

# Spark Dataframe API and Spark SQL

*Spark*  *SQL*



pm jat @ daiict



# Further Higher Abstractions

- RDD are still “low level” to perform analytical tasks!
- Spark provides higher level abstractions
  - “Dataframe” API (more “**structured RDDs**”)
  - Spark-SQL (SQL interface)
- These abstractions makes “cluster programming” amazingly simple, and
- Primarily the reason Spark is becoming popular for big data processing.



# SQL Interface over huge raw files!

- We require SQL interface over raw data file, so that, we are able to run queries without loading data into some database systems
  - For many ETL operations, load time into database systems is forbiddingly high
  - Schema check happens only when we run the query (Read time schema validation)
- This makes running ad-hoc queries on data files quick
- This has been motivation for Pig, Hive, and now Spark-SQL



# Before Spark SQL

- Hive from Facebook
- Pig from Yahoo
  - Not really SQL but a scripting language like perl
- Cloudera Impala
- Google Dremel



# Hive and Pig

- **Hive:** SQL like interface for HDFS files. Initially developed at Facebook (now Apache project).

facebook



```
SELECT count(*) FROM users
```

In reality, 90+% of MR jobs are generated by Hive SQL

- **Pig:** scripting language for various data transformations. Initially developed at Yahoo. Now again apache project

YAHOO!



```
A = load 'foo';  
B = group A all;  
C = foreach B generate COUNT(A);
```

