

# ES集群在雷达大数据平台的演进

# 目录

CONTENTS

- 1 雷达平台简介
- 2 ES集群的演进
- 3 总结





# 雷达大数据平台项目背景

当前社会快消行业，业务的发展、行业的竞争、市场趋势的变化，都非常的快和复杂，这就决定了，品牌需要以更快、更全面的速度去了解、洞察消费者

这中间会遇到各种问题，由于信息或数据的多和复杂，品牌会面临决策难、信息割裂、时间紧迫、信息安全等问题

所以需要建立一个系统化的、体系化的大数据应用平台，来满足业务快速发展的需要



# 雷达大数据平台业务模块

品牌

了解品牌网络舆情、热点话题及消费者画像等方面的信息，为品牌营销管理提供策略

产品

发现新兴趋势，创造新产品，以及监测消费者对本品和竞品的产品体验反馈，发现痛点优化产品

媒介

了解品牌赞助的综艺节目和签约的代言人对品牌影响效果，是否有助于提升品牌的认知度

渠道

监控各电商平台内品牌产品及服务评价，为产品优化、电商平台深度沟通合作提供信息支持

市场  
与消  
费者

对用户画像的刻画，帮助客户发现潜在的目标消费人群



# 雷达大数据平台数据存储概况

## 数据量级

- 十亿+

## 数据类型

- 线下结构化文本数据
- 线上非结构化文本数据
- 图片数据

## 数据种类

- 新闻、论坛、微博、微信、电商、视频、贴吧、博客等

## 机器配置

- CPU: 8 cores
- Mem: 64G
- Disk: 1T

## ES集群

- 节点: 18
- 索引: 4
- 分片: 500+

# 目录

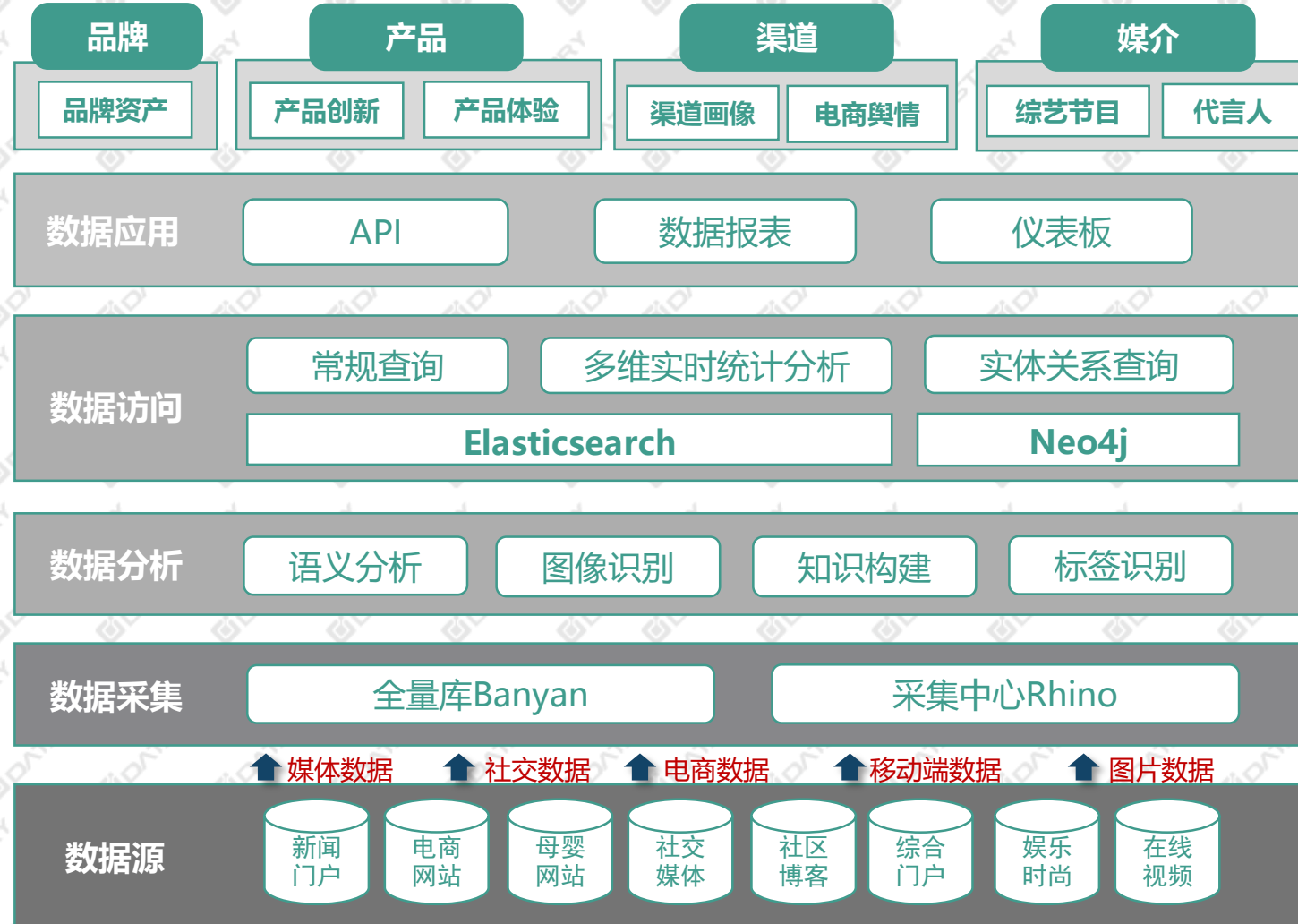
CONTENTS

- 1 雷达平台简介
- 2 **ES**集群的演进
- 3 总结





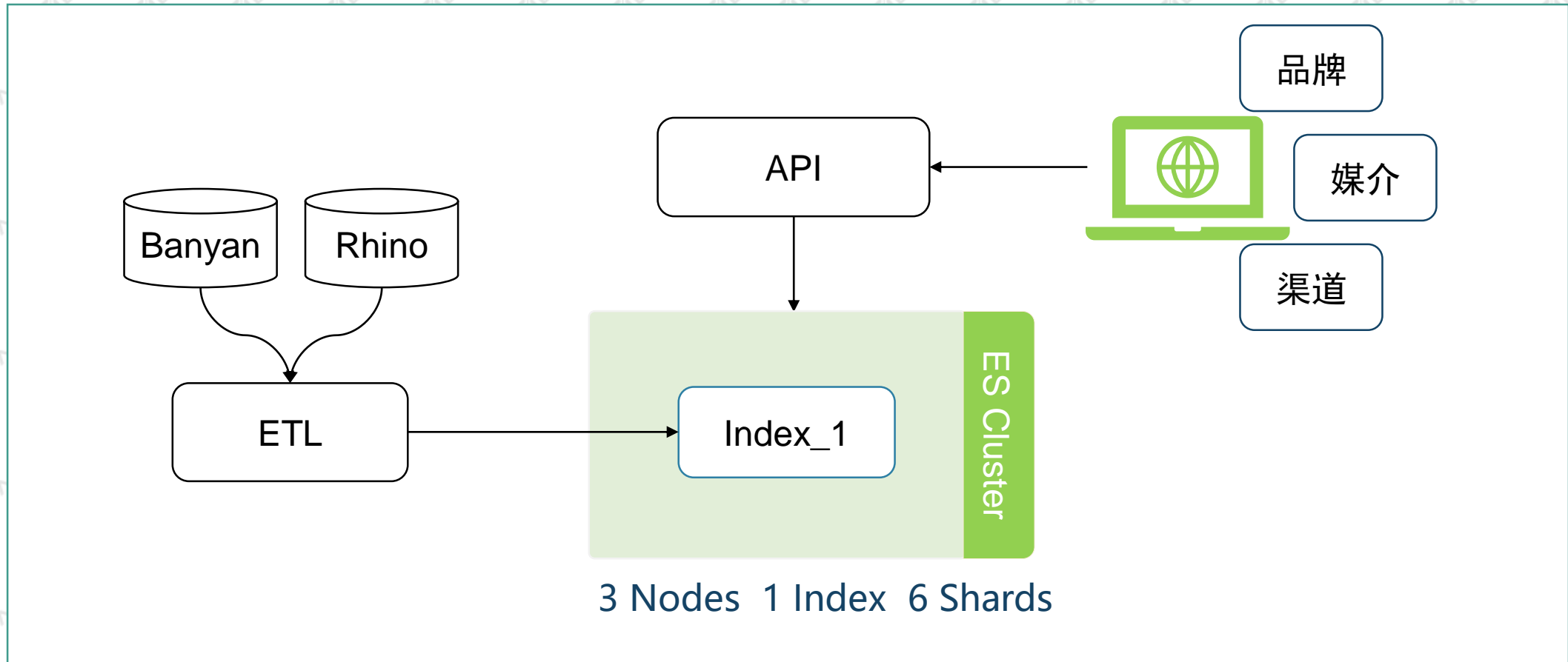
