

A complex network diagram with numerous nodes and connecting lines. Nodes are represented by circles and hexagons, some in red and some in grey. Lines are thin and light red, with some thicker grey lines forming a central structure. Dotted lines also connect some nodes.

BigData

Iniciando meu estudo com Grandes Dados

email: susana.bouchardet@gmail.com



redis



ZooKeeper



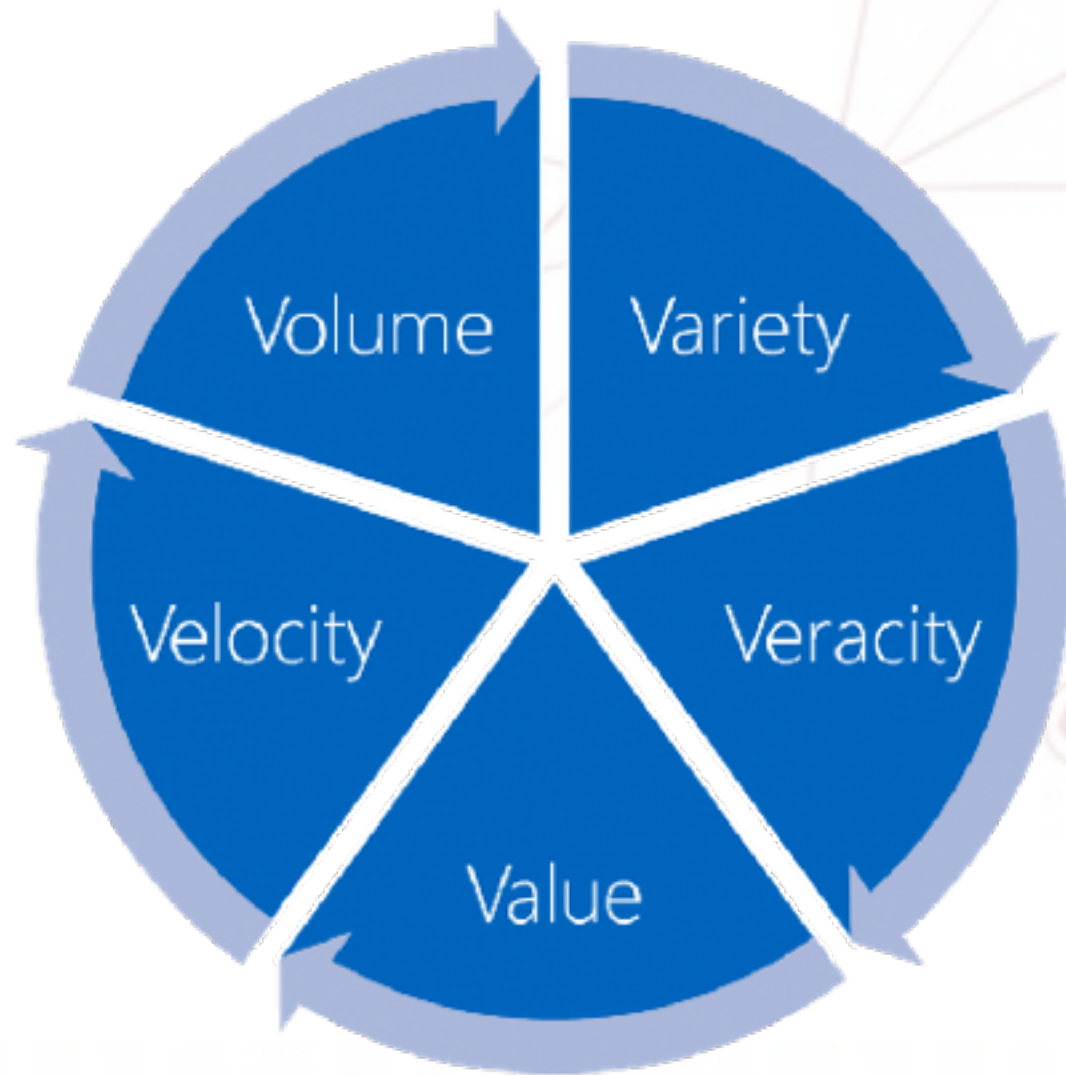
NOT SURE IF ITS TOO MUCH TO LEARN

OR I AM TOO STUPID

Troll.me

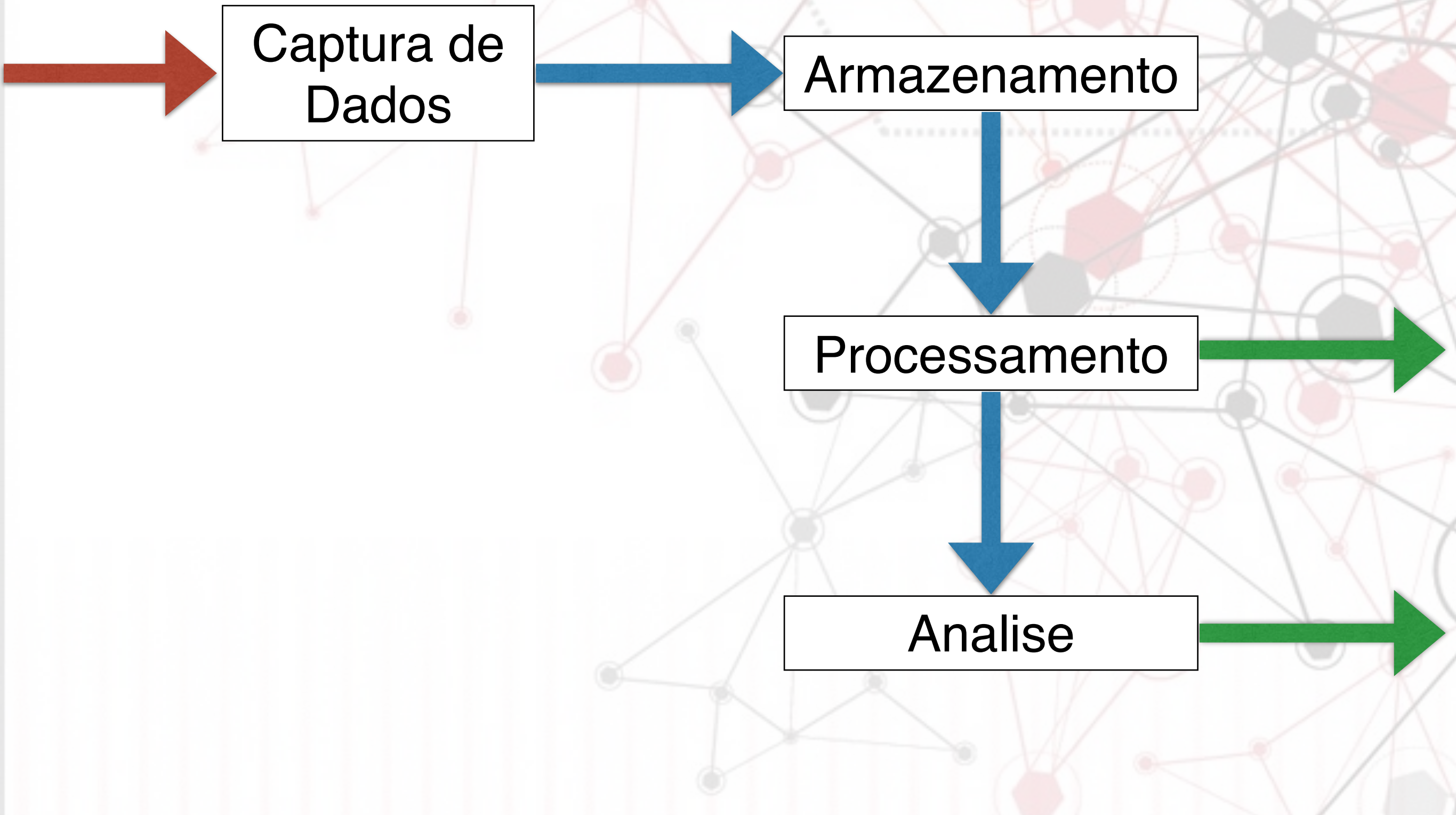
BigData

5V



- **Volume:** Muitos dados, desafio de armazenamento.
- **Variedade:** diferentes formatos de dados estruturados e/ou não.
- **Veracidade:** Acurácia e autenticidade.
- **Valor:** Retorno desses dados para o negócio.
- **Velocidade:** Tempo entre dado gerado e analisado.

