

FusionInsight应用开 发总指导

www.huawei.com



目标

- 学完本课程后，您将能够：
 - 理解**FusionInsight**大数据平台
 - 理解**FusionInsight**应用解决方案
 - 掌握**FusionInsight**应用开发流程
 - 了解方案应用案例



目录

1. 大数据整体介绍

- 大数据平台介绍
- 大数据应用方案和组件介绍

2. 大数据应用开发

3. 经典案例介绍

大数据平台

FusionInsight是华为企业级大数据存储、查询、分析的统一平台，能够帮助企业快速构建海量数据信息处理系统，通过对海量信息数据实时与非实时的分析挖掘，发现全新价值点和企业商机。



目录

1. 大数据整体介绍

- 大数据平台介绍
- 大数据应用方案和组件介绍

2. 大数据应用开发

3. 经典案例介绍

FusionInsight架构：分层解耦开放



FusionInsight 大数据平台



•敏捷

- 完全开放的架构，性能线性扩展
- 强大的**SQL**能力，业务移植便捷
- 丰富的工具支持，开发运维高效

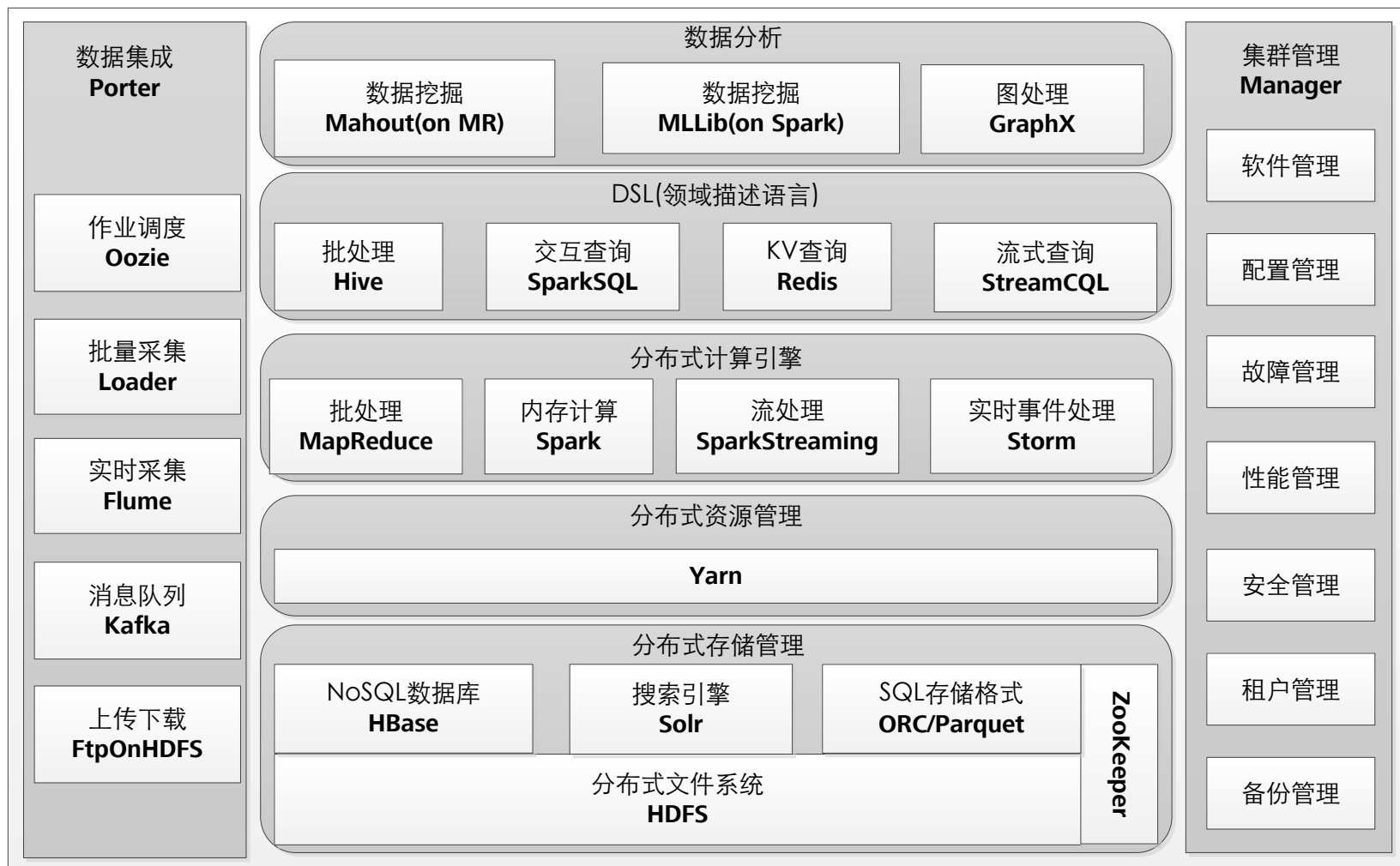
•智慧

- 全量建模，深刻洞察
- 自研算法，高效精准

•可信

- 全组件**HA**、异地容灾
- 开放共赢，可信赖的合作伙伴

FusionInsight HD解决方案架构





本章小结

- **FusionInsight HD** 是一个大数据全栈商用平台，支持各种通用大数据应用场景。



目录

1. 大数据整体介绍
2. 大数据应用开发
 - 技能需求
 - 应用开发流程和介绍
3. 经典案例介绍

技能需求

扎实编程基础

1. 具备**JAVA/Scala**编程能力，熟悉**SQL**
2. 熟悉**JAVA/Scala**编程开发涉及的调试/部署/问题定位等技巧
3. 熟悉和**Linux**常规操作，例如**Shell**命令等