# 雷达大数据平台架构





# 雷达大数据平台ES集群演进 V1

客户需求1：把我司的品牌、赞助的节目、品牌的代言人、电商平台的评论等相关数据实时监控起来

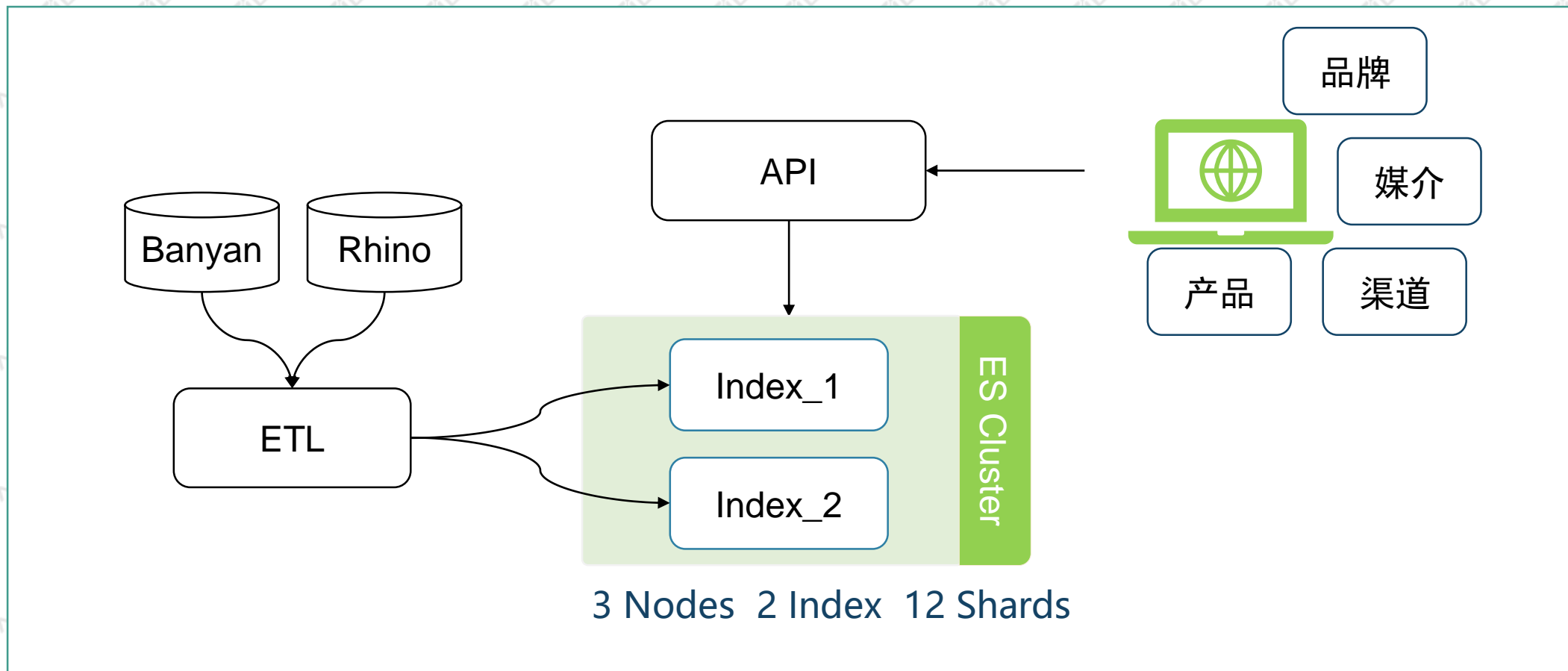






# 雷达大数据平台ES集群演进 V2

客户需求2：把消费者反馈的产品体验问题实时监控起来，后续要对产品做优化改进





# 雷达大数据平台ES集群演进 V2

## V2 存在的问题

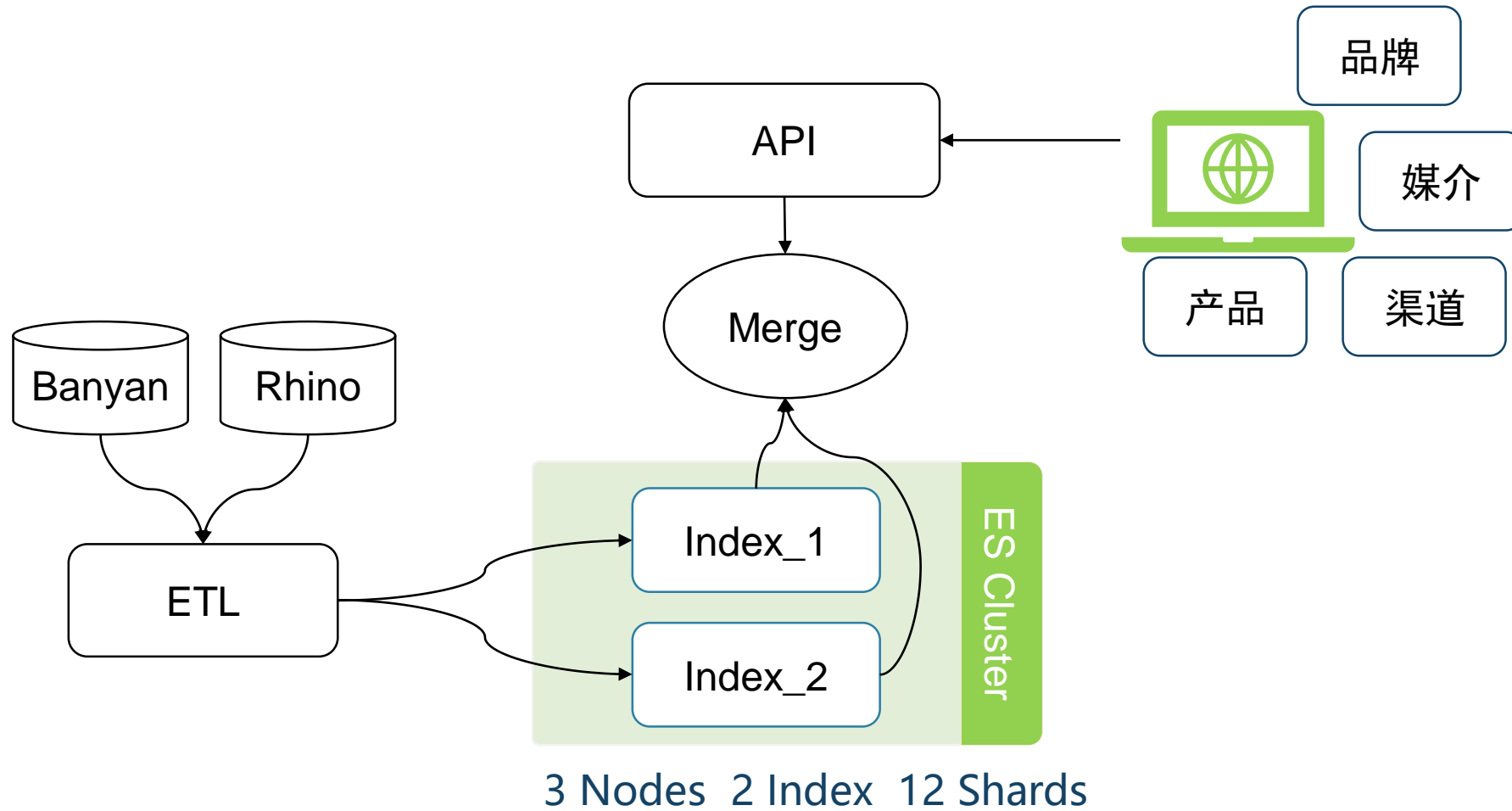
- 跨库查询

客户：我要对Index\_1和Index\_2两个索引中关联的数据做横向查询对比

- 方案一：合并两个索引并成一个更大的宽表，从合并的索引中查询
- 方案二：分别从两个索引查询，把结果Merge展示



# 雷达大数据平台ES集群演进 V2.1

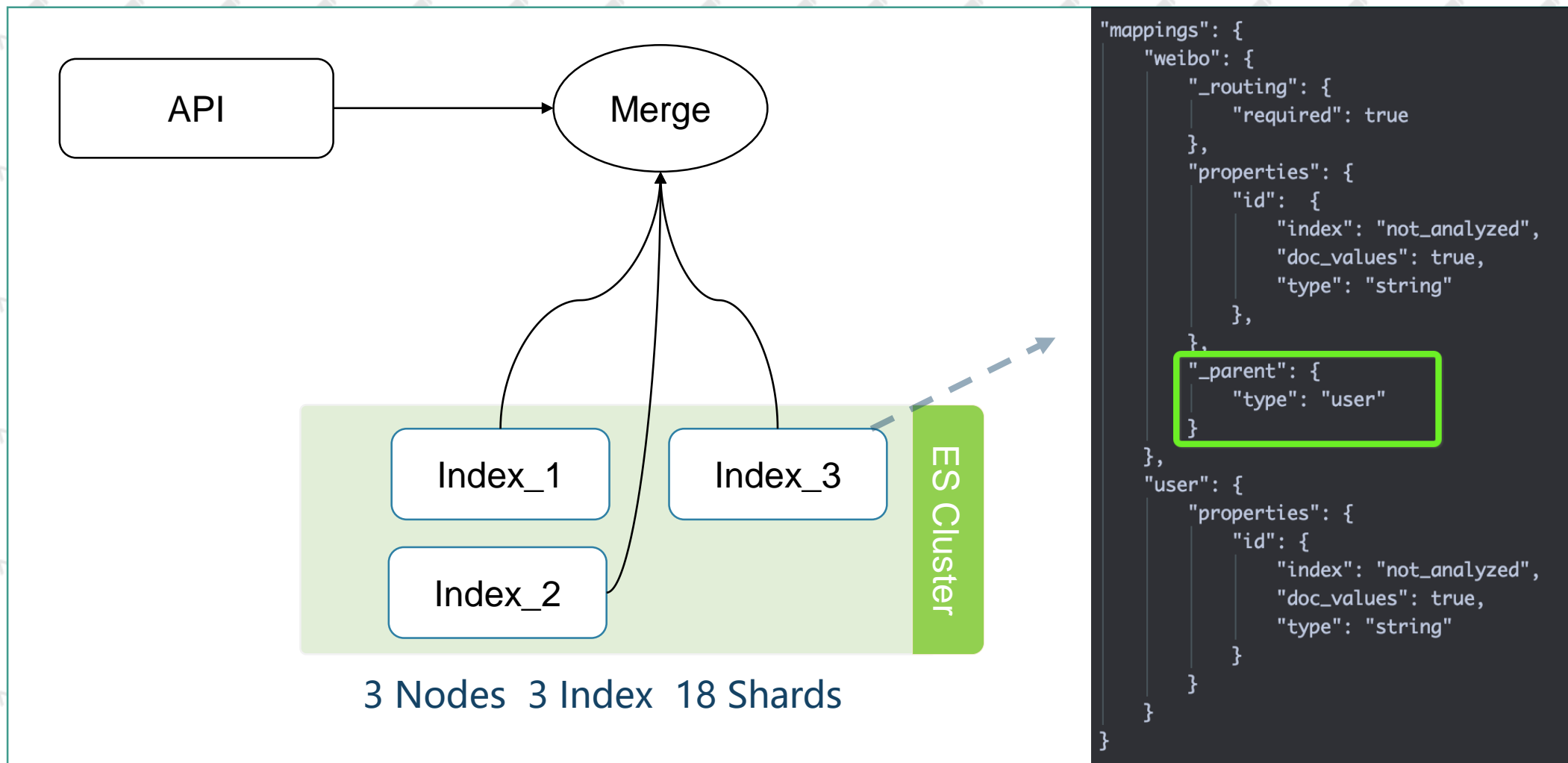