Descrição do Problema



GFS

- Publicado em 2003
- Revolucionário
- arquivos em **chunks** de 64mb
- cada chunk esta no mínimo em 3 máquina



MapReduce



- Publicado em 2004
- modelo para processar grandes volumes de dados.
- Técnica muito usada em linguagens funcionais

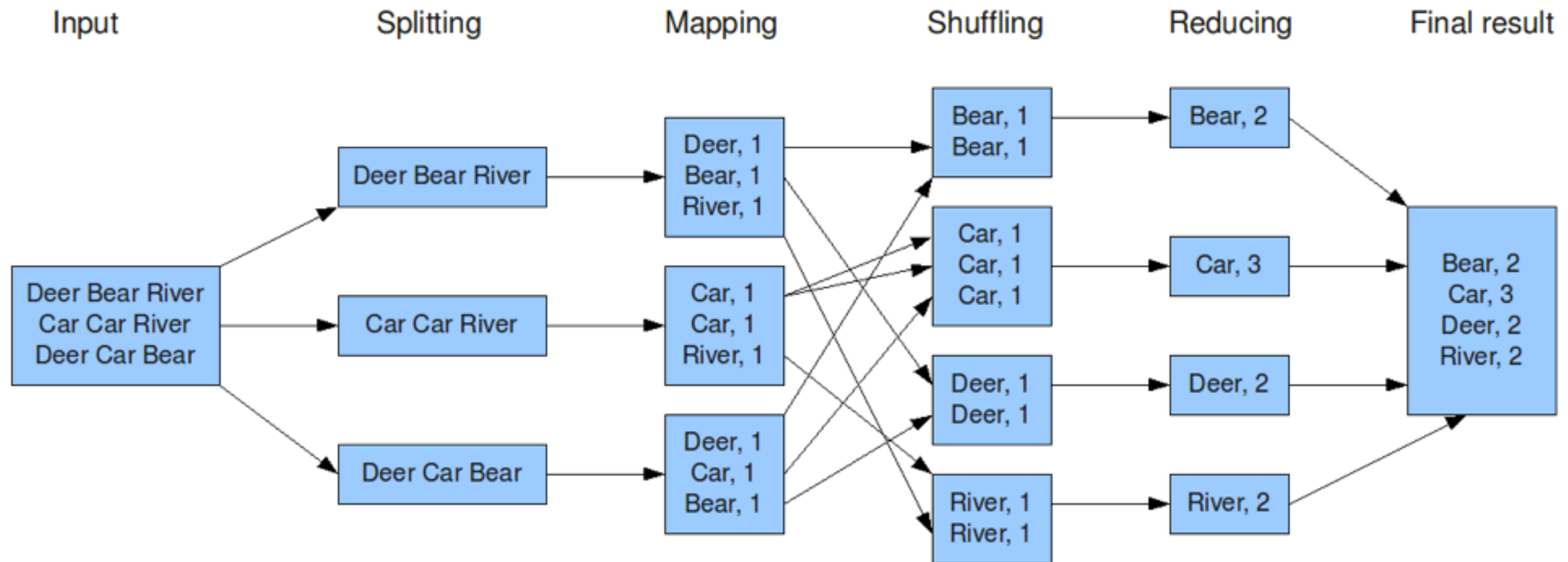
MapReduce

```
1 map(String input_key, String input_value):  
2     // input_key: document name  
3     // input_value: document contents  
4     for each word w in input_value:  
5         EmitIntermediate(w, "1");  
6  
7 reduce(String output_key, Iterator intermediate_values):  
8     // output_key: a word  
9     // output_values: a list of counts  
10    int result = 0;  
11    for each v in intermediate_values:  
12        result += ParseInt(v);  
13    Emit(AsString(result));
```