掌握FusionInsight

1. 经过**FusionInsight** 管理员培训及其认证
2. 熟悉大数据常用案例
3. 掌握大数据应用开发关键点，例如基本安全认证和业务权限管理；

熟悉业务开发

1. 理解研发开发流程
2. 理解本应用业务背景

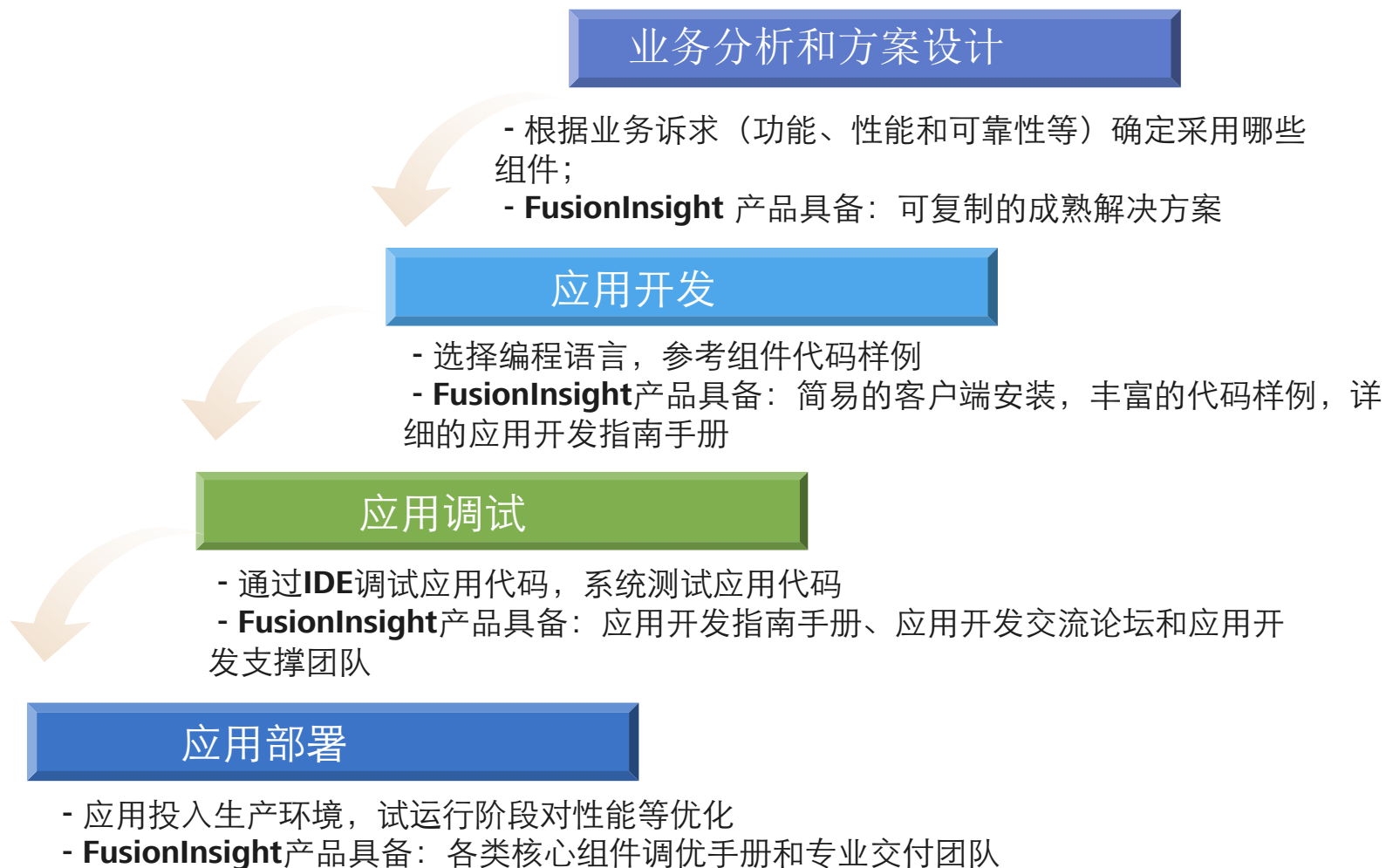
怎样才能做好大数据应用开发



目录

1. 大数据整体介绍
2. 大数据应用开发
 - 技能需求
 - 应用开发流程和介绍
3. 经典案例介绍

大数据应用开发流程



成熟解决方案简介

金融领域

历史明细: **Loader** (数据导入) + **HBase** (数据存储) + **Phoenix** (SQL查询)
在线日志查询: **Flume** (实时采集) + **Kafka** (消息队列) + **Streaming** (流处理)
+ **Redis** (分布式缓存) + **CQL** (日志实时分析) + **HBase/Phoenix** (结果数据存储) + **Solr** (日志搜索、故障定位)

