

Apache Spark

Дмитрий Лахвич(frostball@gmail.com)

Руководитель группы разработки Tele 2

Apache Spark

Apache Spark

Spark is a unified analytics engine for large-scale data processing. It provides high-level APIs in Scala, Java, Python, and R, and an optimized engine that supports general computation graphs for data analysis. It also supports a rich set of higher-level tools including Spark SQL for SQL and DataFrames, MLlib for machine learning, GraphX for graph processing, and Structured Streaming for stream processing.

<https://spark.apache.org/>

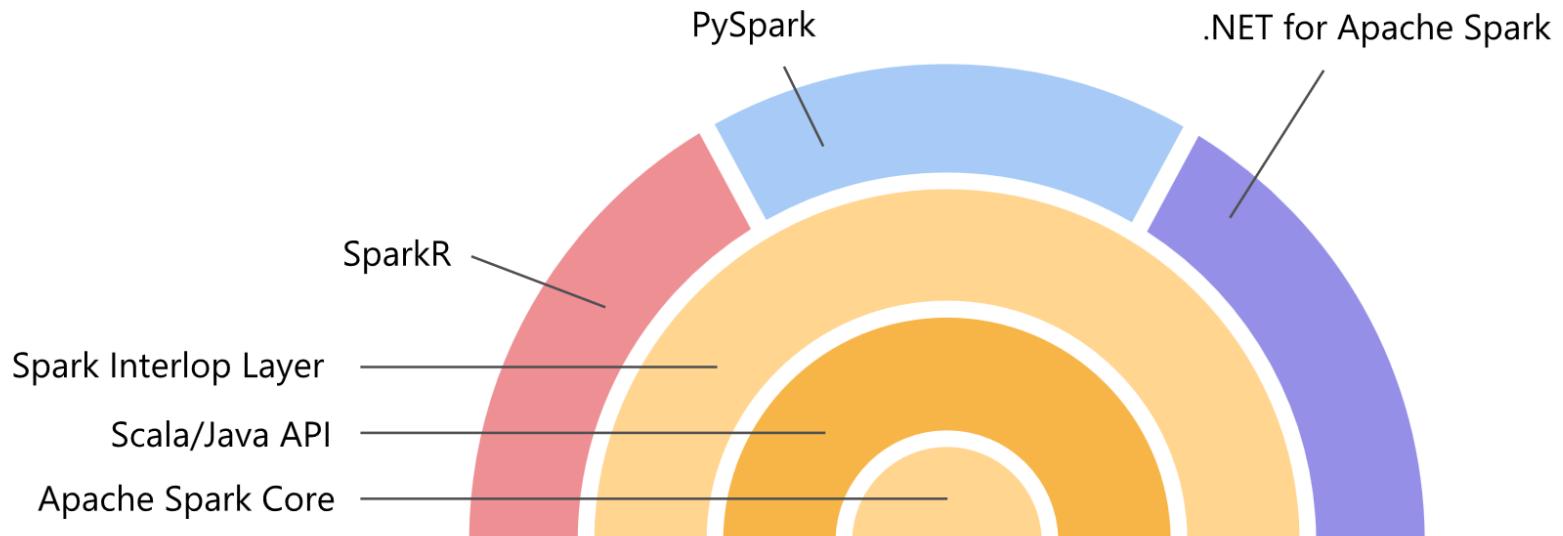
<https://github.com/apache/spark>

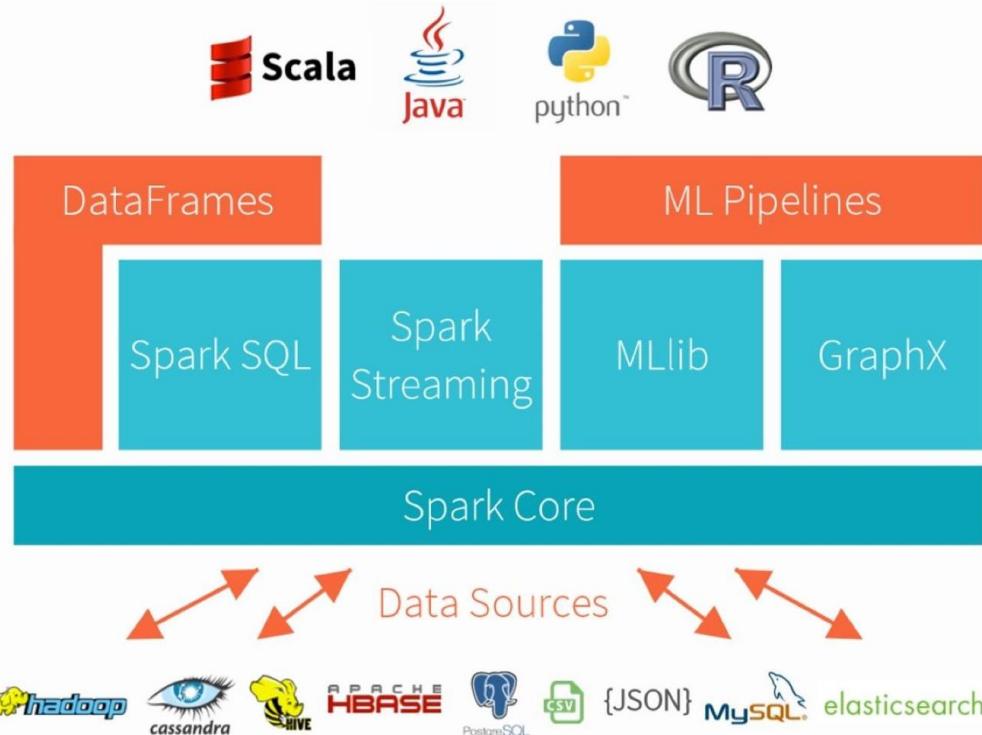
Apache Spark

Spark - фреймворк общего назначения для распределенных вычислений и аналитики над распределенными данными. С поддержкой API на таких языках как Scala, Java, Python, R. Spark поддерживает широкий набор возможностей: поддержка графовых вычислений, распределенных алгоритмов машинного обучения, встроенный Spark SQL, Dataframes, обработку потоковых данных, структурированную обработку потоковых данных, low-latency потоковую обработку данных.

<https://spark.apache.org/>

Spark





External Data Sources API



{ JSON }

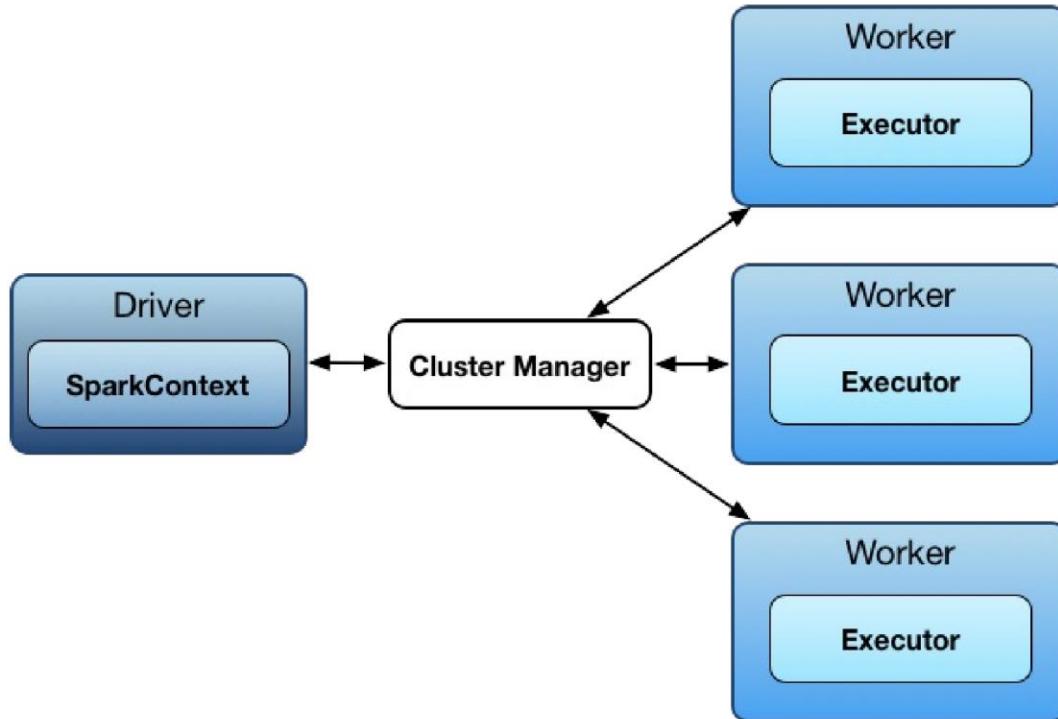


DataFrame

APACHE
HBASE

PostgreSQL





Варианты деплоимента

- Локально
- Amazon EC2
- Standalone Deploy
- YARN
- Kubernetes
- Другое
- Mesos

Пример

Scala - Число Pi

```
object SparkPi {  
    def main(args: Array[String]): Unit = {  
        val spark = SparkSession  
            .builder  
            .appName("Spark Pi")  
            .getOrCreate()  
        val slices = if (args.length > 0) args(0).toInt else 2  
        val n = math.min(100000L * slices, Int.MaxValue).toInt  
        val count = spark.sparkContext.parallelize(1 until n, slices).map { i =>  
            val x = random * 2 - 1  
            val y = random * 2 - 1  
            if (x*x + y*y <= 1) 1 else 0  
        }.reduce(_ + _)  
        println(s"Pi is roughly ${4.0 * count / (n - 1)}")  
        spark.stop()  
    }  
}
```

Python - Число Pi

```
if __name__ == "__main__":
    spark = SparkSession\
        .builder\
        .appName("PythonPi")\
        .getOrCreate()

    partitions = int(sys.argv[1]) if len(sys.argv) > 1 else 2
    n = 100000 * partitions

    def f(_):
        x = random() * 2 - 1
        y = random() * 2 - 1
        return 1 if x ** 2 + y ** 2 <= 1 else 0

    count = spark.sparkContext.parallelize(range(1, n + 1), partitions).map(f).reduce(add)
    print("Pi is roughly %f" % (4.0 * count / n))
```

Python - WordCount

```
if __name__ == "__main__":
    if len(sys.argv) != 2:
        print("Usage: wordcount <file>", file=sys.stderr)
        sys.exit(-1)

    spark = SparkSession\
        .builder\
        .appName("PythonWordCount")\
        .getOrCreate()

    lines = spark.read.text(sys.argv[1]).rdd.map(lambda r: r[0])
    counts = lines.flatMap(lambda x: x.split(' ')) \
        .map(lambda x: (x, 1)) \
        .reduceByKey(add)
    output = counts.collect()
    for (word, count) in output:
        print("%s: %i" % (word, count))

    spark.stop()
```

RDD

RDD — Resilient Distributed Dataset

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

Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley,
Michael J. Franklin, Scott Shenker, Ion Stoica
University of California, Berkeley

https://cs.stanford.edu/~matei/papers/2012/nsdi_spark.pdf

RDD — Resilient Distributed Dataset

- **Resilient** Эластичный, стойкий - отказоустойчивость позволяющая производить пересчёт даже в случае выхода из строя нод кластера
- **Distributed** Распределенный - данные распределены по кластеру, состоящему из многих нод
- **Dataset** Набор данных - коллекцияパーティцированных данных примитивов или комплексных типов

RDD — Resilient Distributed Dataset

Основная абстракция Spark позволяющая производить отказоустойчивые кластерные вычисления в памяти

RDD

```
abstract class RDD[T: ClassTag](  
    @transient private var _sc: SparkContext,  
    @transient private var deps: Seq[Dependency[_]]  
) extends Serializable with Logging {
```

RDD — Основные внутренние свойства

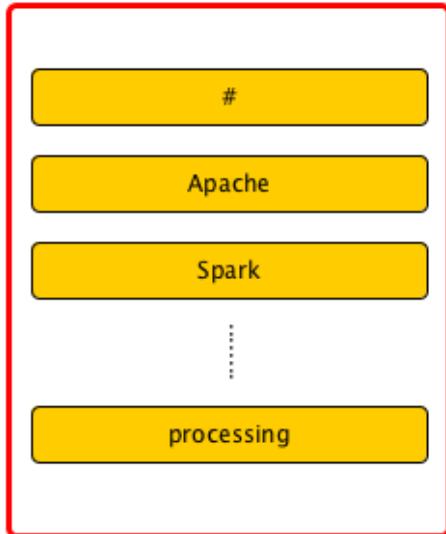
- Список партиций
- Функция которая вычисляет каждый сплит
- Список зависимых RDD
- Partitioner(опционально) - то какパーティционируется key-value RDD
- Список предпочтаемых мест(опционально), где будут происходить вычисление каждого сплита(к примеру местоположение HDFS блока)

RDD — Дополнительные особенности

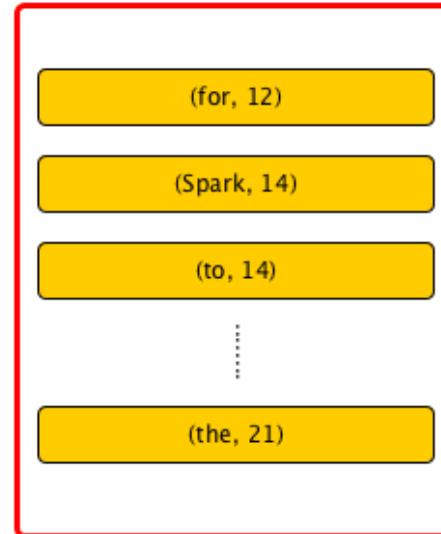
- **In-Memory**, Данные RDD хранятся в памяти до тех пор пока это возможно (по времени, по размеру)
 - **Иммутабельность**, данные не изменяемые, любые трансформации данных получаются только путем создания нового RDD
 - **Ленивые вычисления**, Нет доступа к данным и результатам трансформации, ровно до того момента пока action не вызовет вычисления явно
 - **Кешируемость**, данные можно закешировать на различных носителях (RAM,HDD)
 - **Параллельность**, Данные обрабатываются в параллельно
 - **Типизированы** — RDD Строго типизированы(RDD[Long] ,RDD[(Double,Long)], RDD[User])
 - **Партицированы**— Данные RDD разделены на логические партиции , которые распределены по нодам кластера
 - **Location-Stickiness** — RDD возможность определить наиболее близкий блок данных для партиции
-

RDD — Одна коллекция - много партиций

RDD of Strings



RDD of Pairs



Партиции

Партиция

Логический кусок данных(split) на которые делится RDD. Каждый RDD состоит из какого-то количества партиций. И является базовой единицей над которой происходит параллельно вычисление(Task).

