

PB级数据量背后阿里云Elasticsearch内核优化实践

慕少琼（广富）
阿里巴巴搜索推荐事业部
高级开发工程师
Elastic认证工程师

自研系统

Elasticsearch 增强优化

- 弹性伸缩
- 冷热数据分离
- 高可用
- 向量检索

搜索AI-OS

EYou

基于规则学习的AIOps

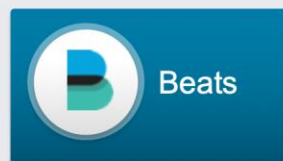
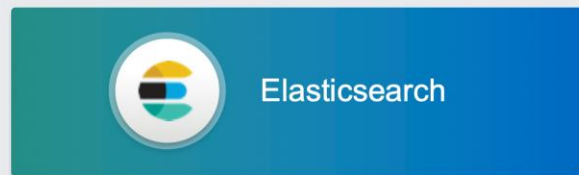


ElasticFlow

数据导入、加工、索引构建



开源系统



X-Pack商业套件

- Security
- Alerting
- Monitoring
- Graph
- Reporting
- ML



阿里云
aliyun.com

ECS/PC/SLB

低成本

开源生态

场景化

源于开源，不止于开源

17个

全球部署

PB 级以上

数据量



【品牌背书】阿里云与社区战略合作



【全球服务】服务覆盖全部阿里云数据中心并且支持本地化专有云交付，和混合云方案