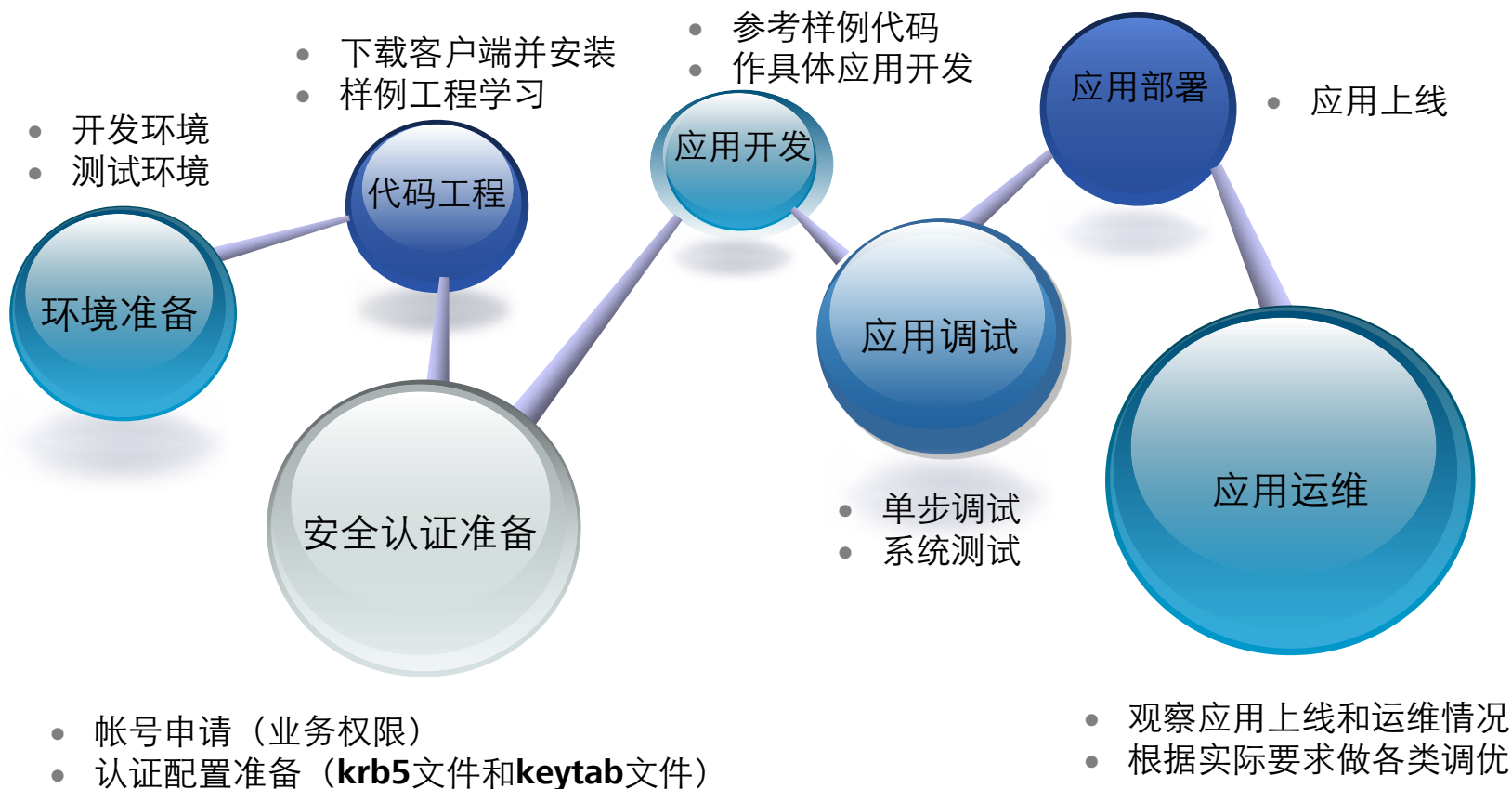
运营商领域

清单查询: **Flume** (实时采集) + **Bulkload** (定时采集) + **HBase** 多实例 (数据存储, 多租户)
实时处理 (位置提醒/实时营销): **Kafka** (消息队列) + **Streaming** (流处理)

公安领域

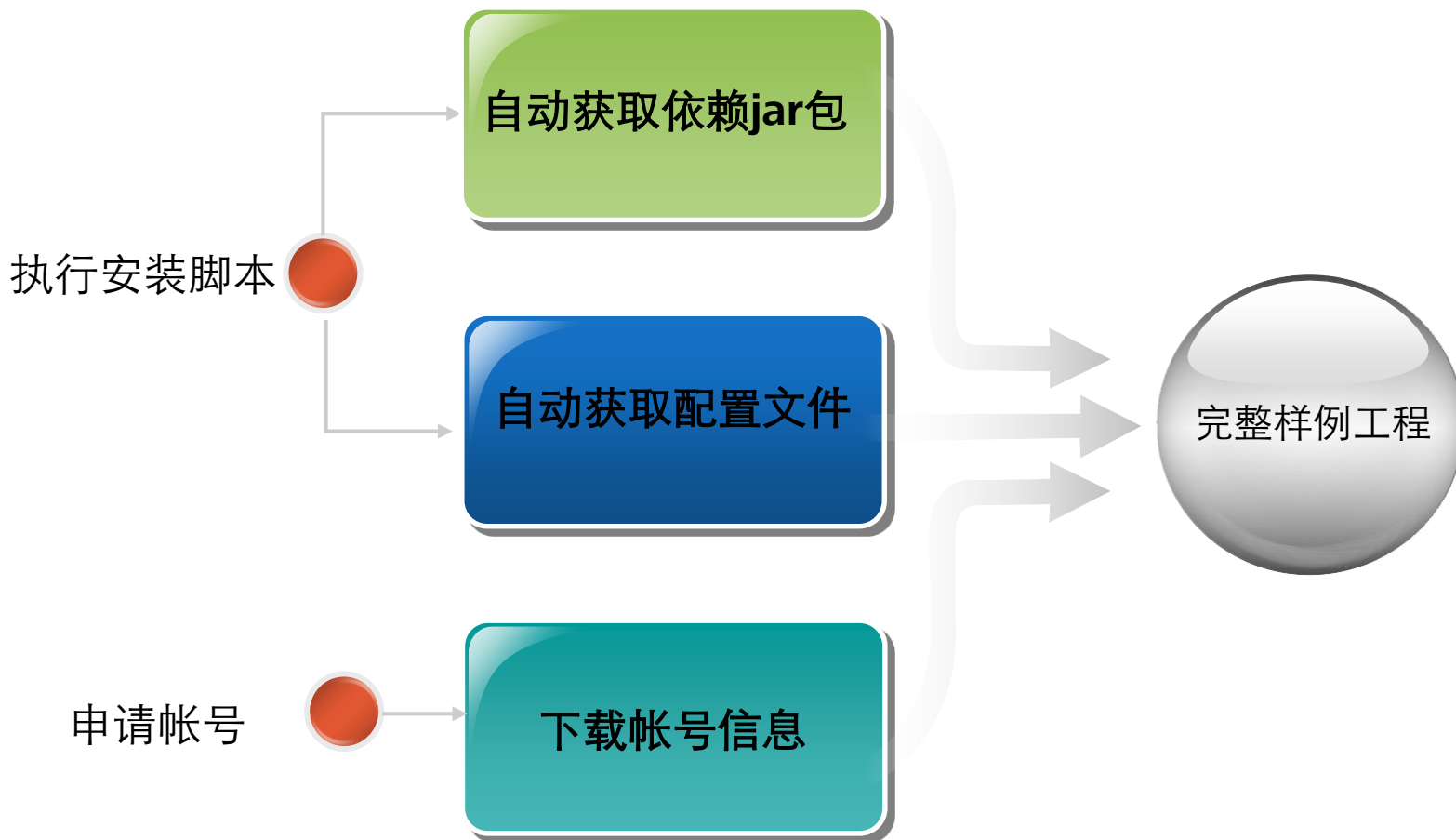
实时监控方案 (套牌车/黑牌车): **Kafka** (消息队列) + **Streaming** (流处理)
+ **Redis** (分布式缓存) + **HBase** (数据存储)
全文检索方案: **Loader** (数据导入) + **HBase** (数据存储) + **Solr** (搜索引擎)

