



Elasticsearch在vivo搜索中台的实践

王文谦

vivo 互联网搜索团队

分享嘉宾



王文谦

vivo互联网技术经理
搜索中台负责人



一、vivo搜索架构演进



二、ElasticSearch在vivo搜索
中台的实践



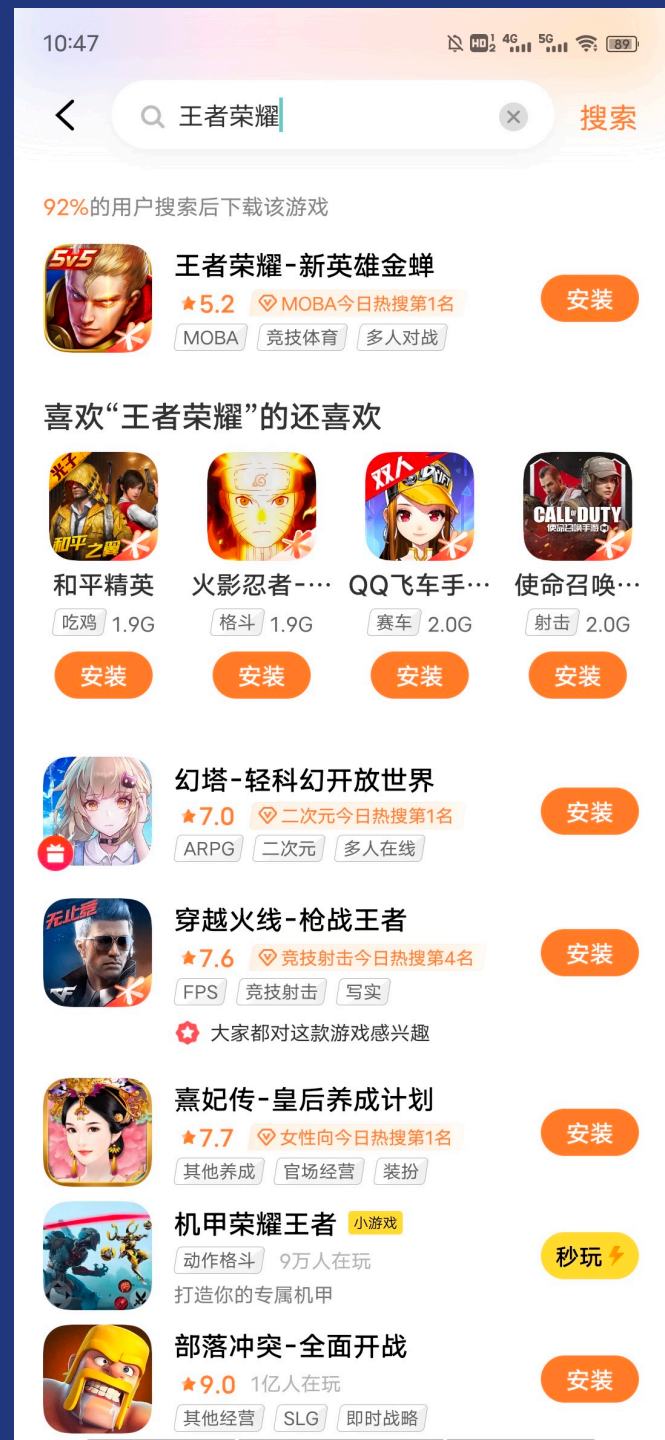
三、搜索中台业务应用总结

一、vivo搜索架构演进

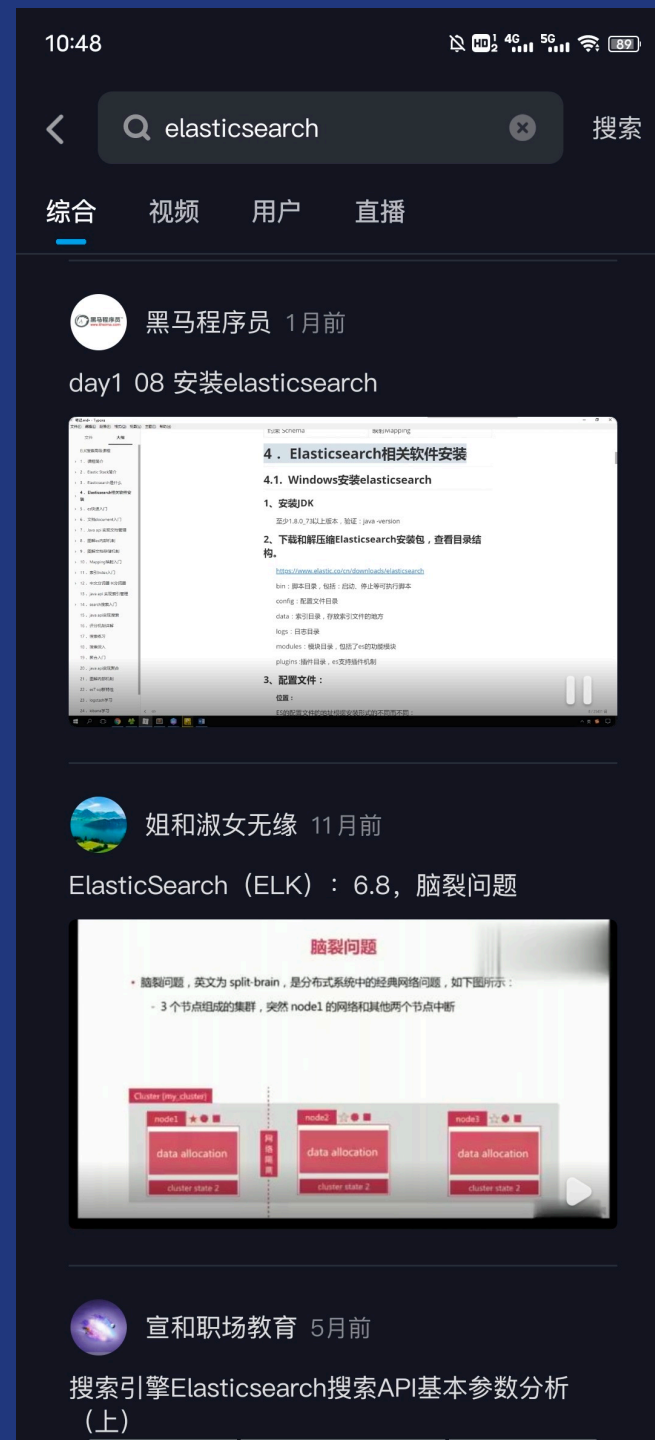
vivo 搜索业务介绍



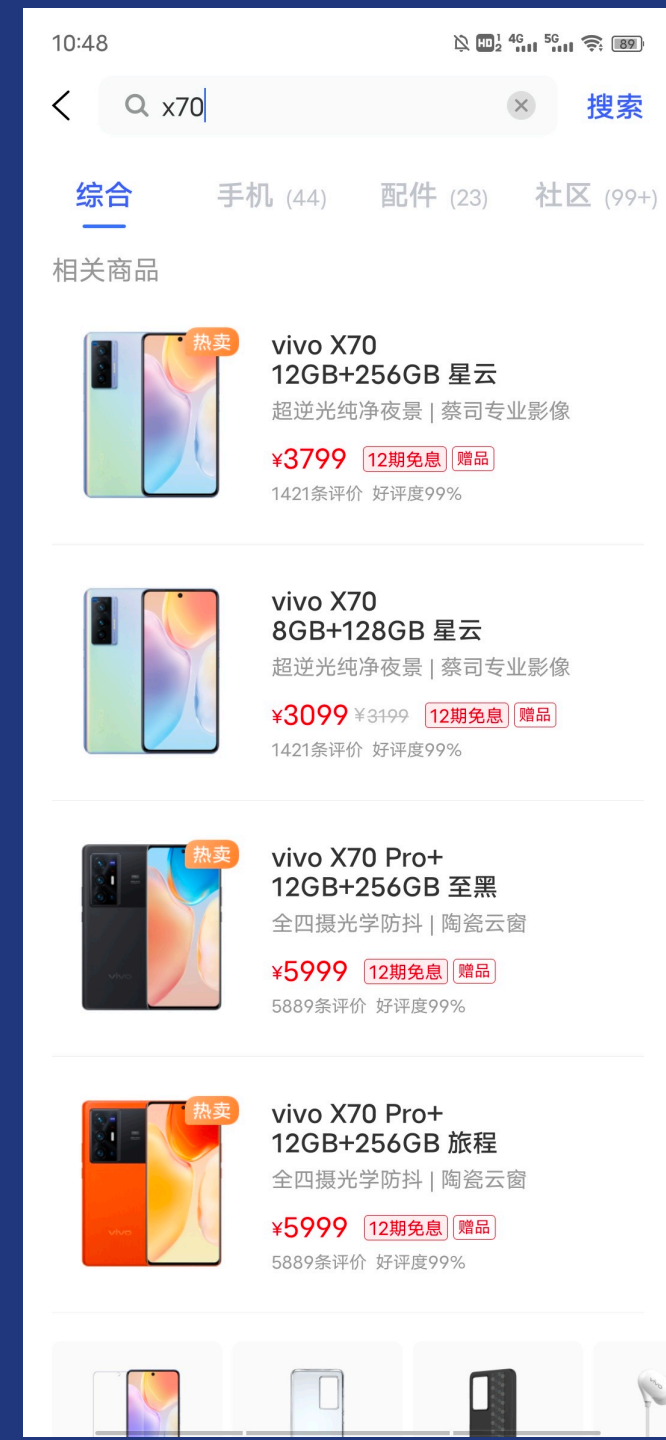
应用搜索
亿级PV



