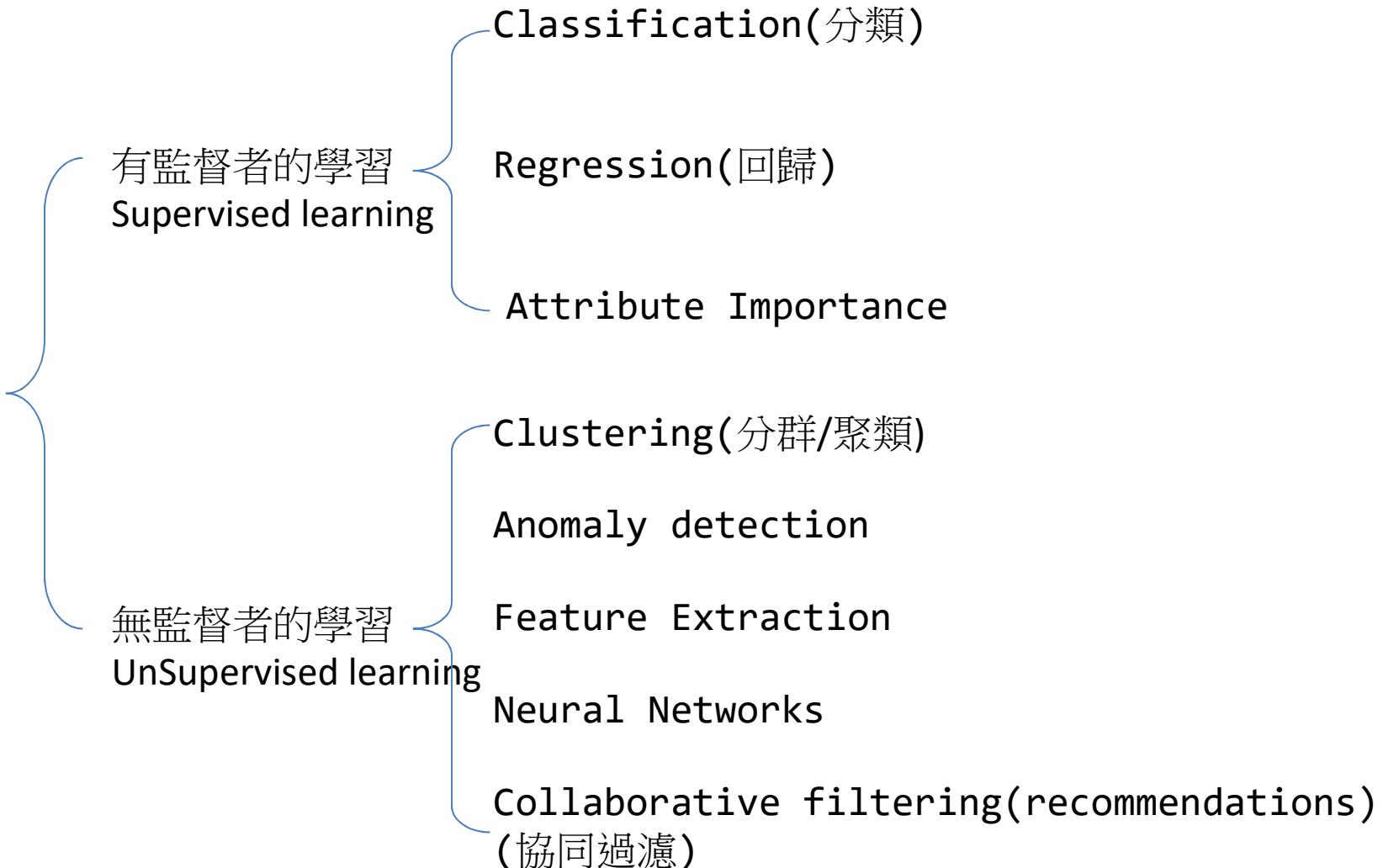




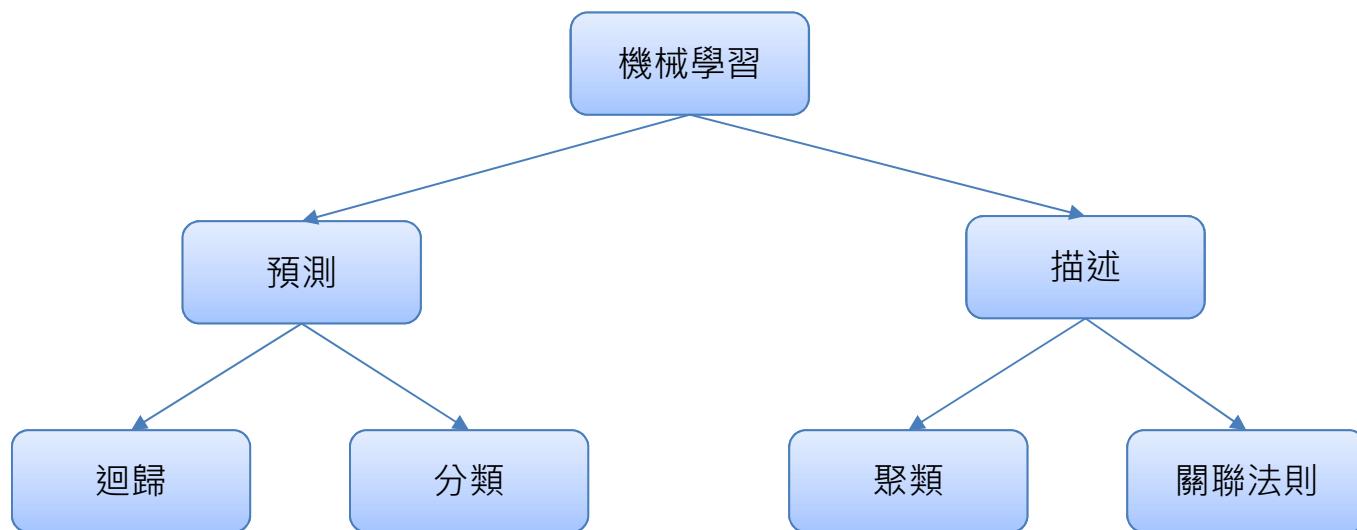








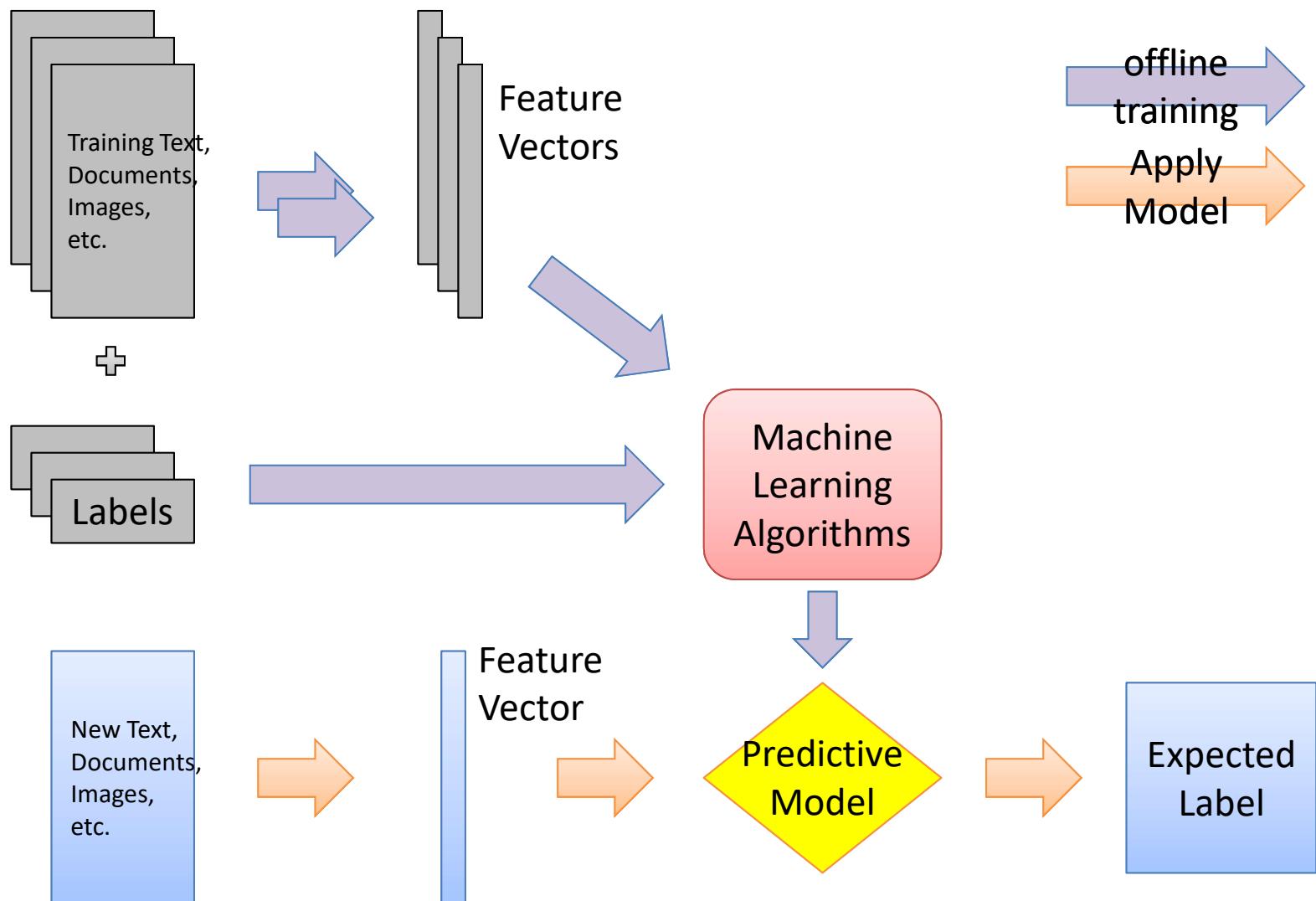
# 選擇適當的機械學習方法



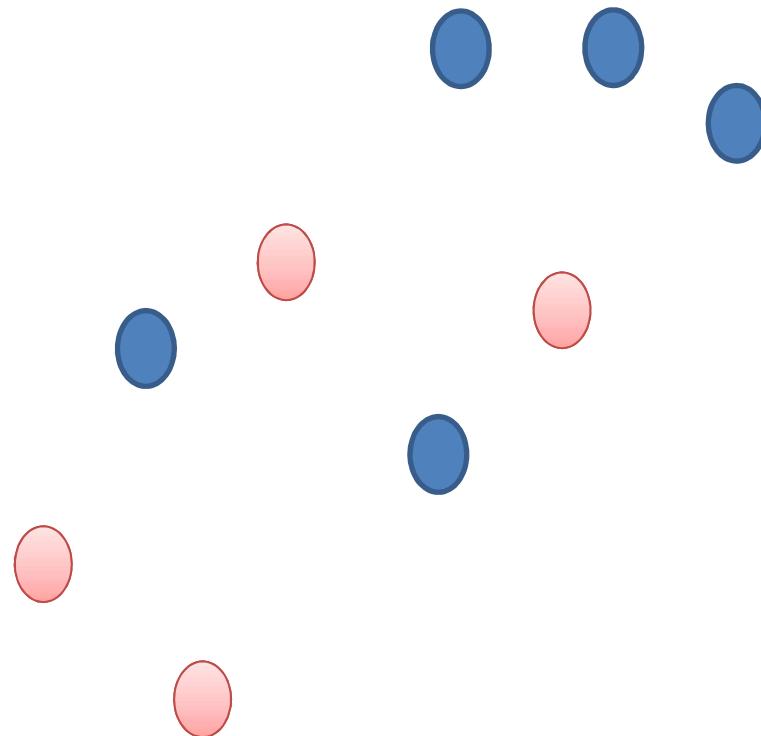
# 監督式學習

- 從給定的訓練資料集中學習出一個函式，當新的資料到來時，可以根據這個函式預測結果。監督學習的訓練集要求是包括輸入和輸出，也可以說是特徵和目標。訓練集中的目標是由人標註的。
  - 迴歸(Regression)-輸出為數值型資料
    - 線性回歸(Linear Regression)
    - 多項式迴歸(Polynomial Regression)
  - 分類(Classification)-輸出為類別型資料
    - 邏輯回歸(Logistic Regression)
    - 決策樹(Decision Tree)、隨機森林(Random Forest)
    - k-最近鄰居(k-Nearest Neighbors, k-NN)
    - 支持向量機(SVM)、相關向量機(RVM)
    - 樸素貝葉斯(Naïve Bayes)
    - 神經網絡(Neural Network)
    - 感知器(Perceptron)