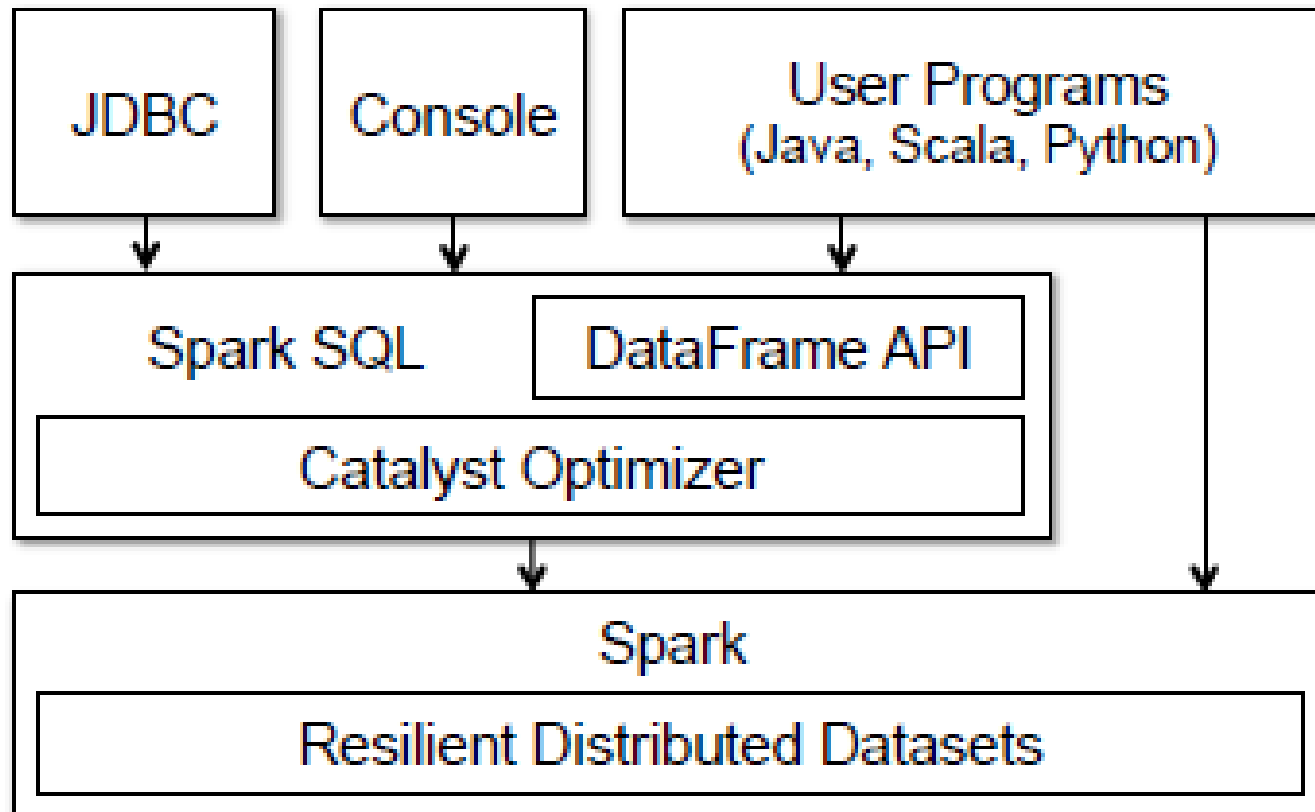


# Spark-SQL

- Spark SQL was introduced in 2015 through paper “[Spark SQL: Relational data processing in spark.](#)” in ACM SIGMOD
- Initial effort to have SQL interface over Spark was **Shark**, and was basically an adaption of “Hive over Spark”.
- Shark, however
  - could not integrate well will with intermediate RDD datasets, and
  - Secondly Hive optimizer could not properly adapted to Spark as it was primarily designed for Map Reduce.
- And, we have [Spark-SQL](#)[1, 2015]



# Spark-SQL – Stack<sup>[1]</sup>



[1] Armbrust, Michael, et al. "Spark sql: Relational data processing in spark." *Proceedings of the 2015 ACM SIGMOD international conference on management of data*. ACM, 2015



# Spark-SQL - key characteristics

- Spark-SQL has following key characteristics-
  - Performance (since underlying engine is Spark)
  - Ability to access “procedures” from SQL statements
    - Opens up gate (for SQL) to interact with ML, Graph Processing algorithms as “User Defined Functions”
  - Query Optimizer
  - `Select regmodel.predictsalary(resume) from employee;`





# Spark SQL - Programming Interface

- Spark SQL runs as a library on top of Spark (refer stack diagram)
- It has become buzzword for “**Declarative Big Data Processing**”
  - Write less code
  - Optimize the execution
- Spark exposes SQL interfaces, which can be accessed through
  - Command Line Interface (Console)
  - DataFrame API integrated into Spark programs
  - JDBC/ODBC access (through Spark Thrift Server)



# Spark “Data Frame API” Source: [1]

- The main abstraction in Spark SQL’s API is a **DataFrame**, a “distributed collection of rows with a **homogeneous schema**”.
- A DataFrame is equivalent to a table in a relational database, and can also be manipulated in similar ways.
- Unlike RDDs, DataFrames keep track of their schema and support various relational operations on them.
- DataFrames can be constructed from data files, tables in a system catalog (based on external data sources) or from existing RDD objects!



# Transformations and Actions on DataFrame

Source: [2]

## Transformation examples

- filter
- select
- drop
- intersect
- join

## Action examples

- count
- collect
- show
- head
- take

- *Note: Dataframe evaluations are lazy. Transformations are just getting expressed (not executed immediately). It is the request of actions that lead to real execution (and this is supported by optimization)*



# DataFrame transformation are Lazy

Source: [1]

- Unlike traditional (like in Panda and R) data frame APIs, Spark Data Frames are lazy,
  - in spark, each DataFrame object represents a logical plan to compute a dataset frame
- It is like, when specified, it adds to execution path, and executed only when some action is to be performed.
- This enables optimization across all operations that were used to build the Data Frame.



# Data Frames and RDDs

- DataFrames are built on top of the RDD and are “structured”, that is have schema attached with them.
- We can call them Schema aware RDDs - element type is “Row with Schema”
- Also immutable
- DataFrames are more structured, provides abstraction similar to a database table. This makes it very simple to program.
- Dataframes are promised to be more efficient, because operations on them **can be optimized through “Catalyst”, the optimizer for Spark-SQL!**



# Dataframe API - summarize

- An Higher level abstraction on Spark
- Dataframes can be think of as “RDD with Schema”
- Abstraction attempts to make working with RDDs like a Relational table.
- Dataframes can be constructed from “data files”, tables in a system catalog (based on external data sources), or from existing RDD objects!
- Dataframes are promised to be more efficient, because operations on them **can be optimized through “Catalyst”, the optimizer for Spark-SQL!**



# Programming Dataframes



# Spark Session

- Prior to Spark 2.0, Spark had three Contexts
  - Spark Context: for RDD operations
  - SQL Context: for Spark SQL operations
  - Hive Context: for Hive queries
- However Spark 2.0 onwards, there is a unified context called Spark Session, all operations can be performed through spark session, therefore
- Spark Session is more recent!