游戏搜索
千万级PV

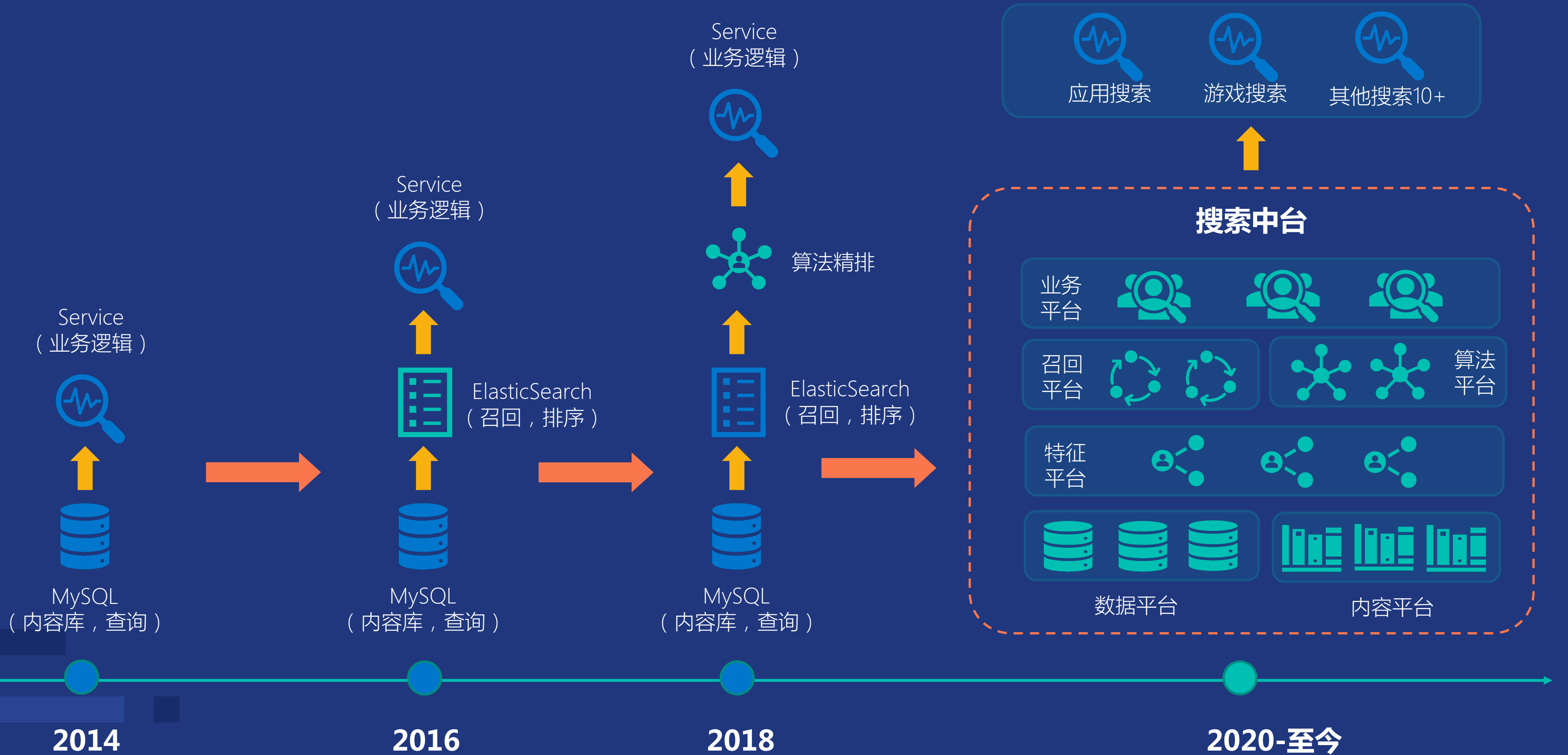


视频搜索
百万级PV

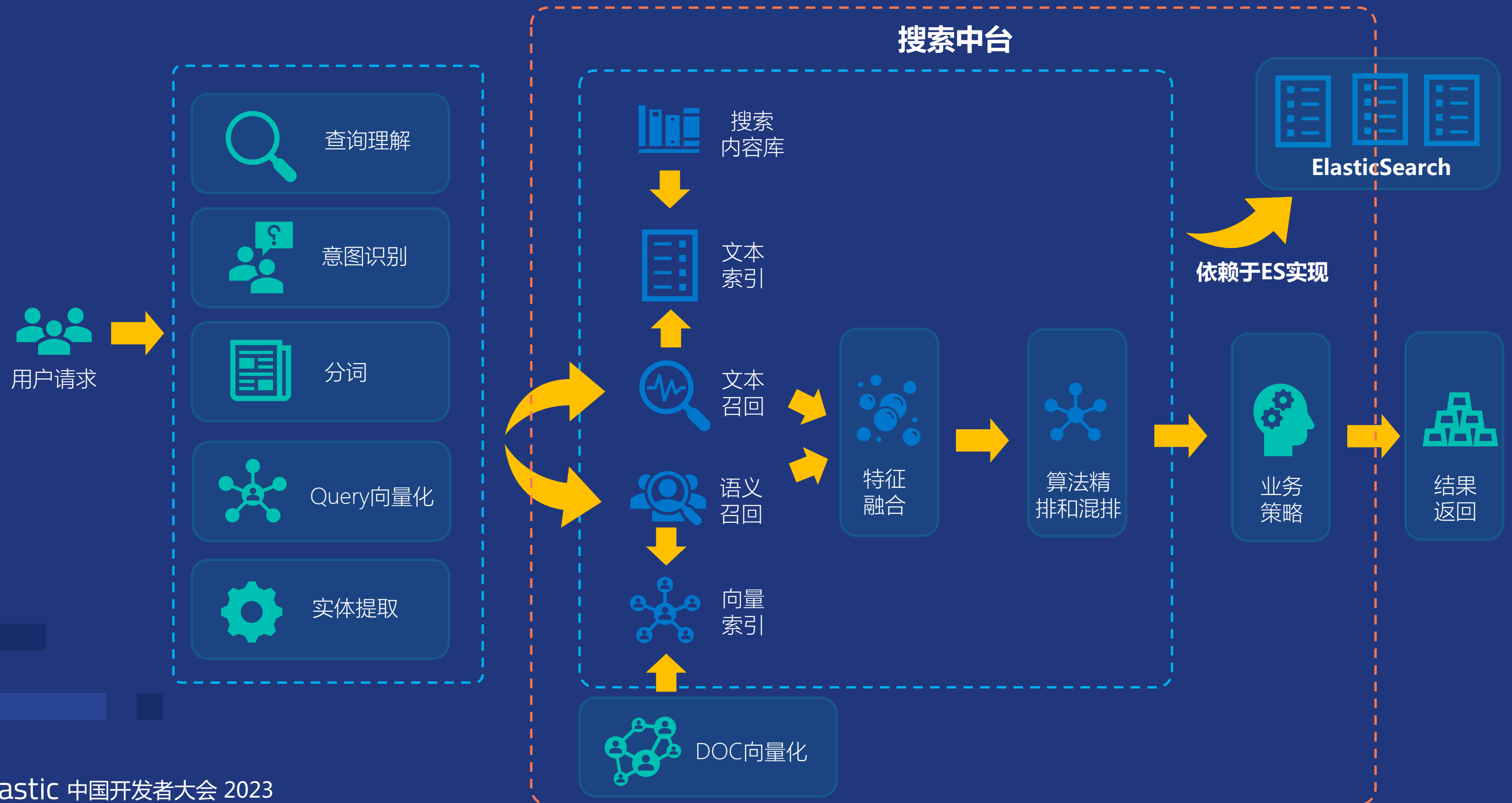


官网搜索
百万级PV

vivo 搜索架构演进



vivo 搜索中台业务流程

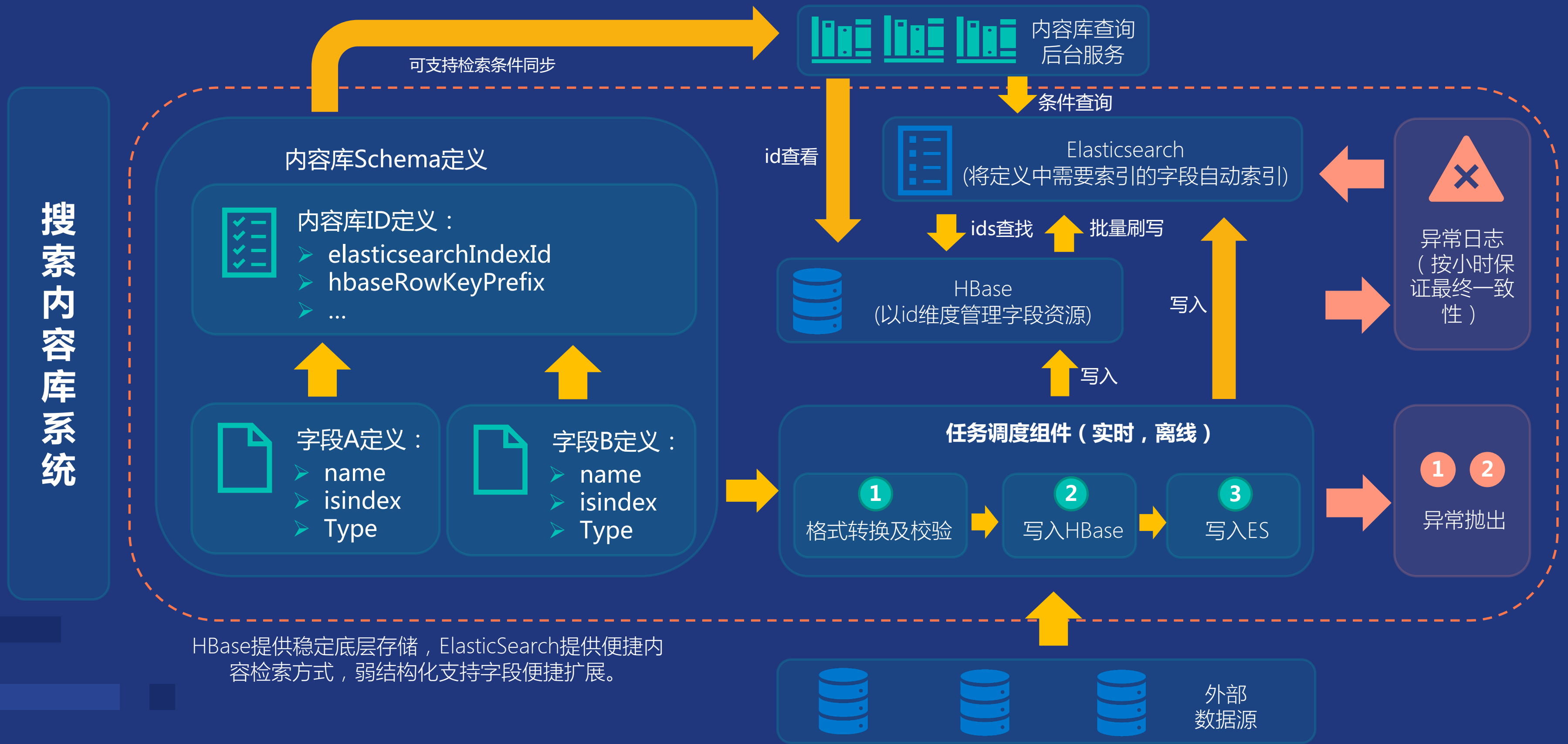


二、Elasticsearch在vivo搜索中台的实践

- Elasticsearch 在搜索中台内容库管理的实践
- Elasticsearch 在搜索中台文本召回的实践
- Elasticsearch 在搜索中台向量召回的实践
- Elasticsearch 在搜索中台精排的实践

Elasticsearch在搜索中台 内容库管理的实践

搜索内容库架构



利用Elasticsearch提供内容库索引的优势

垂类资源名称*

选择一级分类*

资源表名(不含中文)*

是否允许爬取*

相关业务方*

增量更新频率*

合并同类

ES集群id

资源整体定义

搜索内容库索引



字段具体定义

字段名 (field)	字段描述	类型	索引	索引
spuld		lor	是	lo
spuName		str	是	te
spuCode		str	是	te
skuld		lor	是	lo
skuName		str	是	te