奥运会全球指定云服务商

阿里云Elasticsearch内核优化实践

通用增强

性能优化

离线数据处理平台
通用物理复制
带主键写入去重
Translog无锁优化

高可用方案

读写分离与起新下老
数据自动备份
跨可用区高可用

稳定性保障

主节点调度优化
aliyun-QOS
真实内存熔断器

成本优化

数据压缩
弹性伸缩
堆外FST

场景化内核

日志场景

计算存储分离
冷热分离

搜索场景

向量检索
阿里分词
aliyun-sql

时序场景

场景化推荐模板
时序查询剪枝

性能优化——离线数据处理平台



痛点: 全量写入影响现有集群读写性能及稳定性

ElasticBuild全量

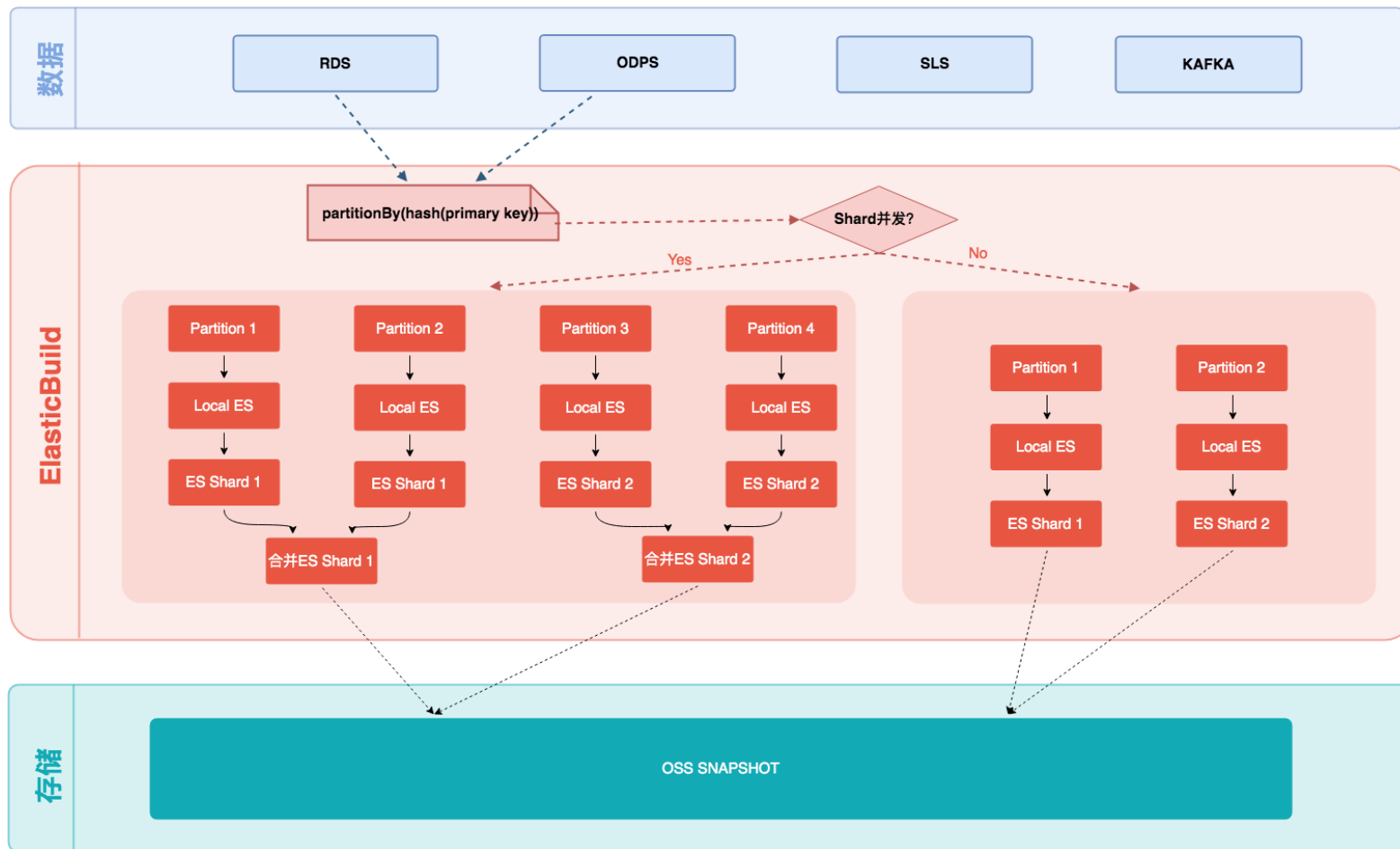
痛点: 全量写入速度慢, 海量数据场景下耗时过长

痛点: 复杂数据处理逻辑实时写入, 影响写入速度

ElasticBuild增量

ElasticBuild – 全量

- 阿里云大数据处理能力
- 内存索引合并避免重复IO
