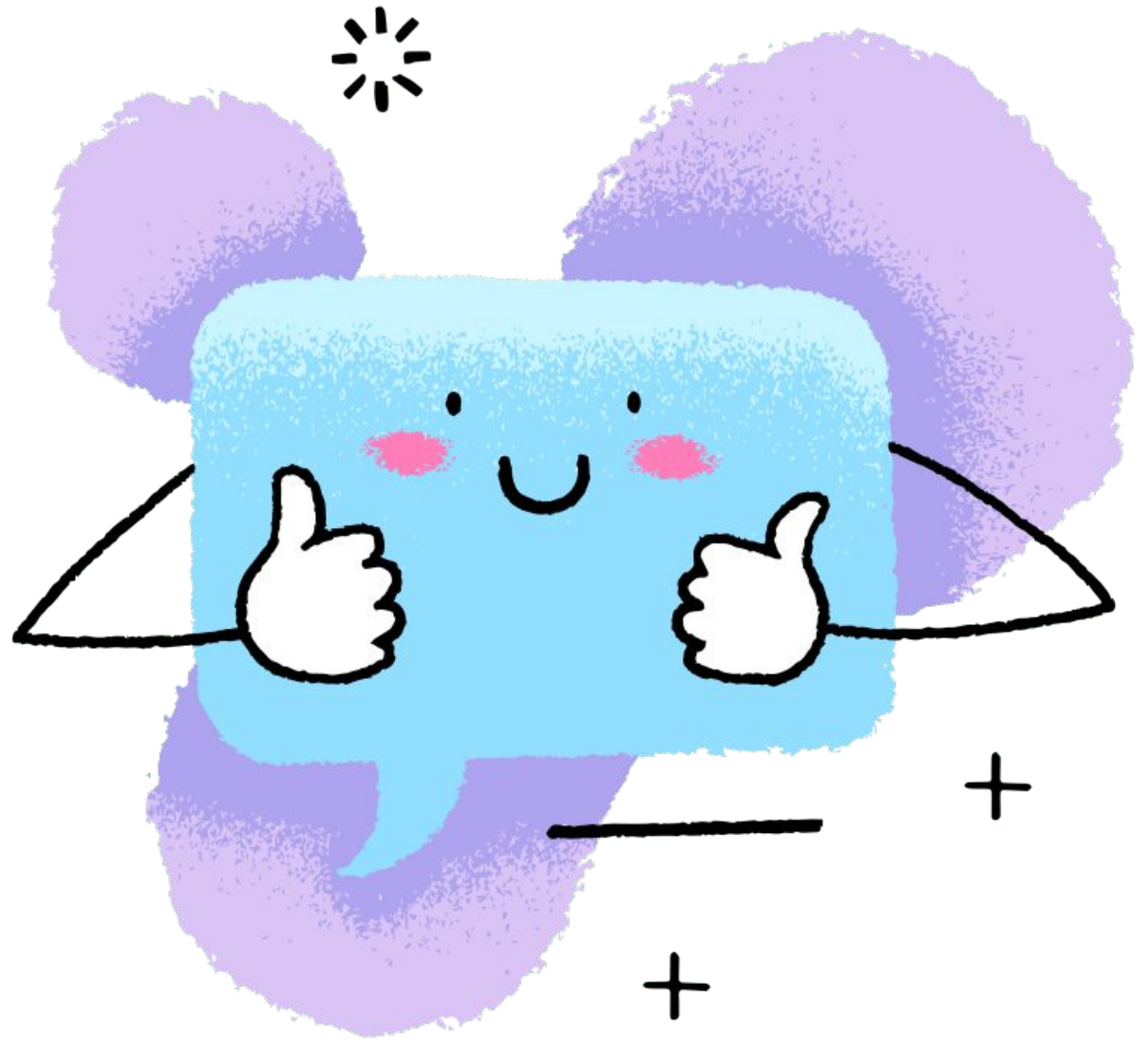


# PySpark





# PySpark DataFrame API

Apache Spark

# Что будет на уроке

1. RDD: Resilient Data Distribution. Класс DataFrame.
2. План запроса. Перемешивание данных (шаффл). Сохранение данных.
3. groupBy, join
4. Получение плана выполнения запроса.
5. Способы оптимизации: push down фильтрация, стратегия исполнения join.
6. Работа со Spark UI
7. Сложные случаи: перекос в объеме данных между экзекьюторами. Функции repartition, coalesce
8. Агрегирующие функции.
9. Агрегирование по синтетическому ключу.
10. Pivot

# RDD

## низкоуровневый API для работы с распределенными данными

```
1 # In Python
2 # Create an RDD of tuples (name, age)
3 dataRDD = sc.parallelize([("Brooke", 20), ("Denny", 31), ("Jules", 30),
4 ("TD", 35), ("Brooke", 25)])
5 # Use map and reduceByKey transformations with their lambda
6 # expressions to aggregate and then compute average
7 agesRDD = (dataRDD
8 .map(lambda x: (x[0], (x[1], 1)))
9 .reduceByKey(lambda x, y: (x[0] + y[0], x[1] + y[1]))
10 .map(lambda x: (x[0], x[1][0]/x[1][1])))
```

# DataFrame

## Начиная со Spark 1.6 обобщенный API для работы с RDD

```
1 from pyspark.sql import SparkSession
2 from pyspark.sql.functions import avg
3 # Create a DataFrame using SparkSession
4 spark = (SparkSession
5 .builder
6 .appName("AuthorsAges")
7 .getOrCreate())
8 # Create a DataFrame
9 data_df = spark.createDataFrame([("Brooke", 20), ("Denny", 31), ("Jules", 30), ("TD",
10 35), ("Brooke", 25)], ["name", "age"])
11 # Group the same names together, aggregate their ages, and compute an average
12 avg_df = data_df.groupBy("name").agg(avg("age"))
13 # Show the results of the final execution
14 avg_df.show()
```

