



日志平台Elasticsearch升级之路

华泰证券-信息技术部 葛宝磊

目录 CONTENTS



项目背景介绍



关键系统设计



Elasticsearch跨版本升级



坑和运维



项目背景介绍

01

项目背景介绍



2016

日志平台项目在2016下半年上线第一版，第一个接入的系统是OA办公，日志量每天不到1G，Elasticsearch选用的版本是2.4.0

2016-2018

随着接入用户的增多，数据量和用户的需求也成指数级的增加，目前已接入的应用系统有近100个，每天的日志量1T+，高峰的Elasticsearch入库速率20W+

2018-2019

ES性能瓶颈明显，通过横向扩展无法解决，需要升级Elasticsearch版本

项目背景介绍

Indexing Rate

Indexing 速率明显跟不上采集端的采集速率，特别是在9:30和下午1:00，非常明显
查询日志无法实时的展现。

CPU告警

业务高峰期个别节点CPU使用率超过90%，影响整个集群的性能

磁盘告警

机器中部分服务器采用的裸盘的方式挂载，经常会出现部分磁盘占用率超过70%情况，而其他盘利用率低于50%

GC问题

分片数达到一定数量后，频繁的full GC，导致集群不响应
索引创建失败导致堆内存无法释放，需要定时重启

数据恢复问题

当某个节点或者集群需要重启时，耗时时间长
