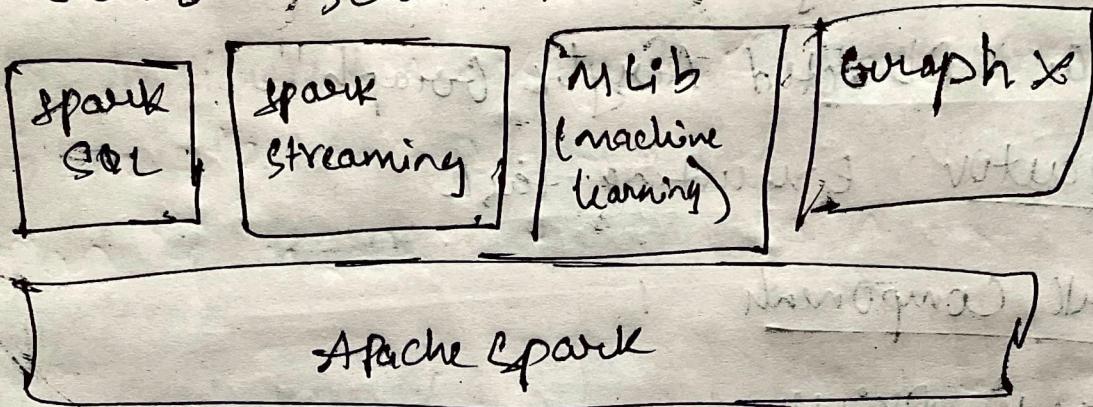


SPARK SQL

- It is a Spark module for structured data processing
- It is a component on top of Spark Core that introduces a new data abstraction called "Schema RDD"



Challenges:

→ Perform ETL to and from various data sources

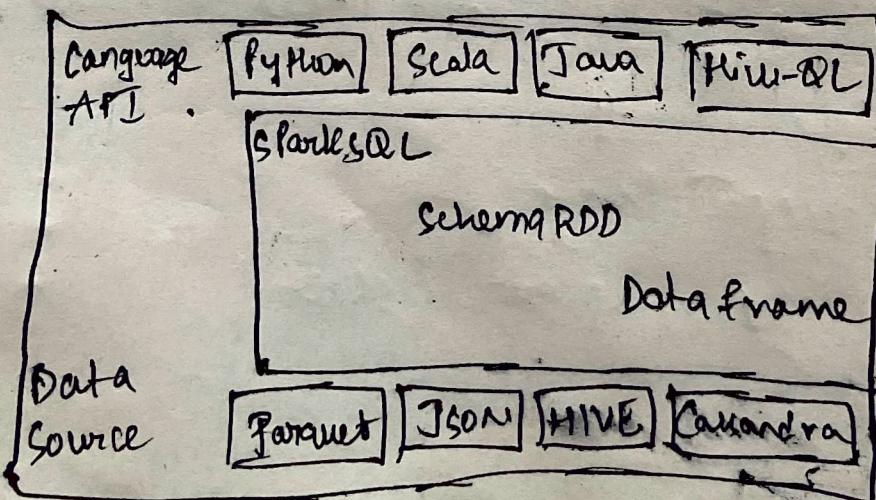
Solutions:

→ A Dataframe API that can perform relational operations.

→ Perform advanced analytics

→ A highly extensible optimizer, catalyst which uses scala to add rules etc.

Architecture:



Features:

→ Integrated:

- ① Integrated APIs in python, Scala, Java
- ② easy to run SQL queries algorithms.

→ Unified Data Access:

- ① Load up every data from various Broads
- ② Schemas RDD provide a single interface for affinity working with structured data.

→ Flink compatibility

- ① Run unmodified engines on workbenches
- ② Simply install it alongside Flink

→ Scalability

- ① Connect through JDBC or ODBC makes the same engine for both engines
- ② takes advantage of RDD Model to support mid-way fault tolerance.

SPARK RDD

→ RDD is a fundamental data structures of spark

→ It is immutable collection of objects that can be stored in memory.

→ RDD is divided into logical partitions

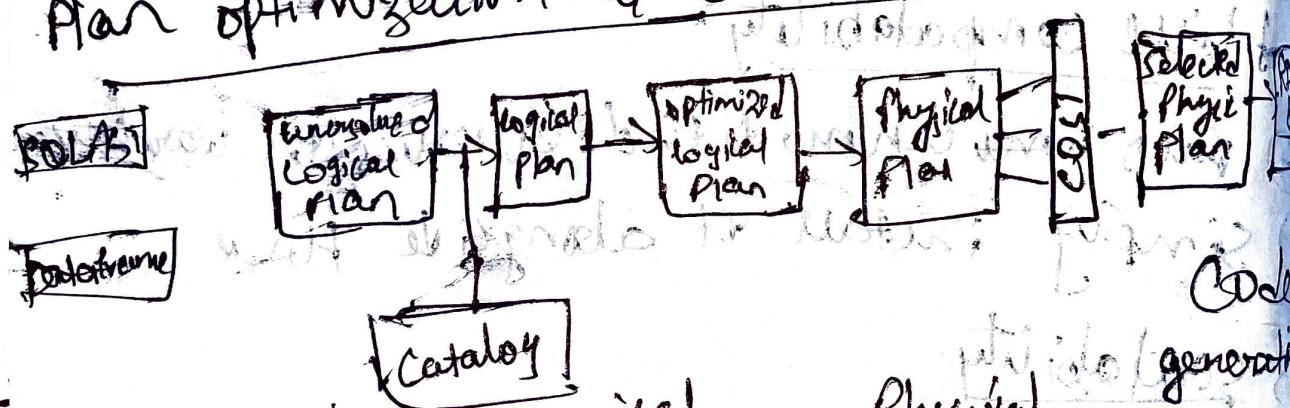
→ Parallel functions transformations

→ Automatically rebuilt on failure

Dataset and Data frame

- A distributed collection of data, into named columns.
- equivalent to relational tables
- The API was designed for modern big data, e.g. Data Science.

Plan optimization & Execution



Analytic logical planning

logical optimization

Physical planning

Code generation