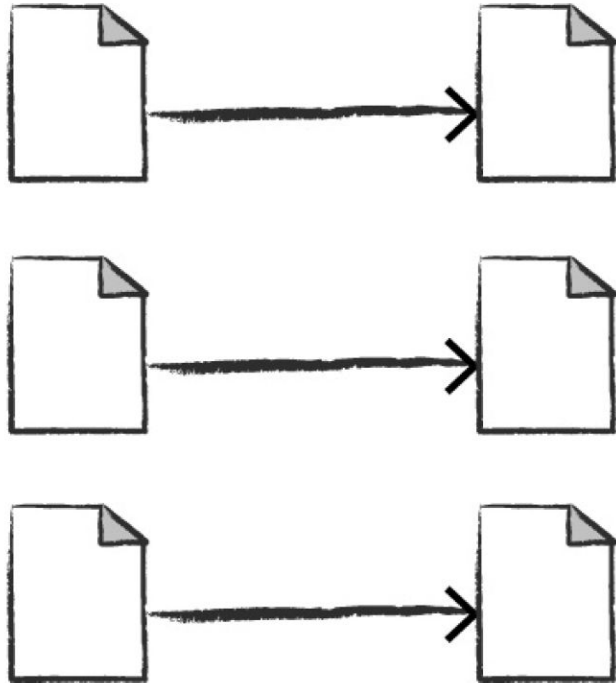
name	avg(age)
Brooke	22.5
Jules	30.0
TD	35.0
Denny	31.0

# Transformation types

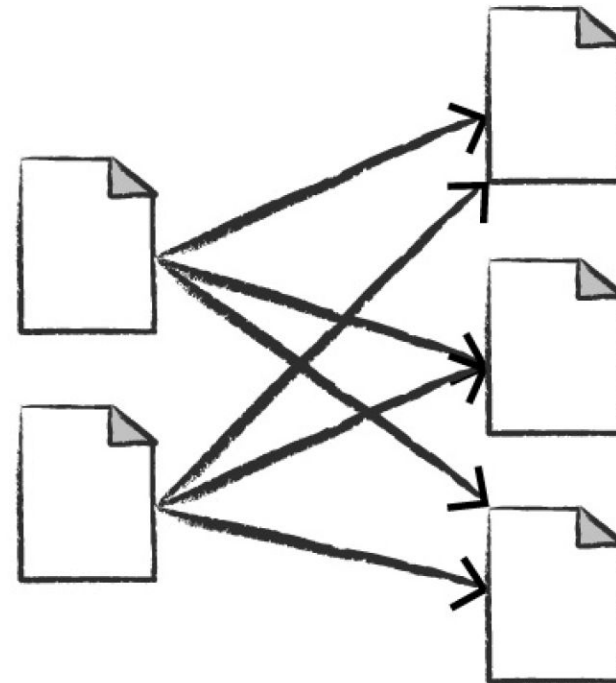
narrow -> in-memory filters

wide -> shuffle, write result on disk

Narrow transformations  
1 to 1



Wide transformations  
(shuffles) 1 to N





# Transformations and actions as Spark operations

Трансформации копируются до вызова action

Transformations	Actions
orderBy()	show()
groupBy()	take()
filter()	count()
select()	collect()
join()	save()

# Spark DataFrame

таблица с типизированными колонками

<b>Id (Int)</b>	<b>First (String)</b>	<b>Last (String)</b>	<b>Url (String)</b>	<b>Published (Date)</b>	<b>Hits (Int)</b>	<b>Campaigns (List[Strings])</b>
1	Jules	Damji	https:// tinyurl.1	1/4/2016	4535	[twitter, LinkedIn]
2	Brooke	Wenig	https:// tinyurl.2	5/5/2018	8908	[twitter, LinkedIn]
3	Denny	Lee	https:// tinyurl.3	6/7/2019	7659	[web, twitter, FB, LinkedIn]
4	Tathagata	Das	https:// tinyurl.4	5/12/2018	10568	[twitter, FB]



# Основные типы данных Python

## значения в колонках

Data type	Value assigned in Python	API to instantiate
ByteType	int	DataTypes.ByteType
ShortType	int	DataTypes.ShortType
IntegerType	int	DataTypes.IntegerType
LongType	int	DataTypes.LongType
FloatType	float	DataTypes.FloatType
DoubleType	float	DataTypes.DoubleType
StringType	str	DataTypes.StringType
BooleanType	bool	DataTypes.BooleanType
DecimalType	decimal.Decimal	DecimalType

# Python structured data types in Spark

Data type	Value assigned in Python	API to instantiate
BinaryType	bytearray	BinaryType()
TimestampType	datetime.datetime	TimestampType()
DateType	datetime.date	DateType()
ArrayType	List, tuple, or array	ArrayType(dataType, [nullable])
MapType	dict	MapType(keyType, valueType, [nullable])
StructType	List or tuple	StructType([fields])
StructField	A value type corresponding to the type of this field	StructField(name, dataType, [nullable])

# Задание schema

## и зачем это нужно

- Избегаем приведение типов
- Не нужно считывать файл только для определения схемы данных (существенно для больших файлов)
- Можно отловить несоответствие данных

```
1 from pyspark.sql.types import *
2 schema = StructType([StructField("author", StringType(), False),
3 StructField("title", StringType(), False),
4 StructField("pages", IntegerType(), False)])
5
6 #the same schema using DDL
7 schema = "author STRING, title STRING, pages INT"
```

# Создаем статический DataFrame

```
1 from pyspark.sql import SparkSession
2 # Define schema for our data using DDL
3 schema = "`Id` INT, `First` STRING, `Last` STRING, `Url` STRING,
4 `Published` STRING, `Hits` INT, `Campaigns` ARRAY<STRING>"
5 # Create our static data
6 data = [
7 [1, "Jules", "Damji", "https://tinyurl.1", "1/4/2016", 4535, ["twitter", "LinkedIn"]],
8 [2, "Brooke", "Wenig", "https://tinyurl.2", "5/5/2018", 8908, ["twitter", "LinkedIn"]]
9 ]
10
11 # Create a DataFrame using the schema defined above
12 blogs_df = spark.createDataFrame(data, schema)
13 # Show the DataFrame; it should reflect our table above
14 blogs_df.show()
15 # Print the schema used by Spark to process the DataFrame
16 print(blogs_df.printSchema())
```