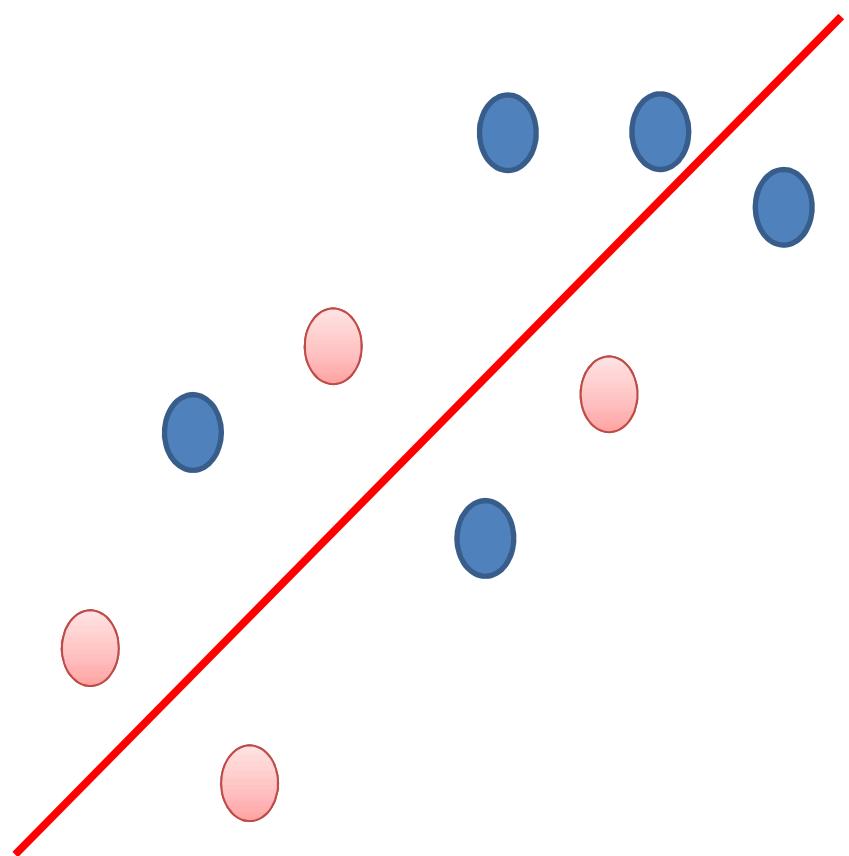
# 監督式學習流程



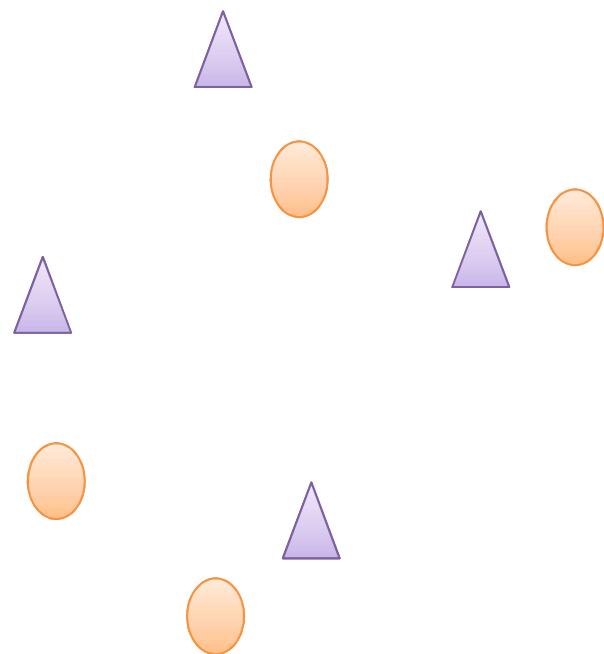
# 線性迴歸



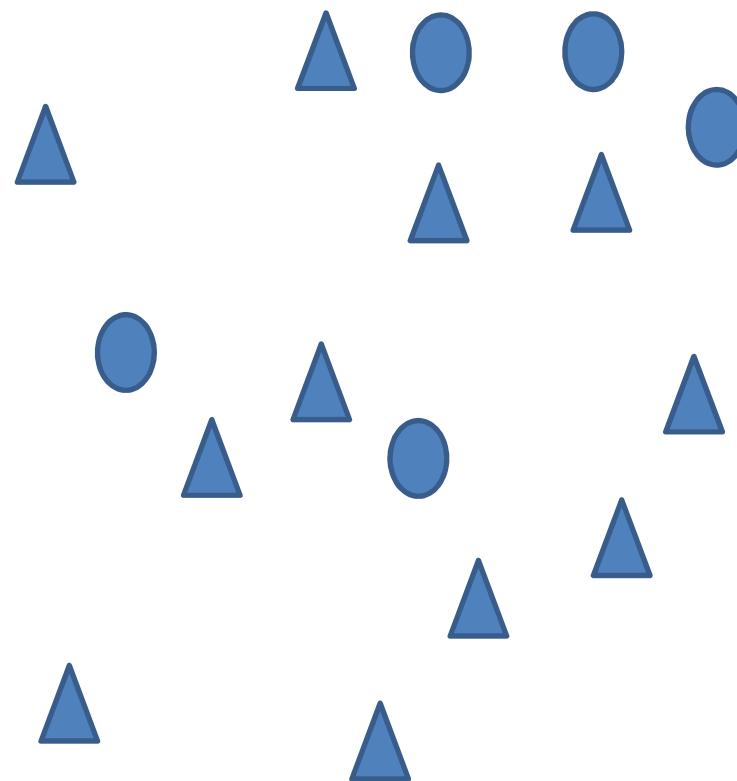
# 線性迴歸



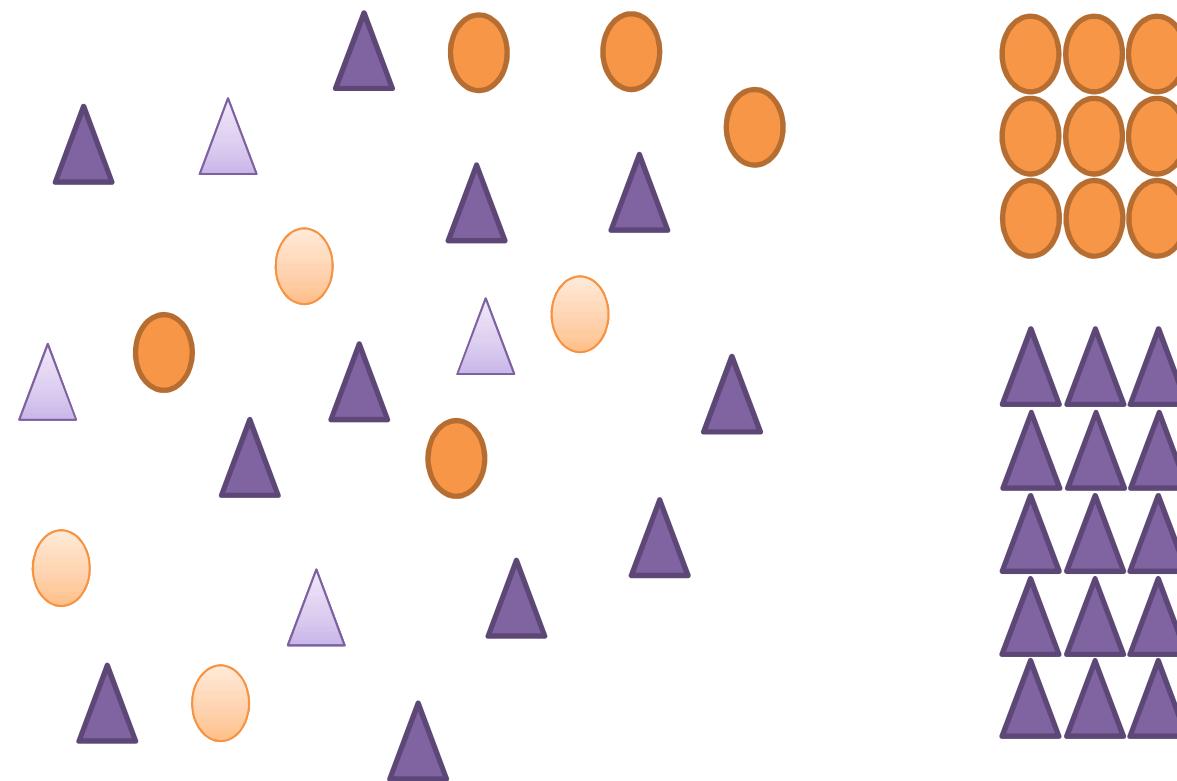
# 分類



# 分類



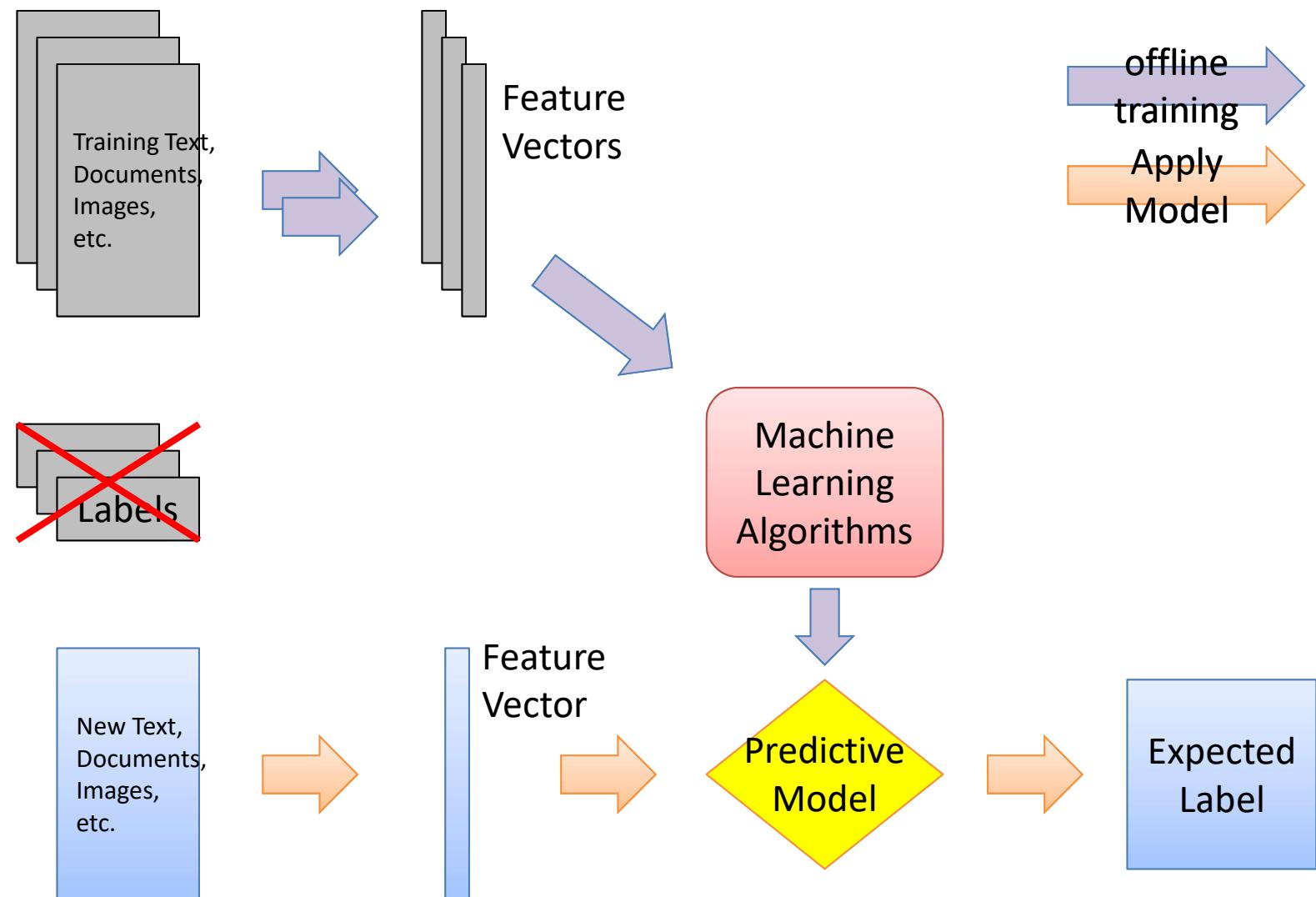
# 分類



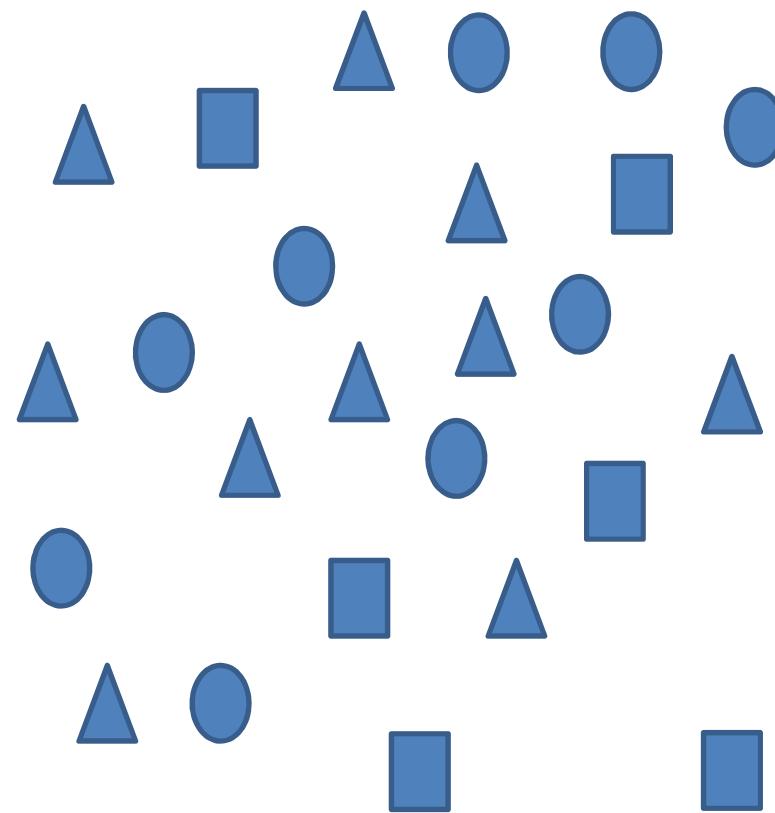
# 無監督式學習

- 與監督學習相比，訓練集沒有人為標註的結果。
  - 聚類/分群(Clustering)
    - hierarchical clustering
    - k-平均數(k-means)
    - Density-based spatial clustering of applications with noise (DBSCAN)
  - 異常偵測(Anomaly Detection)
  - 主要成分分析(Principal Component Analysis , PCA)