应用开发指南--过程



注意：各个组件可能有差异，细节请参考FusionInsight产品文档的应用开发指南

应用开发指南—代码工程



应用开发指南--关键点



帐号

1. 业务应用请申请业务应用帐号，切忌使用原有系统帐号；
2. 分清人机帐号和机机帐号，切忌混用；
3. 新帐号权限遵从最小原则。



安全认证

1. 一个进程里只用一个帐号。
2. 一个进程里只需显式认证一次。



场景约束

如下场景都需遵从上述“安全认证”的要求（设计环节关注）

1. 一个进程里需同时访问多个集群（各自有独立的**KrbServer**和**LdapServer**）；
2. 在容器（例如**Tomcat**）上运行**app**，一个容器里所有**app**属于同一个进程。

FusionInsight 应用开发指南--代码讲解1

开源社区各个组件认证接口不一致

没有对多次登录认证作保护

jaas.conf配置内容跟随**JDK**类型变化

- 1.统一认证接口
- 2.优化**ZooKeeper**配置

应用开发指南--代码讲解2

1.初始化;
2.调用统一认证API

```
public static void main(String[] args) throws Exception{
    init();
    LoginUtil.login(USER_NAME, USER_KEYTAB_FILE, KRB5_FILE, CONF);

    //test HBase
    HBaseSample ts = new HBaseSample(CONF);
    ts.testCreateTable();
    ts.dropTable();
    System.out.println("-----finish HBase -----");
}
```

业务配置初始化

```
public static void init() throws Exception{
    String userdir2 = System.getProperty("user.dir") + File.separator
        + "conf2" + File.separator;
    CONF = new Configuration();
    CONF.addResource(new Path(userdir2 + "core-site.xml"));
    CONF.addResource(new Path(userdir2 + "hdfs-site.xml"));
    CONF.addResource(new Path(userdir2 + "hbase-site.xml"));

    USER_NAME = "tester1";
    USER_KEYTAB_FILE = userdir2 + "user.keytab";
    KRB5_FILE = userdir2 + "krb5.conf";
}
```

Zookeeper初始化

```
/*
 * if need to connect zk, please provide jaas info about zk.
 * of course, you can do it as below:
 * System.setProperty("java.security.auth.login.config", confDirPath + "jaas.conf");
 * but the demo can help you more :
 * Note: if this process will connect more than one zk cluster, the demo may be not proper. you can contact us for
 */
LoginUtil.setJaasConf(ZOOKEEPER_DEFAULT_LOGIN_CONTEXT_NAME, USER_NAME, USER_KEYTAB_FILE);
LoginUtil.setZookeeperServerPrincipal(ZOOKEEPER_SERVER_PRINCIPAL_KEY, ZOOKEEPER_DEFAULT_SERVER_PRINCIPAL);
```

应用开发指南--代码讲解3

统一认证接口



LoginUtil类关键部分

```
public synchronized static void login(String userPrincipal, String userKeytabPat
    throws IOException
{
    // 1.check input parameters
    if ((userPrincipal == null) || (userPrincipal.length() <= 0))
    {
```

判断是否重复登录



```
private static boolean checkNeedLogin(String principal)
    throws IOException
{
    if (!UserGroupInformation.isSecurityEnabled())
    {
        LOG.error("UserGroupInformation is not SecurityEnabled")
        throw new IOException(
```

应用开发指南--代码讲解3

Zookeeper相关:

- 1.环境变量设置
- 2.提供内部类, 生成jaas配置对象

LoginUtil类关键部分

```
public static void setJaasConf(String loginContextName, String principal, String
    throws IOException
{
    public static void setZookeeperServerPrincipal(String zkServerPrincipal
        throws IOException
    {
        System.setProperty(zkServerPrincipalKey, zkServerPrincipal);
    }
    private static class JaasConfiguration extends javax.security.auth.login.Cc
    {
        private static final Map<String, String> BASIC_JAAS_OPTIONS = new HashM
        static
        {
```

应用开发指南—调试

常规手段

1. 利用**JAVA**和**IDE**所带功能完成本地/远程单步调试；
2. 部分组件应用调试需加辅助手段进行远程单步调试（如**Spark/MapReduce/Storm**）；

协助资料

1. 遇到问题，可参考**FAQ**等资料做初步分析；
http://support.huawei.com/ecomunity/bbs/list_1069,1420018021.html?dist=1
2. 交流平台：
http://support.huawei.com/ecomunity/bbs/list_1069,1420018021.html