MapReduce

The overall MapReduce word count process

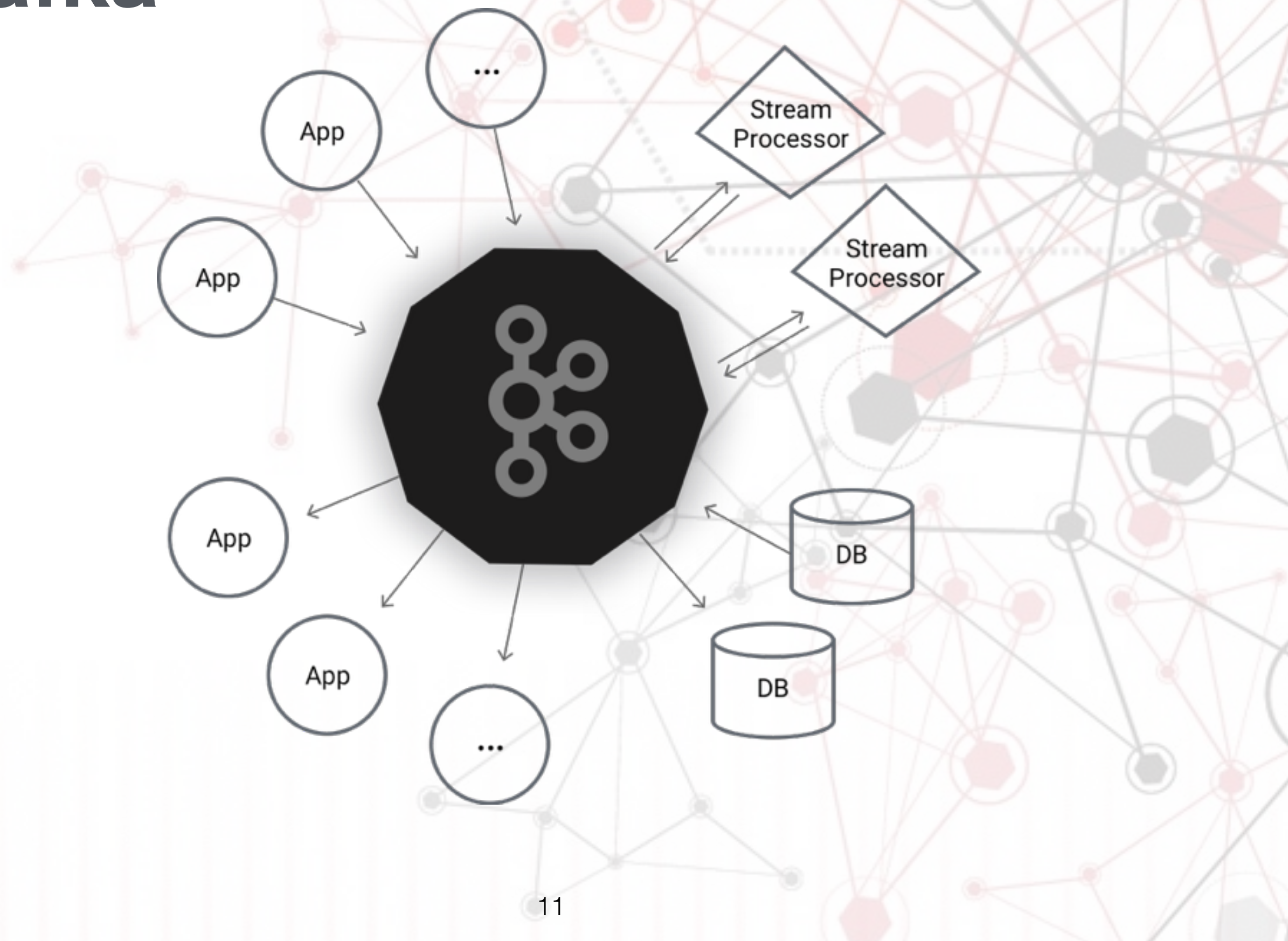


Hadoop

- Nasceu no Yahoo em 2006
- HDFS
- HadoopMapReduce
- HadoopYarn
- Bancos de dados e novos Frameworks surgiram baseados nele.



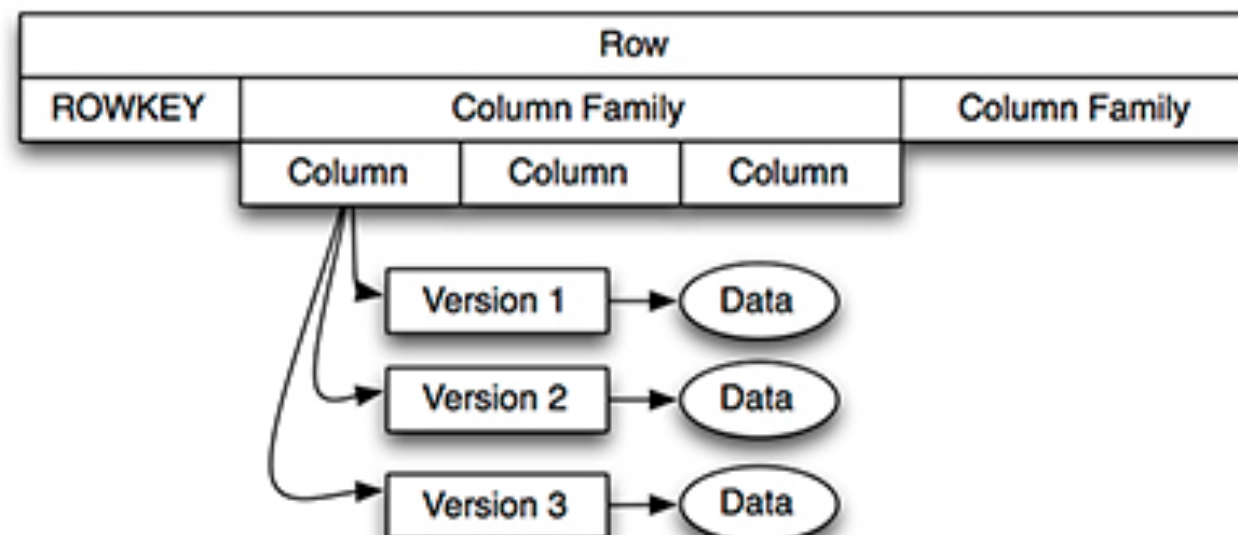
Kafka



HBase



- Foi criado para ser usado com Hadoop
- Banco Distribuído
- Distribuído em Master > RegionServers > Regions
- Orientado por colunas



Airflow

- Criado em 2014 pelo AirBnB
- Conceito de DAG: sequência de tarefas a serem executadas em ordem.
- Agendador de DAG's
- Divide a DAG em tasks



Airflow

DAG: example_qubole_operator

schedule: 1 day, 0:00:00

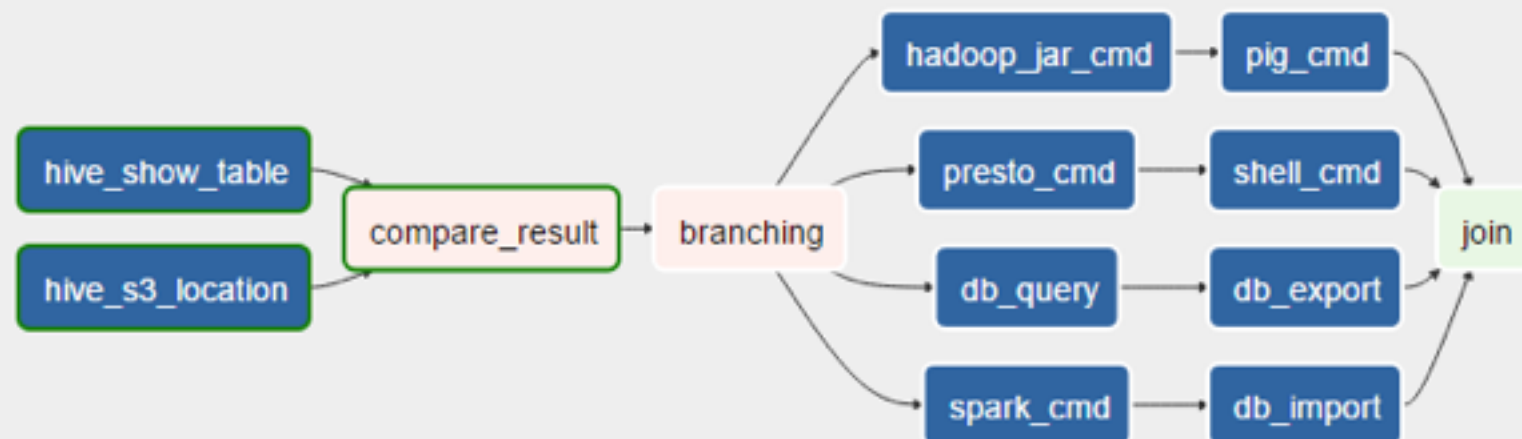
Graph View Tree View Task Duration Landing Times Gantt Details Code

running Run: scheduled__2016-06-20T00:00:00 Layout: Left->Right Go

Search for...

