



十亿级人脸搜索的实践和优化

刘刚, CTO

谱时智能云, 2023/04/08

分享嘉宾

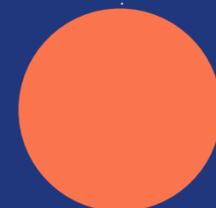


刘刚，谱时智能云 CTO

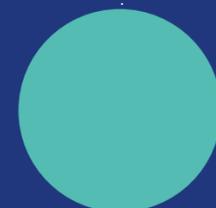
2014年Elasticsearch中国技术分享会分享嘉宾



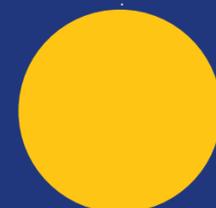
我们的业务



我们的问题



我们的方案



我们的规划

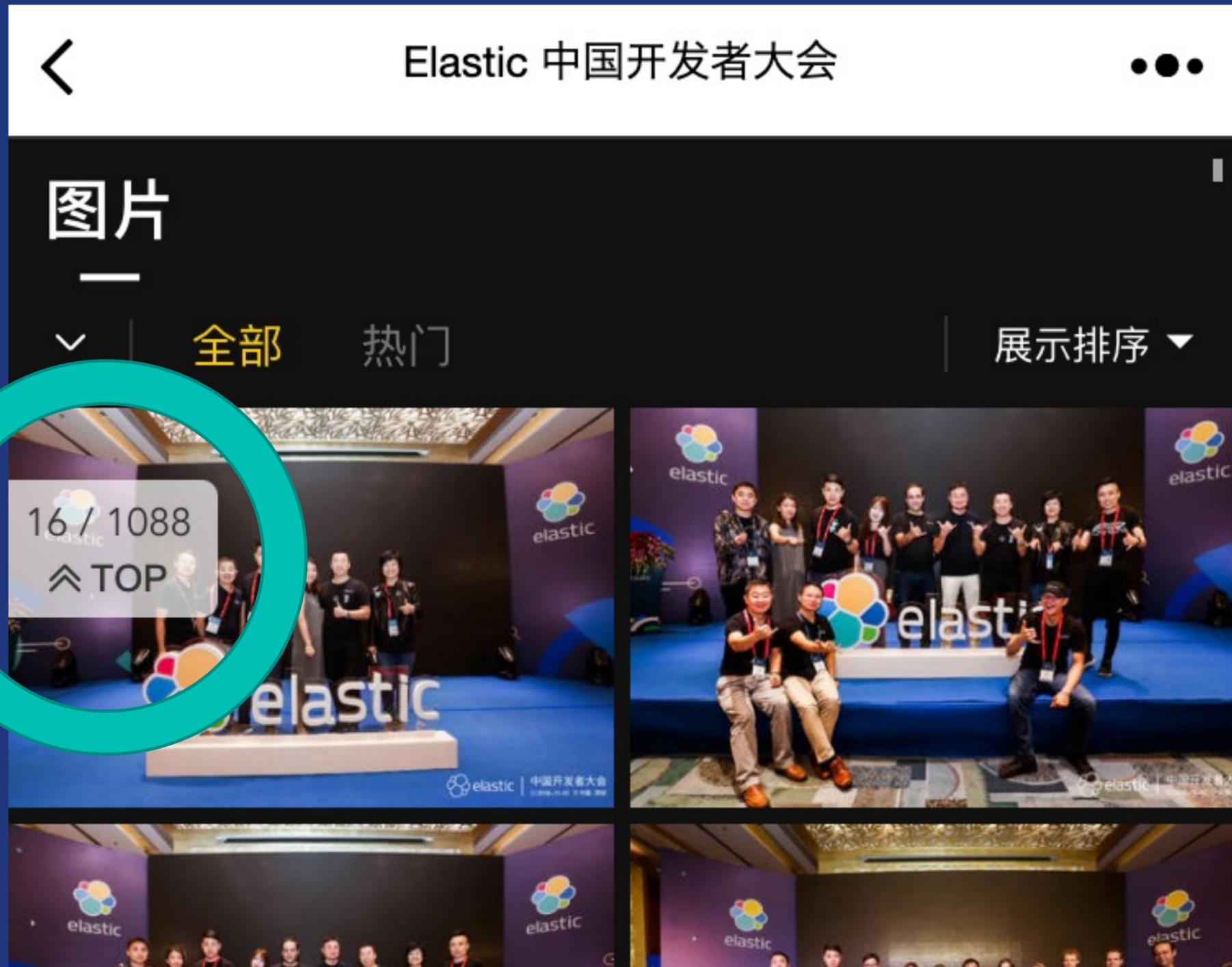
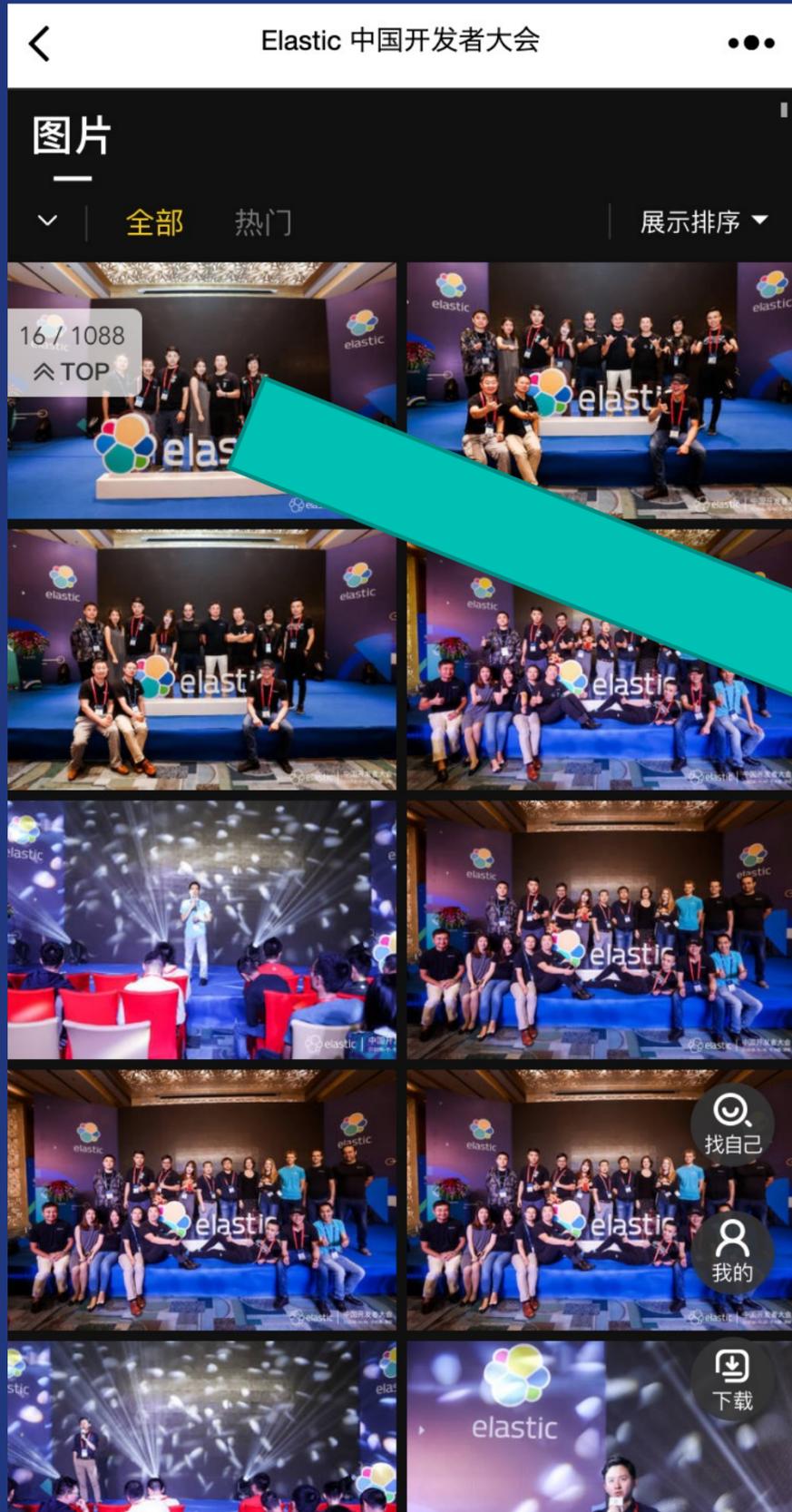
我们的业务

