

Spark项目架构与实战

陈超

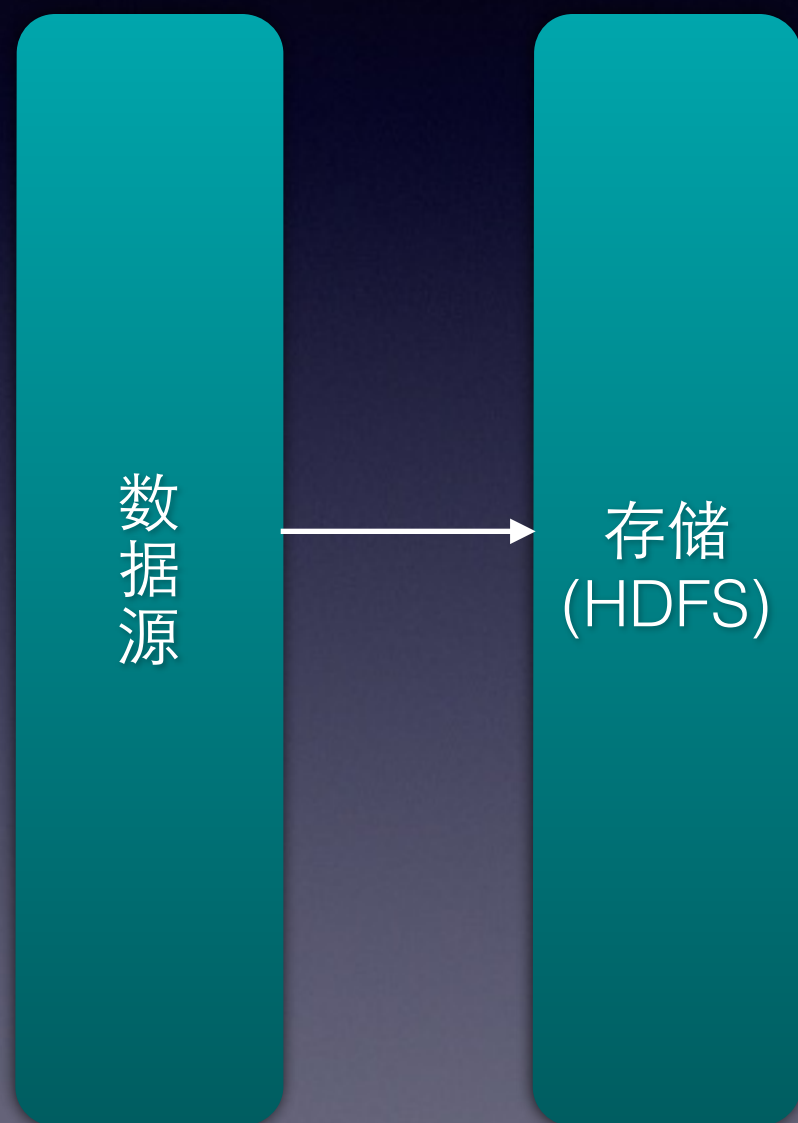
@CrazyJvm



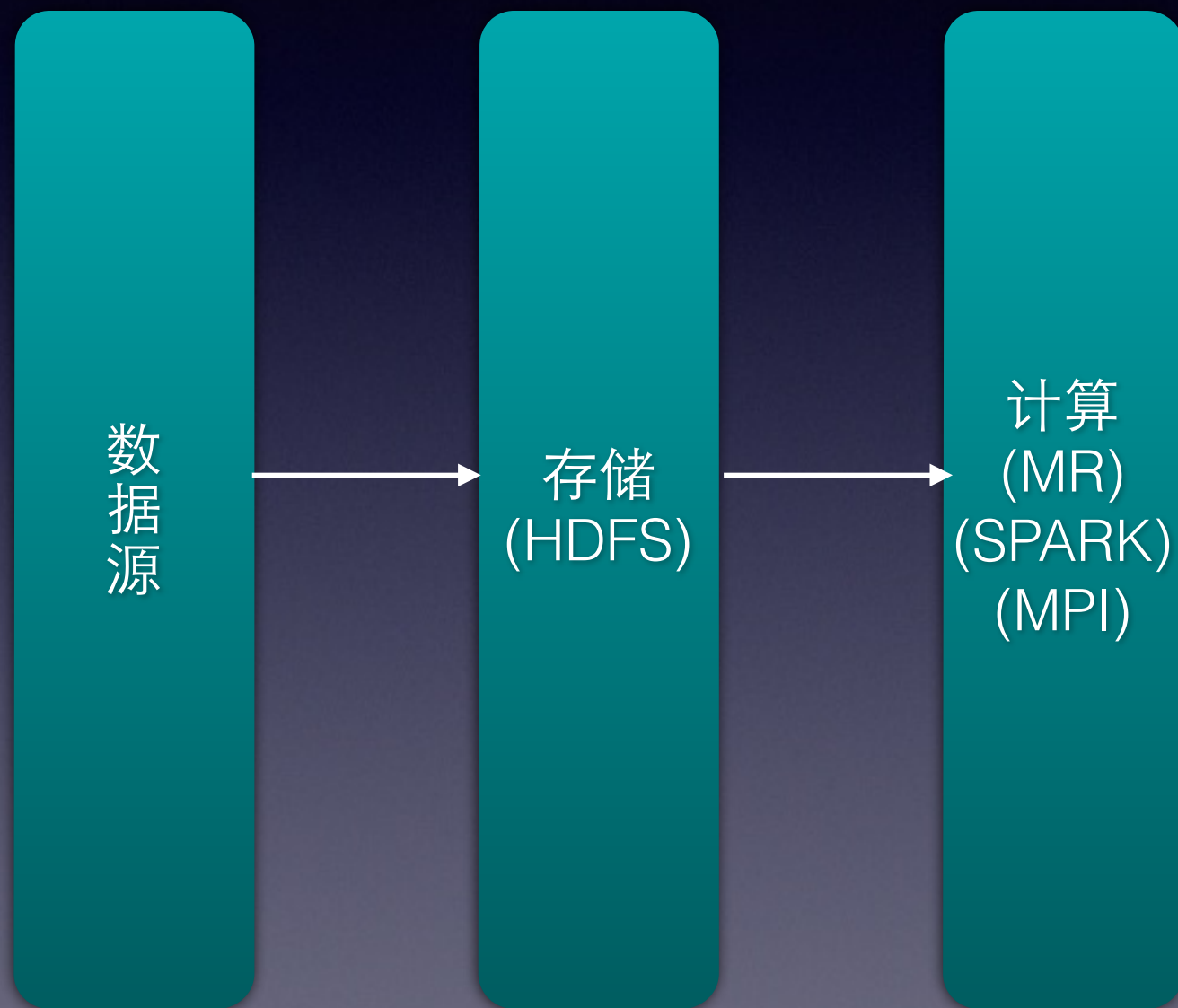
整体流程

数据源

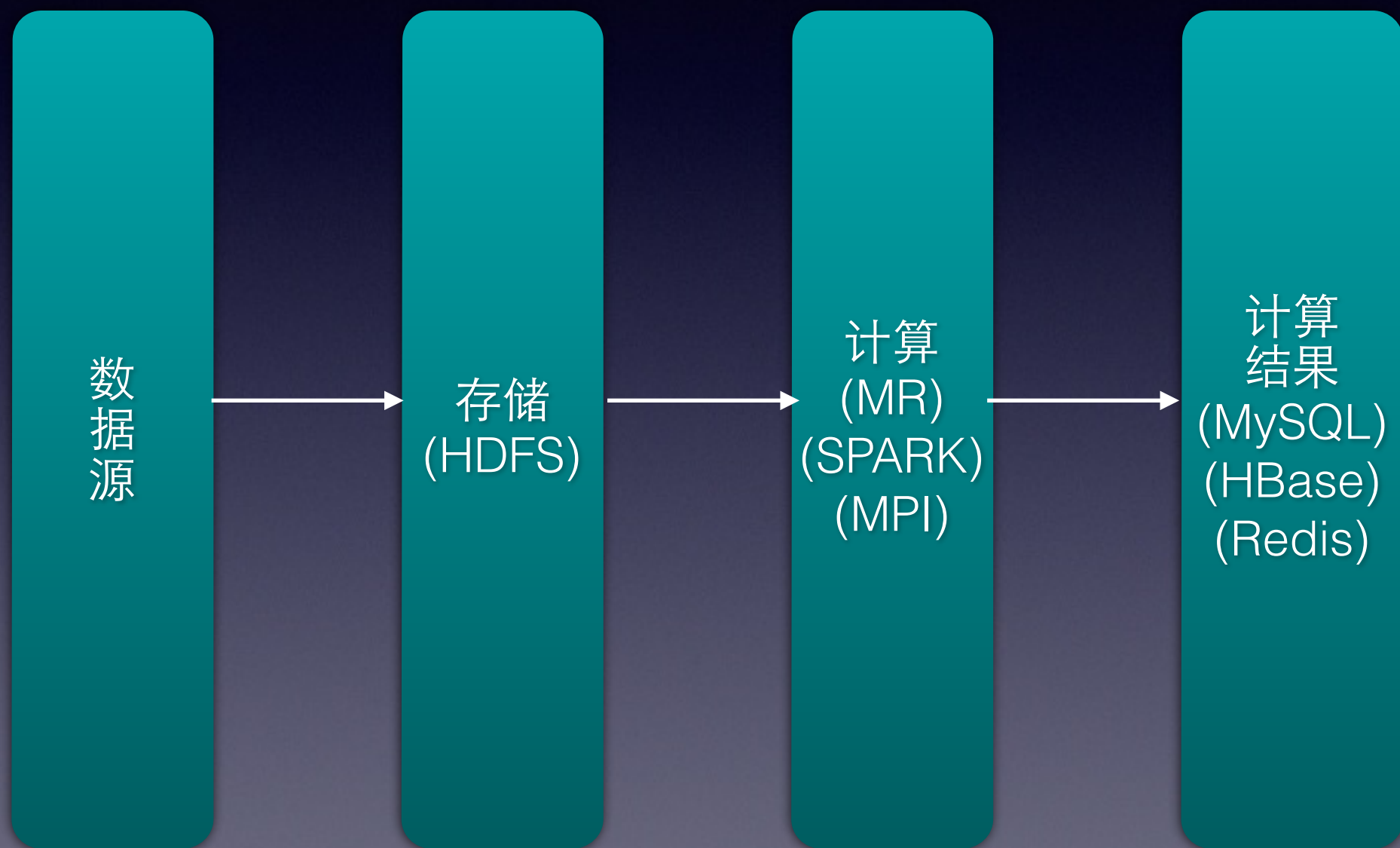