BranchPythonOperator DummyOperator PythonOperator QuboleOperator

success running failed skipped retry queued no status



Será que a Globo.com precisa mesmo?

Números

3 Bilhões de eventos diários

2 Milhões de conexões simultâneas

50 Milhões de usuários únicos por mês

100 Mil novos conteúdos por mês

+20 algoritmos diferentes

100 Mil recomendações por minuto

Pra que tudo isso?



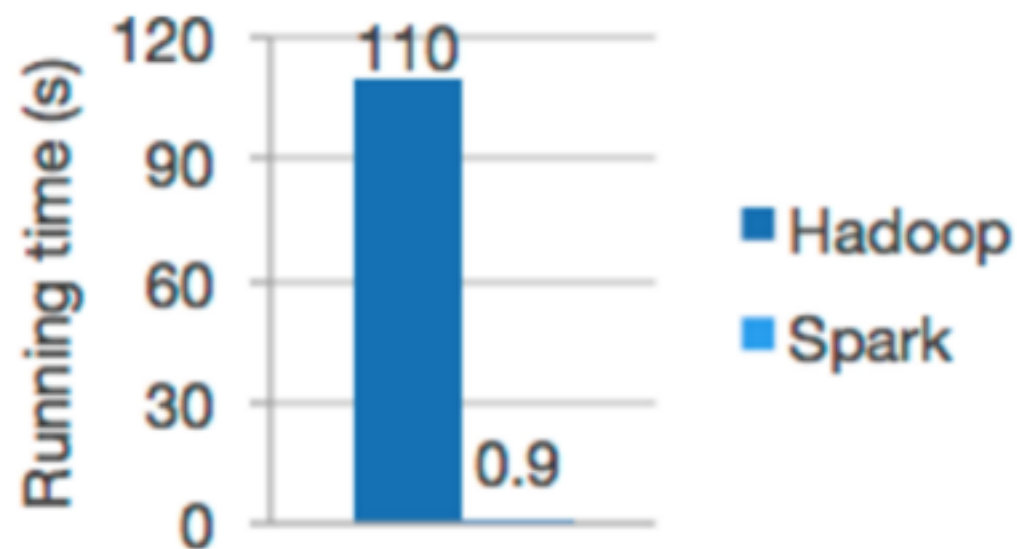
Pra que tudo isso?



- **Predição**
- **Personalização**
- **Saber mais sobre o público**
- **Distribuição mais inteligente de publicidade**

Spark

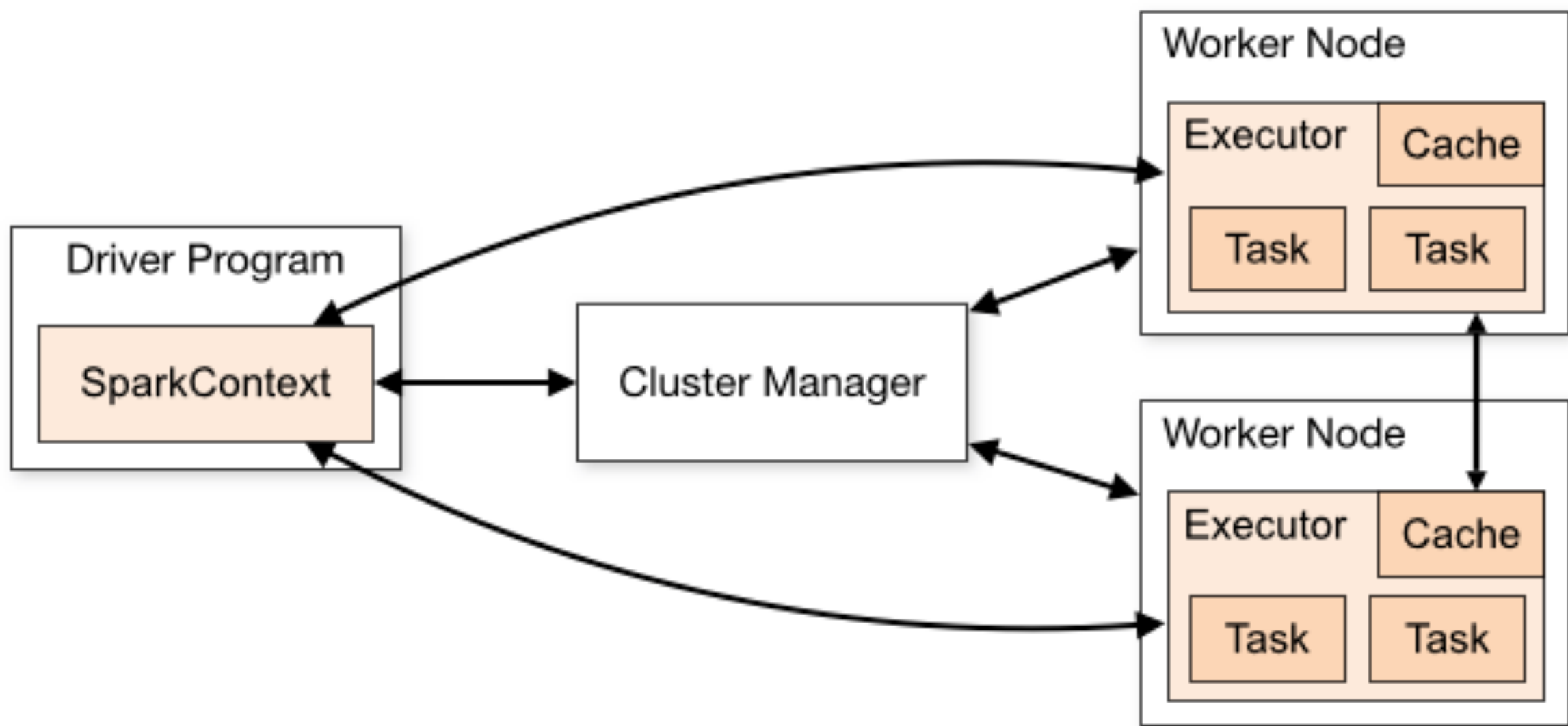
- Nasceu em 2014
- Foco em **velocidade**, facilidade de uso e análises sofisticadas