图片直播

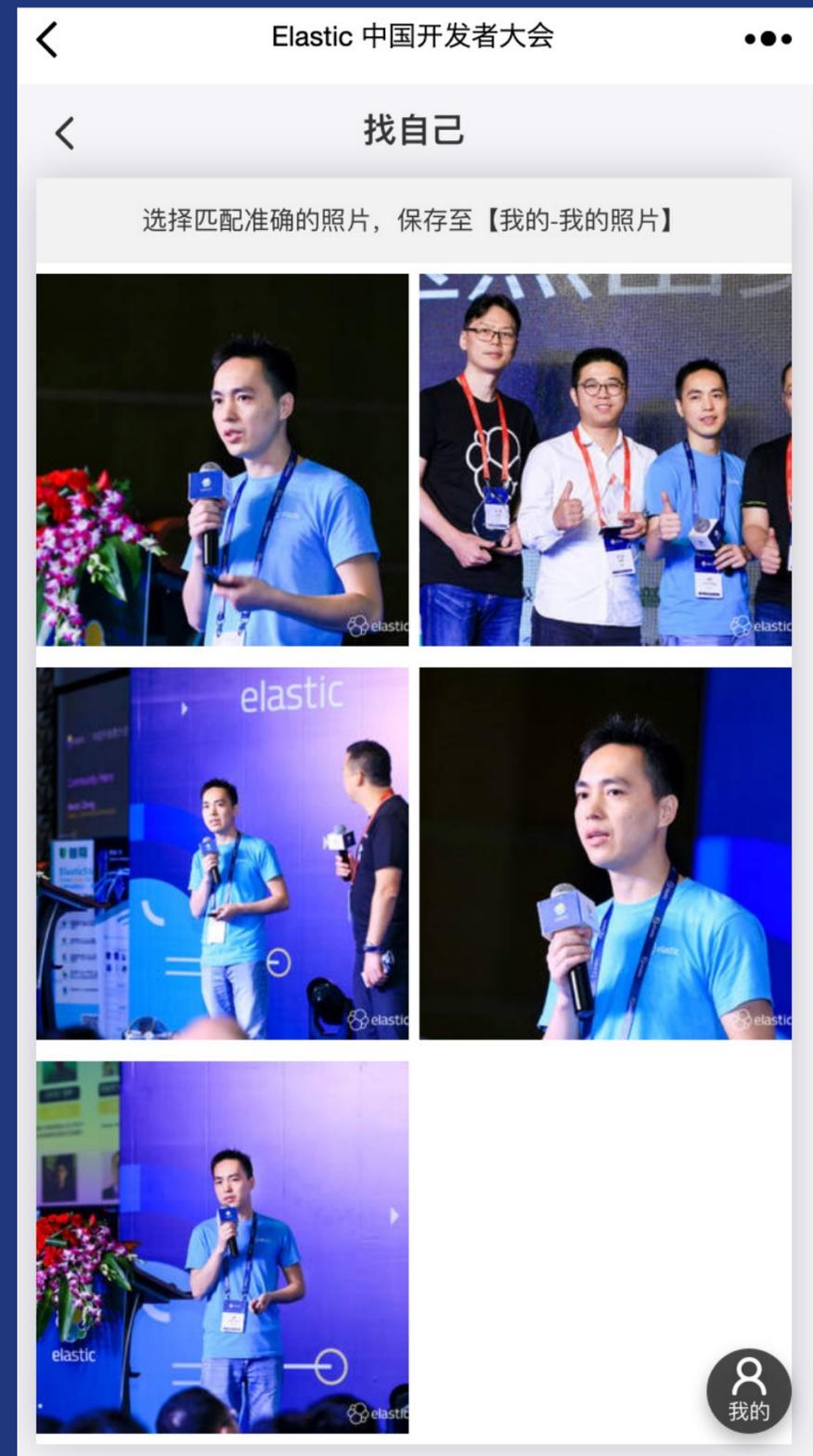
The screenshot shows the Elastic China Developer Conference 2021 photo live stream interface. At the top, there is the Elastic logo and the text "elastic". Below that, the title "中国开发者大会" (China Developer Conference) is displayed, along with the date "2018-11-10" and location "中国-深圳". A search icon and a language selector "En" are also visible. The main content area features the title "Elastic 中国开发者大会" and the number of viewers "36804人已收看". The event details are listed: "时间: 2018.11.10 08:00 - 2018.11.10 19:00" and "地点: 金茂深圳JW万豪酒店". Below the details, there is a "图片" (Photos) section with a dropdown menu showing "全部" (All) and "热门" (Popular), and a "展示排序" (Display Sort) option. The bottom part of the screenshot shows a grid of four photo thumbnails from the event, with interactive icons for "找自己" (Find Me), "我的" (My), and "下载" (Download).