cluster.routing 不支持节点级别的配置

Query内存占用

ES使用方查询海量数据时，容易导致OOM

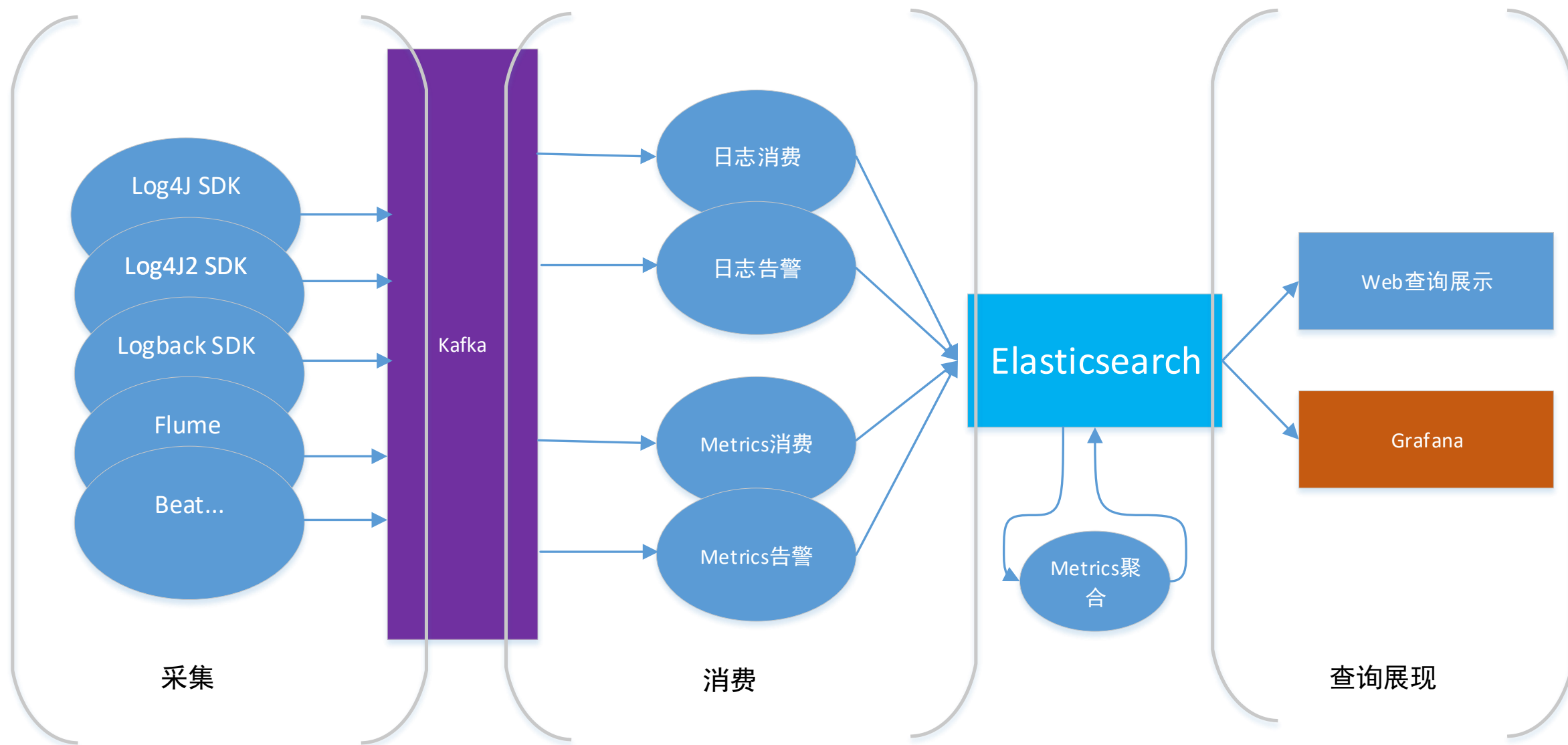




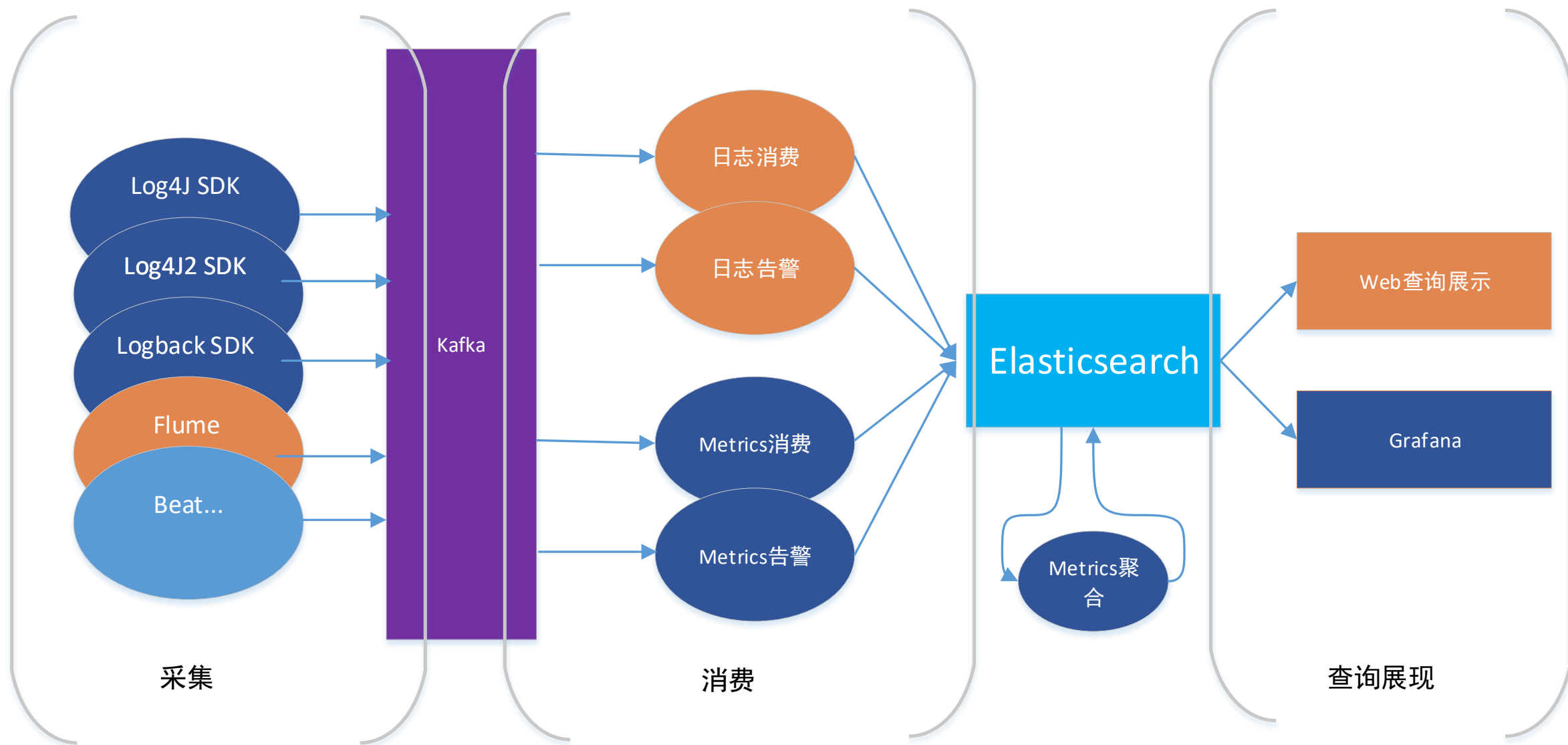
关键系统设计

02

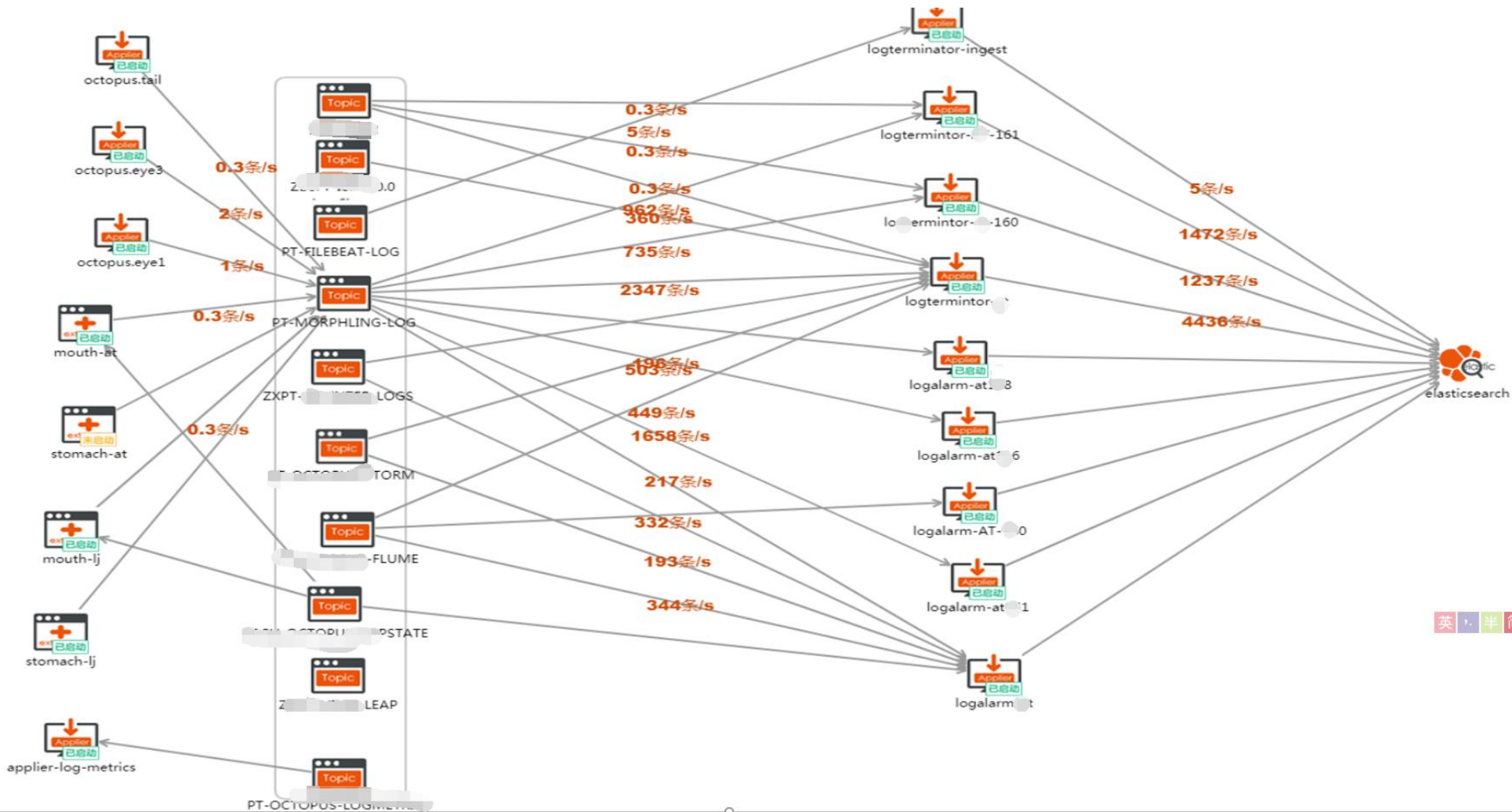
关键系统设计 – 系统架构



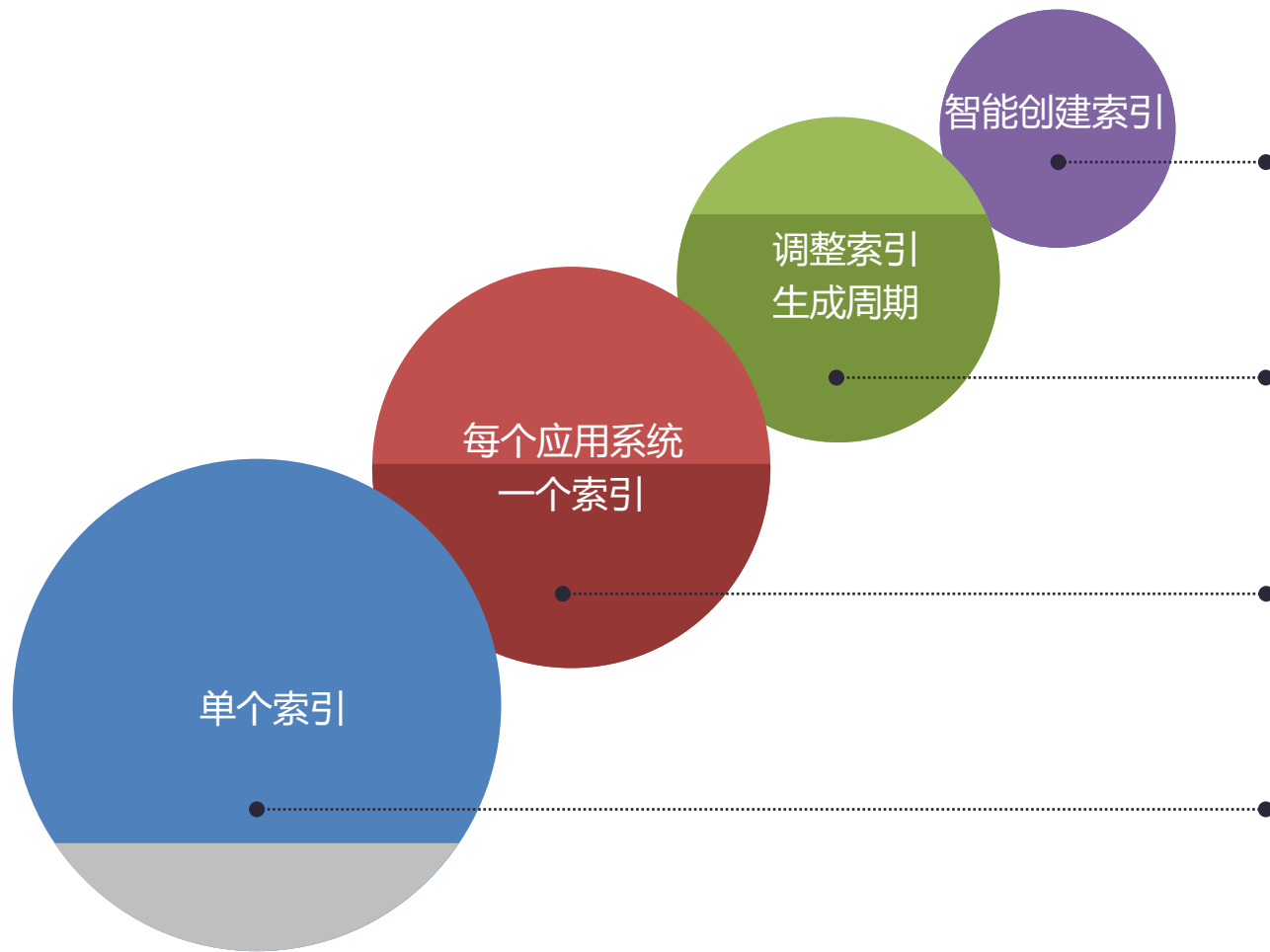
关键系统设计 – 系统架构



关键系统设计 – 系统架构



关键系统设计 – 索引设计



智能创建索引

