



厦门大学研究生课程 《大数据处理技术Spark》

<http://dblab.xmu.edu.cn/post/7659/>

温馨提示：编辑幻灯片母版，可以修改每页PPT的厦大校徽和底部文字

第6章 Spark SQL

(PPT版本号：2017年春季学期)



扫一扫访问班级主页

林子雨

厦门大学计算机科学系

E-mail: ziyulin@xmu.edu.cn ▶▶

主页: <http://www.cs.xmu.edu.cn/linziyu>





提纲

- 6.1 Spark SQL简介
- 6.2 DataFrame与RDD的区别
- 6.3 DataFrame的创建
- 6.4 从RDD转换得到DataFrame
- 6.5 读取和保存数据



厦门大学林子雨



子雨大数据之Spark入门教程

披荆斩棘，在大数据丛林中开辟学习捷径

免费在线教程：<http://dblab.xmu.edu.cn/blog/spark/>





6.1 Spark SQL简介

6.1.1 从Shark说起

6.1.2 Spark SQL设计

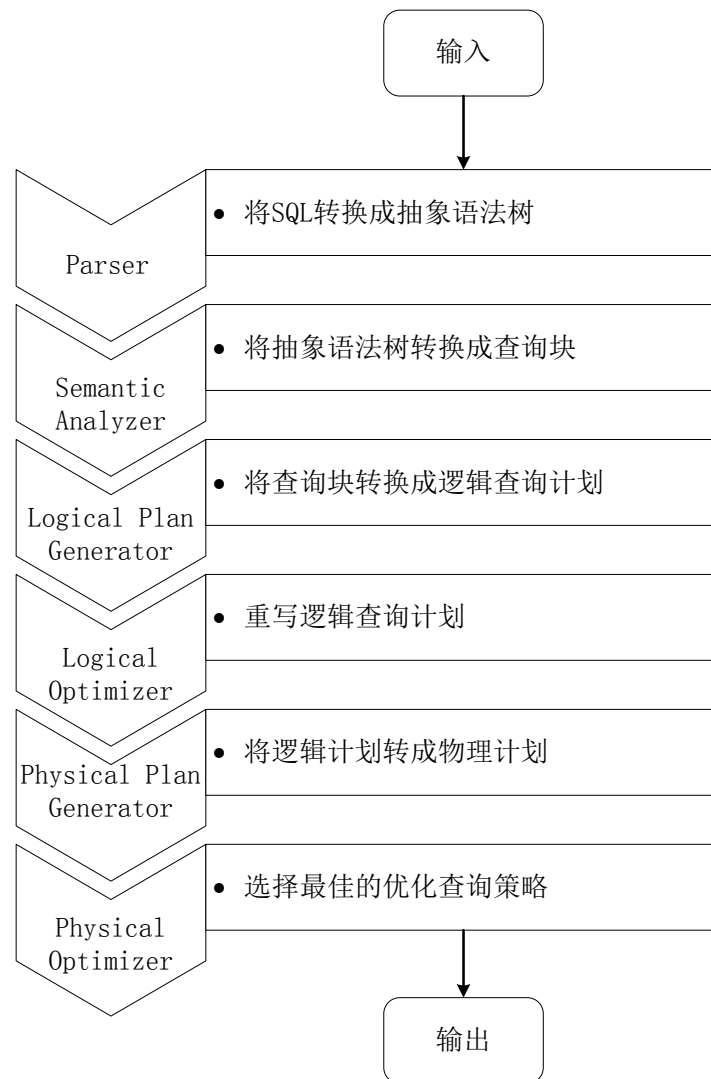
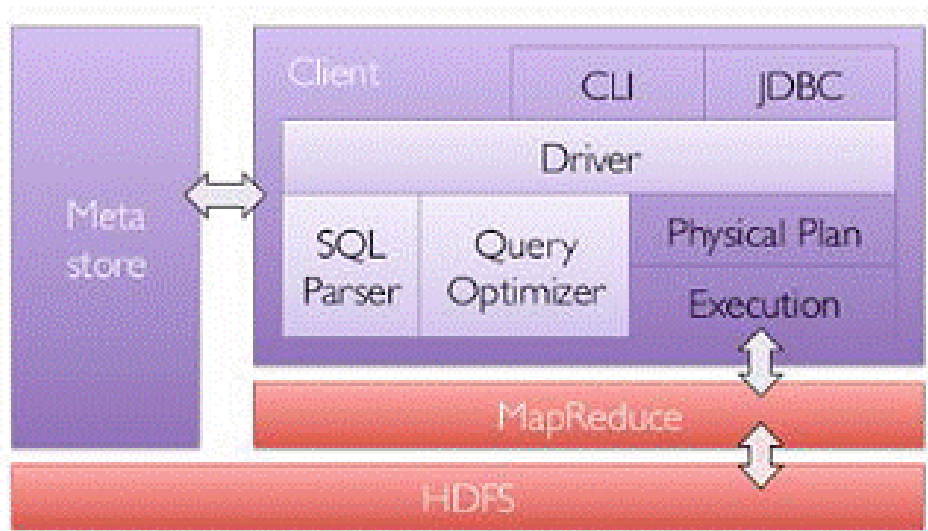
6.1.3 为什么推出Spark SQL