Logistic regression in Hadoop and Spark

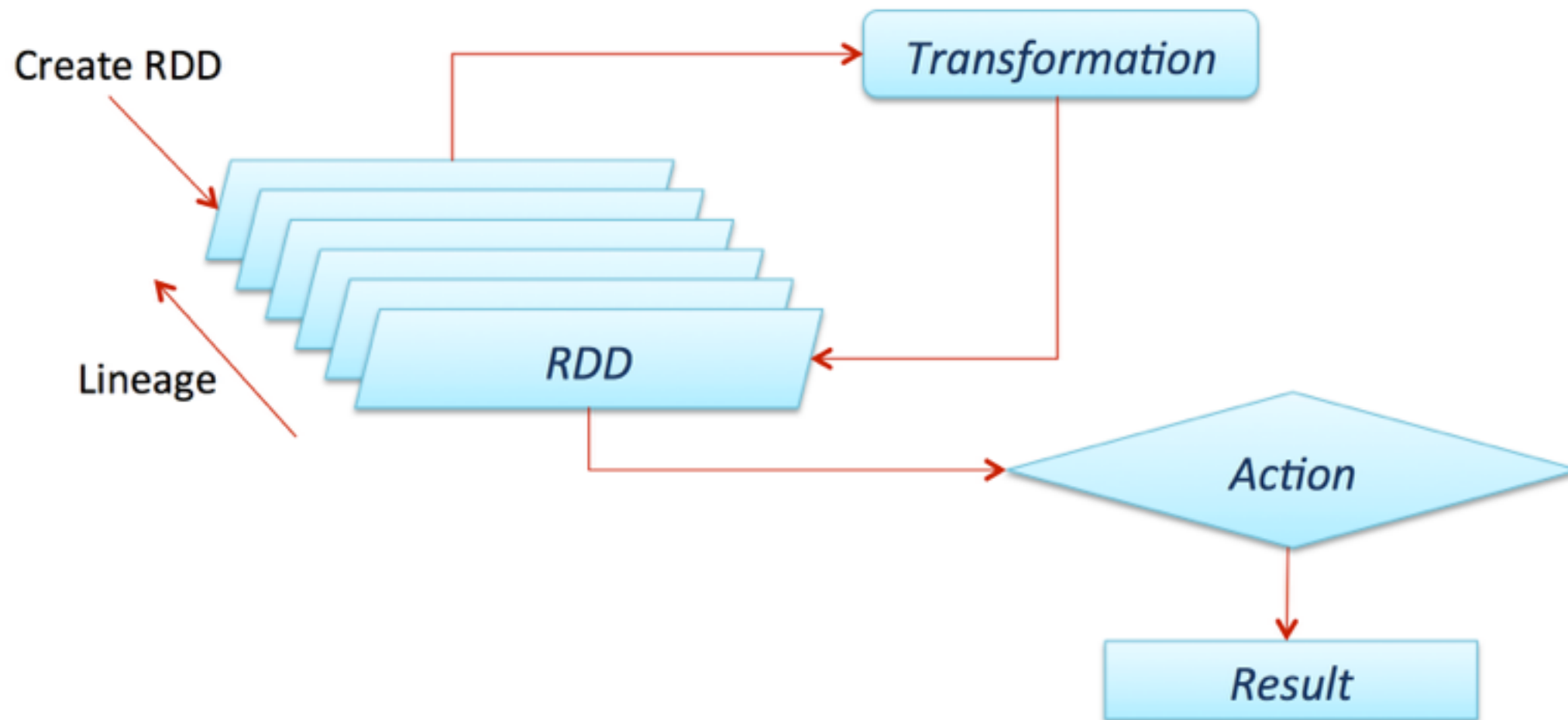
- Aplicações até **100 vezes** mais rápido em memória e até **10 vezes** mais rápido em disco

Spark



Spark

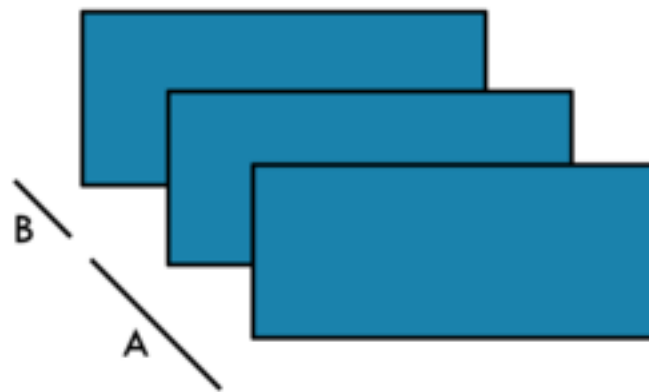
➤ **RDD** (**R**esilient **D**istributed **D**atasets)



Spark

➤ RDD

COLLECT

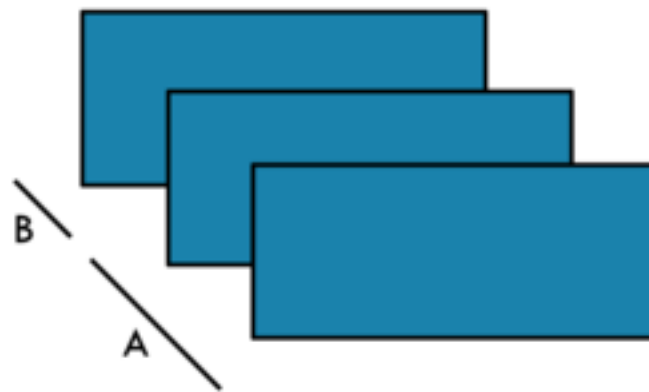


```
11 x = sc.parallelize([1,2,3],2)
12 y = x.collect()
13
14 print x.glom().collect()
15 print y
16
17 sc.stop()
```


Spark

➤ RDD

COLLECT



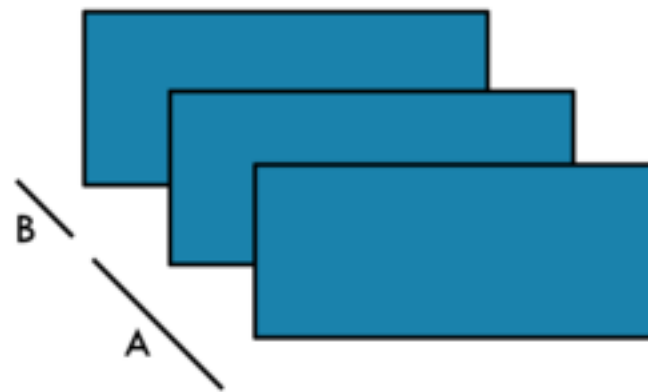
```
11 x = sc.parallelize([1,2,3],2)
12 y = x.collect()
13
14 print x.glom().collect()
15 print y
16
17 sc.stop()
```

x: [[1],[2,3]]
y: [1,2,3]

Spark

➤ RDD

GETNUMPARTITIONS

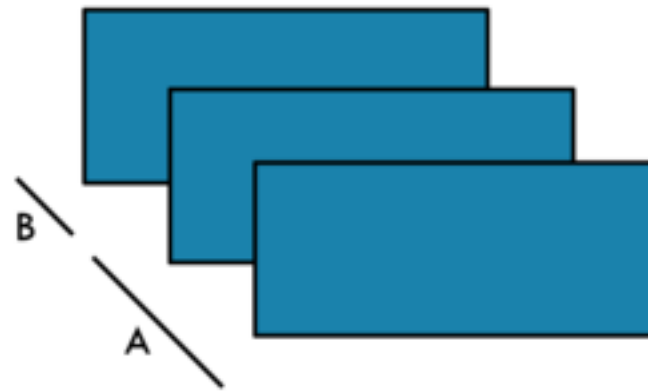


```
11 x = sc.parallelize([1,2,3],2)
12 y = x.getNumPartitions()
13
14 print x.glom().collect()
15 print y
16
```