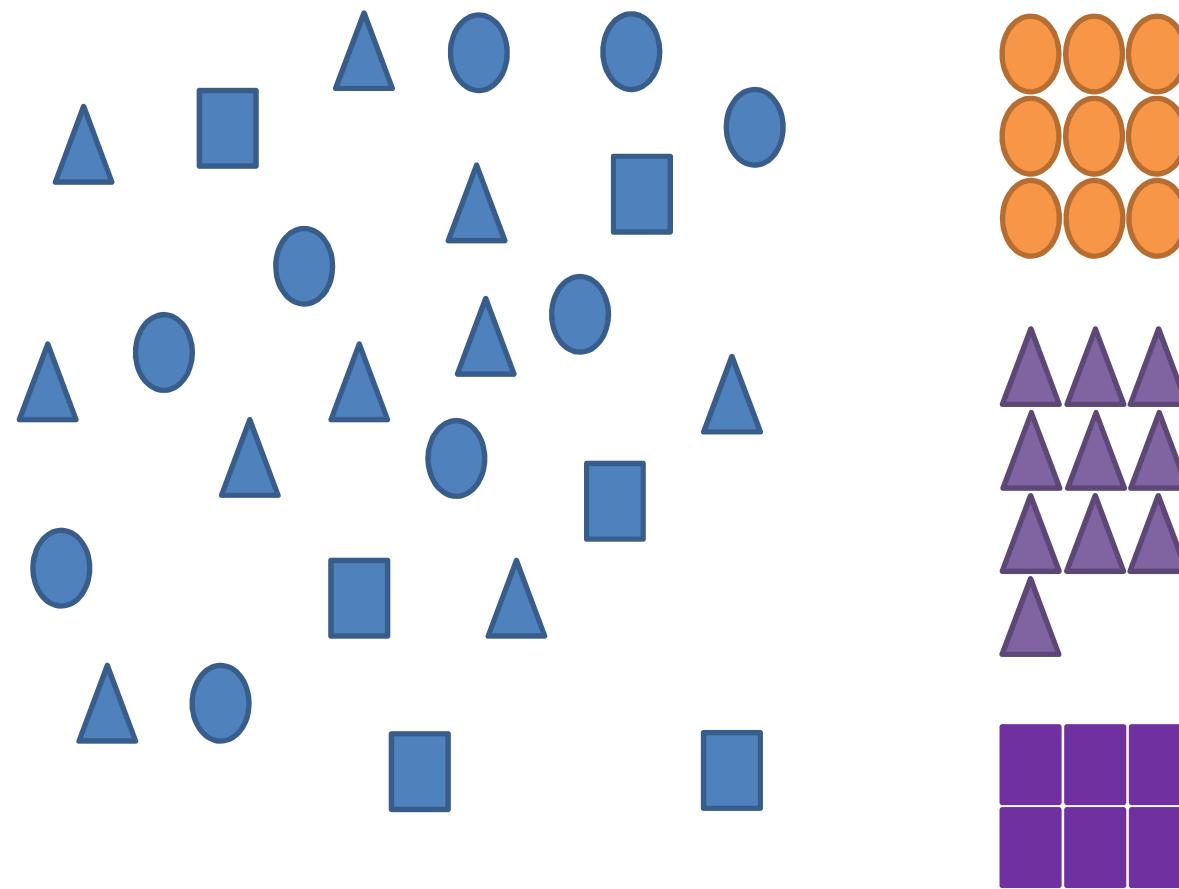
# 無監督式學習流程



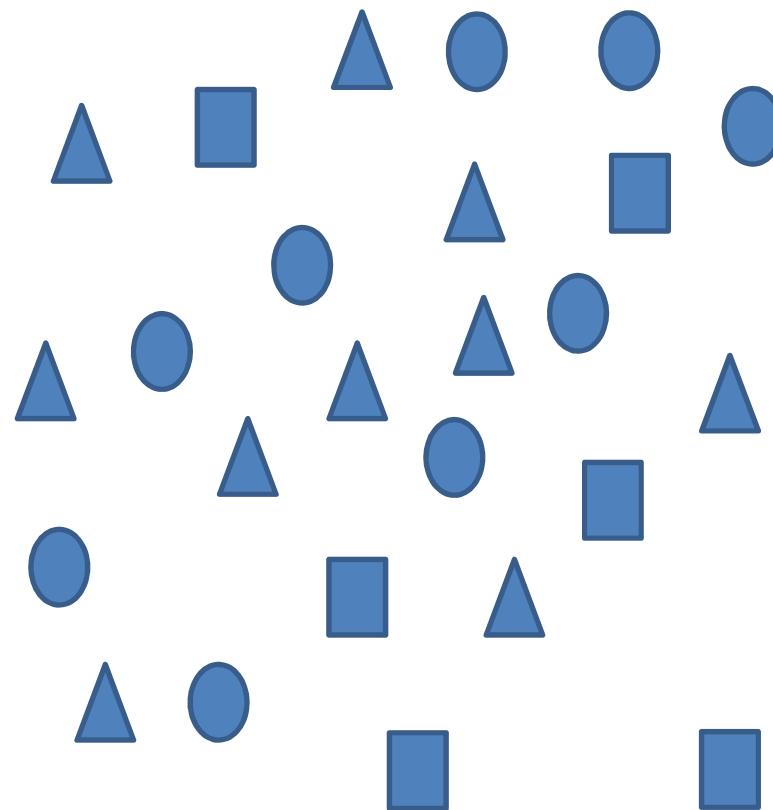
# 分群



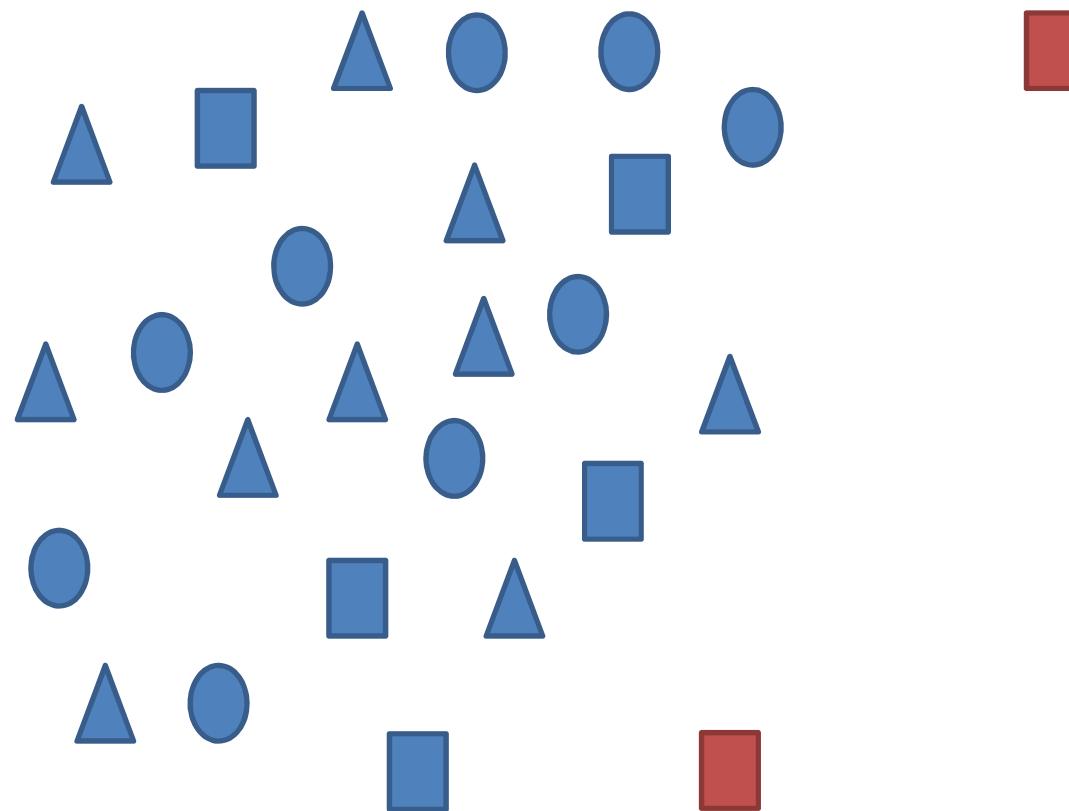
# 分群



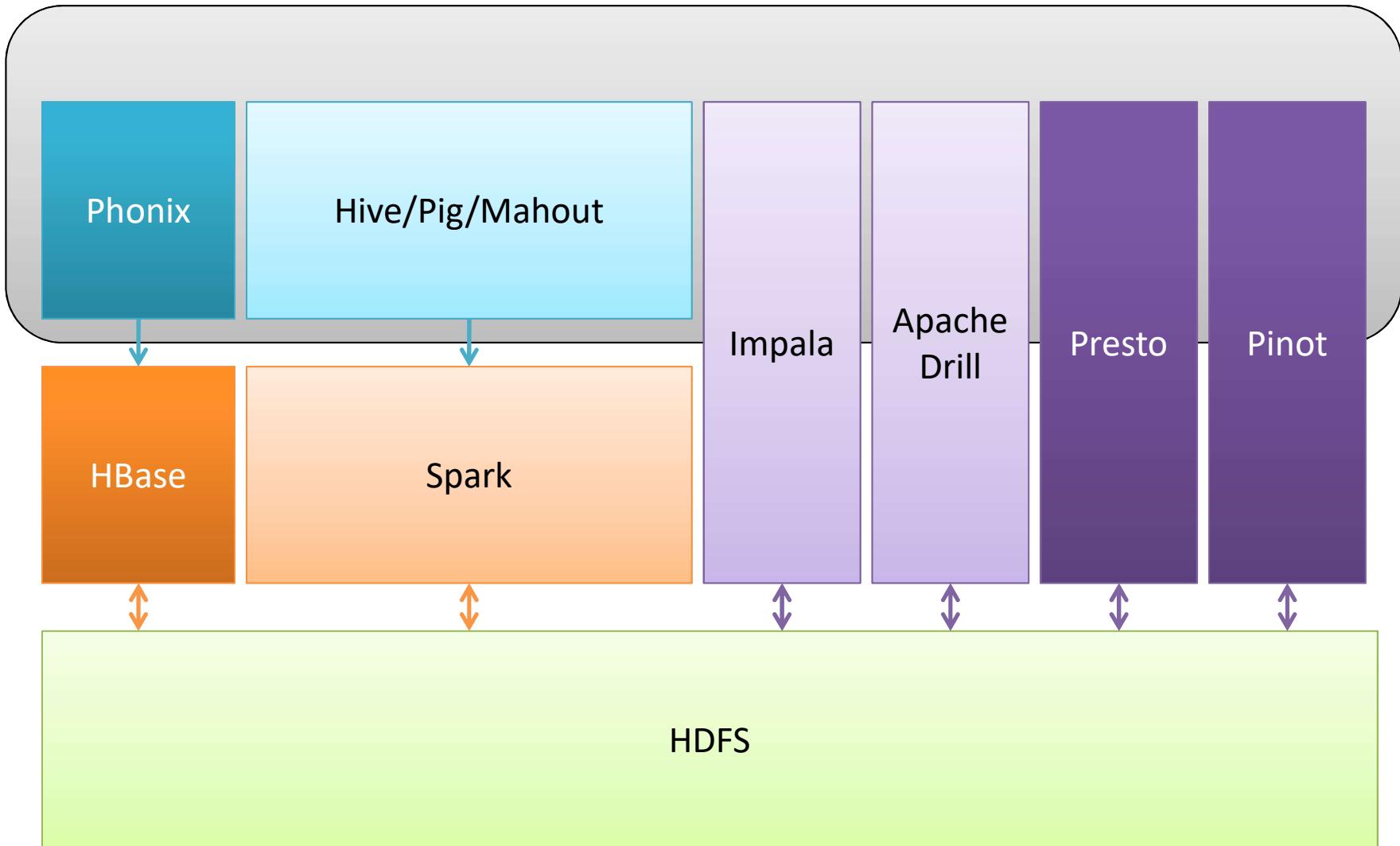
# 異常偵測



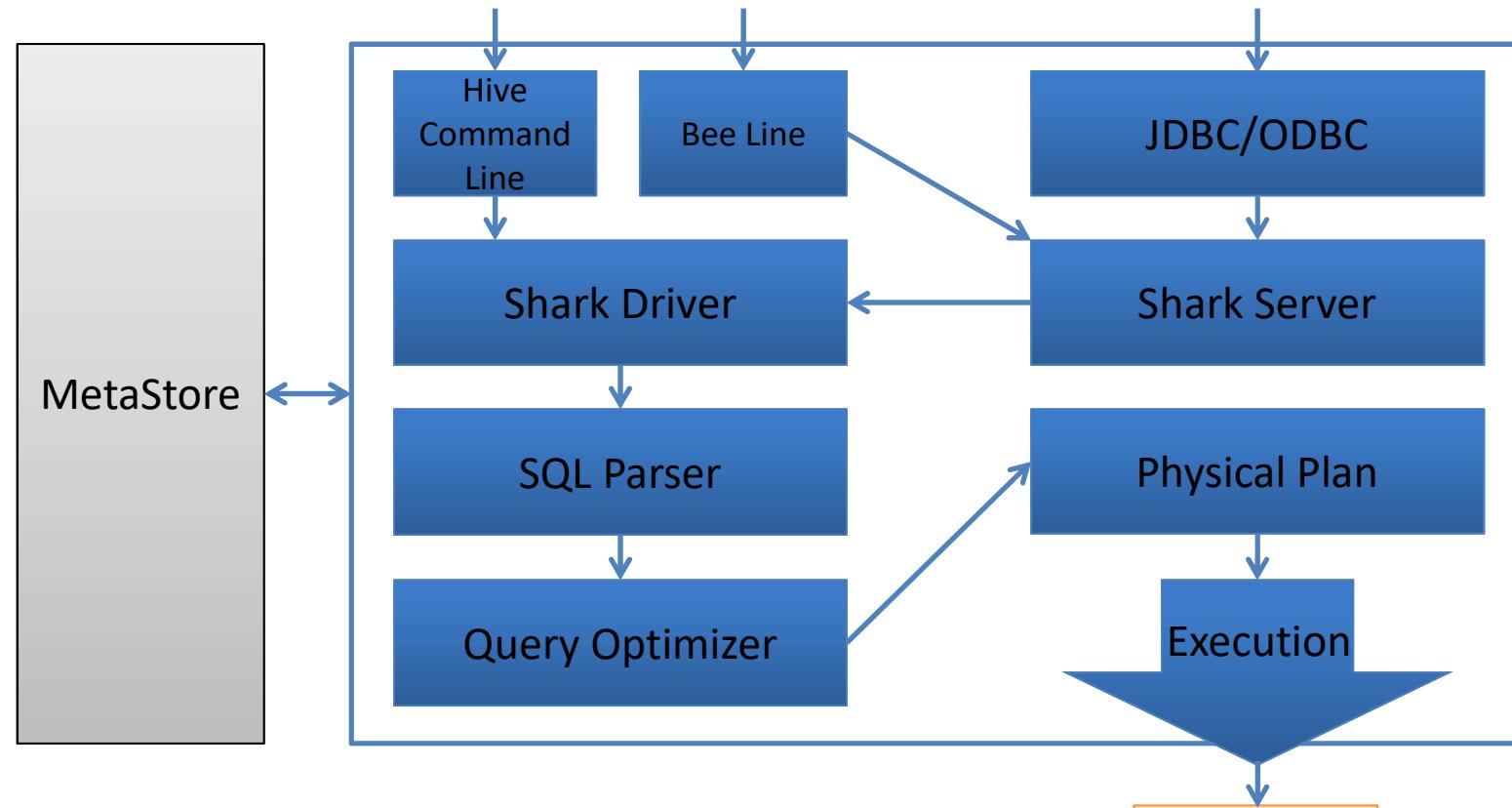
# 異常偵測



## SQL on Hadoop



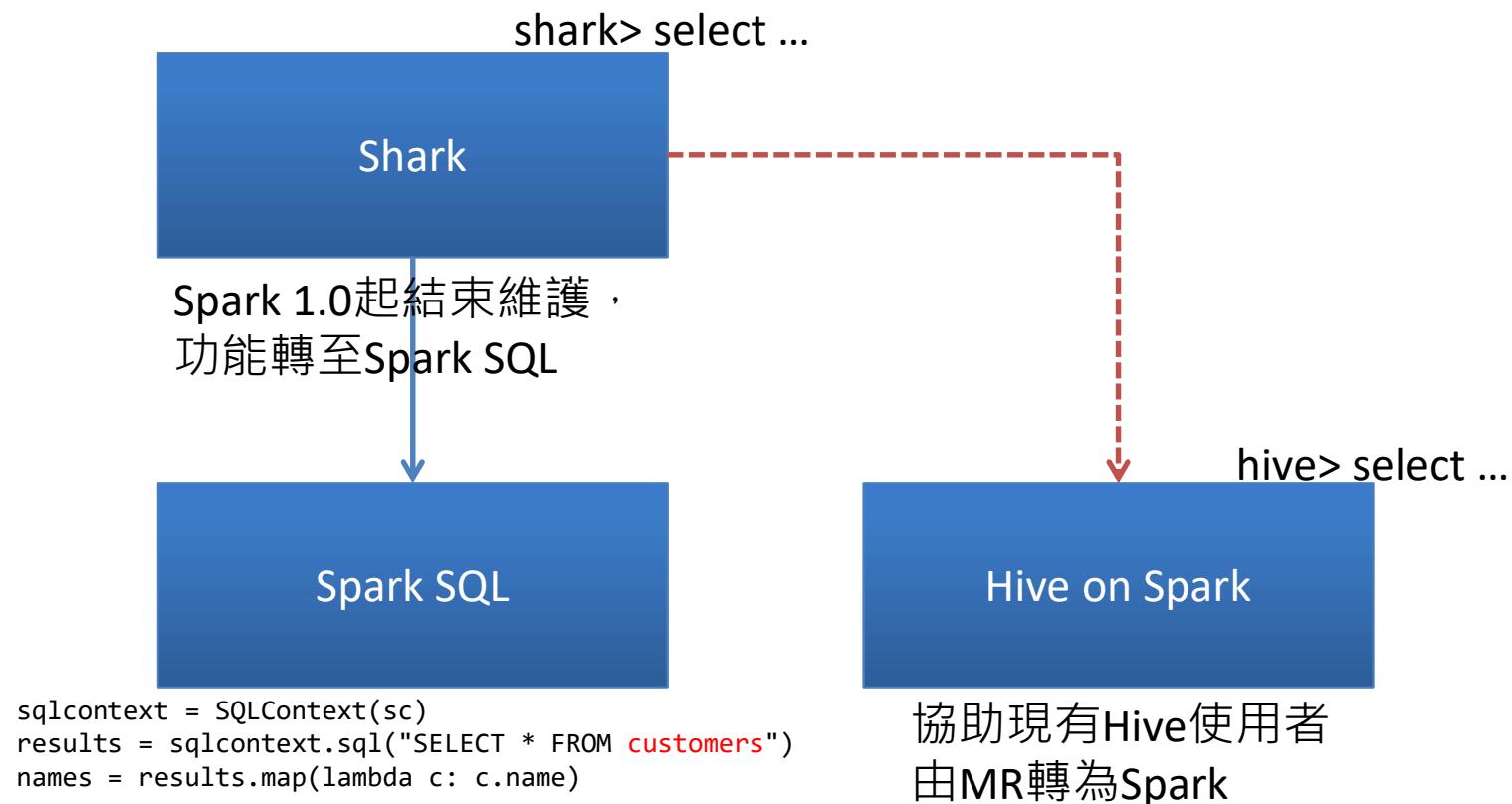