# Count, countDistinct

агрегирующие функции могут  
вызываться по всему набору данных

```
1 from pyspark.sql.functions import count
2 df.select(count("StockCode")).show()
3
4 from pyspark.sql.functions import countDistinct
5 df.select(countDistinct("StockCode")).show()
6
7 from pyspark.sql.functions import approx_count_distinct
8 df.select(approx_count_distinct("StockCode", 0.1)).show() |
```

# Column alias

## МЕНЯЕМ НАЗВАНИЕ КОЛОНКИ

```
1 from pyspark.sql.functions import sum, count, avg, expr
2 df.select(
3     count("Quantity").alias("total_transactions"),
4     sum("Quantity").alias("total_purchases"),
5     avg("Quantity").alias("avg_purchases"),
6     expr("mean(Quantity)").alias("mean_purchases"))\
7     .selectExpr(
8         "total_purchases/total_transactions",
9         "avg_purchases",
10        "mean_purchases")\
11 .show()
```



# Считываем данные для дальнейших упражнений

```
1 # in Python
2 df = spark.read.format("csv")\
3 .option("header", "true")\
4 .option("inferSchema", "true")\
5 .load("/data/retail-data/by-day/2010-12-01.csv")
6 df.printSchema()
7 df.createOrReplaceTempView("dfTable")
```

# Booleans

## Фильтрация по значению в колонке

```
1 # in Python
2 from pyspark.sql.functions import col
3 df.where(col("InvoiceNo") != 536365)\
4 .select("InvoiceNo", "Description")\
5 .show(5, False)|
```

```
1 from pyspark.sql.functions import instr
2
3 priceFilter = col("UnitPrice") > 600
4 descripFilter = instr(df.Description, "POSTAGE") >= 1
5
6 df.where(df.StockCode.isin("DOT"))\
7 .where(priceFilter | descripFilter).show()
```

# SQL для создания колонок и логических выражений через expr

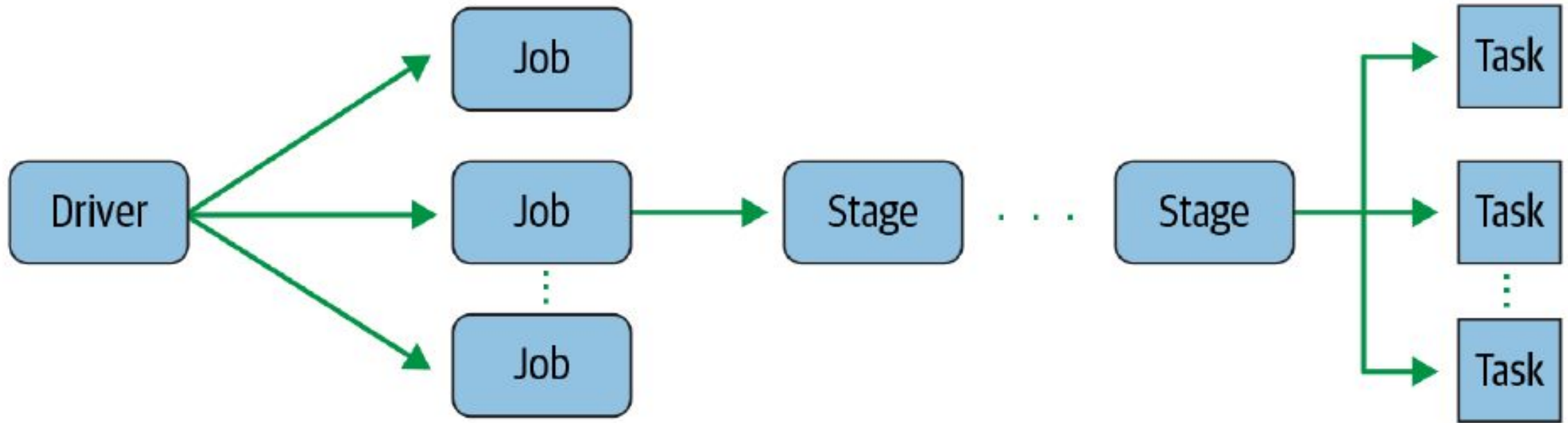
```
1 from pyspark.sql.functions import expr
2 df.withColumn("isExpensive", expr("NOT UnitPrice <= 250"))\
3 .where("isExpensive")\
4 .select("Description", "UnitPrice").show(5)
```

# Численные значения

```
1 from pyspark.sql.functions import expr, pow
2 fabricatedQuantity = pow(col("Quantity") * col("UnitPrice"), 2) + 5
3
4 df.select(expr("CustomerId"), \
5           fabricatedQuantity.alias("realQuantity"))\
6 .show(2)
7
8 # либо через SQL
9 df.selectExpr("CustomerId",
10 "(POWER((Quantity * UnitPrice), 2.0) + 5) as realQuantity")\
11 .show(2)
```