Партиция

```
trait Partition extends Serializable {  
    /**  
     * Get the partition's index within its parent RDD  
     */  
    def index: Int  
  
    // A better default implementation of hashCode  
    override def hashCode(): Int = index  
  
    override def equals(other: Any): Boolean = super.equals(other)  
}
```

Партиция

- Содержит записи RDD
- Можно управлять количеством партиций, с помощью функций coalesce и repartition
- Можно определить как partitionируются данные(по умолчанию это HashPartitioner)

Партиция

- Как правило управление партициями просходит автоматически.
 - Чем больше/меньше партиций тем больше/меньше распределены данные по кластеру
 - Чем меньше/больше партиций тем больше/меньше данных обрабатывается за 1 таск
 - Spark запускает 1 многопоточный таск на партицию
 - Если у вас 50 ядер на кластере, то для их полной утилизации нужно как минимум 50 партиций
 - Количество партиций определяет количество файлов который Spark сгенерирует action при сохранении данных на диск.
 - Размер памяти экзекьютора определяет максимальный размер партиции
-

repartition

```
repartition(numPartitions:Int)(implicit ord:Ordering[T] = null):RDD[T]
```

Возвращает новый RDD с указанным количеством партиций
(выполняет coalesce), при этом возможен шафл (shuffle)

repartition

```
coalesce(numPartitions: Int, shuffle:Boolean = false)(implicit ord: Ordering[T] = null): RDD[T]
```

Трасформирует указанный RDD в новый с указанным количеством партиций

Операции над RDD

Типы операций над RDD

- Transformations
- Actions

Transformation

Ленивые операции которые на вход принимают RDD , делают некоторое преобразование, и возвращают новый RDD, при этом исходный RDD остается в неизменном виде, а непосредственных вычислений не происходит ровно до того момента, пока не будет вызыван action в дальнейшей цепочки вычислений.

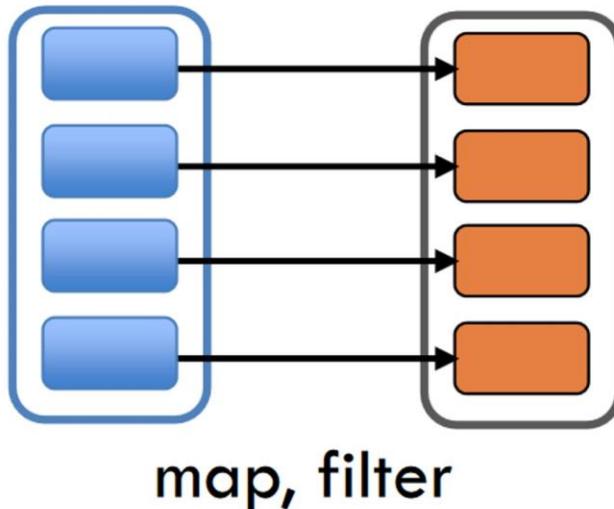
Transformation

- map
- filter
- flatMap
- mapPartitions
- mapPartitionsWithIndex
- groupBy
- sortBy
- sample
- randomSplit
- union
- intersection
- subtract
- distinct
- cartesian
- zip
- keyBy
- zipWithIndex
- zipWithUniqueID
- zipPartitions
- coalesce
- repartition
- repartitionAndSortWithinPartitions
- pipe

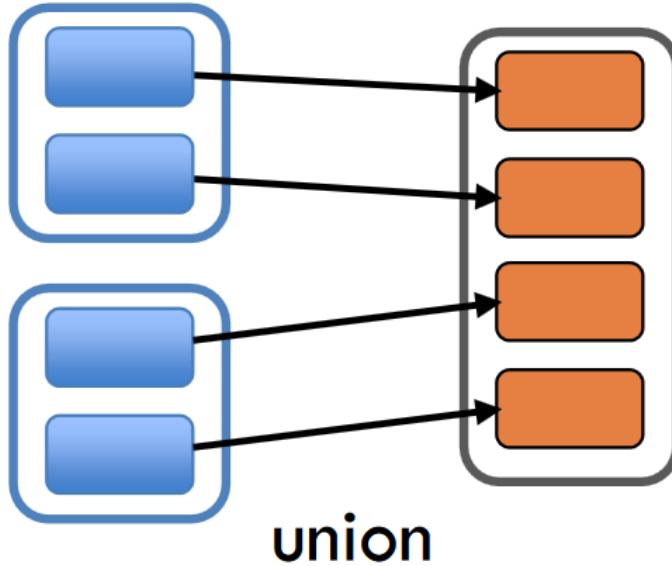
Transformation

- flatMapValues
- groupByKey
- reduceByKey
- reduceByKeyLocally
- foldByKey
- aggregateByKey
- sortByKey
- combineByKey
- sampleByKey
- cogroup (=groupWith)
- join
- subtractByKey
- fullOuterJoin
- leftOuterJoin
- rightOuterJoin
- partitionBy

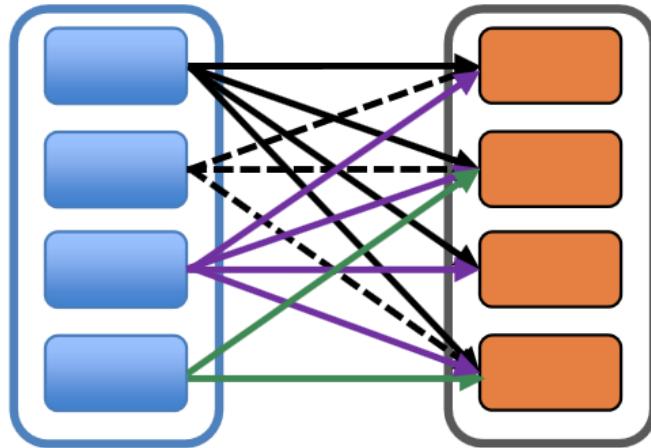
Narrow transformation



Narrow transformation

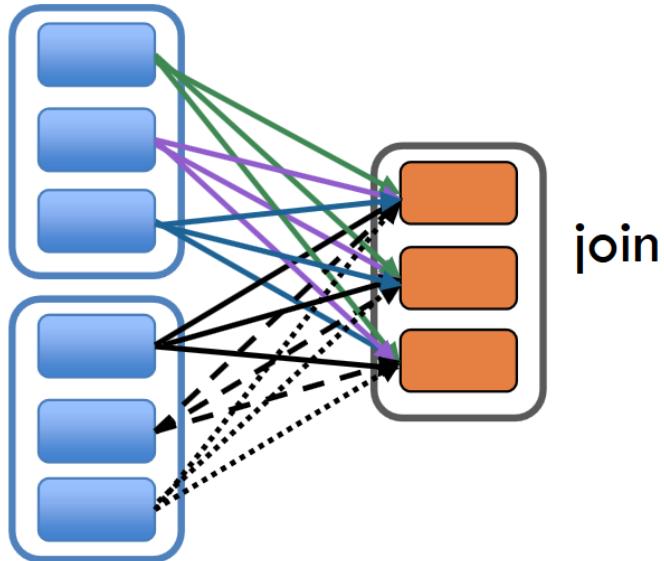


Wide Transformation



`groupByKey`

Wide Transformation



Actions

Операции над RDD которые запускают цепочку вычислений вычисления и возвращают какое-то не RDD значение. Только action может материализовать исчисление

Actions

- reduce
- collect
- aggregate
- fold
- first
- take
- foreach
- top
- treeAggregate
- treeReduce
- foreachPartition
- collectAsMap
- count
- takeSample
- max
- min
- sum
- histogram
- mean
- variance
- stdev
- sampleVariance
- countApprox
- countApproxDistinct
- takeOrdered
- saveAsTextFile
- saveAsSequenceFile
- saveAsObjectFile
- saveAsHadoopDataset
- saveAsHadoopFile
- saveAsNewAPIHadoopDataset
- saveAsNewAPIHadoopFile

Actions

- keys
- values
- countByKey
- countByValue
- countByValueApprox
- countApproxDistinctByKey
- countApproxDistinctByKey
- countByKeyApprox
- sampleByKeyExact

Transformations и Actions

```
val lines = sc.textFile("data.txt")
val lineLengths = lines.map(s => s.length)
val totalLength = lineLengths.reduce((a, b) => a + b)
```

Shuffling

Shuffling

Shuffling(перемешивание) - процесс перераспределения данных по кластеру, который сопровождается при многих операциях.

Shuffling

- Нужно избегать всеми способами

```
scala> val kv = (0 to 5) zip Stream.continually(5)
kv: scala.collection.immutable.IndexedSeq[(Int, Int)] = Vector((0,5), (1,5), (2,5), (3,5), (4,5), (5,5))

scala> val kw = (0 to 5) zip Stream.continually(10)
kw: scala.collection.immutable.IndexedSeq[(Int, Int)] = Vector((0,10), (1,10), (2,10), (3,10), (4,10), (5,10))

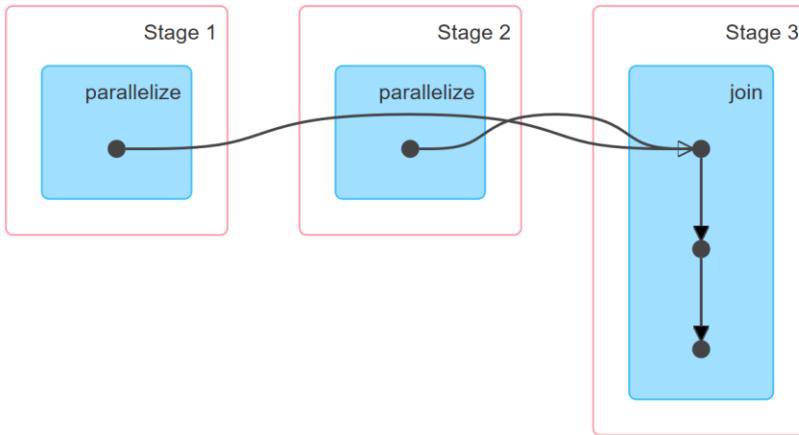
scala> val kvR = sc.parallelize(kv)
kvR: org.apache.spark.rdd.RDD[(Int, Int)] = ParallelCollectionRDD[2] at parallelize at <console>:27

scala> val kwR = sc.parallelize(kw)
kwR: org.apache.spark.rdd.RDD[(Int, Int)] = ParallelCollectionRDD[3] at parallelize at <console>:27

scala> val joined = kvR join kwR
joined: org.apache.spark.rdd.RDD[(Int, (Int, Int))] = MapPartitionsRDD[6] at join at <console>:28

scala> joined.toDebugString
res3: String =
(8) MapPartitionsRDD[6] at join at <console>:28 []
 |  MapPartitionsRDD[5] at join at <console>:28 []
 |  CoGroupedRDD[4] at join at <console>:28 []
 +- (8) ParallelCollectionRDD[2] at parallelize at <console>:27 []
 +- (8) ParallelCollectionRDD[3] at parallelize at <console>:27 []

scala> joined.foreach(println)
(1,(5,10))
(5,(5,10))
(3,(5,10))
(0,(5,10))
(4,(5,10))
(2,(5,10))
```



Completed Stages (3)

Stage Id	Description	Submitted	Duration
3	foreach at <console>:27	+details 2019/11/03 21:36:04	0,1 s
2	parallelize at <console>:27	+details 2019/11/03 21:36:04	55 ms
1	parallelize at <console>:27	+details 2019/11/03 21:36:04	85 ms

Кешировние

Persist и Cache

Оптимизация для итеративных и интерактивных вычислений.