Elastic 中国开发者大会 2021
图片直播二维码



我们的业务





» 我们的问题

依赖第三方人脸搜索

遇到问题



性能限制

2QPS

5个搜索结果



功能限制

多活动聚合

活动人物相册



费用暴增

业务指数增长

不能额外收费

我们的问题



2QPS----->20QPS

依赖第三方人脸搜索

性能限制



5个搜索结果-----> 无解 (没有预期解决时间)

我们的问题

依赖第三方人脸搜索

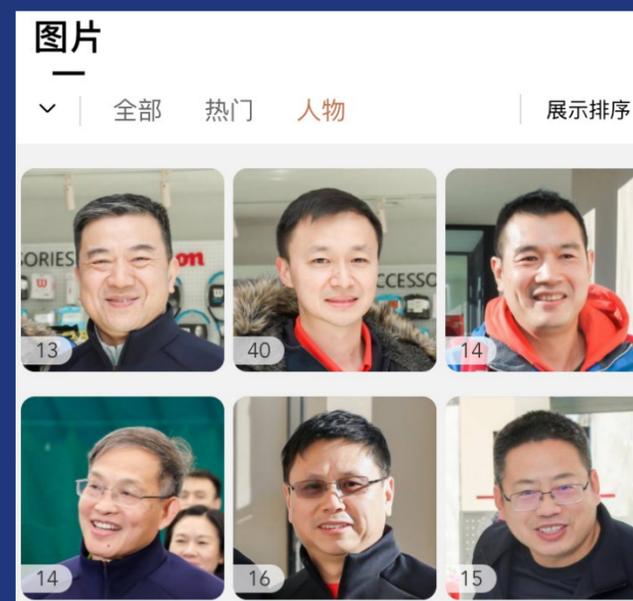
功能限制

多活动聚合

一场马拉松几十万张图

多场马拉松聚合

活动人物相册



» 我们的问题



业务指数增长

依赖第三方人脸搜索

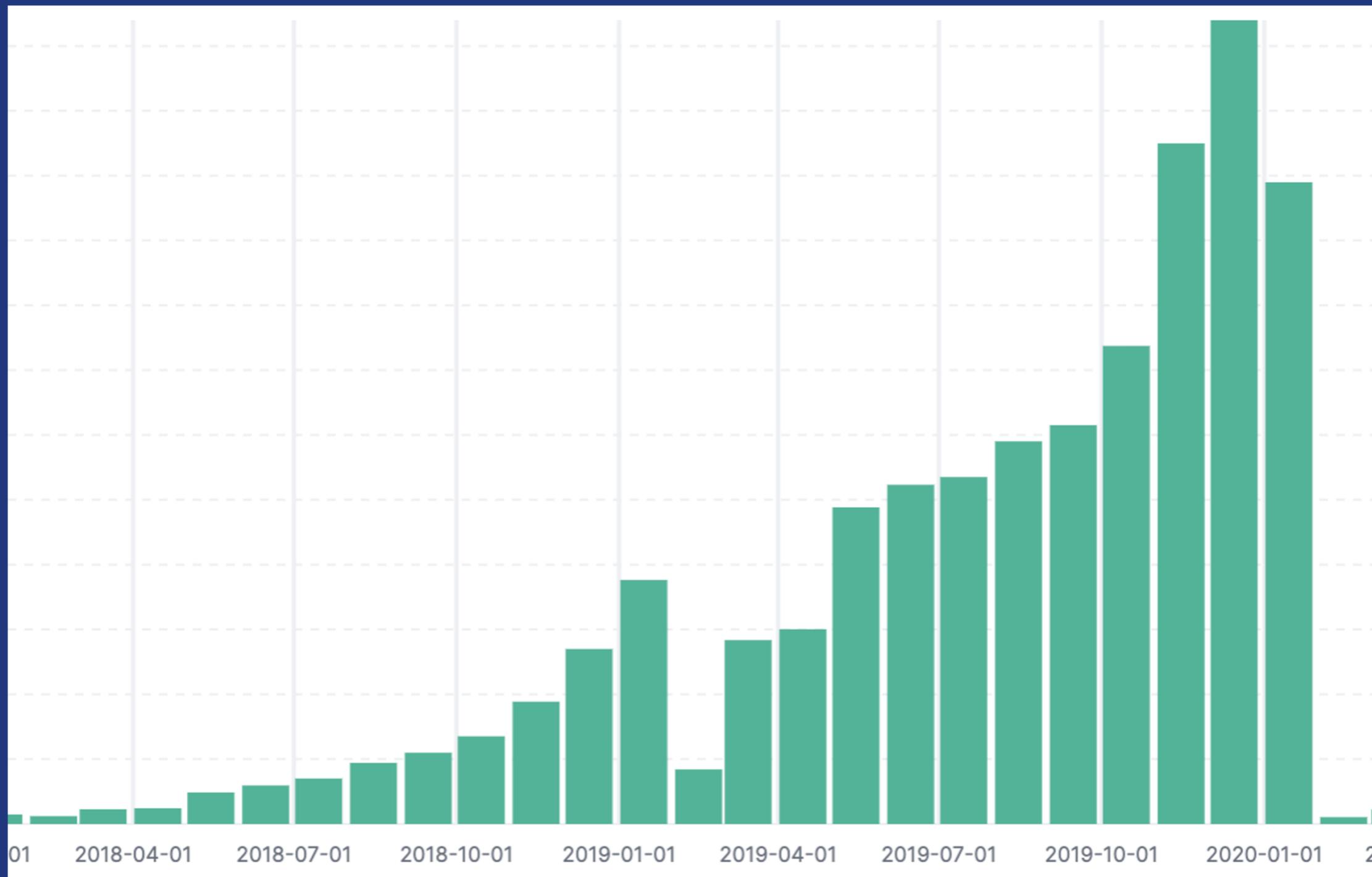
费用暴增



不能额外收费

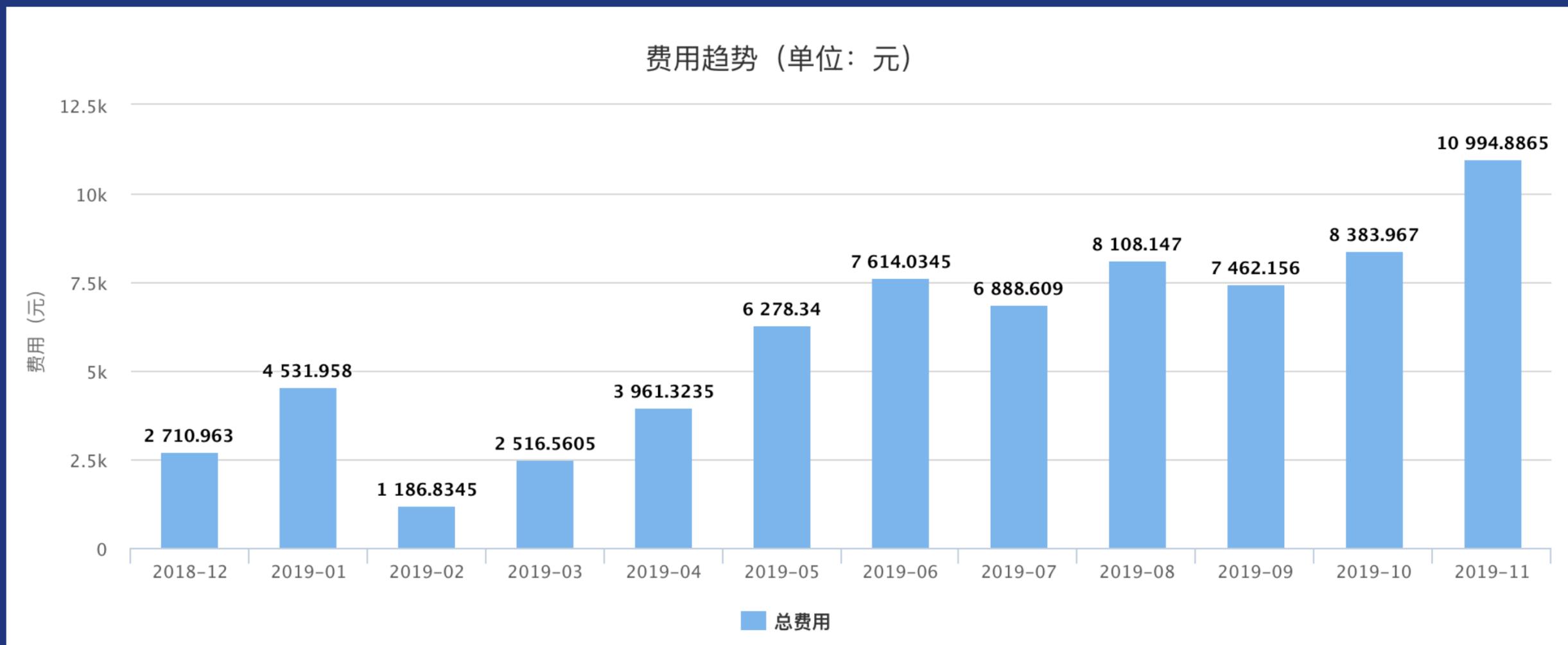
我们的问题

业务指数增长



我们的问题

费用暴增



我们的问题

怎么解决???

小公司



没钱



没人



没时间

“万能的”



elastic

Search. Observe. Protect.