# Initialize Spark Session

- Before doing anything on Spark Dataframe API or Spark SQL we require creating Spark Session
- Typically done as following

```
SparkSession spark = SparkSession.builder()  
    .appName("Spark Demo")  
    .master("local")  
    .getOrCreate();
```



# Reading from a CSV file as Dataframe

- Typically done as following:

```
Dataset<Row> emp = spark.read()  
    .option("header", "true")  
    .option("sep", ",")  
    .csv("data/employee_m.csv");
```

- Peek into read data rows Can from the CSV file.

```
//show all rows
```

```
emp.show(); //select * from emp;
```

```
//show first 3 rows
```

```
emp.show(3); //select * from emp limit 3;
```



# Reading from a CSV file as Dataframe

- By default
  - If No Header Row in data, then columns are named as \_c0, \_c1, \_c2, and so on
  - To indicate that file has header row, we say:  
`.option("header", "true")`
  - Schema remains undefined. By default data type for each column remains string.
- We can ask spark to automatically infer the schema, or we can manually supply the schema before reading a CSV file.



# Inferring Schema

- Set : `.option("inferSchema", "true");` to auto infer the schema (in read options)
- We can dump the inferred schema as following:  
`emp.printSchema();`

```
root
|-- eno: integer (nullable = true)
|-- name: string (nullable = true)
|-- dob: timestamp (nullable = true)
|-- gender: string (nullable = true)
|-- salary: integer (nullable = true)
|-- sup_eno: integer (nullable = true)
|-- dno: integer (nullable = true)
```

- For inferring the schema, processor requires scanning full file. This may be expensive when file is huge, and may not always be successful.



# Manually Specifying the Schema

## (1) Define Schema:

```
StructType schema = DataTypes.createStructType(new StructField[] {  
    //eno,name,dob,gender,salary,sup_eno,dno  
    DataTypes.createStructField("eno", DataTypes.StringType, false),  
    DataTypes.createStructField("name", DataTypes.StringType, true),  
    DataTypes.createStructField("dob", DataTypes.DateType, true),  
    DataTypes.createStructField("gender", DataTypes.StringType, true),  
    DataTypes.createStructField("salary", DataTypes.IntegerType, true),  
    DataTypes.createStructField("sup_eno", DataTypes.StringType, true),  
    DataTypes.createStructField("dno", DataTypes.IntegerType, true)  
});  
//System.out.println(schema.prettyJson());
```

## (2) Specify Schema:

```
Dataset<Row> emp = spark.read()  
    .option("header", "true")  
    .option("sep", ",")  
    .schema(schema)  
    .csv("data/employee_m.csv");
```



# Reading from a JSON file as Dataframe

- Data frame object is constructed and used as following:

```
Dataset<Row> emp1 = spark.read()  
    .json("data/employees.json");  
emp1.printSchema();  
emp1.show();
```

name	salary
Michael	3000
Andy	4500
Justin	3500
Berta	4000

```
root  
 |-- name: string (nullable = true)  
 |-- salary: long (nullable = true)
```



# Examples – Filter/Selection

```
emp.filter( "salary > 30000" ).show();  
// selection * from employee where salary > 30000
```

```
emp.where( "salary > 40000" ).show();  
// selection * from employee where salary > 40000
```

```
emp.filter( emp.col("salary").geq(40000)  
.and(emp.col("dno").equalTo(5))).show(); //More TYPE SAFE  
// selection * from employee where salary >= 40000  
//                                     dno=5;
```



# Examples – Sort/Order By

```
emp.sort(emp.col("dno"),emp.col("salary").desc())  
.show();
```

```
// selection * from employee order by dno, salary desc
```

```
emp.orderBy(emp.col("dno"),emp.col("salary"))  
.show();
```

```
// selection * from employee order by dno, salary
```





# Examples - Project

```
emp.select("eno", "name", "salary")  
  .where("dno == 4")  
  .show();  
//SELECT eno, name, salary from emp where dno=4
```

```
emp.select(  
  emp.col("eno"),  
  emp.col("name"),  
  emp.col("salary").multiply(1.1)  
) .where("dno == 4")  
  .show();  
//SELECT eno, name, salary*1.1 from emp where dno=4
```



# Examples - Aggregation

//on whole Table

```
emp.agg(avg(emp.col("salary")),max(emp.col("salary")))  
.show();
```

```
//select avg(salary), max(salary) from employee;
```

//Group By and Aggregation

```
emp.groupBy(emp.col("dno"))  
.agg(avg(emp.col("salary")), max(emp.col("salary")))  
.show();
```

```
//select avg(salary), max(salary) from employee  
//      group by dno;
```



