

# 搜索平台工程化实践

海拍客 王新博



海拍客是杭州洋驼网络科技有限公司旗下网站，创立于2015年2月，是全国知名的母婴平台。创始团队主要来自阿里巴巴，致力于将海内外最新的品牌、最新的知识、最好的消费理念通过全中国母婴店，带给三线以下城市的消费者，帮助消费者完成消费升级。



- 日志(业务日志，监控日志，tracing log)
- 搜索 (主站搜索，后台搜索，各业务搜索)





## 平台化概述



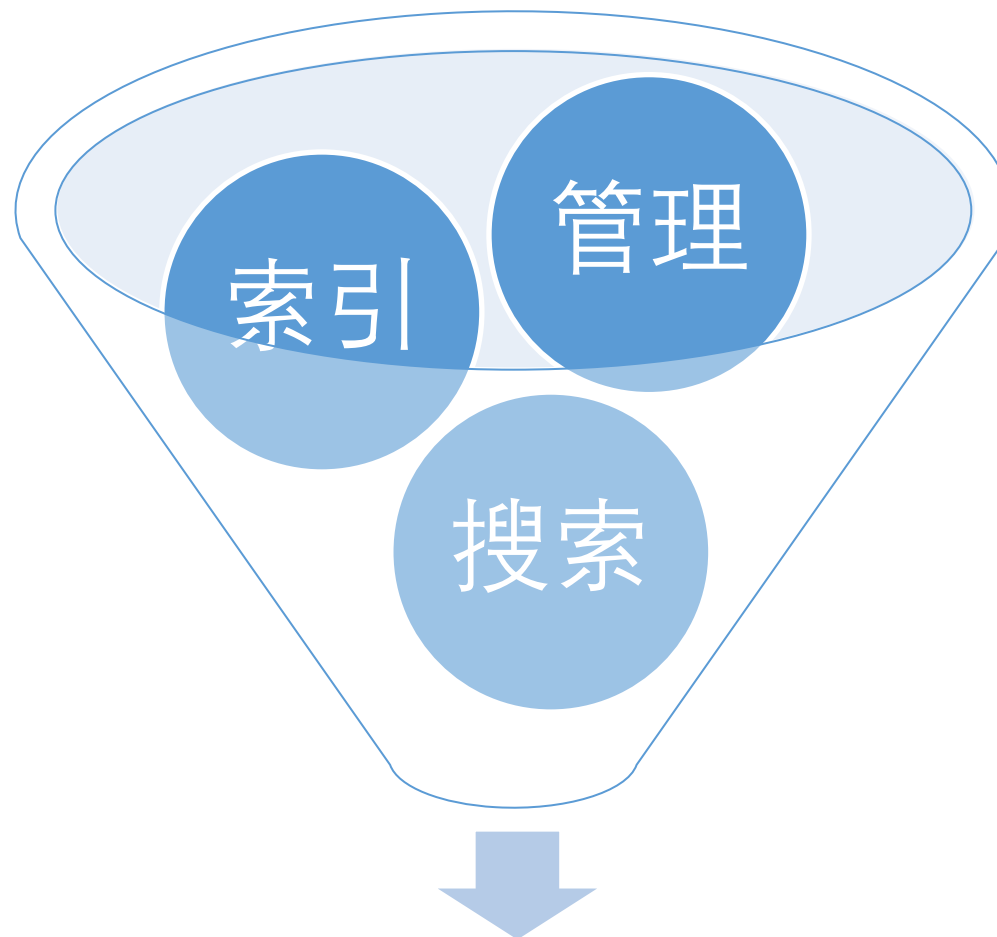
## 搜索管理



## 搜索服务



标准化索引构建过程  