我们的方案

解决问题

第一阶段：

elastic方案探索



理论基础

搜索是计算多维词向量距离

人脸搜索是计算512位向量欧几里得距离



数据类型支持

Elasticsearch 7.0之后新数据类型

Dense_vector

支持最多500维数据



算法支持？

已经支持余弦相似度计算

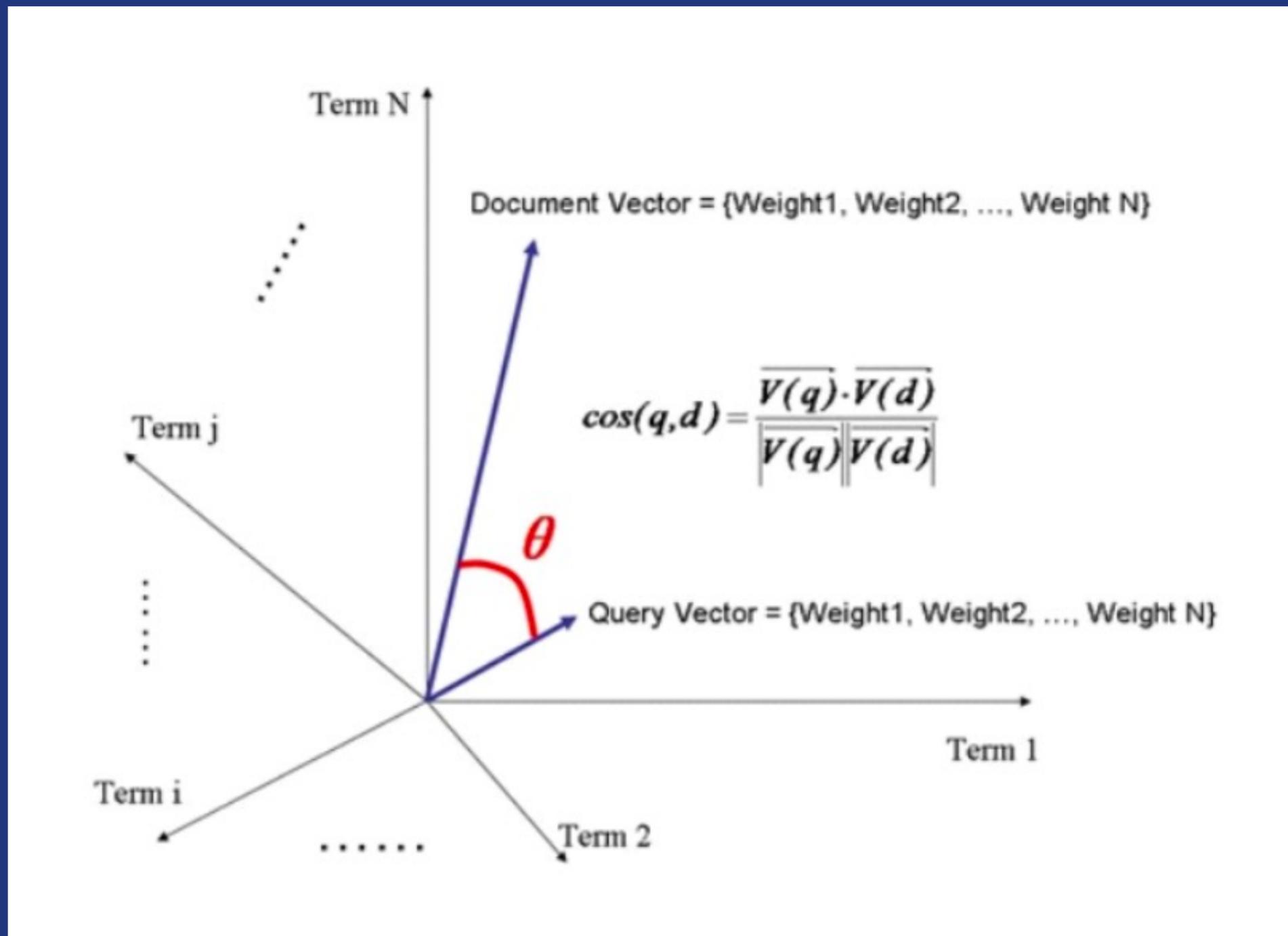
大佬承诺后续版本支持欧几里得距离

我们的方案

理论基础

TF/IDF算法

空间向量算法/
余弦相似度

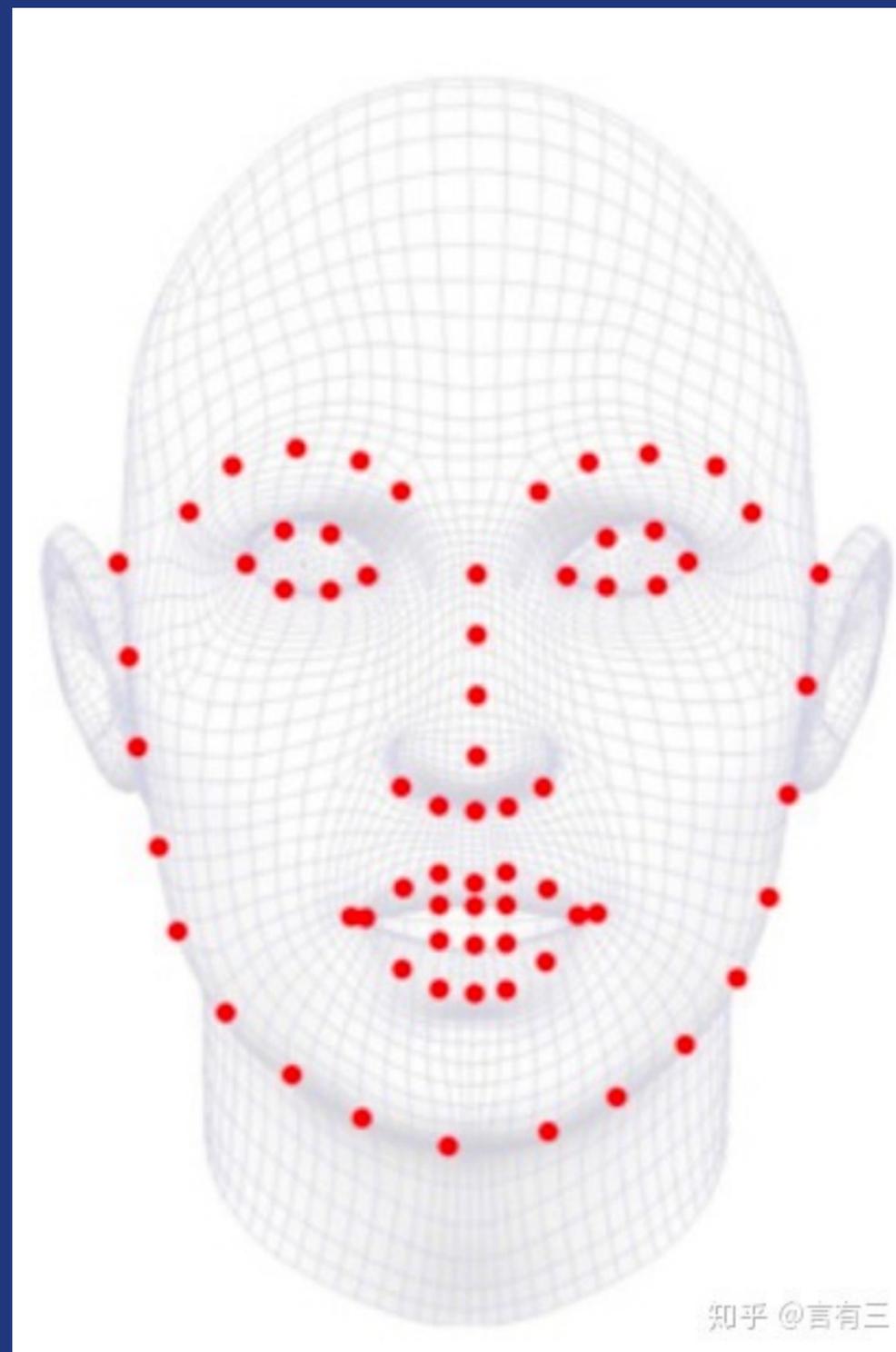


我们的方案

理论基础

欧式空间特征512维

欧几里得距离

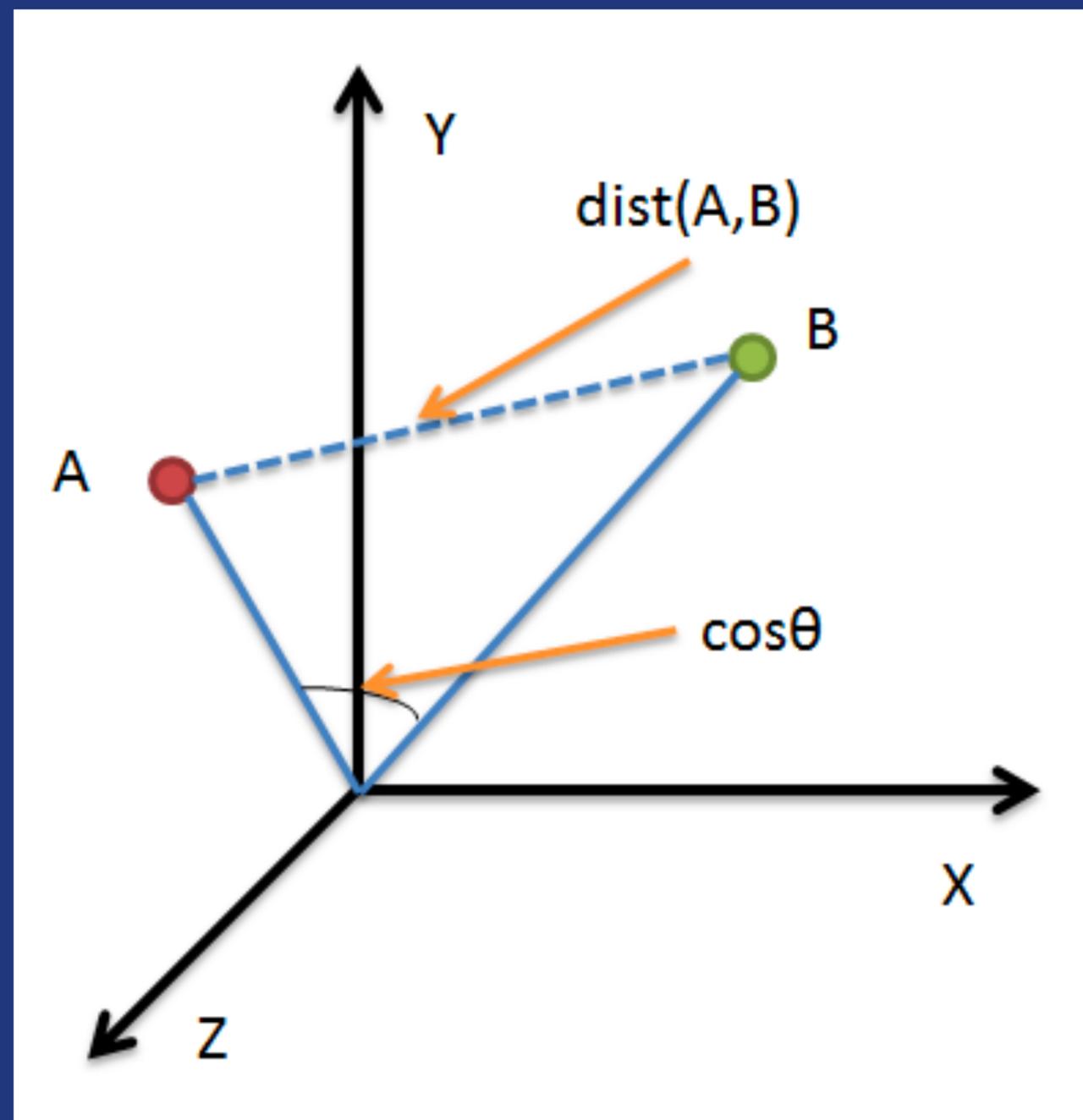


我们的方案

理论基础

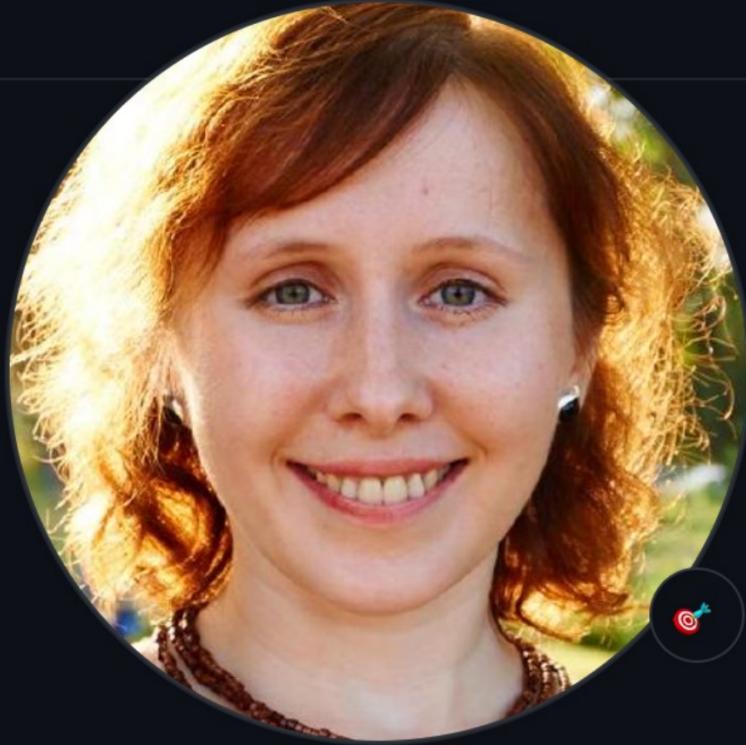
欧式空间特征512维

欧几里得距离



$$d = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2 + \dots + (y_{n_1} - y_{n_2})^2}$$

大佬救命



Mayya Sharipova
mayya-sharipova

Follow

61 followers · 0 following

@elastic

Canada

mayya@apache.org

How to search dense_vector

Elastic Stack Elasticsearch



LiuGangR

3 Jan 2019

Jan 2019

which query is used to search dense_vector. There is no query type to work.

dense_vector is the new feature of es 7.0

<https://www.elastic.co/guide/en/elasticsearch/reference/master/dense-vector.html> 181

1 1

created

Jan 2019

last reply

Feb 2019

2

replies

1.8k

views

1

user

1

like

4

links



LiuGangR

Jan 2019

Feb 2019

there is a way.

github.com/elastic/elasticsearch