# Examples – Join and Project

```
Dataset<Row> empdep =  
emp.join(dep, emp.col("dno").equalTo(dep.col("dno")));  
  
empdep.select(emp.col("name"), dep.col("name"),  
              emp.col("salary"))  
.show();  
  
//select e.name, d.name from employee e inner join  
//      department d on e.dno=d.dno;
```



# Examples – Join and Aggregate

```
emp.join(dep, emp.col("dno").equalTo(dep.col("dno")))
    .groupBy(dep.col("name"))
    .agg(sum(emp.col("salary")))
    .show();
```

```
//select d.name, sum(e.salary) from employee e inner join
//      department d on e.dno=d.dno
//      group by d.name;
```



# [Python] Constructing Data Frame

Source: [2]

```
# Construct a DataFrame from a "users" table in Hive.  
df = sqlContext.table("users")  
  
# Construct a DataFrame from a log file in S3.  
df = sqlContext.load("s3n://someBucket/path/to/data.json", "json")
```





# [Python] Operation on Data Frames

Source: [2]

```
# Create a new DataFrame that contains only "young" users
young = users.filter(users["age"] < 21)

# Alternatively, using a Pandas-like syntax
young = users[users.age < 21]

# Increment everybody's age by 1
young.select(young["name"], young["age"] + 1)

# Count the number of young users by gender
young.groupBy("gender").count()

# Join young users with another DataFrame, logs
young.join(log, logs["userId"] == users["userId"], "left_outer")
```





# [Python] Example: DataFrame API

- Let us say a tab separated data file called “**SalesProduct.txt**”, where attributes Name, and Weight at 2<sup>nd</sup> and 8<sup>th</sup> position respectively.
- Following Spark program (using dataframe API) lists Name and Weight of top 15 products in the descending order weight!

```
sqlContext = SQLContext(sc)
```

```
rdd1 = (content.filter(lambda line: line.split("\t")[7] != "NULL")  
        .map(lambda line: (line.split("\t")[1], float(line.split("\t")[7])))  
        )
```

```
df = sqlContext.createDataFrame(rdd1, schema = ["Name", "Weight"])
```

```
df.orderBy("weight", ascending = False).show(15, truncate = False)
```



# Spark SQL





# Spark SQL

- Spark SQL page (<https://spark.apache.org/sql/>) describes spark-sql with following characteristics:

```
results = spark.sql(  
    "SELECT * FROM people")  
names = results.map(lambda p: p.name)
```

## Integrated:

- Seamlessly mix SQL queries with Spark programs.
- Spark SQL lets you query structured data inside Spark programs, using either SQL or a familiar DataFrame API.
- Usable in Java, Scala, Python and R.



# Spark SQL

## Uniform Data Access

- Connect to any data source the same way.
- DataFrames and SQL provide a common way to access a variety of data sources, including Hive, Avro, Parquet, ORC, JSON, and JDBC. You can even join data across these sources.

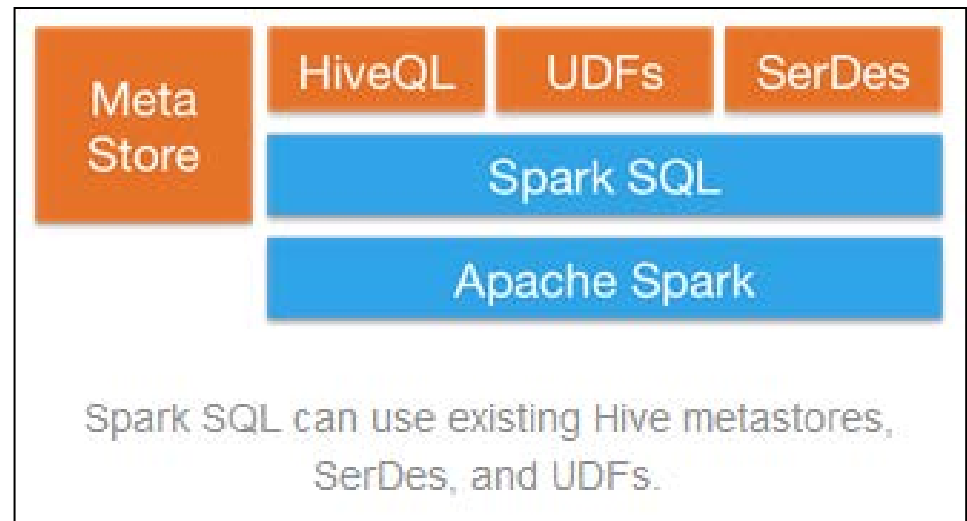
```
spark.read.json("s3n://...")  
  .registerTempTable("json")  
results = spark.sql(  
  """SELECT *  
    FROM people  
    JOIN json ...""")
```



# Spark SQL

## Hive Integration

- Run SQL or HiveQL queries on existing warehouses.
- Spark SQL supports the HiveQL syntax as well as Hive SerDes (Serializer/Deserializer) and User Defined Functions (UDFs), allowing you to access existing Hive warehouses.

