# Spark SQL

Шугаепов Ильнур
VK.com
ilnur.shug@gmail.com

февраль 2020 г.

## История

2006 MapReduce

2009 Hive

2009 Pig

2010 Spark

2013 Shark[1]

2015 Impala[2]

2015 Spark SQL[3]

---

[1] Reynold S Xin и др. "Shark: SQL and rich analytics at scale". B: *Proceedings of the 2013 ACM SIGMOD International Conference on Management of data*. 2013, с. 13—24.

[2] Marcel Kornacker и др. "Impala: A Modern, Open-Source SQL Engine for Hadoop.". B: *Cidr*. Т. 1. 2015, с. 9.

[3] Michael Armbrust и др. "Spark sql: Relational data processing in spark". B: *Proceedings of the 2015 ACM SIGMOD international conference on management of data*. 2015, с. 1383—1394.

Spark
Limitations

- Low-level procedural code
- No optimizations

## Shark

- First effort to build a relational interface on Spark
- Shark modified the Apache Hive system to run on Spark

## Shark
### Limitations

- Shark could only be used to query external data stored in the Hive catalog
- The only way to call Shark from Spark programs was to put together a SQL string
- Hive optimizer was tailored for MapReduce and difficult to extend

### Наблюдение

*Most data pipelines are combination of relational and procedural algorithms.*

Background and Goals
Programming Interface
Catalyst
Evaluation
Список литературы

Ограничения существующих систем
Наблюдение
Предложенное решение
Spark overview
Goals

Spark SQL — new module in Apache Spark

- DataFrame API
- Catalyst (optimizer)

### Определение

*DataFrames are collections of structured records that can be manipulated using Spark's underline{procedural} API, or using new underline{relational} APIs that allow richer optimizations*

They can be created directly from Spark's RDDs, enabling relational processing in existing Spark programs.

## Spark
### RDDs

### Определение

*RDD:*

- *Resilient — отказоустойчивый*
- *Distributed — разбитый на партиции*
- *Dataset*

*read-only, partitioned collection of records*

RDDs can be manipulated through operations like `map`, `filter`, and reduce, which take functions and ship them to nodes on the cluster.

## Spark
### Fault-tolerance

- Запомним граф вычислений (linage)
- Тогда если часть данных будет потеряна, то их легко можно восстановить

Spark
Lazy evaluation

- Each RDD represents a "logical plan" to compute a dataset
- Spark waits until <u>action</u> to launch a computation
- Allows to do some simple query optimization, such as pipelining operations (narrow dependencies)
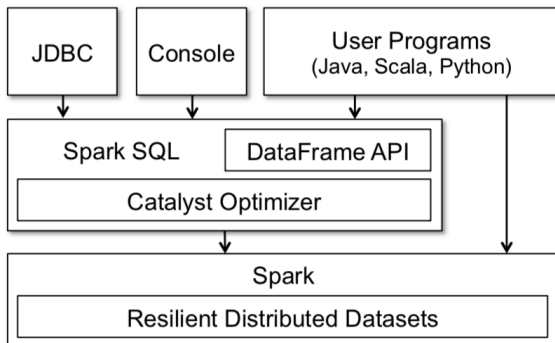
## Spark
Optimizations

Spark engine does not understand the structure of the data in RDDs (which is arbitrary Java/Python objects) or the semantics of user functions (which contain arbitrary code)

Background and Goals
Programming Interface
Catalyst
Evaluation
Список литературы

Ограничения существующих систем
Наблюдение
Предложенное решение
Spark overview
Goals

Цели

1. Support relational processing both within Spark programs (on native RDDs) and on external data sources using a programmer-friendly API.