Issue: Introduce vector field, vector query and rescoring based on them 71

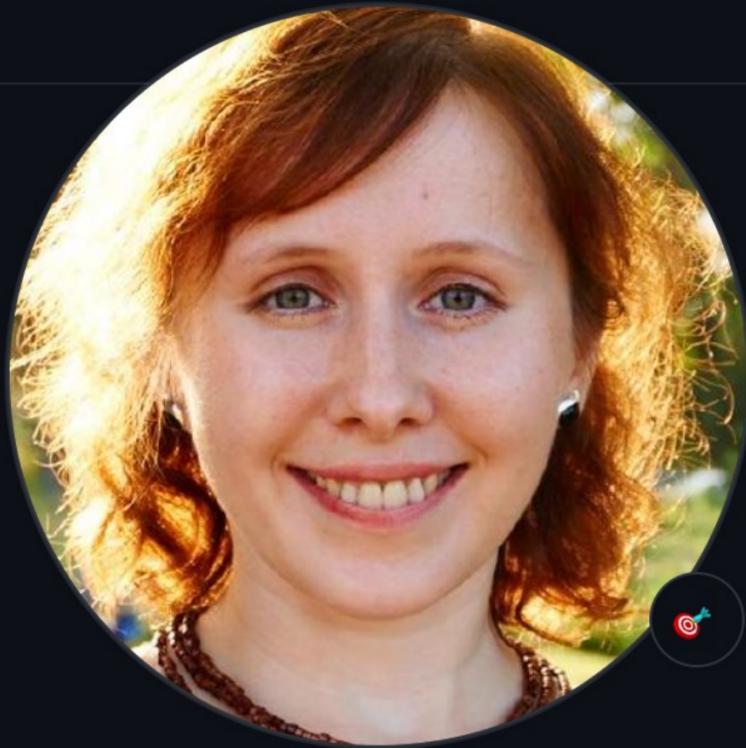
opened by [mayya-sharipova](#) on 2018-06-27

Introduce a new field of type vector on which vector calculations can be done during rescoring phase

```
PUT my_index
{
  "mappings": {
    "_doc": ....
```

[:Search/Ranking](#)

大佬救命



Mayya Sharipova
mayya-sharipova

Follow

61 followers · 0 following

@elastic

Canada

mayya@apache.org

 LiuGangR commented on 1 Feb 2019

@jpountz
That is cool. And you have any plan to support that in which version ?

 jpountz commented on 1 Feb 2019 Contributor

@LiuGangR Hopefully 7.1.

 LiuGangR commented on 2 Feb 2019

@jpountz
another question. If I what to search 'dense_vector' field, the 'cosineSimilarity' is the only why. And is there a default vector query?
Thanks!

 mayya-sharipova commented on 6 Feb 2019 Contributor Author

@LiuGangR yes, the only way to use `dense_vector` or `sparse_vector` in queries is through `cosineSimilarity` and `dotProduct` functions

 1

» 我们的方案

第一阶段结论：

es作为人脸搜索引擎

第一阶段验证：

es7.x (未发布) + script

```
"script": {
  "lang": "painless",
  "source": ""
  double total = 0;
