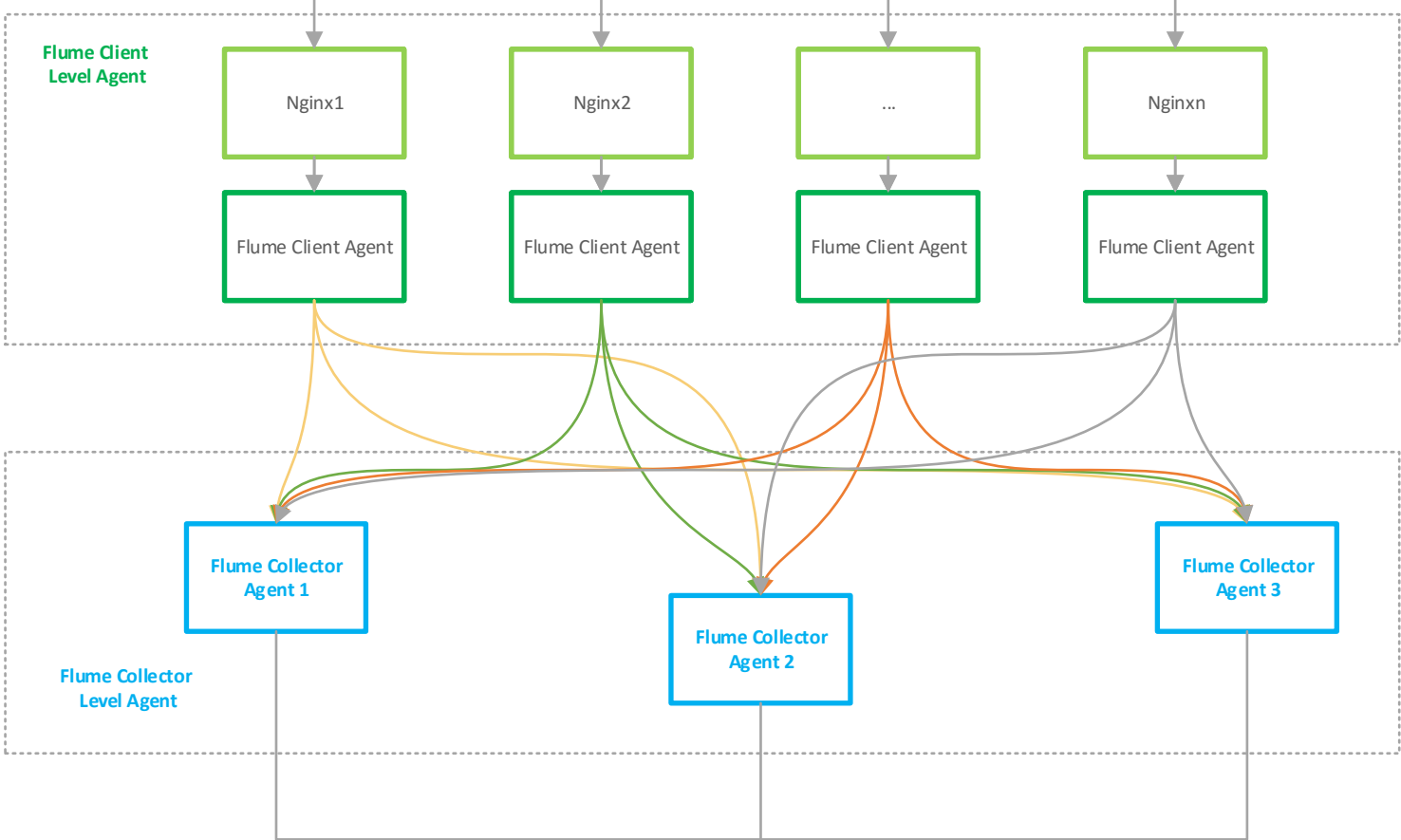




Nginx入口服务器  
(运维负责)