- 多阶段并发提高处理速度
- Blink checkpoint保证failover



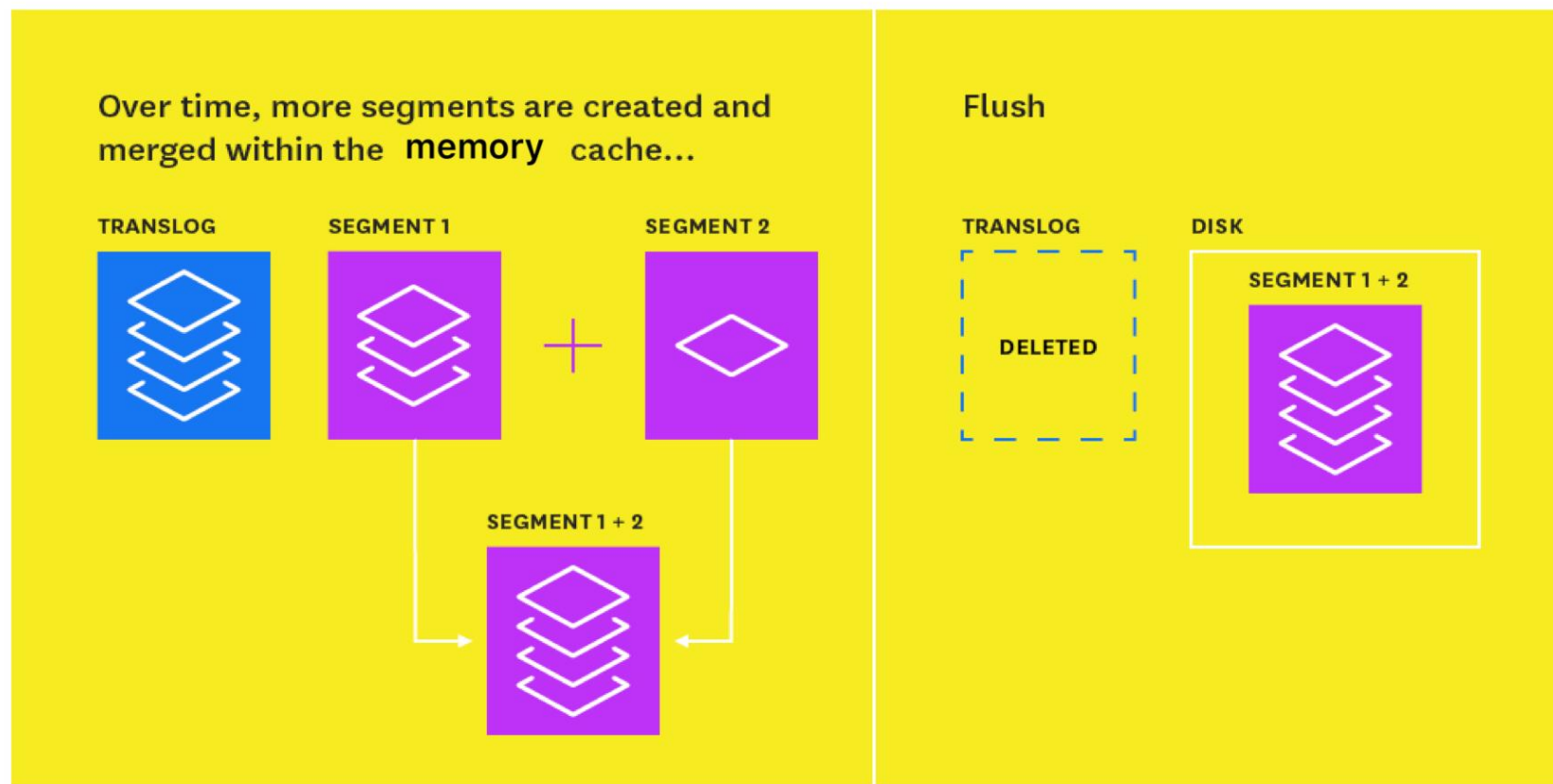
ElasticBuild – 全量

索引内存合并

瓶颈: Segment写入磁盘后, 被反复读到内存中合并, 接着又写入磁盘。这个过程会有大量重复IO开销

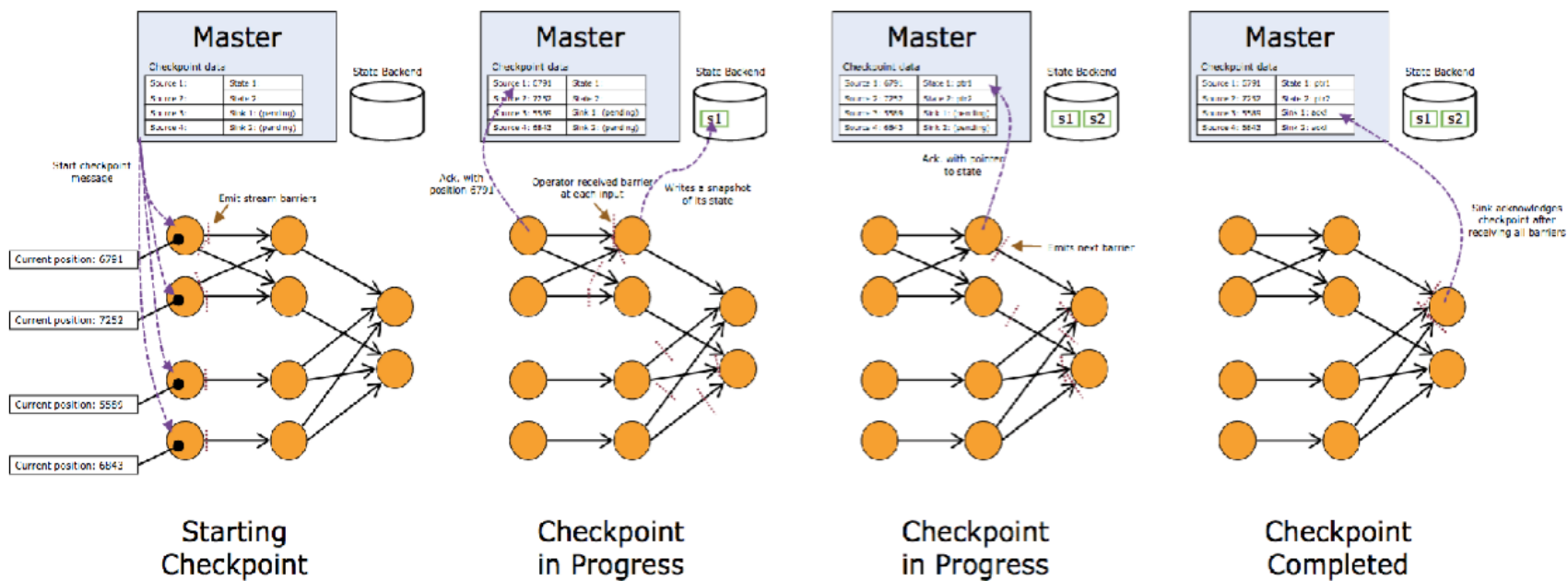
方案: 将索引合并过程放在内存中如图, 多个 segment会在内存中达到一定大小后, 才Flush到磁盘, 减少IO开销

Inside a Shard



ElasticBuild – 全量

Blink Checkpoints取代Translog



优势: 相对于Translog降低了近一倍的IO开销

优势: At Least Once保证数据准确

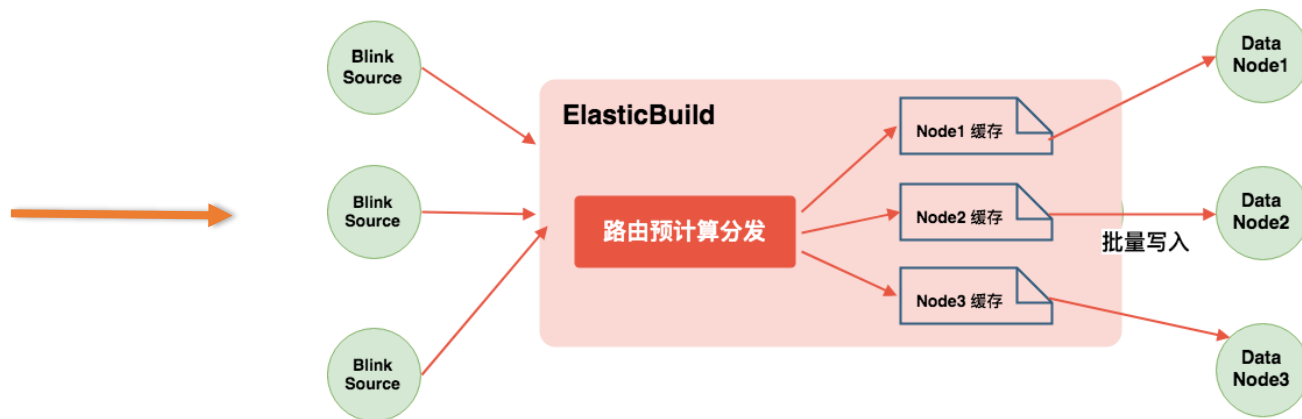
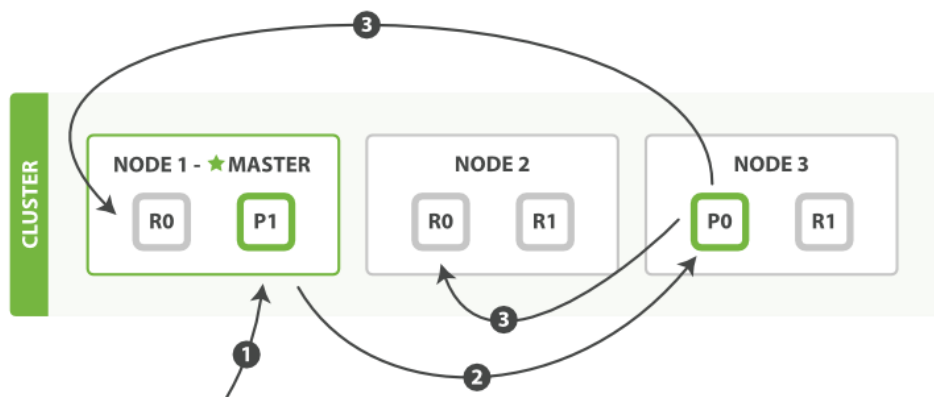
优势: 秒级Failover恢复