整体流程



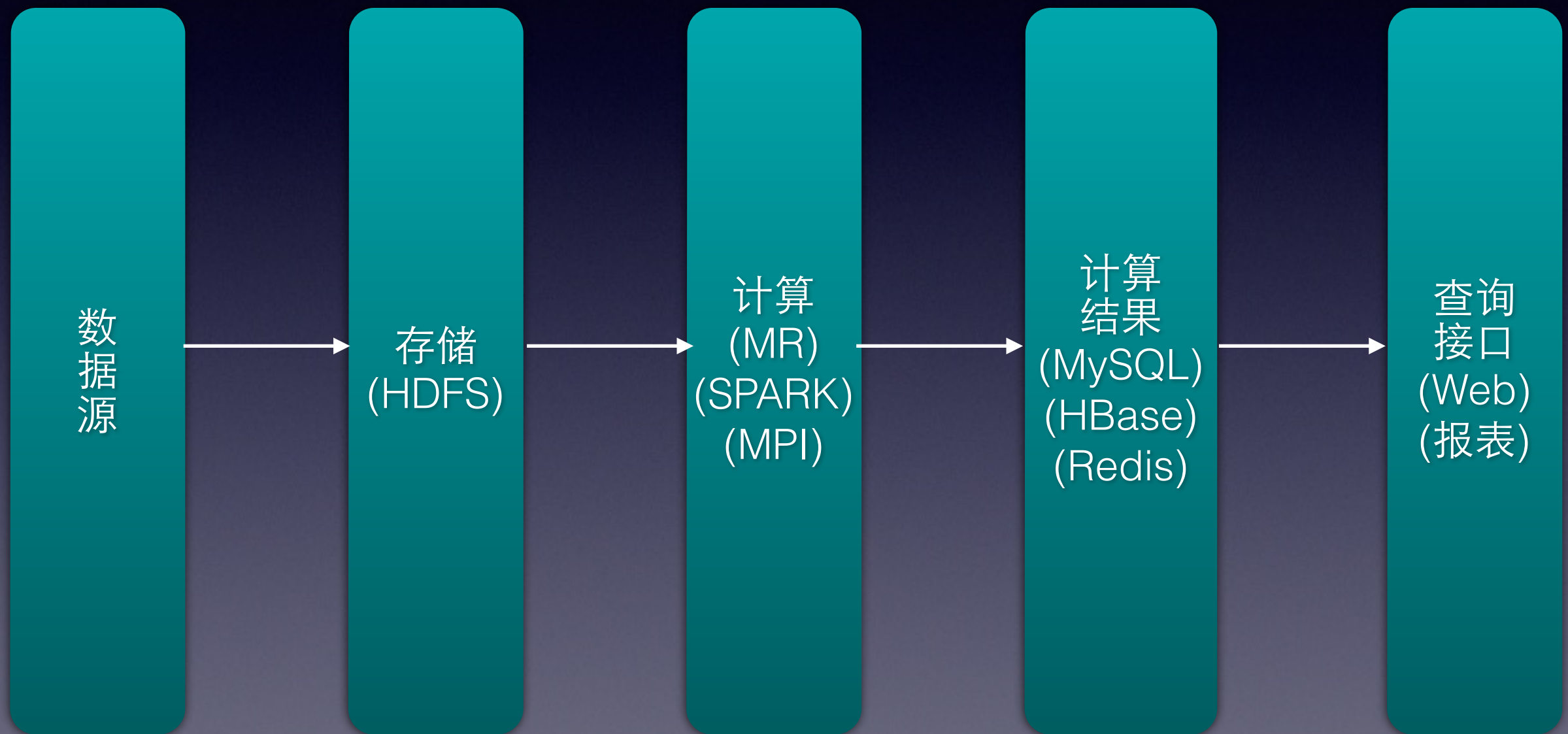
整体流程



整体流程



整体流程



数据源

MySQL(slave)

Oracle(slave)

日志

数据源

MySQL(slave)

Oracle(slave)

日志

数据同步层

实时同步

非实时同步

数据源

MySQL(slave)

Oracle(slave)

日志

数据同步层

实时同步

非实时同步

HDFS

MR Job

Hive Job

Spark Job

Shark Job

流计算

数据源

MySQL(slave)

Oracle(slave)

日志

数据同步层

实时同步

非实时同步

HDFS

MR Job

Hive Job

Spark Job

Shark Job

流计算

计算结果

MySQL

Oracle

HBase

Redis

GemFire

数据源

MySQL(slave)

Oracle(slave)