保证数据到索引通道畅通

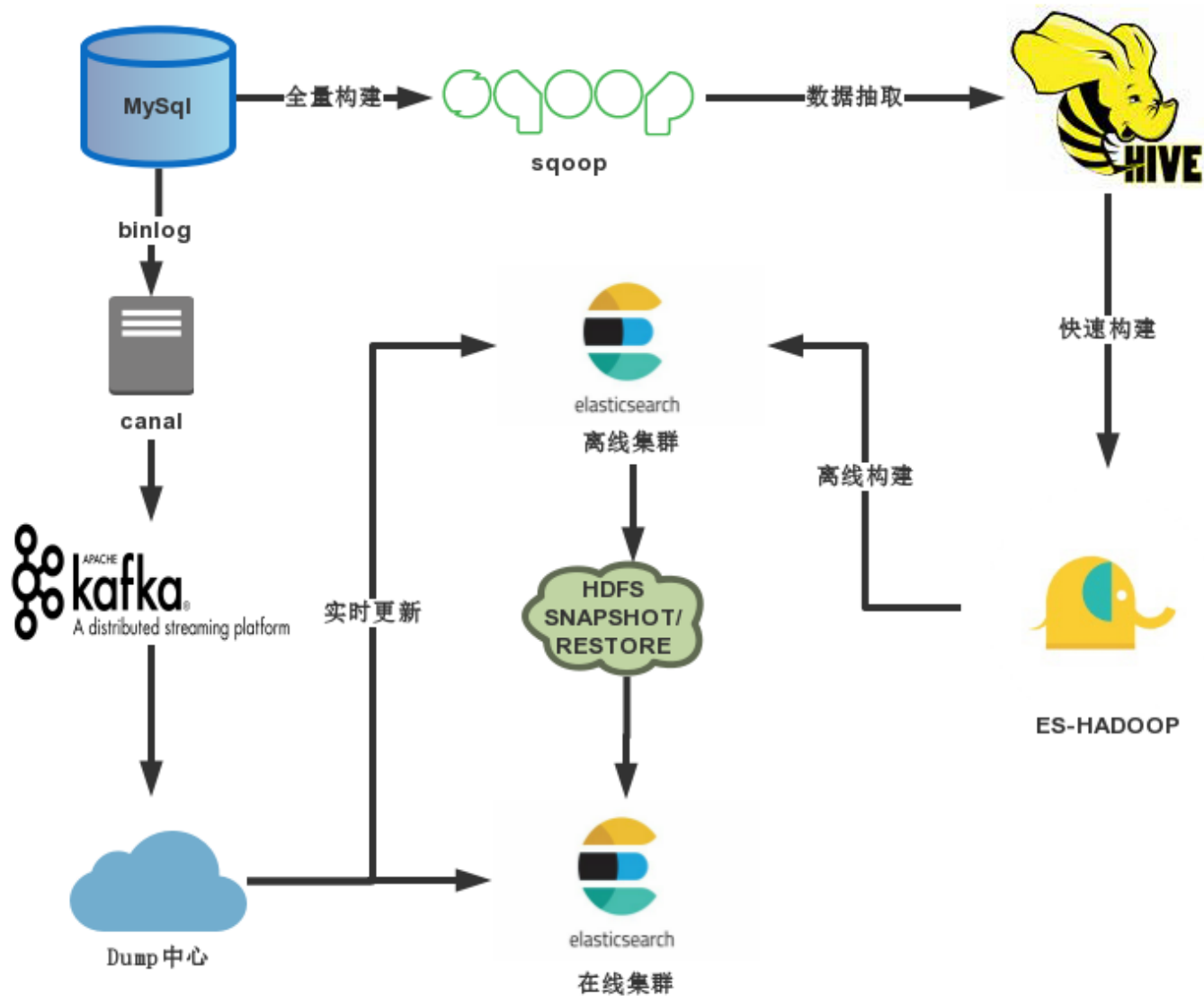


方便创建实例  
管理实例信息

统一的搜索接口  
索引对业务透明

搜索平台

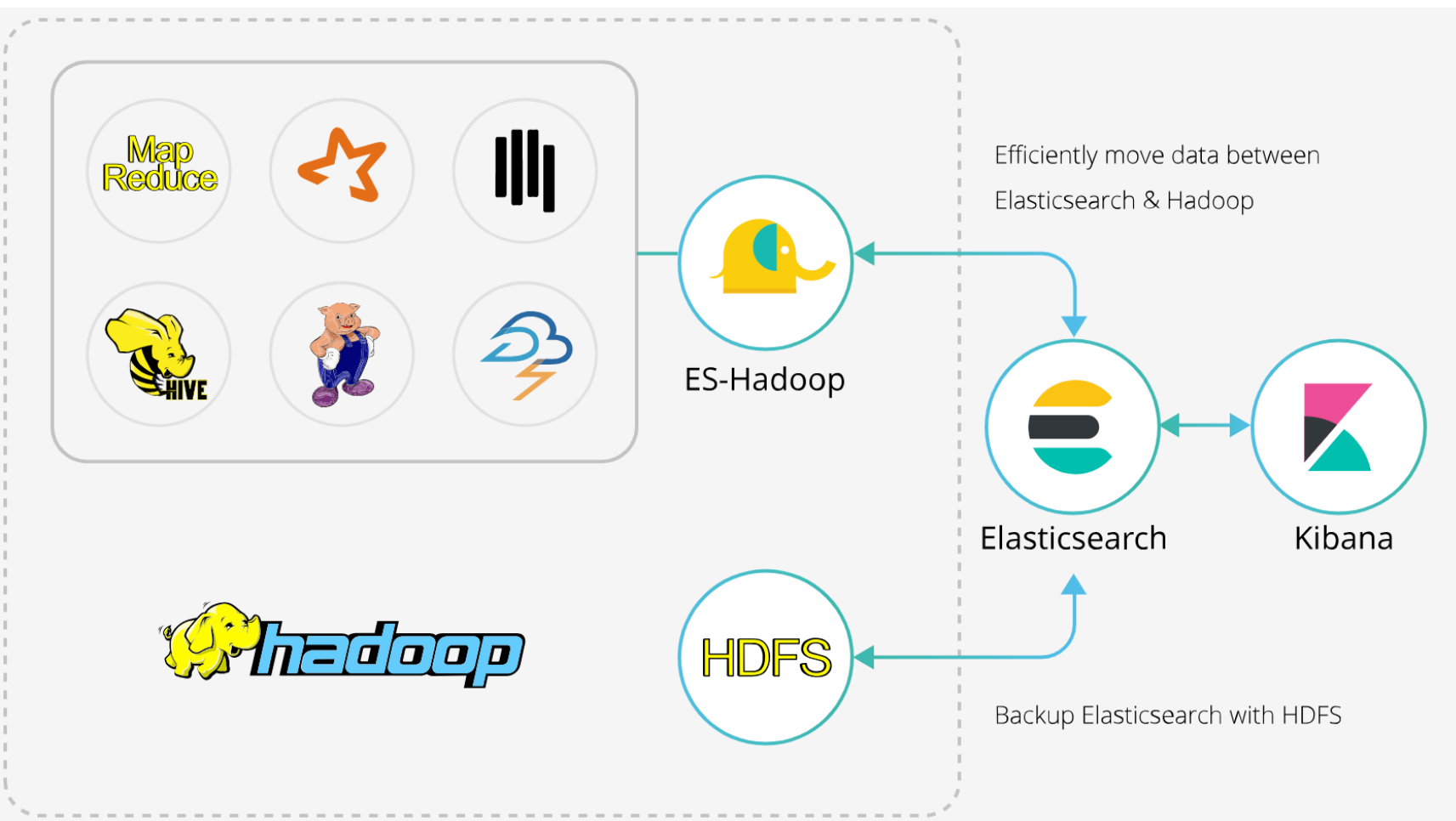




- ◆ 全量构建由sqoop从slave db抽取数据到hive
- ◆ 快速构建直接从hive通过es-hadoop写入es
- ◆ 大索引或重要业务通过离线构建，不影响线上集群