查询条件组合

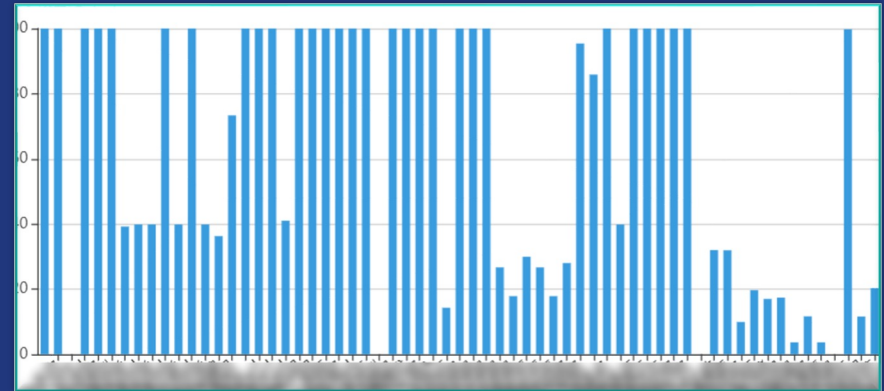
数据表格 | 质量分析

数据范围 已审核表 资源id 排序依据 无

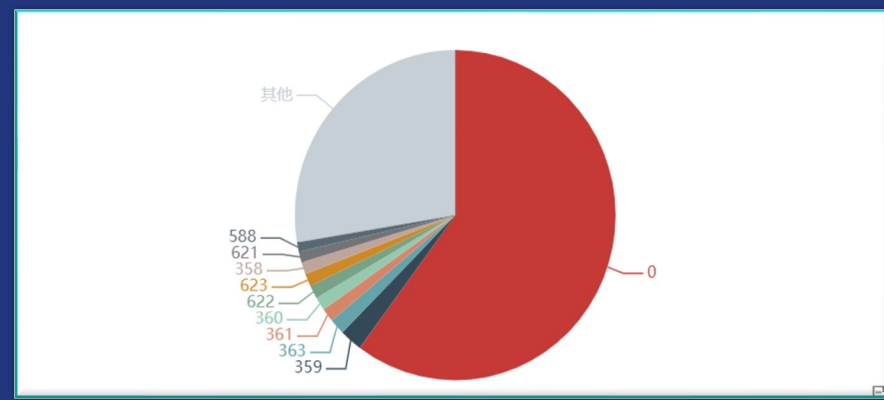
must 商品类型 分词匹配 电源

must 最低价格 分词匹配 100

字段存在性统计



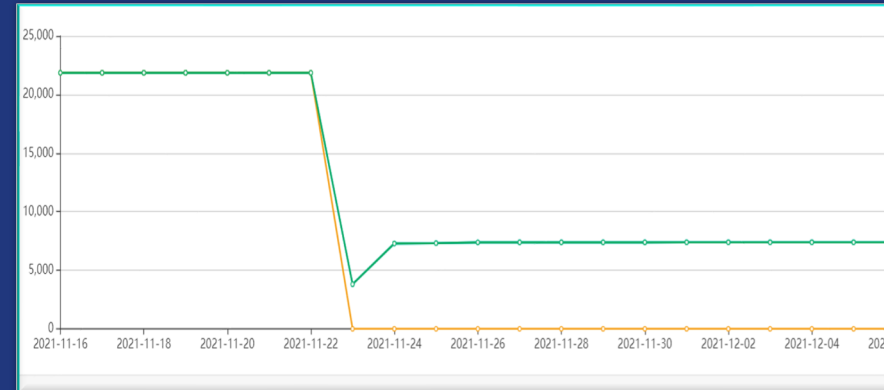
字段top值分布



数据条目查看和管理

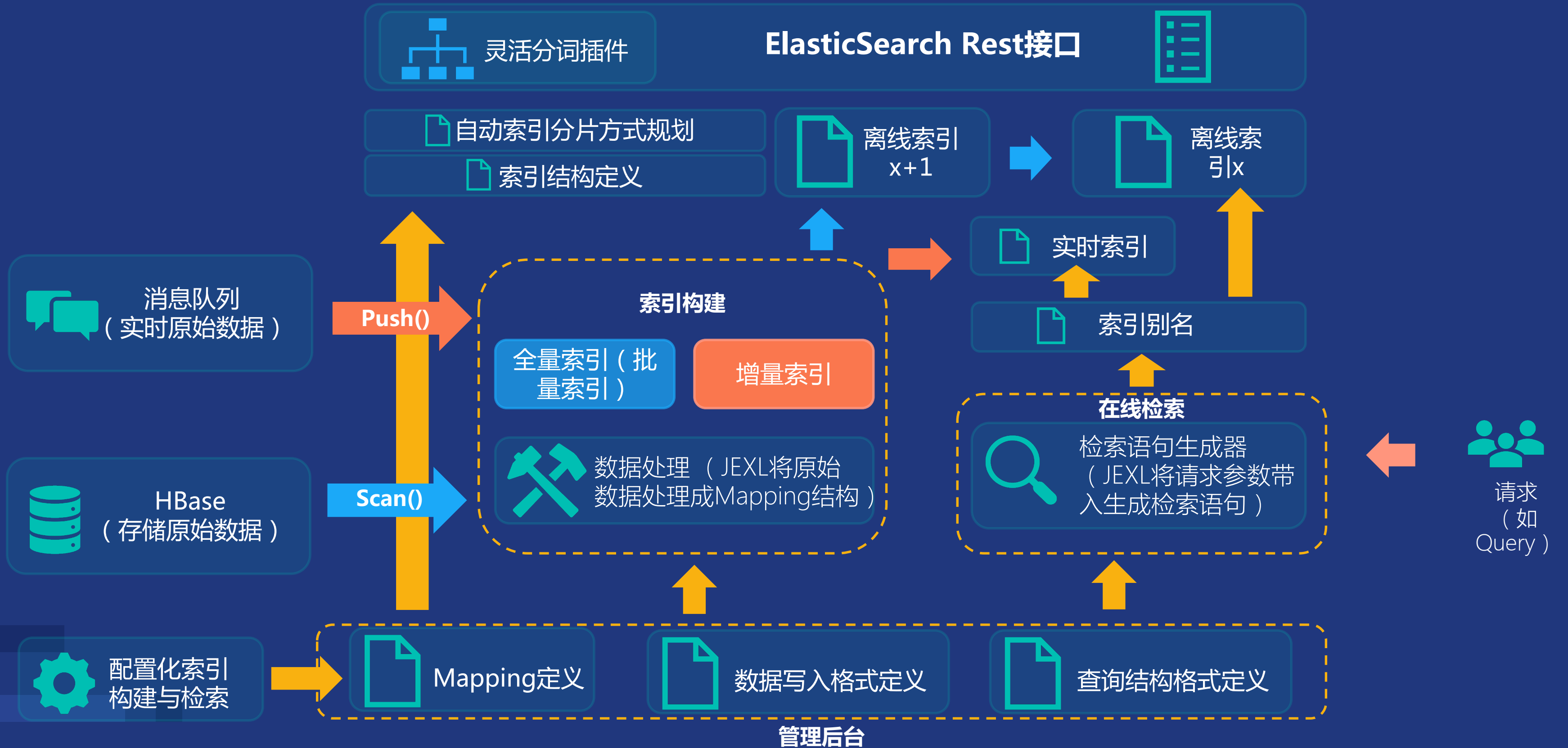
...	0	编辑 图透编辑 删除 加入黑名单 索引下架 操作日志
...	0	编辑 图透编辑 删除 加入黑名单 索引下架 操作日志
...	0	编辑 图透编辑 删除 加入黑名单 索引下架 操作日志
...	0	编辑 图透编辑 删除 加入黑名单 索引下架 操作日志
...	0	编辑 图透编辑 删除 加入黑名单 索引下架 操作日志
...	0	编辑 图透编辑 删除 加入黑名单 索引下架 操作日志
...	0	编辑 图透编辑 删除 加入黑名单 索引下架 操作日志
...	0	编辑 图透编辑 删除 加入黑名单 索引下架 操作日志

数据量变化 趋势监控



Elasticsearch在搜索中台 文本召回的实践

索引管理架构



配置化的索引构建与检索

请求：

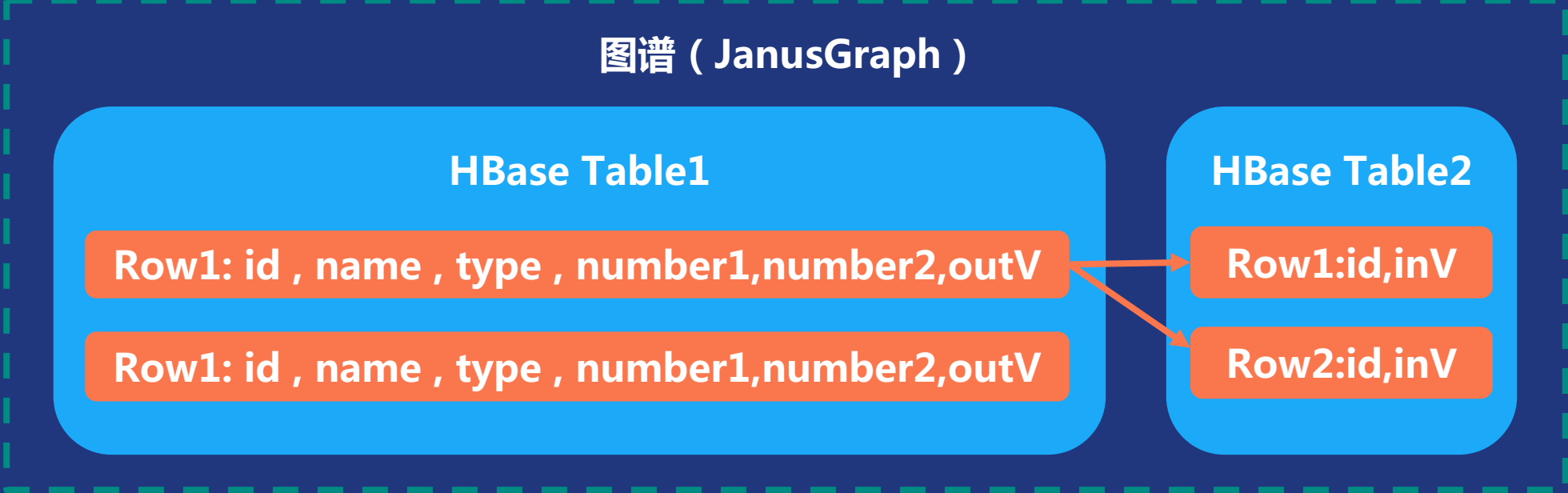
Query= "测试测试"

在线请求解析器 (JEXL)

```
{
  "query": {
    "bool": {
      "must": [
        {"match": {"name": "${obj.getQuery()}"}} ,
        {"term": {"type": "1"}}
      ],
      "should": []
    }
  },
  "from": 0,
  "size": 100,
  "aggs": {}
}
```

[{"match":{"name":"测试测试"}, {"term":{"type":"1"}}]

图谱 (JanusGraph)



Scan 批量变更

单条变更

Hook变更

离线入库数据解析器 (JEXL)

```
{
  "id": "${obj.targetSource().property('id')}",
  "name": "${obj.targetSource().property('name')}",
  "type": "${obj.targetSource().property('type')}",
  "totalNum": "${obj.targetSource().property('number1') + obj.targetSource().property('number2')}"
}
```

