



elastic

基于Elasticsearch的 斗鱼搜索服务实现

白凡@斗鱼

2016-12-10, QQ:285698756

Elastic{ON} DevChina – Dec 10, 2016

关于我

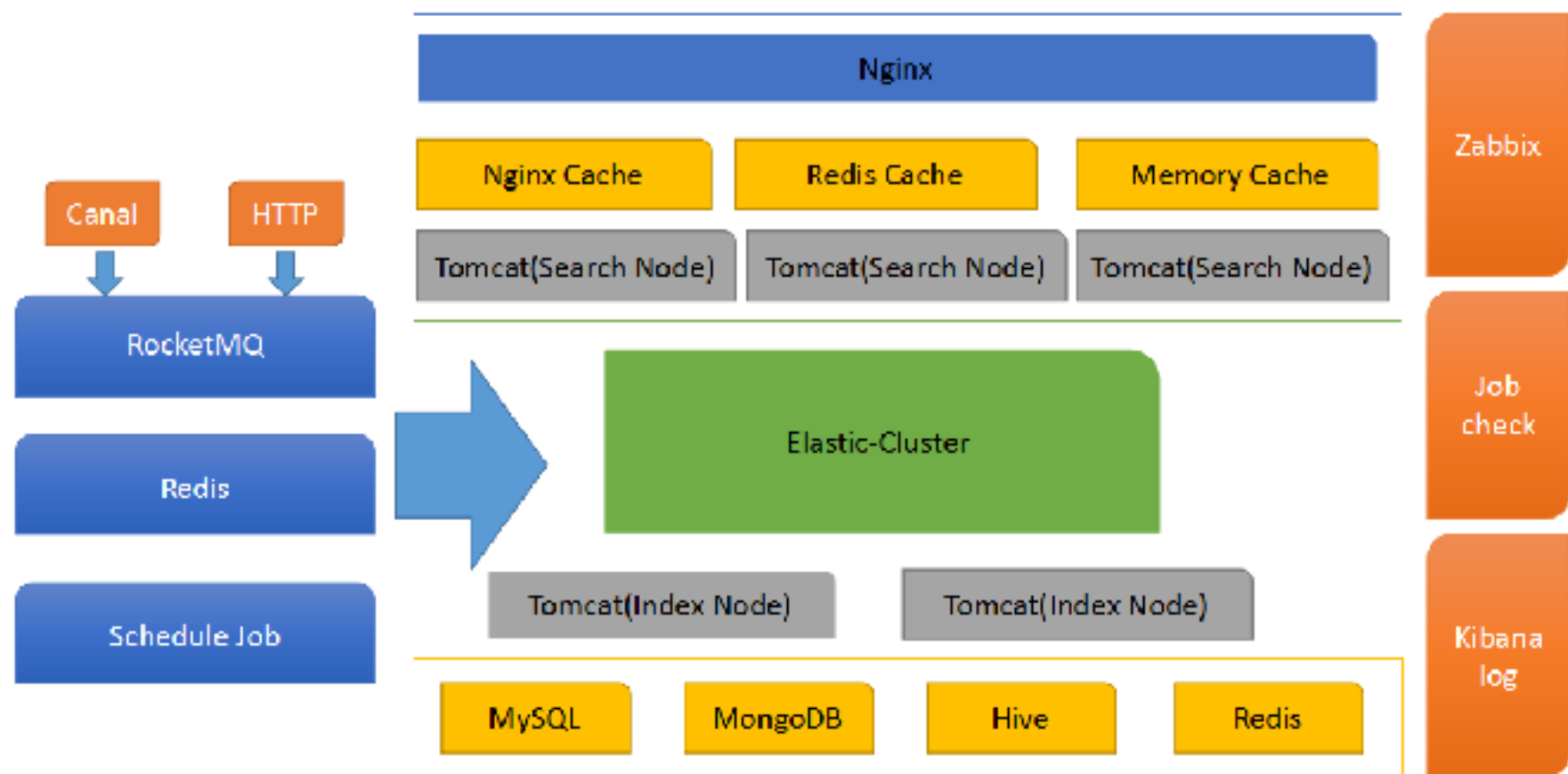
2015年1月进入斗鱼，一直从事斗鱼的搜索服务搭建。

目前主要项目包含，斗鱼主站，移动端、鱼秀、鱼吧、点播站等搜索服务。



主要内容

- 技术架构
- 搜索服务
- 增量索引
- 监控
- 容灾
- 心得体会





搜索服务

IK自定义词库来源，主要来源于用户搜索词打点，以及弹幕分词

搜索与索引分离

HTTP协议，仅作为内部接口，提供给前端

Nginx负载均衡，做好缓存

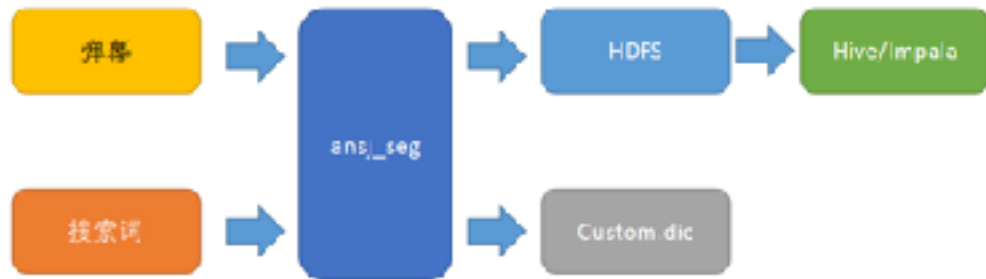
如何丰富自定义词库?

主要来源于斗鱼“弹幕”以及搜索词

采用Ansj中文分词器

HDFS，离线统计分析

Custom.dic，供索引分词



自定义分词器

场景:

斗鱼直播间有房间描述，如“五五开”的直播间，其描述可能就会有“五五开，电竞卢本伟，大神”等。且对于描述的权重比例很高。