- ◆ 实时数据通过canal发送binlog到kafka,dump中心监听kafka消息完成增量数据更新
- ◆ Dump中心双写数据, 离线集群作为备份索引



# Why ES-HADOOP?

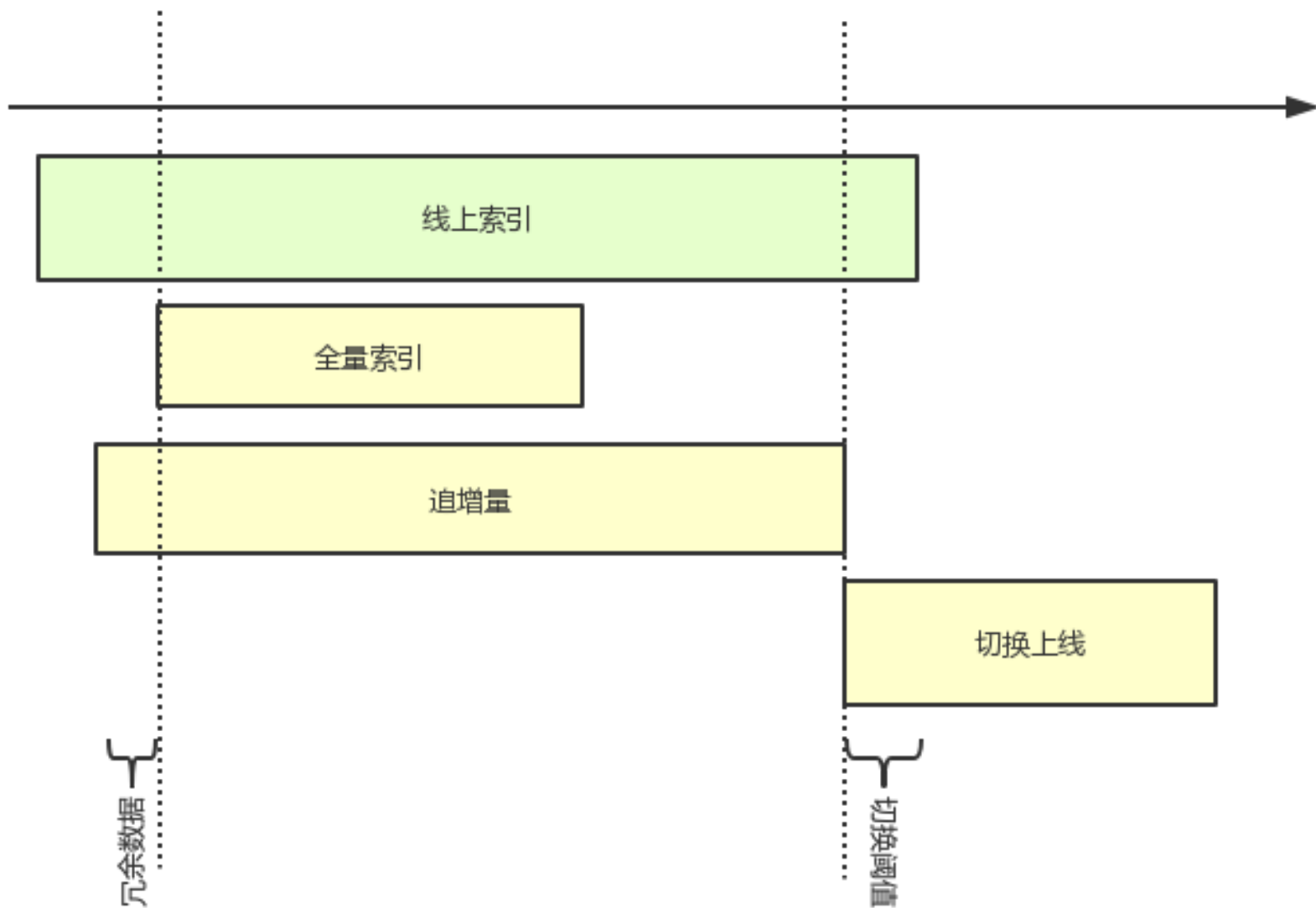


- 凭借现有 Hadoop API 的动态扩展，ES-Hadoop 让您能够在 Elasticsearch 和 Hadoop 之间轻松地双向移动数据，同时借助 HDFS 作为存储库，进行长期存档。分区感知、故障处理、类型转换和数据共置均可透明地完成
- 开发快速，便于部署，快速完成索引全量构建功能