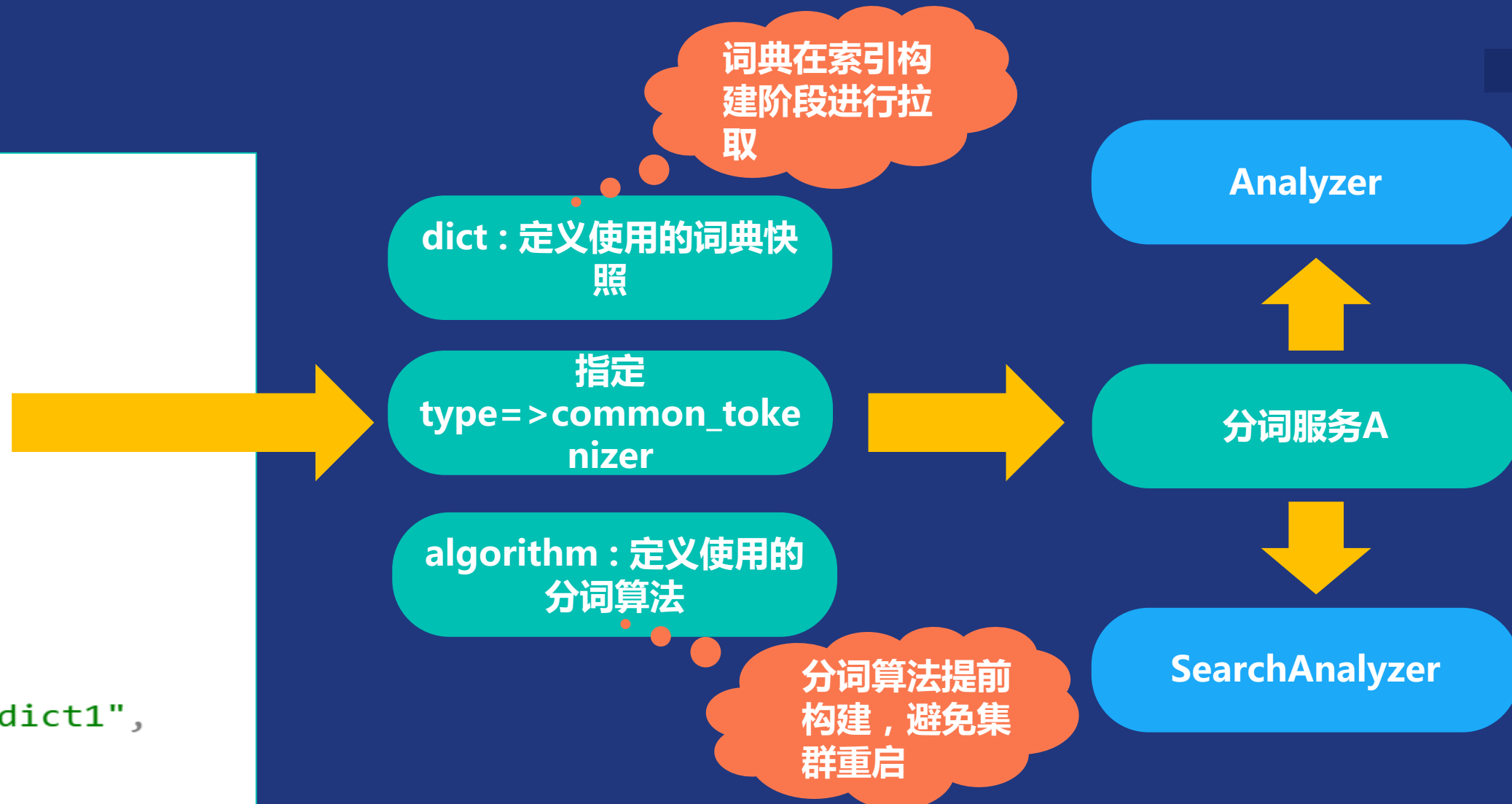
{"id":"xxx","name":"xxx","type":"xxx","totalNum":"xxx"}

灵活的分词

```

"settings": {
  "analysis": {
    "tokenizer": {
      "tokenizer_ik_smart_dict1": {
        "type": "common_tokenizer",
        "dict": "test_dict1",
        "algorithm": "ik_smart"
      }
    },
    "analyzer": {
      "analyzer_ik_smart_dict1": {
        "type": "custom",
        "char_filter": [],
        "tokenizer": "tokenizer_ik_smart_dict1",
        "filter": []
      }
    }
  }
}

```



- 问题1 : 分词后的term无法灵活更改 (比如干预分词term) 。
- 问题2 : query的分词term难以附加词权重
- 问题3 : match操作隐藏具体匹配细节, 问题不易定位