6.1.1 从Shark说起

Hive: SQL-on-Hadoop

Hive Architecture

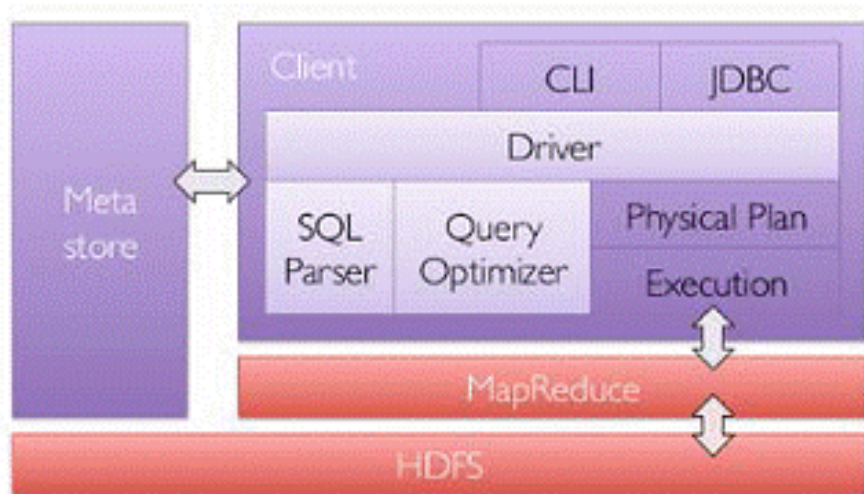


Hive中SQL查询的MapReduce作业转化过程

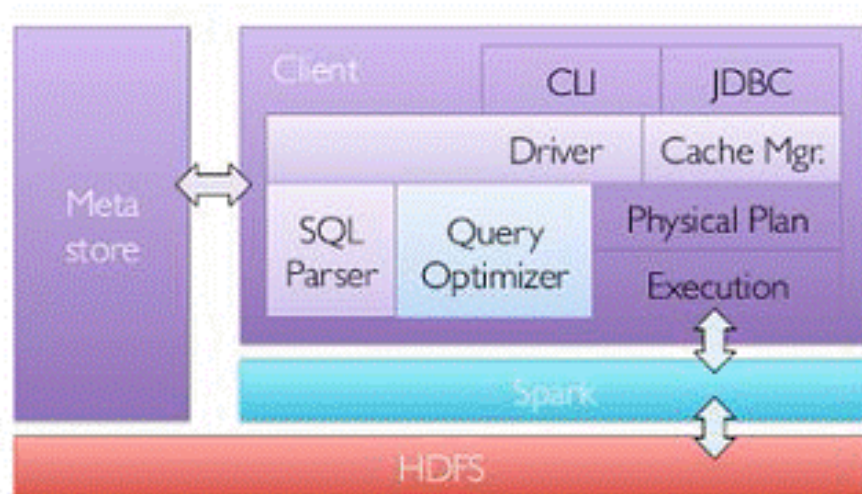


6.1.1 从Shark说起

Hive Architecture



Shark Architecture

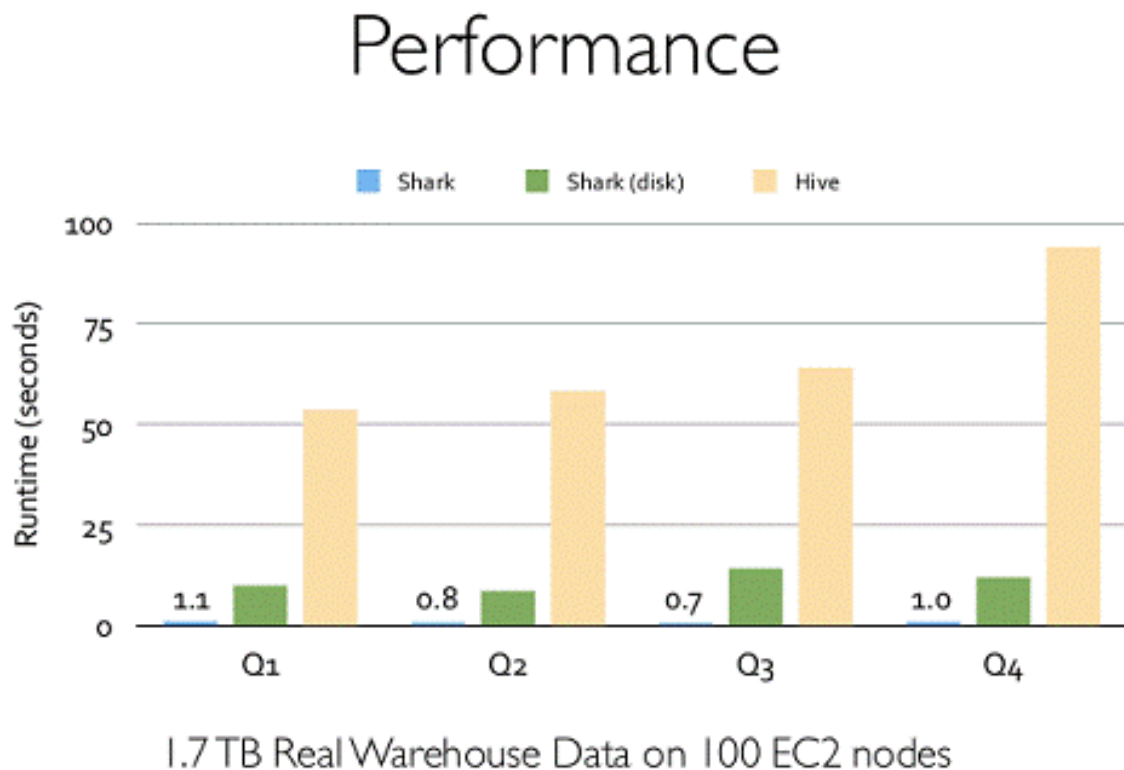


•Shark即Hive on Spark，为了实现与Hive兼容，Shark在HiveQL方面重用了Hive中HiveQL的解析、逻辑执行计划翻译、执行计划优化等逻辑，可以近似认为仅将物理执行计划从MapReduce作业替换成了Spark作业，通过Hive的HiveQL解析，把HiveQL翻译成Spark上的RDD操作



6.1.1 从Shark说起

Shark的出现，使得SQL-on-Hadoop的性能比Hive有了10-100倍的提高





6.1.1 从Shark说起

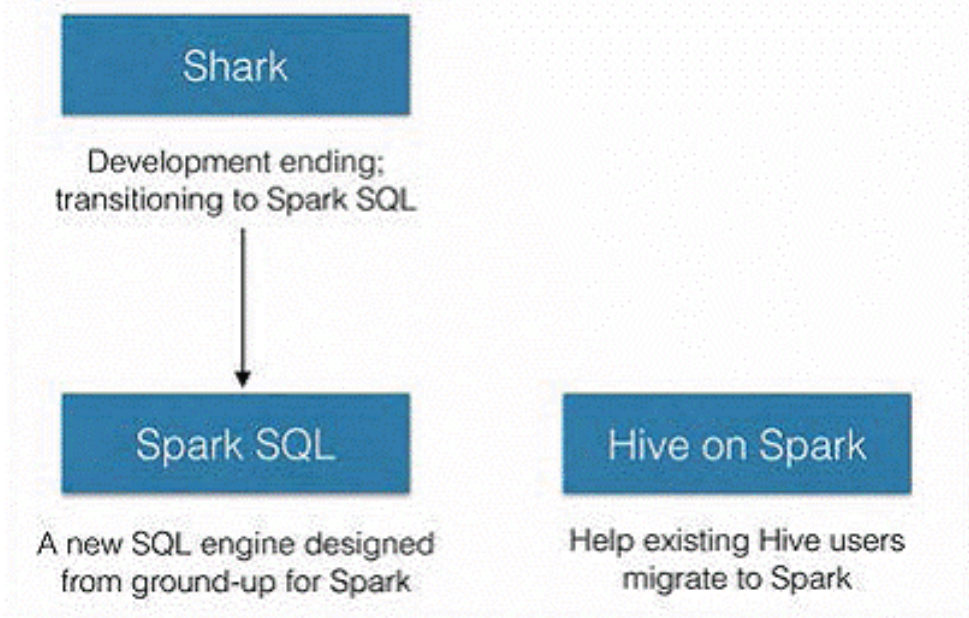
Shark的设计导致了两个问题：

- 一是执行计划优化完全依赖于Hive，不方便添加新的优化策略
- 二是因为Spark是线程级并行，而MapReduce是进程级并行，因此，Spark在兼容Hive的实现上存在线程安全问题，导致Shark不得不使用另外一套独立维护的打了补丁的Hive源码分支



6.1.1 从Shark说起

2014年6月1日Shark项目和SparkSQL项目的主持人Reynold Xin宣布：停止对Shark的开发，团队将所有资源放SparkSQL项目上，至此，Shark的发展画上了句号，但也因此发展出两个直线：SparkSQL和Hive on Spark



- Spark SQL作为Spark生态的一员继续发展，而不再受限于Hive，只是兼容Hive
- Hive on Spark是一个Hive的发展计划，该计划将Spark作为Hive的底层引擎之一，也就是说，Hive将不再受限于一个引擎，可以采用Map-Reduce、Tez、Spark等引擎



6.1.2 Spark SQL设计

Spark SQL在Hive兼容层面仅依赖HiveQL解析、Hive元数据，也就是说，从HQL被解析成抽象语法树（AST）起，就全部由Spark SQL接管了。Spark SQL执行计划生成和优化都由Catalyst（函数式关系查询优化框架）负责