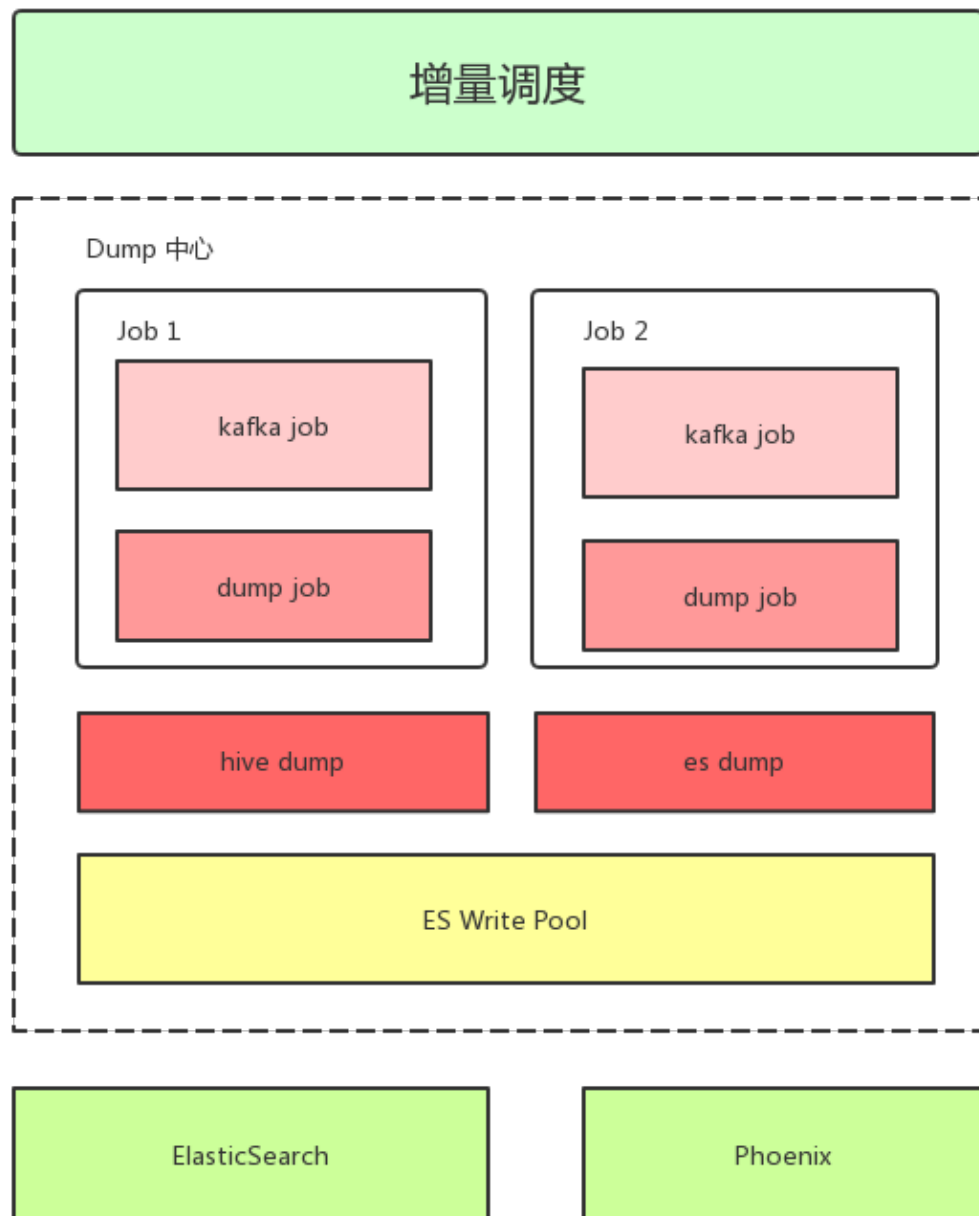


ID	类型	mysql表		过滤条件		操作					
<div>-</div> 22	主表	yt_icp.t_item				<a href="#">修改表</a>   <a href="#">增加字段</a>   <a href="#">删除</a>					
es字段名称	es字段类型	mysql字段	mysql字段类型	hive字段类型	是否主键	关联	是否索引	分词	是否展示	操作	
id	keyword	id	int	INT	主键		索引	不分词	显示	<a href="#">操作</a> √	
item_name	text	item_name	varchar	String			索引	不分词	显示	<a href="#">操作</a> √	
brand	keyword	brand	bigint	BIGINT		23	索引	不分词	显示	<a href="#">操作</a> √	
<div>-</div> 23	辅表	yt_icp.t_brand				<a href="#">修改表</a>   <a href="#">增加字段</a>   <a href="#">删除</a>					
es字段名称	es字段类型	mysql字段	mysql字段类型	hive字段类型	是否主键	关联	是否索引	分词	是否展示	操作	
brand_id	keyword	id	bigint	BIGINT	主键			不分词		<a href="#">操作</a> √	
brand_name	keyword	name	varchar	String			索引	不分词	显示	<a href="#">操作</a> √	