Spark

➤ RDD

GETNUMPARTITIONS



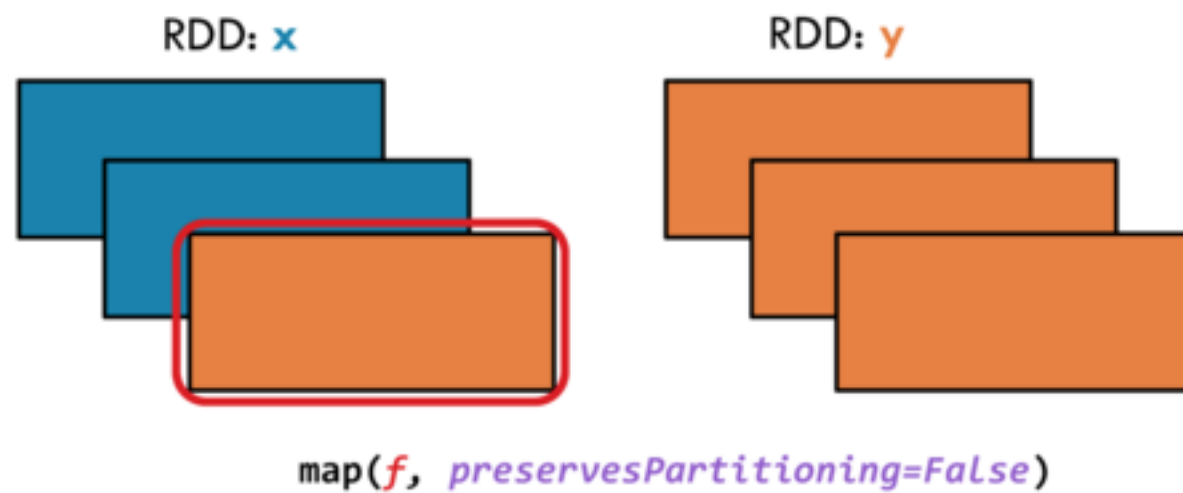
```
11 x = sc.parallelize([1,2,3],2)
12 y = x.getNumPartitions()
13
14 print x.glom().collect()
15 print y
16
```

x: [[1],[2,3]]
y: 2

Spark

➤ RDD

MAP



Spark

➤ RDD

```
x = sc.parallelize(['a','b','c'])  
y = x.map(lambda z:(z,1))  
  
print x.collect()  
print y.collect()
```

Spark

➤ RDD

```
x = sc.parallelize(['a','b','c'])
y = x.map(lambda z:(z,1))

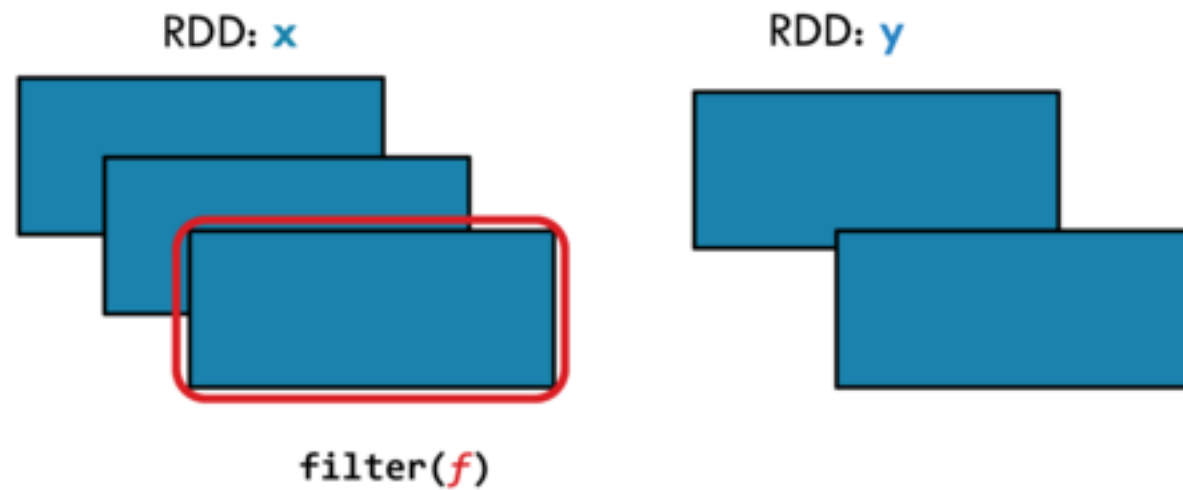
print x.collect()
print y.collect()
```

x: ['a', 'b', 'c']
y: [('a', 1),
 ('b', 1),
 ('c', 1)]

Spark

➤ RDD

FILTER

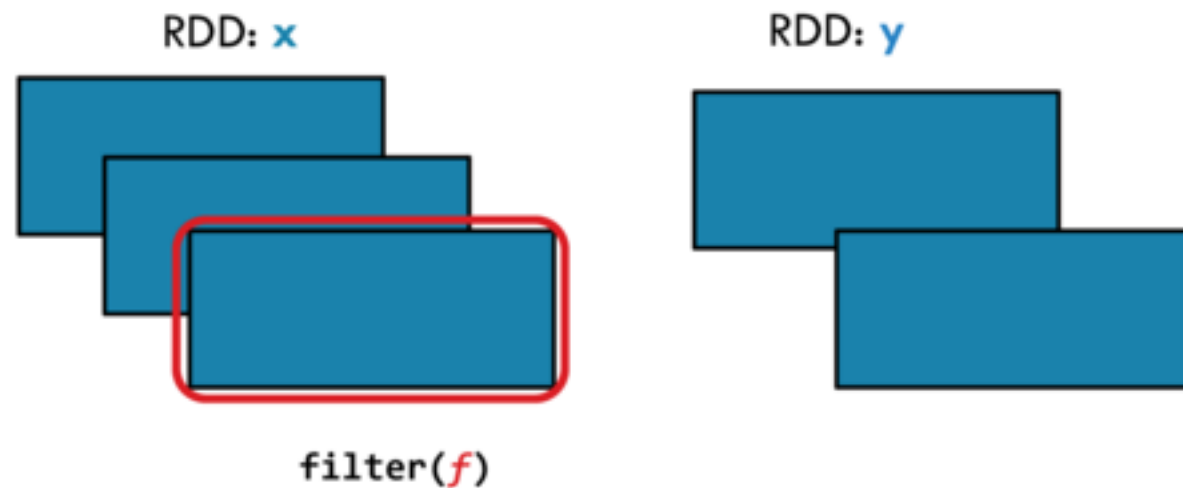


```
11 x = sc.parallelize([1,2,3])
12 y = x.filter(lambda x: x > 1)
13
14 print x.collect()
15 print y.collect()
16
```