场景1 : 需要term的位置信息时

提前分词, 利用空格分开, 再利用空格分词器进行分词

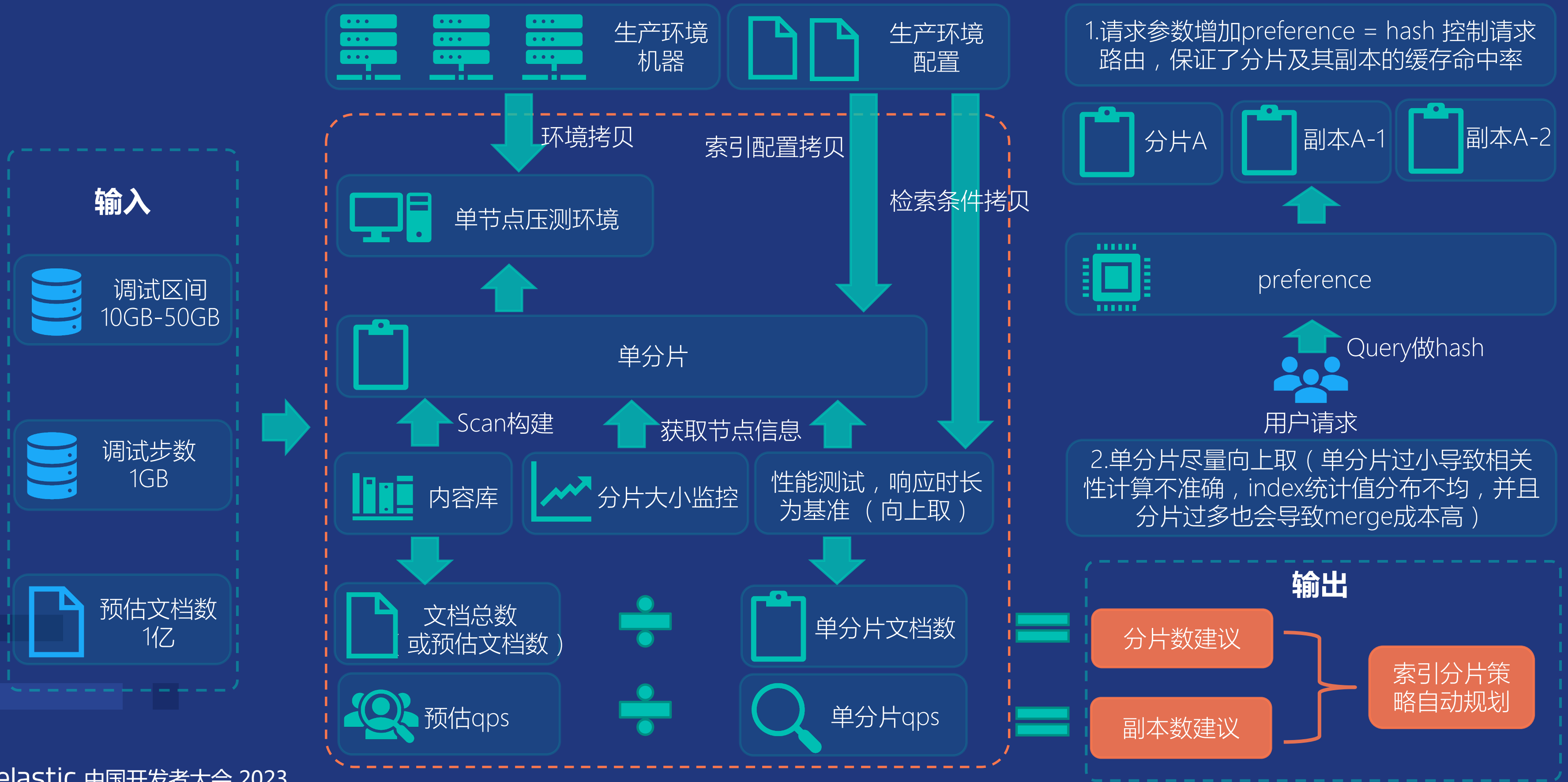
场景2 : 无需用到位置信息

提前分词, 使用keyword, 并将分词结果创建成数组的形式

场景3 : 需要对分词term附加权重

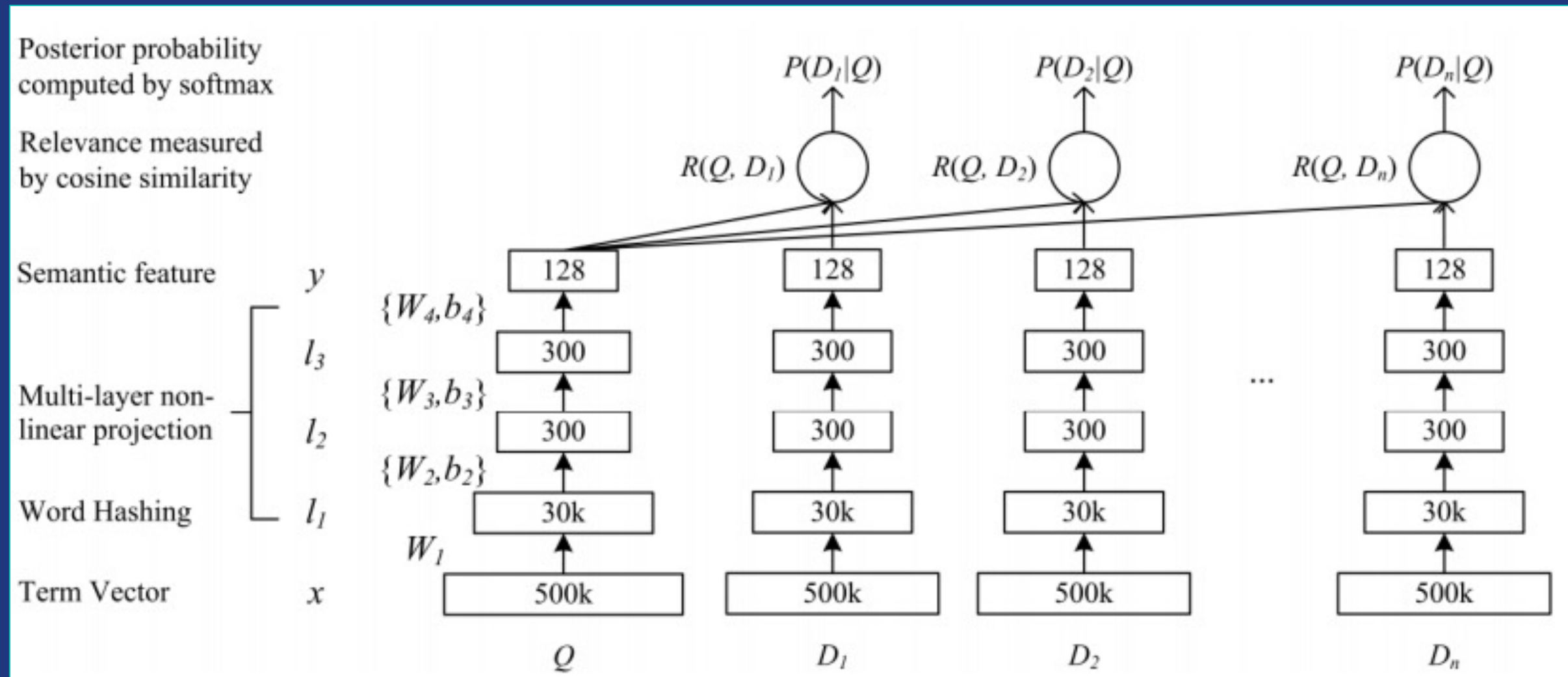
利用bool表达式+should+term查询替换match

自动索引分片规划



Elasticsearch在搜索中台 向量召回的实践

向量索引的构建与检索基本架构

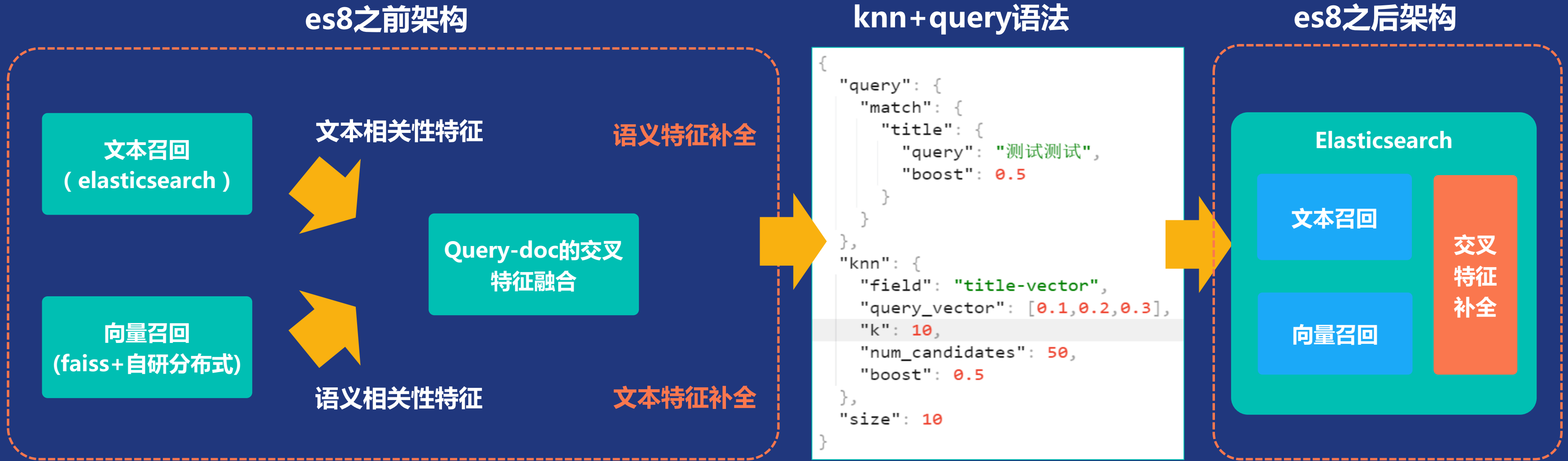


DSSM架构

- Query->Doc的转换概率 (用户行为)
- Embedding过程
- Query和doc(比如title或其他可向量化表示的特征)



