



基于ES的企业搜索中台

付国庆

2019-12-7, 总经理, 北京顶尖时代科技发展有限公司

目录

1 关于 顶尖时代 TOPTIME

2 关于中台

3 搜索中台，了解一下

4 基于ES企业搜索中台的解决方案

5 案例分享

背景

- 2005年成立
- 技术团队来自百度
- 十几年专注于搜索技术领域

产品服务

- 产品：企业搜索中台
- 产品：企业情报系统
- 服务：互联网云采集
- 服务：互联网云搜索

客户 500+

- 企业
- 政府
- 金融机构
- 研究机构

目录

1 关于 顶尖时代 TOPTIME

2 关于中台

3 企业搜索中台

4 基于ES企业搜索中台的解决方案

5 案例分享

中台概念

战斗单元

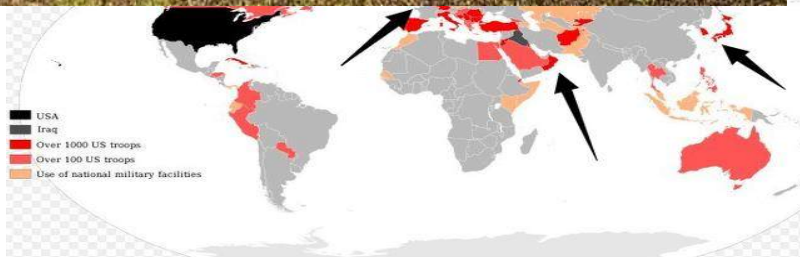


前台



中台

军事基地



后台



互联网公司的中台战略



提出中台战略：**大中台** 小前台；

中台运营数据能力、产品技术能力，帮助前台适应市场、业务支撑、业务创新。



腾讯副总裁：腾讯将进一步开放数据中台和技术中台



京东集团黎科峰：作为一家技术企业 京东将打造产业互联网中台



企业中台

应用形式多
需求变化快
灵活、机动

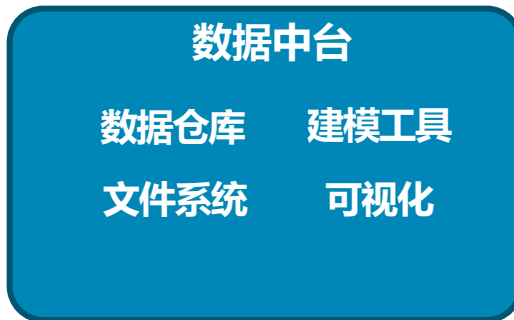
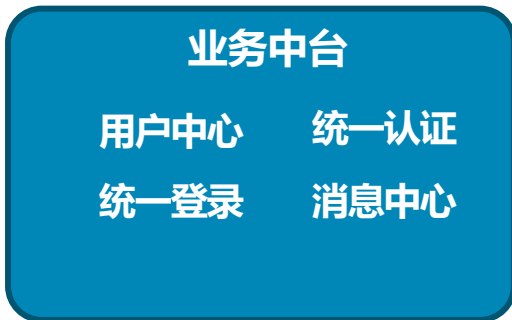
前台
(业务)



营销部门
服务部门
客户
合作伙伴
供应商

资源、能力、数据
复用、赋能平台

中台
(能力)



中台建设团队


稳定、安全、
审计、合规

后台
(管理)