# 雷达大数据平台ES集群演进 V3

客户需求3：建设人群库，通过分析用户画像找出潜在的消费者





# 雷达大数据平台ES集群演进 V3

## V3 存在的问题

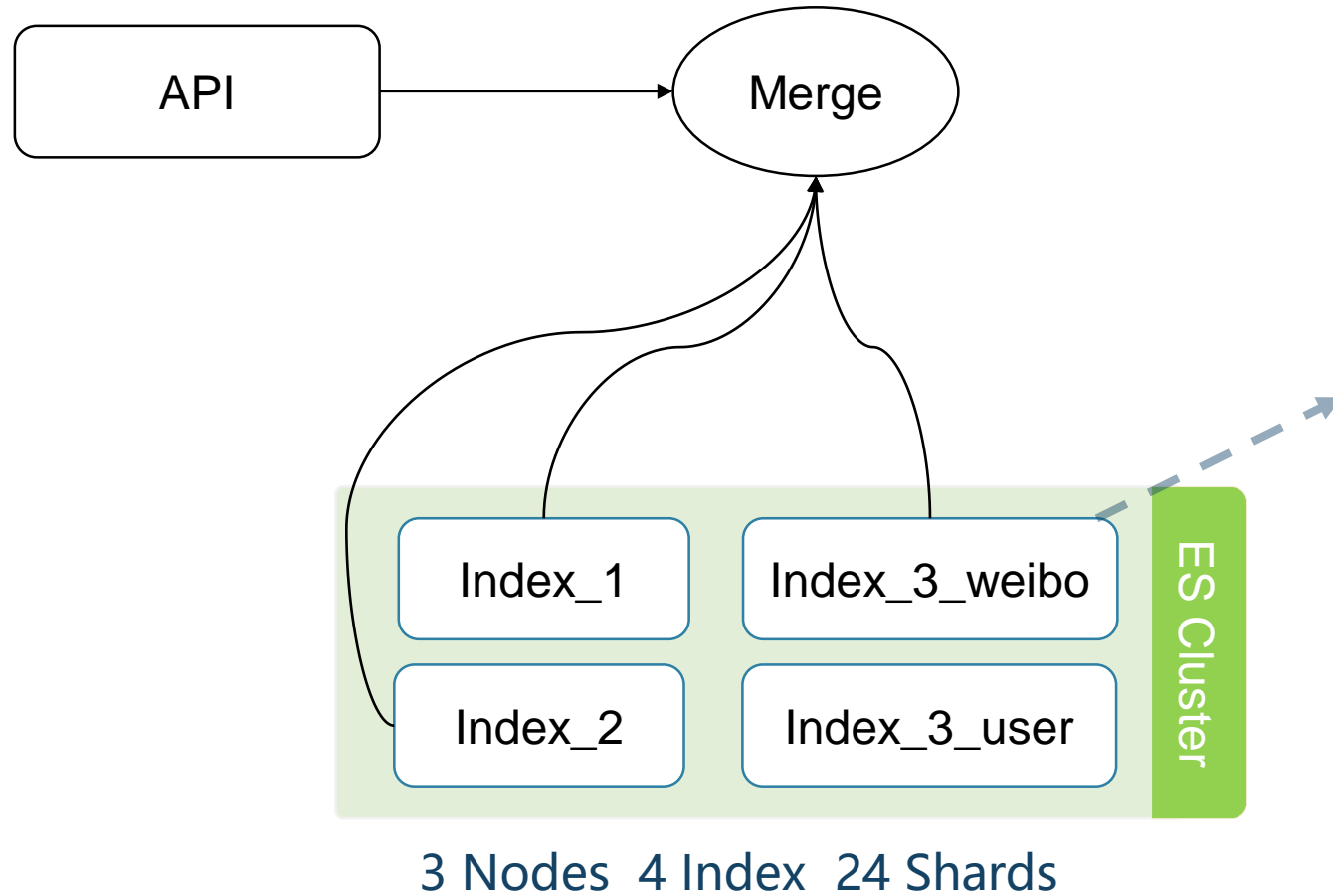
- 父子文档查询性能低
- 相同名称的字段类型不一致

- 方案：

- 拆分父子文档为两个索引，在子文档中冗余存储必要的查询字段；



# 雷达大数据平台ES集群演进 V3.1



```
"mappings": {
  "weibo": {
    "properties": {
      "id": {
        "index": "not_analyzed",
        "doc_values": true,
        "type": "string"
      },
      "uid": {
        "index": "not_analyzed",
        "doc_values": true,
        "type": "string"
      }
    }
  }
}

"mappings": {
  "user": {
    "properties": {
      "id": {
        "index": "not_analyzed",
        "doc_values": true,
        "type": "string"
      }
    }
  }
}
```





# 雷达大数据平台ES集群演进 V3.1

## V3.1 还可能存在的坑

- 当业务需要对子文档做查询，涉及到父文档的字段，而这个字段冗余到子文档中
- 方案一：冗余全部的父文档字段到子文档中
  - 引来另一个问题：如果某个字段内容特别大，会占用很大的存储空间
- 方案二：不依赖ES查询，考虑其他外部替代方案