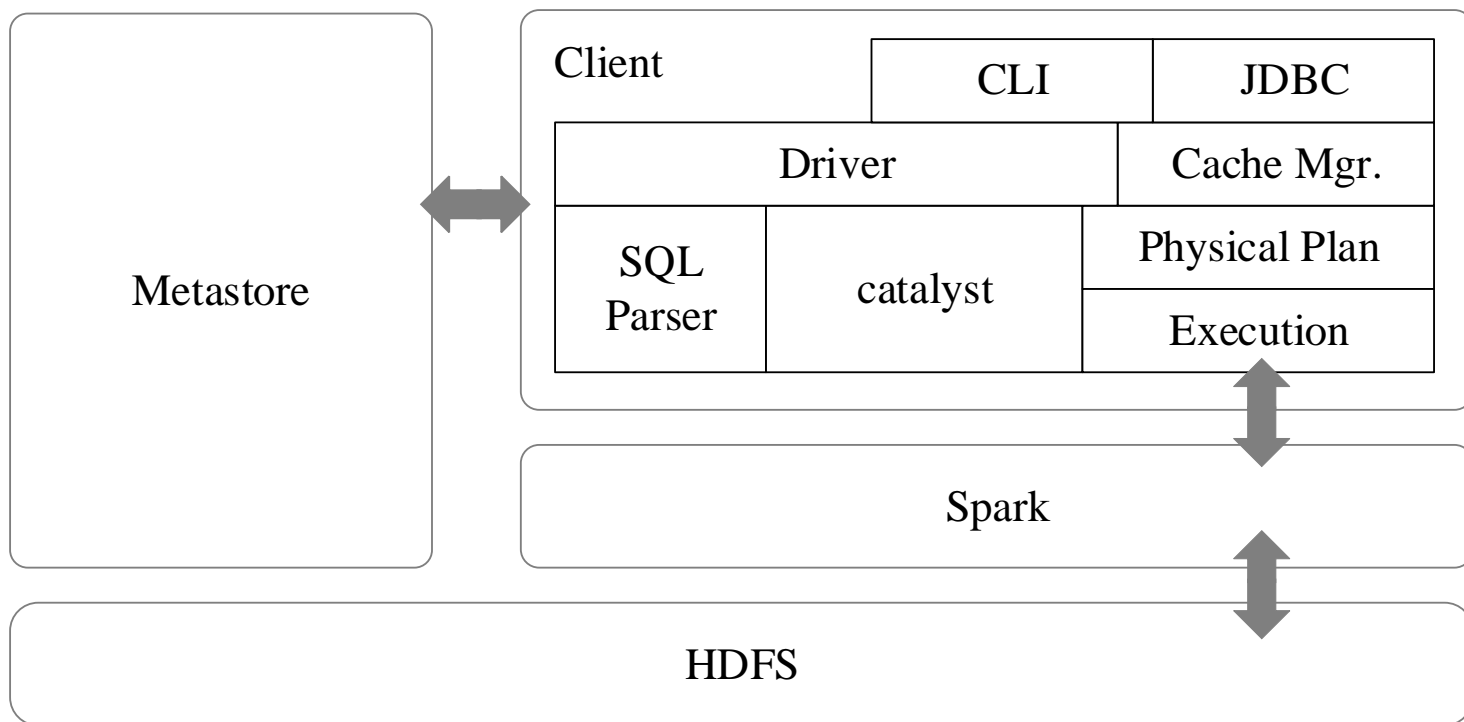


图 Spark SQL架构



6.1.2 Spark SQL设计

- Spark SQL增加了DataFrame（即带有Schema信息的RDD），使用户可以在Spark SQL中执行SQL语句，数据既可以来自RDD，也可以是Hive、HDFS、Cassandra等外部数据源，还可以是JSON格式的数据
- Spark SQL目前支持Scala、Java、Python三种语言，支持SQL-92规范

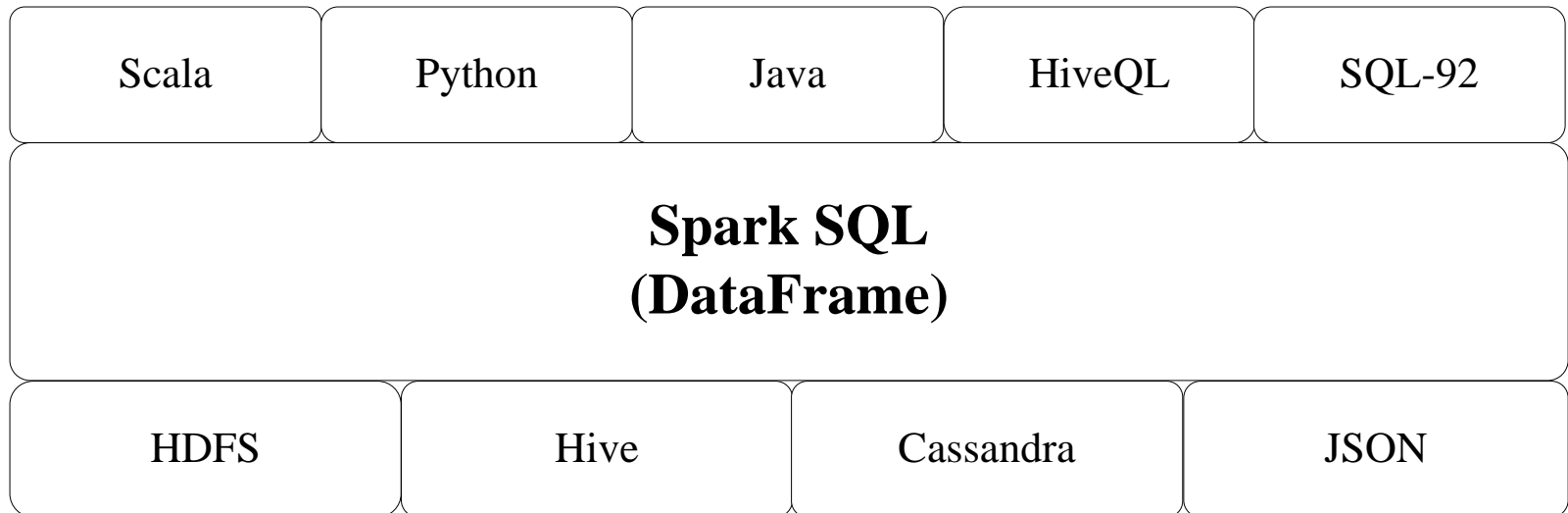


图 Spark SQL支持的数据格式和编程语言



6.1.3 为什么推出Spark SQL

Spark SQL: Relational Data Processing in Spark

Michael Armbrust¹, Reynold S. Xin¹, Cheng Lian¹, Yin Huai¹, Davies Liu¹, Joseph K. Bradley¹,
Xiangrui Meng¹, Tomer Kaftan¹, Michael J. Franklin^{1*}, Ali Ghodsi¹, Matei Zaharia^{1*}

¹Databricks Inc. ^{*}MIT CSAIL [†]AMPLab, UC Berkeley

While the popularity of relational systems shows that users often prefer writing declarative queries, the relational approach is insufficient for many big data applications. First, users want to perform ETL to and from various data sources that might be semi- or unstructured, requiring custom code. Second, users want to perform advanced analytics, such as machine learning and graph processing, that are challenging to express in relational systems. In practice, we have observed that most data pipelines would ideally be expressed with a combination of both relational queries and complex procedural algorithms. Unfortunately, these two classes of systems—relational and procedural—have until now remained largely disjoint,



6.1.3 为什么推出Spark SQL

Spark SQL: Relational Data Processing in Spark

Michael Armbrust¹, Reynold S. Xin¹, Cheng Lian¹, Yin Huai¹, Davies Liu¹, Joseph K. Bradley¹,
Xiangrui Meng¹, Tomer Kaftan¹, Michael J. Franklin^{2*}, Ali Ghodsi¹, Matei Zaharia^{1*}

¹Databricks Inc. ²MIT CSAIL ³AMPLab, UC Berkeley

Spark SQL bridges the gap between the two models through two contributions. First, Spark SQL provides a *DataFrame API* that can perform relational operations on both external data sources and Spark's built-in distributed collections. This API is similar to the widely used data frame concept in R [32], but evaluates operations lazily so that it can perform relational optimizations. Second, to support the wide range of data sources and algorithms in big data, Spark SQL introduces a novel extensible optimizer called *Catalyst*. Catalyst makes it easy to add data sources, optimization rules, and



6.1.3 为什么推出Spark SQL

- 关系数据库已经很流行
- 关系数据库在大数据时代已经不能满足要求
 - 首先，用户需要从不同数据源执行各种操作，包括结构化和非结构化数据
 - 其次，用户需要执行高级分析，比如机器学习和图像处理
- 在实际大数据应用中，经常需要融合关系查询和复杂分析算法（比如机器学习或图像处理），但是，缺少这样的系统

Spark SQL填补了这个鸿沟：

- 首先，可以提供**DataFrame API**，可以对内部和外部各种数据源执行各种关系操作
- 其次，可以支持大量的数据源和数据分析算法

Spark SQL可以融合：传统关系数据库的结构化数据管理能力和机器学习算法的数据处理能力



6.2 DataFrame与RDD的区别

- DataFrame的推出，让Spark具备了处理大规模结构化数据的能力，不仅比原有的RDD转化方式更加简单易用，而且获得了更高的计算性能
- Spark能够轻松实现从MySQL到DataFrame的转化，并且支持SQL查询