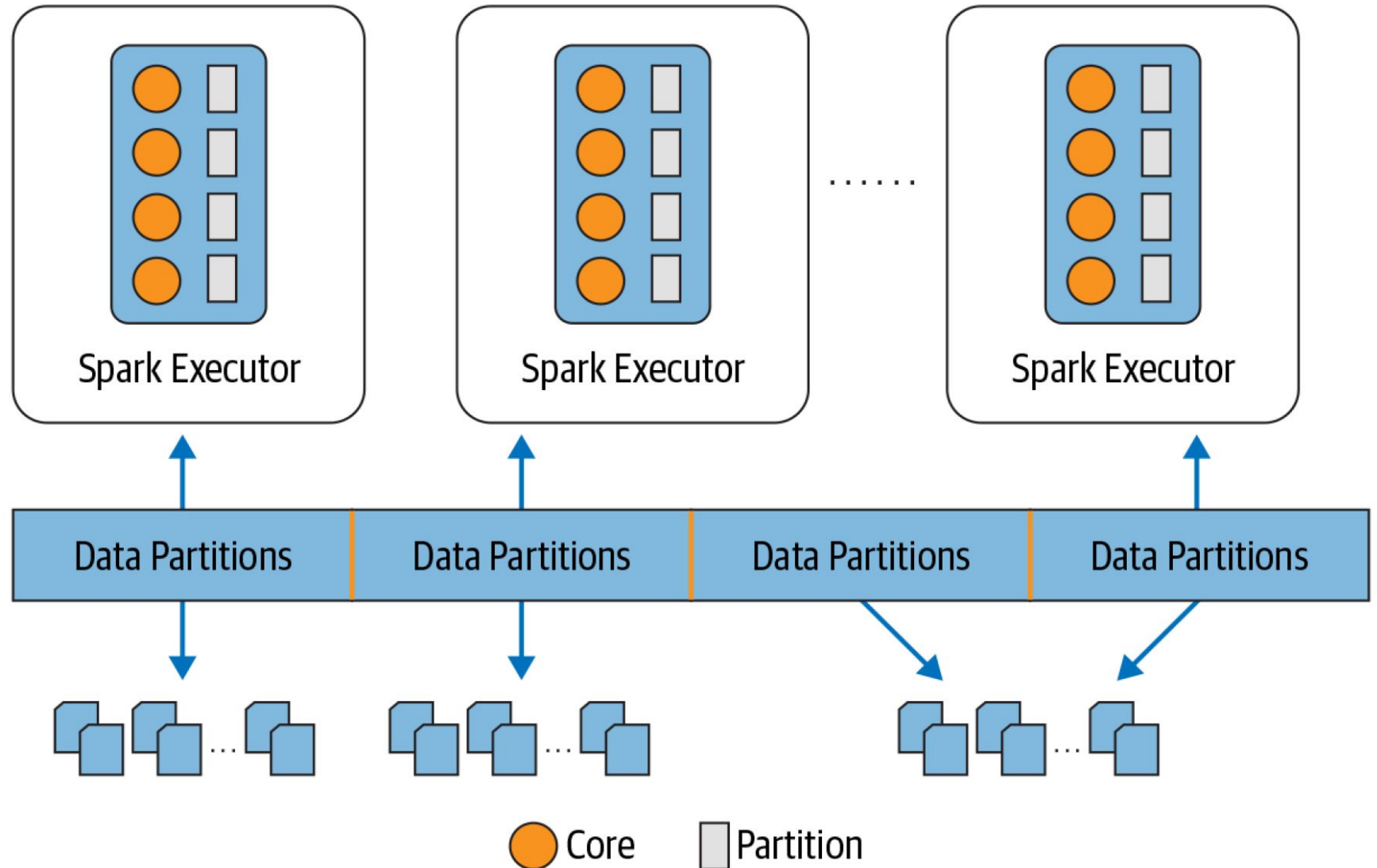
# Spark jobs, stage, task

shuffle occurs between every two stages



# Spark executors, partitions

## Relationship of Spark tasks, cores, partitions





# **df.cache()**

## **store as many of the partitions in memory across Spark executors as memory allows**

When to Cache and Persist:

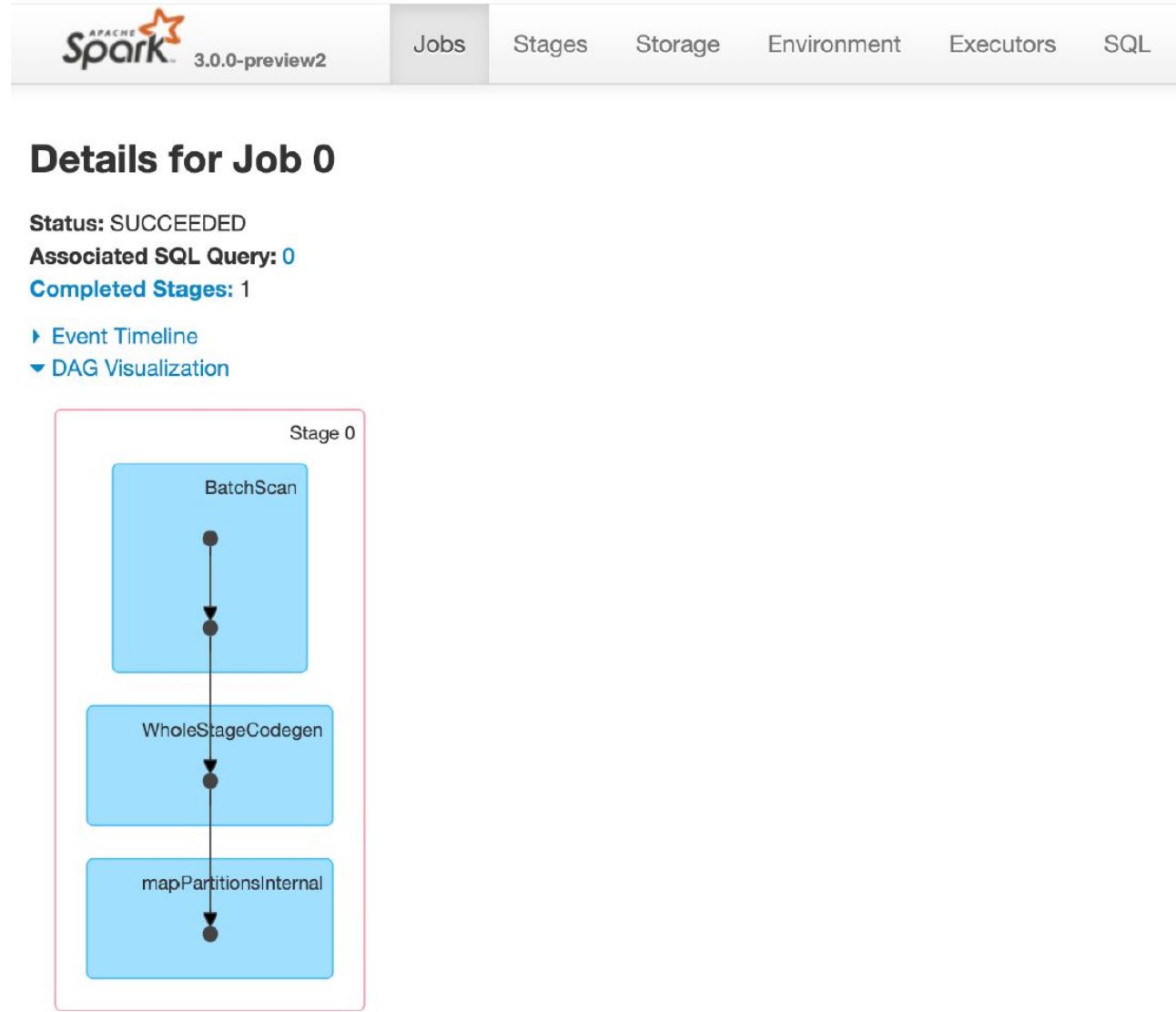
- DataFrames commonly used during iterative machine learning training
- DataFrames accessed commonly for doing frequent transformations during ETL or building data pipelines