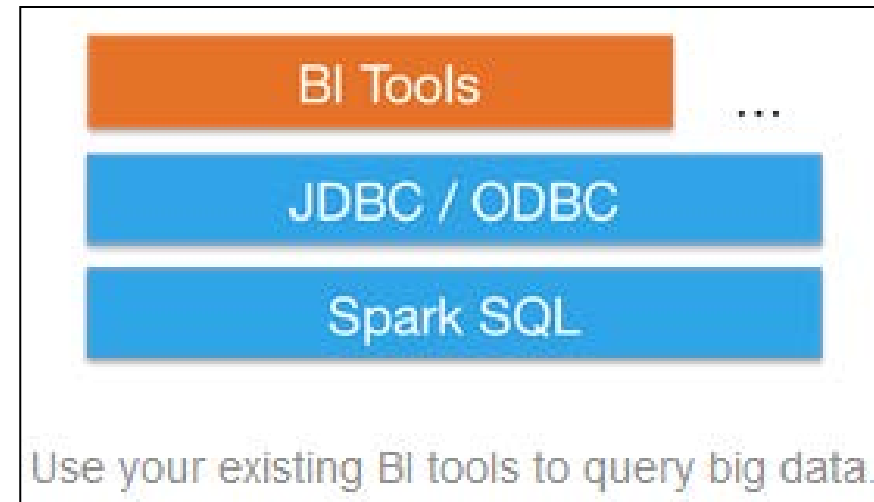




# Spark SQL

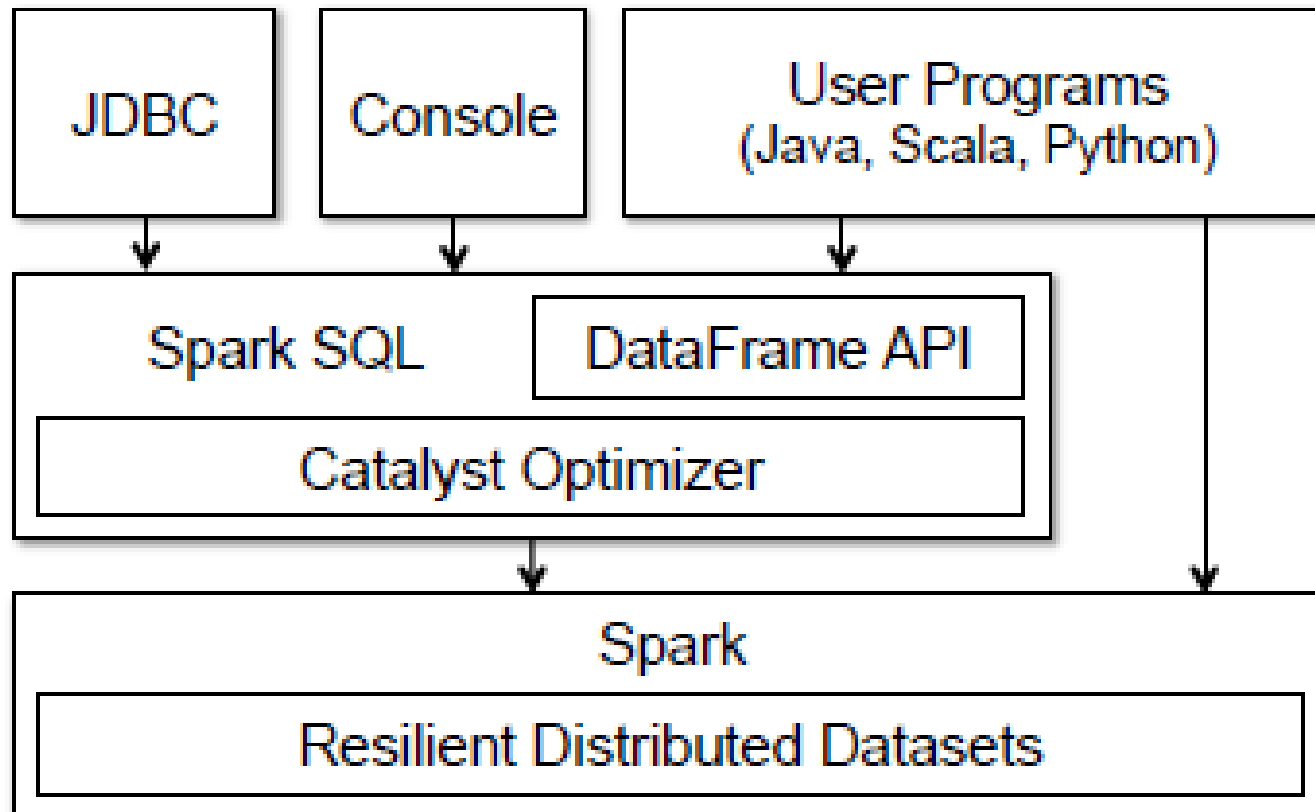
## Standard Connectivity

- Connect through JDBC or ODBC.
- A server mode provides industry standard JDBC and ODBC connectivity for business intelligence tools.





# Spark-SQL - Stack



[1] Armbrust, Michael, et al. "Spark sql: Relational data processing in spark." *Proceedings of the 2015 ACM SIGMOD international conference on management of data*. ACM, 2015



# Spark SQL

- Concepts are very similar to relational “Table” and “SQL”

## Tables

- In programmer’s perspective table is identical to relational table.
- Here tables can be temporary or be stored. Stored tables are typically managed by Spark based data repository.
- We typically have one data file for one table.
- Spark data repository may store tables in “storage formats” that helps in efficient execution of queries.



# Spark SQL

- Spark SQL provides same set of commands **DDL + DML**
- However commands do have some additional parameters that are specific to distributed data and clustered computing
  - For instance partitioning, shuffling, and sorting strategy, etc.
- Spark SQL run-time system provides “Spark SQL Engine” that “efficiently” executes Spark-SQL operations.
- SQL statements can be executed from CLI (command level interface) or host programs.
- CLI allows executing SQL commands interactively.



# Spark SQL statements

- CREATE TABLE.  
Also ALTER TABLE, DROP TABLE, etc
- INSERT, UPDATE, LOAD data,
- CREATE, DROP FUNCTION