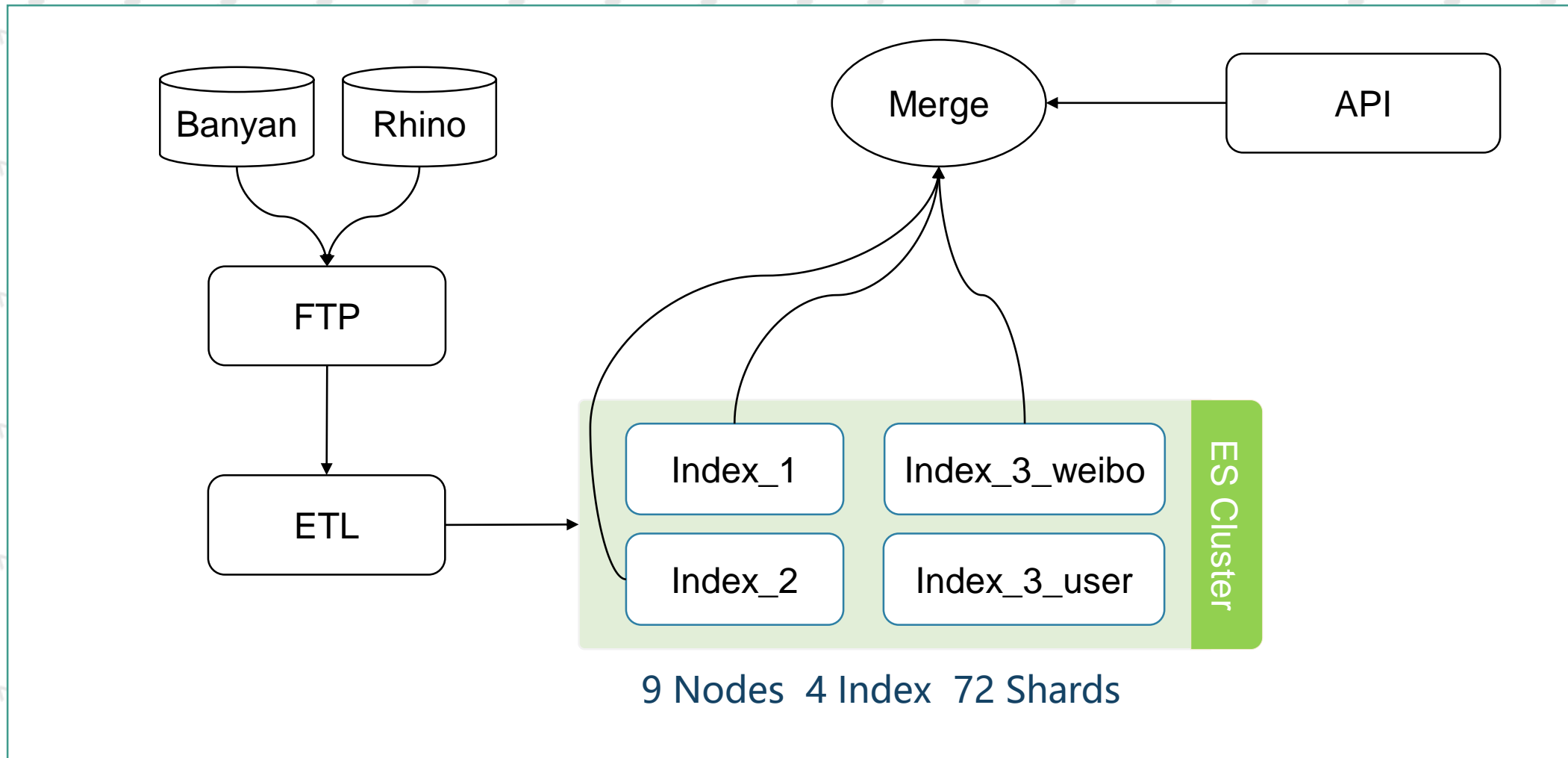
# 雷达大数据平台ES集群演进 V3.1

- 实施过程中遇到的另一个坑：  
对用户标签字段聚合，排查发现这个字段的 `doc_values = false`，  
直接导致ES集群Heap负载过高
- 方案：对于需要做 Agg 的字段，设置 `doc_values = true`



# 雷达大数据平台ES集群演进 V4

客户需求4：私有化部署，数据放到客户内部平台作为资产







部署完成，查询速度挺快的.....



但是，几个月之后.....