When Not to Cache and Persist:

- DataFrames that are too big to fit in memory
- An inexpensive transformation on a DataFrame not requiring frequent use, regardless of size

# Spark UI

## DAG, resources

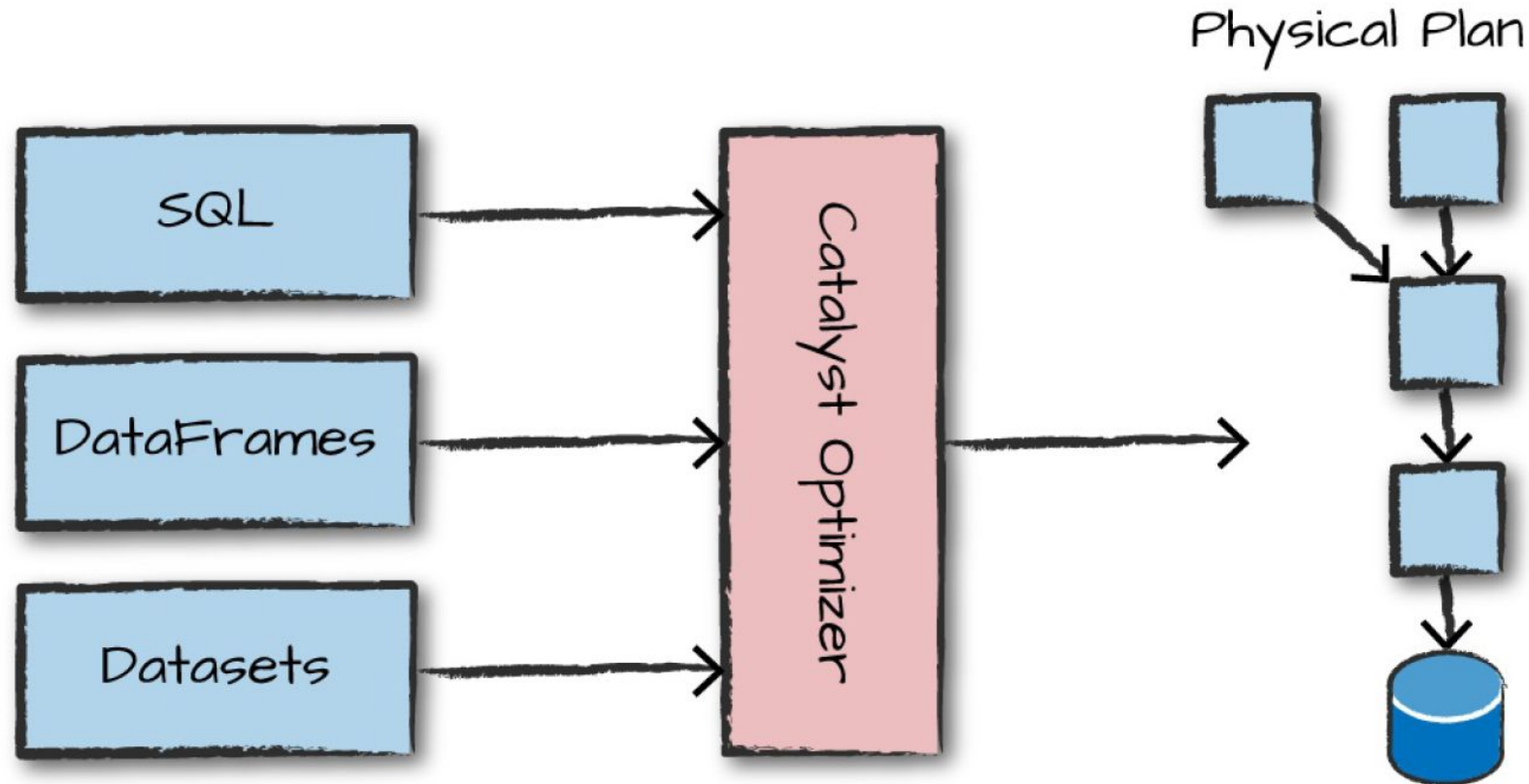


# Практика

## Zeppelin -> Lecture 2

# Оптимизатор запросов

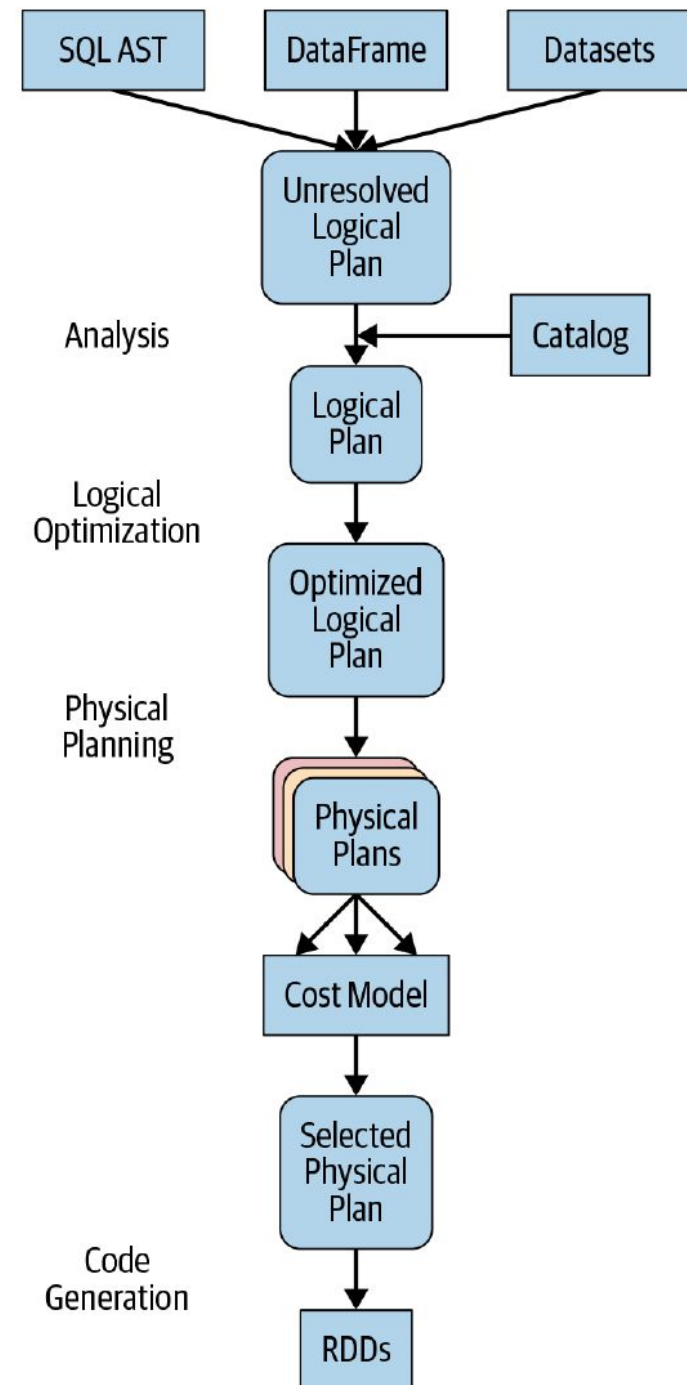