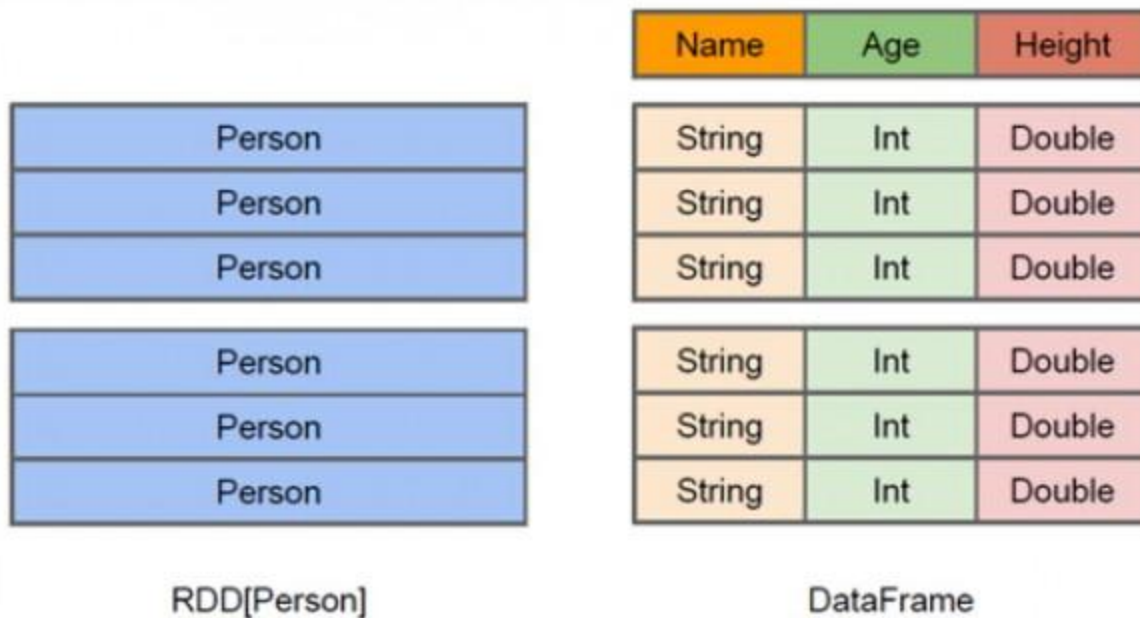


图 DataFrame与RDD的区别

- RDD是分布式的 Java对象的集合，但是，对象内部结构对于RDD而言却是不可知的
- DataFrame是一种以RDD为基础的分布式数据集，提供了详细的结构信息



6.3 DataFrame的创建

- 从Spark2.0以上版本开始，Spark使用全新的SparkSession接口替代Spark1.6中的SQLContext及HiveContext接口来实现其对数据加载、转换、处理等功能。SparkSession实现了SQLContext及HiveContext所有功能
- SparkSession支持从不同的数据源加载数据，并把数据转换成DataFrame，并且支持把DataFrame转换成SQLContext自身中的表，然后使用SQL语句来操作数据。SparkSession亦提供了HiveQL以及其他依赖于Hive的功能的支持



6.3 DataFrame的创建

- 在“/usr/local/spark/examples/src/main/resources/”这个目录下，这个目录下有两个样例数据people.json和people.txt。
people.json文件的内容如下：

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

people.txt文件的内容如下：

```
Michael, 29  
Andy, 30  
Justin, 19
```



6.3 DataFrame的创建

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

scala> val spark=SparkSession.builder().getOrCreate()
spark: org.apache.spark.sql.SparkSession = org.apache.spark.sql.SparkSession@2bdab835

//使支持RDDs转换为DataFrames及后续sql操作
scala> import spark.implicits._
import spark.implicits._

scala> val df = spark.read.json("file:///usr/local/spark/examples/src/main/resources/people.json")
df: org.apache.spark.sql.DataFrame = [age: bigint, name: string]

scala> df.show()
+----+-----+
| age|  name|
+----+-----+
| null|Michael|
|  30|  Andy|
|  19| Justin|
+----+-----+
```



6.3 DataFrame的创建

可以执行一些常用的DataFrame操作

```
// 打印模式信息
scala> df.printSchema()
root
 |-- age: long (nullable = true)
 |-- name: string (nullable = true)

// 选择多列
scala> df.select(df("name"),df("age")+1).show()
+-----+-----+
|  name|(age + 1)|
+-----+-----+
|Michael|      null|
|  Andy|       31|
| Justin|       20|
+-----+-----+

// 条件过滤
scala> df.filter(df("age") > 20 ).show()
+---+-----+
|age|name|
+---+-----+
| 30|Andy|
+---+-----+
```




6.3 DataFrame的创建

```
// 分组聚合
scala> df.groupBy("age").count().show()
+----+-----+
| age|count|
+----+-----+
|  19|    1|
| null|    1|
|  30|    1|
+----+-----+

// 排序
scala> df.sort(df("age").desc).show()
+----+-----+
| age|  name|
+----+-----+
|  30|  Andy|
|  19| Justin|
| null|Michael|
+----+-----+
```



6.3 DataFrame的创建

//多列排序

```
scala> df.sort(df("age").desc, df("name").asc).show()
```

```
+----+-----+
| age|   name|
+----+-----+
|  30|   Andy|
|  19|  Justin|
| null|Michael|
+----+-----+
```

//对列进行重命名

```
scala> df.select(df("name").as("username"),df("age")).show()
```

```
+-----+-----+
|username| age|
+-----+-----+
| Michael| null|
|   Andy|  30|
|  Justin|  19|
+-----+-----+
```



6.4 从RDD转换得到DataFrame

- 6.4.1 利用反射机制推断RDD模式
- 6.4.2 使用编程方式定义RDD模式
- 6.4.3 把RDD保存成文件



6.4.1 利用反射机制推断RDD模式

在利用反射机制推断RDD模式时，需要首先定义一个case class，因为，只有case class才能被Spark隐式地转换为DataFrame

```
scala> import org.apache.spark.sql.catalyst.encoders.ExpressionEncoder
import org.apache.spark.sql.catalyst.encoders.ExpressionEncoder

scala> import org.apache.spark.sql.Encoder
import org.apache.spark.sql.Encoder

scala> import spark.implicits._ //导入包，支持把一个RDD隐式转换为一个
DataFrame
import spark.implicits._
```

剩余代码见下一页



6.4.1 利用反射机制推断RDD模式

```
scala> case class Person(name: String, age: Long) #定义一
↑case class
defined class Person
scala> val peopleDF =
spark.sparkContext.textFile("file:///usr/local/spark/examples/s
rc/main/resources/people.txt").map(_._split(",")).map(attributes
=> Person(attributes(0), attributes(1).trim.toInt)).toDF()
peopleDF: org.apache.spark.sql.DataFrame = [name: string,
age: bigint]
```

剩余代码见下一页



6.4.1 利用反射机制推断RDD模式

```
scala> peopleDF.createOrReplaceTempView("people") #必须注册为临时表才能供下面的查询使用
```

```
scala> val personsRDD = spark.sql("select name,age from people where age > 20")
```

#最终生成一个DataFrame

```
personsRDD: org.apache.spark.sql.DataFrame = [name: string, age: bigint]
```

```
scala> personsRDD.map(t => "Name:"+t(0)+","+"Age:"+t(1)).show()
```

#DataFrame中的每个元素都是一行记录，包含name和age两个字段，分别用t(0)和t(1)来获取值

```
+-----+
```

```
| value|
```

```
+-----+
```

```
|Name:Michael,Age:29|
```

```
| Name:Andy,Age:30|
```

```
+-----+
```



6.4.2 使用编程方式定义RDD模式

当无法提前定义case class时，就需要采用编程方式定义RDD模式

```
scala> import org.apache.spark.sql.types._
import org.apache.spark.sql.types._
scala> import org.apache.spark.sql.Row
import org.apache.spark.sql.Row
//生成 RDD
scala> val peopleRDD =
spark.sparkContext.textFile("file:///usr/local/spark/examples/
src/main/resources/people.txt")
peopleRDD: org.apache.spark.rdd.RDD[String] =
file:///usr/local/spark/examples/src/main/resources/people.txt
MapPartitionsRDD[1] at textFile at <console>:26
```

剩余代码见下一页



6.4.2 使用编程方式定义RDD模式

//定义一个模式字符串

```
scala> val schemaString = "name age"
```

```
schemaString: String = name age
```

//根据模式字符串生成模式

```
scala> val fields = schemaString.split(" ").map(fieldName =>  
StructField(fieldName, StringType, nullable = true))
```

```
fields: Array[org.apache.spark.sql.types.StructField] =  
Array(StructField(name,StringType,true),  
StructField(age,StringType,true))
```

剩余代码见下一页



6.4.2 使用编程方式定义RDD模式

```
scala> val schema = StructType(fields)
schema: org.apache.spark.sql.types.StructType =
StructType(StructField(name,StringType,true),
StructField(age,StringType,true))
```

*//从上面信息可以看出，schema描述了模式信息，模式中
包含name和age两个字段*

剩余代码见下一页



6.4.2 使用编程方式定义RDD模式

//对peopleRDD 这个RDD中的每一行元素都进行解析