  for (int i = 0; i < doc['face_detail'].length; ++i) {
    total +=(params['_source']['face_detail'][i] - params.queryVector[l])
            *(params['_source']['face_detail'][i] - params.queryVector[i]);
    if (total>params.confidence){
      return total;
    }
  }
  return total;
},
"params": {
  "confidence":0.16,
  "queryVector": [
    0.07027599215507507,
    0.09286174178123474,
    0.03654364496469498,
```

我们的方案

第二阶段：

elastic方案上线



7.2版本之前

没有正式版，是等还是用开发板？

人脸向量512位，dense_vector最多支持500位

不支持欧几里得距离计算



7.2版本

人脸向量512位，dense_vector支持1024位

不支持欧几里得距离计算

集群管理引入x-pack



7.4版本之后

支持欧几里得距离计算

搜索优化

异步搜索

我们的方案

自研人脸向量512维

es7.2版本之前

20kb/人脸

dense_vector只支持500位



非es方案一

数据存数据库
计算用java



非es方案二

数据存redis
计算用python

我们的方案

es7.2版本



7.2版本dense_vector支持1024位

功能测试

性能测试

建表

测试+集群整体升级



7.1版本x-pack免费

更安全的集群管理

身份管理

我们的方案

es7.2版本

测试



功能测试

- 1.数据类型，欧几里得距离计算脚本
- 2.几十万样本计算



性能测试

- 1.预估正式数据量和压力
- 2.不同硬件，索引结构
- 3.确认集群硬件和索引结构



信心

- 1.能不能替代行业独角兽大厂的服务？
- 2.能不能一直稳定？
- 3.能不能省钱？

我们的方案

场景

1场活动

500-50万张图

500-100万张人脸

硬件

cpu 内存 硬盘

索引结构和分布

几台物理机？几个索引？几个分片？几个备份？

» 我们的方案

场景

1场活动

500-50万张图

500-100万张人脸

3*8c32G3T高性能云盘

性价比

10个index(face_01)