而对于搜索业务，我们不希望搜索“五”，“开”，“卢本伟”等词能匹配到，仅在搜索“五五开”，“电竞卢本伟”时返回结果。

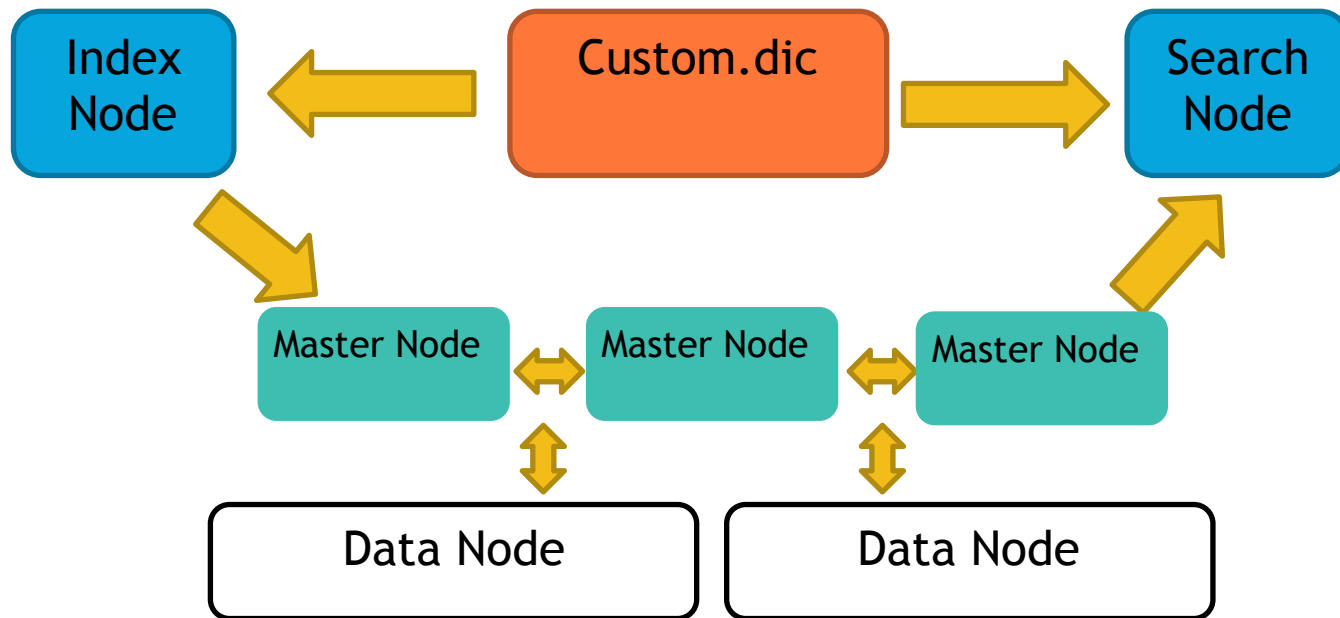
因此，whitespace analyzer加入基于“，”的pattern tokenizer。实现空格和逗号分词。

创建mapping时，analyzer采用 whitespace,search_analyzer仍然使用ik，搜索使用 termQuery

```
index:  
  analysis:  
    analyzer:  
      whitespace:  
        type: "pattern"  
        pattern: "，"
```

```
.startObject("describe")  
  .field("type", "string")  
  .field("analyzer", "whitespace")  
  .field("search_analyzer", "ik")  
  .endObject()
```