```
scala> val rowRDD =  
peopleRDD.map(_.split(",")).map(attributes =>  
Row(attributes(0), attributes(1).trim))  
rowRDD:  
org.apache.spark.rdd.RDD[org.apache.spark.sql.Row] =  
MapPartitionsRDD[3] at map at <console>:29  
  
scala> val peopleDF = spark.createDataFrame(rowRDD,  
schema)  
peopleDF: org.apache.spark.sql.DataFrame = [name:  
string, age: string]
```

剩余代码见下一页



6.4.2 使用编程方式定义RDD模式

//必须注册为临时表才能供下面查询使用

```
scala> peopleDF.createOrReplaceTempView("people")
```

```
scala> val results = spark.sql("SELECT name,age FROM people")
results: org.apache.spark.sql.DataFrame = [name: string, age: string]
```

```
scala> results.map(attributes => "name: " +
attributes(0)+", "+ "age:" +attributes(1)).show()
```

```
+-----+
```

```
| value|
```

```
+-----+
```

```
|name: Michael,age:29|
```

```
| name: Andy,age:30|
```

```
| name: Justin,age:19|
```

```
+-----+
```



6.4.3 把RDD保存成文件

第1种保存方法:

```
scala> val peopleDF =  
spark.read.format("json").load("file:///usr/local/spark/examples/src/main/resources/people.json")  
peopleDF: org.apache.spark.sql.DataFrame = [age: bigint, name: string]  
scala> peopleDF.select("name",  
"age").write.format("csv").save("file:///usr/local/spark/mycode/newpeople.csv")
```

write.format()支持输出 json,parquet, jdbc, orc, libsvm, csv, text等格式文件

可以看到/usr/local/spark/mycode/这个目录下面有个newpeople.csv文件夹, 包含如下两个文件:

```
part-r-00000-33184449-cb15-454c-a30f-9bb43faccac1.csv  
_SUCCESS
```

如果要再次把newpeople.csv中的数据加载到RDD中, 可以直接使用newpeople.csv目录名称, 而不需要使用part-r-00000-33184449-cb15-454c-a30f-9bb43faccac1.csv 文件

```
scala> val textFile = sc.textFile("file:///usr/local/spark/mycode/newpeople.csv")
```



6.4.3 把RDD保存成文件

第2种保存方法:

```
scala> val peopleDF =  
spark.read.format("json").load("file:///usr/local/spark/examples/src/main/resources/people.json")  
peopleDF: org.apache.spark.sql.DataFrame = [age: bigint, name: string]  
scala> df.rdd.saveAsTextFile("file:///usr/local/spark/mycode/newpeople.txt")
```

可以看到/usr/local/spark/mycode/这个目录下面有个newpeople.txt文件夹，内容如下:

```
part-00000  
_SUCCESS
```

如果要再次把newpeople.txt中的数据加载到RDD中，可以直接使用newpeople.txt目录名称，而不需要使用part-00000文件

```
scala> val textFile = sc.textFile("file:///usr/local/spark/mycode/newpeople.txt")
```



6.5 读取和保存数据

Spark SQL可以支持Parquet、JSON、Hive等数据源，并且可以通过JDBC连接外部数据源

6.5.1 读写Parquet

6.5.2 通过JDBC连接数据库

6.5.3 连接Hive读写数据



6.5.1 读写Parquet

Parquet是一种流行的列式存储格式，可以高效地存储具有嵌套字段的记录。**Parquet**是语言无关的，而且不与任何一种数据处理框架绑定在一起，适配多种语言和组件，能够与**Parquet**配合的组件有：

- 查询引擎: Hive, Impala, Pig, Presto, Drill, Tajo, HAWQ, IBM Big SQL
- 计算框架: MapReduce, Spark, Cascading, Crunch, Scalding, Kite
- 数据模型: Avro, Thrift, Protocol Buffers, POJOs



6.5.1 读写Parquet

从parquet文件中加载数据生成DataFrame

```
scala> import spark.implicits._
import spark.implicits._

scala> val parquetFileDF =
spark.read.parquet("file:///usr/local/spark/examples/src/main/resources/
users.parquet")
SLF4J: Failed to load class "org.slf4j.impl.StaticLoggerBinder".
SLF4J: Defaulting to no-operation (NOP) logger implementation
SLF4J: See http://www.slf4j.org/codes.html#StaticLoggerBinder for
further details.
parquetFileDF: org.apache.spark.sql.DataFrame = [name: string,
favorite_color: string ... 1 more field]
```

剩余代码见下一页



6.5.1 读写Parquet

```
scala> parquetFileDF.createOrReplaceTempView("parquetFile")
```

```
scala> val namesDF = spark.sql("SELECT * FROM parquetFile")
namesDF: org.apache.spark.sql.DataFrame = [name: string,
favorite_color: string ... 1 more field]
```

```
scala> namesDF.foreach(attributes =>println("Name: " + attributes(0)+"
favorite color:"+attributes(1)))
```

```
16/12/02 10:18:49 WARN hadoop.ParquetRecordReader: Can not
initialize counter due to context is not a instance of
```

```
TaskInputOutputContext, but is
```

```
org.apache.hadoop.mapreduce.task.TaskAttemptContextImpl
```

```
Name: Alyssa favorite color:null
```

```
Name: Ben favorite color:red
```




6.5.1 读写Parquet

将DataFrame保存成parquet文件

```
scala> import spark.implicits._
import spark.implicits._

scala> val peopleDF =
spark.read.json("file:///usr/local/spark/examples/src/main/resources/people.json")
peopleDF: org.apache.spark.sql.DataFrame = [age: bigint, name: string]

scala>
peopleDF.write.parquet("file:///usr/local/spark/mycode/newpeople.parquet")
```

可以看到” /usr/local/spark/myCode/”这个目录下多了一个newpeople.parquet

```
part-r-00000-8d3a120f-b3b5-4582-b26b-f3693df80d45.snappy.parquet
_SUCCESS
```

```
scala> val users =
spark.read.parquet("file:///usr/local/spark/myCode/people.parquet")
```



6.5.2 通过JDBC连接数据库

6.5.2.1 准备工作

6.5.2.2 读取MySQL数据库中的数据

6.5.2.3 向MySQL数据库写入数据



6.5.2.1 准备工作

- 请参考厦门大学数据库实验室博客教程《Ubuntu安装MySQL》，在Linux系统中安装好MySQL数据库

教程地址：<http://dblab.xmu.edu.cn/blog/install-mysql/>

- 在Linux中启动MySQL数据库

```
$ service mysql start  
$ mysql -u root -p  
#屏幕会提示你输入密码
```

输入下面SQL语句完成数据库和表的创建：

```
mysql> create database spark;  
mysql> use spark;  
mysql> create table student (id int(4), name char(20), gender char(4), age int(4));  
mysql> insert into student values(1,'Xueqian','F',23);  
mysql> insert into student values(2,'Weiliang','M',24);  
mysql> select * from student;
```



6.5.2.1 准备工作

- 下载MySQL的JDBC驱动程序，比如mysql-connector-java-5.1.40.tar.gz
- 把该驱动程序拷贝到spark的安装目录“/usr/local/spark/jars”下
- 启动一个spark-shell，启动Spark Shell时，必须指定mysql连接驱动jar包

```
$ cd /usr/local/spark
$ ./bin/spark-shell \
--jars /usr/local/spark/jars/mysql-connector-java-5.1.40/mysql-connector-java-
5.1.40-bin.jar \
--driver-class-path /usr/local/spark/jars/mysql-connector-java-5.1.40/mysql-
connector-java-5.1.40-bin.jar
```



6.5.2.2 读取MySQL数据库中的数据

执行以下命令连接数据库，读取数据，并显示：

```
scala> val jdbcDF = spark.read.format("jdbc").option("url",  
"jdbc:mysql://localhost:3306/spark").option("driver", "com.mysql.jdbc.Driv  
er").option("dbtable", "student").option("user", "root").option("password",  
"hadoop").load()
```

```
scala> jdbcDF.show()  
+---+-----+-----+---+  
| id| name|gender|age|  
+---+-----+-----+---+  
| 1| Xueqian| F| 23|  
| 2| Weiliang| M| 24|  
+---+-----+-----+---+
```



6.5.2.3 向MySQL数据库写入数据

在MySQL数据库中创建了一个名称为spark的数据库，并创建了一个名称为student的表。创建后，查看一下数据库内容：