Lazy исполнение трансформаций  
позволяет оптимизировать план



# Catalyst Optimizer

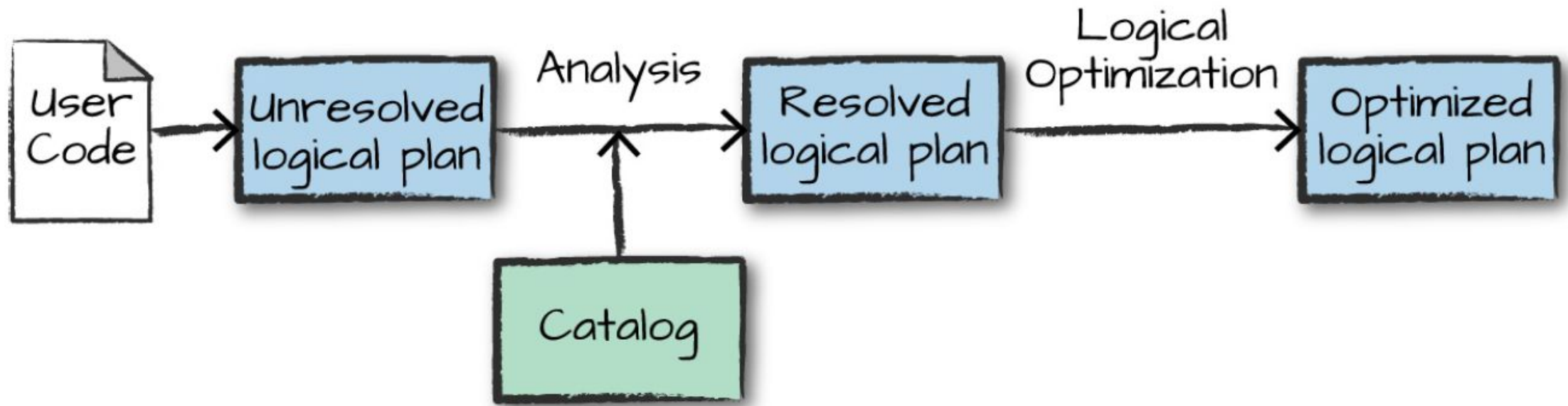
## transformational phases

1. Analysis
2. Logical optimization
3. Physical planning
4. Code generation



# Logical plan

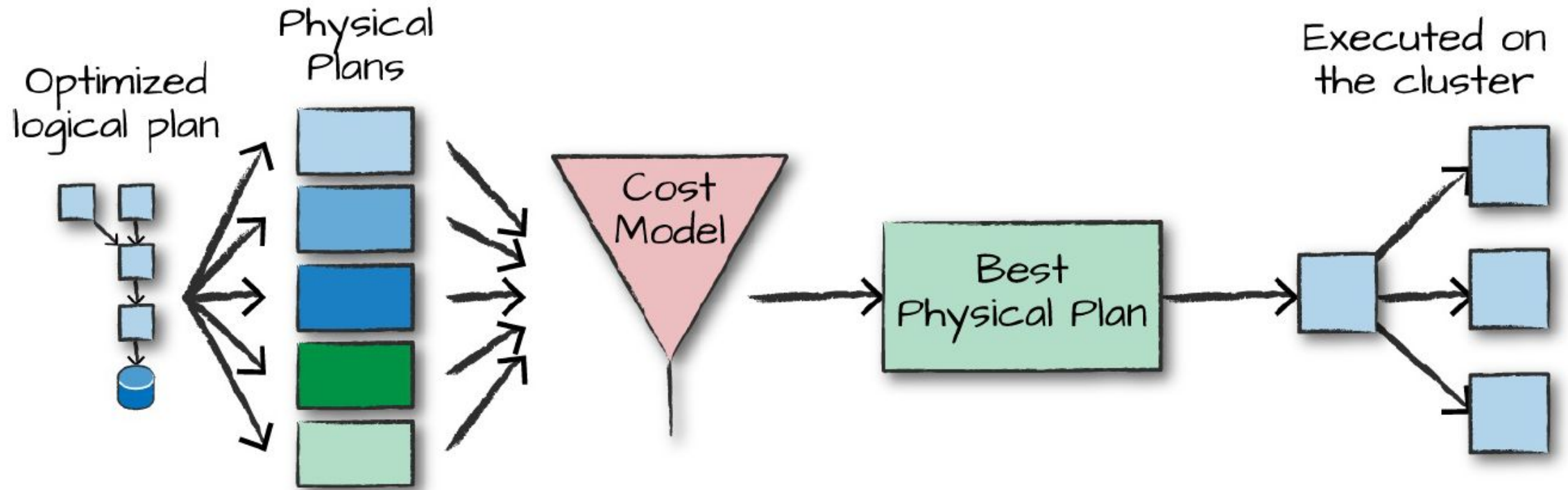
- 1) проверка наличия таблиц, колонок
- 2) pushing down predicates or selections



