

# Spark 其他开源系统



董西成  
2017年5月

# 主要内容

1

Alluxio

2

Spark JobServer

3

Zeppelin

4

总结

# 主要内容

1

Alluxio

2

Spark JobServer

3

Zeppelin

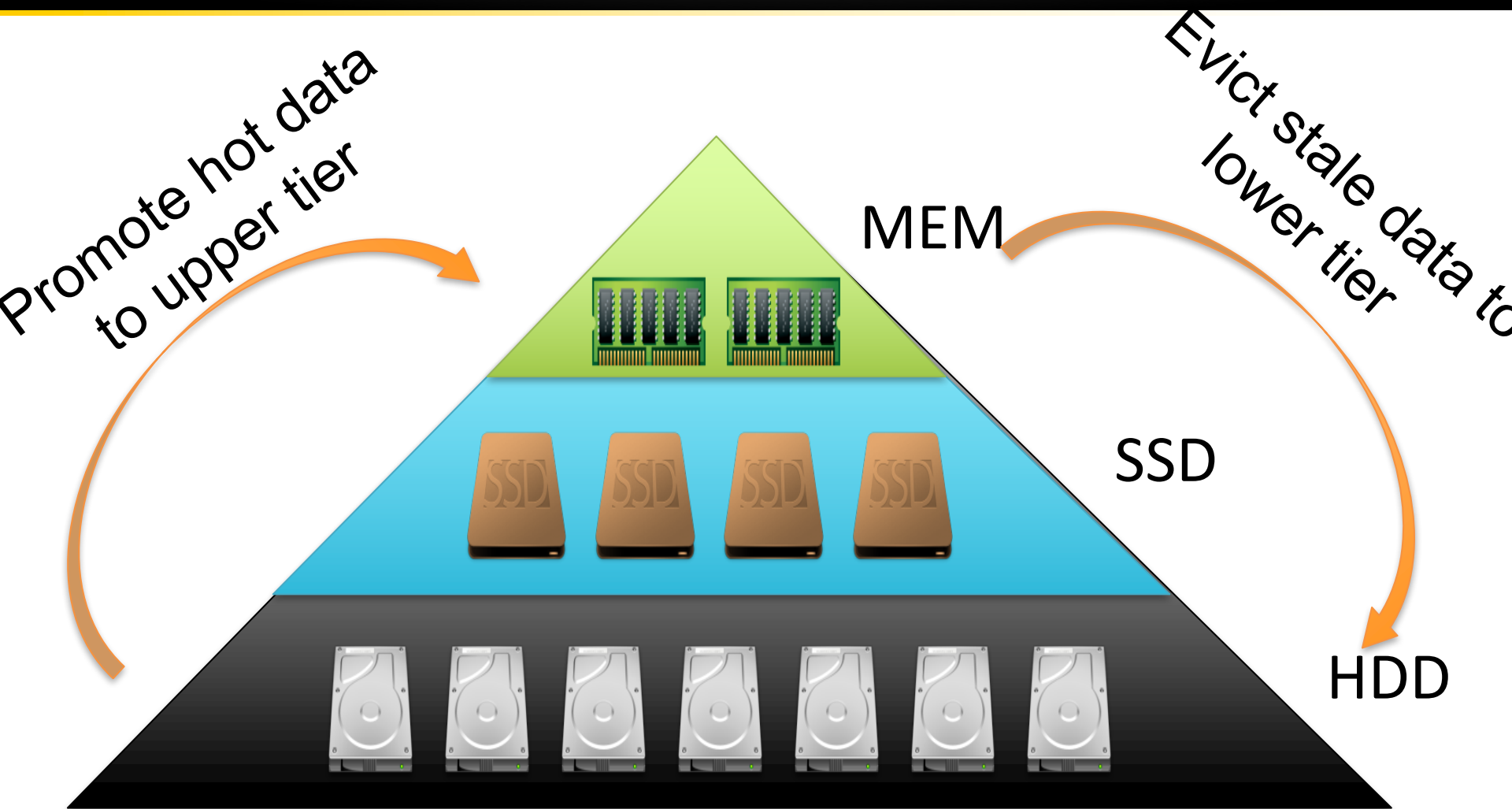
4

总结

# Alluxio是什么？

- 一个虚拟分布式文件系统；
- 架构在底层分布式文件系统和上层分布式计算框架之间的一个中间件；
- 以文件形式在内存或其它存储设施中提供数据的存取服务。
- 官网链接：<http://www.alluxio.org/>