Позволяют сохранить результат промежуточного вычисления в памяти/другом носителе, для последующего переиспользования.

Отличие persist от cache лишь синтаксическое cache = persist(MEMORY_ONLY)

Также есть возможность удалить данные из кеша с помощью функции unpersist()

Persist Cache

```
persist(): this.type
```

```
persist(newLevel: StorageLevel): this.type
```

Указать StorageLevel для заданного персиста можно однократно.

Persist Cache

- Оптимизация для итеративных и интерактивных вычислений.
Позволяют сохранить результат промежуточного вычисления в памяти/другом носителе, для последующего переиспользования
 - Отличие persist от cache лишь синтаксическое cache = persist(MEMORY_ONLY)
 - Также есть возможность удалить данные из кеша с помощью функции unpersist()
-

Уровни хранения

- NONE
 - DISK_ONLY
 - DISK_ONLY_2
 - MEMORY_ONLY
 - MEMORY_ONLY_2
 - MEMORY_ONLY_SER
 - MEMORY_ONLY_SER_2
 - MEMORY_AND_DISK
 - MEMORY_AND_DISK_2
 - MEMORY_AND_DISK_SER
 - MEMORY_AND_DISK_SER_2
 - OFF_HEAP
-

Стейджи

Стейджи

[Jobs](#)[Stages](#)[Storage](#)[Environment](#)[Executors](#)

Stages for All Jobs

[Completed Stages: 1](#)

Completed Stages (1)

Stage Id ▾	Description	Submitted	Duration	Tasks: Succeeded/Total
0	reduce at <console>:32 +details	2019/11/03 16:55:49	2,8 min	2000/2000

Аккумуляторы

Аккумулятор

Аккумулятор - переменная доступная только на запись для экзекьютора, при этом запись выражается функцией += (добавить), и на чтение для драйвера.

Аккумуляторы

Accumulators										
Accumulable	Value									
counter	45									
Tasks										
Index ▲	ID	Attempt	Status	Locality Level	Executor ID / Host	Launch Time	Duration	GC Time	Accumulators	Errors
0	0	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2016/04/21 10:10:41	17 ms			
1	1	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2016/04/21 10:10:41	17 ms		counter: 1	
2	2	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2016/04/21 10:10:41	17 ms		counter: 2	
3	3	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2016/04/21 10:10:41	17 ms		counter: 7	
4	4	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2016/04/21 10:10:41	17 ms		counter: 5	
5	5	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2016/04/21 10:10:41	17 ms		counter: 6	
6	6	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2016/04/21 10:10:41	17 ms		counter: 7	
7	7	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2016/04/21 10:10:41	17 ms		counter: 17	

Аккумуляторы - Scala

```
scala> val accum = sc.longAccumulator("My Accumulator")
accum: org.apache.spark.util.LongAccumulator = LongAccumulator(id: 0, name: Some(My Accumulator), value: 0)

scala> sc.parallelize(Array(1, 2, 3, 4)).foreach(x => accum.add(x))

scala> accum.value
res2: Long = 10
```

Аккумуляторы - Python

```
>>> accum = sc.accumulator(0)
```

```
>>> accum
```

```
Accumulator<id=0, value=0>
```

```
>>> sc.parallelize([1, 2, 3, 4]).foreach(lambda x: accum.add(x))
```

```
>>> accum.value
```

```
10
```

Нестандартный аккумулятор - Scala

```
class VectorAccumulatorV2 extends AccumulatorV2[MyVector, MyVector] {  
    private val myVector: MyVector = MyVector.createZeroVector  
  
    def reset(): Unit = {  
        myVector.reset()  
    }  
  
    def add(v: MyVector): Unit = {  
        myVector.add(v)  
    }  
    ...  
}  
val myVectorAcc = new VectorAccumulatorV2  
sc.register(myVectorAcc, "MyVectorAcc1")
```

Нестандартный аккумулятор - Python

```
class VectorAccumulatorParam(AccumulatorParam):  
    def zero(self, initialValue):  
        return Vector.zeros(initialValue.size)  
  
    def addInPlace(self, v1, v2):  
        v1 += v2  
        return v1  
  
vecAccum = sc.accumulator(Vector(...), VectorAccumulatorParam())
```

Аккумуляторы - Особенности

- Значение `accum.value` доступно только на драйвере
 - Операция `+=` допустима только в `action`
 - Никак не влияет на модель ленивых вычислений
-

Аккумуляторы - особенности

```
val accum = sc.longAccumulator  
data.map { x => accum.add(x); x }  
// Здесь accum.value все еще 0 ,так как не было никаких action
```

Broadcast variables

Broadcast переменная

Для оптимизации выполнения программы, можно заранее переслать экземпляр заранее вычисленной коллекции с драйвера на каждый экзекьютор

Broadcast переменная

```
scala> val broadcastVar = sc.broadcast(Array(1, 2, 3))
broadcastVar: org.apache.spark.broadcast.Broadcast[Array[Int]] = Broadcast@0
```

```
scala> broadcastVar.value
res0: Array[Int] = Array(1, 2, 3)
```

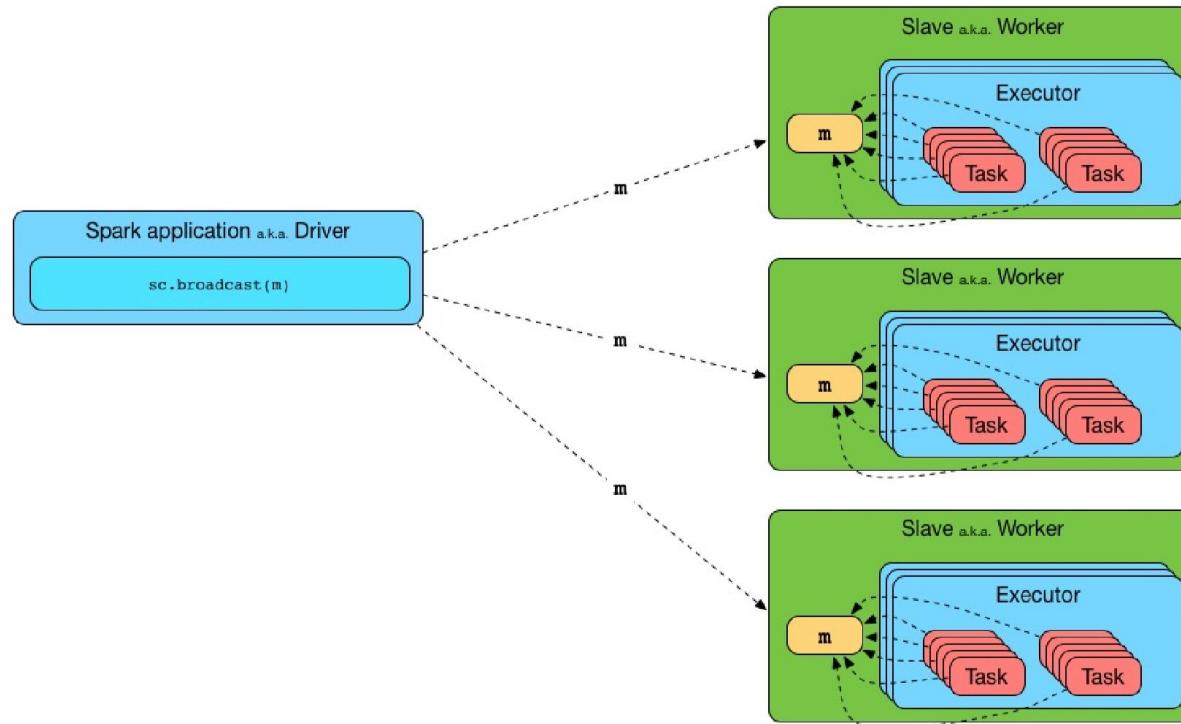
Аккумуляторы- Особенности

```
val acMap = sc.broadcast(myRDD.map { case (a,b,c,b) => (a, c) }.collectAsMap)
```

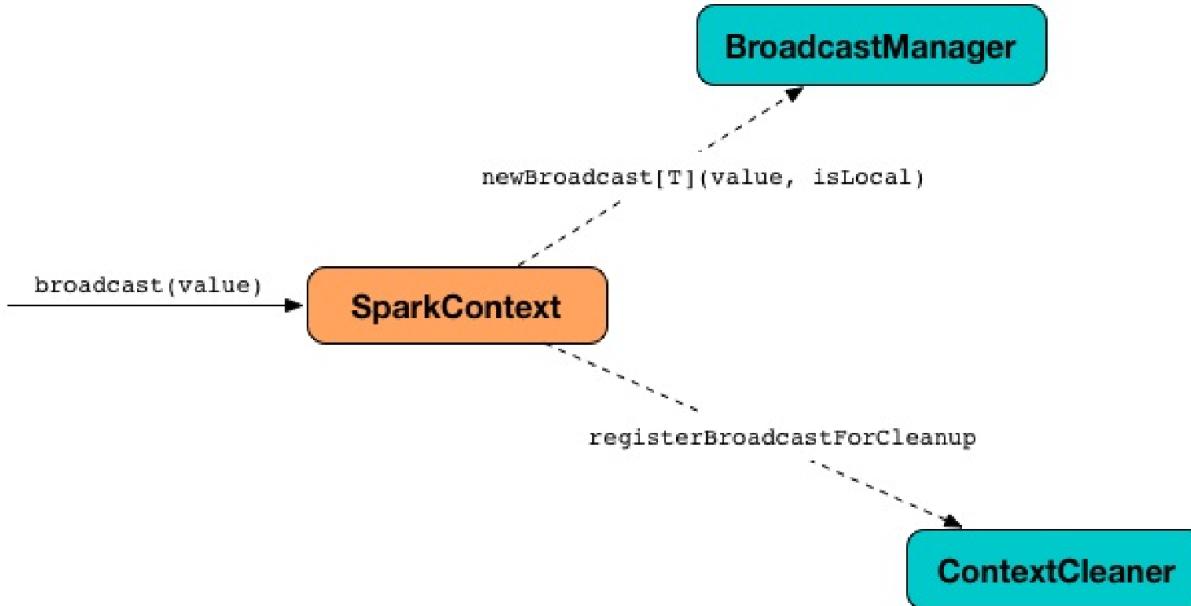
```
val otherMap = sc.broadcast(myOtherRDD.collectAsMap)
```

```
myBigRDD.map { case (a, b, c, d) =>  
    (acMap.value.get(a).get, otherMap.value.get(c).get)  
}.collect
```

Broadcast переменная



Broadcast переменная



Тестовый запуск

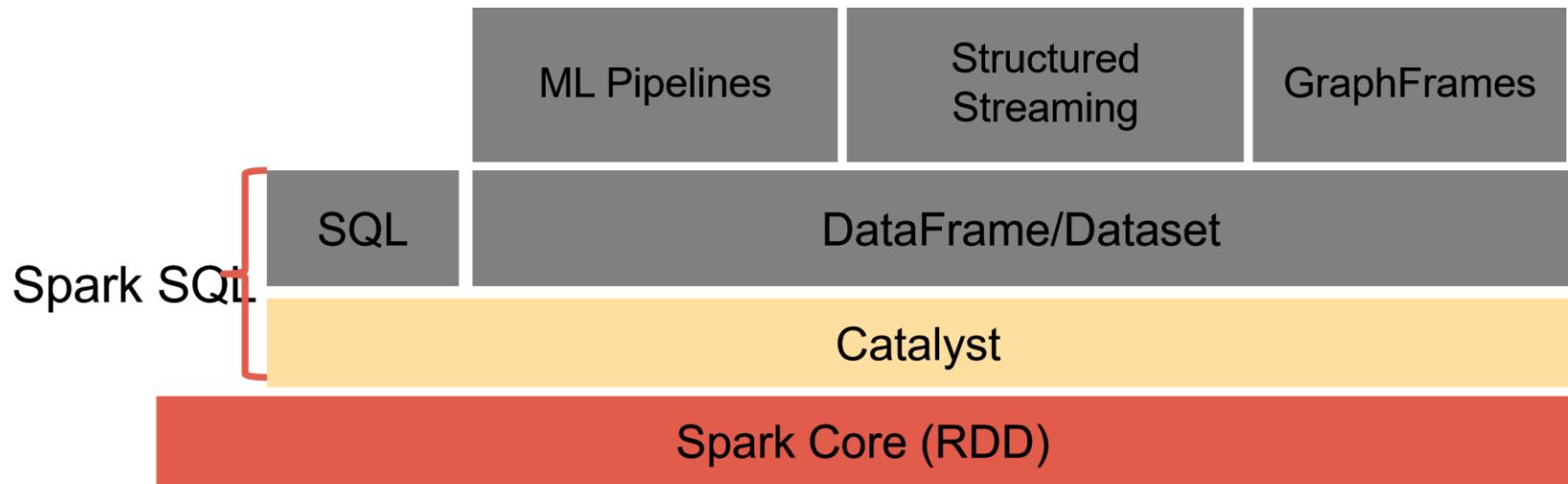
Jupyter ноутбук

```
docker run -p 8888:8888 -p 4040:4040 --rm -v "/Users/reireirei/Documents/data_lab:/home/jovyan/work"  
jupyter/pyspark-notebook
```