每周统计应用系统的日志量，根据上周的日志量大小决定下周写入索引生成周期。

调整索引生成周期

根据应用的预估日志量按不同的时间间隔创建索引，支持按天，按周，按月，按年创建索引

优点：索引量适中，分片数适中，查询效率高

缺点：对于突发情况或者应用日志量发生较大变更无法自适应。

每个应用系统每天创建一个索引

每个应用系统每天创建一个索引

优点：索引管理比较方便，查询效率较高

缺点：索引数量多，分片数太多，jvm使用率高，恢复慢

每天创建一个索引

按天创建索引，接入的应用系统通过Type来区分

优点：索引量小，管理方便

缺点：单索引数据量太大，查询效率低，各应用系统日志的留存时间必须保持一致。

关键系统设计 – 索引设计

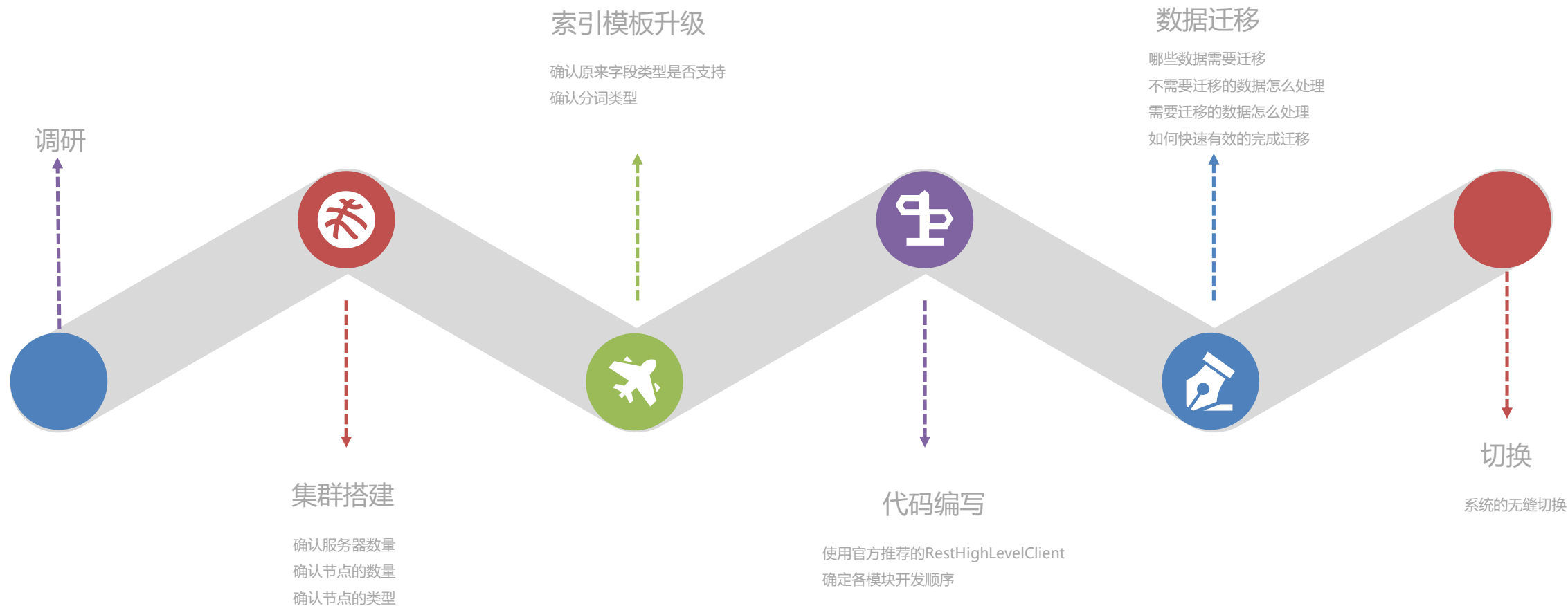




Elasticsearch
跨版本升级

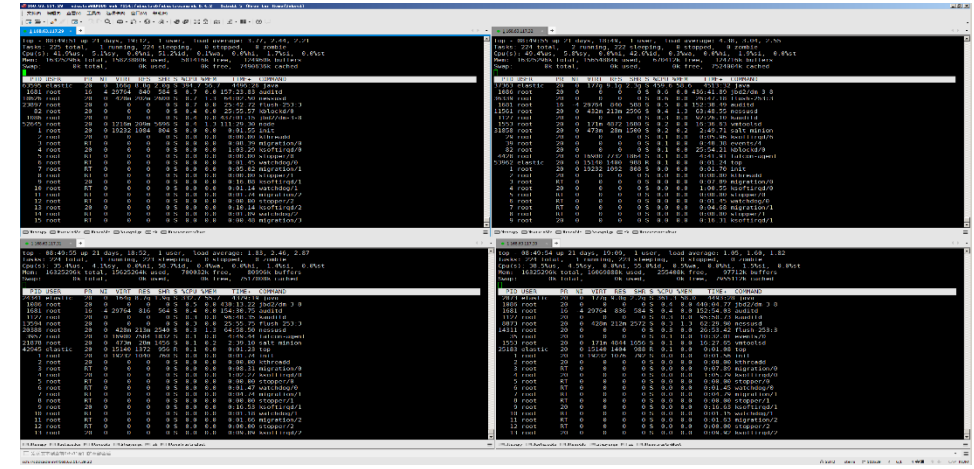
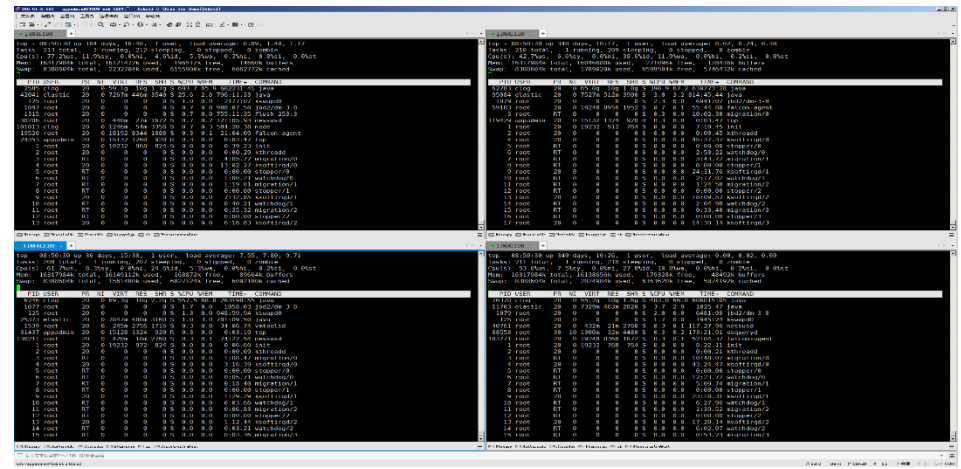
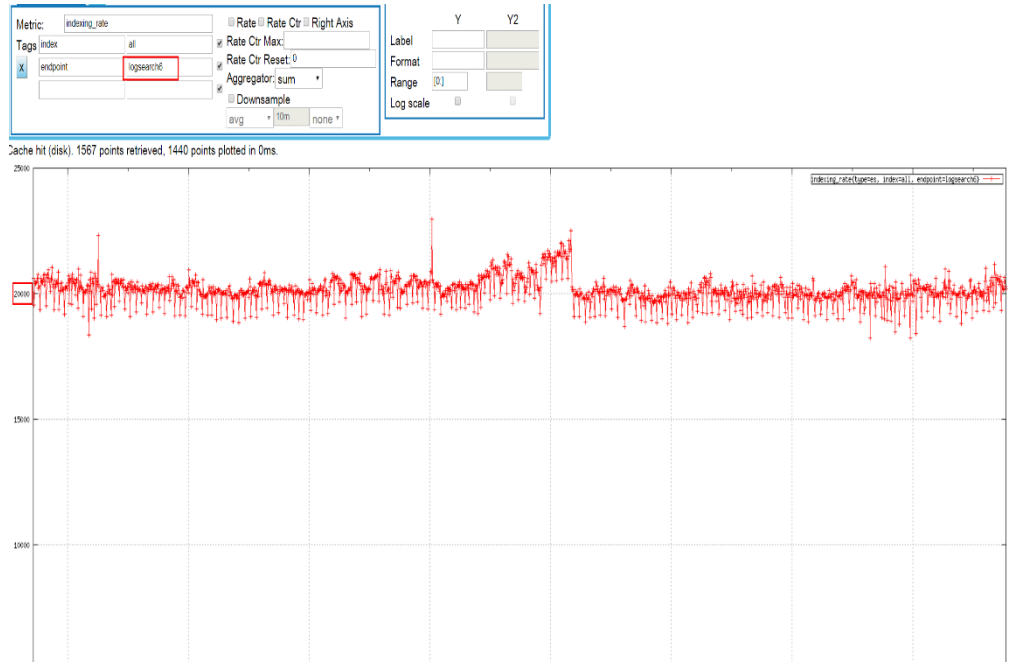
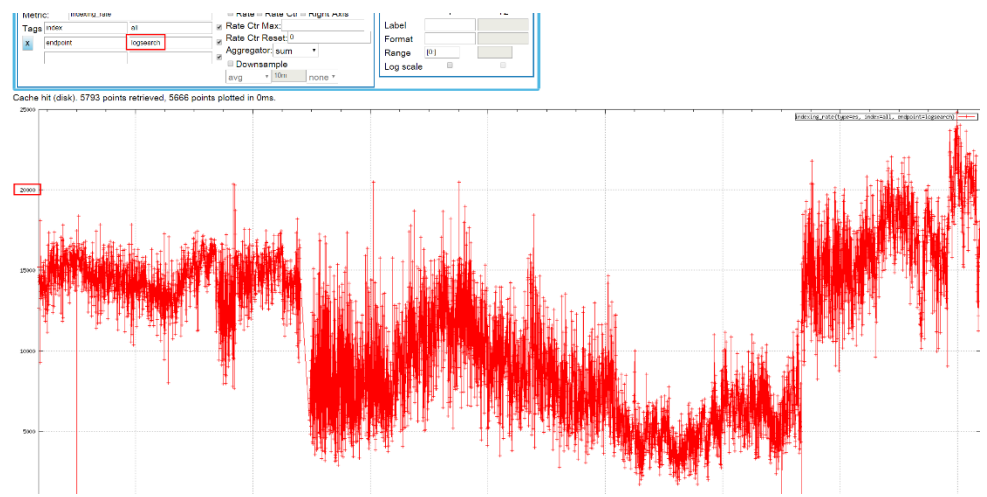
03

跨版本升级 – 升级步骤



跨版本升级 - 集群搭建

性能对比



跨版本升级 – 集群搭建

磁盘占用对比

http://16...:19200/

连接

uat_clog_es6

集群健康值: green (2119 of 2119)

概览

索引

数据浏览

基本查询 [+]

复合查询 [+]

indices

View Aliases

logstash-2018-11-21

logstash-2018-11-21

size: 69.9Gi (69.9Gi)
docs: 109,265,025 (109,265,025)

- es6-1 3
- es6-2 0
- ★ es6-3 2
- es6-4 1

http://16...:19200/_plugin/head/

连接

uat_clog_es

集群健康值: green (2450 of 2450)