全量构建通过配置生成宽表SQL，将宽表数据写入ES，增量更新根据宽表配置组装数据。搜索平台动态组装宽表，加快了业务的接入速度，不必预先准备宽表逻辑。



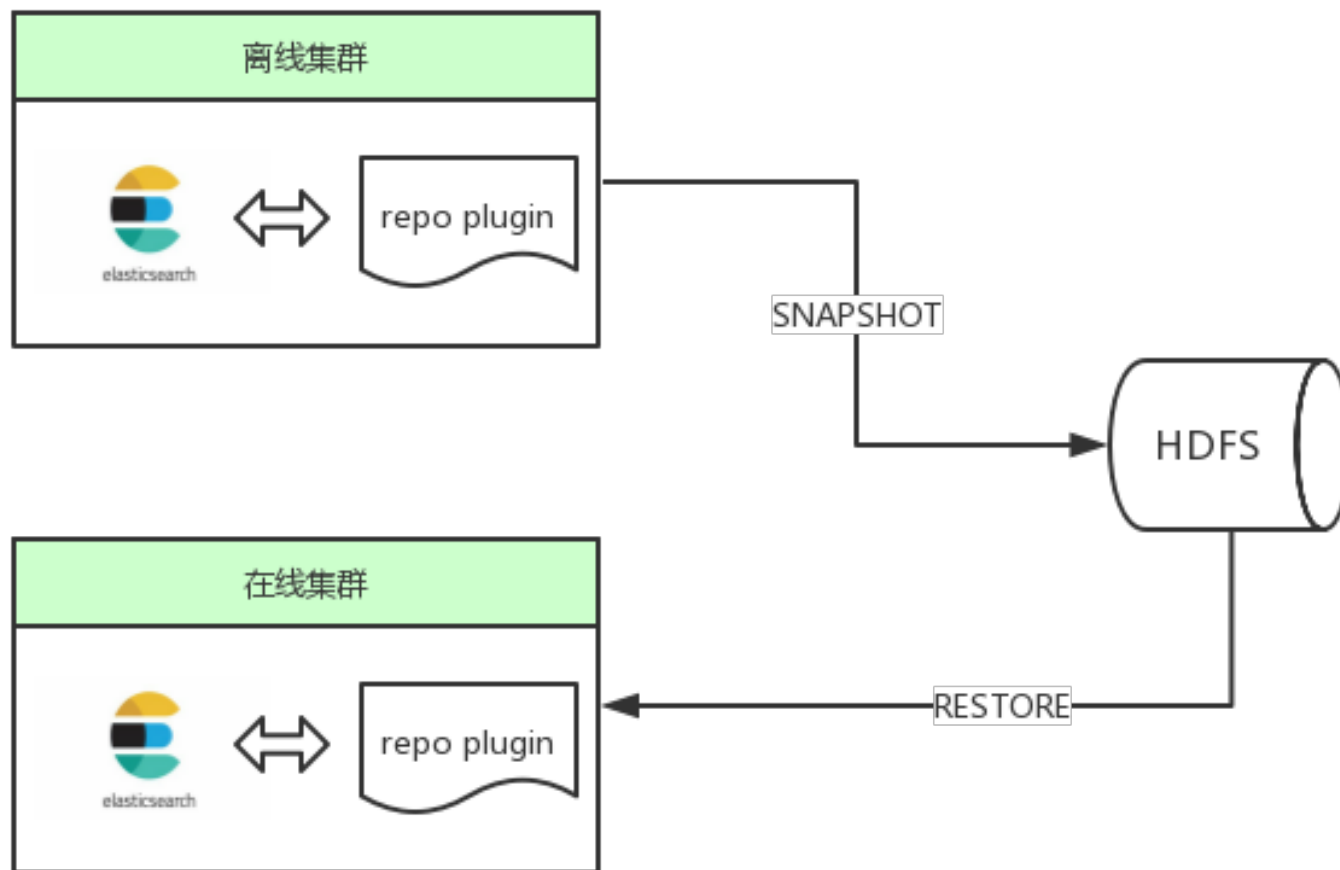
- 全量构建完成后进行追增量
- 追增量指保证线上或邻近当前时间的数据进入索引
- 追增量保证切换到线上后不出现数据不同步情况



- 增量调度控制dump任务的开启及关闭
- Dump任务接受消息通知，将实时更新映射到索引
- Es Write Pool 利用es bulk定时定量的写es从而提高写入性能



- 保证核心业务存在离线备份,大索引构建不影响线上集群
- 利用hdfs repo插件进行索引的snapshot/restore
- 离线构建首先将数据写入离线集群
- 离线集群向hdfs写入快照
- 在线集群恢复快照