# Shark



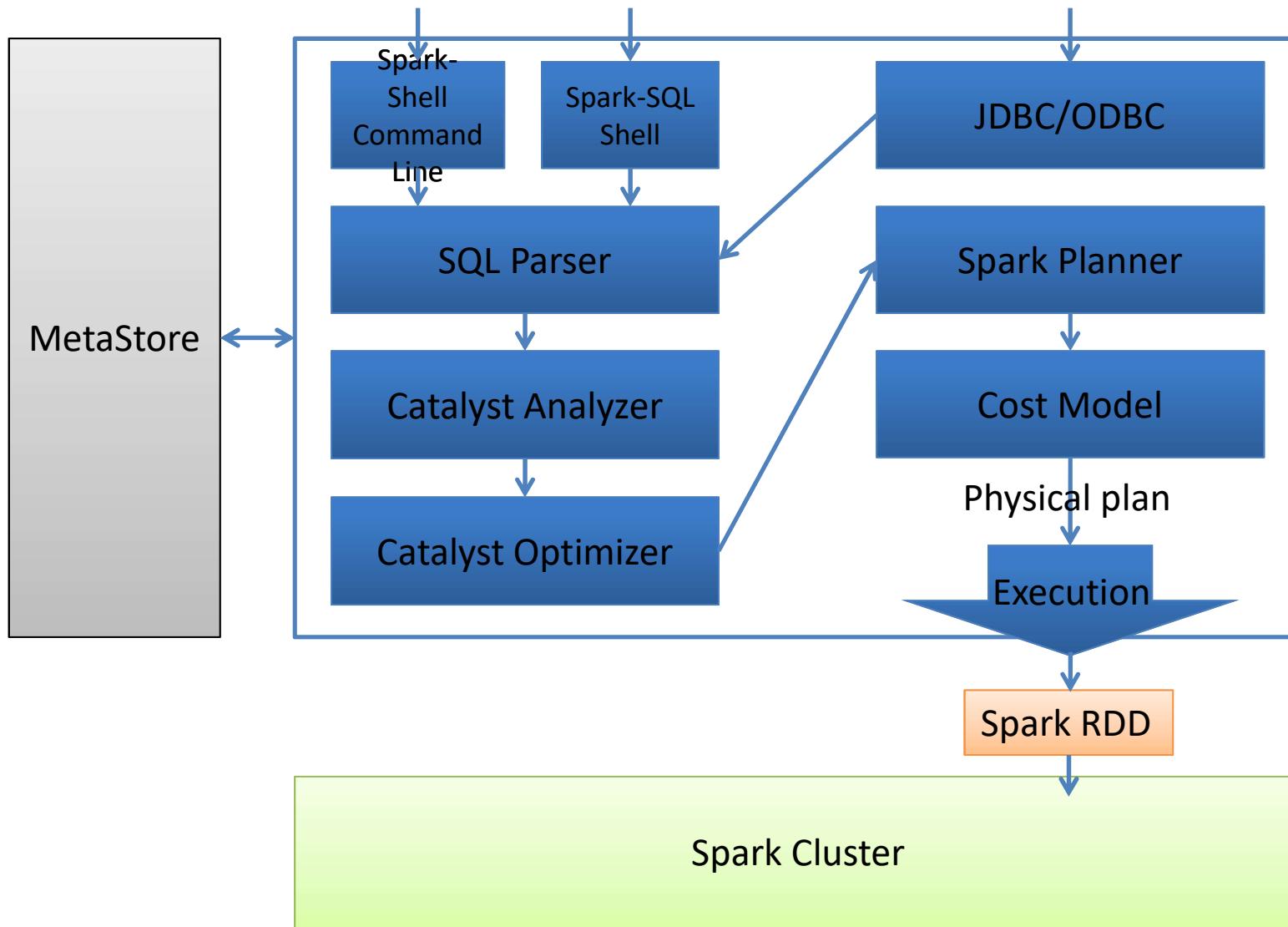
Spark 1.0後停止維護，  
相關功能轉由Spark SQL提供。

Spark Cluster

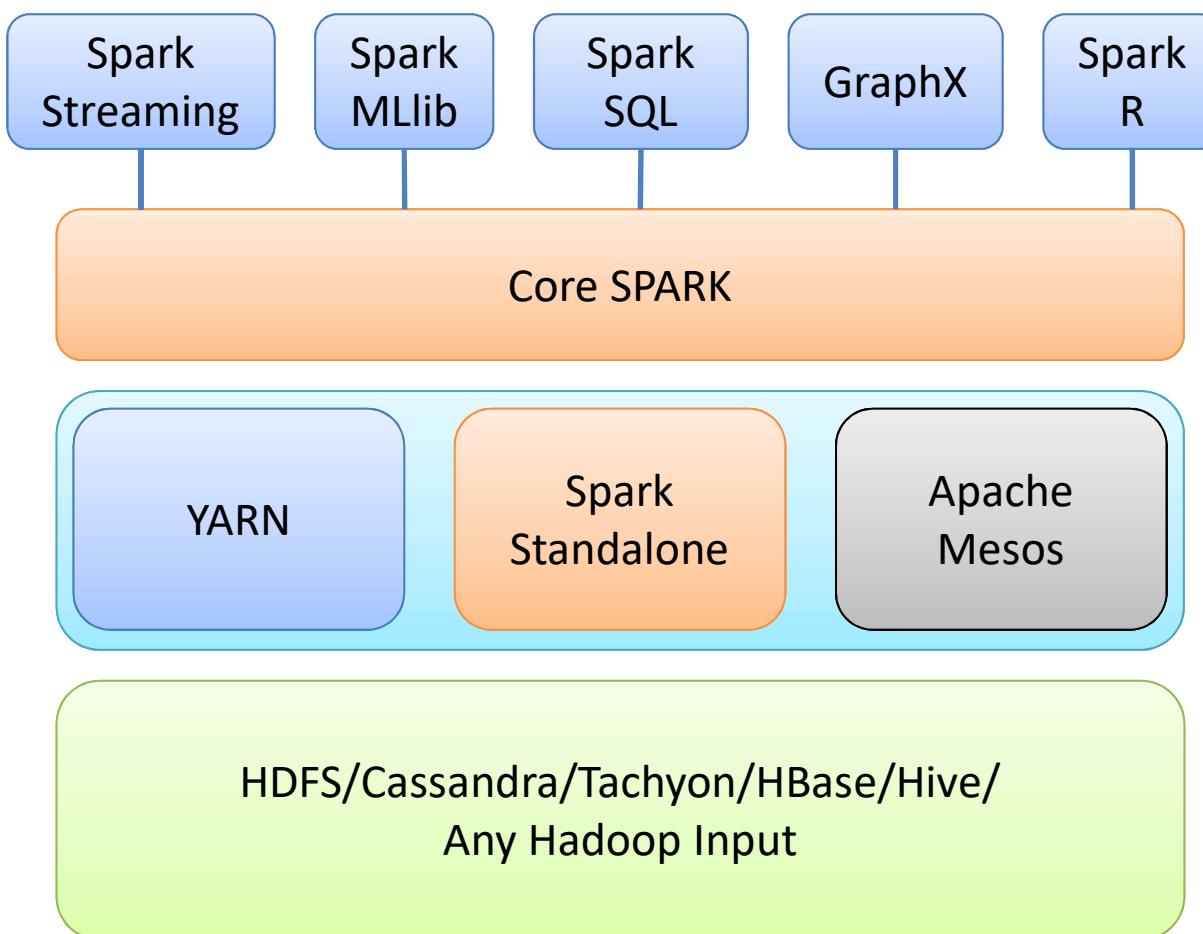
# SQL on Spark



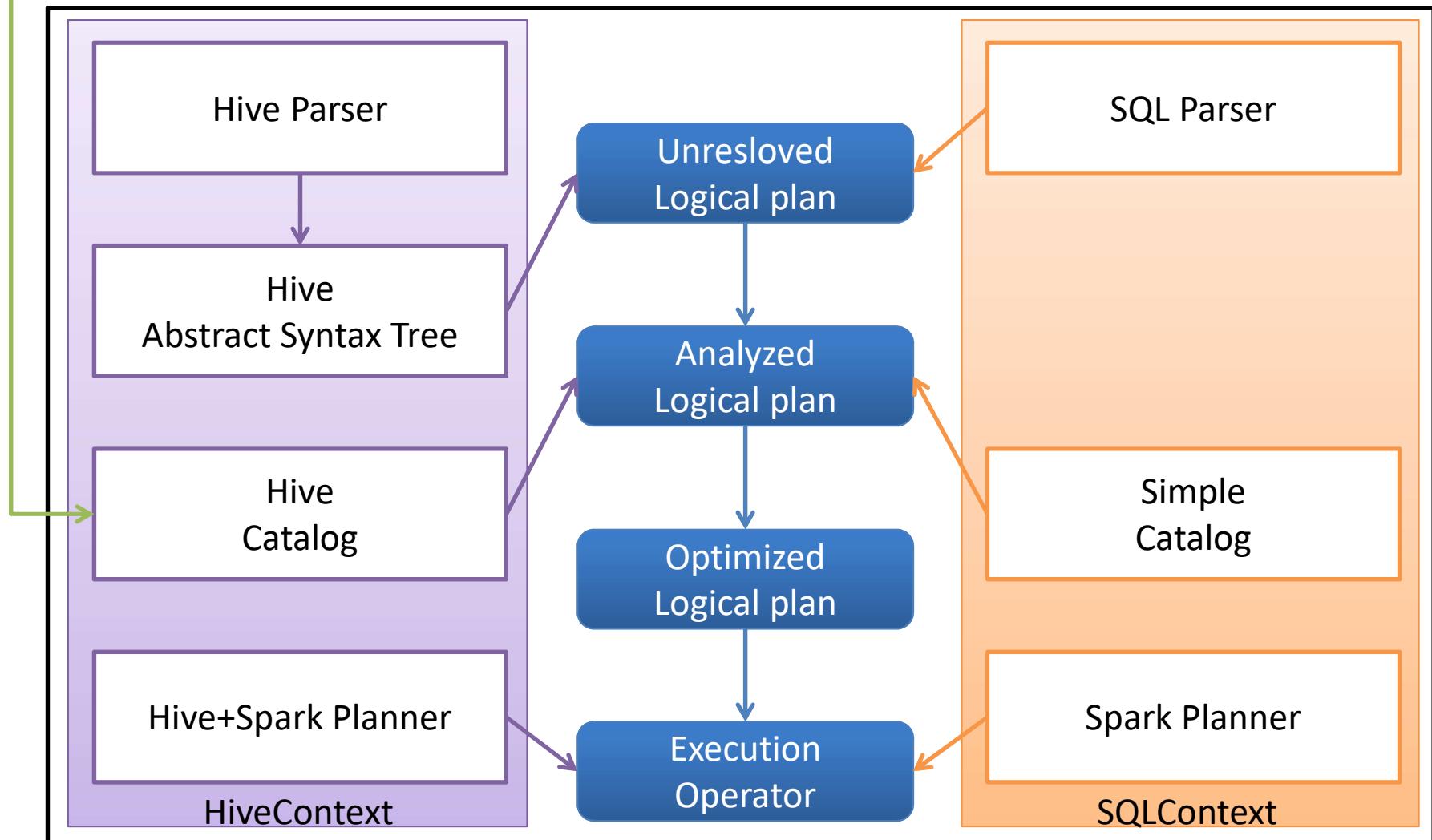
# SparkSQL

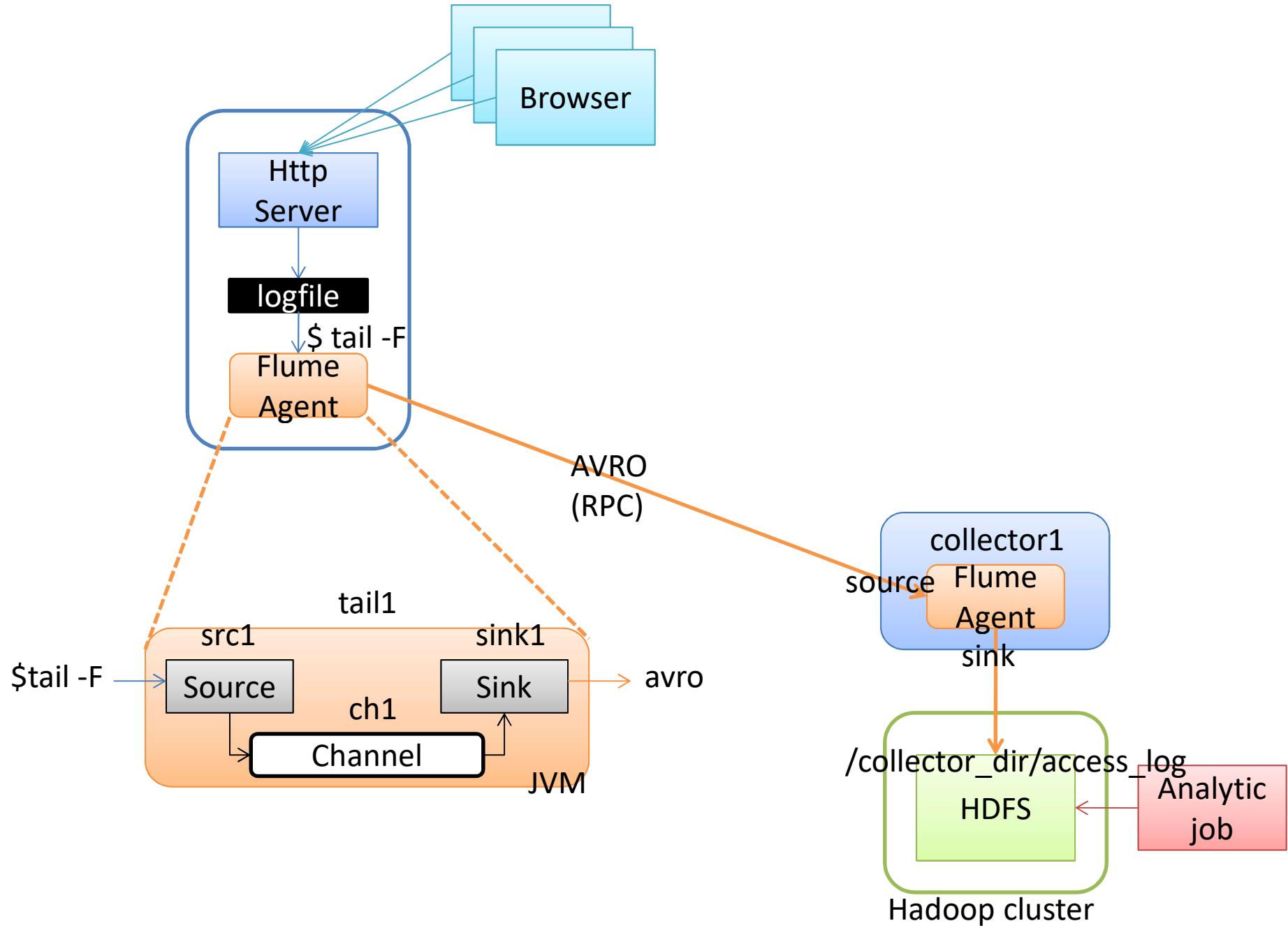