Spark

➤ RDD

FILTER



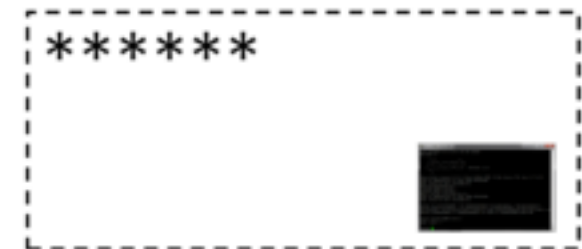
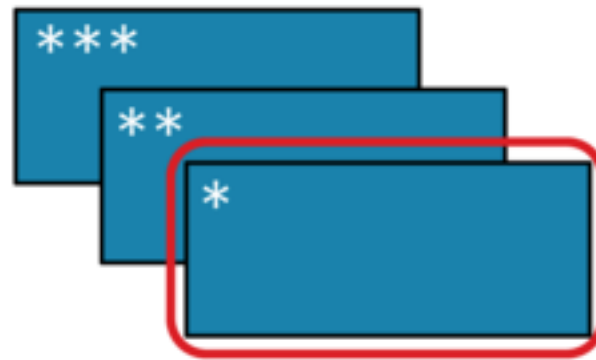
```
11 x = sc.parallelize([1,2,3])
12 y = x.filter(lambda x: x > 1)
13
14 print x.collect()
15 print y.collect()
16
```

x: [1,2,3]
y: [2,3]

Spark

➤ RDD

REDUCE

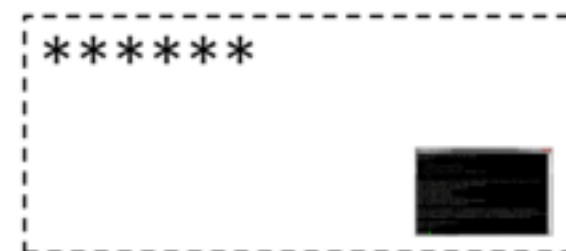
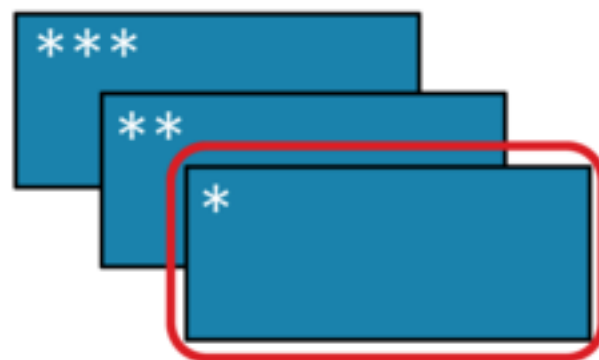


```
11 x = sc.parallelize([1,2,3,4])
12 y = x.reduce(lambda x,y: x + y)
13
14 print x.glom().collect()
15 print y
```

Spark

➤ RDD

REDUCE



```
11 x = sc.parallelize([1,2,3,4])
12 y = x.reduce(lambda x,y: x + y)
13
14 print x.glom().collect()
15 print y
```

```
x: [[1],[2],
    [3],[4]]
y: 10
```


Spark

➤ RDD

REDUCE

```
x = sc.parallelize([1,2,3],2)  
y = x.reduce(lambda x,y: x - y)
```

==

```
x = sc.parallelize([1,2,3],3)  
y = x.reduce(lambda x,y: x - y)
```

?

Spark

➤ RDD

REDUCE

```
x = sc.parallelize([1,2,3],2)  
y = x.reduce(lambda x,y: x - y)
```

==

```
x = sc.parallelize([1,2,3],3)  
y = x.reduce(lambda x,y: x - y)
```

?

y: 2

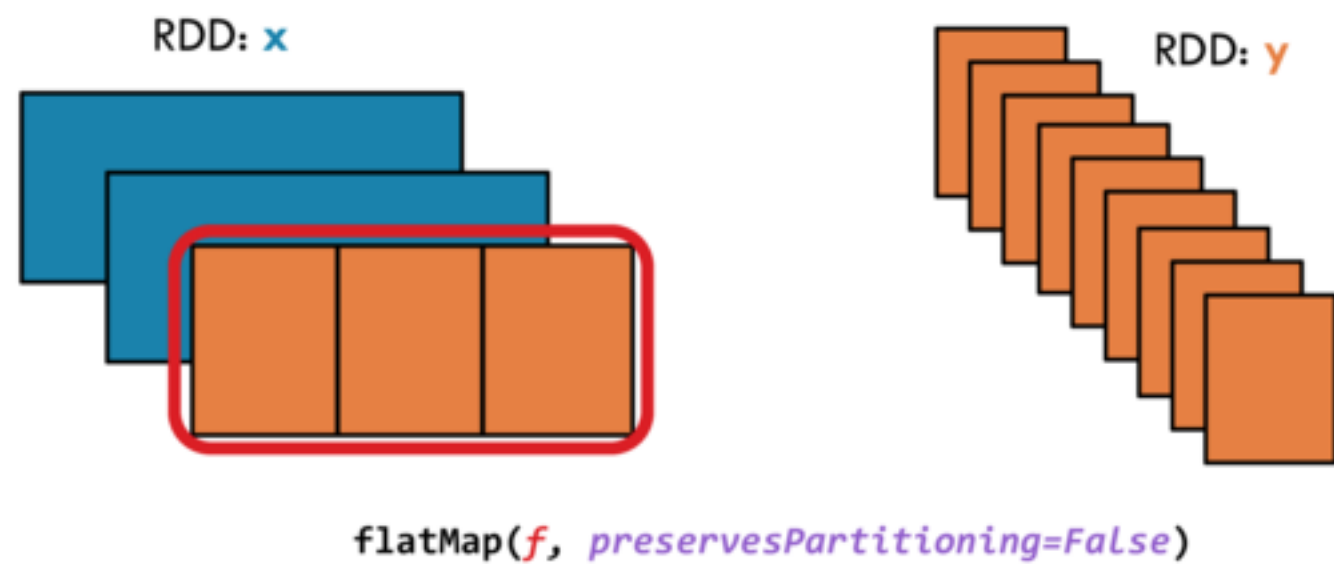
y: -4

Porque?