- Various commands to add, merge, delete tables in “data lakes”. Data lake is basically large scale data repository. One can be spark based data lake with SQL interface.
- Here is SQL reference manual from Databrick  
<https://docs.databricks.com/spark/latest/spark-sql/index.html>





# CREATE TABLE

```
CREATE [TEMPORARY] TABLE [IF NOT EXISTS] [db_name.]table_name
  [(col_name1 col_type1 [COMMENT col_comment1], ...)]
  USING datasource
  [OPTIONS (key1=val1, key2=val2, ...)]
  [PARTITIONED BY (col_name1, col_name2, ...)]
  [CLUSTERED BY (col_name3, col_name4, ...) INTO num_buckets BUCKETS]
  [LOCATION path]
  [COMMENT table_comment]
  [TBLPROPERTIES (key1=val1, key2=val2, ...)]
  [AS select_statement]
```



# CREATE TABLE “USING datasource”

- Specifies file format to use for the table.
- “**data\_source**” must be one of TEXT, CSV, JSON, JDBC, PARQUET, ORC, HIVE, DELTA, or LIBSVM, or a user defined class
- HIVE is supported to create a Hive SerDe table.



# CREATE TABLE options

**PARTITIONED BY (col\_name1, col\_name2, ...)**

- Partition the created table by the specified columns. A directory is created for each partition.

**CLUSTERED BY col\_name3, col\_name4, ...)**

- Each partition in the created table will be split into a fixed number of buckets by the specified columns.
- This is typically used with partitioning to read and shuffle less data.

**LOCATION path**

- The directory to store the table data. This clause automatically implies EXTERNAL.



# CREATE TABLE options

## AS select\_statement

- Populate the table with input data from the SELECT statement.