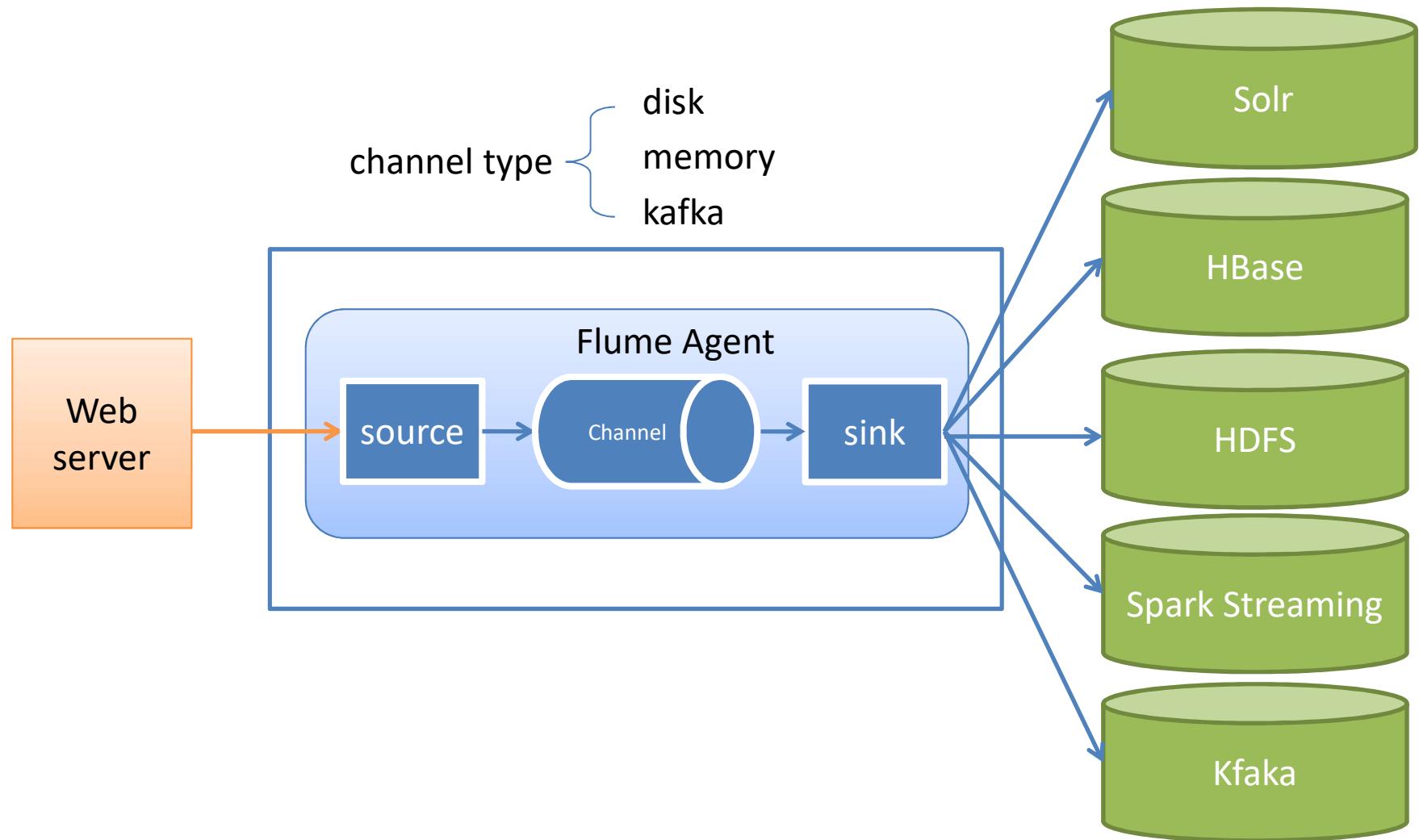
# SPARK



# Catalyst Context







# Spark Streaming

- Spark Streaming是core Spark API的延伸，對即時的資料流提供縮放性、高吞吐量、容錯的串流處理。
- Spark Streaming可以處理來自多種資料源的串流資料
  - Kafka、Flume、HDFS/S3、Kinesis、Twitter
- Spark Streaming接收到即時的資料串流輸入，將串流資料切分為批次資料，然後交由Spark引擎處理產生最終結果。