ElasticBuild – 增量

用户痛点：整体实时写入QPS为200W/s，单机QPS 3.3W/s，实时写入需优化

瓶颈分析：CPU瓶颈，将路由计算转移至blink中，减少client node的计算和网络开销

结果：ElasticBuild实时部分加入路由预计算批量发送优化，实时写入性能提升**30%**

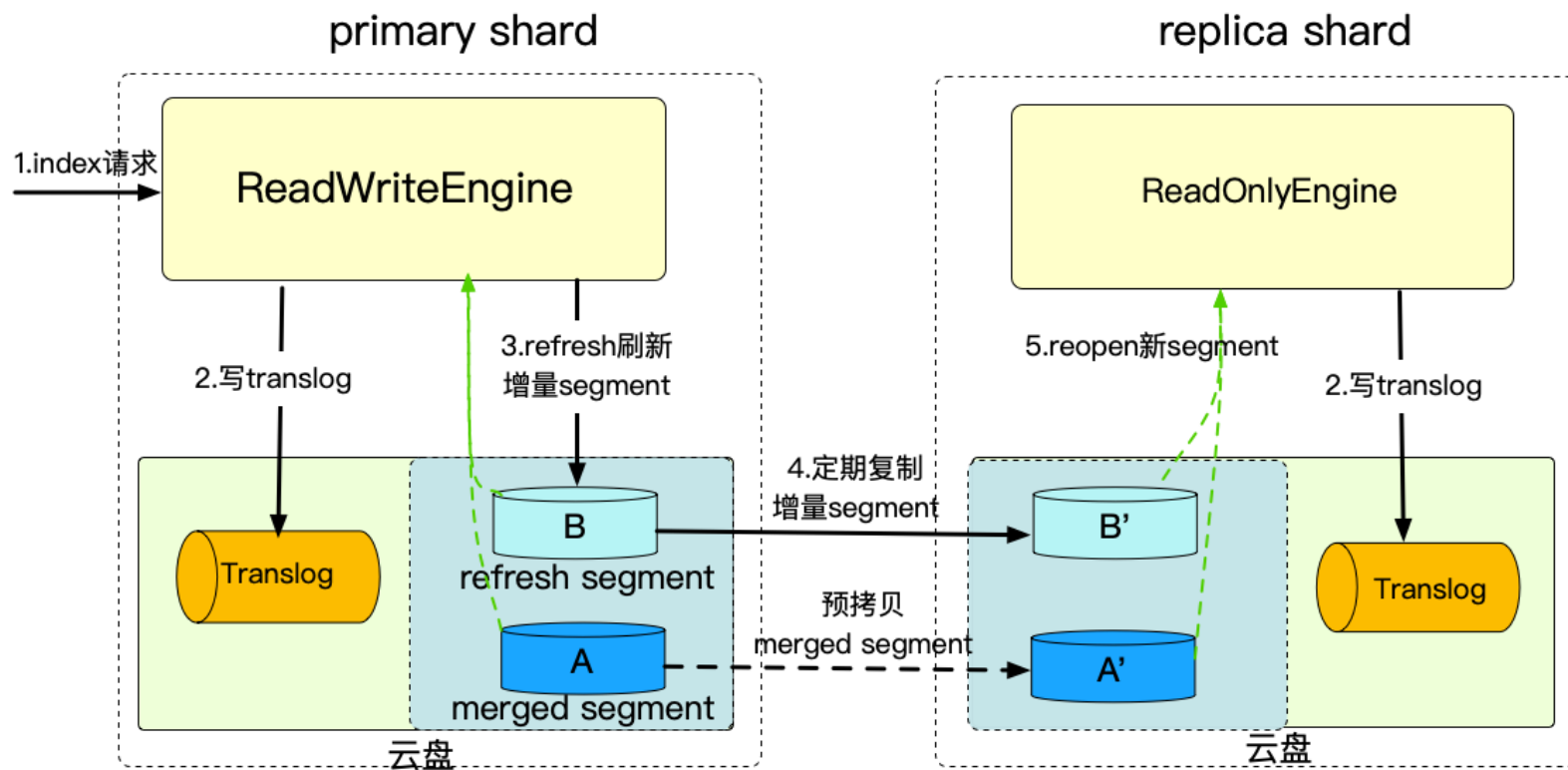


ElasticFlow产品化输出

全量数据build相对于在线build, 性能提升**3倍**

增量数据build, 实时写入性能提升**30%**

通用增强——通用物理复制



主分片保持与原生一致

副本分片只写translog

Refresh时拷贝增量segment

主副分片切换

写入性能提升 **50%**

高可用方案——一起新下老与读写分离

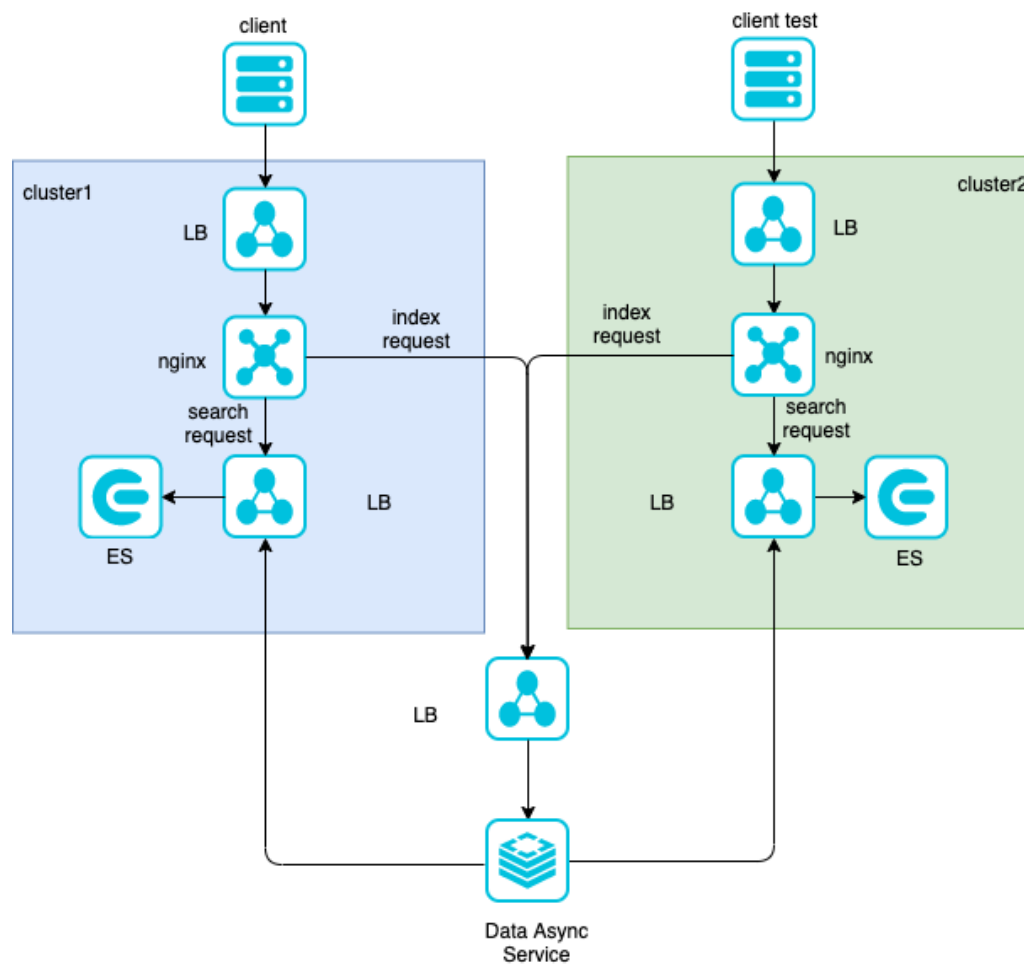
集群升级业务痛点：

- 节点兼容性问题
- 升级回滚
- 保证服务可用性

引入nginx实现平滑读写分流

引入消息队列实现数据高可用

数据双写保证新老数据一致性



通用增强——高可用方案

数据备份

数据备份 (免费试用) 修改配置

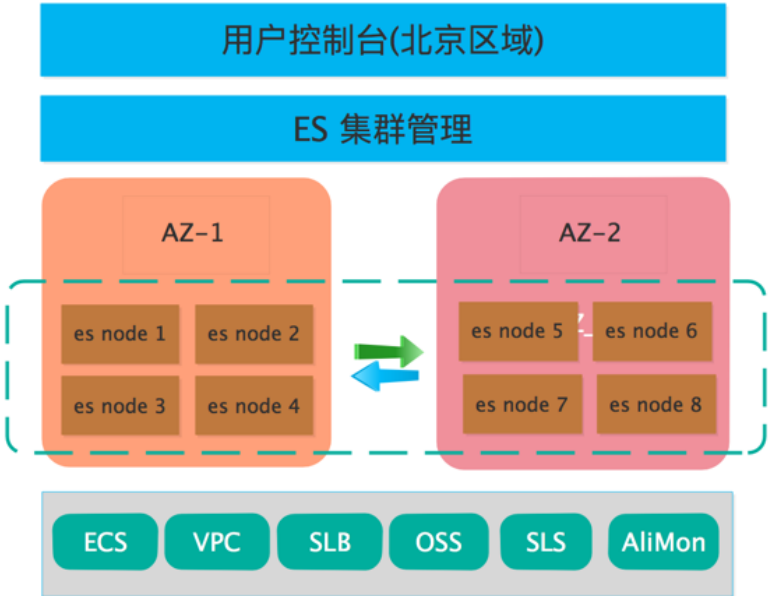
开启自动备份: 自动备份开始时间: 每天 02:00 ?

备份恢复: [点击查看教程](#) 备份状态: [点击查看教程](#)

跨集群OSS仓库设置

引用实例ID	引用仓库名称	引用仓库状态	操作
未创建其他集群OSS仓库引用 立即添加			

跨可用区部署



- 自动异常检测
- 快速故障可用区隔离
- 故障期间可用性保证
- 快速故障恢复
- 业务0影响

稳定性保障——master调度优化

用户痛点: 集群有3个专有主节点、10个热节点、2个冷节点，超过5万个 shard，创建索引和删除耗时超过1分钟。

解决方案: reroute的调度算法复杂度 $O(n^2)$ 降为 $O(n)$ ，大规模集群索引创建和删除耗时降到1s。