```
mysql> use spark;
Database changed

mysql> select * from student;
//上面命令执行后返回下面结果
+-----+-----+-----+-----+
| id  | name    | gender | age  |
+-----+-----+-----+-----+
|  1  | Xueqian | F      |  23  |
|  2  | weiliang | M      |  24  |
+-----+-----+-----+-----+
2 rows in set (0.00 sec)
```



6.5.2.3 向MySQL数据库写入数据

现在开始在spark-shell中编写程序，往spark.student表中插入两条记录

```
import java.util.Properties
import org.apache.spark.sql.types._
import org.apache.spark.sql.Row

//下面我们设置两条数据表示两个学生信息
val studentRDD = spark.sparkContext.parallelize(Array("3 Rongcheng M 26", "4
Guanhua M 27")).map(_.split(" "))

//下面要设置模式信息
val schema = StructType(List(StructField("id", IntegerType,
true),StructField("name", StringType, true),StructField("gender", StringType,
true),StructField("age", IntegerType, true)))
```

剩余代码见下一页



6.5.2.3 向MySQL数据库写入数据

//下面创建Row对象，每个Row对象都是rowRDD中的一行

```
val rowRDD = studentRDD.map(p => Row(p(0).toInt, p(1).trim, p(2).trim, p(3).toInt))
```

//建立起Row对象和模式之间的对应关系，也就是把数据和模式对应起来

```
val studentDF = spark.createDataFrame(rowRDD, schema)
```

//下面创建一个prop变量用来保存JDBC连接参数

```
val prop = new Properties()
```

```
prop.put("user", "root") //表示用户名是root
```

```
prop.put("password", "hadoop") //表示密码是hadoop
```

```
prop.put("driver", "com.mysql.jdbc.Driver") //表示驱动程序是com.mysql.jdbc.Driver
```

//下面就可以连接数据库，采用append模式，表示追加记录到数据库spark的student表中

```
studentDF.write.mode("append").jdbc("jdbc:mysql://localhost:3306/spark",  
"spark.student", prop)
```



6.5.2.3 向MySQL数据库写入数据

可以看一下效果，看看MySQL数据库中的spark.student表发生了什么变化

```
mysql> select * from student;
+-----+-----+-----+-----+
| id | name | gender | age |
+-----+-----+-----+-----+
| 1 | Xueqian | F | 23 |
| 2 | Weiliang | M | 24 |
| 3 | Rongcheng | M | 26 |
| 4 | Guanhua | M | 27 |
+-----+-----+-----+-----+
4 rows in set (0.00 sec)
```



6.5.3 连接Hive读写数据

6.5.3.1 Hive简介和安装

6.5.3.2 让Spark包含Hive支持

6.5.3.3 在Hive中创建数据库和表

6.5.3.4 连接Hive读写数据



6.5.3.1 Hive简介和安装

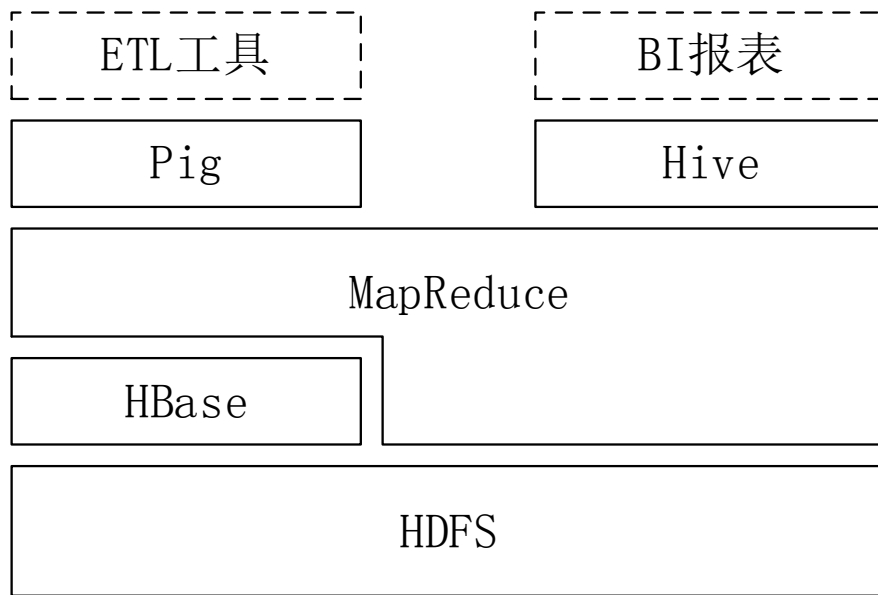
- **Hive**是一个构建于**Hadoop**顶层的数据仓库工具
- 支持大规模数据存储、分析，具有良好的可扩展性
- 某种程度上可以看作是用户编程接口，本身不存储和处理数据
- 依赖分布式文件系统**HDFS**存储数据
- 依赖分布式并行计算模型**MapReduce**处理数据
- 定义了简单的类似**SQL**的查询语言——**HiveQL**
- 用户可以通过编写的**HiveQL**语句运行**MapReduce**任务
- 可以很容易把原来构建在关系数据库上的数据仓库应用程序移植到**Hadoop**平台上
- 是一个可以提供有效、合理、直观组织和使用数据的分析工具



6.5.3.1 Hive简介和安装

- Hive依赖于HDFS 存储数据
- Hive依赖于MapReduce 处理数据

Hadoop生态系统

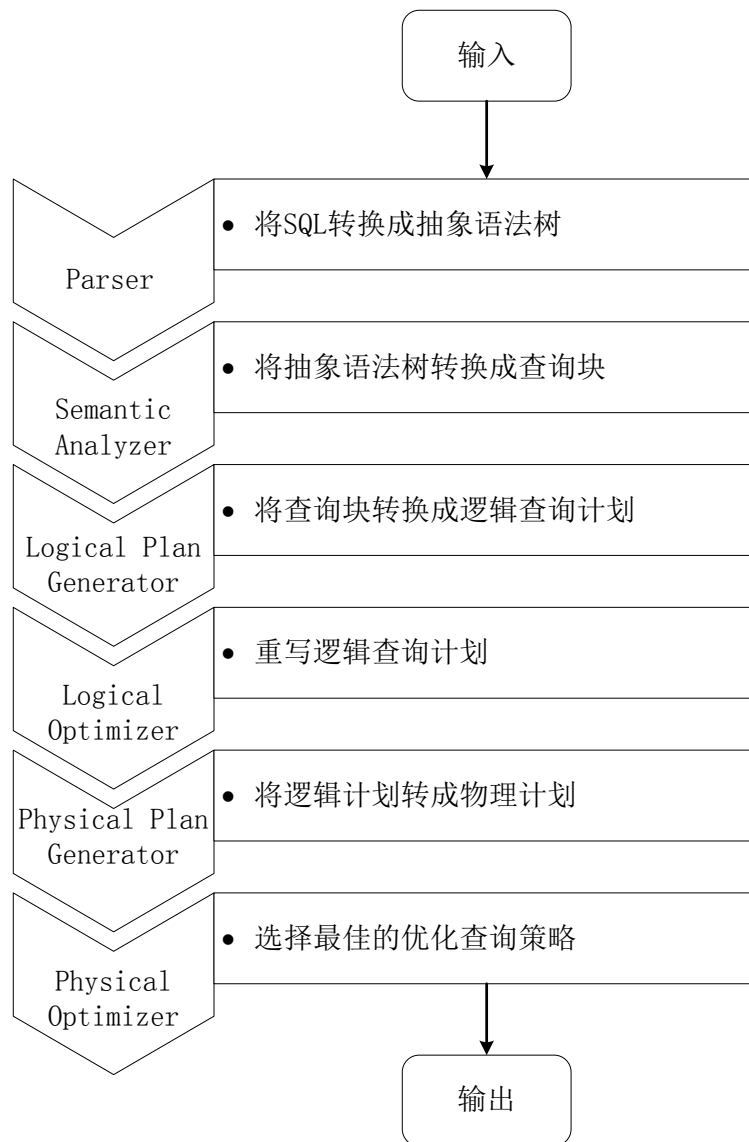
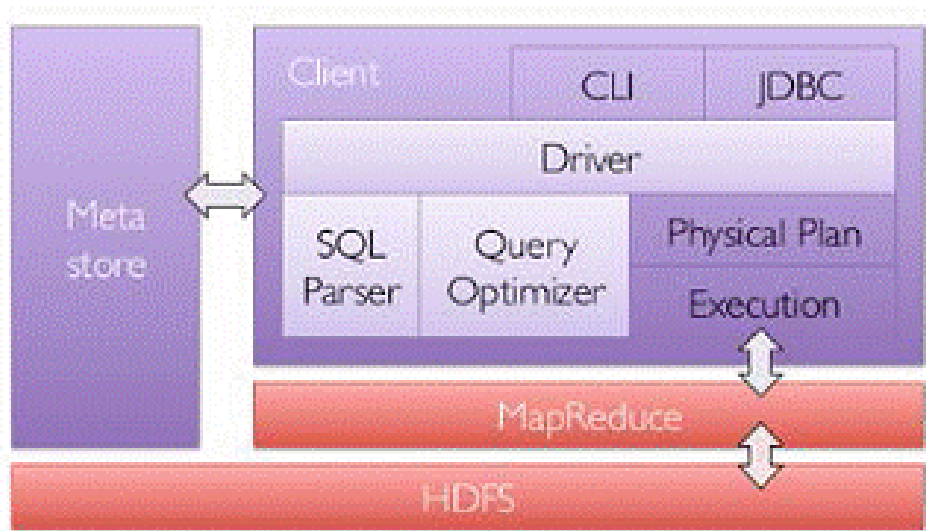


Hive与Hadoop生态系统中其他组件的关系



6.5.3.1 Hive简介和安装

Hive Architecture





6.5.3.1 Hive简介和安装

- Hive的安装，请参考厦门大学数据库实验室建设的中国高校大数据课程公共服务平台上的技术博客：
- 《Ubuntu安装hive，并配置mysql作为元数据库》

<http://dblab.xmu.edu.cn/blog/install-hive/>



中国高校大数据课程
公 共 服 务 平 台



6.5.3.2 让Spark包含Hive支持

- 为了让Spark能够访问Hive，必须为Spark添加Hive支持
- Spark官方提供的预编译版本，通常是不包含Hive支持的，需要采用源码编译，编译得到一个包含Hive支持的Spark版本

(1) 测试一下自己电脑上已经安装的Spark版本是否支持Hive

启动进入了spark-shell，如果不支持Hive，会显示如下信息：