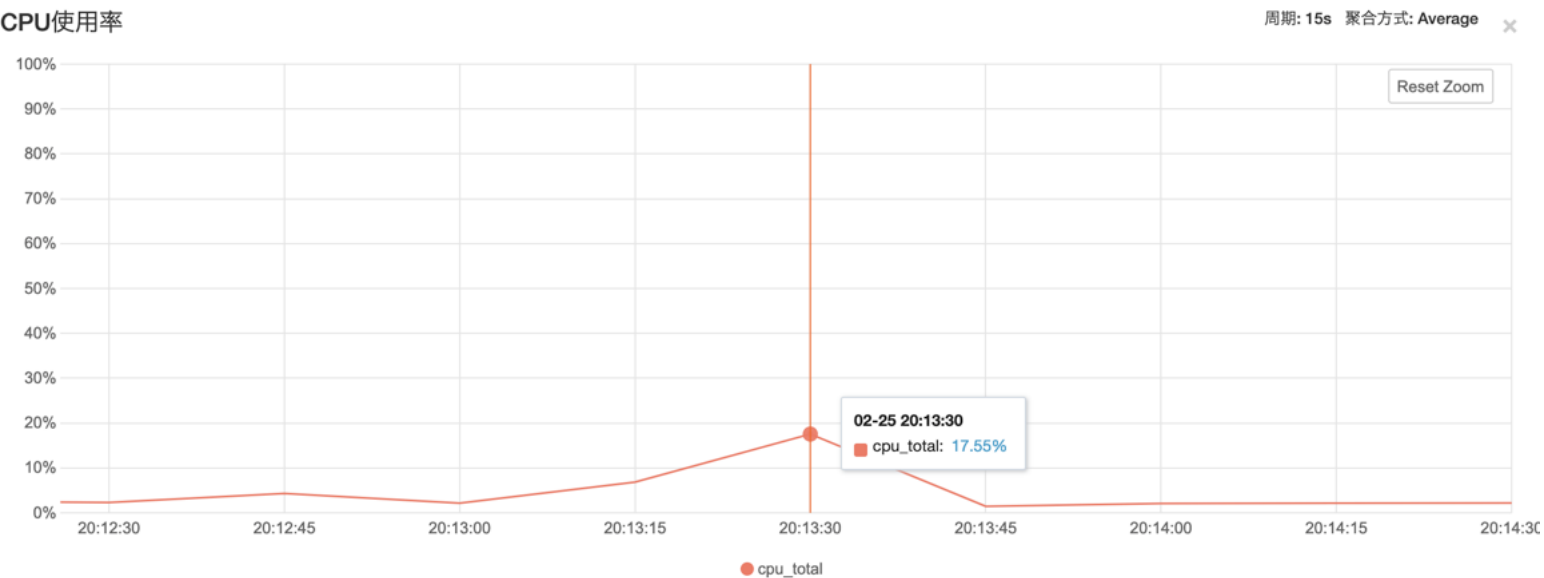
CPU使用率

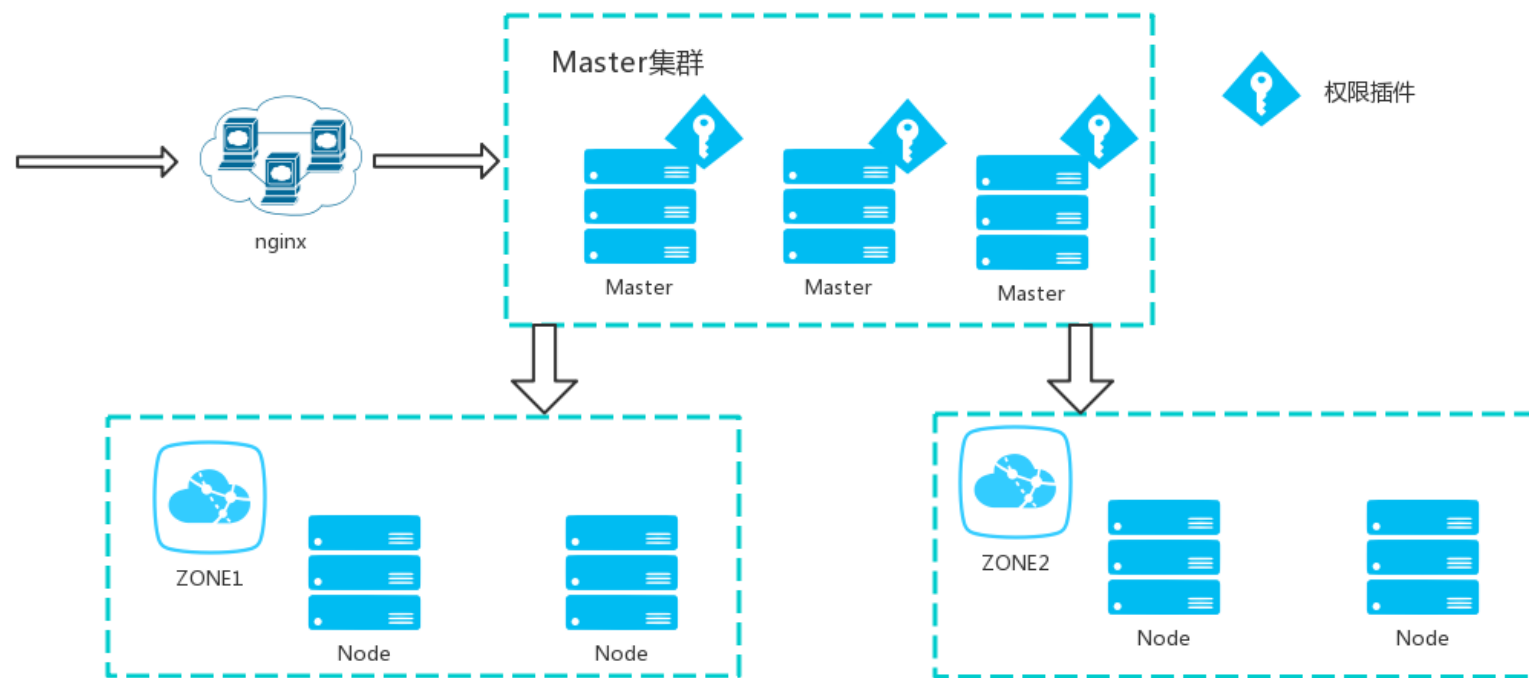


大索引直接线上构建

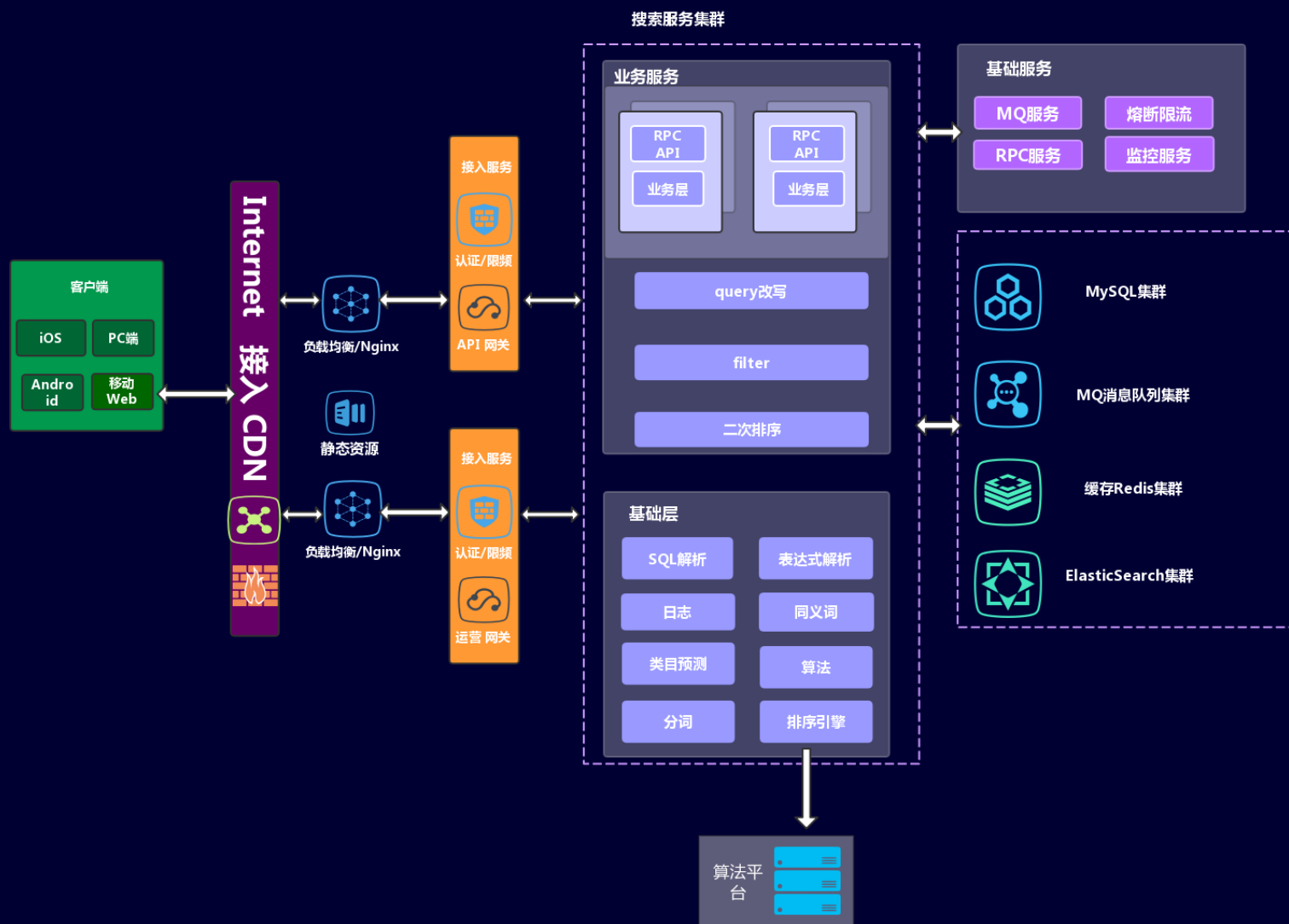
大索引离线构建

CPU使用率

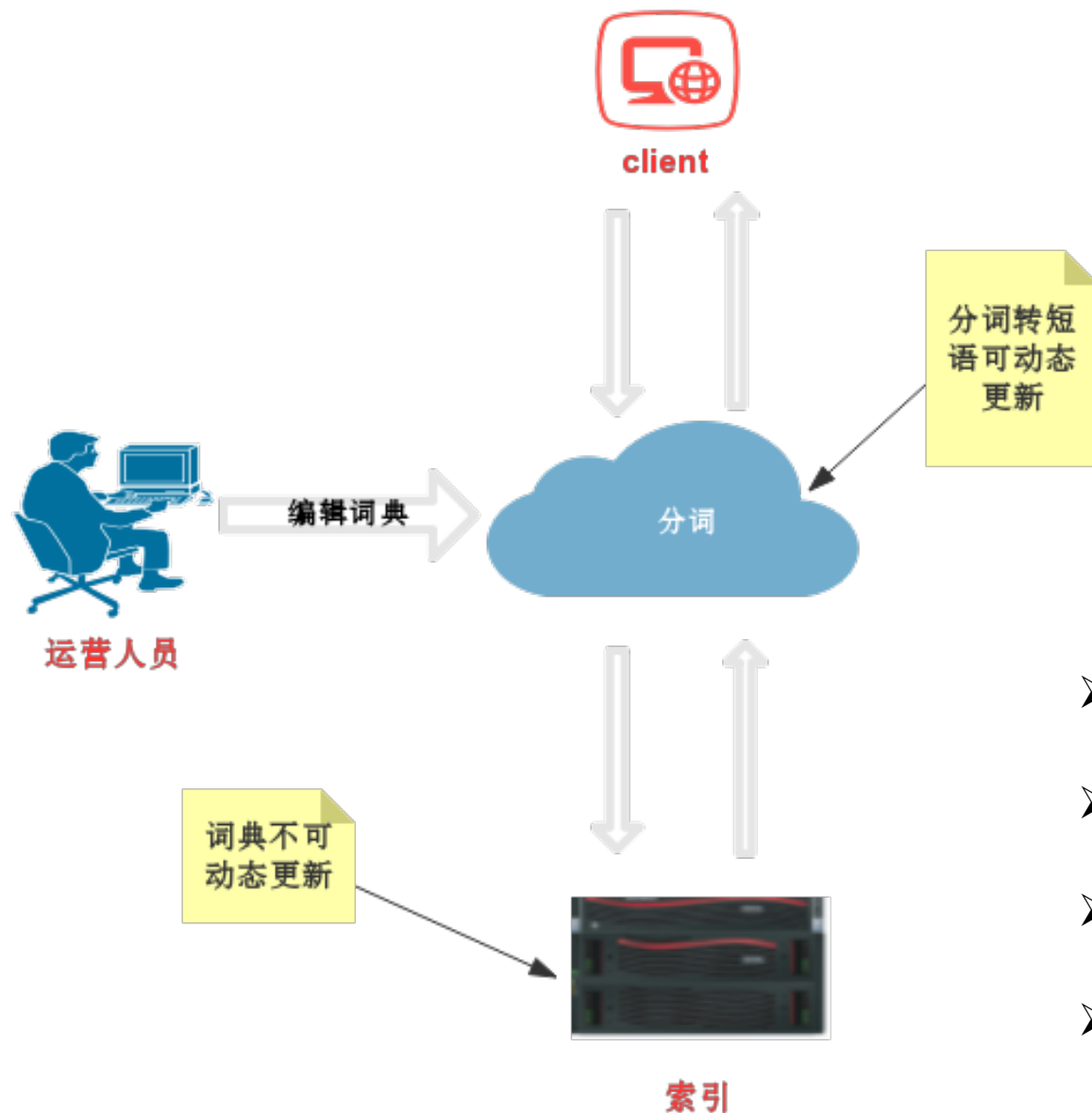