<pre> @@ -144,51 +156,34 @@ public int numberOfShardsWithState(ShardRoutingState... states) { 144 * @return List of shards 145 */ 146 public List<ShardRouting> shardsWithState(ShardRoutingState... states) { 147 - List<ShardRouting> shards = new ArrayList<>(); 148 - for (ShardRouting shardEntry : this) { 149 - for (ShardRoutingState state : states) { 150 - if (shardEntry.state() == state) { 151 - shards.add(shardEntry); 152 - } 153 - } 154 - } 155 - return shards; 156 } </pre>	<pre> 156 * @return List of shards 157 */ 158 public List<ShardRouting> shardsWithState(ShardRoutingState... states) { 159 + return Arrays.stream(states) 160 + .map(state -> statesToShards.get(state)) 161 + .flatMap(Collection::stream) 162 + .collect(Collectors.toList()); 163 } </pre>
--	--

<https://github.com/elastic/elasticsearch/pull/48579>

稳定性保障——阿里云QOS

插件配置

集群监控

日志查询

安全配置

数据备份

可视化控制

智能运维

集群概况

系统默认插件列表

自定义插件列表

刷新

请输入插件名称

<input type="checkbox"/>	插件名称	类型	状态	描述	操作
<input type="checkbox"/>	aliyun-knn	系统默认	● 未安装	Elasticsearch向量检索插件	安装
<input type="checkbox"/>	aliyun-qos	系统默认	● 已安装	search rate and bulk size limit	卸载
<input type="checkbox"/>	analysis-aliws	系统默认	● 未安装	Elasticsearch Aliws分析插件	安装 词库配置