日志

数据同步层

实时同步

非实时同步

HDFS

MR Job

Hive Job

Spark Job

Shark Job

流计算

结算结果

MySQL

Oracle

HBase

Redis

GemFire

各数据需求方

产品1

产品2

产品3

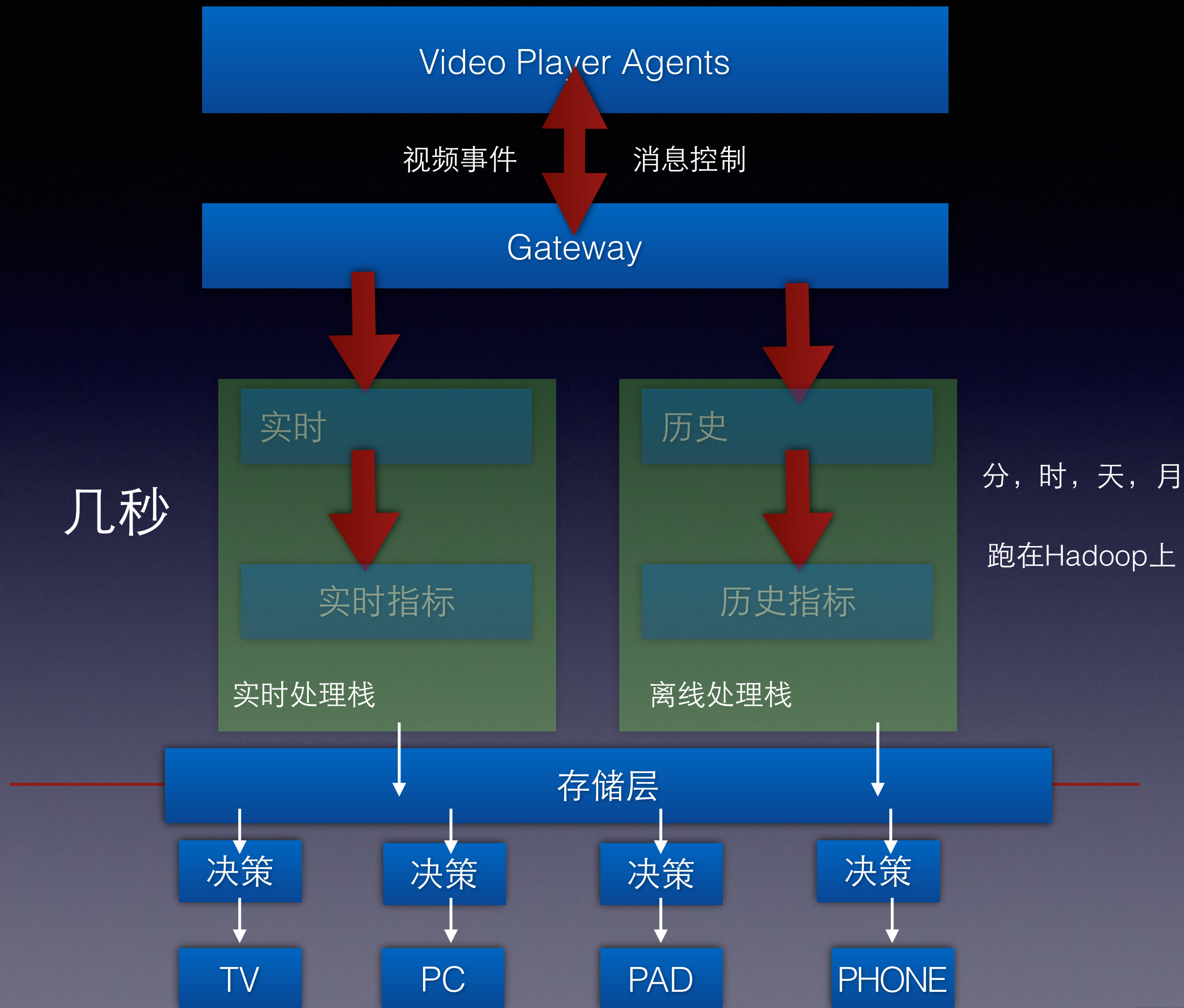
.....