- Nginx 进行IP白名单及 Action 过滤
- Master集群安装自研权限控制插件
- 将Node节点划分到不同的zone，按业务类型分配到对应的zone内
- 未来考虑使用openresty代替插件功能，可获得更灵活的能力



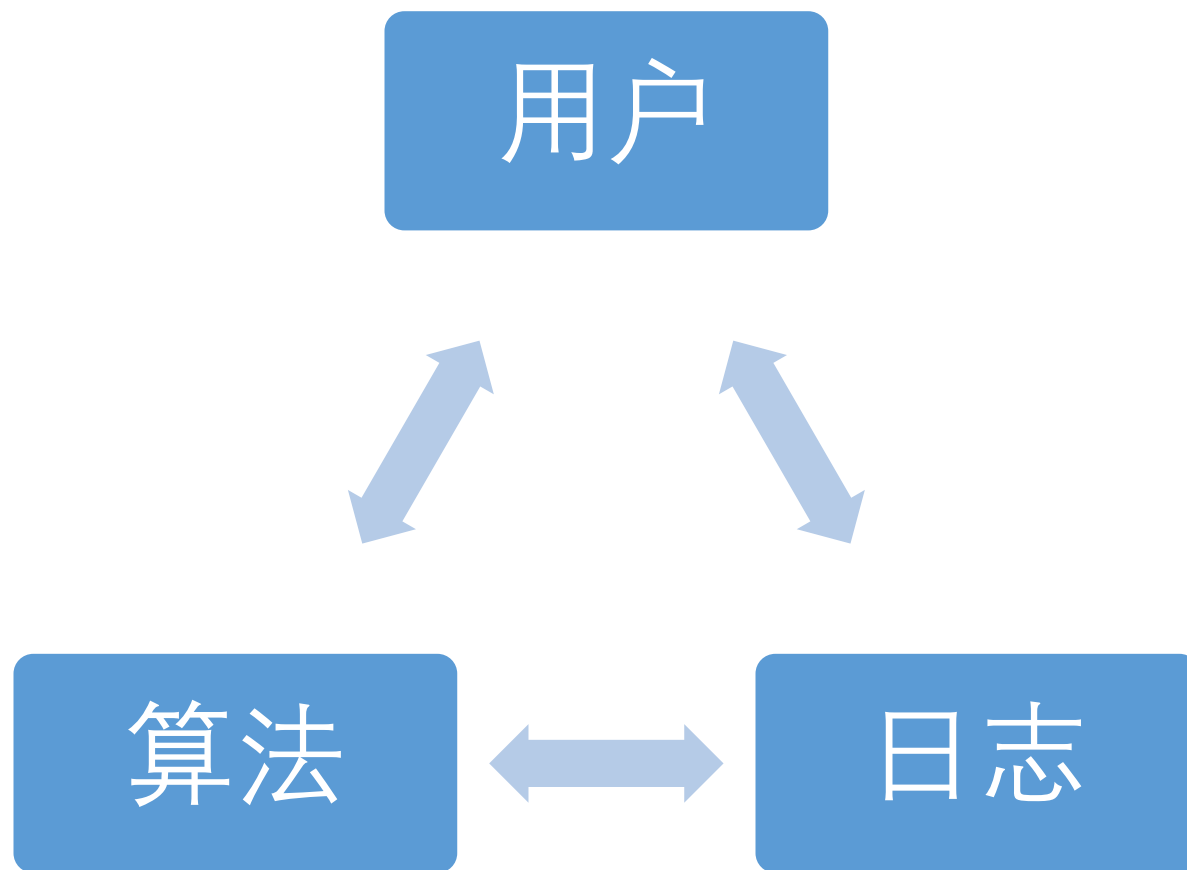
- 封装业务API,同时也提供通用API
- 业务插件为业务提供特定改写, 同义词, 分词,类目推荐等功能
- 支持sql
- 对外提供中间表达式, 插件层将中间表达式转换为ES DSL



- ES 集群分词插件将中文分成单字
- 搜索服务分词后变成短语查询
- 可快速加入新词而不需要重新构建词典
- 索引会变大，增加CPU占用，但是性能完全可接受



- 回馈算法，更懂用户
- 评判排序算法优劣
- 为自己的算法设计合理的日志格式



谢谢！