# Tied Storage



# Alluxio生态系统

Spark

MapReduce

Flink

Zeppelin

HBase

Presto

Native File System  
Interface

Hadoop Compatible  
FS Interface

Native Key-Value  
Interface

FUSE Compatible  
FS Interface



Amazon S3

Swift

Alibaba OSS

HDFS

GlusterFS

NFS

# Case Study: Baidu

## Baidu Queries Data 30 Times Faster with Alluxio

- Application: Spark
- Alluxio Storage: MEM + HDD
- Backend Storage: Baidu's File System
- 200+ nodes deployment, 2PB+ managed space
- Result: Speeding up data querying by 30x

# 主要内容

1

Alluxio

2

Spark JobServer

3

Zeppelin

4

总结



# 什么是Spark JobServer

- “Spark as Service”：针对 job 和 contexts 的各个方面提供了 REST 风格的 api 接口进行管理（用REST提交、杀死以及管理spark作业）
- 通过长期运行的job contexts支持亚秒级别低延迟的任务
- 可以通过结束 context 来停止运行的作业(job)
- 分割 jar 上传步骤以提高 job 的启动速度
- 异步和同步的 job API，其中同步 API 对低延时作业非常有效
- 支持 Standalone Spark 和 Mesos、yarn
- 对RDD或DataFrame对象命名并缓存，通过该名称获取RDD或DataFrame。这样可以提高对象在作业间的共享和重用
- 官网链接：<https://github.com/spark-jobserver/spark-jobserver>

# Spark JobServer Job实例

```
/**
 * A super-simple Spark job example that implements the SparkJob trait and
 * can be submitted to the job server.
 */
object WordCountExample extends SparkJob {
  override def validate(sc: SparkContext, config: Config): SparkJobValidation = {
    Try(config.getString("input.string"))
      .map(x => SparkJobValid)
      .getOrElse(SparkJobInvalid("No input.string"))
  }

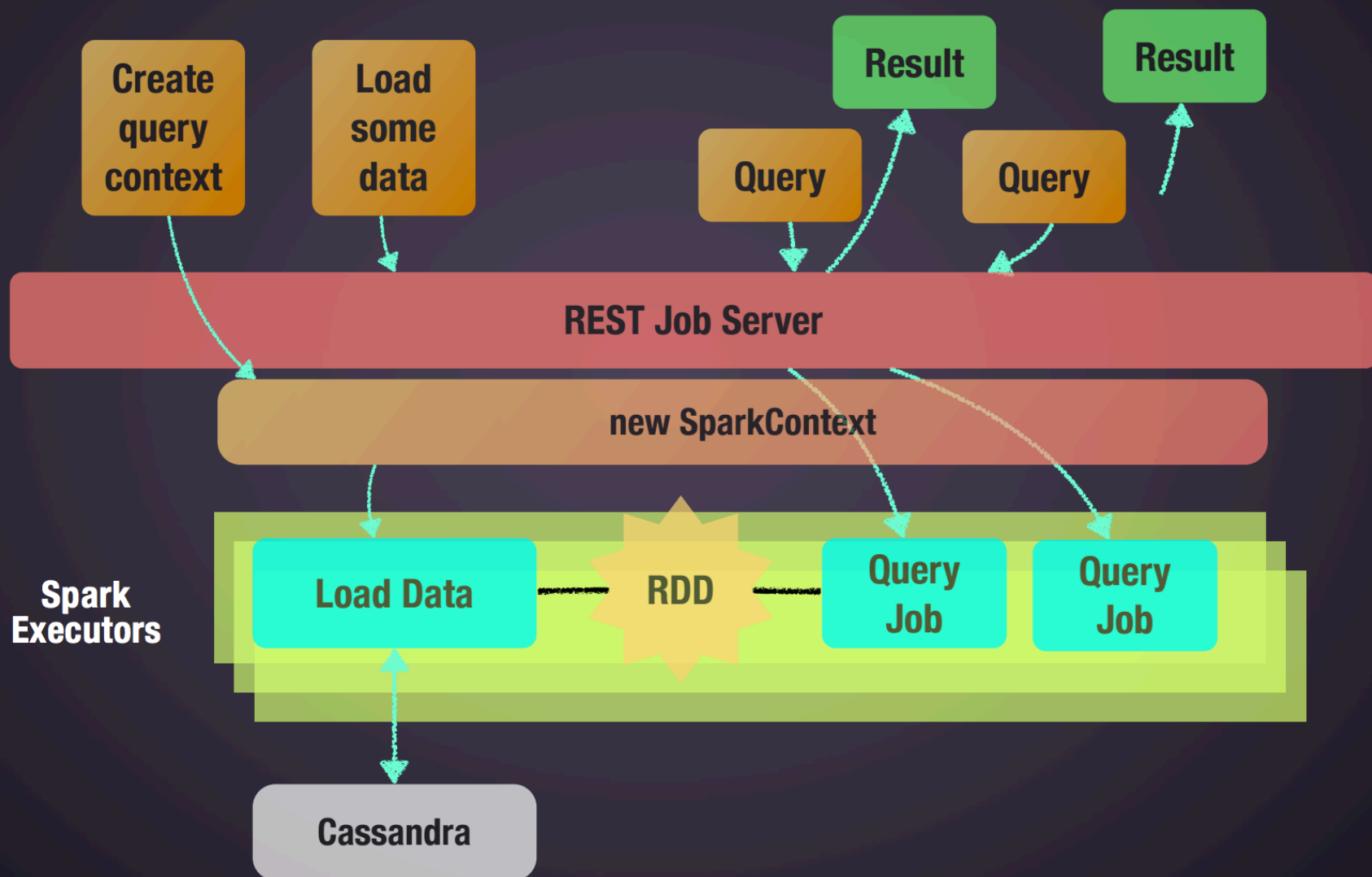
  override def runJob(sc: SparkContext, config: Config): Any = {
    val dd = sc.parallelize(config.getString("input.string").split(" ").toSeq)
    dd.map((_, 1)).reduceByKey(_ + _).collect().toMap
  }
}
```