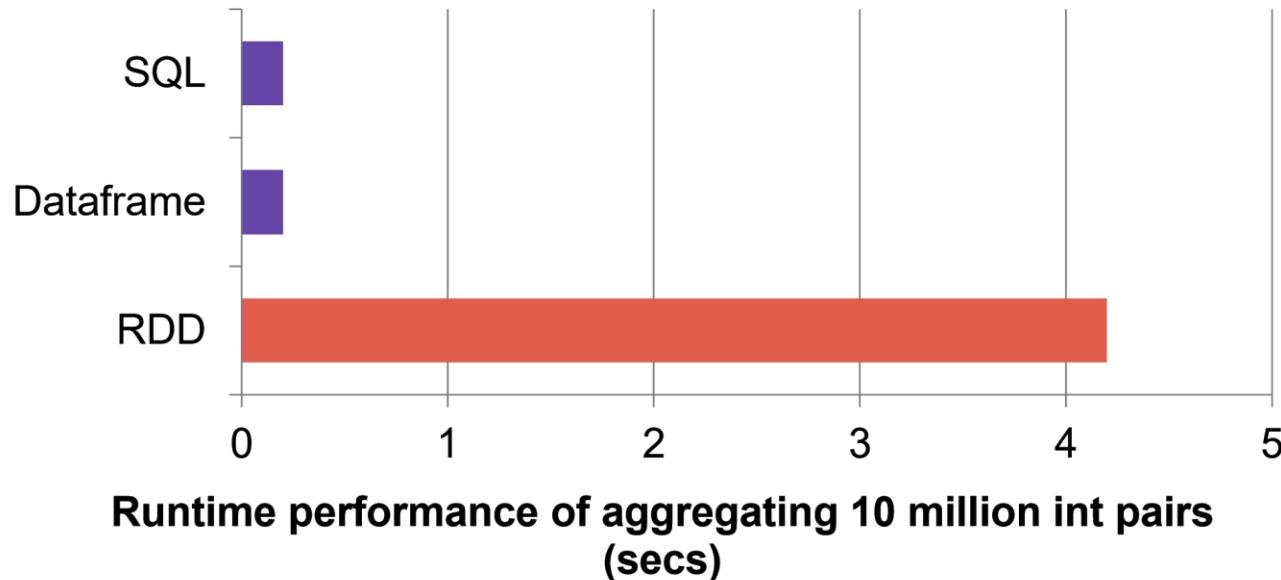
Demo

Dataframes

SparkSQL



Преимущества нового API



Новое API - зачем ?

Dataframe

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

SQL

```
select dept, avg(age) from data group by 1
```

RDD

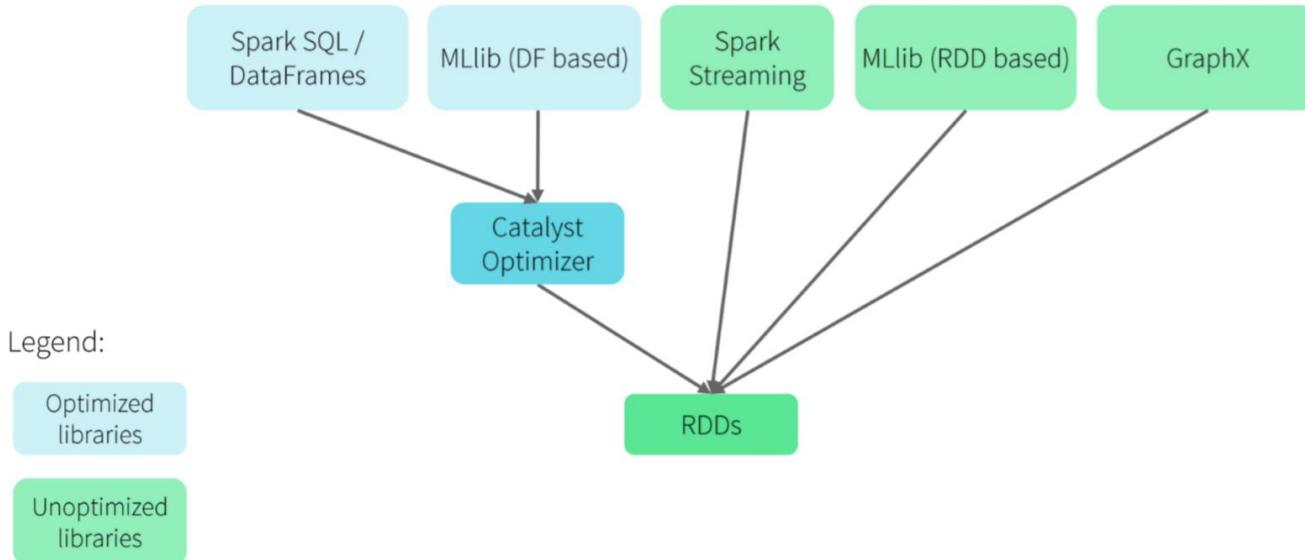
```
data.map { case (dept, age) => dept -> (age, 1) }
    .reduceByKey { case ((a1, c1), (a2, c2)) => (a1 + a2, c1 + c2)}
    .map { case (dept, (age, c)) => dept -> age / c }
```

Structured API

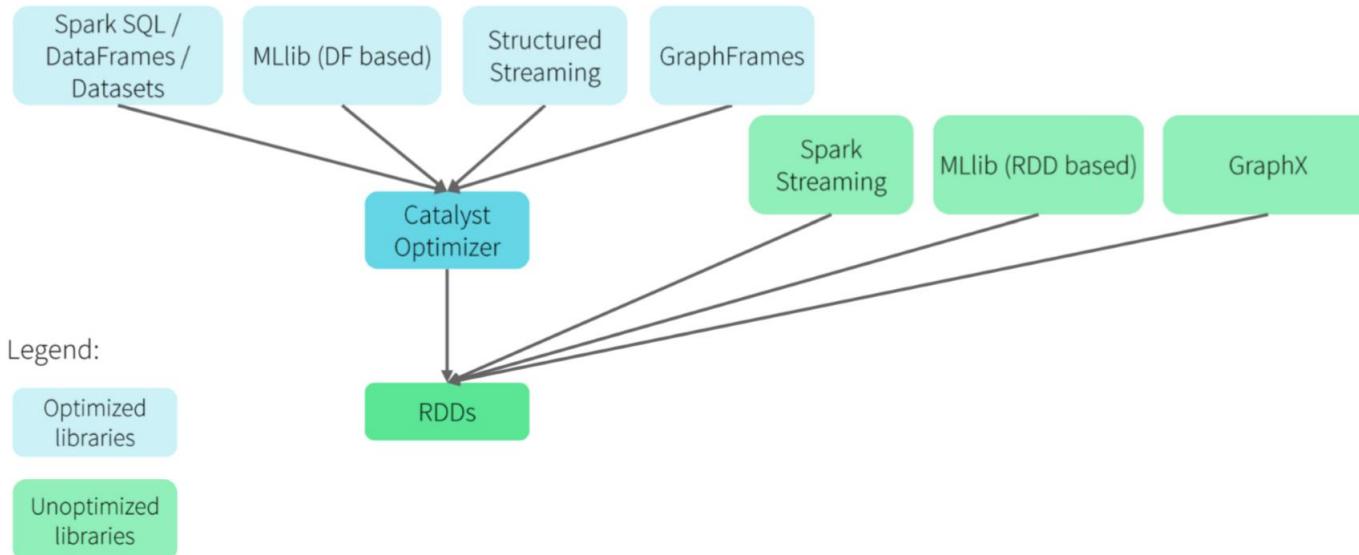
- Структура - уменьшает пространство того, что может быть выражено
- Большая часть вычислений может быть выражена

Уменьшение гибкости API приводит к тому, что вычисления можно сильно оптимизировать исходя из тех ограничений, что мы наложили.

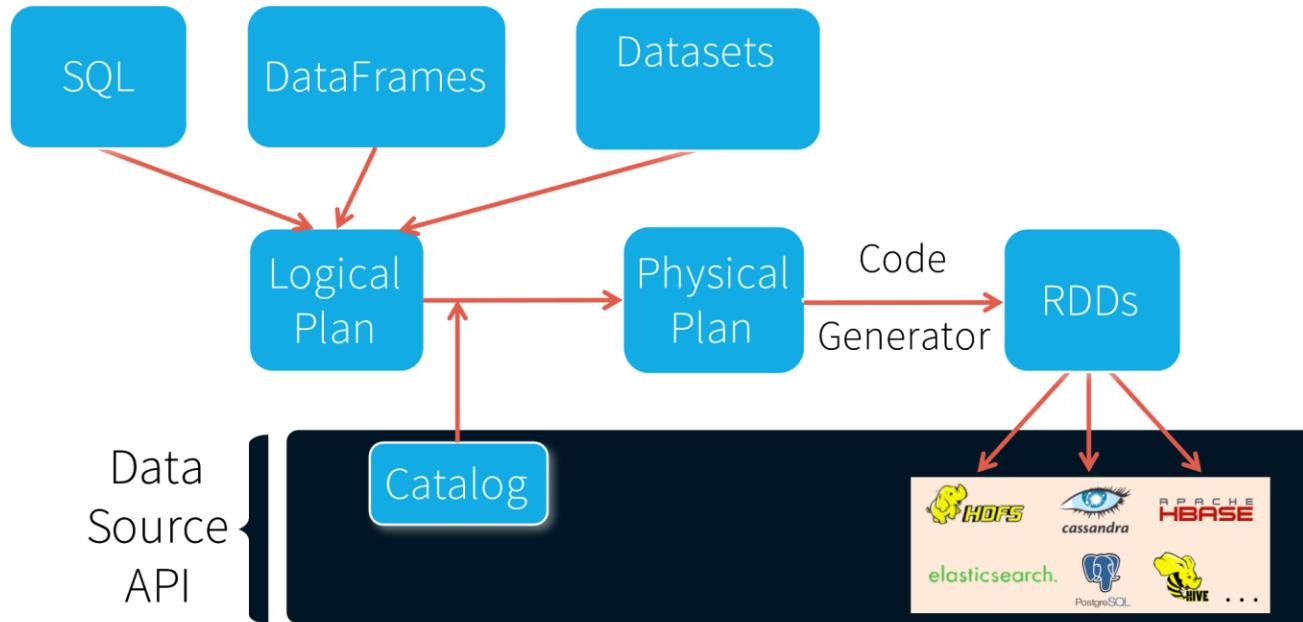
Spark 1.6.x



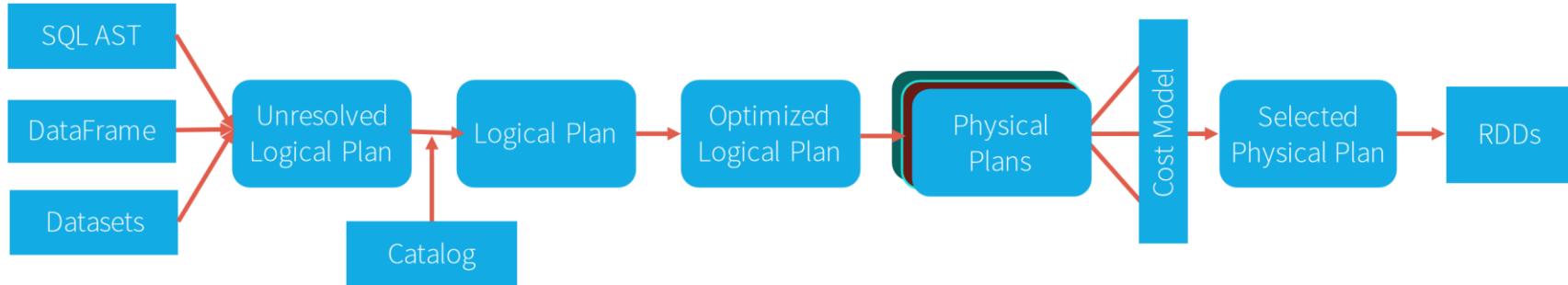
Spark 2.x.x



Новое API - архитектура



Новое API - архитектура



Откуда прирост?