保障团队

1. 疑难杂症可向**FusionInsight** 交付团队求助；

应用开发指南—部署运行

应用程序运行方式：

- 借助**FusionInsight** 客户端命令：例如，提交**Spark**应用时，在**FusionInsight**客户端下执行命令：**spark-submit**。
- 借助**JDK**自带命令：例如，启动一个**HBase**应用程序，此程序根据**JDK**规范被打包为一个可执行**jar**包，以如下命令启动：**java -jar** 。
- 借助第三方平台：例如，应用被部署在**Tomcat**容器或**OSGI**平台上。

应用开发指南—观察运行结果

观察应用程序运行结果的方式：

- 借助**FusionInsight** 客户端命令：例如，启动一个**HBase**应用程序，往某个**HBase**表插入一条记录，可以通过**HBase**的**shell**操作窗口查看（执行**get**命令）此记录是否插入成功。
- 注：需先执行**kinit**命令完成**kerberos**认证。**Kinit**命令支持人机账号和机机账号。例如；
 - 人机账号：**kinit tester1** （执行此命令后，会弹出输入密码的界面）
 - 机机账号：**kinit -kt /opt/tester1.keytab tester1** （执行此命令后，进入业务操作界面）
- 借助应用程序运行返回结果：例如应用程序打印执行结果信息，或通过单步调试跟踪运行结果。
- 借助组件原生态**Web**页面：例如使用**YARN Web**页面查看**YARN**应用执行结果，使用**HUE Web**页面查看**HDFS**组件的数据存储情况。



思考题

若存在如下场景，有什么解决方法？

- 问题：刚部署了一个安全模式的**FusionInsight**集群，需在一个**Tomcat**容器上部署两个**app**，一个**app**访问此集群的**HBase**组件，另外一个**app**访问此集群的**Hive**组件。这种场景下，**kerberos**认证代码逻辑该如何处理？需放在哪个**app**上？或者你是否有更好的解决办法？



本章小结

- 认证是应用开发的关键点，要根据业务需要，申请合适帐号，完成安全认证。
- 华为**FusionInsight**易集成开发，提供了包括应用开发指南、样例代码和支撑团队等各类协助。