Spark

➤ RDD

FLATMAP



Spark

➤ RDD

```
11 x = sc.parallelize([1,2,3])
12 y = x.flatMap(lambda z:(z,z**2,42))
13
14 print x.collect()
15 print y.collect()
16
```


Spark

➤ RDD

```
11 x = sc.parallelize([1,2,3])
12 y = x.flatMap(lambda z:(z,z**2,42))
13
14 print x.collect()
15 print y.collect()
16
```

x: [1,2,3]
y: [1,1,42,
2,4,42,
3,9,42]

Spark

➤ RDD

REDUCEBYKEY

```
11 x = sc.parallelize(['a','a','b','b'])
12 y = x.map(lambda z:(z,1)).reduceByKey(lambda x,y:x+y)
13
14 print x.collect()
15 print y.collect()
16
```


Spark

➤ RDD

REDUCEBYKEY

```
11 x = sc.parallelize(['a','a','b','b'])
12 y = x.map(lambda z:(z,1)).reduceByKey(lambda x,y:x+y)
13
14 print x.collect()
15 print y.collect()
16
```

```
x: ['a', 'a', 'b', 'b']
y: [('a', 2),
    ('b', 2)]
```

Spark

➤ RDD

GROUPBYKEY

```
11 x = sc.parallelize(['a','a','b','b']).map(lambda x:(x,1))
12 y = x.groupByKey().map(lambda x:(x[0],list(x[1])))
13
14 print x.collect()
15 print y.collect()
```


Spark

➤ RDD

GROUPBYKEY

```
11 x = sc.parallelize(['a','a','b','b']).map(lambda x:(x,1))
12 y = x.groupByKey().map(lambda x:(x[0],list(x[1])))
13
14 print x.collect()
15 print y.collect()
```

```
x: [('a', 1),
     ('a', 1),
     ('b', 1),
     ('b', 1)]
y: [('a', [1, 1]),
     ('b', [1, 1])]
```

Spark

Exemplo




```
1  import time
2  import pandas as pd
3  words_list = []
4
5  start_time= time.time()
6  with open('really_big_text.txt') as f:
7      for i in f.readlines():
8          words = i.split()
9          words_list = words_list + map(lambda x:(x,1),words)
10
11  df = pd.DataFrame(words_list,columns=['word', 'count'])\
12      .groupby('word').sum()\
13      .sort('count',ascending=False)
14
15
16  print "TIME: %s"%(time.time()-start_time)
17  print df
18
```

× ..ython_sudeste (zsh)

```
(.env) python_sudeste % python only_python.py  
only_python.py:11: FutureWarning: sort(columns=...  
    df = pd.DataFrame(words_list, columns=['word', 'co  
TIME: 1506.24237919
```

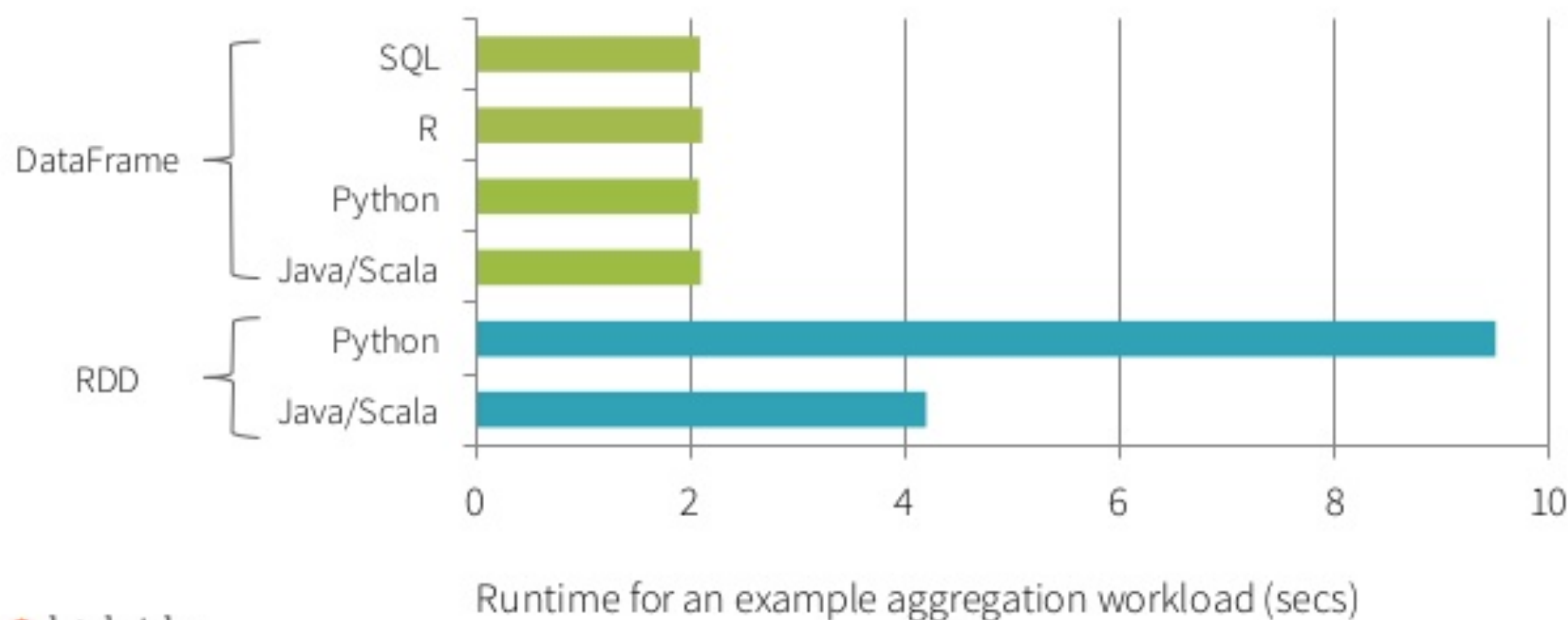
	count
word	
the	72271
and	46446
of	40346
to	16539
that	15553
in	14504
he	11377
shall	10872
	10056


```
1  import time
2  from pyspark.sql import SQLContext
3  from pyspark import SparkConf, SparkContext
4
5  conf = (SparkConf()
6          .setAppName("MyFirstApp")
7          .set("spark.executor.memory", "12g"))
8
9  sc = SparkContext(conf = conf)
10 sqlContext = SQLContext(sc)
11
12 start_time = time.time()
13
14 text = sqlContext.read.text('really_big_text.txt').rdd
15 text_words = text.flatMap(lambda x: x['value'].split())
16 text_words.cache()
17
18 count_words = text_words\
19     .map(lambda x: (x,1))\
20     .reduceByKey(lambda x,y:x+y)\
21     .sortBy(lambda x:x[1],False)\
22     .toDF(schema=['word', 'count'])
23
24 print "TIME: %s"%(time.time()-start_time)
25 print "Top 100"
26 count_words.show(10)
27 print "Total : %s"%len(text_words.collect())
28 sc.stop()
29
```

```
(.env) python_sudeste % ./submit_job.sh
TIME: 10.1741509438
Top 100
+-----+-----+
| word|count|
+-----+-----+
| the|72271|
| and|46446|
| of|40346|
| to|16539|
| that|15553|
| in|14504|
| he|11377|
| shall|10872|
| unto|10056|
| his|10017|
+-----+-----+
only showing top 10 rows

Total : 948432
```


Benefit of Logical Plan: Performance Parity Across Languages



OBRIGADA!



susana bouchardet



@s_bouchardet