集团内部
运营部门

IAAS 企业云



目录

1 关于 顶尖时代 TOPTIME

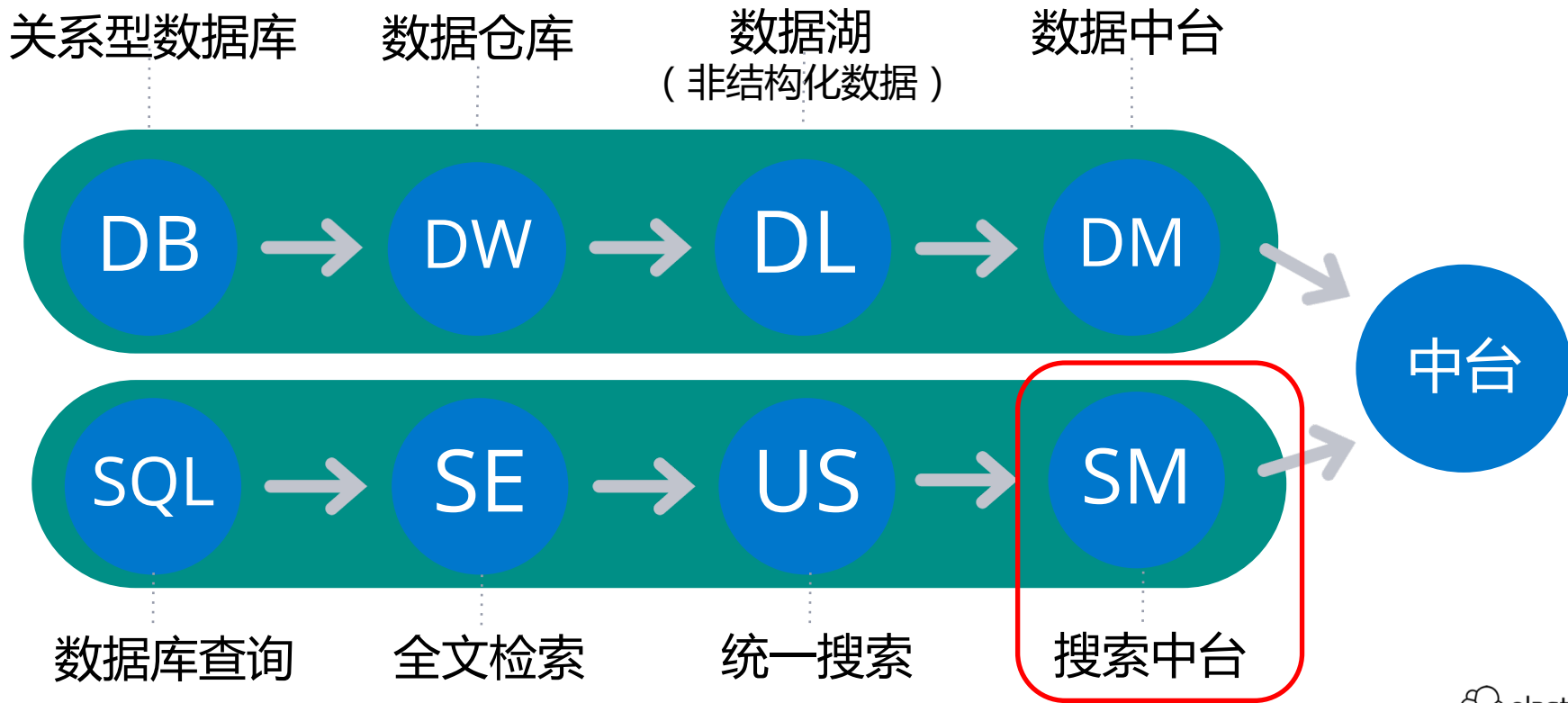
2 关于中台

3 搜索中台，了解一下

4 基于ES企业搜索中台的解决方案

5 案例分享

数据与搜索的交汇



企业搜索中台

应用形式多
需求变化快
灵活、机动

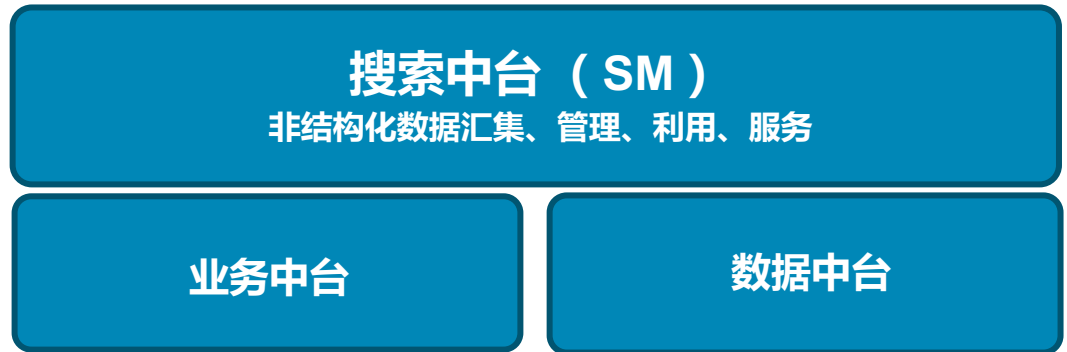
前台
(业务)



营销部门
服务部门
客户
合作伙伴
供应商

资源、能力、数据
复用、赋能平台

中台
(能力)



中台研发团队

稳定、安全、
审计、合规

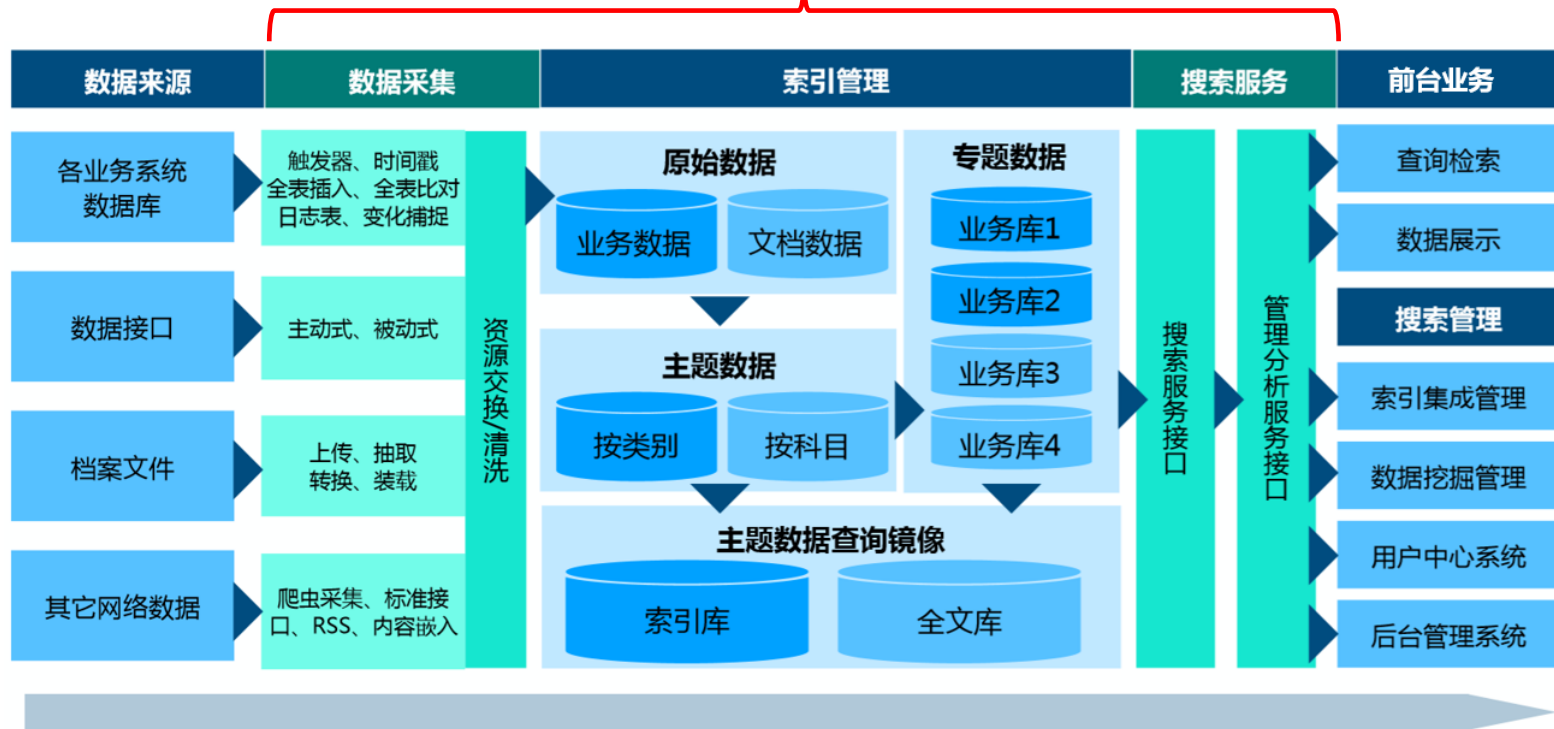
后台
(管理)