Оптимизатор Spark'a пытается построить наиболее оптимальный план вычисления конкретной пользовательской программы, понимая с какими данными он работает, какие операции он производит, а также какие данные ожидаются на выходе

SparkSession - Scala

```
import org.apache.spark.sql.SparkSession

val spark = SparkSession
    .builder()
    .appName("Spark SQL basic example")
    .config("spark.some.config.option", "some-value")
    .getOrCreate()

// For implicit conversions like converting RDDs to DataFrames
import spark.implicits._
```

SparkSession - Scala

```
from pyspark.sql import SparkSession

spark = SparkSession \
    .builder \
    .appName("Python Spark SQL basic example") \
    .config("spark.some.config.option", "some-value") \
    .getOrCreate()
```

Создание DataFrame – people.json

```
{"name": "Michael"}  
{"name": "Andy", "age": 30}  
{"name": "Justin", "age": 19}
```



Создание DataFrame - Scala

```
val df = spark.read.json("examples/src/main/resources/people.json")

// Displays the content of the DataFrame to stdout
df.show()
// +---+-----+
// | age|    name|
// +---+-----+
// |null|Michael|
// |  30|    Andy|
// |  19| Justin|
// +---+-----+
```

Создание DataFrame - Python

```
# spark is an existing SparkSession
df = spark.read.json("examples/src/main/resources/people.json")
# Displays the content of the DataFrame to stdout
df.show()
# +---+---+
# | age| name|
# +---+---+
# | null|Michael|
# | 30| Andy|
# | 19| Justin|
# +---+---+
```

Создание DataFrame - Python

```
# spark is an existing SparkSession
df = spark.read.json("examples/src/main/resources/people.json")
# Displays the content of the DataFrame to stdout
df.show()
# +---+---+
# | age| name|
# +---+---+
# | null|Michael|
# | 30| Andy|
# | 19| Justin|
# +---+---+
```

Операции над DataFrame - Scala

```
// This import is needed to use the $-notation
import spark.implicits._

// Print the schema in a tree format
df.printSchema()
// root
// |-- age: long (nullable = true)
// |-- name: string (nullable = true)

// Select only the "name" column
df.select("name").show()
// +---+
// |  name|
// +---+
// |Michael|
// |  Andy|
// | Justin|
// +---+


// Select everybody, but increment the age by 1
df.select($"name", $"age" + 1).show()
// +-----+-----+
// |  name|(age + 1)|
// +-----+-----+
// |Michael|      null|
// |  Andy|       31|
// | Justin|      20|
// +-----+-----+


// Select people older than 21
df.filter($"age" > 21).show()
// +---+
// |age|name|
// +---+
// | 30|Andy|
// +---+
```

Операции над DataFrame - Scala

```
// Select people older than 21
df.filter($"age" > 21).show()
// +---+---+
// |age|name|
// +---+---+
// | 30|Andy|
// +---+---+

// Count people by age
df.groupBy("age").count().show()
// +---+---+
// | age|count|
// +---+---+
// | 19|    1|
// |null|    1|
// | 30|    1|
// +---+---+
```

Операции над DataFrame - Python

```
# spark, df are from the previous example
# Print the schema in a tree format
df.printSchema()
# root
# |-- age: long (nullable = true)
# |-- name: string (nullable = true)

# Select only the "name" column
df.select("name").show()

# +---+
# |   name|
# +---+
# |Michael|
# |   Andy|
# | Justin|
# +---+
```

```
# Select everybody, but increment the age by 1
df.select(df['name'], df['age'] + 1).show()
# +-----+-----+
# |    name|(age + 1)|
# +-----+-----+
# |Michael|      null|
# |   Andy|       31|
# | Justin|       20|
# +-----+-----+
# Select people older than 21
df.filter(df['age'] > 21).show()
# +---+
# |age|name|
# +---+
# | 30|Andy|
# +---+
```

Операции над DataFrame - Python

```
# Count people by age  
df.groupBy("age").count().show()  
# +---+---+  
# | age|count|  
# +---+---+  
# | 19|    1|  
# |null|    1|  
# | 30|    1|  
# +---+---+
```

Больше операций

[https://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.functions\\$](https://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.functions$)

SparkSQL - простой select

```
// Register the DataFrame as a SQL temporary view  
df.createOrReplaceTempView("people")  
  
val sqlDF = spark.sql("SELECT * FROM people")  
sqlDF.show()  
// +---+---+  
// | age| name|  
// +---+---+  
// |null|Michael|  
// | 30| Andy|  
// | 19| Justin|  
// +---+---+
```

Dataset – Scala

```
case class Person(name: String, age: Long)

// Encoders are created for case classes
val caseClassDS = Seq(Person("Andy", 32)).toDS()
caseClassDS.show()
// +----+----+
// |name|age|
// +----+----+
// |Andy| 32|
// +----+----+

// Encoders for most common types are automatically provided by importing spark.implicits._
val primitiveDS = Seq(1, 2, 3).toDS()
primitiveDS.map(_ + 1).collect() // Returns: Array(2, 3, 4)
```

Dataset – Scala

```
case class Person(name: String, age: Long)

// Encoders are created for case classes
val caseClassDS = Seq(Person("Andy", 32)).toDS()
caseClassDS.show()
// +----+----+
// |name|age|
// +----+----+
// |Andy| 32|
// +----+----+

// Encoders for most common types are automatically provided by importing spark.implicits._
val primitiveDS = Seq(1, 2, 3).toDS()
primitiveDS.map(_ + 1).collect() // Returns: Array(2, 3, 4)
```

Dataset – Scala

```
// DataFrames can be converted to a Dataset by providing a class. Mapping will be done by name
val path = "examples/src/main/resources/people.json"
val peopleDS = spark.read.json(path).as[Person]
peopleDS.show()
// +-----+
// | age|  name|
// +-----+
// |null|Michael|
// |  30|  Andy|
// |  19| Justin|
// +-----+
```

Dataset - Больше операций

<https://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.Dataset>

Взаимодействие с RDD – Scala

```
// For implicit conversions from RDDs to DataFrames
import spark.implicits._

// Create an RDD of Person objects from a text file, convert it to a Dataframe
val peopleDF = spark.sparkContext
    .textFile("examples/src/main/resources/people.txt")
    .map(_.split(","))
    .map(attributes => Person(attributes(0), attributes(1).trim.toInt))
    .toDF()

// Register the DataFrame as a temporary view
peopleDF.createOrReplaceTempView("people")
```

Взаимодействие с RDD – Scala

```
// SQL statements can be run by using the sql methods provided by Spark
val teenagersDF = spark.sql("SELECT name, age FROM people WHERE age BETWEEN 13 AND 19")

// The columns of a row in the result can be accessed by field index
teenagersDF.map(teenager => "Name: " + teenager(0)).show()
// +-----+
// |      value|
// +-----+
// |Name: Justin|
// +-----+
```

Взаимодействие с RDD – Scala

```
// SQL statements can be run by using the sql methods provided by Spark
val teenagersDF = spark.sql("SELECT name, age FROM people WHERE age BETWEEN 13 AND 19")

// The columns of a row in the result can be accessed by field index
teenagersDF.map(teenager => "Name: " + teenager(0)).show()
// +-----+
// |      value|
// +-----+
// |Name: Justin|
// +-----+
```

Взаимодействие с RDD – Scala

```
// or by field name
teenagersDF.map(teenager => "Name: " + teenager.getAs[String]("name")).show()
// +-----+
// |      value|
// +-----+
// |Name: Justin|
// +-----+

// No pre-defined encoders for Dataset[Map[K,V]], define explicitly
implicit val mapEncoder = org.apache.spark.sql.Encoders.kryo[Map[String, Any]]
// Primitive types and case classes can be also defined as
// implicit val stringIntMapEncoder: Encoder[Map[String, Any]] = ExpressionEncoder()

// row.getValuesMap[T] retrieves multiple columns at once into a Map[String, T]
teenagersDF.map(teenager => teenager.getValuesMap[Any](List("name", "age"))).collect()
// Array(Map("name" -> "Justin", "age" -> 19))
```

Взаимодействие с RDD – Scala

```
// or by field name
teenagersDF.map(teenager => "Name: " + teenager.getAs[String]("name")).show()
// +-----+
// |      value|
// +-----+
// |Name: Justin|
// +-----+

// No pre-defined encoders for Dataset[Map[K,V]], define explicitly
implicit val mapEncoder = org.apache.spark.sql.Encoders.kryo[Map[String, Any]]
// Primitive types and case classes can be also defined as
// implicit val stringIntMapEncoder: Encoder[Map[String, Any]] = ExpressionEncoder()

// row.getValuesMap[T] retrieves multiple columns at once into a Map[String, T]
teenagersDF.map(teenager => teenager.getValuesMap[Any](List("name", "age"))).collect()
// Array(Map("name" -> "Justin", "age" -> 19))
```

Взаимодействие с RDD – people.txt

Michael, 29

Andy, 30

Justin, 19

Взаимодействие с RDD – Python

```
from pyspark.sql import Row

sc = spark.sparkContext

# Load a text file and convert each line to a Row.
lines = sc.textFile("examples/src/main/resources/people.txt")
parts = lines.map(lambda l: l.split(","))
people = parts.map(lambda p: Row(name=p[0], age=int(p[1])))

# Infer the schema, and register the DataFrame as a table.
schemaPeople = spark.createDataFrame(people)
schemaPeople.createOrReplaceTempView("people")
```

Взаимодействие с RDD – Python

```
# SQL can be run over DataFrames that have been registered as a table.  
teenagers = spark.sql("SELECT name FROM people WHERE age >= 13 AND age <= 19")  
  
# The results of SQL queries are Dataframe objects.  
# rdd returns the content as an :class:`pyspark.RDD` of :class:`Row`.  
teenNames = teenagers.rdd.map(lambda p: "Name: " + p.name).collect()  
for name in teenNames:  
    print(name)  
# Name: Justin
```

Явное указание схемы данных - Scala

```
import org.apache.spark.sql.types._

// Create an RDD
val peopleRDD = spark.sparkContext.textFile("examples/src/main/resources/people.txt")
// The schema is encoded in a string
val schemaString = "name age"
// Generate the schema based on the string of schema
val fields = schemaString.split(" ")
    .map(fieldName => StructField(fieldName, StringType, nullable = true))
val schema = StructType(fields)
// Convert records of the RDD (people) to Rows
val rowRDD = peopleRDD
    .map(_.split(","))
    .map(attributes => Row(attributes(0), attributes(1).trim))
```

Явное указание схемы данных - Scala

```
// Apply the schema to the RDD
val peopleDF = spark.createDataFrame(rowRDD, schema)

// Creates a temporary view using the DataFrame
peopleDF.createOrReplaceTempView("people")

// SQL can be run over a temporary view created using DataFrames
val results = spark.sql("SELECT name FROM people")

// The results of SQL queries are DataFrames and support all the normal RDD operations
// The columns of a row in the result can be accessed by field index or by field name
results.map(attributes => "Name: " + attributes(0)).show()

// +-----+
// |      value|
// +-----+
// |Name: Michael|
// |  Name: Andy|
// | Name: Justin|
// +-----+
```

Явное указание схемы данных - Python

```
# Import data types
from pyspark.sql.types import *

sc = spark.sparkContext

# Load a text file and convert each line to a Row.
lines = sc.textFile("examples/src/main/resources/people.txt")
parts = lines.map(lambda l: l.split(","))
# Each line is converted to a tuple.
people = parts.map(lambda p: (p[0], p[1].strip()))
# The schema is encoded in a string.
schemaString = "name age"
fields = [StructField(field_name, StringType(), True) for field_name in schemaString.split()]
schema = StructType(fields)
```

Явное указание схемы данных - Python

```
# Apply the schema to the RDD.  
schemaPeople = spark.createDataFrame(people, schema)  
  
# Creates a temporary view using the DataFrame  
schemaPeople.createOrReplaceTempView("people")  
  
# SQL can be run over DataFrames that have been registered as a table.  
results = spark.sql("SELECT name FROM people")  
results.show()  
  
# +-----+  
# | name |  
# +-----+  
# |Michael|  
# | Andy |  
# | Justin|  
# +-----+
```

Пользовательские employees.json

```
{"name": "Michael", "salary": 3000}  
 {"name": "Andy", "salary": 4500}  
 {"name": "Justin", "salary": 3500}  
 {"name": "Berta", "salary": 4000}
```

Пользовательские агрегаторы(UDAF)

```
import org.apache.spark.sql.{Row, SparkSession}  
import org.apache.spark.sql.expressions.MutableAggregationBuffer  
import org.apache.spark.sql.expressions.UserDefinedAggregateFunction  
import org.apache.spark.sql.types._
```

Пользовательские агрегаторы(UDAF)

```
object MyAverage extends UserDefinedAggregateFunction {
```

Пользовательские агрегаторы(UDAF)

```
// Data types of input arguments of this aggregate function
def inputSchema: StructType = StructType(StructField("inputColumn", LongType) :: Nil)
// Data types of values in the aggregation buffer
def bufferSchema: StructType = {
    StructType(StructField("sum", LongType) :: StructField("count", LongType) :: Nil)
}
```

Пользовательские агрегаторы(UDAF)

```
def dataType: DataType = DoubleType
// Whether this function always returns the same output on the identical input
def deterministic: Boolean = true
// Initializes the given aggregation buffer. The buffer itself is a `Row` that in addition to
// standard methods like retrieving a value at an index (e.g., get(), getBoolean()), provides
// the opportunity to update its values. Note that arrays and maps inside the buffer are still
// immutable.
def initialize(buffer: MutableAggregationBuffer): Unit = {
    buffer(0) = 0L
    buffer(1) = 0L
}
```