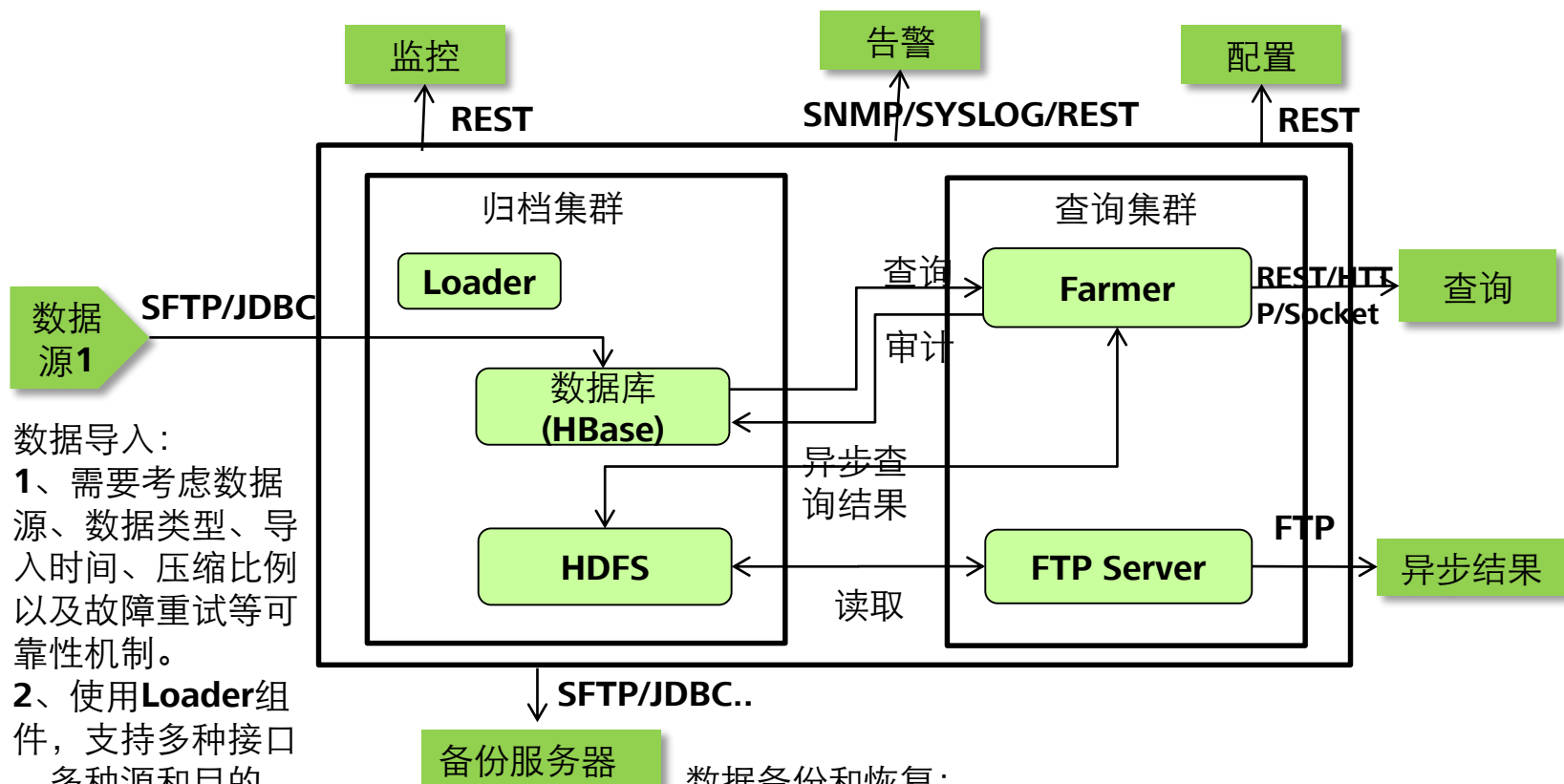
目录

1. 大数据整体介绍
2. 大数据应用开发
3. 经典案例介绍

金融领域（历史明细）

	历史明细查询
客户问题	<p>1、银行历史数据存放分为：主机（1~2年内）、数据仓库（5年内）、归档库（5年以上），未集中，比较散。</p> <p>2、数据仓库和归档库，因性能压力，原则上不对外提供实时查询服务。</p>
场景（要求）	<p>1、内部客户历史查询：比如柜员等发起对用户的历史明细查询。</p> <p>2、互联网用户的历史查询：借记卡历史明细查询、信用卡历史明细查询等。</p>
解决方案 案例	<p>特点：离线历史数据查询类业务，结构化TB级数据大融合，构建一个超大集群，提供毫秒级~几秒级查询业务：</p> <p>Loader（数据导入）+ HBase（数据存储）+ Phoenix（SQL查询）+ Farmer（数据服务）</p>

金融领域（历史明细）



数据导入：

1、需要考虑数据源、数据类型、导入时间、压缩比例以及故障重试等可靠性机制。

2、使用Loader组件，支持多种接口，多种源和目的，支持定时调度，无需二次开发。

数据备份和恢复：

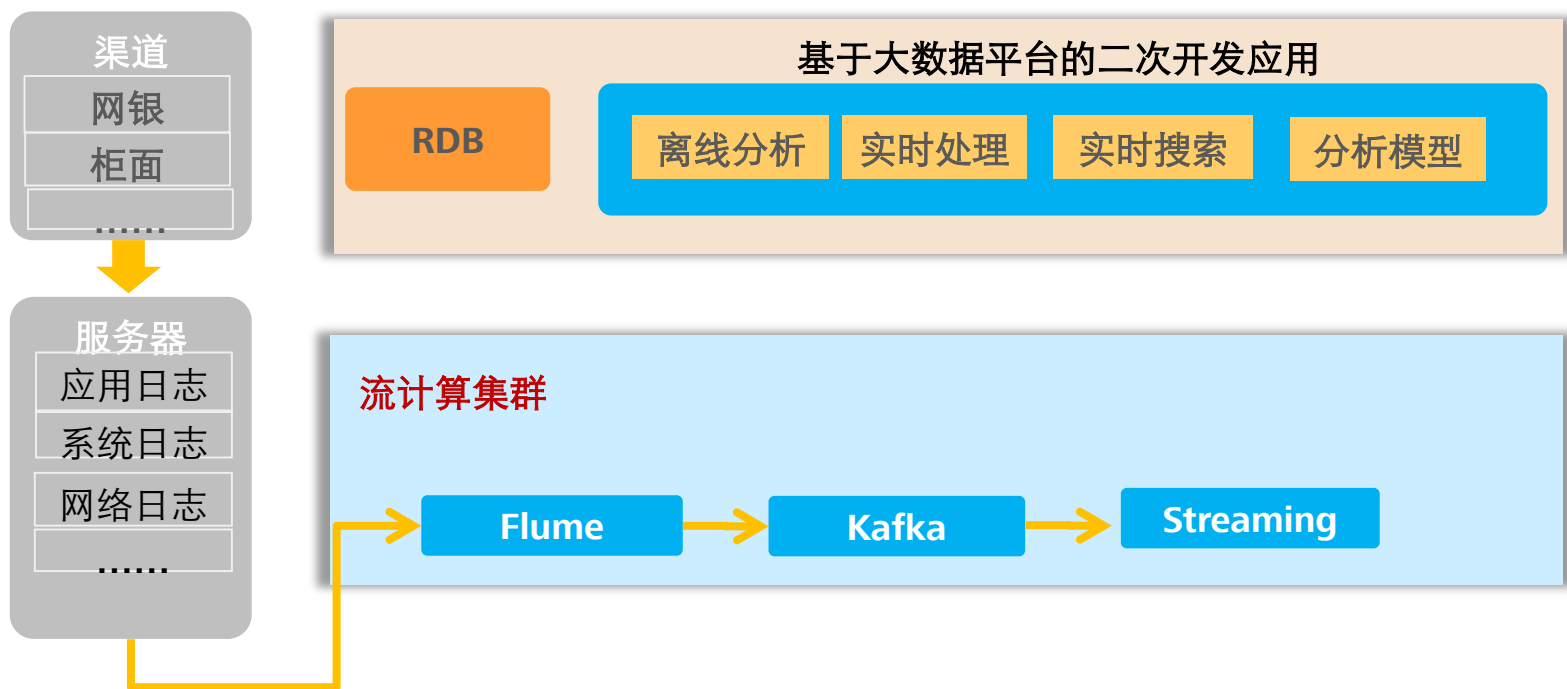
1、需要和客户指定备份恢复策略，包括：备份时间，恢复时间，不同故障如何恢复等

2、使用Loader组件，无需二次开发

金融领域（在线日志分析）

	在线日志分析
客户问题	<p>1、银行系统当前大部分仅对告警部分进行实时处理和上报监控，缺乏对监控KPI、运行日志的实时处理。</p> <p>2、银行系统因磁盘空间有限，当前运维数据仅能存放几个月，过期全部删除掉，无法利用历史数据进行故障风险，以及故障预测等职能运维方面的事情。</p>
场景（要求）	<p>1、对现有应用系统日志进行实时采集和分析，为TOP应用提供秒级的监控能力，包括交易情况、资源消耗情况。</p> <p>2、针对故障或者异常情况，提供实时日志搜索能力。</p> <p>3、针对历史日志，通过数据挖掘和分析，实现故障预测的能力。</p>
解决方案案例	<p>特点：采用流式框架，对日志进行实时采集、分析和展现；未来扩展到数据预测。</p> <p>方案实现：Flume（实时采集） + kafka（消息队列） + Streaming（流处理）。</p>