集团内部
运营部门



搜索中台的数据与服务流程



目录

1 关于 顶尖时代 TOPTIME

2 关于中台

3 企业搜索中台

4 基于ES企业搜索中台的解决方案

5 案例分享

统一搜索

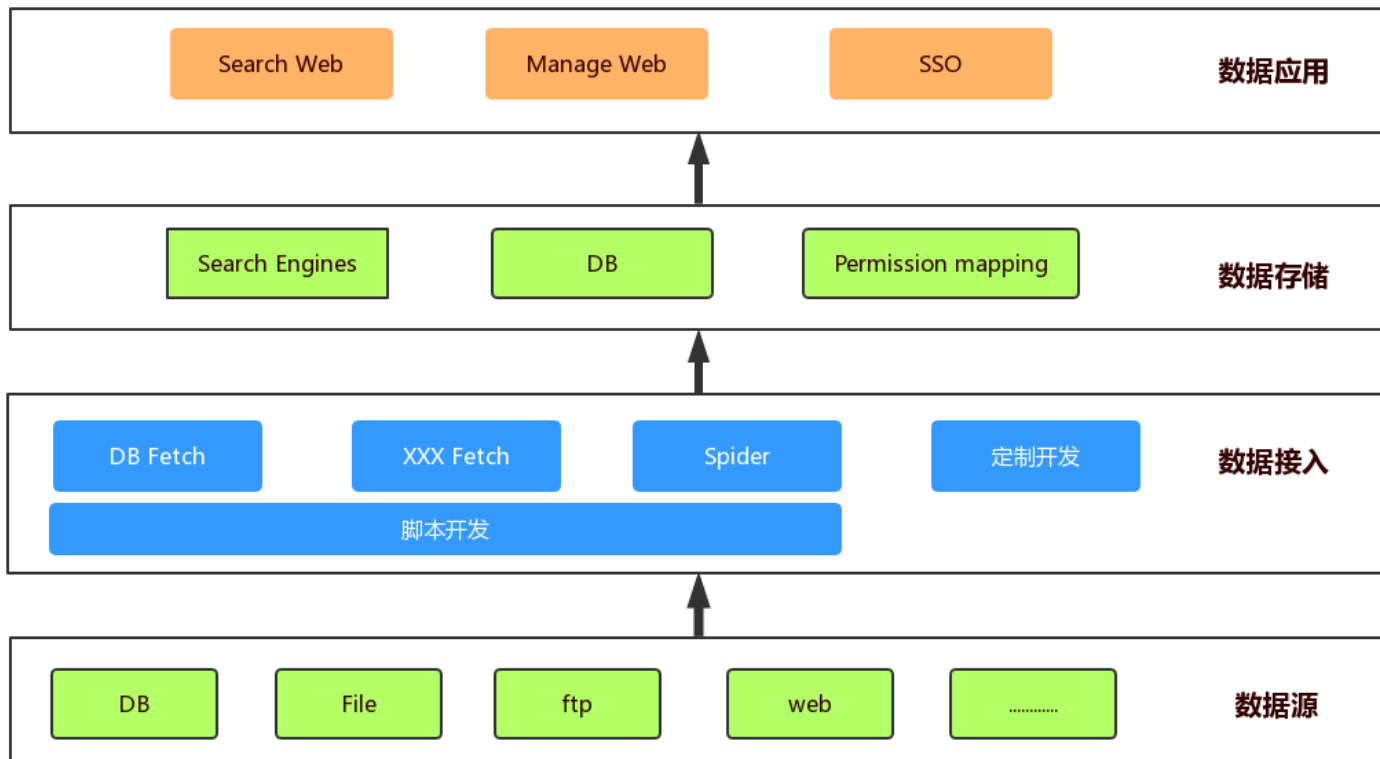


搜索中台

发展演变

统一搜索流程


统一搜索团队



统一搜索方案

优点

实现了企业内部资源的一键搜索，解决信息孤岛
统一登录，继承业务系统权限进行搜索

缺点

开发、实施周期长，工作量大，
搜索技术团队，需适应各业务系统进行开发
缺少复用性
后期各系统升级、调整后，搜索功能需要调整



企业搜索中台优势 专业的事情交给专业团队完成

快速开发

接入、搜索实现了可视化配置

统一解决方案

不直接面向应用,提供统一解决方案,支持二次开发,各系统变化时工作量大幅减少

灵活复用

功能、数据多次利用,解决方案不局限在某一领域



搜索中台

数据汇交

平台数据同步接口(云端数据落地)

内部数据汇交

REST

FILE

MQ

DB VIEW

Spider (本地网页采集)

数据治理

数据加工

数据标签

数据分类

关联分析

用户画像

实体提取

索引管理

数据利用

专题分析

人物分析

主题分析

信息汇聚

集约化搜索

标准搜索

场景化

搜索统计

搜索服务接口

数据可视化

智能问答

数据共享

资源同步接口(对外提供数据)

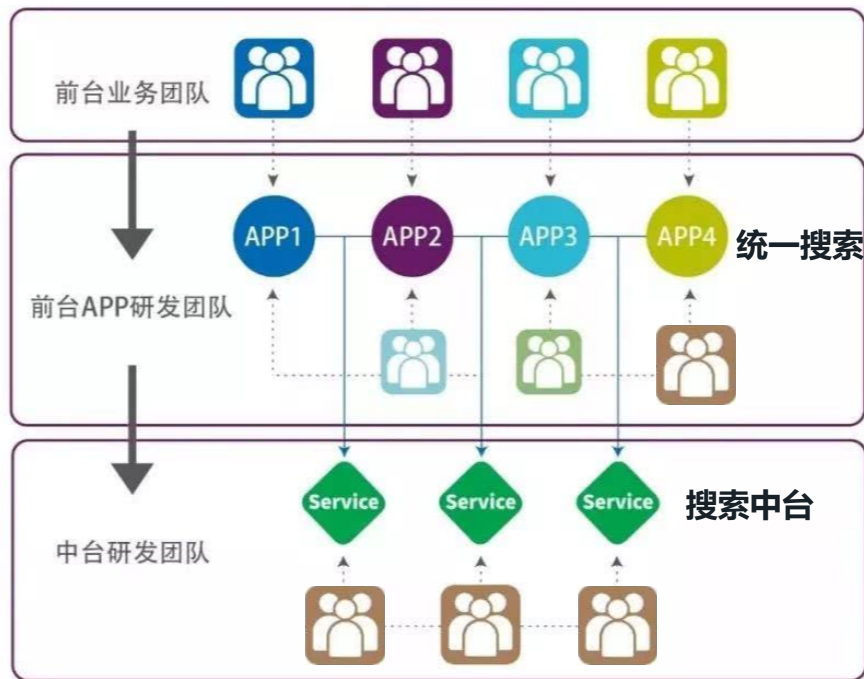
应用同步接口(系统应用外围调用)