Пользовательские агрегаторы(UDAF)

```
// Updates the given aggregation buffer `buffer` with new input data from `input`
def update(buffer: MutableAggregationBuffer, input: Row): Unit = {
    if (!input.isNullAt(0)) {
        buffer(0) = buffer.getLong(0) + input.getLong(0)
        buffer(1) = buffer.getLong(1) + 1
    }
}

// Merges two aggregation buffers and stores the updated buffer values back to `buffer1`
def merge(buffer1: MutableAggregationBuffer, buffer2: Row): Unit = {
    buffer1(0) = buffer1.getLong(0) + buffer2.getLong(0)
    buffer1(1) = buffer1.getLong(1) + buffer2.getLong(1)
}

// Calculates the final result
def evaluate(buffer: Row): Double = buffer.getLong(0).toDouble / buffer.getLong(1)
}
```

Пользовательские агрегаторы(UDAF)

```
// Register the function to access it
spark.udf.register("myAverage", MyAverage)

val df = spark.read.json("examples/src/main/resources/employees.json")
df.createOrReplaceTempView("employees")
df.show()
// +-----+-----+
// |    name|salary|
// +-----+-----+
// |Michael|  3000|
// |   Andy|  4500|
// | Justin|  3500|
// |  Berta|  4000|
// +-----+-----+
```

Пользовательские агрегаторы(UDAF)

```
val result = spark.sql("SELECT myAverage(salary) as average_salary FROM employees")
result.show()
// +-----+
// |average_salary|
// +-----+
// |      3750.0|
// +-----+
```

Типизированные пользовательские агрегаторы(Typed UDAF)

```
import org.apache.spark.sql.{Encoder, Encoders, SparkSession}  
import org.apache.spark.sql.expressions.Aggregator  
  
case class Employee(name: String, salary: Long)  
case class Average(var sum: Long, var count: Long)
```

Типизированные пользовательские агрегаторы(Typed UDAF)

```
import org.apache.spark.sql.{Encoder, Encoders, SparkSession}  
import org.apache.spark.sql.expressions.Aggregator  
  
case class Employee(name: String, salary: Long)  
case class Average(var sum: Long, var count: Long)
```

Типизированные пользовательские агрегаторы(Typed UDAF)

```
object MyAverage extends Aggregator[Employee, Average, Double] {
```

Типизированные пользовательские агрегаторы(Typed UDAF)

```
// A zero value for this aggregation. Should satisfy the property that any b + zero = b
def zero: Average = Average(0L, 0L)

// Combine two values to produce a new value. For performance, the function may modify `buffer`
// and return it instead of constructing a new object
def reduce(buffer: Average, employee: Employee): Average = {
    buffer.sum += employee.salary
    buffer.count += 1
    buffer
}

// Merge two intermediate values
def merge(b1: Average, b2: Average): Average = {
    b1.sum += b2.sum
    b1.count += b2.count
    b1
}
```

Типизированные пользовательские агрегаторы(Typed UDAF)

```
// Transform the output of the reduction
def finish(reduction: Average): Double = reduction.sum.toDouble / reduction.count
// Specifies the Encoder for the intermediate value type
def bufferEncoder: Encoder[Average] = Encoders.product
// Specifies the Encoder for the final output value type
def outputEncoder: Encoder[Double] = Encoders.scalaDouble
```

Типизированные пользовательские агрегаторы(Typed UDAF)

```
val ds = spark.read.json("examples/src/main/resources/employees.json").as[Employee]
ds.show()
// +-----+-----+
// |    name|salary|
// +-----+-----+
// |Michael|  3000|
// |   Andy|  4500|
// | Justin|  3500|
// |  Berta|  4000|
// +-----+-----+

// Convert the function to a `TypedColumn` and give it a name
val averageSalary = MyAverage.toColumn.name("average_salary")
val result = ds.select(averageSalary)
result.show()
// +-----+
// |average_salary|
// +-----+
// |          3750.0|
// +-----+
```

Источники данных - Python

```
df = spark.read.load("examples/src/main/resources/users.parquet")
df.select("name", "favorite_color").write.save("namesAndFavColors.parquet")
```

Источники данных - Scala

```
val peopleDF = spark.read.format("json").load("examples/src/main/resources/people.json")
peopleDF.select("name", "age").write.format("parquet").save("namesAndAges.parquet")

val peopleDFCsv = spark.read.format("csv")
.option("sep", ";")
.option("inferSchema", "true")
.option("header", "true")
.load("examples/src/main/resources/people.csv")

usersDF.write.format("orc")
.option("orc.bloom.filter.columns", "favorite_color")
.option("orc.dictionary.key.threshold", "1.0")
.save("users_with_options.orc")
```

Sql над файлами - Python

```
val df = spark.sql("SELECT * FROM parquet.`examples/src/main/resources/users.parquet`")
```

Режимы сохранения

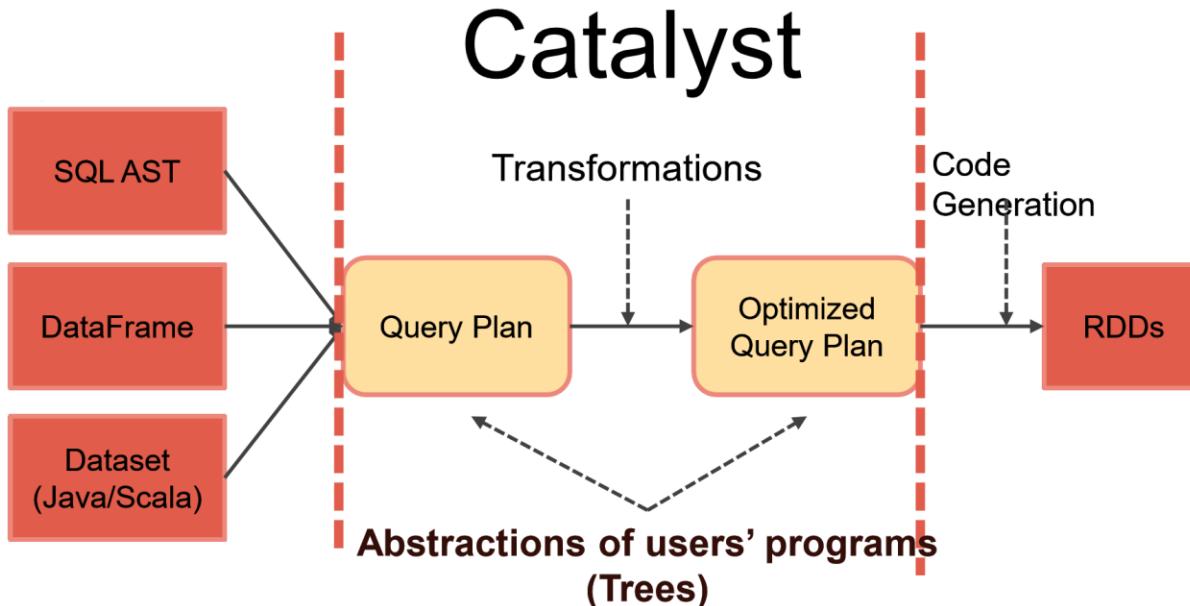
Scala/Java		Определение
SaveMode.ErrorIfExists (default)	"error" or "errorIfExists" (default)	Запись прерывается если данные уже есть
SaveMode.Append	"append"	Добавляет данные в конец
SaveMode.Overwrite	"overwrite"	Переписывает данные
SaveMode.Ignore	"ignore"	CREATE TABLE IF NOT EXISTS

Бакетирование, сортировка,パーティционирование

```
peopleDF.write.bucketBy(42, "name").sortBy("age").saveAsTable("people_bucketed")
usersDF.write.partitionBy("favorite_color").format("parquet").save("namesPartByColor.parquet")
usersDF
    .write
    .partitionBy("favorite_color")
    .bucketBy(42, "name")
    .saveAsTable("users_partitioned_bucketed")
```

Catalyst optimizator

Catalyst



Catalyst - Expression

```
SELECT sum(v)
FROM (
  SELECT
    t1.id,
    1 + 2 + t1.value AS v
  FROM t1 JOIN t2
  WHERE
    t1.id = t2.id AND
    t2.id > 50000) tmp
```

Catalyst - Expression

```
SELECT sum(v)
FROM (
  SELECT
    t1.id,
    1 + 2 + t1.value AS v
  FROM t1 JOIN t2
  WHERE
    t1.id = t2.id AND
    t2.id > 50000) tmp
```

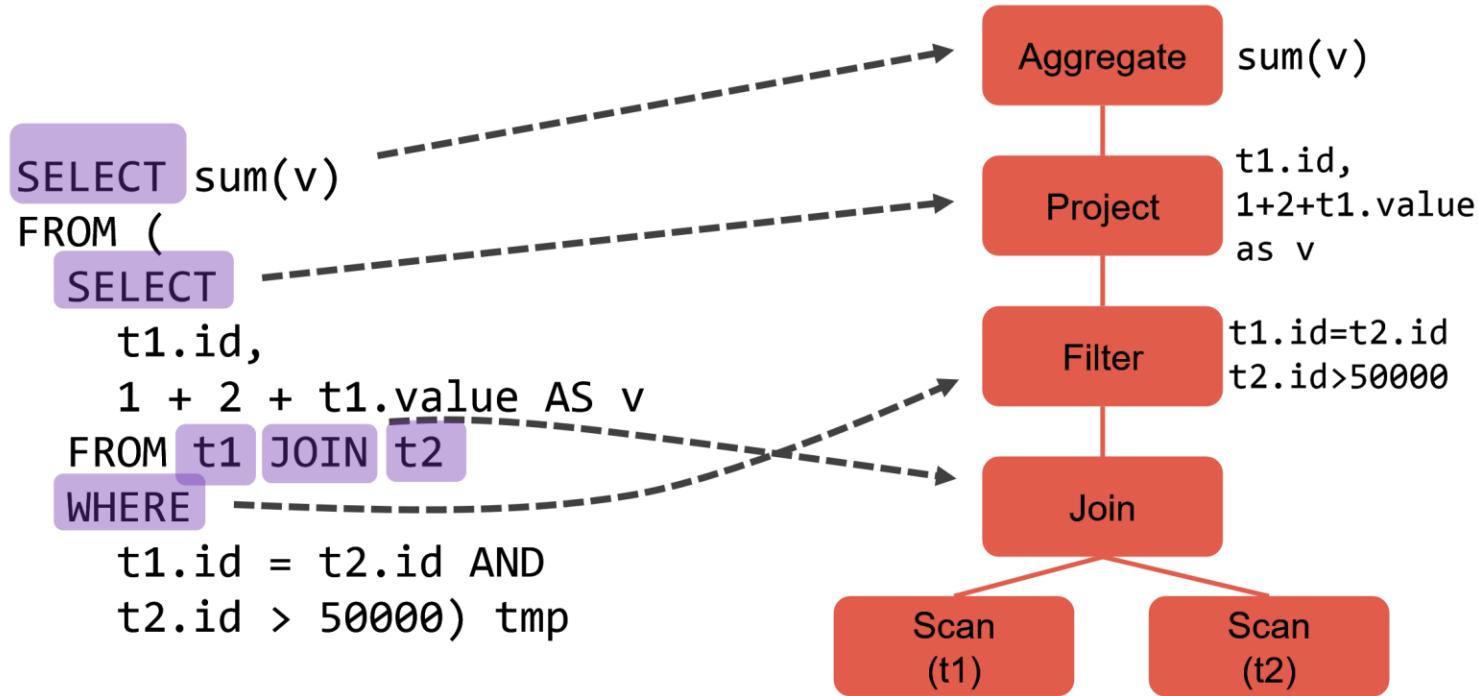
- Выражение представляет собой новое значение, которое получается вычислением входных данных($1 + 2 + t1.value$)
- Атрибут - колонка датасета или результирующая колонка какой либо функции

Catalyst - Expression

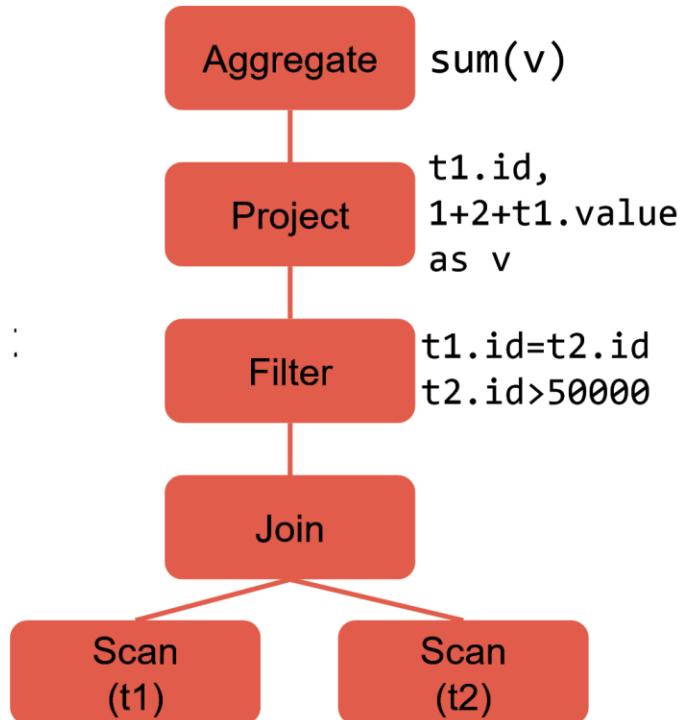
```
SELECT sum(v)
FROM (
  SELECT
    t1.id,
    1 + 2 + t1.value AS v
  FROM t1 JOIN t2
  WHERE
    t1.id = t2.id AND
    t2.id > 50000) tmp
```