資料來源：<http://spark.apache.org/docs/latest/streaming-programming-guide.html>

# DStream(Discretized Stream)

- DStream用來表示一個連續的資料串流
- DStream可以看作是RDD的順序
- 建立Dstream有兩種方式：
  - 由資料源接受資料串流
    - Kafka、Flume、Kinesis
  - 原本的Dstream經過Spark Streaming的高階操作而生成
    - Scala、Java、Python

# StreamingContext()

- StreamingContext是Spark Streaming的入口點
- StreamingContext用來與Spark Cluster通信，  
以及建立Dstream。
- 建立StreamingContext需要通過SparkContext

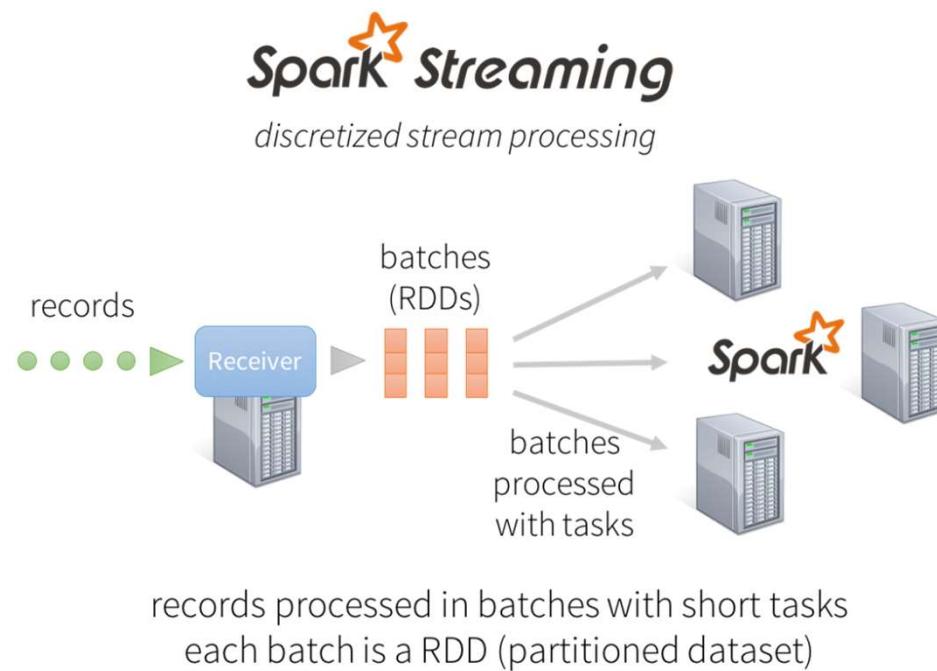
```
from pyspark import SparkContext  
from pyspark.streaming import StreamingContext
```

```
sc = SparkContext()  
ssc = StreamingContext(sc,1)
```

# Processing Time and Event Time

- Event Time
  - 資料生成時間，通常內含在資料本身。
    - 例如ETC記錄：車號ABC-123的車輛於12:01:50通過閘門1
    - 12:01:50即為Event Time
- Processing Time
  - 資料被處理的時間
    - 閘門1的資料於12:01:59秒送到Spark Streaming進行處理
    - 12:01:59則是Processing Time
- Spark Streaming基於micro-batch模式，將串流資料依照接受時間每幾秒切分為一個batch資料。所以Spark Streaming是基於Processing Time進行切分。

# Spark Streaming 架構



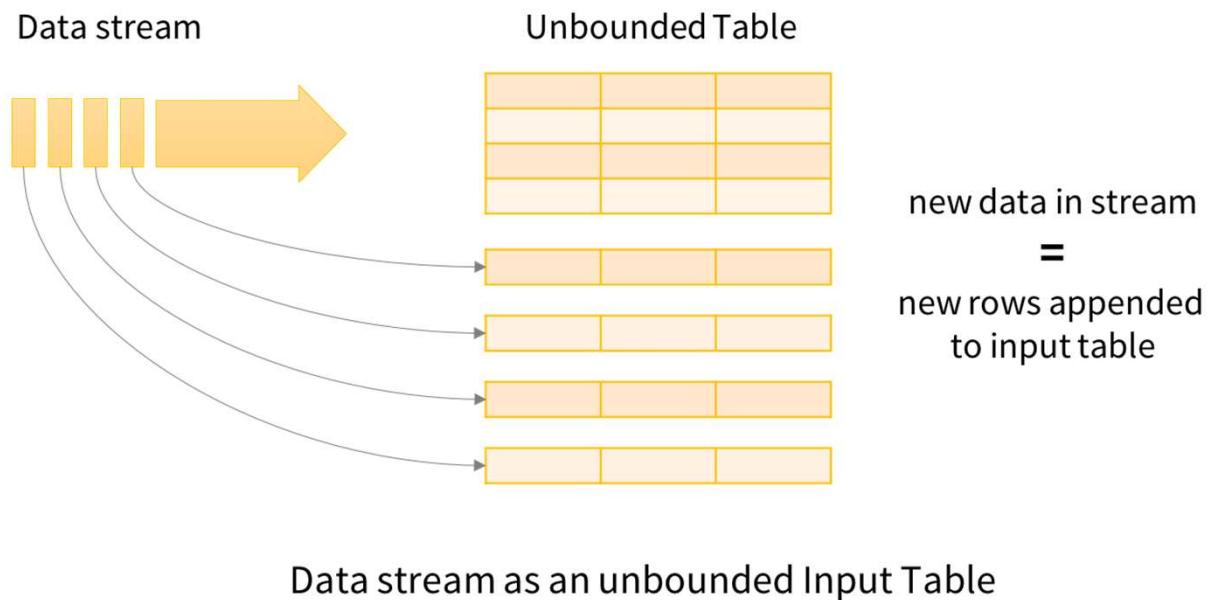
# Structured Streaming



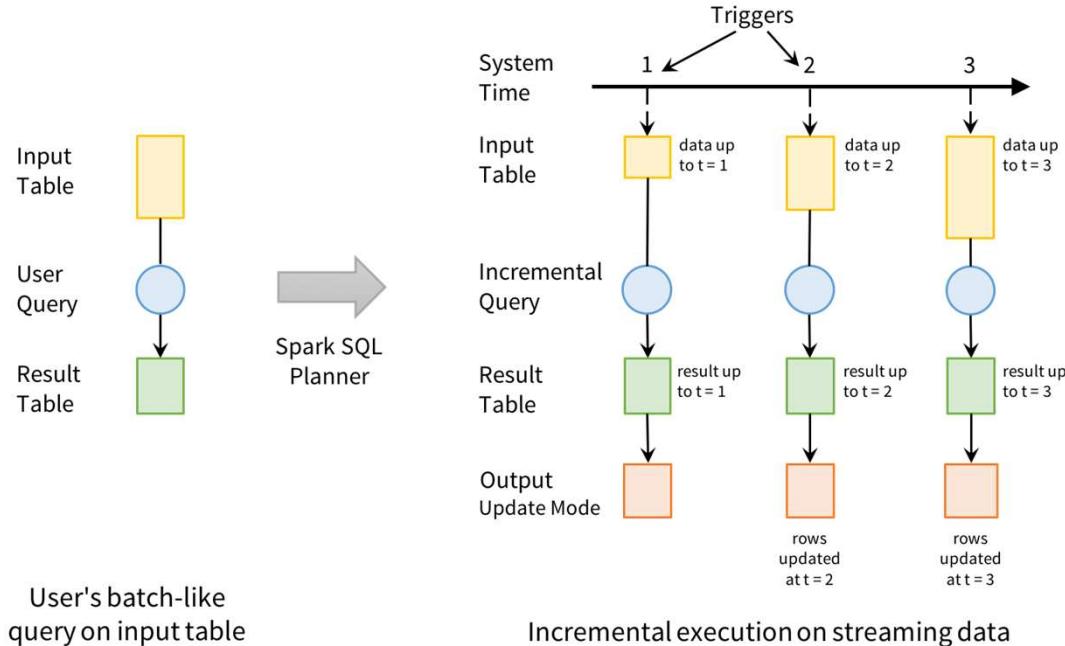
- 2016年於Spark 2.0開始引入
  - 可以區分processing time與event time
  - 使用relational execution engine提升效能
- Declarative(宣告式) API
  - Structured Streaming
- Imperative(命令式) API
  - Spark Streaming

# Structured Streaming的優點

- Incremental query model
  - 可以執行增量查詢。且程式碼與批次操作 (RDD) API 完全一樣 (基於 DataFrame 與 DataSet API)
- Support for end-to-end application
  - Structured Streaming 與內建的 connector 使得 end-to-end 程式非常容易撰寫，而且“correct by default”。
- Reuse Spark SQL engine
  - 通過 DataFrame / DataSet API 讓 Structured Streaming 可以直接利用 Spark SQL 之前所做的最佳化努力，如執行計畫優化、codegen、記憶體管理等。



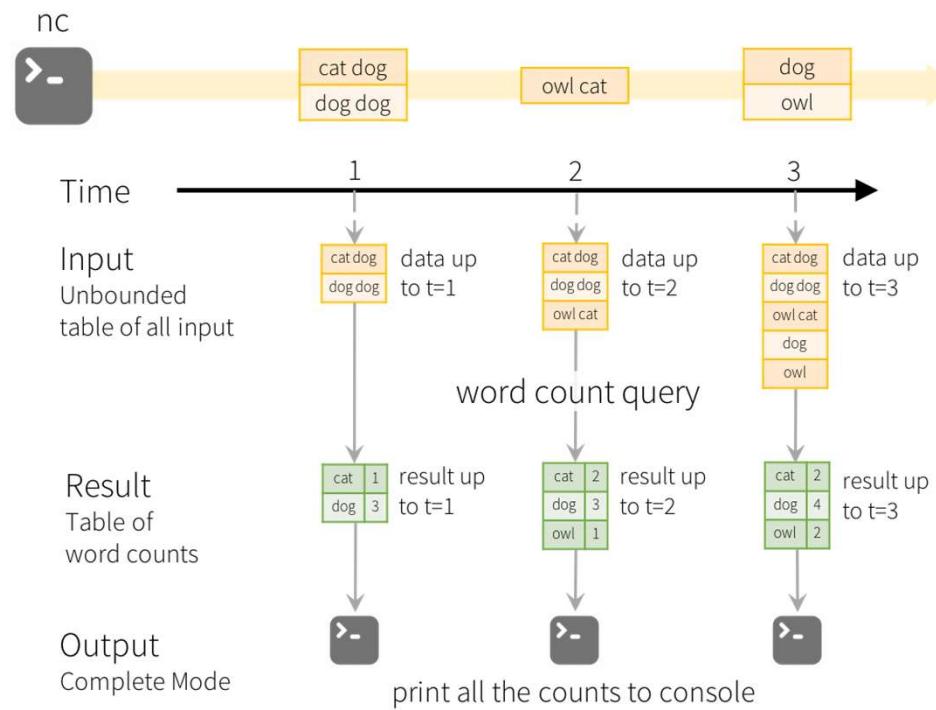
資料來源：<https://spark.apache.org/docs/latest/structured-streaming-programming-guide.html>



### Structured Streaming Processing Model

Users express queries using a batch API; Spark incrementalizes them to run on streams

資料來源：<https://spark.apache.org/docs/latest/structured-streaming-programming-guide.html>



Model of the Quick Example

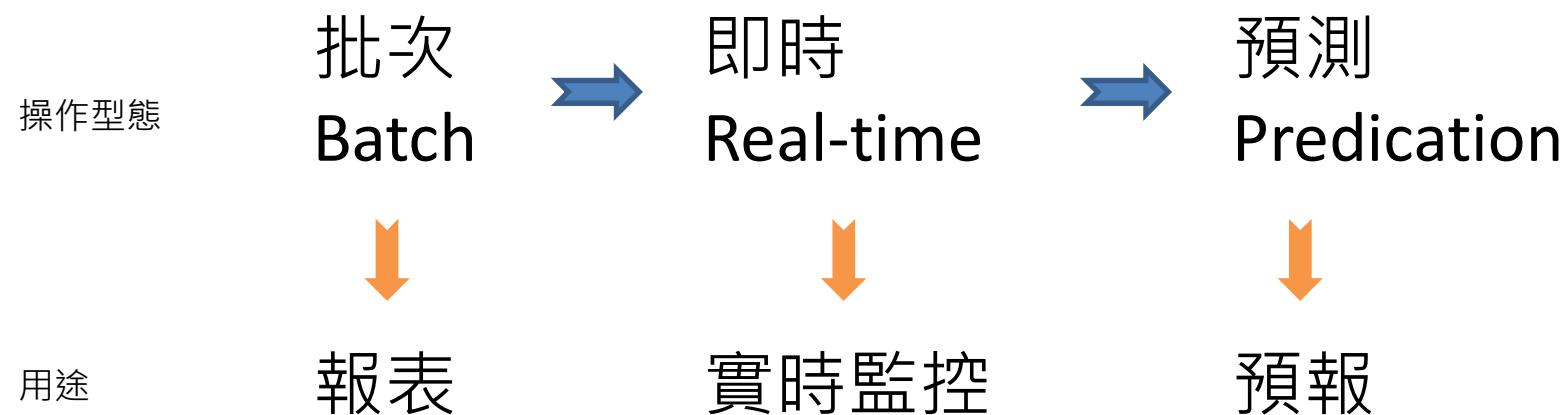
資料來源：<https://spark.apache.org/docs/latest/structured-streaming-programming-guide.html>