配置搜索接口(数据汇交搜索服务)

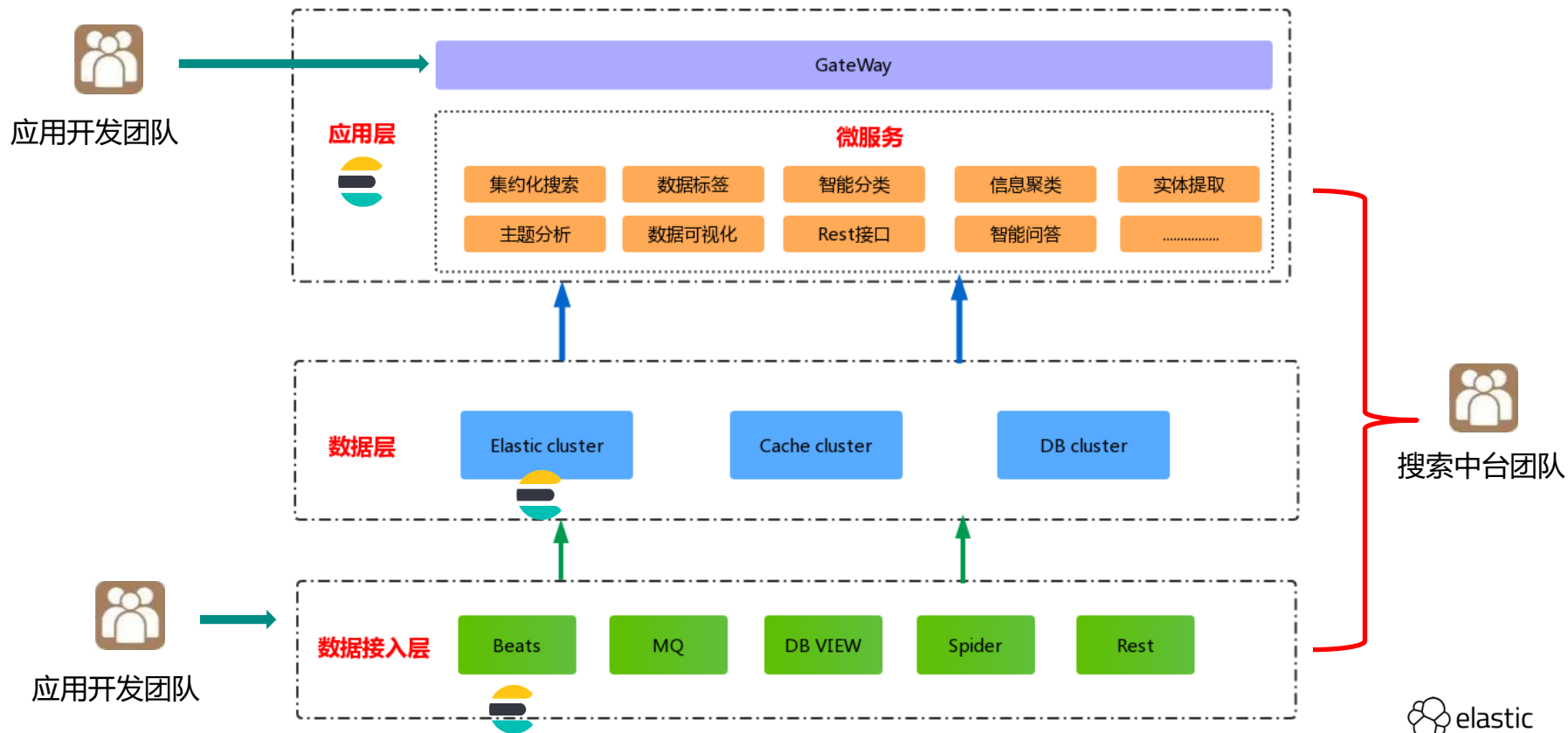
原生搜索接口

问答服务接口

采用搜索中台架构后，研发团队定位变化



基于ES搜索中台架构



数据汇交

REST接口

支持ElasticSearch原始脚本进行数据提交

MQ方式

支持Kafka、RabbitMQ、RocketMQ方式进行增量数据提交

DB视图方式

支持ORACLE、MYSQL、SQLSERVER等数据库接口视图进行可视化配置

文件方式

支持csv文件方式进行初始化、增量数据提交

编辑导入任务

任务名称 *	任务类别 *	任务状态 *
<input type="text" value="同步客户房间数据"/>	<input type="text" value="db view"/>	<input type="text"/>
索引库 *	任务关联映射 *	
<input type="text" value="crmss"/>	<input type="text" value="CRM_AccountRoom X"/>	
JDBC连接字符串 *	用户名 *	密码 *
<input type="text" value="jdbc:sqlserver://...;databaseName=E."/>	<input type="text"/>	<input type="text"/>
数据库类型 *	视图名称 *	
<input type="text" value="SQL SERVER"/>	<input type="text" value="AccountRoomView"/>	
时间戳字段 *	状态字段 *	
<input type="text" value="ModifiedOn"/>	<input type="text" value="StateCode"/>	
是否包含附件 *		
<input type="text" value="不包含"/>		



数据汇交

REST接入样例

接口参数说明

调用地址: `http://[redacted]rest/bulk`

请求方式: POST

返回类型: JSON

请求头:

名称	类型	是否必须	描述
X-Gaia-API-Key	STRING	必填	网关校验信息, 请在网关管理处申请

请求参数:

名称	类型	是否必须	描述
indexDb	STRING	必填	索引库别名
requestBody	STRING	必填	批量数据

测试样例:

```
indexDb=resttestdb0
requestBody=
{ "index" : { "_id" : "1" } }
{ "field1" : "value1" }
```

接口请求测试

indexDb *:

requestBody
*:



发送请求

测试返回结果

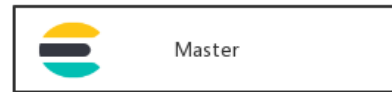
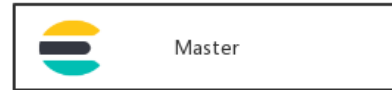
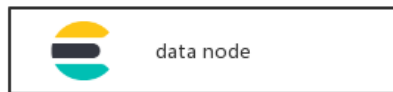
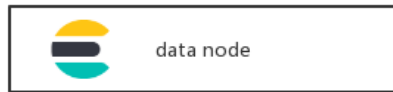
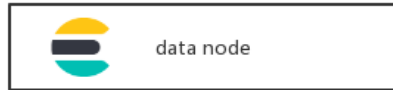
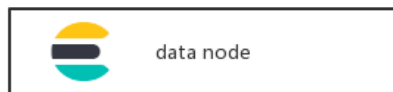
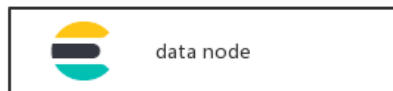
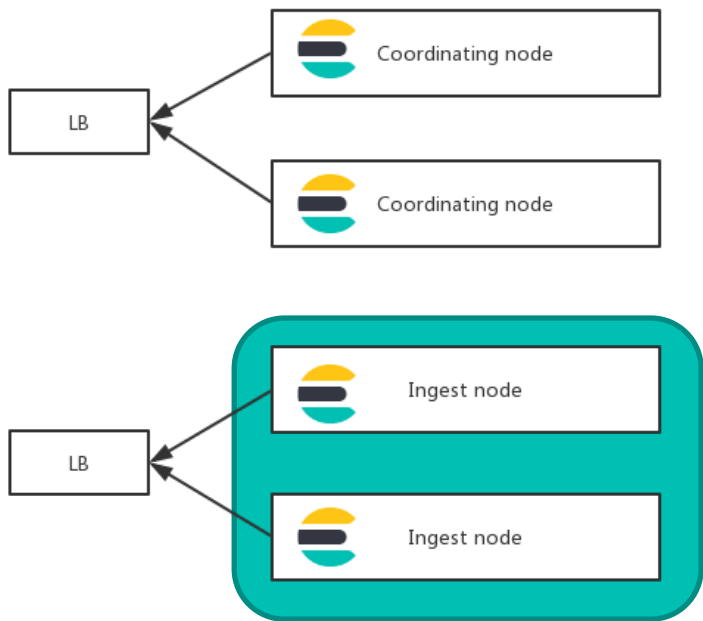
数据汇交

MQ接入样例

```
{
  "customerId": "56431", //租户ID
  "indexName": "oadb", //索引
  "mappingName": "cps", //索引结构
  "type": 0, //操作类别 (增删改)
  "dataList": [ //数据详情
    {
      "DRREFERENCE": "",
      "DRETITLE": "",
      "URL": ""
    }
  ]
}
```



数据治理



数据治理

Pipeline

```
curl -XPUT "http://192.168.0.7:9200/_ingest/pipeline/toptime-default-pipeline" -u
'elastic:toptime' -H 'Content-Type: application/json' -d'
{
  "description": "toptime default pipeline",
  "processors": [
    {
      "set": {
        "field": "DREDBNAME",
        "value": "{{_index}}"
      }
    },
    {
      "script": {
        "source": "if(ctx._index==~/(.*)[\\d-]*?$/){ctx.DREDBNAME=/(.*)[\\d-
_]*/$.matcher(ctx._index).replaceAll(\"$1\")}"
      }
    }
  ]
}'
```

数据利用

基于ES搜索中台微服务

搜索功能	智能化服务	场景化服务	索引管理 搜索统计服务	接口服务
关键词搜索	下拉提示	主题场景	搜索次数统计	互联网云采集接口 Cloud Spider
分类导航搜索	智能纠错	应用场景	搜索热词统计	
搜索结果分类统计	拼音搜索	机构名片	搜索点击统计	数据库同步接口 DB View
复合条件搜索	关键词推荐	人物名片	搜索转化率统计	
标签过滤搜索	通俗语言	知识图谱	搜索改进率统计	文档提交与解析接口 FILE
地域过滤搜索	敏感词屏蔽	图形化	索引库管理	
图片搜索	框计算	智能问答	索引数据管理	API提交接口 (REST)
附件文档搜索	智能排序	地图搜索	词汇管理	API搜索接口 (REST)
微信、微博搜索	搜索结果置顶	用户画像		API删除接口 (REST)
视频搜索	智能推荐			MQ
	智能语义识别			

数据利用

默认搜索条件

默认排序方式：

相关度

默认搜索位置：

标题+全文

去重字段：

关闭

规则配置

搜索规则：

智能匹配

智能匹配

精准匹配

完全匹配优先

普通搜索

自定义规则

信息

修改

删除

加权语法

时间加权：

关闭

范围天数：

90

加数值：

5

自定义加权：

关闭

加权方式：

boost

value :

加数值：

5

添加

- Query DSL

Query and filter context

+ Compound queries

+ Full text queries

+ Geo queries

+ Joining queries

Match all

+ Span queries

+ Specialized queries

+ Term-level queries

`minimum_should_match` parameter

`rewrite` parameter

Regular expression syntax

数据利用

搜索代码解耦-Search Template



```
POST _scripts/template1
{
  "script": {
    "lang": "mustache",
    "source": "{{from}}: {{{^from}}}0{{/from}}, "size":
    {{{size}}}{^size}}20{{/size}}, "query": {"bool": {"should":
    [{"bool": {"must": [{"bool": {"should": [{"query_string":
    "query":
    "{{queryText}}", "fields": [{"{{queryField1}}"{{#queryField2}}", "{{
    queryField2}}"{{/queryField2}}], "default_operator":
    "and", "auto_generate_phrase_queries": true, "boost":
    {{{completeBoost}}}{^completeBoost}}2{{/completeBoost}}}}]"}}
  }
}
```

数据共享

搜索接口配置

+ 新建 编辑 删除 请输入搜索词

接口ID	接口名称	状态
B82DAA5C395E47CFF2BAF624225EBF8C	门牌号搜索	正常

显示第 1 至 1 项结果, 共 1 项

上页 1 下页

接口参数说明

调用地址: <https://.../custom>

请求方式: POST

返回类型: JSON

请求头:

名称	类型	是否必	描述
----	----	-----	----

接口请求测试

dbName:

pageSize:

page:

搜索参数:

发送请求

基本配置

自定义搜索参数

参数名称:

索引库字段 (多个字段请用英文半角逗号,分割):

高亮字段:

请选择

默认值:

数据类别:

STRING

高级

添加

参数名称	索引库字段	高亮字段	默认值	数据类别	查询类别	分词方式	命中最小分词比例	操作
a	DRETITLE			STRING	QUERY	细分	1	<input type="text"/>

系统安全

ES访问安全

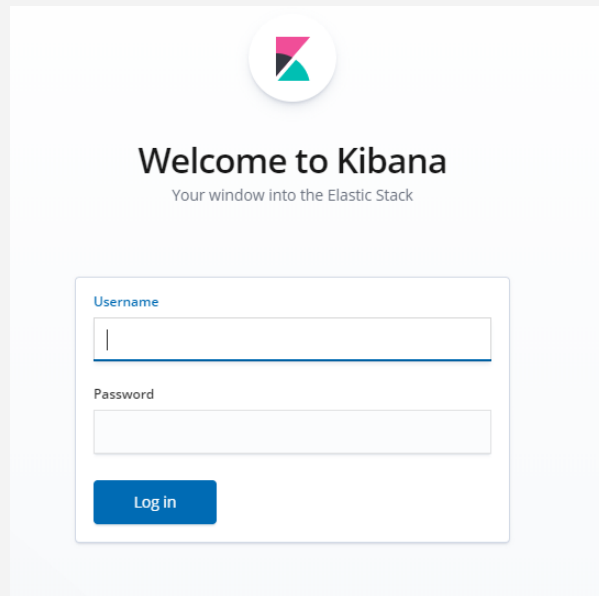
文件和原生 Realm，可用于创建和管理用户TLS功能，可对通信进行加密
每个Elasticsearch节点生成私钥和X.509证书

租户权限分离

租户单独进行授权，禁止跨库操作

公私钥分发

租户单独分发密钥，通过RSA方式进行接口通信



Elasticsearch 安全功能现免费提供
(从 6.8.0 和 7.1.0 版本开始)