产品N

项目实例

- 看看在不同行业大家都是怎么用的

Conviva



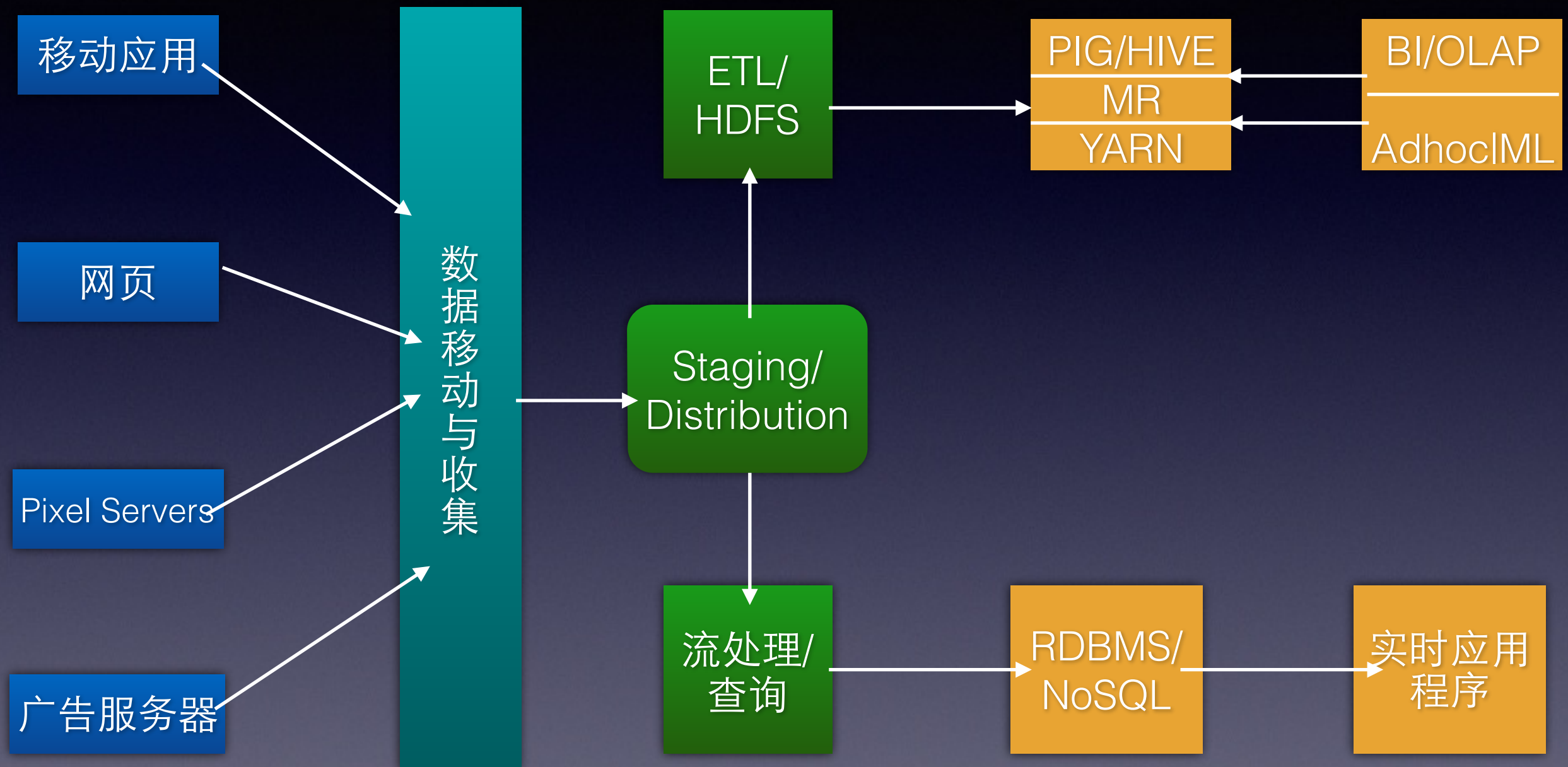
good enough?

切记!!!