分表--扩展性

7个shards

3,5,7,11,13测试后确认

2个replicas

写多读少

» 我们的方案

es7.2版本

信心哪里来？

是骡子是马，遛出来



生产数据在双方同时生成



同时搜索互为补充

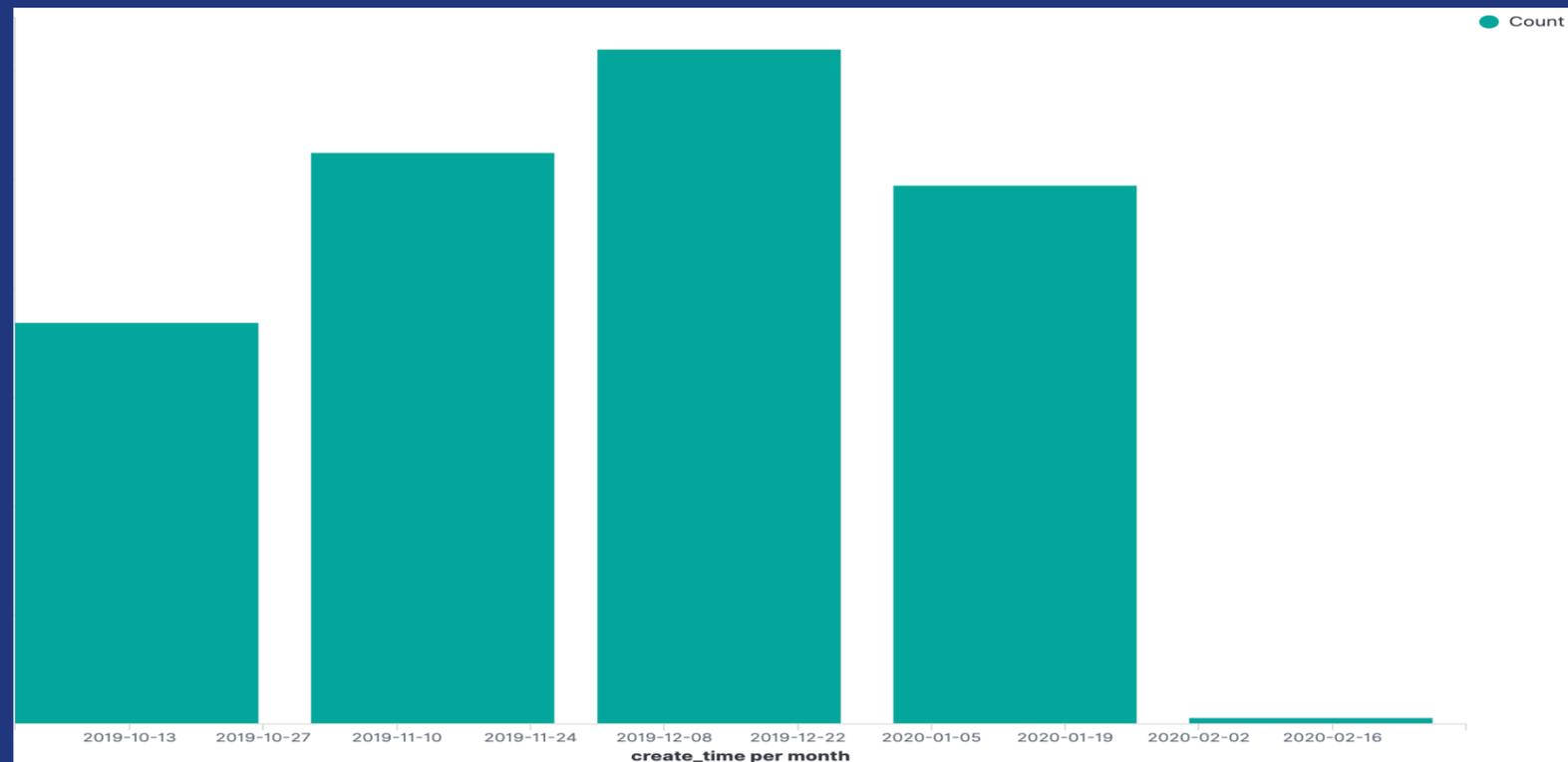


对比结果的数量，准确率，失败率，召回率

我们的方案

es7.2版本

正式上线



我们的方案

es7.2版本



7.2版本dense_vector支持1024位

功能测试

性能测试

建表

测试+集群整体升级



7.1版本x-pack免费

更安全的集群管理

身份管理

» 我们的方案

es7.4版本之后

优化



支持欧几里得距离计算



搜索优化



异步搜索

我们的方案

es7.4版本之后

优化

```
"sort": [  
  {  
    "_score": {  
      "order": "desc"  
    }  
  }  
],  
"size": 100,  
"from": 0,  
"query": {  
  "script_score": {  
    "min_score": 10 - confidence,  
    "script": {  
      "source": "10-l2norm(params.face, 'face_detail')",  
      "params": {  
        "face": [-0.0016885467, -0.08441351, -0.020592572,
```