系统运维&统计

DevOps & Count

日志分析

- 中间件日志
- 应用日志

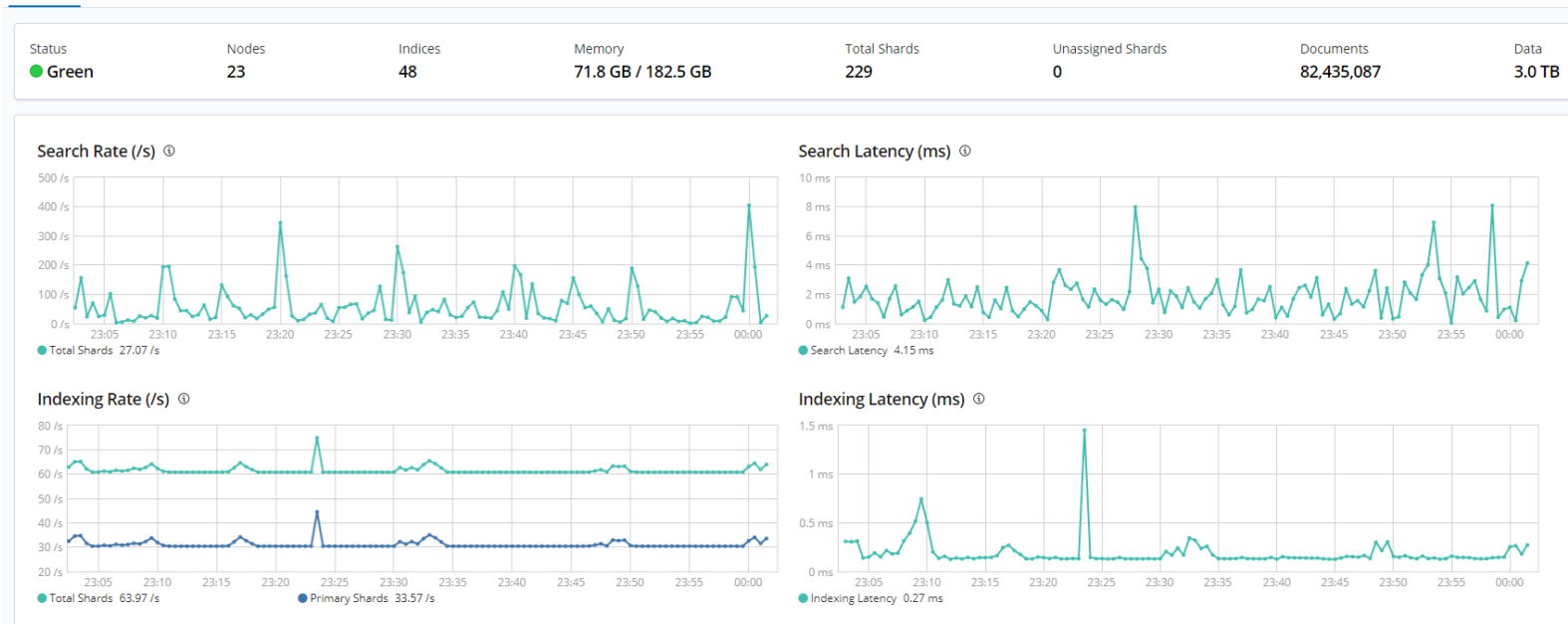
系统指标

- CPU
- RAM
- HDD

安全报警

- 机器人攻击
- 恶意刷榜

系统监控



系统监控



Name	Status	CPU Usage	Load Average	JVM Memory	Disk Free Space ↑	Shards
★ master-01 192.168.0.7:9300	● Online	0% ↑ 1% max 0% min	0.16 ↓ 0.23 max 0 min	44% ↑ 45% max 40% min	52.0 GB ↓ 52.0 GB max 52.0 GB min	0
☰ master-03 192.168.1.99:9300	● Online	0% ↑ 1% max 0% min	0 ↓ 0.27 max 0 min	22% ↑ 26% max 20% min	52.3 GB ↓ 52.3 GB max 52.3 GB min	0
☰ master-02 192.168.1.98:9300	● Online	0% ↑ 1% max 0% min	0.01 ↑ 0.17 max 0 min	21% ↑ 26% max 20% min	52.3 GB ↓ 52.3 GB max 52.3 GB min	0
☰ data-04 192.168.0.12:9300	● Online	0% ↓ 5% max 0% min	0.03 ↓ 0.14 max 0 min	61% ↑ 63% max 58% min	101.5 GB ↑ 101.5 GB max 101.3 GB min	12
☰ data-05 192.168.0.18:9300	● Online	0% ↑ 1% max 0% min	0.08 ↑ 0.11 max 0 min	44% ↑ 46% max 41% min	102.8 GB ↓ 102.8 GB max 102.7 GB min	11
☰ data-11 192.168.1.106:9300	● Online	0% ↑ 1% max 0% min	0.03 ↓ 0.15 max 0 min	43% ↑ 47% max 41% min	105.3 GB ↓ 105.3 GB max 105.3 GB min	12
☰ data-01 192.168.0.8:9300	● Online	0% ↑ 2% max 0% min	0.15 ↑ 0.34 max 0 min	24% ↑ 24% max 19% min	108.2 GB ↓ 108.2 GB max 108.1 GB min	11
☰ data-08 192.168.0.15:9300	● Online	0% ↑ 1% max 0% min	0.05 ↑ 0.43 max 0 min	38% ↑ 39% max 34% min	111.3 GB ↓ 111.3 GB max 111.2 GB min	11
☰ data-06 192.168.0.14:9300	● Online	0% ↑ 1% max 0% min	0.01 ↑ 0.21 max 0 min	59% ↑ 60% max 55% min	111.3 GB ↓ 111.3 GB max 111.3 GB min	11
☰ data-19 192.168.1.102:9300	● Online	0% ↑ 3% max 0% min	0.08 ↑ 0.72 max 0 min	57% ↑ 60% max 54% min	113.7 GB ↓ 113.7 GB max 113.7 GB min	12

日志分析预警



```
- siteCode: 6101000031, tab: all, qt:2017 统计公报, 搜索总耗时: 48 ms
- siteCode: 1101110035, tab: xxgk, qt:,,用途 辅助功能耗时: 9 ms
- siteCode: 1101110035, tab: xxgk, qt:,,用途 es搜索耗时: 106 ms
- siteCode: 1101110035, tab: xxgk, qt:,,用途, scene搜索耗时: 1 ms
- siteCode: 1101110035, tab: xxgk, qt:,,用途, 搜索总耗时: 116 ms
- siteCode: 1101110035, tab: xxgk, qt:用途分类 辅助功能耗时: 9 ms
- siteCode: 1101110035, tab: xxgk, qt:用途分类 辅助功能耗时: 9 ms
- siteCode: 1101110035, tab: xxgk, qt:用途分类 es搜索耗时: 35 ms
```

IP访问信息

IP	Num
121.9.216.126	474
218.19.215.148	195
222.95.171.222	180
103.85.144.198	151

有效控制企业构建的成本

ES服务器硬件要求

(PC server | 虚拟化主机)

- CPU 8核
- RAM 16G/32G
- HDD SSD
- OS Centos 7.X
- 参数配置

Ulimit

Swappiness

JVM等



未来扩展 2020

未来扩展

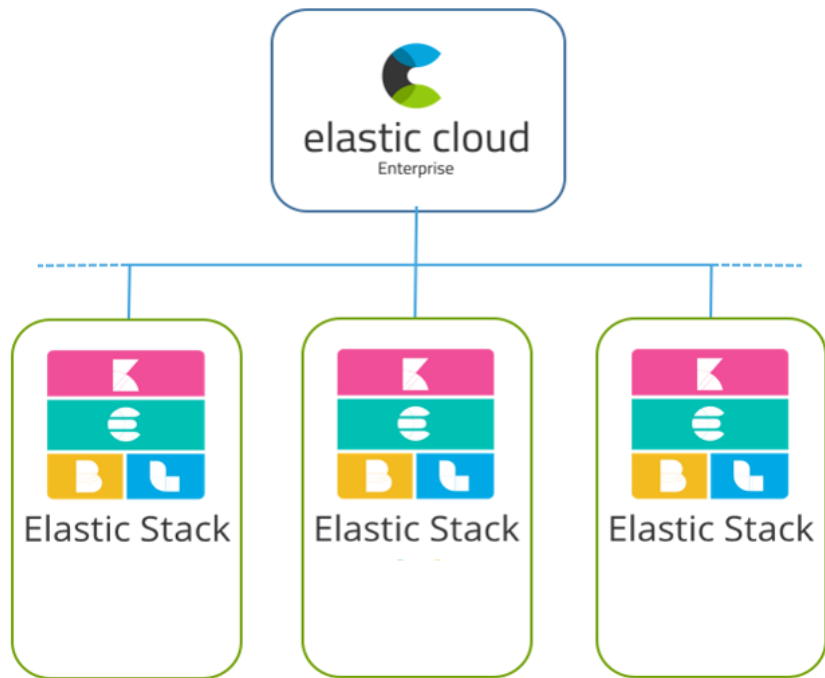
集成ECE

扩展多集群应用场景，简化Elastic Stack集群维护成本

租户安全分离

租户在应用端、后台管理端均进行隔离，有效保障数据安全性

Elastic Cloud Enterprise (ECE)