# Spark JobServer Job提交

```
◆ curl --data-binary @../target/mydemo.jar localhost:8090/jars/demo
OK[11:32 PM] ~

◆ curl -d "input.string = A lazy dog jumped mean dog" 'localhost:8090/jobs?
appName=demo&classPath=WordCountExample&sync=true'
{
  "status": "OK",
  "RESULT": {
    "lazy": 1,
    "jumped": 1,
    "A": 1,
    "mean": 1,
    "dog": 2
  }
}
```

# Spark JobServer加速Job运行



# 主要内容

1

Alluxio

2

Spark JobServer

3

Zeppelin

4

总结

# Zeppelin: 基于web的一站式数据科学平台

- 数据探索和发现
- 官网链接: <http://zeppelin.apache.org/>
- 与Spark和Hadoop深度集成
- 插件化解释器

The screenshot displays the Zeppelin Notebook interface with a blue header bar containing the Zeppelin logo and tabs for 'Notebook' and 'Interpreter'. The main area is divided into three sections:

- Code Editor:** Contains Scala code for loading data from a CSV file into a Spark DataFrame and registering it as a temporary table named 'bank'.

```
val bankText = sc.textFile(s"/user/cloudera/zpdata/bank-full.csv")

case class Bank(age: Integer, job: String, marital: String, education: String, balance: Integer)

val bank = bankText.map(s => s.split(";")).filter(s => s(0) != "").map(
  s => Bank(s(0).toInt,
    s(1).replaceAll("\\s+", ""),
    s(2).replaceAll("\\s+", ""),
    s(3).replaceAll("\\s+", ""),
    s(5).replaceAll("\\s+", "").toInt)
)

bank.toDF().registerTempTable("bank")

bankText: org.apache.spark.rdd.RDD[String] = /user/cloudera/zpdata/bank-full.csv MapPartitionsRDD[55] at textFile at <console>:2
defined class Bank
bank: org.apache.spark.rdd.RDD[Bank] = MapPartitionsRDD[58] at map at <console>:28
Took 1 seconds
```
- SQL Query and Visualization 1:** A query filtering for ages less than 30, visualized as a grouped bar chart.

```
%sql
select age, count(1) value
from bank
where age < 30
group by age
order by age
```
- SQL Query and Visualization 2:** A query filtering for ages less than or equal to a user-defined 'maxAge' (30), visualized as a pie chart.

```
%sql
select age, count(1) value
from bank
where age <= ${maxAge=30}
group by age
order by age
```
- SQL Query and Visualization 3:** A query filtering by marital status, visualized as a table with columns for age and value.