# 與其他串流引擎的比較

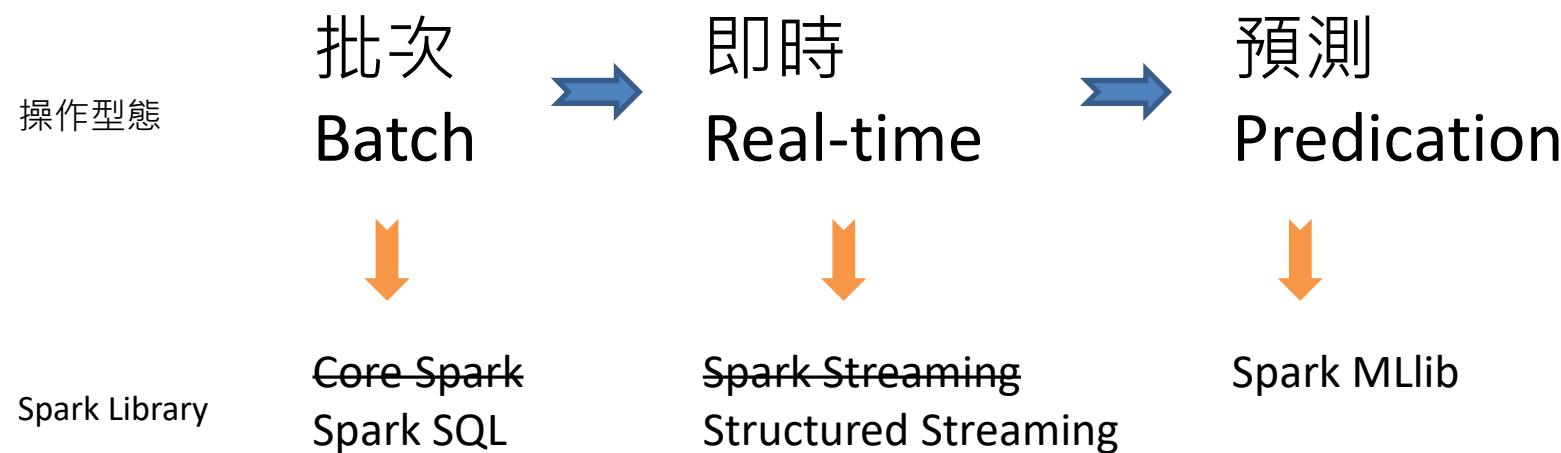
Property	Structured Streaming	Spark Streaming	Apache Storm	Apache Flink	Kafka Streams	Google Dataflow
Streaming API	incrementalize batch queries	integrates with batch	separate from batch	separate from batch	separate from batch	integrates with batch
Prefix Integrity Guarantee	✓	✓	✗	✗	✗	✗
Internal Processing	exactly once	exactly once	at least once	exactly once	at least once	exactly once
Transactional Sources/Sinks	✓	some	some	some	✗	✗
Interactive Queries	✓	✓	✗	✗	✗	✗
Joins with Static Data	✓	✓	✗	✗	✗	✗

資料來源：<https://databricks.com/blog/2016/07/28/structured-streaming-in-apache-spark.html>

# 資料用途的演化



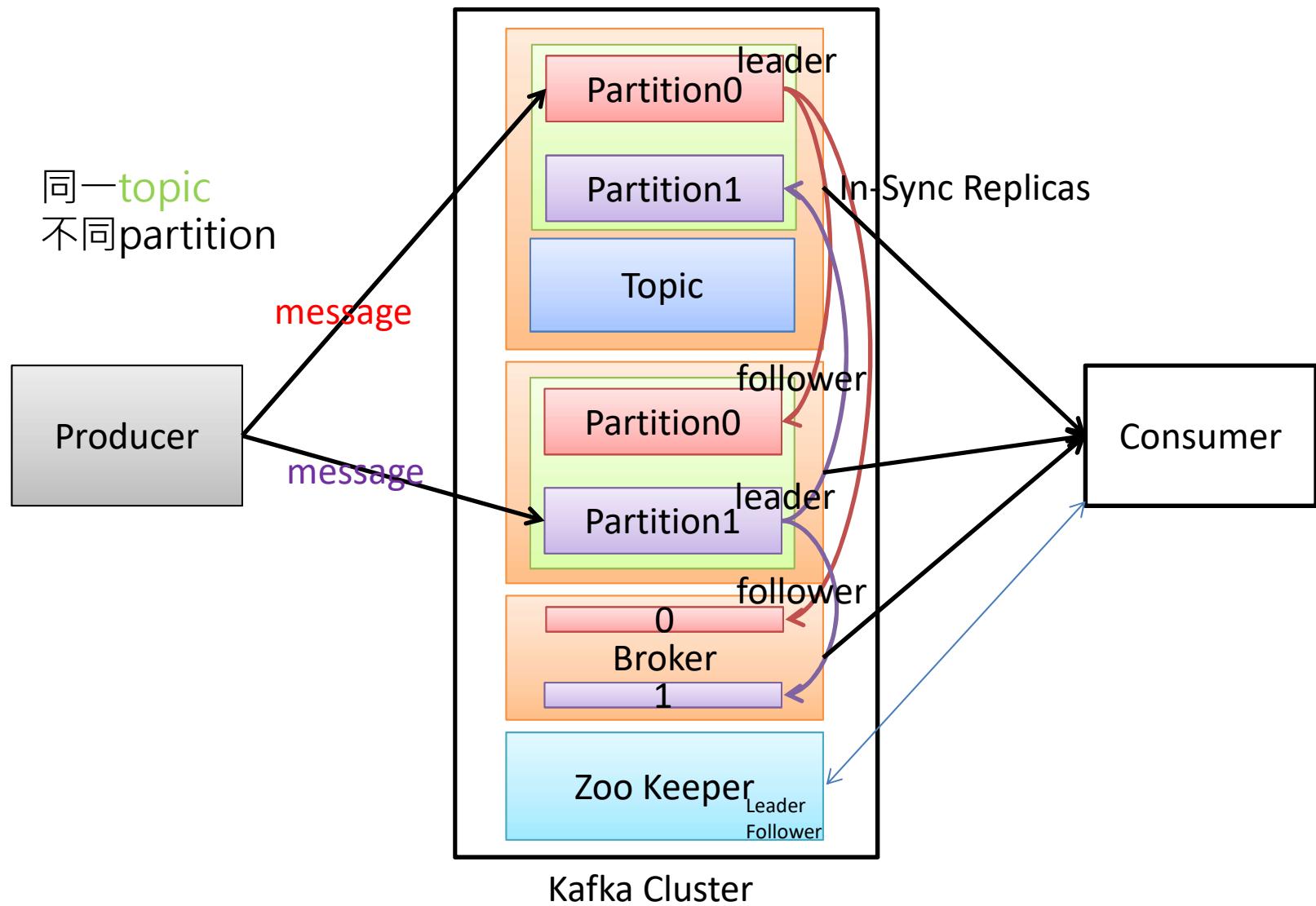
# Spark程式庫與資料用途的對應



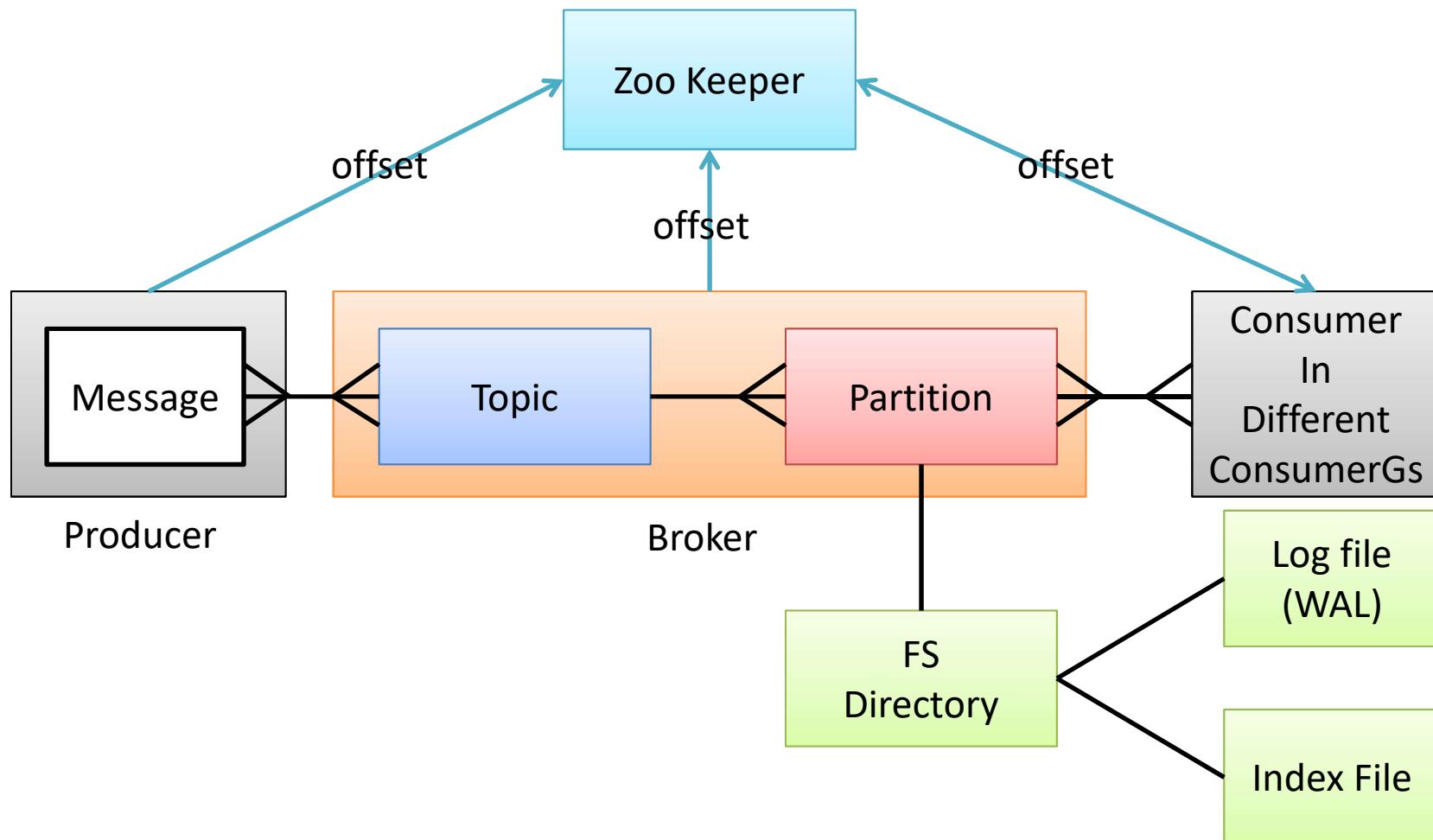
# Structured Streaming

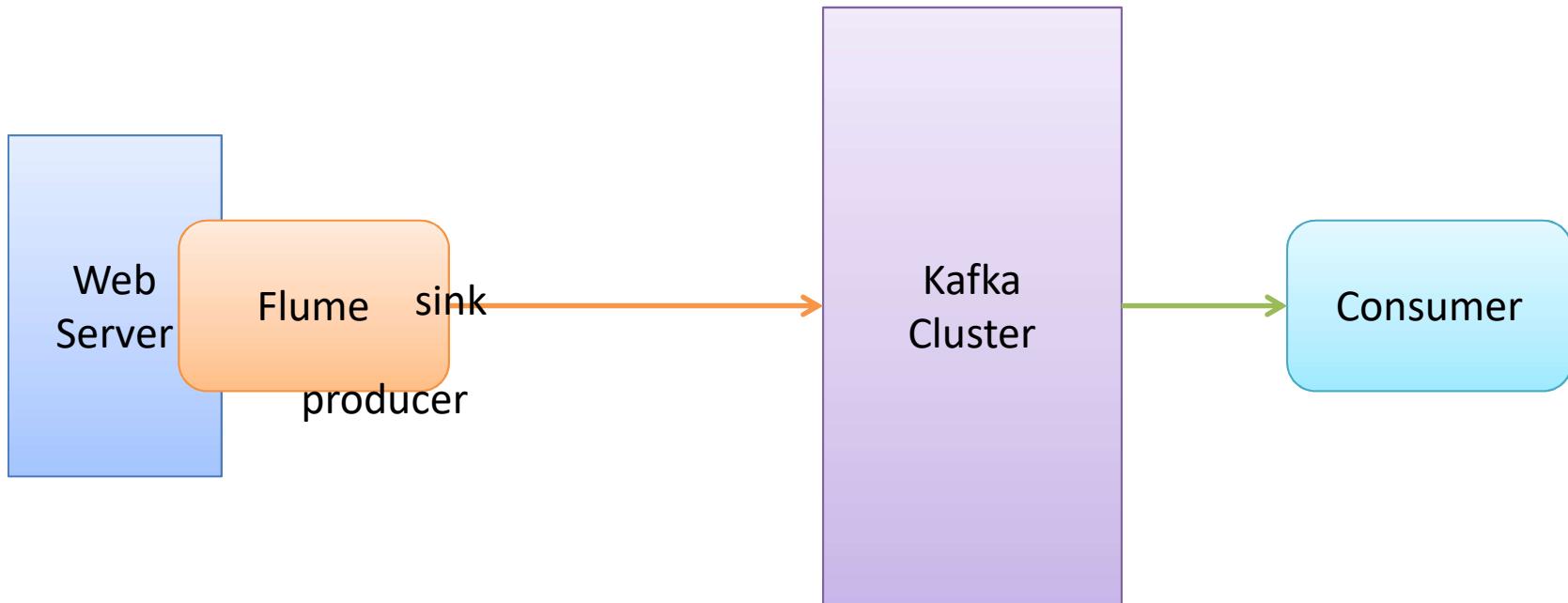
- Spark 2.0開始引入，用來處理stream資料。
  - 底層基於SparkSQL與Catalyst Optimizer
  - 讓使用者可以向使用靜態Dataset開發方式處理stream資料
- 未來將取代Spark Streaming用來處理stream資料