节点及索引级别限流

动态开关及限流规则调整



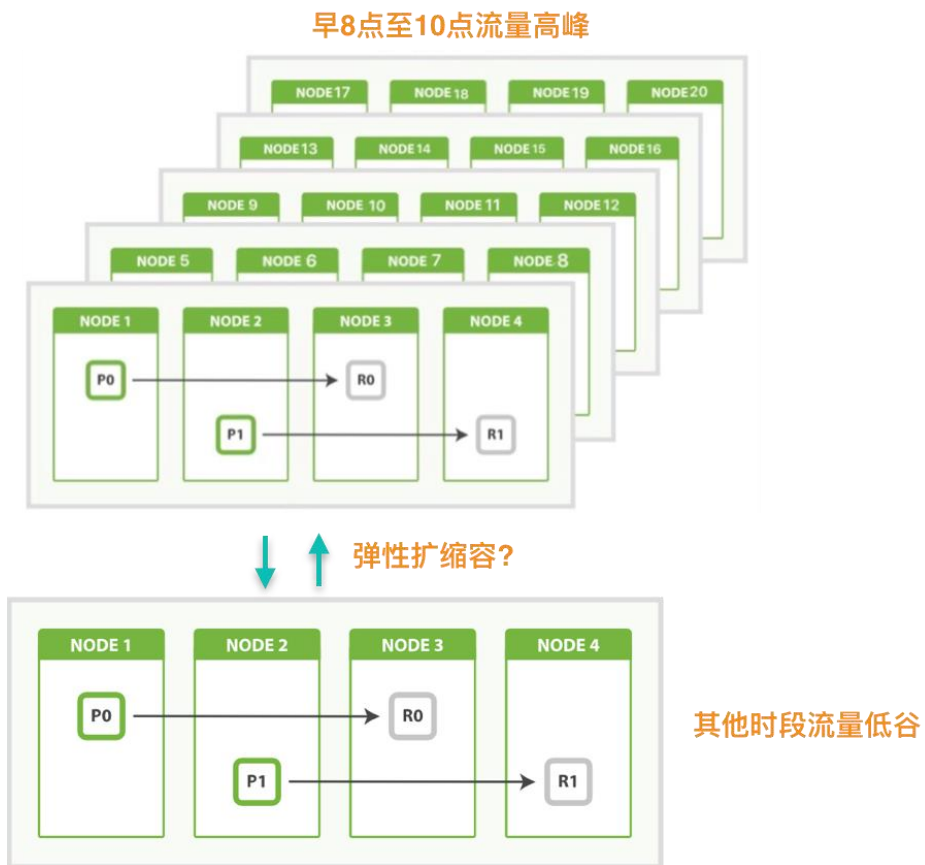
brotili <https://github.com/google/brotli>



zstd <https://github.com/facebook/zstd>

压缩算法	索引大小(GB)	写入TPS(doc/s)
ES默认压缩算法(lz4)	35.5	202682
best-compression(deflate)	26.4	181686
brotili	24.4	182593
zstd	24.6	181393

成本优化——弹性伸缩



针对有明显高/低峰期规律的集群，降低低峰期的使用成本，对数据库加速场景收益更明显

- 触发条件
 - 定时扩缩容
 - 根据业务流量动态扩缩容

- 弹性资源
 - 创建时预估高峰和低峰时的节点个数
 - 后续根据集群运行情况推荐合理值

场景化内核 — 场景化推荐模板

通用场景

日志监控场景

数据库加速场景

搜索场景

经验积淀

场景丰富

考虑因素全面

可扩展性强

登录名 elastic
用来访问Elasticsearch和登录Kibana

登录密码
大写、小写、数字、特殊字符占三种，长度为8-32位；特殊字符为!@#%&*()_+==

场景初始化配置 通用场景

- 通用场景
- 数据分析场景
- 数据库加速场景
- 搜索场景
- 不启用

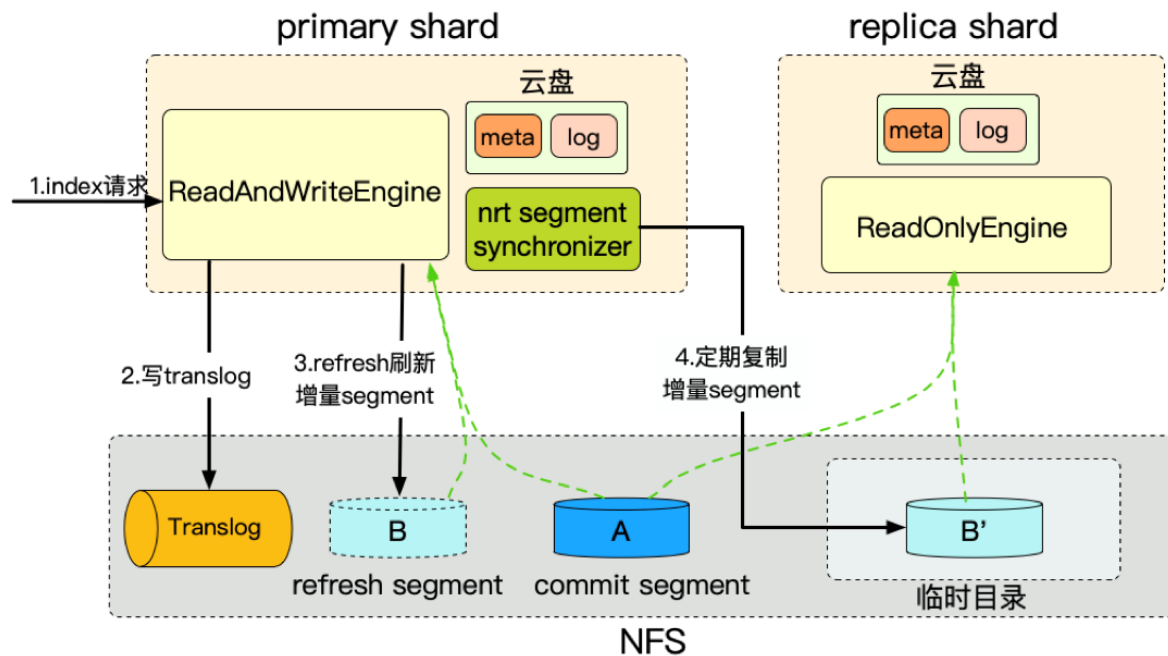
索引模板配置

当前配置	场景化推荐模板 (通用场景)
1 {	1 {
2 "order": -2147483647,	2 "order": -2147483647,
3 "index_patterns": [3 "index_patterns": [
4 "q"	4 "q"
5],	5],
6 "settings": {	6 "settings": {
7 "index": {	7 "index": {
8 "search": {	8 "search": {
9 "slowlog": {	9 "slowlog": {
10 "level": "info",	10 "level": "info",
11 "threshold": {	11 "threshold": {
12 "fetch": {	12 "fetch": {
13 "warn": "200ms",	13 "warn": "200ms",
14 "trace": "50ms",	14 "trace": "50ms",
15 "debug": "80ms",	15 "debug": "80ms",
16 "info": "100ms"	16 "info": "100ms"
17 },	17 },
18 "query": {	18 "query": {
19 "warn": "500ms",	19 "warn": "500ms",
20 "trace": "50ms",	20 "trace": "50ms",
21 "debug": "100ms",	21 "debug": "100ms",
22 "info": "200ms"	22 "info": "200ms"
23 }	23 }
24 }	24 }
25 }	25 }