我们的方案

第二阶段

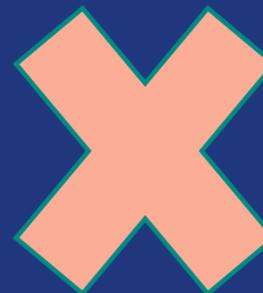
回顾



性能限制

2QPS

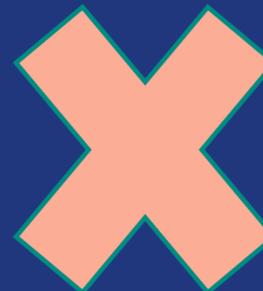
5个搜索结果



功能限制

多活动聚合

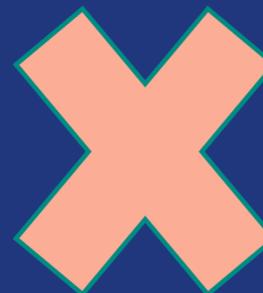
活动人物相册



费用暴增

业务指数增长

不能额外收费



我们的方案

第三阶段

集群优化+省钱



疫情影响

行业是重灾区

降本增效



业务特点

相册有时效性，可以分层

统计，监控，业务都在一个集群



硬件降价

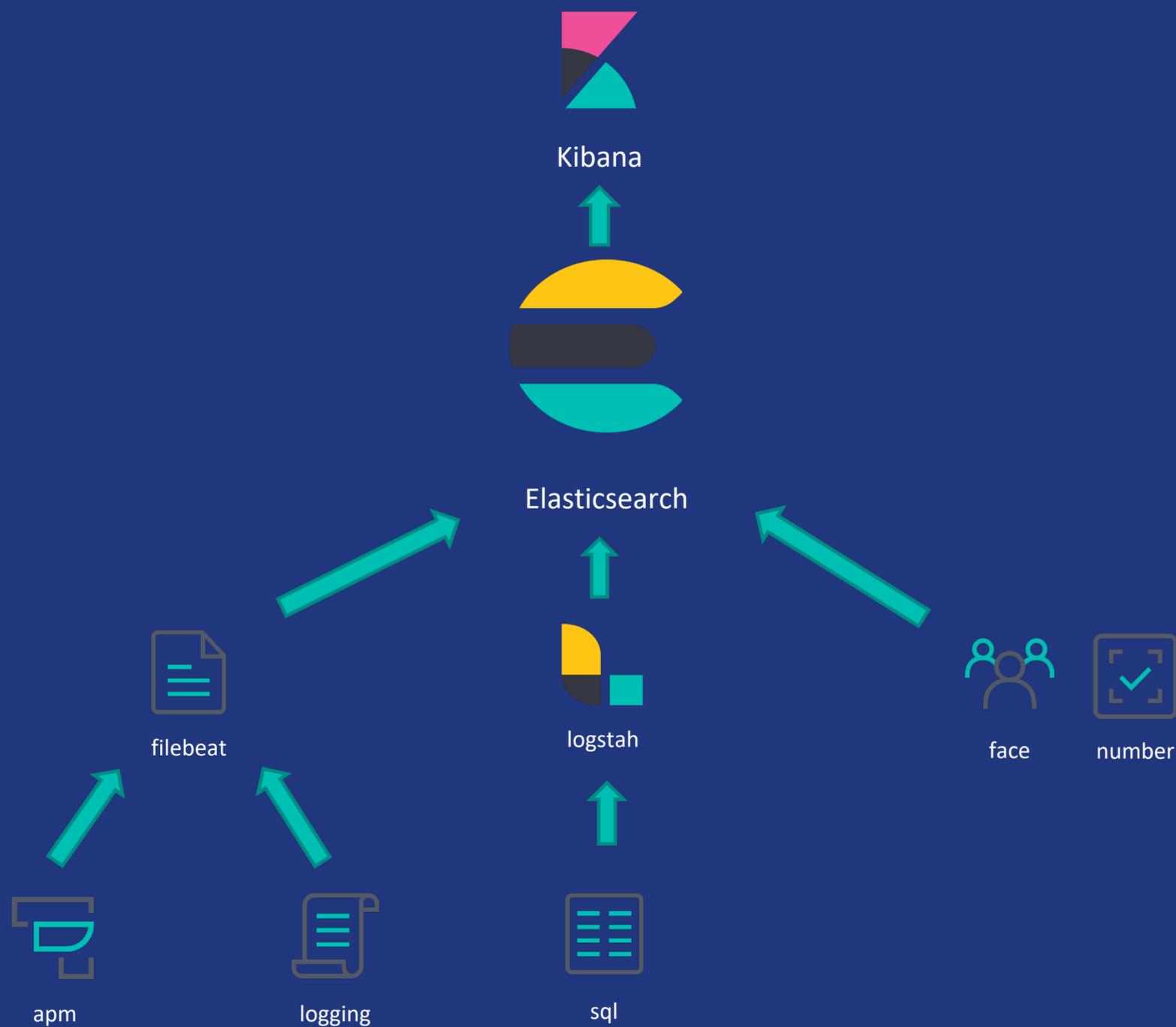
大数据类型服务器出现

HHD和SSD都可选

我们的方案

第三阶段

业务特点



我们的方案

第三阶段

方案



生命周期方案



硬件分层方案



索引汇总作业 (rollup jobs)

我们的方案

第三阶段

方案

硬件分层+索引汇总



硬件分层

服务器分为两部分-----HHD和SSD

热数据放SSD，冷数据放HHD

定期手动配置



索引汇总 (rollup jobs)

桶进数据按照业务需求定期汇总到新索引

调整业务模式

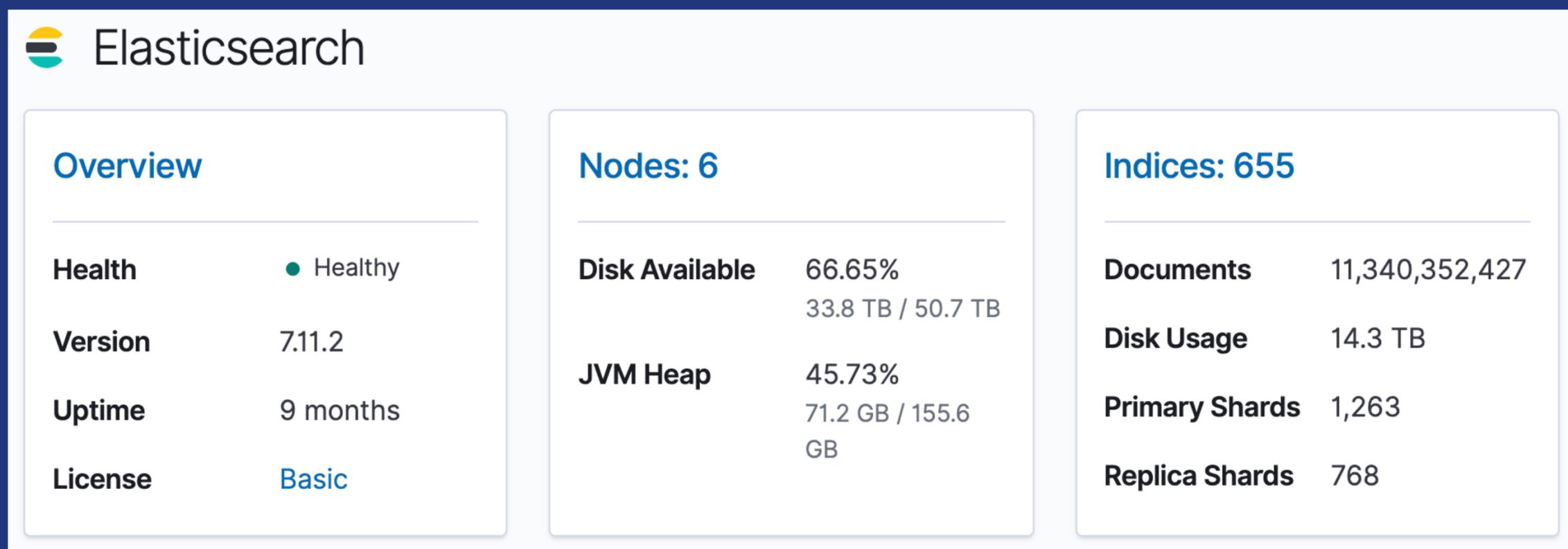
我们的方案

第三阶段

2*8c32G22T(HHD)

4*8c64G2T(SSD)

回顾



Elasticsearch

Overview

Health	● Healthy
Version	7.11.2
Uptime	9 months
License	Basic

Nodes: 6

Disk Available	66.65%
	33.8 TB / 50.7 TB
JVM Heap	45.73%
	71.2 GB / 155.6 GB

Indices: 655

Documents	11,340,352,427
Disk Usage	14.3 TB
Primary Shards	1,263
Replica Shards	768

后续计划

我们的规划



百亿级人脸搜索规划，引入es8.x特性

后续计划



数据迁移到其他云服务商---省钱



视频上人脸搜搜的方案探索

特别感谢



Medcl 黎霖 林瑞凯 徐奥平 王金雷

十亿级人脸搜索的实践和优化

刘刚 CTO @谱时智能云

2023年4月8日



感谢观看



专业、垂直、纯粹的 Elastic 开源技术交流社区

<https://elasticsearch.cn/>