概览

索引

数据浏览

基本查询 [+]

复合查询 [+]

indices

View Aliases

logstash-2018-11-21

logstash-2018-11-21

size: 107Gi (215Gi)
docs: 109,265,025 (218,639,426)

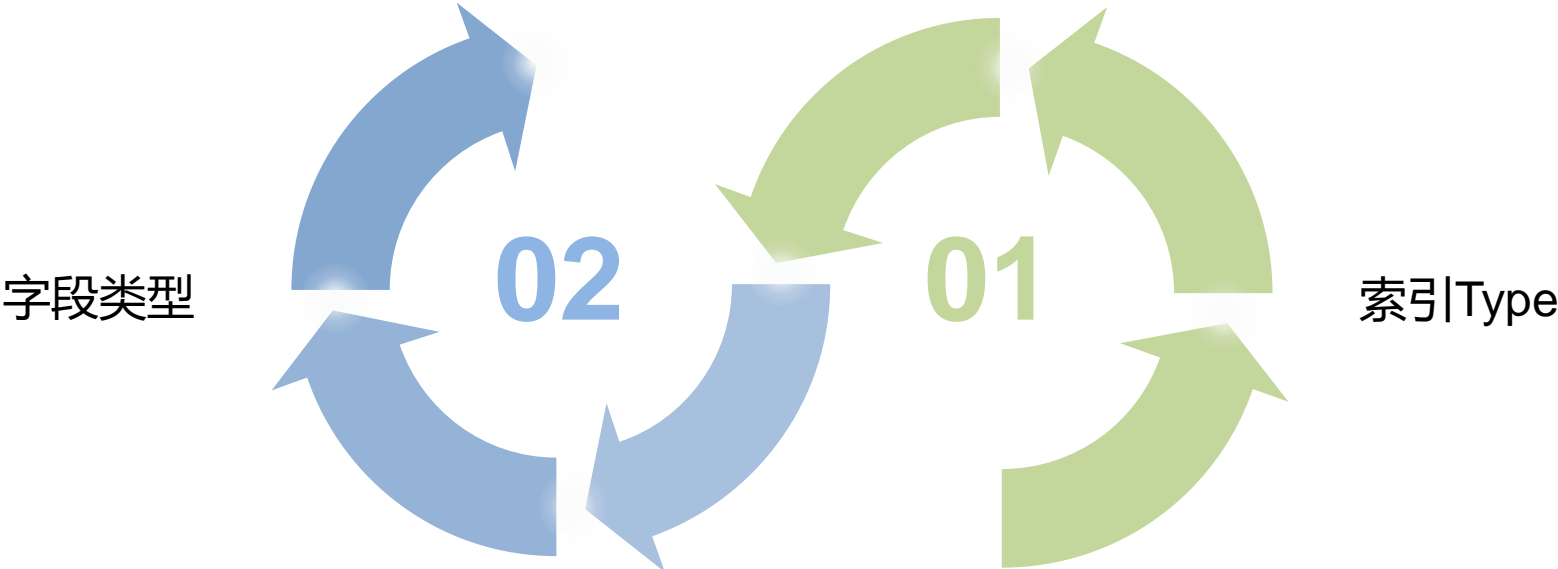
- ★ uat-es-1 1 2
- uat-es-2 0 1 4
- uat-es-3 2 3
- uat-es-4 0 3 4