搜索接口



Java API开发

BoolQuery多维度查询

维度权重比例不一，权重值配置中心统一调度

```
List<BooleanQueryBuilder> booleanList = new ArrayList<BooleanQueryBuilder>();
//
QueryBuilder queryBuilder = QueryBuilders.termQuery(
    " ", sk).boost(
    (float) OceanDataCommonConfig
        .getValueInteger(" "));
//
QueryBuilder queryRoomNameuzzy = QueryBuilders.prefixQuery(
    " ", sk);
//
QueryBuilder query = QueryBuilders.matchQuery(" ", sk);
//
QueryBuilder queryBuilder = QueryBuilders.termQuery(" ", sk)
    .boost(((float) OceanDataCommonConfig
        .getValueInteger(" ")));
//
QueryBuilder queryBuilder = QueryBuilders.termQuery(" ", sk)
    .boost(((float) OceanDataCommonConfig
        .getValueInteger(" ")));
BoolQueryBuilder shouldQuery = QueryBuilders.boolQuery()
    .should( )
    .should( )
    .should( )
    .should( )
    .should( );
//
QueryBuilder queryStatus = QueryBuilders.termQuery(" ", 1);
BoolQueryBuilder mustQuery = QueryBuilders.boolQuery()
    .must(shouldQuery).must(queryStatus);
```

数据索引采用upsert

以主键作为_id，如房间id，用户id，视频id

便于索引更新

```
IndexRequest indexRequest = new IndexRequest(index, type, id).source(data.toString());  
UpdateRequest updateRequest = new UpdateRequest(index, type, id).doc(data).upsert(indexRequest);  
client.update(updateRequest).get();
```

增量索引

- 由于斗鱼语言的异构性，PHP,C++,Java...，所以索引同步显得非常麻烦，在这里，分享目前斗鱼搜索的几个增量索引的方法。
- RocketMQ，实时性、稳定性最好，可复用，需要中间插件，比较推荐。数据源为Canal解析binlog，亦或前端直接通过HTTP请求POST增量数据。只要设置不同的Group，可共同消费。
- Redis队列，实时性、稳定性好，不可复用，需要中间插件。采用Redis队列进行消费，但不可复用，只能一个服务端消费，且如果未做持久化、一旦Redis挂掉，数据不可恢复。
- 定时任务，无需中间件，但实时性、稳定性不好。通过更新时间定时更新增量索引，一旦任务挂掉，时间点一过，则不可恢复。且对更新时间非常敏感，必须为写入数据库时间。

监控

- 监控，监控，监控
- 对于斗鱼直播搜索这个场景，需要运用到大量数据更新，比如直播状态，房间封面等。
- 一旦出现问题，对排序，显示，索引影响很大。



监控

- 目前斗鱼搜索服务，总体的监控任务包含以下三个点：
 - 1.线上服务器状态监控，基于Zabbix对服务器Nginx连接数、CPU，负载的监控。
 - 2.对于定时任务、搜索接口以及MQ，Redis消费等服务的监控，消费监控定位到源头，接口监控如卡顿以及接口状态
 - 3.基于搜索服务的日志监控功能，可直接套用ELK系列，对AOP级日志进行监控分析。

监控

- 出现过线上CPU几近打满，服务器濒临挂掉的情况
- 对应措施：
 - 1.加数据节点分摊压力
 - 2.去掉慢查询wildcard '*xxx*'
 - 3.后续整体将ES集群从1.x升级到了2.x

容灾

- 万一挂掉了怎么办?
- 索引更新失败了怎么办?
- 1.搜索接口和增量索引分离，对于ES的读取能力还是非常强大，所以将搜索接口与增量索引模块隔离，即使增量索引挂掉，也不会影响线上业务。并且基于监控功能，可很快发现问题，并及时fix。
- 2.定时的snapshot也很重要，可迅速恢复索引。
- 3.索引更新失败或者消息队列消费失败，保留手动同步索引接口。

心得体会

- 分片数并不是越多越好，太多会影响数据索引速度。
- ES集群master节点最好为奇数，防止“脑裂”，当然最好不要为1，否则没法重启master节点。
- 对于不需要分词的String类型索引，一定not_analyzed，因为默认String，会采用标准分词，而导致你会发现wildcard或者prefix查询搜索不到。
- 尽量避免wildcard的“*xxx*”，慢查询，请求量上来，CPU，负载有打满风险。
- mapping创建一定记得留有备用字段，谁也不想反复重建索引，影响线上功能。
- index node和search node隔离，search node挂掉可能性不大，index node挂掉至少不会影响线上业务。
- head、kopf插件是个好工具，当然自己开发也可以。



Docker?

斗鱼目前搜索包含，主站搜索、移动端搜索、鱼吧搜索、点播系统搜索等，

而这其中又涉及到了大量的index node,search node,监控，nginx代理...

采用Docker对接口，业务的隔离是非常有必要的。

Questions?!

Thanks!