Nginx分流容错



Kafka集群

读取日志数据

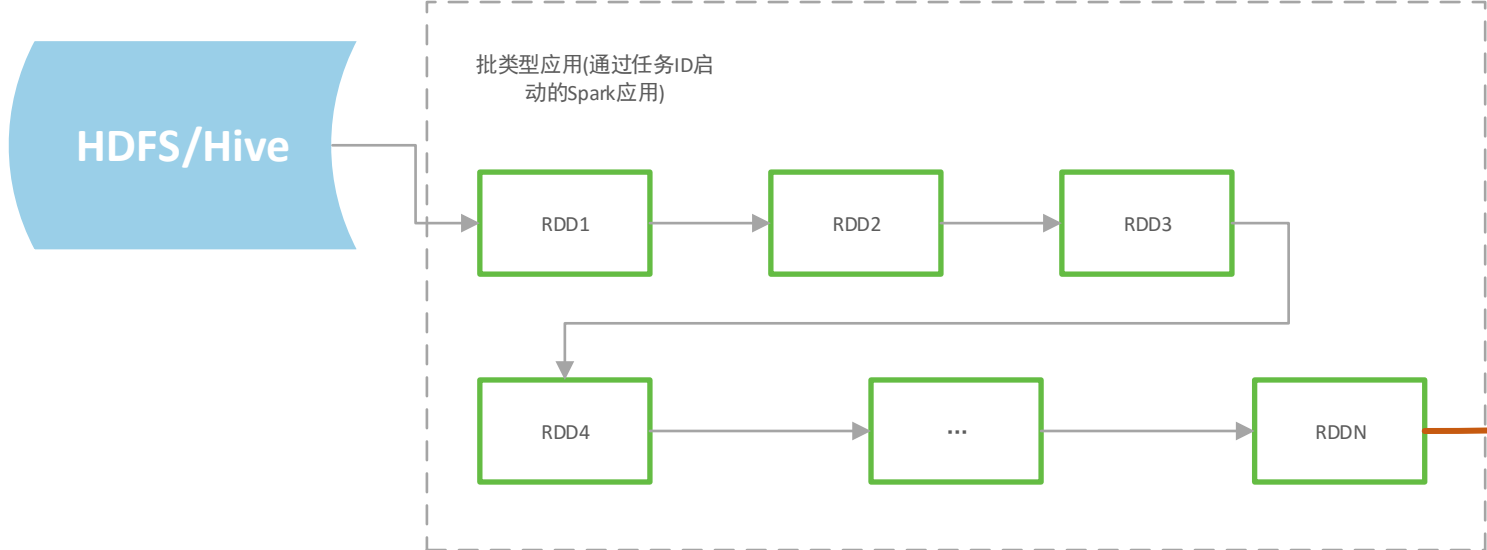
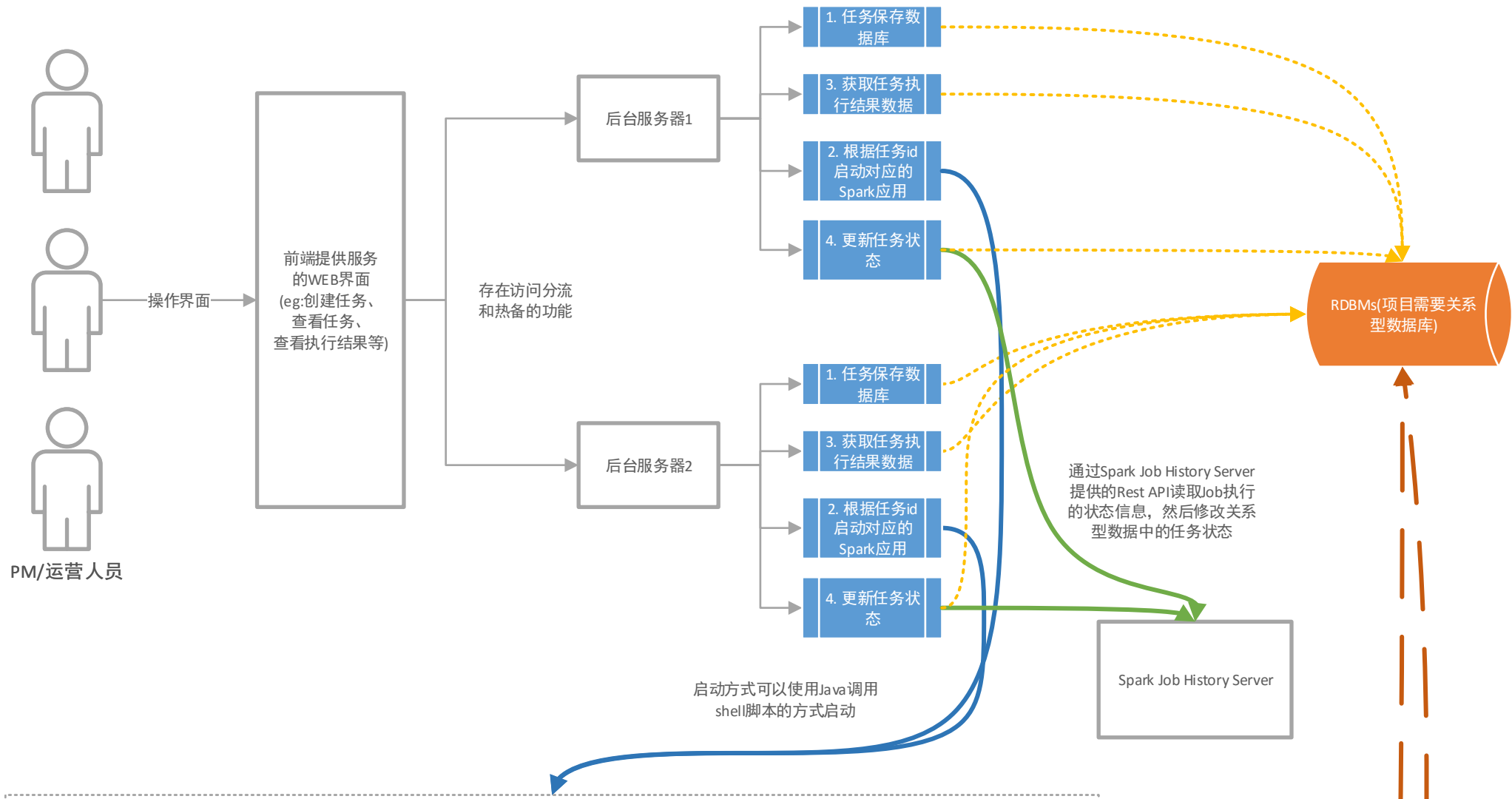
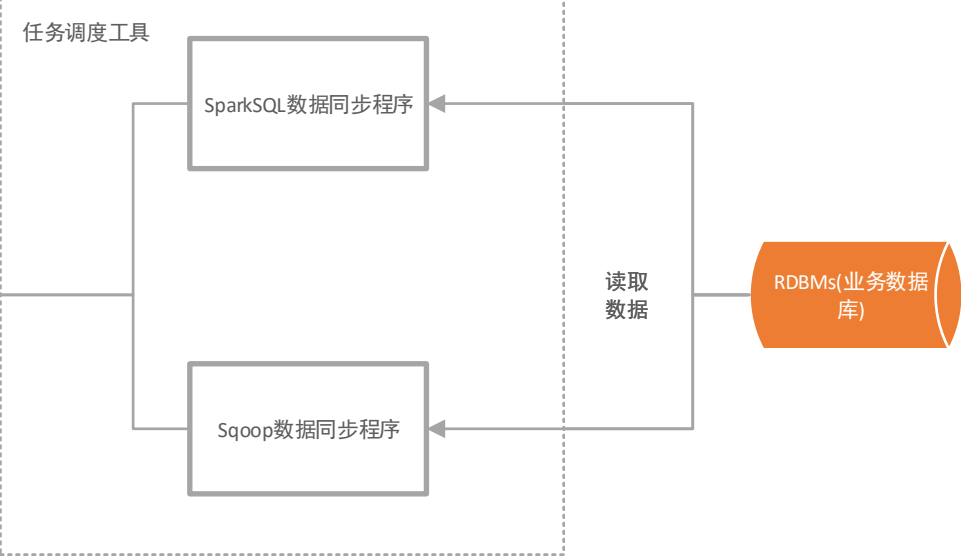
ETL处理后需要继续  
实时处理的数据保存Kafka

SparkStreaming(实时ETL程序)

ETL所有结果数据保存Hive

HDFS/Hive

- 一、这里的Flume client Level中的Agent可以采用以下几种方式：
1. 使用Nginx + Nginx Log File + Exec Source + File/Memory Channel + Avro Sink的方式，类似Hadoop项目
  2. 使用Nginx + Tomcat + Web项目 + Avro Source + File/Memory Channel + Avro Sink的方式，其中Nginx、Tomcat和Flume client agent可以在一台机器上，Web项目在tomcat中，Web项目中将接受的数据以avro的形式写出(Flume API)，参考用户画像项目数据收集讲解
  - 二、Flume Client Agent到Flume Controller Agent层采用sink group设置，增加高可用性
  - 三、Kafka配置多个Partition，增加数据的并行处理能力
  - 四、Kafka对Partition配置三个备份因子，增加数据的健壮性
  - 五、RDBMs和HDFS/Hive之间的数据同步可选技术方案：SparkSQL/SQoop
  - 六、数据同步调度系统可选技术方案：OOzie、Linux Crontab等
  - 七、SparkStreaming的实时ETL程序执行间隔时间为30秒，异步发送数据



任务执行结果数据保存MySQL(根据需求有可能保存HDFS)