- Выражение представляет собой новое значение, которое получается вычислением входных данных($1 + 2 + t1.value$)
- Атрибут - колонка датасета или результирующая колонка какой либо функции

Catalyst - Expression

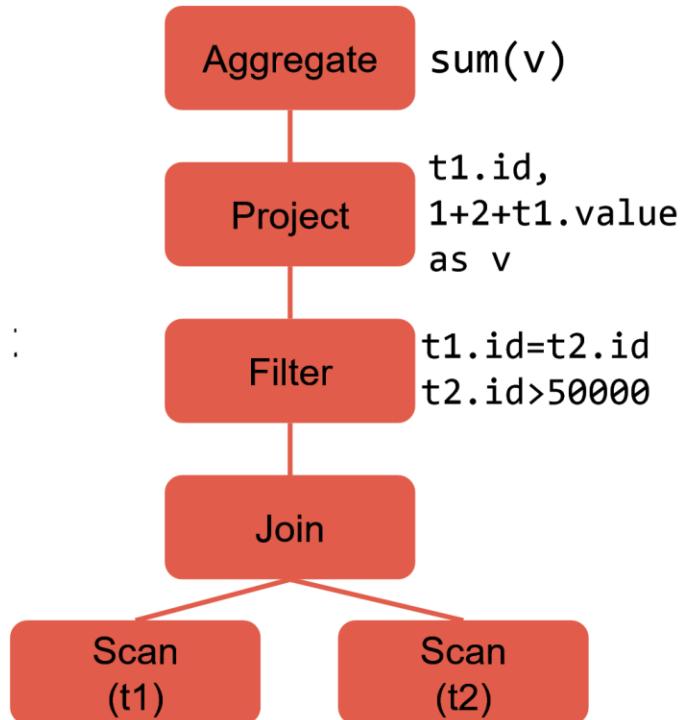


Catalyst – Logical Plan



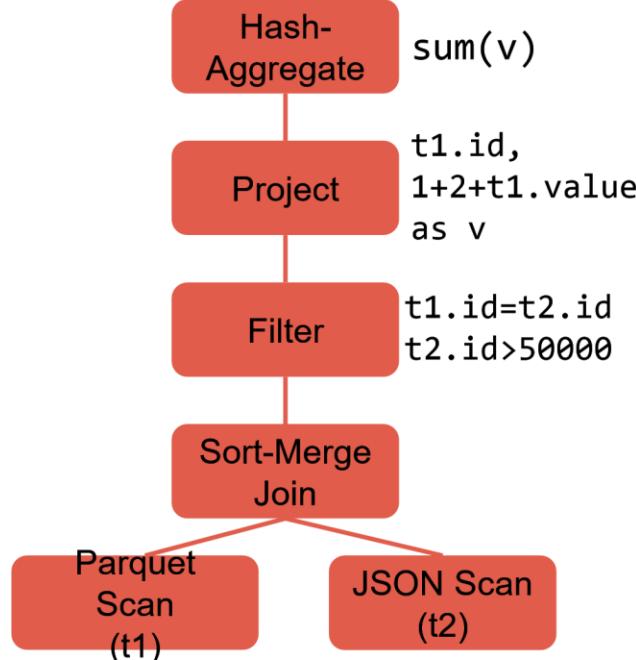
- Логический план описывает вычисление датасета, без описания того, как именно будут произведены вычисления
- Output: список атрибутов которые генерирует логический план [v,id]
- Ограничения(constraints): множество инвариантов для строк полученные из логического плана [t2.id>5000]
- Статистика: различные статистики(размер плана, max/min/null колонок и т.д.)

Catalyst – Логический план



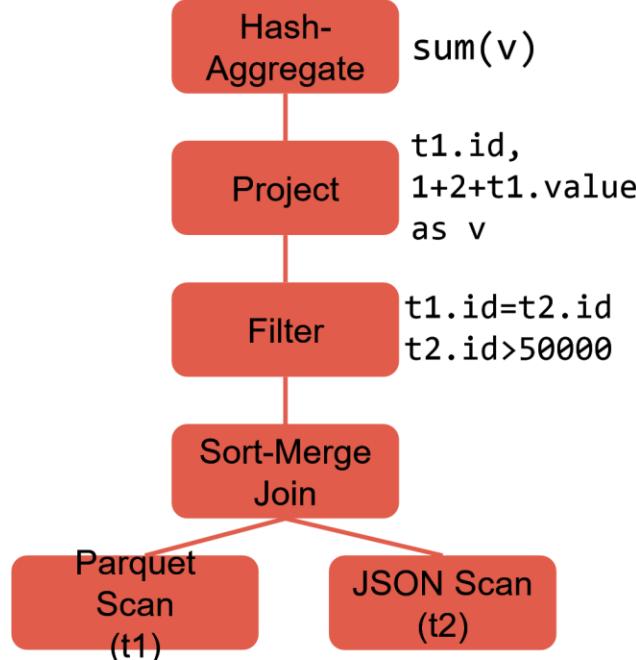
- **Логический план** описывает вычисление датасета, без описания того, как именно будут произведены вычисления
- **Output:** список атрибутов которые генерирует логический план $[v, id]$
- **Ограничения (constraints):** множество инвариантов для строк полученные из логического плана $[t2.id > 5000]$
- **Статистика:** различные статистики (размер плана, max/min/null колонок и т.д.)

Catalyst – Физический план



- Физический план описывает вычисление датасета, с указанием того, как именно будут получены и вычислены те или иные данные
- Можно запустить

Catalyst – Физический план



- Физический план описывает вычисление датасета, с указанием того, как именно будут получены и вычислены те или иные данные
- Можно запустить

Catalyst – Transformation

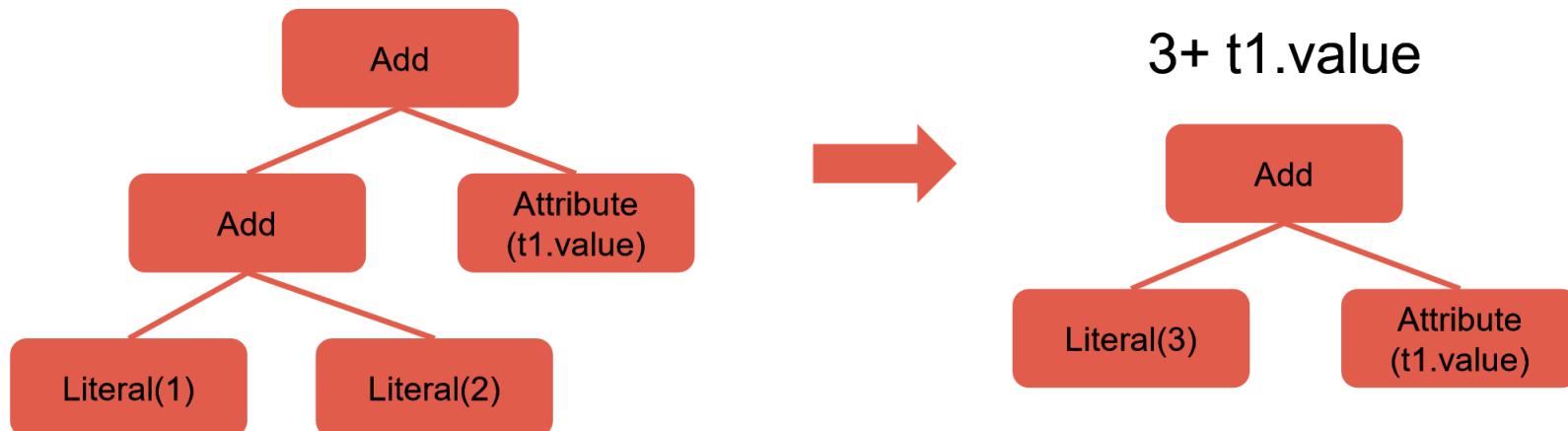
Трасформации без изменения типа(Transform и Rule Executor)

- Expression => Expression
- Logical Plan => Logical Plan
- Physical Plan => Physical Plan