## Examples

```
CREATE TABLE boxes (width INT, length INT, height INT) USING CSV
```

```
CREATE TABLE rectangles  
  USING PARQUET  
  PARTITIONED BY (width)  
  CLUSTERED BY (length) INTO 8 buckets  
  AS SELECT * FROM boxes
```



# SELECT statement

```
SELECT [hints, ...] [ALL|DISTINCT] named_expression[, named_expression, ...]  
  FROM relation[, relation, ...]  
  [lateral_view[, lateral_view, ...]]  
  [WHERE boolean_expression]  
  [aggregation [HAVING boolean_expression]]  
  [ORDER BY sort_expressions]  
  [CLUSTER BY expressions]  
  [DISTRIBUTE BY expressions]  
  [SORT BY sort_expressions]  
  [WINDOW named_window[, WINDOW named_window, ...]]  
  [LIMIT num_rows]
```



# SELECT options

## ALL

- Select all matching rows from the relation. Enabled by default.

## DISTRIBUTE BY

- Repartition rows in the table based on a set of expressions.
- Rows with the same expression values will be hashed to the same worker.
- You cannot use this with ORDER BY or CLUSTER BY.



# SELECT options

## SORT BY

- Impose ordering on a set of expressions “within each partition”. Default sort direction is ascending.
- You cannot use this with ORDER BY or CLUSTER BY.

## CLUSTER BY

- Repartition rows in the relation based on a set of expressions and sort the rows in ascending order based on the expressions.
- In other words, this is a shorthand for DISTRIBUTE BY and SORT BY
- You cannot use this with ORDER BY, DISTRIBUTE BY, or SORT BY.



# SELECT - examples

```
SELECT * FROM boxes
```

```
SELECT width, length FROM boxes WHERE height=3
```

```
SELECT DISTINCT width, length FROM boxes WHERE height=3 LIMIT 2
```

```
SELECT * FROM VALUES (1, 2, 3) AS (width, length, height)
```

```
SELECT * FROM VALUES (1, 2, 3), (2, 3, 4) AS (width, length, height)
```

```
SELECT * FROM boxes ORDER BY width
```

```
SELECT * FROM boxes DISTRIBUTE BY width SORT BY width
```

```
SELECT * FROM boxes CLUSTER BY length
```





# SELECT – DW (Data Cube options)

- Common data cube (a common task in DW applications) options are also available:
  - ROLLUP
  - CUBE
  - GROUPING SETS



# Some example in host programming environment



# Executing Spark-SQL statements

- Spark SQL allows us manipulating “data-frame objects” with SQL queries.
  - That is we can execute an SQL statement on a Dataframe object.
- Doing this, requires us registering a dataframe as View using “`createOrReplaceTempView()`” “~~`registerAsTempTable()`~~”
- Spark Session provides SQL method for executing SQL queries. This method always returns a dataframe object.
- We can mix DataFrame methods and SQL queries in the same code.



# Example: Spark-SQL

- Suppose, we already have a dataset object `emp` constructed.
- Below is how we can execute a SQL statement on this. Details of result dataset are printed, and output is shown here.

```
emp.createOrReplaceTempView("employee");  
String sql = "SELECT eno, name, salary FROM employee WHERE dno=4";  
Dataset<Row> emp_dno4 = spark.sql( sql );  
emp_dno4.show();  
emp_dno4.printSchema();
```

eno	name	salary
106	Jennifer	43000
107	Ahmad	25000
108	Alicia	25000

```
root  
|-- eno: string (nullable = true)  
|-- name: string (nullable = true)  
|-- salary: integer (nullable = true)
```



# Example: Registering Tables

```
SparkSession spark = SparkSession.builder()  
    .appName("Spark Demo").master("local").getOrCreate();
```

```
Dataset<Row> emp = spark.read()  
    .option("header", "true")  
    .option("sep", ",")  
    .option("inferSchema", "true")  
    .csv("data/employee_m.csv");
```

```
Dataset<Row> dep = spark.read()  
    .option("header", "true")  
    .option("sep", ",")  
    .option("inferSchema", "true")  
    .csv("data/department_m.csv");
```

```
emp.createOrReplaceTempView("employee");  
dep.createOrReplaceTempView("department");
```



# Spark SQL Examples

- Various SQL operations on Employee and Department tables.  
Code should be self explanatory!

```
String sql = "SELECT d.name, e.name FROM employee e JOIN "  
            + "department d ON (e.dno=d.dno)";  
spark.sql( sql ).show();
```

```
sql = "SELECT e.name, s.name FROM employee e LEFT JOIN "  
      + "employee s ON (e.sup_eno=s.eno)";  
spark.sql( sql ).show();
```

```
sql = "SELECT dno, sum(salary) as TotalSal FROM employee "  
      + "group by dno order by sum(salary) desc ";  
spark.sql( sql ).show();
```