金融领域（在线日志分析）



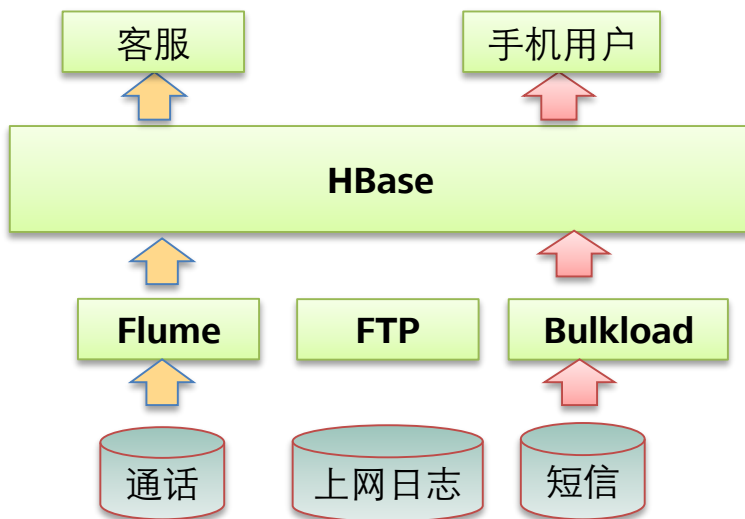
运营商领域（清单查询）

场景介绍

基于**HBase**。对通话清单和上网日志进行查询，主要需求有：

- 1、高效的**Key-Value**查询
- 2、支持灵活的数据加载方式
- 3、支持简单的数据预处理
- 4、支持**HBase**多实例

方案示意图



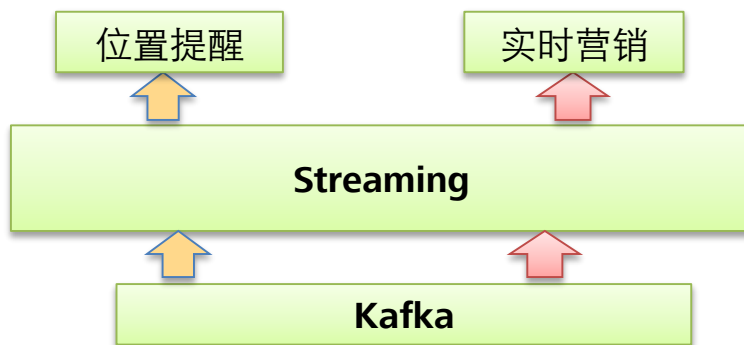
运营商领域（实时处理）

场景介绍

数据源将实时数据丢入**Kafka**，**Streaming**实时订阅，触发实时应用，此场景在运营商很常见。主要需求有：

- 1、**Kafka**支持权限管理
- 2、**Streaming**支持多租户和资源隔离
- 3、支持多种数据进入**Kafka**的方式
- 4、实时处理性能

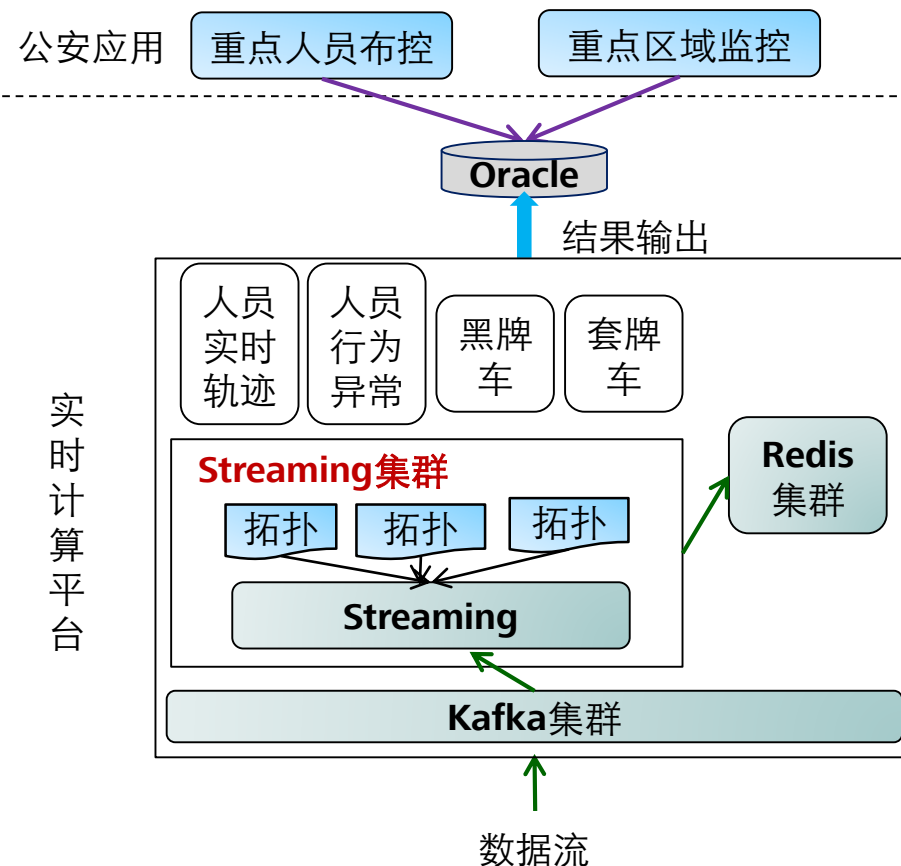
方案示意图



公安领域（实时监控）

	实时监控
客户问题	XX 市局为例，每天 2000W 卡口记录数据流入，数据量大。 Oracle 数据库无法支持实时流处理。
场景（要求）	卡口监控实时提取车辆信息，如车牌，车型，颜色等。海量信息汇入，发现套牌车、黑牌车时实时报警，辅助警员及时应对。
解决方案 案例	Kafka （消息队列） + Streaming （流处理） + Redis （分布式缓存） + HBase （数据存储） / Oracle 。