- 看起来很美好!!! BUT, 细节是魔鬼!!!
- NEXT ——> 全部由Spark(& Spark Streaming)来完成

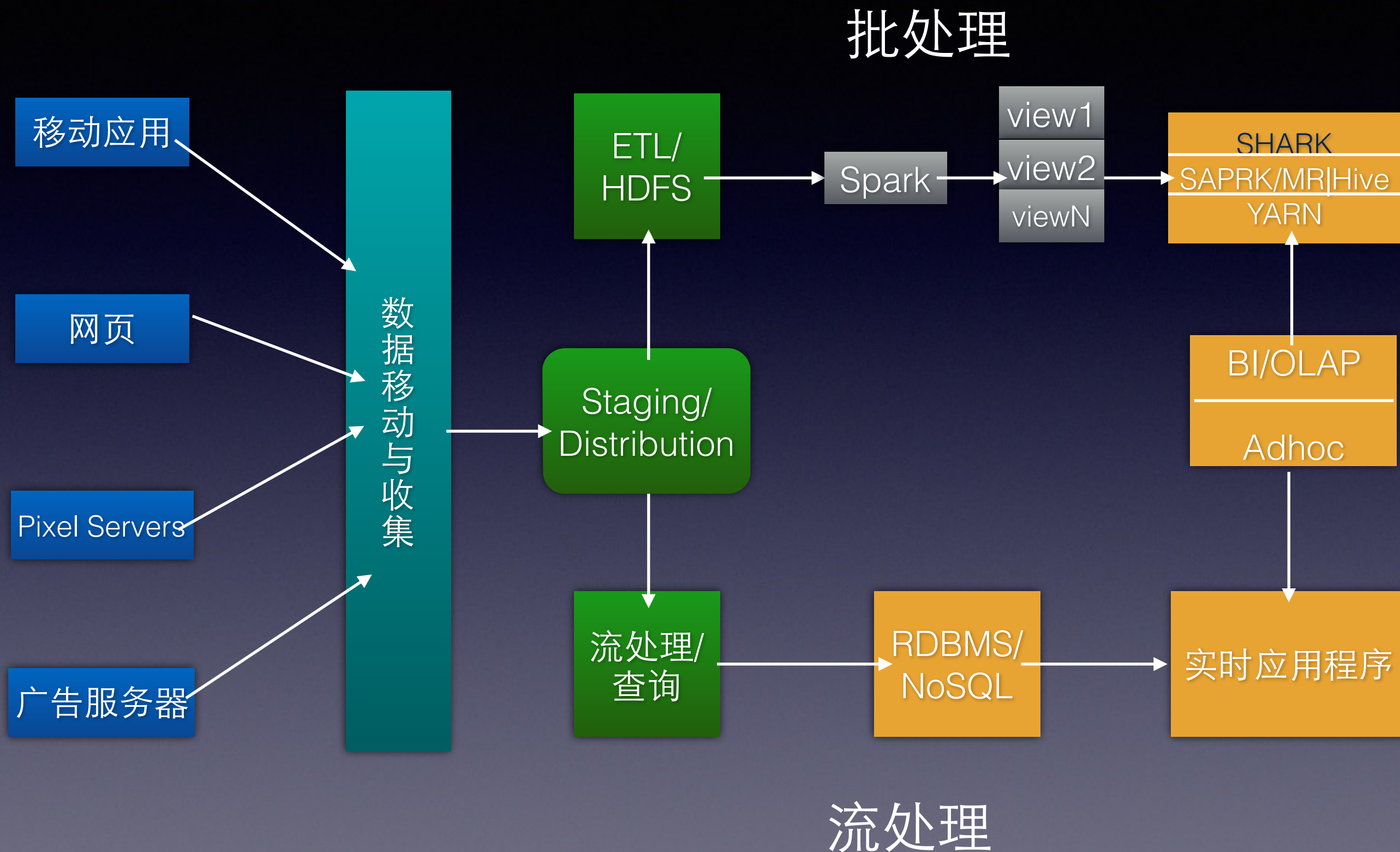
Yahoo!

批处理



流处理

- 慢！ 慢！ 慢！ 解决方案你应该已经想到~



- 很自然的替换/合作!

- 谢谢！ ~~~~~或许，还需要多说几句

- Hadoop
- SPARK

流处理

- Storm
- Spark Streaming
- S4

DB

- Cassandra
- HBase
- MongoDB
- Terrastore
- Redis
- SSDB
- MySQL

SQL on Hadoop(Spark)

- Hive
- Shark(Catalyst)
- Impala

日志(数据)收集

- Sqoop
- Flume
- Chukwa
- Kafka
- DataX
- Dbsync
- TimeTunnel

ML

- Mahout
- MLlib

Others

- Zookeeper
- Oozie
- Hue

- 没有最好的架构，只有最合适的架构！ ！ ！

谢 谢

欢迎关注我的公众微信号，经常会聊`Spark`相关的话题

ChinaScala