» Elasticsearch 向量检索优化特征融合链路

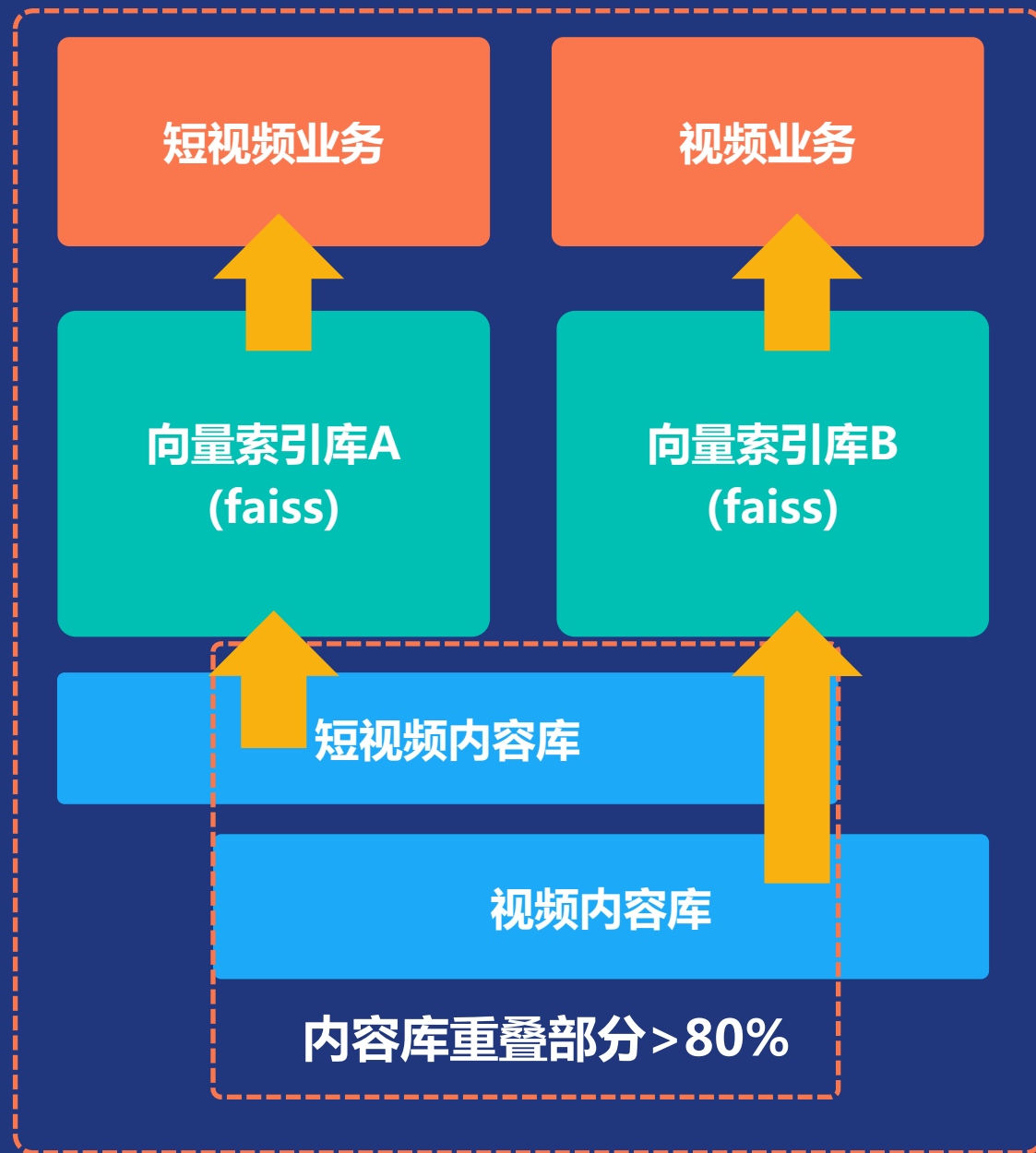


Query-doc的相关性交叉特征
需要再重新回捞计算实现成本
高，一致性难保证

Query-doc的相关性特
征都交给es计算，实现
简单，一致性容易保证

» Elasticsearch 向量检索预过滤实践

利用传统向量引擎



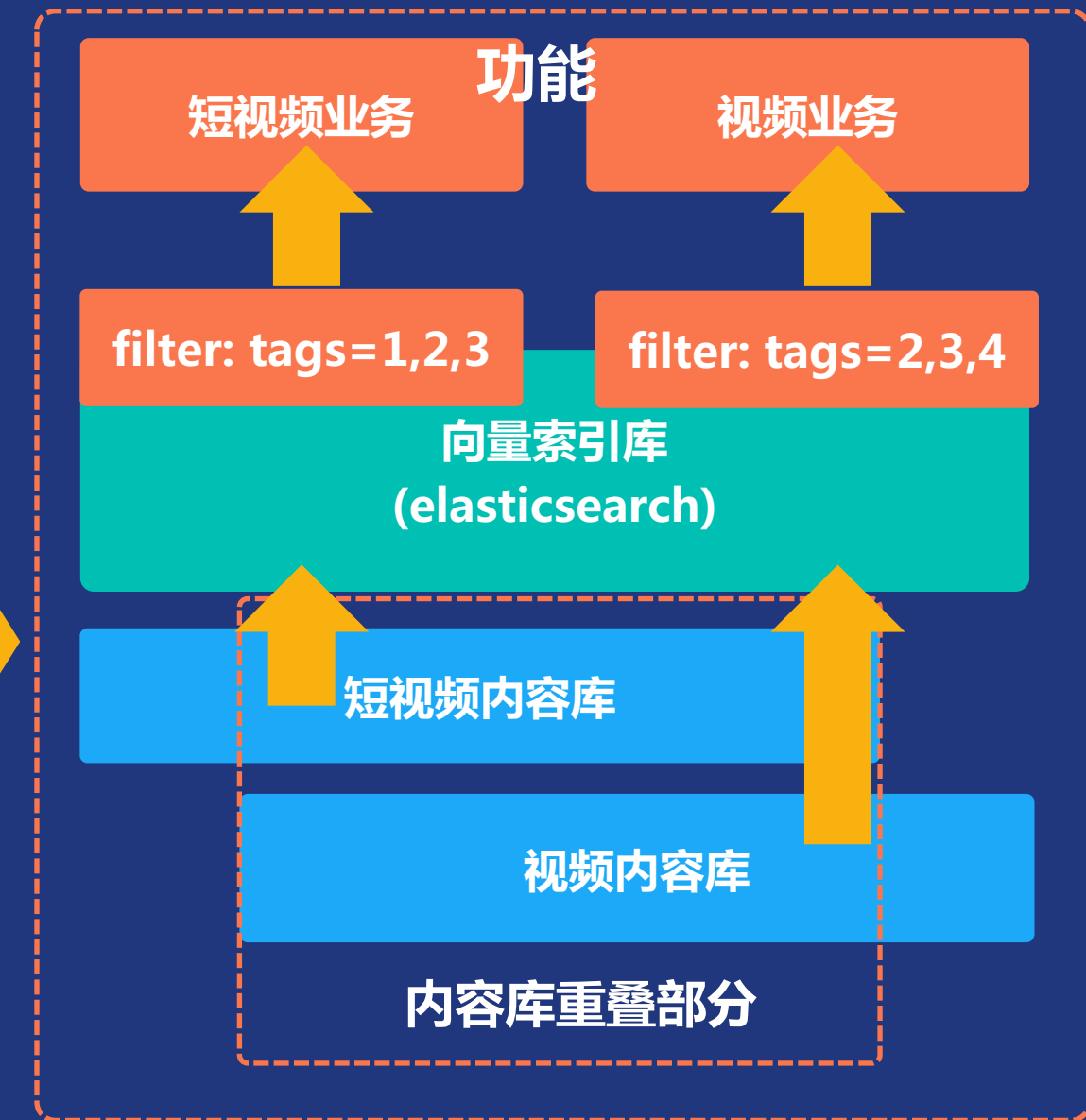
传统向量引擎，不保存正排信息也不支持过滤，如果后置过滤则导致召回结果变少甚至出现无结果的情况，尽管模型一致且内容接近，仍然部署了两个索引库

Knn过滤配置

```
{
  "knn": {
    "field": "title-vector",
    "query_vector": [0.1, 0.2, 0.3],
    "k": 10,
    "num_candidates": 100,
    "filter": {
      "terms": {
        "tags": [1, 2, 3]
      }
    }
  }
}
```

es向量索引支持在检索过程中过滤（存放正排的天然优势）

利用es向量检索过滤

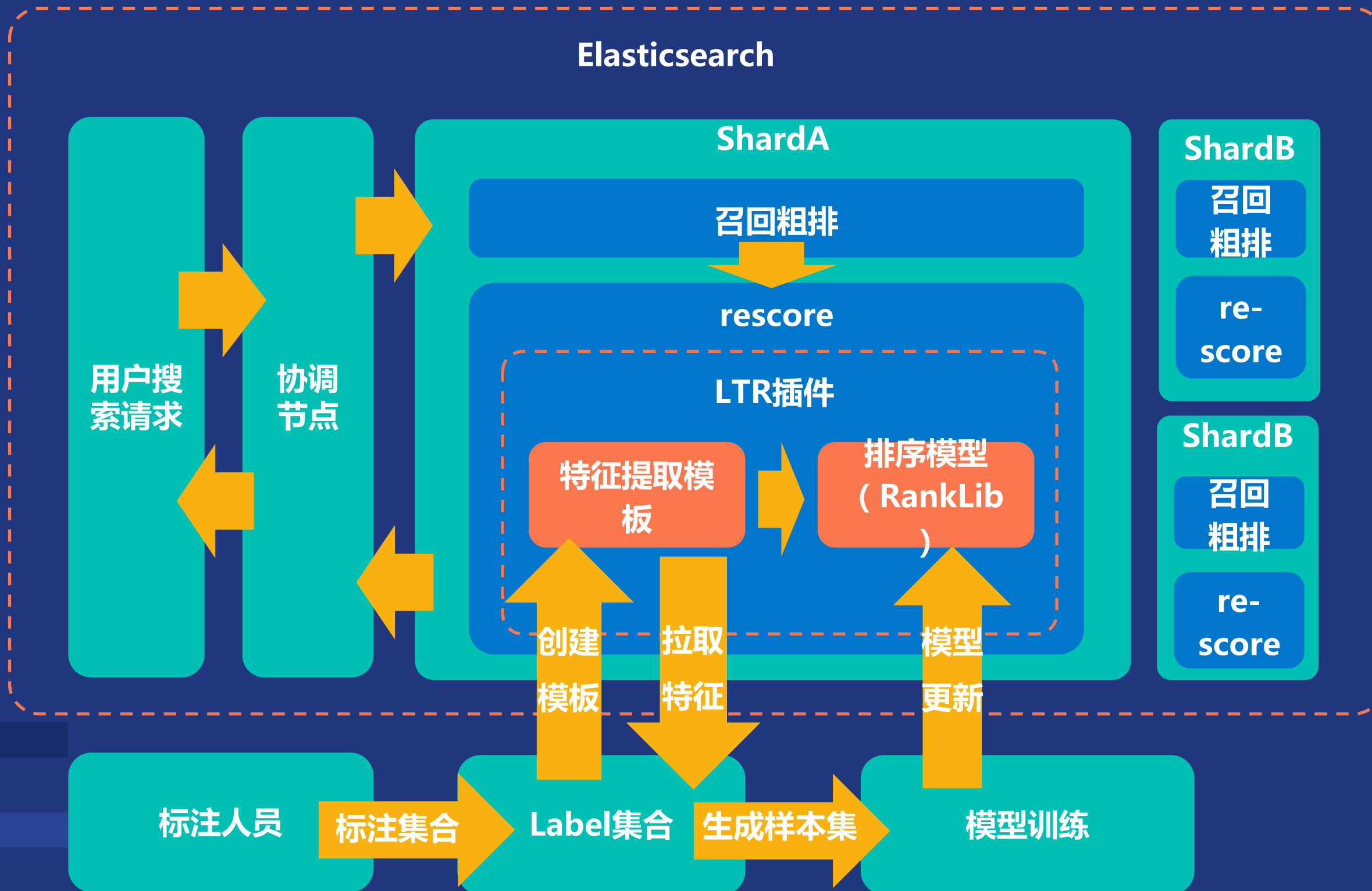


统一索引库，在上层进行过滤操作。索引库统一后节省了大量的维护成本

Elasticsearch在搜索中台 精排的实践

利用rescore插件完成精排

使用LTR插件解决算法排序问题(传统机器学习模型)