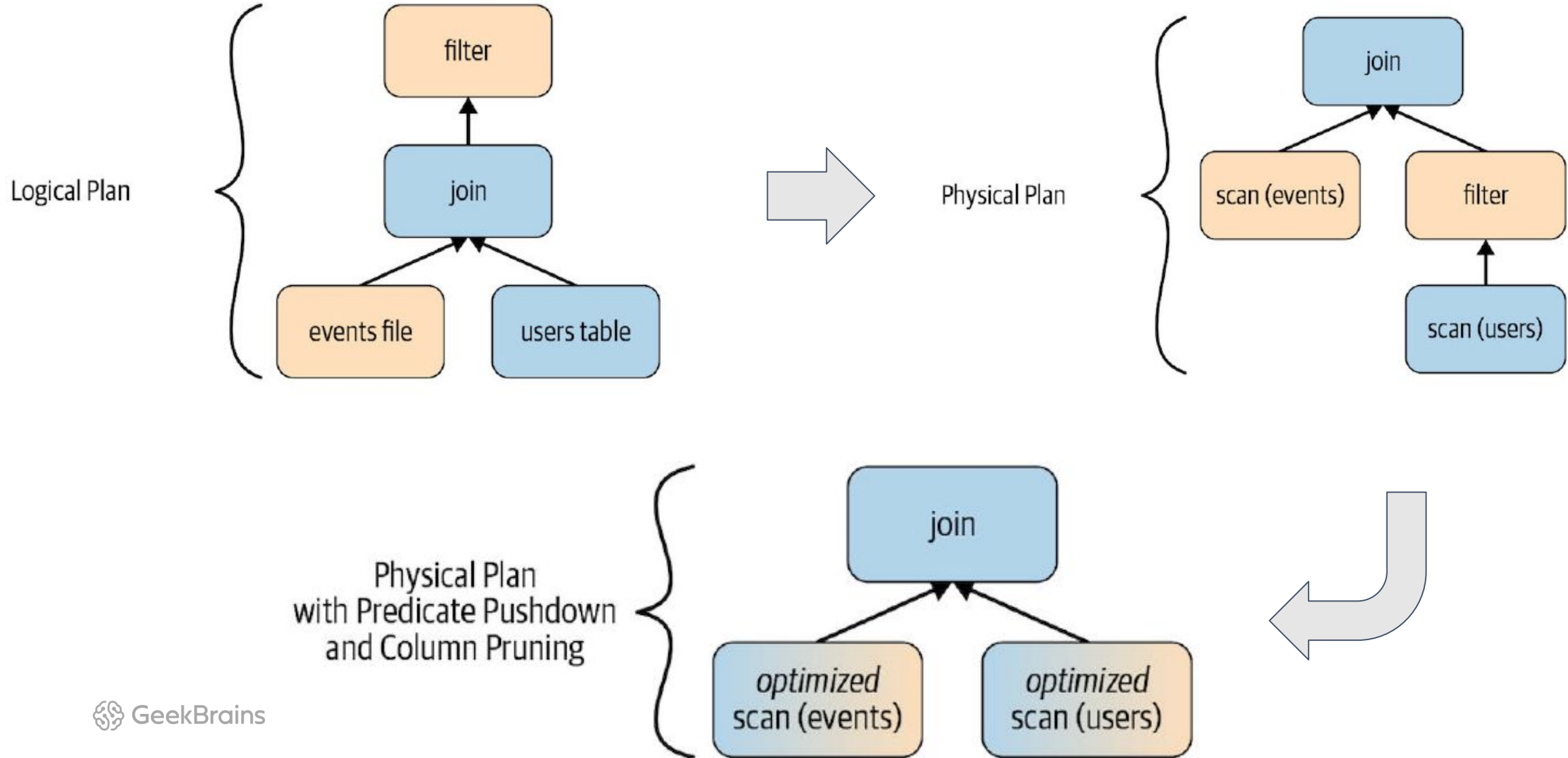


# Physical Planning

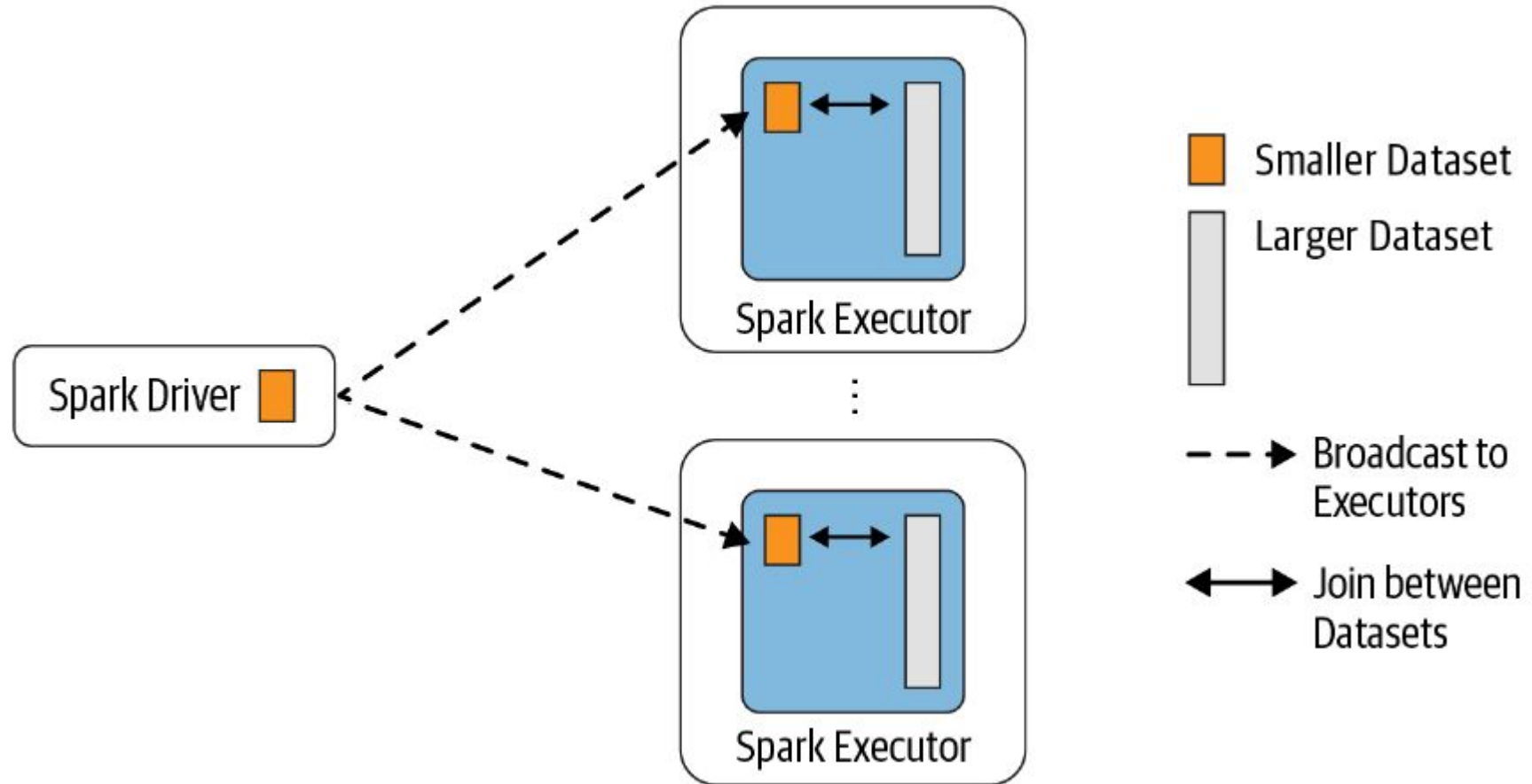
Использует знание о размере и  
распределении партиций



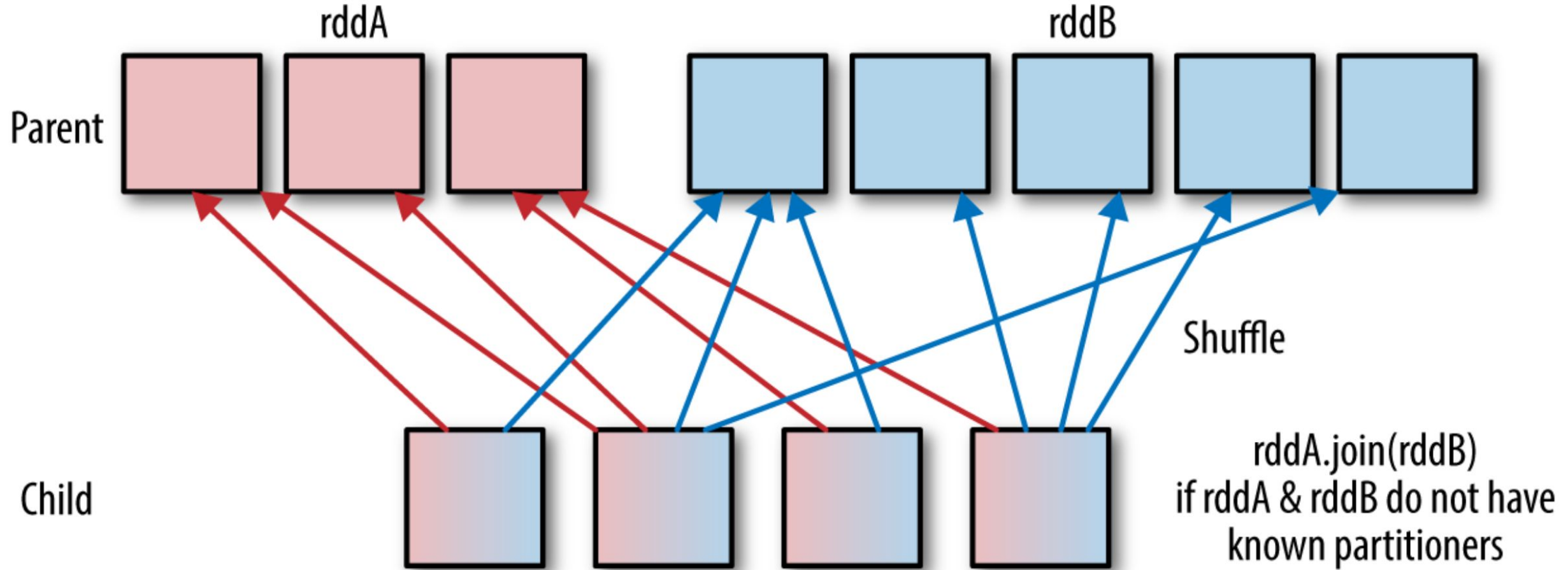
# Планы запроса из примера



# Broadcast Hash Join



# Shuffle Sort Merge Join



# Ссылки

[catalyst-optimizer](#)

[deep-dive-into-spark-sqls-catalyst-optimizer](#)

[high-performance-spark/ch04.html](#)

[towardsdatascience.com/the-art-of-joining-in-spark](#)

[https://itnext.io/handling-data-skew-in-apache-spark-9f56343e58e8](#)

[https://blog.clairvoyantsoft.com/optimize-the-skew-in-spark-e523c6ee18ac](#)

**Спасибо!**  
**Каждый день**  
**вы становитесь**  
**лучше :)**