提交 取消

日志场景——计算存储分离

- 阿里云高速网络环境
- 索引分片一写多读
- 依赖云存储保证数据可靠性
- 状态与索引分离
- IO fence机制保证数据一致性
- 内存物理复制降低主备延迟



写入性能提升100%

计算上避免了副本写入的cpu开销



存储成本倍数级降低

业务数据只存一份

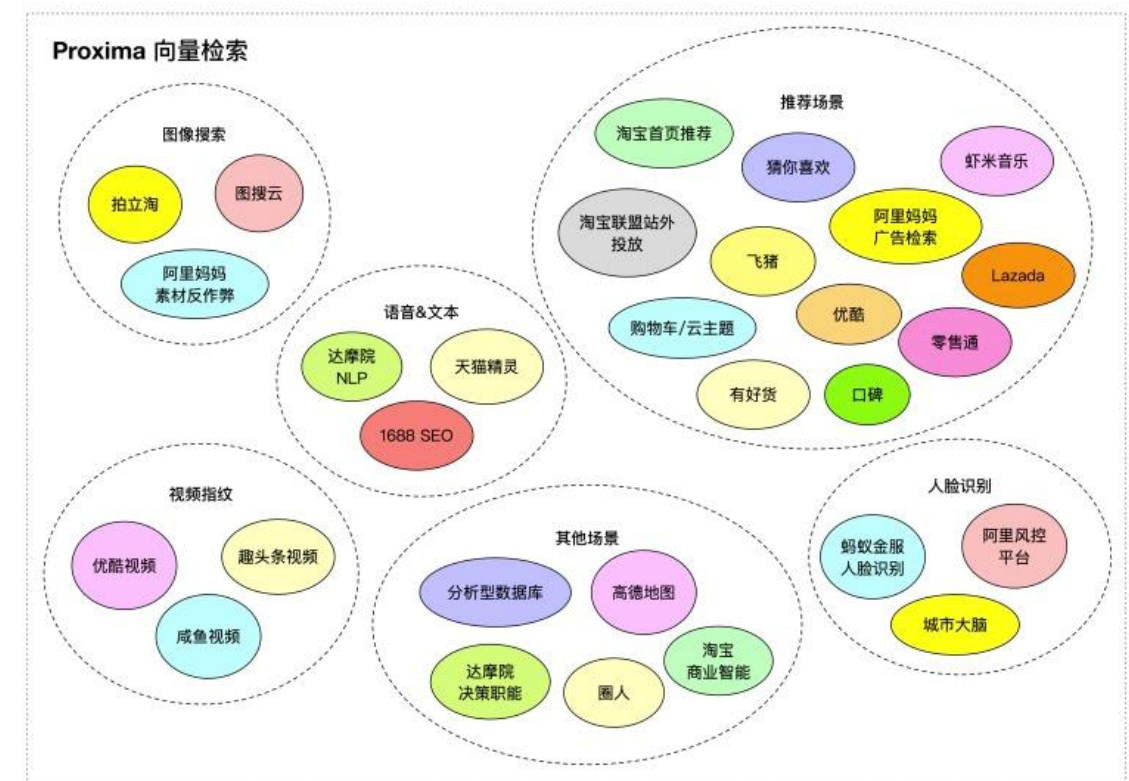


秒级弹性扩缩容

副本秒级快速扩缩容和迁移，轻松应对高峰流量

搜索场景——向量检索

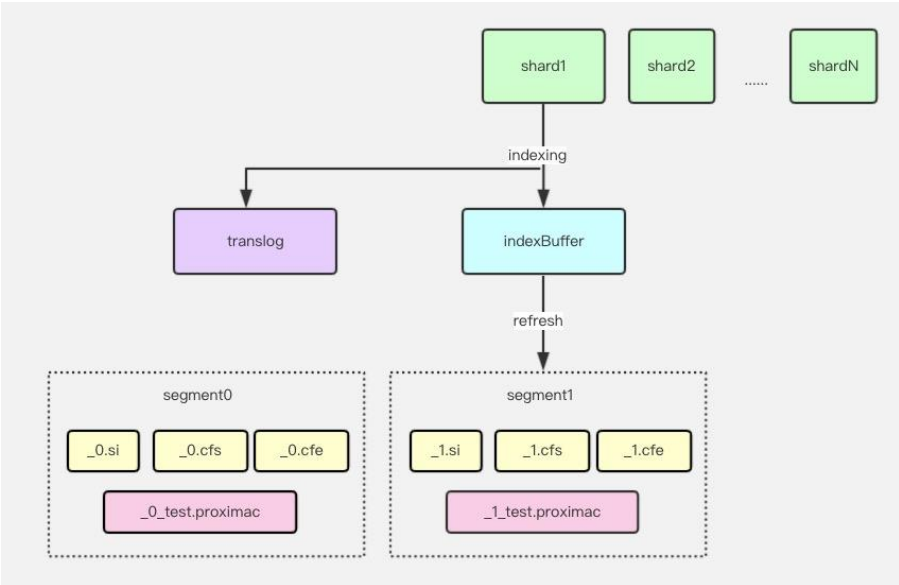
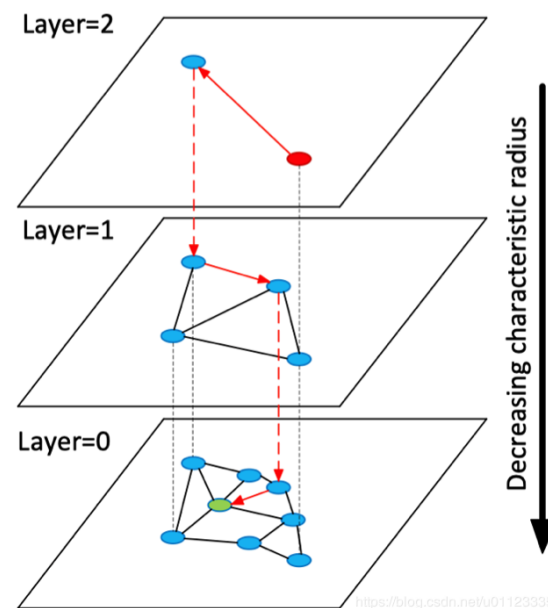
- 基于达摩院高性能向量检索库
- 成熟应用于手淘猜你喜欢、拍立淘、优酷视频、指纹识别等大规模生产应用场景
- 基于Codec机制扩展，完美兼容ES分布式能力
- 基于Hnsw算法，无需训练，查询速度快，召回率高