```
scala> import org.apache.spark.sql.hive.HiveContext
<console>:25: error: object hive is not a member of package org.apache.spark.sql
import org.apache.spark.sql.hive.HiveContext
                              ^
```

如果你当前电脑上的Spark版本包含Hive支持，那么应该显示下面的正确信息：

```
scala> import org.apache.spark.sql.hive.HiveContext
import org.apache.spark.sql.hive.HiveContext
```



6.5.3.2 让Spark包含Hive支持

(2) 采用源码编译方法得到支持Hive的Spark版本

到Spark官网下载源码

<http://spark.apache.org/downloads.html>

Download Apache Spark™

1. Choose a Spark release:
2. Choose a package type:
3. Choose a download type:
4. Download Spark: [spark-2.1.0.tgz](#)
5. Verify this release using the [2.1.0 signatures and checksums](#) and [project release KEYS](#).

Note: Starting version 2.0, Spark is built with Scala 2.11 by default. Scala 2.10 users should download the Spark source package and build with Scala 2.10 support.



6.5.3.2 让Spark包含Hive支持

解压文件

```
$ cd /home/hadoop/下载 //spark-2.1.0.tgz就在这个目录下面  
$ ls #可以看到刚才下载的spark-2.1.0.tgz文件  
$ sudo tar -zxf ./spark-2.1.0.tgz -C /home/hadoop/  
$ cd /home/hadoop  
$ ls #这时可以看到解压得到的文件夹spark-2.1.0
```

在编译时，需要给出电脑上之前已经安装好的Hadoop的版本

```
$ hadoop version
```

运行编译命令，对Spark源码进行编译

```
$ cd /home/hadoop/spark-2.1.0  
$ ./dev/make-distribution.sh --tgz --name h27hive -Pyarn -Phadoop-2.7 -  
Dhadoop.version=2.7.1 -Phive -Phive-thriftserver -DskipTests
```

编译成功后会得到文件名“spark-2.1.0-bin-h27hive.tgz”，这个就是包含Hive支持的Spark安装文件



6.5.3.2 让Spark包含Hive支持

(3) 安装支持Hive的Spark版本

Spark的安装详细过程，请参考厦门大学数据库实验室建设的中国高校大数据课程公共服务平台上的技术博客：

《**Spark2.1.0**入门：Spark的安装和使用》

博客地址：<http://dmlab.xmu.edu.cn/blog/1307-2/>



中国高校大数据课程
公共服务平台

启动进入了spark-shell，由于已经可以支持Hive，会显示如下信息：

```
scala> import org.apache.spark.sql.hive.HiveContext
import org.apache.spark.sql.hive.HiveContext
```



6.5.3.3 在Hive中创建数据库和表

假设已经完成了Hive的安装，并且使用的是MySQL数据库来存放Hive的元数据需要借助于MySQL保存Hive的元数据，首先启动MySQL数据库：

```
$ service mysql start
```

由于Hive是基于Hadoop的数据仓库，使用HiveQL语言撰写的查询语句，最终都会被Hive自动解析成MapReduce任务由Hadoop去具体执行，因此，需要启动Hadoop，然后再启动Hive

启动Hadoop：

```
$ cd /usr/local/hadoop  
$ ./sbin/start-all.sh
```

Hadoop启动成功以后，可以再启动Hive：

```
$ cd /usr/local/hive  
$ ./bin/hive
```



6.5.3.3 在Hive中创建数据库和表

进入Hive，新建一个数据库sparktest，并在这个数据库下面创建一个表student，并录入两条数据

```
hive> create database if not exists sparktest; //创建数据库sparktest
hive> show databases; //显示一下是否创建出了sparktest数据库
//下面在数据库sparktest中创建一个表student
hive> create table if not exists sparktest.student(
> id int,
> name string,
> gender string,
> age int);
hive> use sparktest; //切换到sparktest
hive> show tables; //显示sparktest数据库下面有哪些表
hive> insert into student values(1,'Xueqian','F',23); //插入一条记录
hive> insert into student values(2,'Weiliang','M',24); //再插入一条记录
hive> select * from student; //显示student表中的记录
```



6.5.3.4 连接Hive读写数据

需要修改 “/usr/local/sparkwithhive/conf/spark-env.sh”这个配置文件：

```
export SPARK_DIST_CLASSPATH=$(/usr/local/hadoop/bin/hadoop classpath)
export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64
export CLASSPATH=$CLASSPATH:/usr/local/hive/lib
export SCALA_HOME=/usr/local/scala
export HADOOP_CONF_DIR=/usr/local/hadoop/etc/hadoop
export HIVE_CONF_DIR=/usr/local/hive/conf
export SPARK_CLASSPATH=$SPARK_CLASSPATH:/usr/local/hive/lib/mysql-connector-java-5.1.40-bin.jar
```



6.5.3.4 连接Hive读写数据

请在spark-shell（包含Hive支持）中执行以下命令从Hive中读取数据：

```
Scala> import org.apache.spark.sql.Row
Scala> import org.apache.spark.sql.Session
Scala> case class Record(key: Int, value: String)
// warehouseLocation points to the default location for managed databases and tables
Scala> val warehouseLocation = "spark-warehouse"
Scala> val spark = Session.builder().appName("Spark Hive
Example").config("spark.sql.warehouse.dir",
warehouseLocation).enableHiveSupport().getOrCreate()
Scala> import spark.implicits._
Scala> import spark.sql
//下面是运行结果
scala> sql("SELECT * FROM sparktest.student").show()
+---+-----+-----+---+
| id| name|gender|age|
+---+-----+-----+---+
| 1| Xueqian| F| 23|
| 2| Weiliang| M| 24|
+---+-----+-----+---+
```




6.5.3.4 连接Hive读写数据

编写程序向Hive数据库的sparktest.student表中插入两条数据
在插入数据之前，先查看一下已有的2条数据

```
hive> use sparktest;
OK
Time taken: 0.016 seconds

hive> select * from student;
OK
1  Xueqian F    23
2  Weiliang  M    24
Time taken: 0.05 seconds, Fetched: 2 row(s)
```



6.5.3.4 连接Hive读写数据

编写程序向Hive数据库的sparktest.student表中插入两条数据:

```
scala> import java.util.Properties
scala> import org.apache.spark.sql.types._
scala> import org.apache.spark.sql.Row
//下面我们设置两条数据表示两个学生信息
scala> val studentRDD = spark.sparkContext.parallelize(Array("3 Rongcheng M
26", "4 Guanhua M 27")).map(_.split(" "))
//下面要设置模式信息
scala> val schema = StructType(List(StructField("id", IntegerType,
true),StructField("name", StringType, true),StructField("gender", StringType,
true),StructField("age", IntegerType, true)))
//下面创建Row对象，每个Row对象都是rowRDD中的一行
scala> val rowRDD = studentRDD.map(p => Row(p(0).toInt, p(1).trim, p(2).trim,
p(3).toInt))
//建立起Row对象和模式之间的对应关系，也就是把数据和模式对应起来
scala> val studentDF = spark.createDataFrame(rowRDD, schema)
```

剩余代码见下一页



6.5.3.4 连接Hive读写数据

```
//查看studentDF
scala> studentDF.show()
+---+-----+-----+---+
| id| name|gender|age|
+---+-----+-----+---+
| 3|Rongcheng| M| 26|
| 4| Guanhua| M| 27|
+---+-----+-----+---+
//下面注册临时表
scala> studentDF.registerTempTable("tempTable")

scala> sql("insert into sparktest.student select * from tempTable")
```



6.5.3.4 连接Hive读写数据

输入以下命令查看Hive数据库内容的变化:

```
hive> select * from student;
OK
1  Xueqian F    23
2  Weiliang  M    24
3  Rongcheng M    26
4  Guanhua M    27
Time taken: 0.049 seconds, Fetched: 4 row(s)
```

可以看到，插入数据操作执行成功了!



附录：主讲教师林子雨简介



主讲教师：林子雨

单位：厦门大学计算机科学系

E-mail: ziyulin@xmu.edu.cn

个人网页: <http://www.cs.xmu.edu.cn/linziyu>

数据库实验室网站: <http://dblab.xmu.edu.cn>



扫一扫访问个人主页

林子雨，男，1978年出生，博士（毕业于北京大学），现为厦门大学计算机科学系助理教授（讲师），曾任厦门大学信息科学与技术学院院长助理、晋江市发展和改革局副局长。中国计算机学会数据库专业委员会委员，中国计算机学会信息系统专业委员会委员，荣获“2016中国大数据创新百人”称号。中国高校首个“数字教师”提出者和建设者，厦门大学数据库实验室负责人，厦门大学云计算与大数据研究中心主要建设者和骨干成员，2013年度厦门大学奖教金获得者。主要研究方向为数据库、数据仓库、数据挖掘、大数据、云计算和物联网，并以第一作者身份在《软件学报》《计算机学报》和《计算机研究与发展》等国家重点期刊以及国际学术会议上发表多篇学术论文。作为项目负责人主持的科研项目包括1项国家自然科学基金青年基金项目(No.61303004)、1项福建省自然科学基金项目(No.2013J05099)和1项中央高校基本科研业务费项目(No.2011121049)，同时，作为课题负责人完成了国家发改委城市信息化重大课题、国家物联网重大应用示范工程区域试点泉州市工作方案、2015泉州市互联网经济调研等课题。中国高校首个“数字教师”提出者和建设者，2009年至今，“数字教师”大平台累计向网络免费发布超过100万字高价值的研究和教学资料，累计网络访问量超过100万次。打造了中国高校大数据教学知名品牌，编著出版了中国高校第一本系统介绍大数据知识的专业教材《大数据技术原理与应用》，并成为京东、当当网等网店畅销书籍；建设了国内高校首个大数据课程公共服务平台，为教师教学和学生学习大数据课程提供全方位、一站式服务，年访问量超过50万次。具有丰富的政府和企业信息化培训经验，厦门大学管理学院EDP中心、浙江大学管理学院EDP中心、厦门大学继续教育学院、泉州市科技培训中心特邀培训讲师，曾给中国移动通信集团公司、福州马尾区政府、福建龙岩卷烟厂、福建省物联网科学研究院、石狮市物流协会、厦门市物流协会、浙江省中小企业家、四川泸州企业家、江苏沛县企业家等开展信息化培训，累计培训人数达3000人以上。



附录：林子雨编著《Spark入门教程》

厦门大学林子雨编著《Spark入门教程》
教程内容包括Scala语言、Spark简介、安装、运行架构、RDD的设计与运行原理、部署模式、RDD编程、键值对RDD、数据读写、Spark SQL、Spark Streaming、MLlib等



厦门大学林子雨



子雨大数据之Spark入门教程

披荆斩棘，在大数据丛林中开辟学习捷径



免费在线教程：<http://dbllab.xmu.edu.cn/blog/spark/>



附录：《大数据技术原理与应用》教材



扫一扫访问教材官网

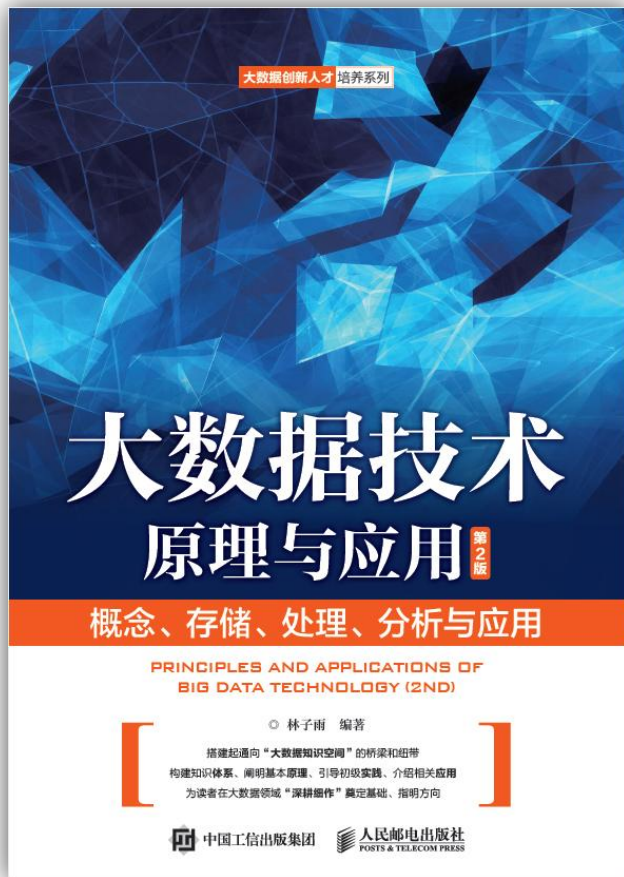
《大数据技术原理与应用——概念、存储、处理、分析与应用（第2版）》，由厦门大学计算机科学系林子雨老师编著，是中国高校第一本系统介绍大数据知识的专业教材。

全书共有15章，系统地论述了大数据的基本概念、大数据处理架构Hadoop、分布式文件系统HDFS、分布式数据库HBase、NoSQL数据库、云数据库、分布式并行编程模型MapReduce、Spark、流计算、图计算、数据可视化以及大数据在互联网、生物医学和物流等各个领域的应用。在Hadoop、HDFS、HBase、MapReduce和Spark等重要章节，安排了入门级的实践操作，让读者更好地学习和掌握大数据关键技术。

本书可以作为高等院校计算机专业、信息管理等相关专业的大数据课程教材，也可供相关技术人员参考、学习、培训之用。

欢迎访问《大数据技术原理与应用——概念、存储、处理、分析与应用（第2版）》教材官方网站：

<http://dbllab.xmu.edu.cn/post/bigdata>





附录：中国高校大数据课程公共服务平台



中国高校大数据课程 公共服务平台

<http://dblab.xmu.edu.cn/post/bigdata-teaching-platform/>



扫一扫访问平台主页



扫一扫观看3分钟FLASH动画宣传片

The background of the slide features a blue gradient with several white silhouettes of people. At the top, there are two groups of people standing and holding hands. On the right side, a person is shown in profile, looking towards the center. On the left side, two people are shown in profile, one appearing to be speaking or gesturing towards the other. The overall theme is one of community and collaboration.

Thank You!

Department of Computer Science, Xiamen University, 2017