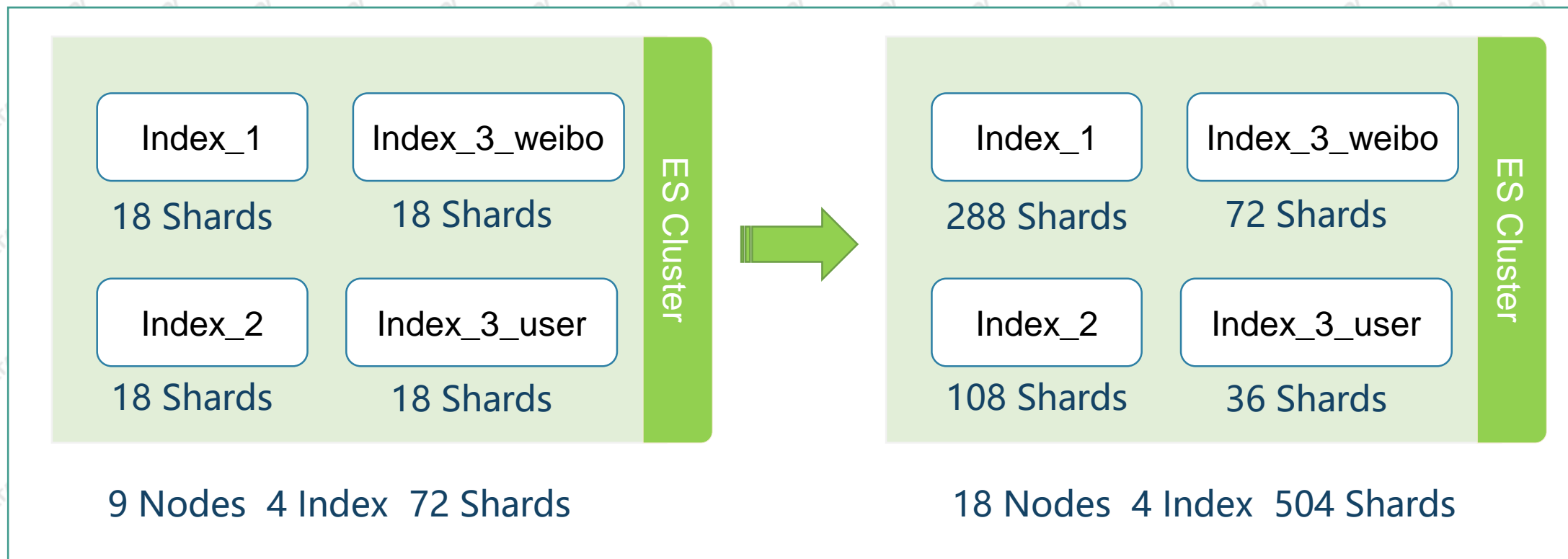


# 雷达大数据平台ES集群演进 V5

故障需求：生产平台图表查询速度降低（而且当时磁盘使用率达到70%）

原因：分片过大，单个分片达200G（远大于官方推荐的30G）

方案：增加分片数，扩展节点数







# 雷达大数据平台ES集群演进 V5

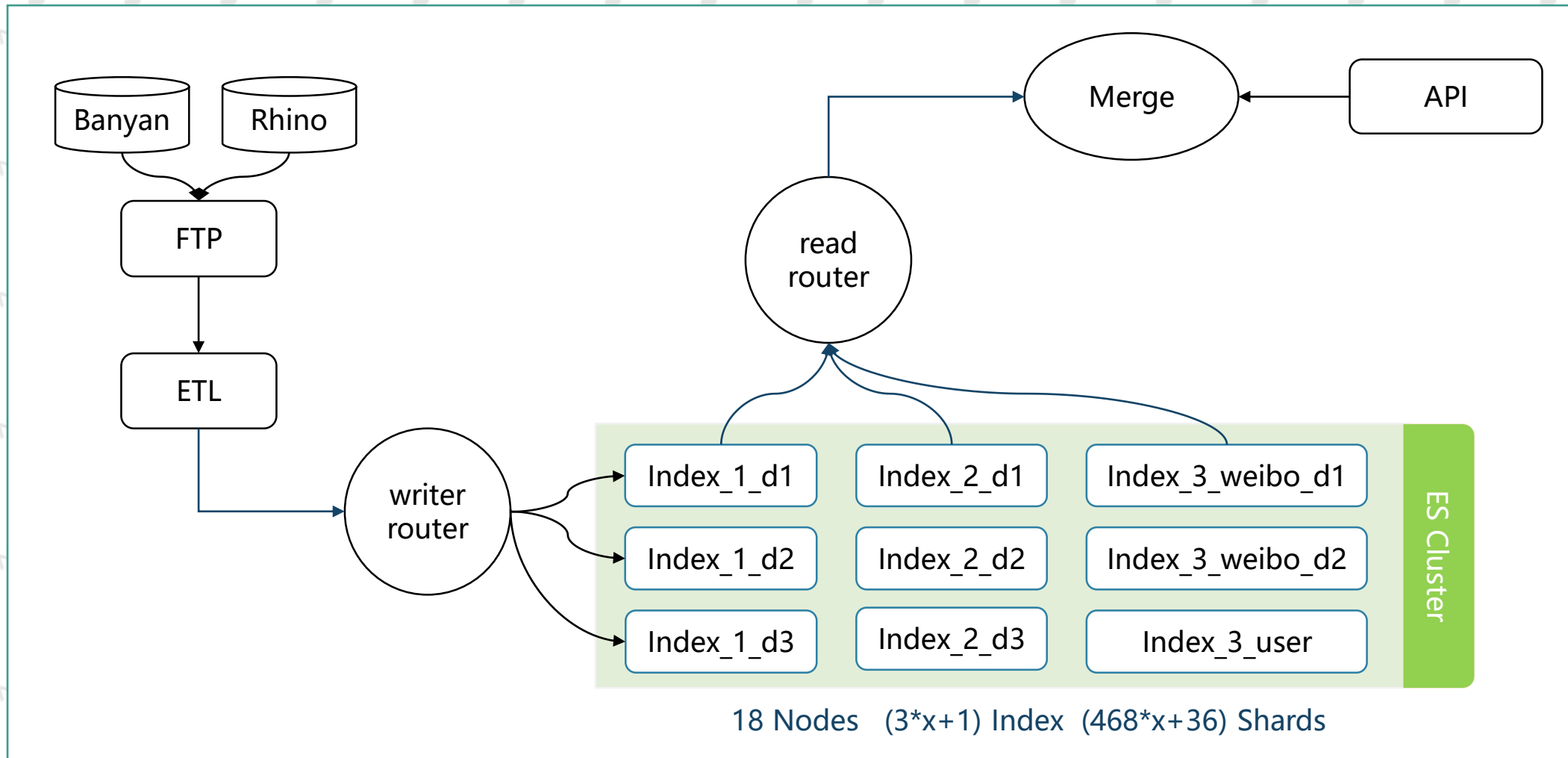
## V5 存在的问题

- 单个索引过大（目前最大的索引大小1.8T，不包括副本），导致单个索引分片过多
- 方案：按时间水平扩展切分索引



# 雷达大数据平台ES集群演进 V6

优化需求：水平扩展切分索引



# 目录

CONTENTS

- 1 雷达平台简介
- 2 ES集群的演进
- 3 总结



# 总结

- 当业务系统的查询性能需求远高于索引性能时，尽量减少父子文档，使用Nested嵌套类型或者冗余存储来实现父子关系效果；
- 对于需要做聚合操作的字段，设置doc\_values=true；
- 官方推荐一个索引只设置一个type，如果业务要设置两个及以上type，则相同名称的字段类型必须一致。
- 索引分片大小尽量在30G左右；
- 索引分片数设置不能过多，根据实际的业务数据量和机器配置评估；
- 单个索引过大后，按照一定的规则分库存储和查询。



# Q&A