优点：

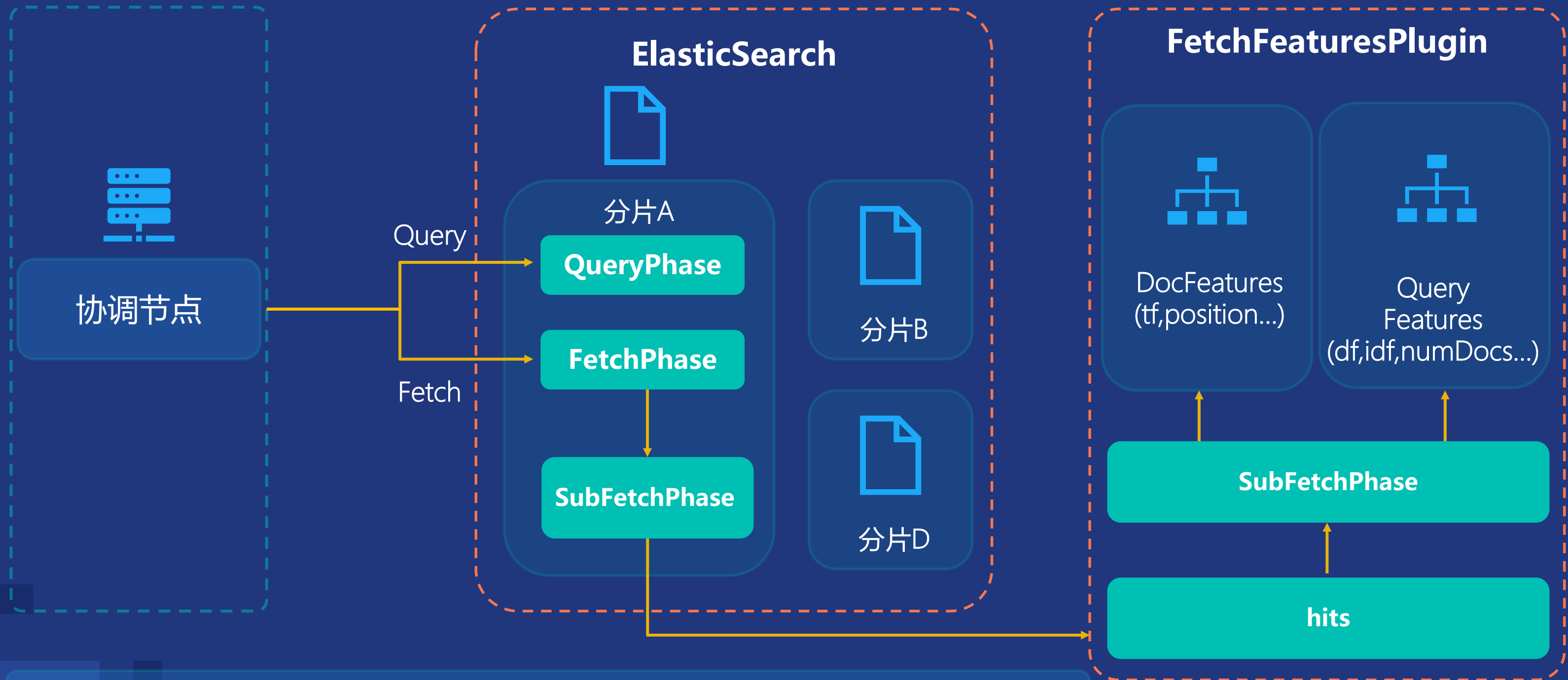
1. 提供基础特征 (痛点)
2. 增加并行度，减少数据传输量

缺点：

1. 绑定太死，多路召回不便同排
2. 增加多余的性能负载 (rescore)
3. 接入非java预估服务实现成本高

提取相关性特征

开发特征提取插件，解决基础相关性特征无法透出的痛点。



放在Fetch阶段，减少机器负载，和网络传输量，将结果传回调用服务再进行算法调用，提升灵活性。

特征提取效果演示

```
"query": {
  "bool": {
    "should": [
      {
        "bool": {
          "_name": "title_terms",
          "should": []
        }
      }
    ]
  },
  "ext": {
    "feature_log": {
      "items": [
        {
          "feature_name": "title_query_features",
          "named_query": "title_terms"
        },
        {
          "feature_name": "title_field_features",
          "field_name": "title"
        }
      ]
    }
  }
},
"from": 0,
"size": 80,
```

```
"title": "分析搜索引擎对网站的....",
"_score": 64.90771,
"_fields": {
  "log_features": [
    {
      "title_query_features": {
        "搜索引擎_df": 2047,
        "搜索引擎_tf": 1,
        "搜索引擎_numDocs": 2938397,
        "搜索引擎_position": [
          1
        ],
        "origin": 15.411291,
        "搜索引擎_idf": 8.268756
      },
      "title_field_features": {
        "adtn": 16.142025,
        "dtn_dr": 11,
        "adtn_dr": 14.749546,
        "dtn": 11
      }
    }
  ]
},
"status": 1
```

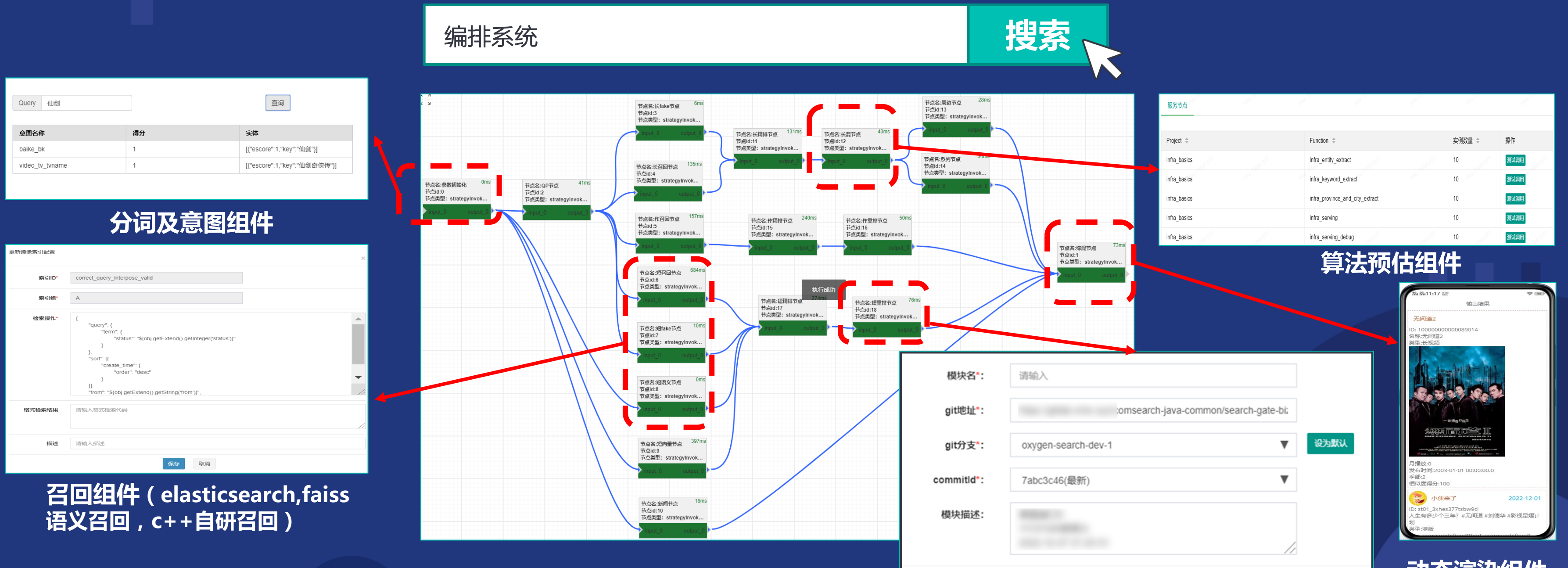
查询特征

文档字段特征，比如dtn
(document_term_num)
文档该字段的term总数

三、搜索中台业务应用总结

搜索中台编排系统

支撑若干业务的搜索中台编排系统





感谢观看





专业、垂直、纯粹的 Elastic 开源技术交流社区

<https://elasticsearch.cn/>