搜索场景 — 向量检索

图算法: 无须训练 查询速度快 召回率高

- 每个向量节点保存到n个邻居节点的距离
- 插入新节点时，从已构建的任一节点出发
- 计算该节点及n个友点与新节点的距离，选取最近点为新节点的友点
- 为加速查找，添加了类似跳表的分层“高速公路”机制



- Lucene索引数据结构的抽象接口可以自定义倒排/正排等索引的实现机制
- 当refresh刷出segment时，调用自定义Codec，构建segment级别的向量索引
- 查询时，通过实现自定义Weight和Scorer，返回当前segment的近似knn文档id和分数

搜索场景 — 向量检索

阿里云ES 6.7.0版本

机器配置：数据节点16c64g*2 + 100G SSD云盘

数据集：sift128维float向量 (<http://corpus-texmex.irisa.fr/>)

数据量：2千万

Top10召回率	98.6%
Top50召回率	97.9%
Top100召回率	97.4%
延迟(p99)	0.093s
延迟(p90)	0.018s



基于阿里巴巴alinlp分词技术

支持多种模型和分词算法包括CRF、结合词典的CRF、MMSEG等，应用于多种业务场景包括淘宝搜索、优酷、口碑等，提供近1G的海量词库



支持热更新alinlp词典

通过控制台上传新词典干预分词效果

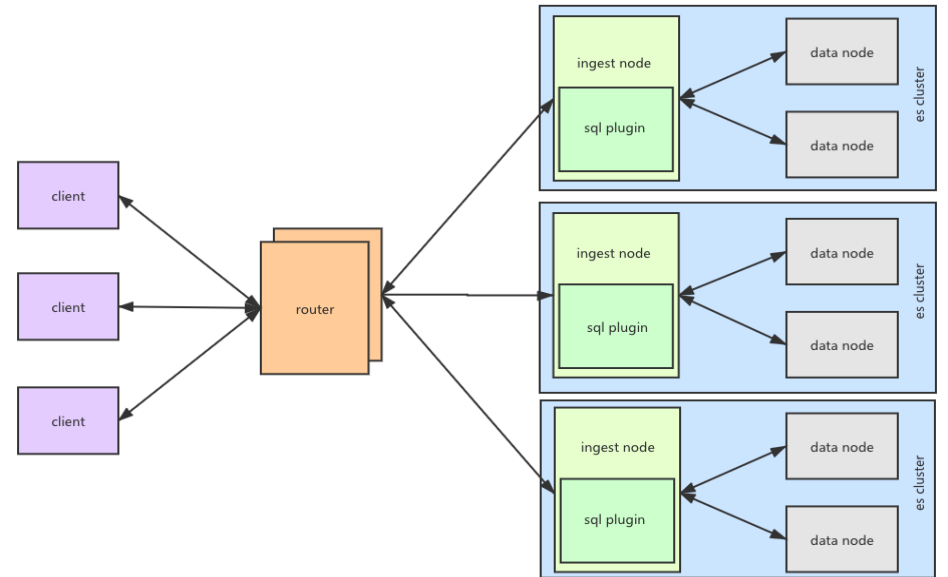
分词效果比较

分词器	分词结果	查全率	查准率	查询性能
标准分词器standard	南/京/市/长/江/大/桥	高	低	低
中日韩文分词器cjk	南京/京市/ 市长 /长江/江大/大桥	高	低	高
IK中文分词器ik_max_word	南京市/南京/ 市长 /长江大桥/长江/大桥	高	低	高
IK中文分词器ik_smart	南京市/长江大桥	低	高	高
阿里中文分词器aliws	南京/市/长江/大桥	高	高	高

搜索场景 — aliyun-sql

aliyun-sql插件是基于Apache Calcite开发的部署在服务端的SQL解析插件，使用此插件可以像使用普通数据库一样使用SQL语句查询Elasticsearch中的数据，从而极大地降低学习和使用ES的成本。

SQL插件	sql解析器	分页查询	Join	Nested	常用Function	Case Function	扩展UDF	执行计划优化
x-pack-sql (6.x版本)	antlr	支持	不支持	支持 (正常语法为a.b)	支持的Function较丰富	不支持	不支持	执行计划优化规则相对较多
opendistro-for-elasticsearch	druid	不支持 (最大查询数量受ES的max_result_window参数限制)	支持	支持 (nested(message.info))	支持的Function较少	不支持	不支持	有少量的执行计划优化规则
aliyun-sql	javacc	支持	支持。具有截断功能，可动态配置单表查询数量，详情请参见 语法介绍 。	支持 (正常语法为a.b)	支持的Function较丰富，详情请参见 Function和表达式 。	支持	支持。详情请参见 自定义UDF函数 。	执行计划优化规则相对较多，并且使用Calcite优化执行计划



Elasticsearch 中文技术... 

3652人



 扫一扫群二维码，立刻加入该群。

欢迎加入Elasticsearch技术交流钉钉群!