公安领域（实时监控）



数据源: 主要包含视频监控数据, 卡口拍照数据, 通信话单, 及上网记录。

Kafka集群: 消息系统。

Streaming集群: 分布式流处理引擎。

Redis集群: 提供高速key/value存储查询能力, 用于流处理中间数据的缓存。

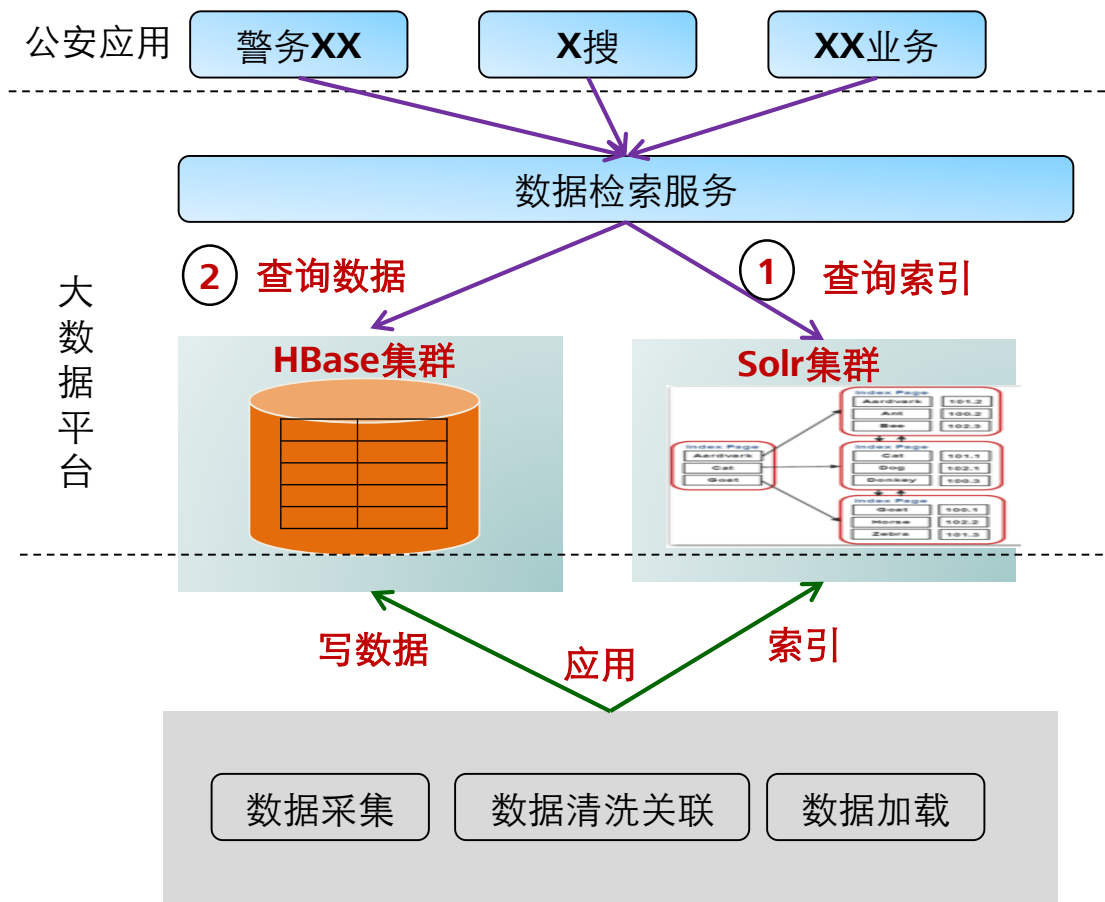
Oracle: 传统关系数据库, 流处理结果由拓扑写入。

公安应用: 直接面向公安客户的业务, 针对实时性要求较高的诉求, 如重点人员布控、重点区域监控。

公安领域（全文检索）

	实时监控
客户问题	现状是数据分散在多个系统中，无法集中存放和查询。 Oracle 数据库管理百亿数据需分表分库，维护麻烦。 Oracle 数据库性能无法支撑百亿记录秒级查询。
场景（要求）	案件发生时，相关的历史案件和嫌疑人员等信息快速获取，辅助警员快速判断和应对，数百亿记录秒级完成精确查询。海量数据 PB 级以上，查询时间秒级返回；上万条数据 3 秒返回； 110 多个字段， 60 多个字段要查询。
解决方案案例	HBase （数据存储）+ Solr （搜索引擎）。

公安领域（全文检索）





思考题

对于上文提到公安领域的全文检索方案，你有什么优化方案？

提示：

- 1.一份数据同时写往**HBase**和**Solr**这两个不同地方，就会存在常见的数据一致性问题；
- 2.应用要同时访问**HBase**接口和**Solr**接口，对接内容较多。



本章总结

- 在金融、运营商和公安三大领域上积累了各种主流方案。
- 设计方案需要考虑安全认证的影响。例如安全模式下，一个应用同时访问两个集群，那么这两个集群需要做互信。



习题

- 判断题
 1. **Loader**适合做实时数据采集。 (T or F)
 2. 新申请帐号的业务权限越大越好，方便。 (T or F)
 3. **Solr**组件在某种程度上也可以当做是非关系数据库使用。 (T or F)



习题

- 多选题-1

如下这几个分析平台可以运行在**YARN**框架上的有？（ ）

A.Spark

B.MapReducec

C.Streaming



习题

- 多选题-2

如下哪些组件可以对外提供**SQL**接口？（ ）

A.HBase

B.Hive

C.SparkSQL

D.Solr



习题

- 多选题-3

如下哪些组件适用于实时分析？（ ）

A.MapReduce

B.Streaming

C.SparkStreaming

D.Hive



习题

- 多选题-4

Kerberos安全认证必须有下列哪些？（ ）

A.krb5.conf配置文件

B.用户名

C.用户的**keytab**文件

D.jaas.conf配置文件

Thank you

www.huawei.com