2. Easily support new data sources, including semi-structured data and external databases amenable to query federation.

3. Enable extension with advanced analytics algorithms such as graph processing and machine learning.

Background and Goals
Programming Interface
Catalyst
Evaluation
Список литературы

DataFrame API
DataFrame Operations
In Memory Caching
User Defined Functions

PROGRAMMING INTERFACE

Background and Goals    DataFrame API
**Programming Interface**    DataFrame Operations
Catalyst    In Memory Caching
Evaluation    User Defined Functions
Список литературы

Background and Goals
Programming Interface
Catalyst
Evaluation
Список литературы

DataFrame API
DataFrame Operations
In Memory Caching
User Defined Functions

- DataFrame — distributed collection of rows with the same *schema*
- DataFrame is equivalent to a table in a relational database
- DataFrames can be manipulated in similar ways to the RDDs
- Schema leads to a more optimized execution

Background and Goals
Programming Interface
Catalyst
Evaluation
Список литературы

DataFrame API
DataFrame Operations
In Memory Caching
User Defined Functions

DataFrame
Construction

- From external data sources (HDFS, Hive)
- From existing RDDs (schema inference algorithm)

### Замечание

*DataFrame can be viewd as an RDD of Row objects, allowing user to call procedural Spark APIs such as map*

Background and Goals
Programming Interface
Catalyst
Evaluation
Список литературы

DataFrame API
DataFrame Operations
In Memory Caching
User Defined Functions

## DataFrame
Execution

- Spark DataFrame are *lazy*
- Spark build *logical plan* before execution
- Laziness enables rich optimization
- Logical plan  $\rightarrow$  Physical plan

DSL

Users can perform relational operations on DataFrames using a
domain-specific language (DSL) similar to Python Pandas

Background and Goals
Programming Interface
Catalyst
Evaluation
Список литературы

DataFrame API
DataFrame Operations
In Memory Caching
User Defined Functions

## Example
Фильмы с наибольшим средним рейтингом

```
1  ratings_df \
2      .groupby('movie_id') \
3      .agg(F.mean('rating').alias('mean_rating'),
4          F.count('rating').alias('ratings_count')) \
5      .join(movies_df, ratings_df['movie_id'] == movies_df['movieId'],
6          how='inner') \
7      .sort(F.col('mean_rating').desc()) \
```

## Difference with native Spark API

- All of these operators build up an abstract syntax tree (AST) of the expression, which is then passed to Catalyst for optimization.
- This is unlike the native Spark API that takes functions containing arbitrary Scala/Java/Python code, which are then opaque to the runtime engine.

Background and Goals
**Programming Interface**
Catalyst
Evaluation
Список литературы

DataFrame API
**DataFrame Operations**
In Memory Caching
User Defined Functions

# DataFrame construction
## Schema inference

- While building DataFrame from RDD user can manually define schema
- Spark SQL can automatically infer the schema of the dataset using reflection
- In Python, Spark SQL samples the dataset to perform schema inference due to the dynamic type system

Background and Goals
Programming Interface
Catalyst
Evaluation
Список литературы

DataFrame API
DataFrame Operations
In Memory Caching
User Defined Functions

.cache()

- Method cache of DataFrame does the same thing as method persist of RDD
- Caching is particularly useful for interactive queries and for the iterative algorithms common in machine learning

UDF
Example

```
1 def get_release_year(title):
2     result = re.match(r'.*(\(\d+\))', title)
3     return int(result.group(1)[1:-1]) if result is not None else None
4
5 get_release_year_udf = F.udf(get_release_year, IntegerType())
6
7 movies_df \
8     .withColumn('year', get_release_year_udf('title')) \
```

Background and Goals
Programming Interface
Catalyst
Evaluation
Список литературы

DataFrame API
DataFrame Operations
In Memory Caching
User Defined Functions

Резюме

- The DataFrame API lets developers seamlessly mix procedural and relational methods.

Background and Goals
Programming Interface
Catalyst
Evaluation
Список литературы

Trees & Rules
Catalyst in Spark SQL

CATALYST

Background and Goals
Programming Interface
Catalyst
Evaluation
Список литературы

Trees & Rules
Catalyst in Spark SQL

Catalyst contains a general library for representing <u>trees</u> and applying <u>rules</u> to manipulate them

Background and Goals
Programming Interface
Catalyst
Evaluation
Список литературы

Trees & Rules
Catalyst in Spark SQL

## Trees



Рис.: Catalyst tree for expression x+(1+2)

Background and Goals
Programming Interface
Catalyst
Evaluation
Список литературы

Trees & Rules
Catalyst in Spark SQL

## Rules

### Определение

$Rule: T \mapsto T'$ — rule maps tree to another tree.