目录

1 关于 顶尖时代 TOPTIME

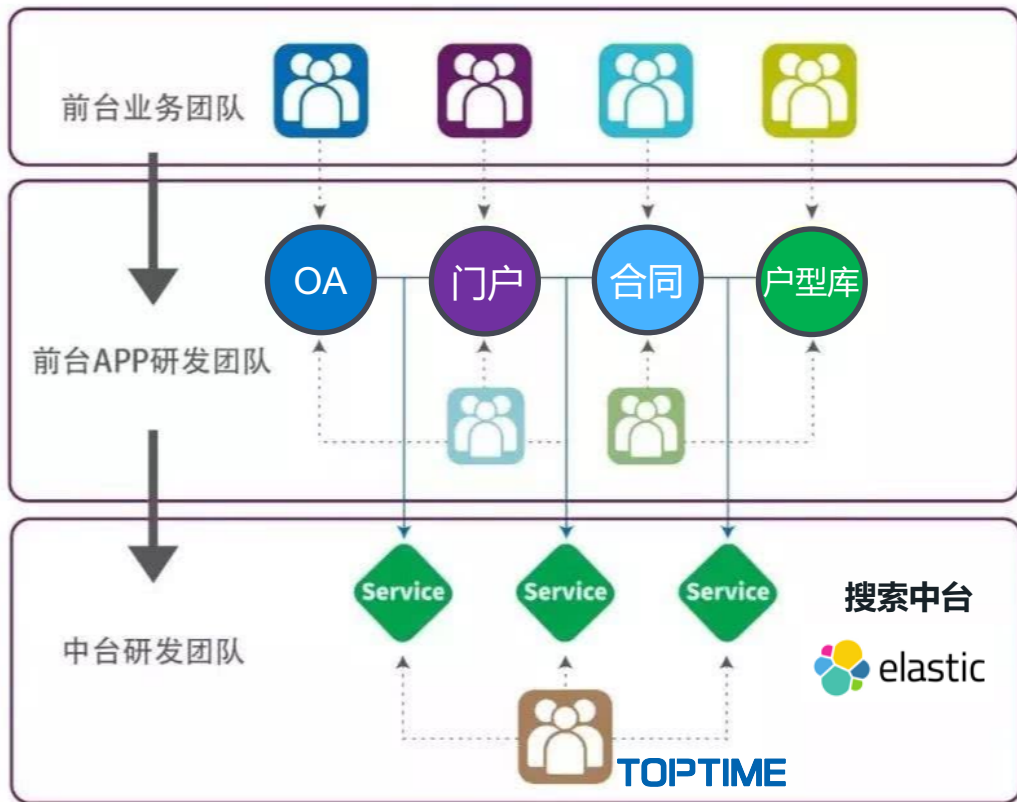
2 关于中台

3 企业搜索中台

4 基于ES企业搜索中台的解决方案

5 案例分享

搜索中台支撑了企业若干系统的搜索功能建设



搜索中台-灵活构建复合条件搜索应用

户型库搜索

输入案名进行搜索，支持多个案名，用空格隔开



封装类型

方案

全套

城市 (公司)

北京

城市 (地区)

北京

业态

请选择

户型计容面积(m²)

0-80

80-85

85-90

90-95

95-100

100-105

105-115

135-145

145-160

160-175

175-190

190-195

195以上

m²

m²

梯户比

T2

T3

T4

T5

T6

房型(交付) - 室

请选择

房型(交付) - 卫

请选择

房型(交付) - 卫

请选择

是否精装

精装

毛坯

配置标准

P1

P2

P3

P4

住宅-小高层

住宅-高层

住宅-超高层

住宅-洋房

住宅-独栋别墅

住宅-联排别墅

住宅-叠拼别墅

展开更多筛选

城市 (公司)

北京

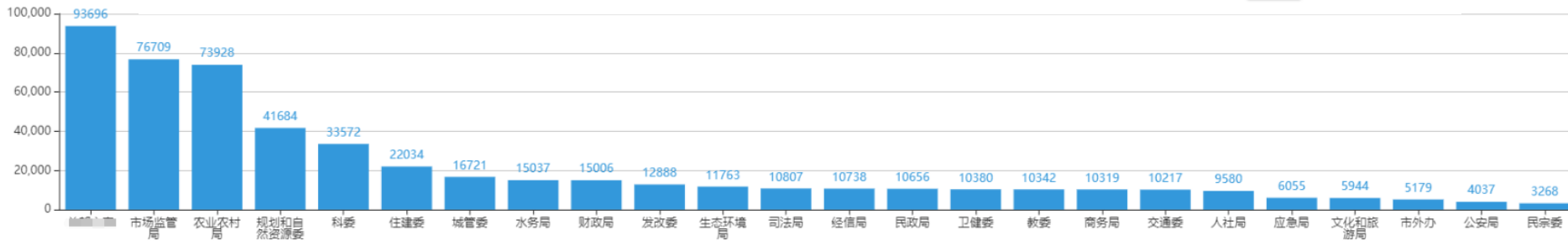
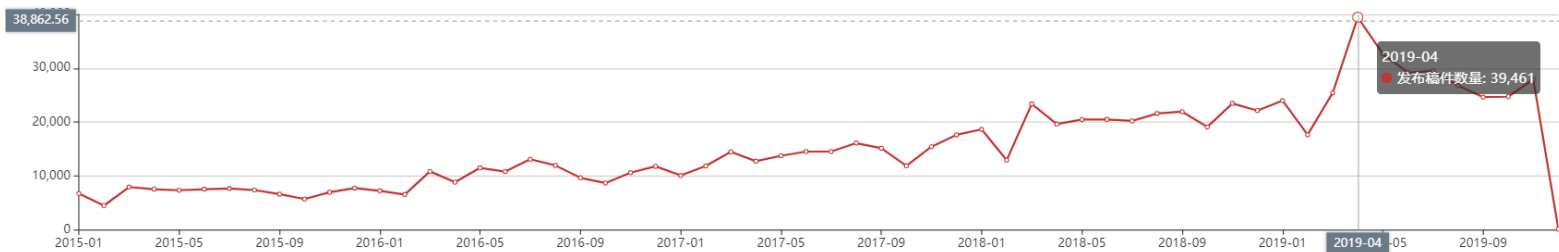
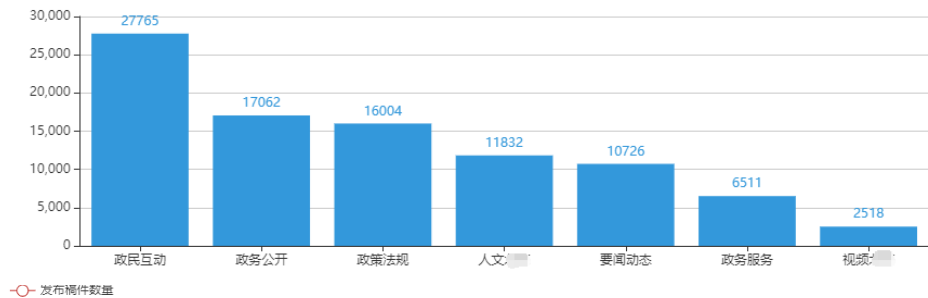
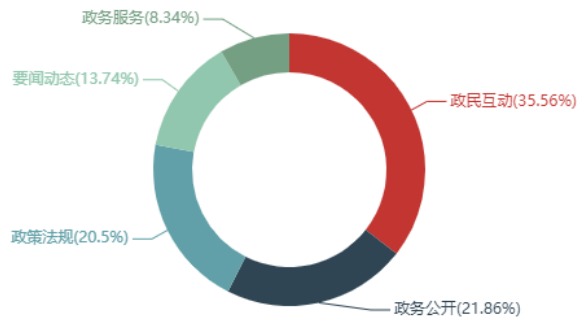


城市 (地区)

北京



搜索中台-快速构建数据可视化分析专题：大数据统计



总结

支撑1000+ 网站搜索服务



谢谢