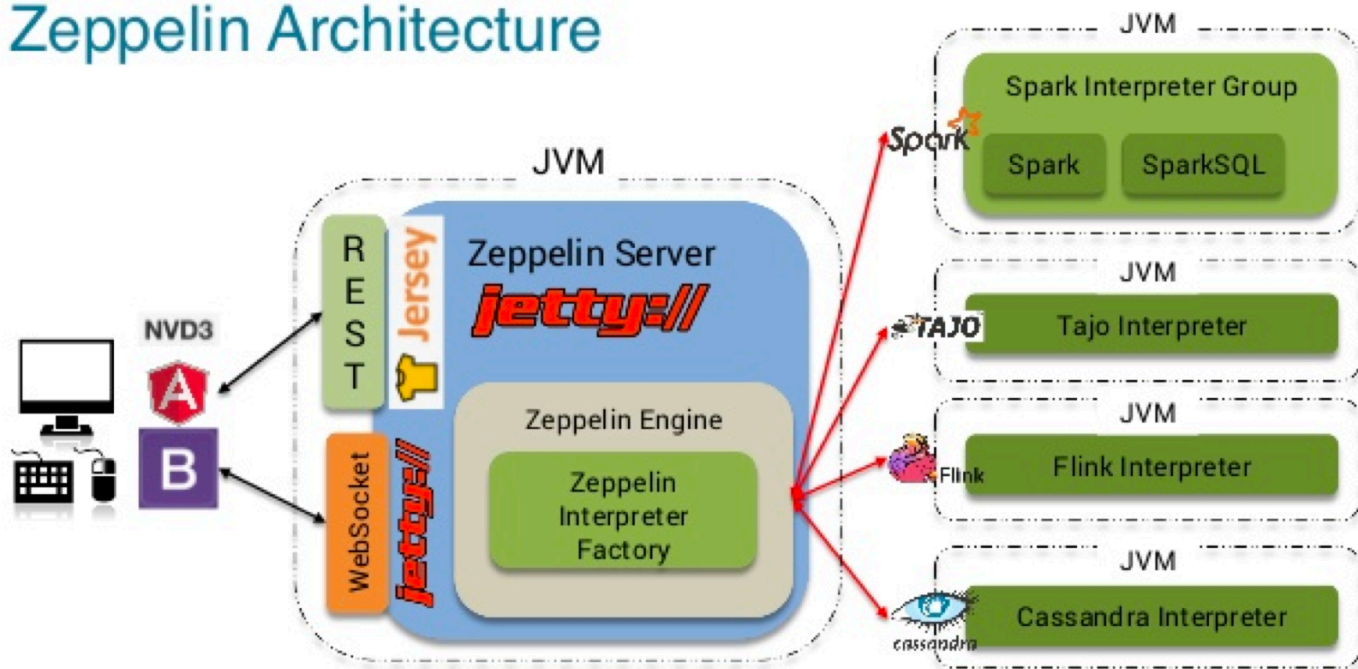
```
%sql
select age, count(1) value
from bank
where marital="${marital=single,singleIddivorcedIdmarried}"
group by age
order by age
```

Annotations with red arrows point to specific features:

- Type here:** Points to the code editor.
- Expose the DataFrame as a SQL Table:** Points to the SQL query section.
- Share notebooks:** Points to the top right corner.
- Use the exposed DataFrame in queries and leverage the built-in visualizations:** Points to the SQL query and visualization sections.

# Zeppelin: 插件化架构

## Zeppelin Architecture



# hadoop123 : 董西成的微信公众号

专注于Hadoop/spark等大数据相关技术的分享







Q&A?

Thank You !