The most common approach is to use a set of pattern matching functions that find and replace subtrees with a specific structure.

Background and Goals
Programming Interface
Catalyst
Evaluation
Список литературы

Trees & Rules
Catalyst in Spark SQL

# Rules
## Example

### Constant folding

```
1 tree.transform {
2     case Add(Literal(c1), Literal(c2)) => Literal(c1+c2)
3     case Add(left, Literal(0)) => left
4     case Add(Literal(0), right) => right
5 }
```

Background and Goals
Programming Interface
Catalyst
Evaluation
Список литературы

Trees & Rules
Catalyst in Spark SQL

# Rules
Fixed point

Catalyst groups rules into batches, and executes each batch until it reaches a *fixed point*, that is, until the tree stops changing after applying its rules.

Background and Goals
Programming Interface
**Catalyst**
Evaluation
Список литературы

Trees & Rules
Catalyst in Spark SQL

Рис.: Phases of query planning in Spark SQL. Rounded rectangles represent Catalyst trees

Background and Goals
Programming Interface
**Catalyst**
Evaluation
Список литературы

Trees & Rules
Catalyst in Spark SQL

## Example
### Query

```
1 query = """
2     SELECT movie_id, COUNT(*), first(title) as title
3     FROM ratings INNER JOIN movies ON ratings.movie_id == movies.movieId
4     WHERE movies.title LIKE '%(1994)%'
5     GROUP BY movie_id
6     ORDER BY COUNT(*) DESC
7 """
8
9 spark.sql(query).explain(True)
```

Background and Goals
Programming Interface
Catalyst
Evaluation
Список литературы

Trees & Rules
Catalyst in Spark SQL

## Example
Unresolved Logical Plan

```
1 == Parsed Logical Plan ==
2 'Sort ['COUNT(1) DESC NULLS LAST], true
3 +— 'Aggregate ['movie_id], ['movie_id, unresolvedalias('COUNT(1), None),
      first('title, false) AS title#53]
4   +— 'Filter 'movies.title LIKE %(1994)%
5     +— 'Join Inner, ('ratings.movie_id = 'movies.movieId)
6       :— 'UnresolvedRelation 'ratings'
7       +— 'UnresolvedRelation 'movies'
```

Background and Goals
Programming Interface
Catalyst
Evaluation
Список литературы

Trees & Rules
Catalyst in Spark SQL

## Example
Logical Plan

```
 1  == Analyzed Logical Plan ==
 2  movie_id: int, count(1): bigint, title: string
 3  Project [movie_id#1, count(1)#56L, title#53]
 4  +- Sort [count(1)#56L DESC NULLS LAST], true
 5     +- Aggregate [movie_id#1], [movie_id#1, count(1) AS count(1)#56L,
            first(title#19, false) AS title#53]
 6        +- Filter title#19 LIKE %(1994)%
 7           +- Join Inner, (movie_id#1 = movieId#18)
 8              :- SubqueryAlias 'ratings'
 9              : +- Relation[user_id#0,movie_id#1,rating#2,timestamp#3] csv
10              +- SubqueryAlias 'movies'
11                 +- Relation[movieId#18,title#19,genres#20] csv
```

Background and Goals
Programming Interface
**Catalyst**
Evaluation
Список литературы