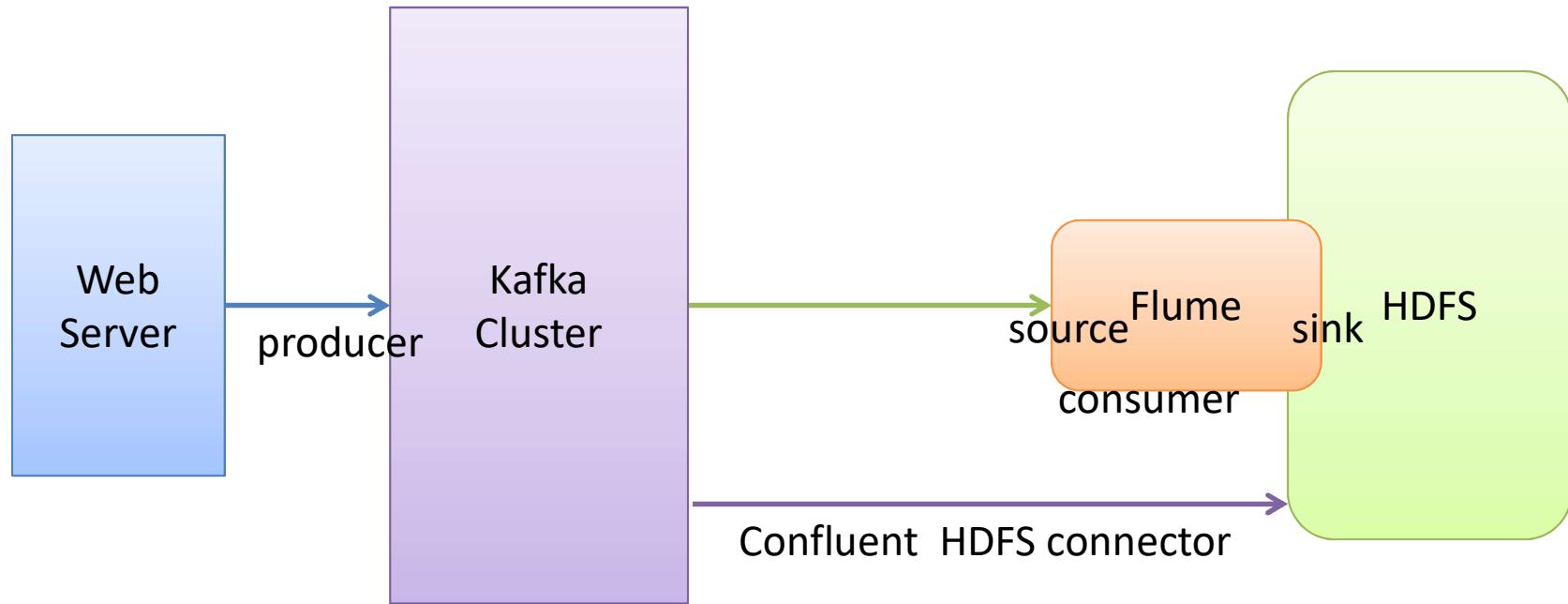
# Kafka Architecture



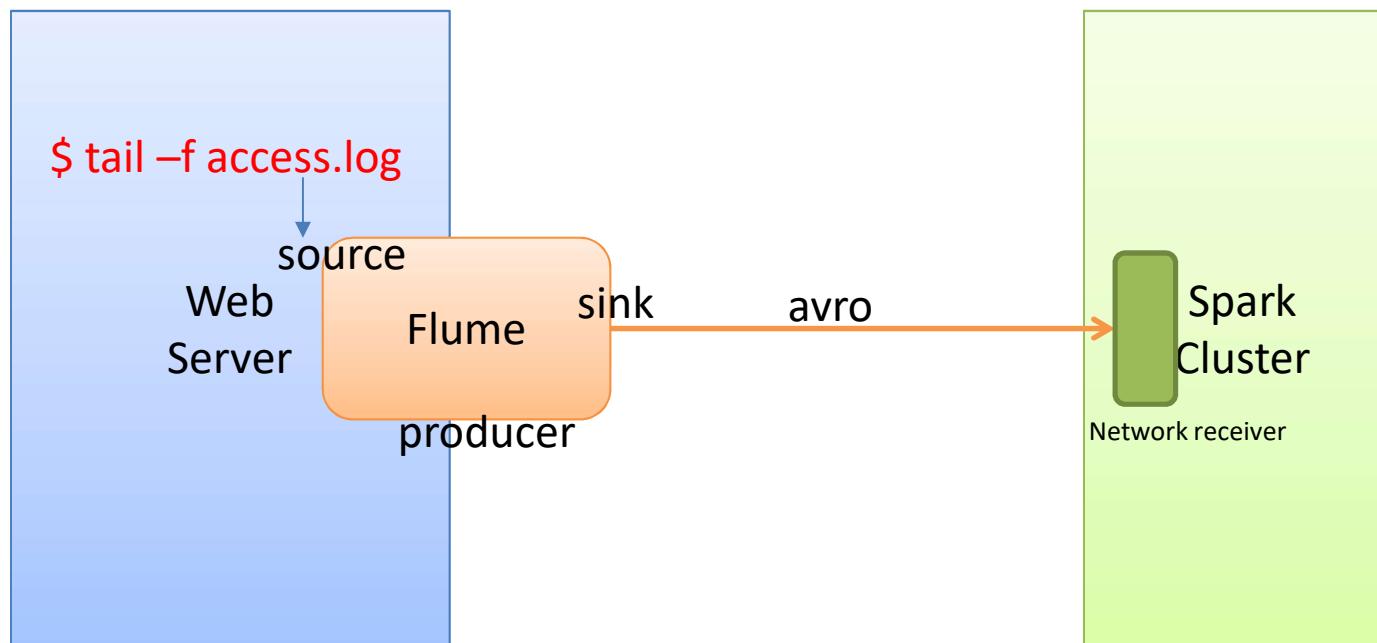
# Kafka terminology



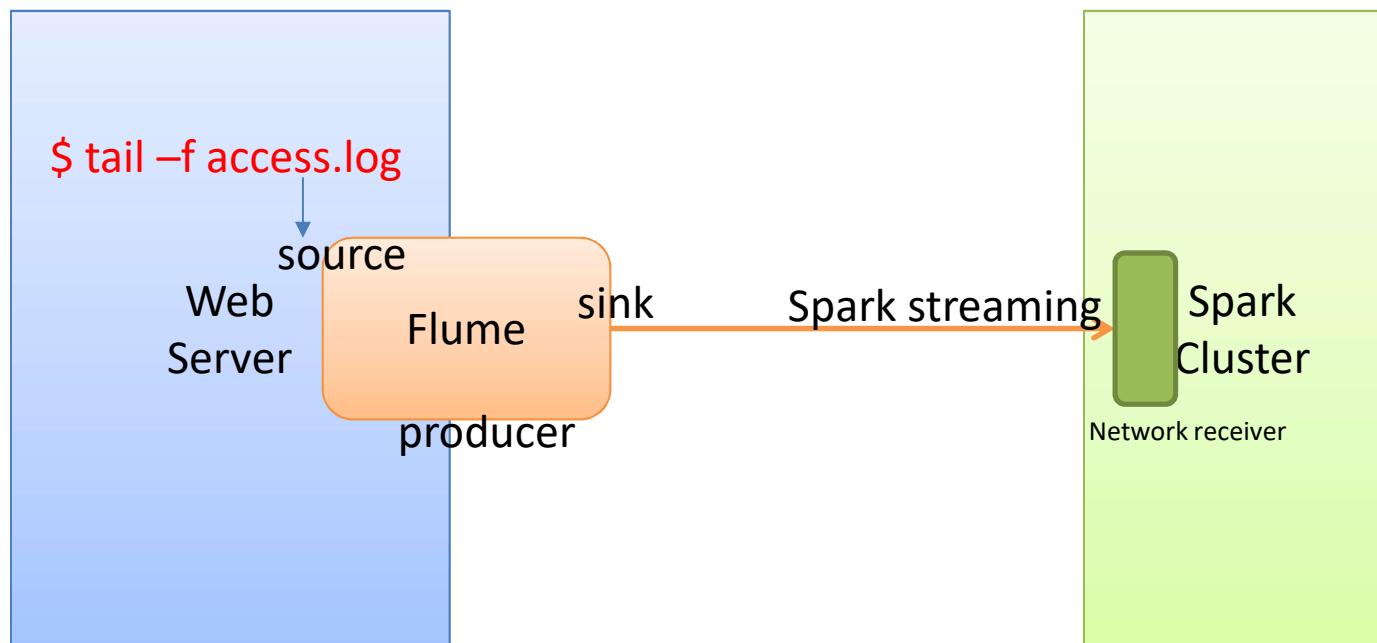




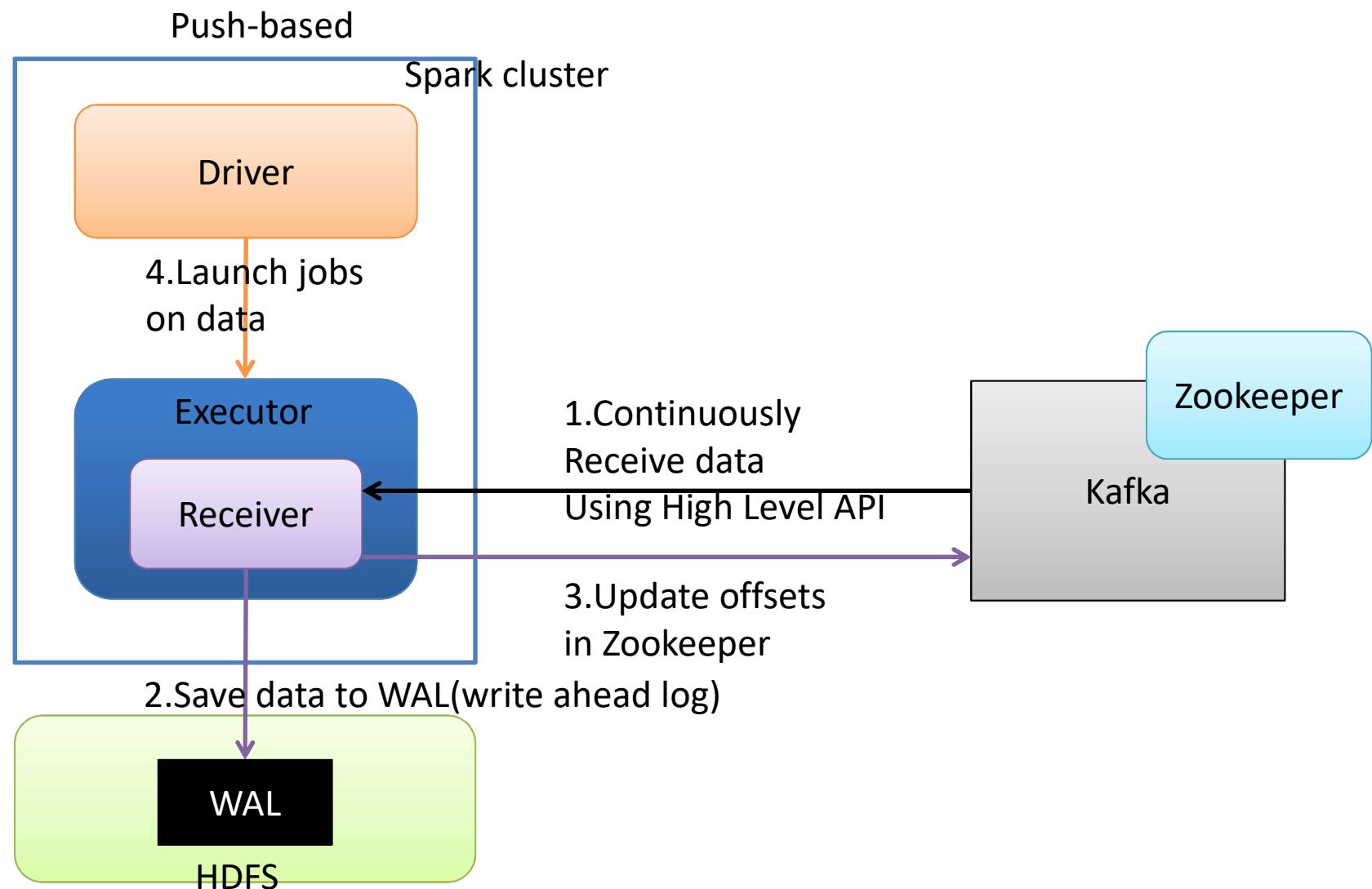
# PUSH



# PUSH



# Kafka with Receiver and WAL



# Kafka w/o Receiver and WAL