Трансформации с изменением типа дерева
Logical Plan => Physical Plan

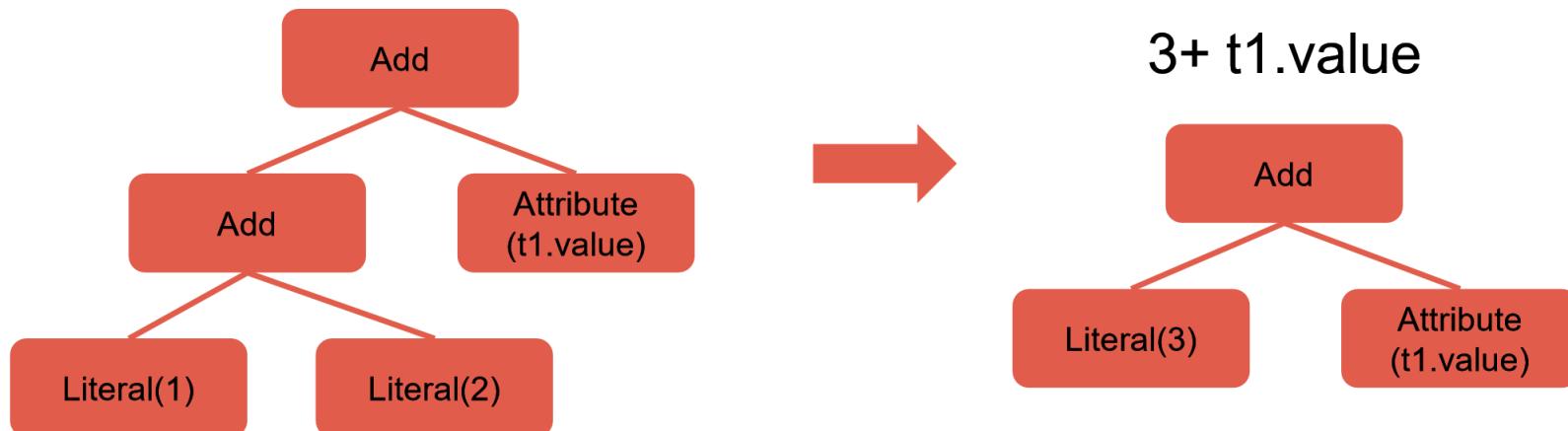
Catalyst – Transformation

$1 + 2 + t1.value$



Catalyst – Transformation

$1 + 2 + t1.value$

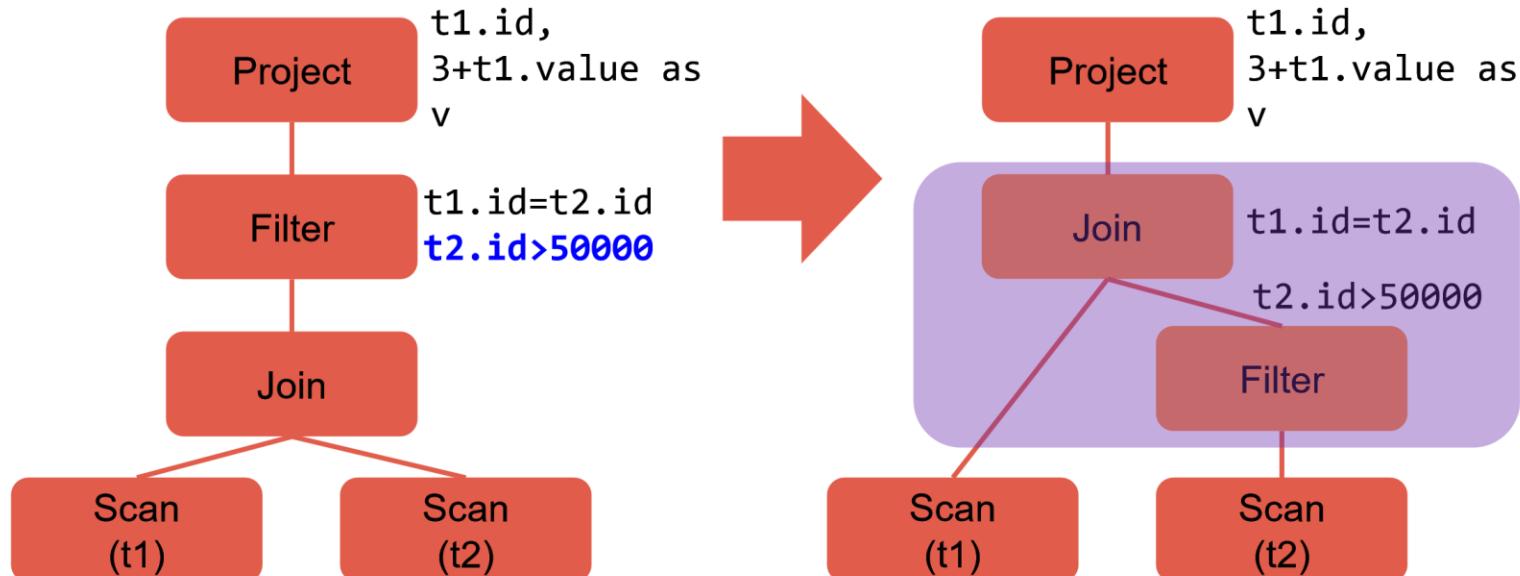


Catalyst – Transformation

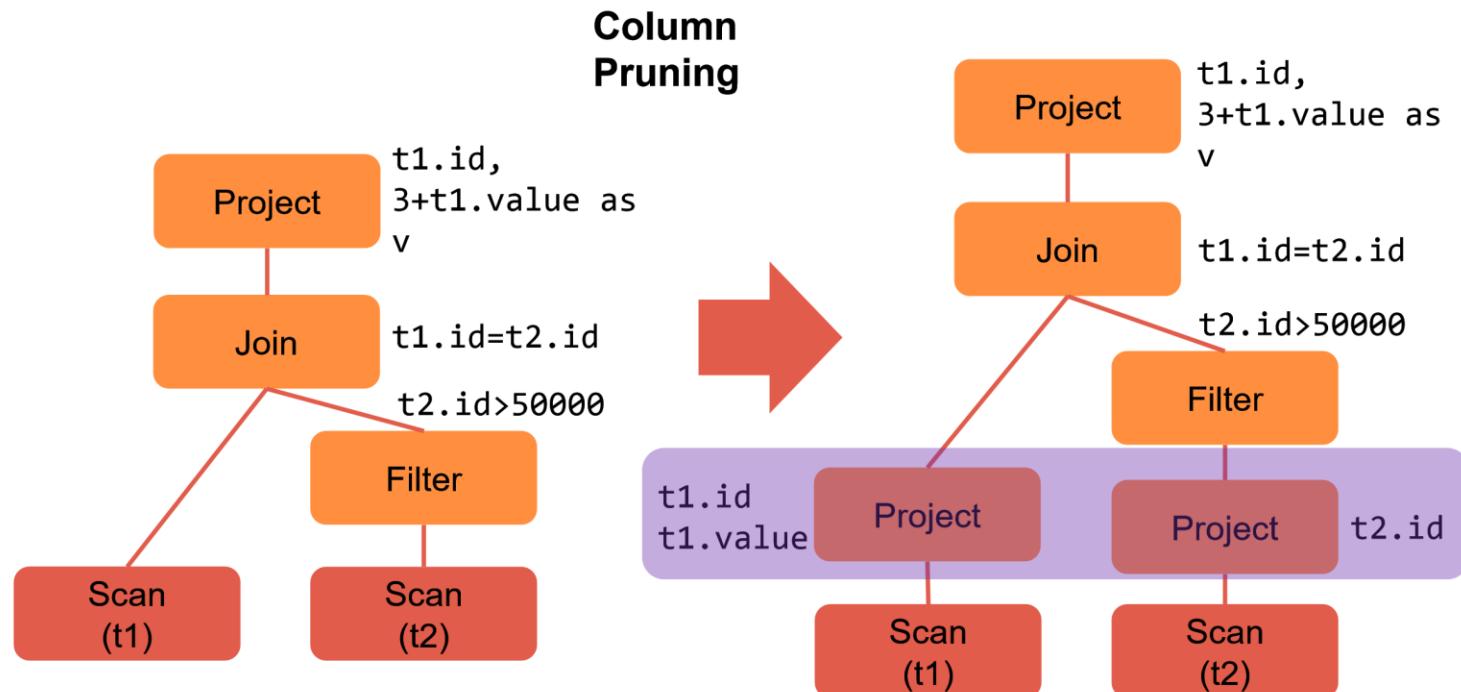
```
val expression: Expression = ...
expression.transform {
    case Add(Literal(x, IntegerType), Literal(y, IntegerType)) =>
        Literal(x + y)
}
```

Catalyst – Комбинирование нескольких правил

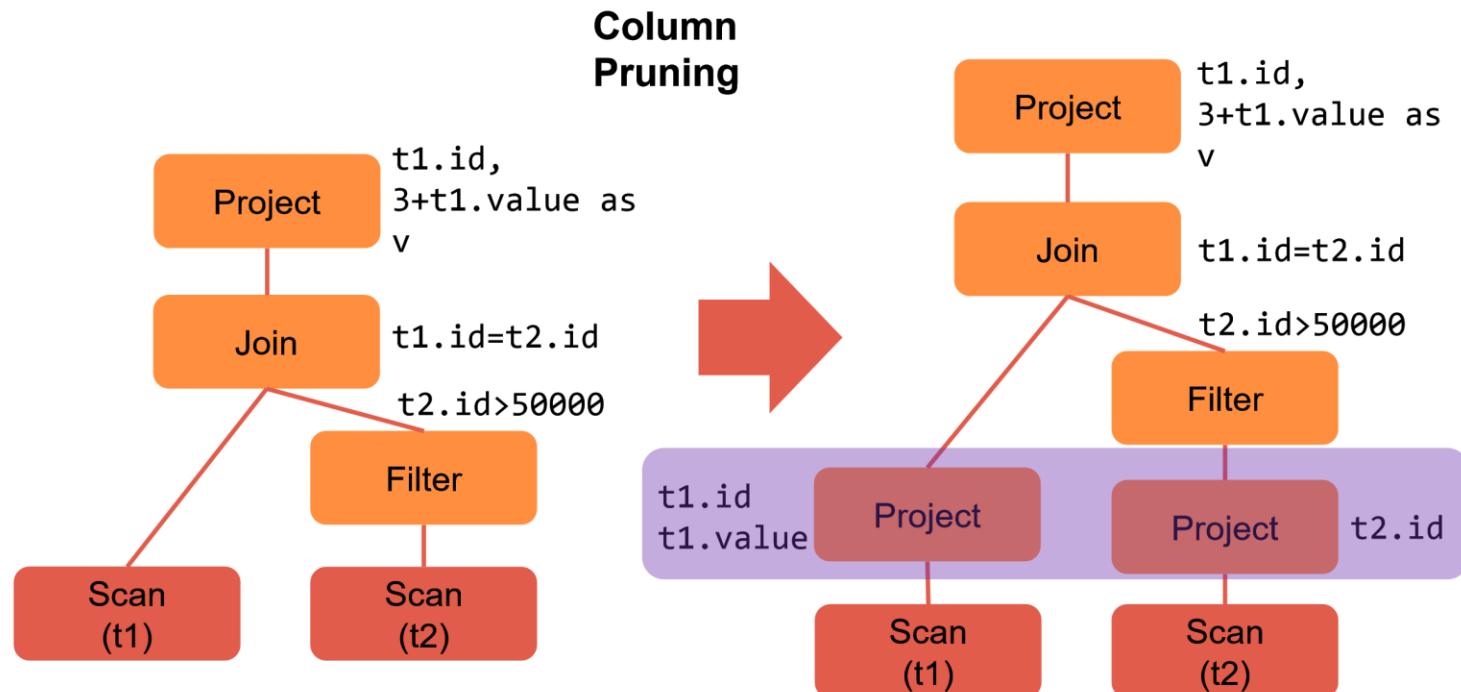
Predicate Pushdown

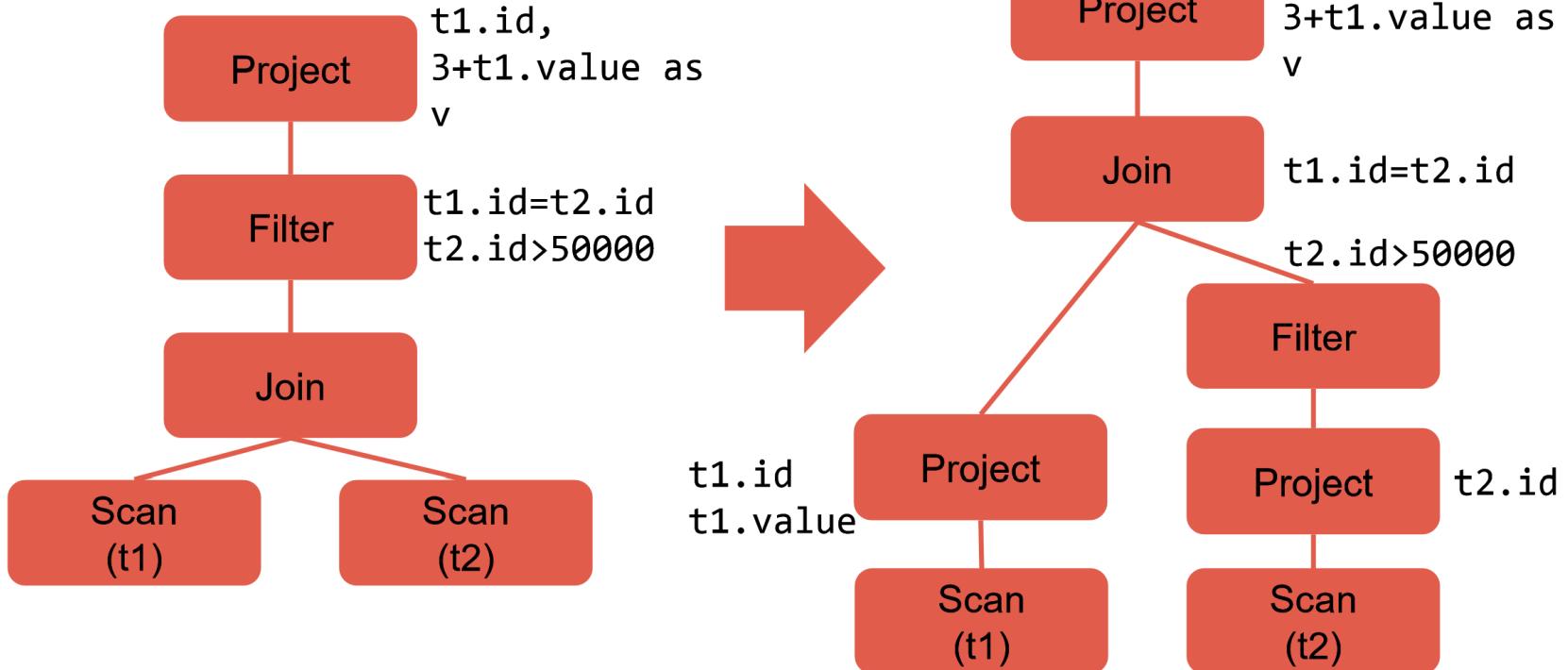


Catalyst – Комбинирование нескольких правил



Catalyst – Комбинирование нескольких правил





Catalyst: Logical Plan => Physical Plan

```
object BasicOperators extends Strategy {  
    def apply(plan: LogicalPlan): Seq[SparkPlan] = plan match {  
        ...  
        case logical.Project(projectList, child) =>  
            execution.ProjectExec(projectList, planLater(child)) :: Nil  
        case logical.Filter(condition, child) =>  
            execution.FilterExec(condition, planLater(child)) :: Nil  
        ...  
    }  
}
```

Catalyst: Logical Plan => Physical Plan

```
object BasicOperators extends Strategy {  
    def apply(plan: LogicalPlan): Seq[SparkPlan] = plan match {  
        ...  
        case logical.Project(projectList, child) =>  
            execution.ProjectExec(projectList, planLater(child)) :: Nil  
        case logical.Filter(condition, child) =>  
            execution.FilterExec(condition, planLater(child)) :: Nil  
        ...  
    }  
}
```

Catalyst: Logical Plan => Physical Plan

- Анализ(Rule Executor): трансформирует неразрешенный логический план в разрешенный логический план
- Логическая оптимизация(Rule Executor): Трасформирует разрешенный логический план в оптимизированный логический план
- Физическое планирование(Стратегии + Rule Executor):
 1. Трансформация оптимизированного логический план в физический план
 2. Применение rule executor и подготовка плана к исполнению

Спасибо

Dmitry Lakhvich

-  LinkedIn: <https://www.linkedin.com/in/dmitry-lakhvich-a610706b/>
-  Twitter: [@KrivdaTheTriewe](https://twitter.com/KrivdaTheTriewe)
-  Telegram: [@KrivdaTheTriewe](https://t.me/KrivdaTheTriewe)
- Email: frostball@gmail.com