# (Python) Example: Spark SQL

```
sqlContext = SQLContext(sc)
```

```
rdd1 = (content.filter(lambda line: line.split("\t")[7] != "NULL")  
        .map(lambda line: (line.split("\t")[1], float(line.split("\t")[7])))  
        )
```

```
df = sqlContext.createDataFrame(rdd1, schema = ["Name", "Weight"])
```

```
df.createOrReplaceTempView("df_table")
```

```
sqlContext.sql(" SELECT * FROM df_table  ORDER BY Weight DESC limit 15").show()
```



# Spark SQL – Procedural integration<sup>[1]</sup>

- Spark SQL can seamlessly integrate with procedures through its concept of User Defined Functions!
- Below is Scala example

```
val model: LogisticRegressionModel = ...  
  
ctx.udf.register("predict",  
  (x: Float, y: Float) => model.predict(Vector(x, y)))  
  
ctx.sql("SELECT predict(age, weight) FROM users")
```

[2] Armbrust, Michael, et al. "Spark sql: Relational data processing in spark." *Proceedings of the 2015 ACM SIGMOD international conference on management of data*. ACM, 2015





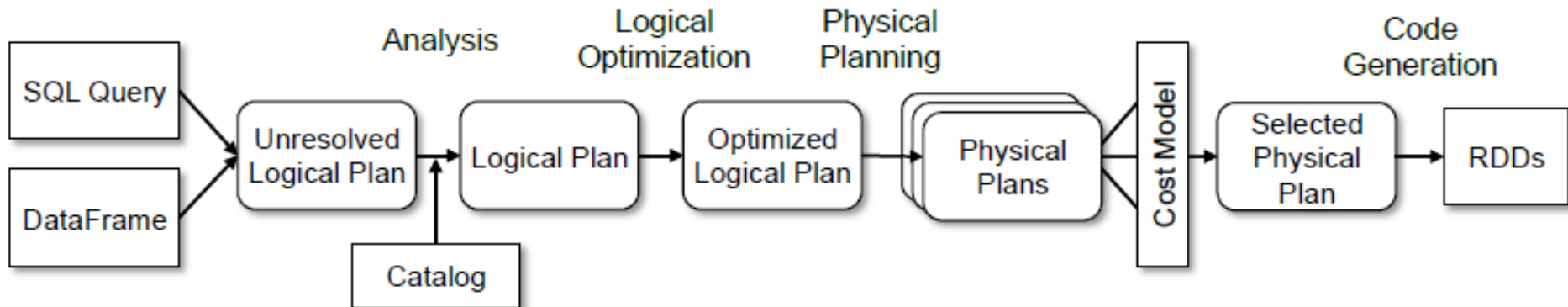
# Query Optimization in RDB

1. Parse the query and convert into RA tree
2. Logical Optimization: Reorganize the RA-Tree into “better evaluation tree” using certain “algebraic equivalence” and “heuristics rules” rules.
  1. For example: performing selection as early in the sequence is better in terms of query execution time
  2. Join as sequence of cross-product followed by Selection (selection-cond) (Emp cross dep )
3. Physical Optimization: choose optimal physical plan, which is sequence of “physical operations” (like file scan, index scan, or so)



# Spark SQL Optimizer Catalyst<sup>[1]</sup>

- All the statements are cached as Abstract Syntax Tree (AST)
- Lazy evaluation of AST enables optimization of expressed operations.
- Diagram here depicts optimization pipeline



[2] Armbrust, Michael, et al. "Spark sql: Relational data processing in spark." *Proceedings of the 2015 ACM SIGMOD international conference on management of data*. ACM, 2015



# Spark SQL Optimizer Catalyst<sup>[1]</sup>

- Optimizer can work based on rules and use some cost based model for choosing an execution plan.
- Catalyst is designed to be Extensible, and workload specific optimizers can be created.

[2] Armbrust, Michael, et al. "Spark sql: Relational data processing in spark." *Proceedings of the 2015 ACM SIGMOD international conference on management of data*. ACM, 2015



# References/Further Reading

- [1] Armbrust, Michael, et al. "Spark SQL: Relational data processing in spark." *Proceedings of the 2015 ACM SIGMOD international conference on management of data*. ACM, 2015.
- [2] Spark SQL, DataFrames and Datasets Guide  
<https://spark.apache.org/docs/latest/sql-programming-guide.html>
- [3] Spark SQL language reference  
<https://docs.databricks.com/spark/2.x/spark-sql/language-manual/index.html>
- [4] Documentation pages (version 2.4)  
<https://spark.apache.org/docs/2.4.0/>