Trees & Rules
Catalyst in Spark SQL

## Example
Optimized Logical Plan

```
1  == Optimized Logical Plan ==
2  Sort [count(1)#56L DESC NULLS LAST], true
3  +- Aggregate [movie_id#1], [movie_id#1, count(1) AS count(1)#56L, first(
      title#19, false) AS title#53]
4     +- Project [movie_id#1, title#19]
5        +- Join Inner, (movie_id#1 = movieId#18)
6           :- Project [movie_id#1]
7           :  +- Filter isnotnull(movie_id#1)
8           :     +- Relation[user_id#0,movie_id#1,rating#2,timestamp#3] csv
9           +- Project [movieId#18, title#19]
10             +- Filter ((isnotnull(title#19) && Contains(title#19, (1994))
                  ) && isnotnull(movieId#18))
11                +- Relation[movieId#18,title#19,genres#20] csv
```
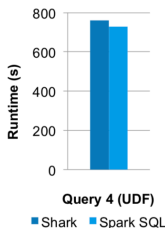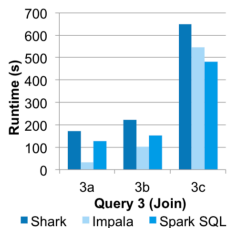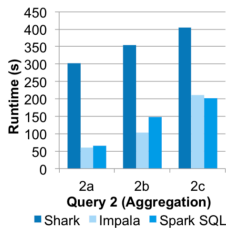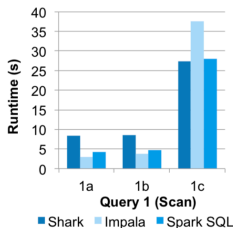
! Filter push down rule !

Background and Goals
Programming Interface
Catalyst
**Evaluation**
Список литературы

SQL Performance
DataFrames vs. Native Spark

EVALUATION

Background and Goals
Programming Interface
Catalyst
**Evaluation**
Список литературы

SQL Performance
DataFrames vs. Native Spark

Evaluation of the performance of Spark SQL on two dimensions:

1. SQL query processing performance
2. Spark program performance

Background and Goals
Programming Interface
Catalyst
Evaluation
Список литературы

SQL Performance
DataFrames vs. Native Spark

# Benchmark[4]

[4] Andrew Pavlo. "A comparison of approaches to large-scale data analysis". В: *Proceedings of the 2009 ACM SIGMOD international conference on management of data*. 2009, с. 165—178.

Background and Goals
Programming Interface
Catalyst
Evaluation
Список литературы

SQL Performance
DataFrames vs. Native Spark

## Distributed Aggregation
Dataset and Task

Dataset 1 billion integer pairs, $(a, b)$ with $100000$ distinct values of $a$

Task compute the average of $b$ for each value of $a$

Background and Goals
Programming Interface
Catalyst
**Evaluation**
Список литературы

SQL Performance
DataFrames vs. Native Spark

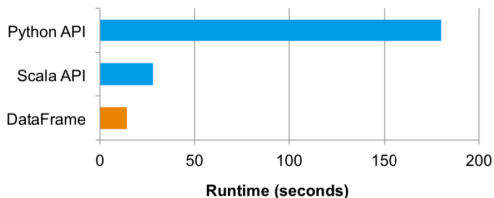# Distributed Aggregation
Solutions

### Native Spark

```
1 sum_and_count = data \
2     .map(lambda x: (x.a, (x.b, 1))) \
3     .reduceByKey(lambda x, y: (x[0]+y[0], x[1]+y[1])) \
4     .collect()
5
6 [(x[0], x[1][0] / x[1][1]) for x in sum_and_count]
```

### Spark SQL

```
1 df.groupBy("a").avg("b")
```

Background and Goals
Programming Interface
Catalyst
Evaluation
Список литературы

SQL Performance
DataFrames vs. Native Spark

## Distributed Aggregation
Performance



In the DataFrame API, only the <u>logical plan</u> is constructed in Python, all <u>physical execution</u> is compiled into native Spark code as JVM bytecode.

📄 Armbrust, Michael и др. "Spark sql: Relational data processing in spark". В: *Proceedings of the 2015 ACM SIGMOD international conference on management of data*. 2015, с. 1383—1394.

📄 Karau, Holden и Rachel Warren. *High performance Spark: best practices for scaling and optimizing Apache Spark*. " O'Reilly Media, Inc.", 2017.

📄 Kornacker, Marcel и др. "Impala: A Modern, Open-Source SQL Engine for Hadoop.". В: *Cidr*. Т. 1. 2015, с. 9.

📄 Pavlo, Andrew. "A comparison of approaches to large-scale data analysis". В: *Proceedings of the 2009 ACM SIGMOD international conference on management of data*. 2009, с. 165—178.

📄 Xin, Reynold S и др. "Shark: SQL and rich analytics at scale". В: *Proceedings of the 2013 ACM SIGMOD International Conference on Management of data*. 2013, с. 13—24.