# 实时数据流处理

**采用框架：kafka+spark streaming**

1. 环境搭建
   - kafka 环境搭建
   - spark streaming 环境搭建（Mesos）
     - mesos环境搭建，参考文档Mesos.md
     - mesos上部署spark
       1. 下载spark

          ```
          wget http://d3kbcqa49mib13.cloudfront.net/spark-2.0.2-bin-hadoop2.7.tgz
          tar zvxf spark-2.0.2-bin-hadoop2.7.tgz
          ```

       2. 配置

          ```
          cd spark-2.0.2-bin-hadoop2.7/conf
          cat spark-env.sh.template spark-env.sh
          vim spark-env.sh #参照示例，添加内容
          ```

          无hdfs版:

          

          hdfs版：
          修改上述示例的最后一行，即export
          SPARK_EXECUTOR_URI=hdfs://spark.tar.gz

          ```
          cat spark-defaults.conf.template spark-defaults.conf
          vim spark-defaults.conf #添加如下内容
          ```

          spark.io.compression.codec lzf
       3. 打包分发

          ```
          cd ../..
          tar zcvf spark.tar.gz spark-2.0.2-bin-hadoop2.7 #压缩文件名和
          位置与第2步骤SPARK_EXECUTOR_URI一致
          #无hdfs
          scp spark.tar.gz 192.168.125.172:/home/null
          scp spark.tar.gz 192.168.125.173:/home/null
          ```

4. 在master主机运行测试：

```
cd spark-2.0.2-bin-hadoop2.7
bin/spark-shell --master mesos://192.168.125.171:8081
```

效果如下：



2. 上传程序

```
./sbin/start-mesos-dispatcher.sh mesos://192.168.125.171:5050
./bin/spark-submit \
  --class <main-class> \
  --master <master-url> \
  --deploy-mode <deploy-mode> \
  --conf <key>=<value> \
  ... # other options
  <application-jar> \
  [application-arguments]
```

- 客户端模式
  上面选项中的−deploy-mode 指定为client
- 集群模式
  上面选项中的−deploy-mode 指定为cluster
  例子：

```
./bin/spark-submit --class *** --master mesos://192.168.125.171:7077 --d
eploy-mode client ***.jar 1000
```

注意，master地址是MesosClusterDispatcher地址，默认是7077

3.