跨版本升级 – 集群搭建

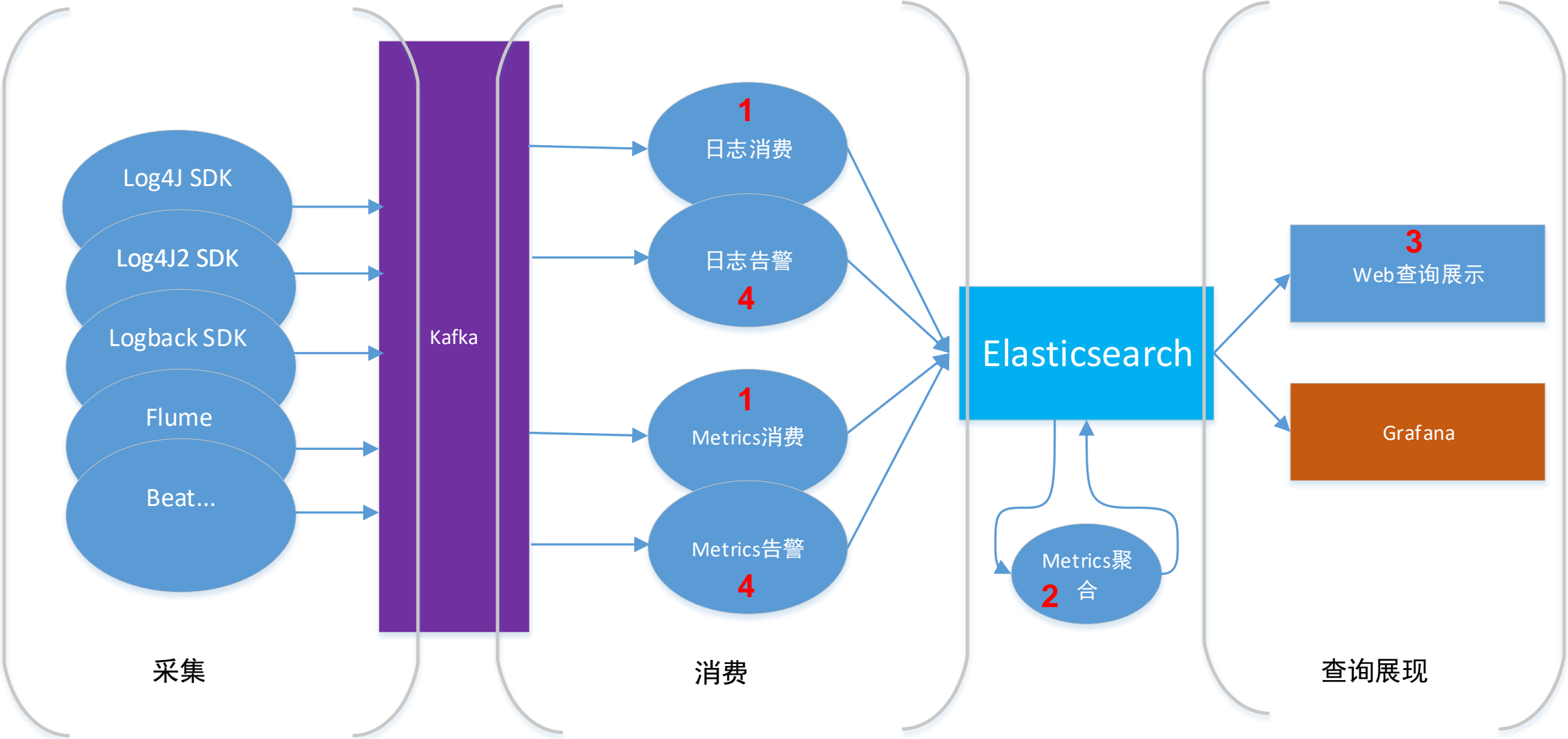
```
curl -X GET http://100.97.0.113:9200/_cat/nodes?v
```

ip	heap.percent	ram.percent	cpu	load_1m	load_5m	load_15m	node.role	master	name	
16	107	56	99	4	0.40	0.75	0.84	di	-	es7-1
16	146	60	99	2	0.69	0.85	0.90	m	*	es2
1	112	64	95	7	1.17	1.39	1.51	di	-	es4-1
1	113	53	97	6	1.60	1.37	1.40	di	-	es3-1
1	147	30	98	6	1.24	1.14	1.13	m	-	es1
1	108	67	99	4	0.56	0.64	0.74	di	-	es6-1
1	107	48	99	4	0.40	0.75	0.84	di	-	es7-2
1	147	46	98	6	1.24	1.14	1.13	di	-	es1-1
16	112	44	95	7	1.17	1.39	1.51	m	-	es4
16	146	65	99	3	0.69	0.85	0.90	di	-	es2-1
16	113	47	97	7	1.60	1.37	1.40	m	-	es3
16	108	63	99	4	0.56	0.64	0.74	di	-	es6-2
16	106	66	99	5	0.50	0.88	0.95	di	-	es8-1
16	111	72	96	6	1.90	1.98	1.72	di	-	es5-1
16	111	42	96	4	1.90	1.98	1.72	m	-	es5
1	106	52	99	4	0.50	0.88	0.95	di	-	es8-2

原来： 6台服务器， 3个master节点， 18的data节点， 70T存储空间
升级： 5台服务器， 5个master节点， 5(11)个data节点， 60T存储空间



跨版本升级 – 代码改造





跨版本升级 – 数据迁移



Reindex API



IMPORTANT

Reindex requires `_source` to be enabled for all documents in the source index.



IMPORTANT

Reindex does not attempt to set up the destination index. It does not copy the settings of the source index. You should set up the destination index prior to running a `_reindex` action, including setting up mappings, shard counts, replicas, etc.

The most basic form of `_reindex` just copies documents from one index to another. This will copy documents from the `twitter` index into the `new_twitter` index:

```
POST _reindex
{
  "source": {
    "index": "twitter"
  },
  "dest": {
    "index": "new_twitter"
  }
}
```

[COPY AS CURL](#) [VIEW IN CONSOLE](#)

Reindex from Remote

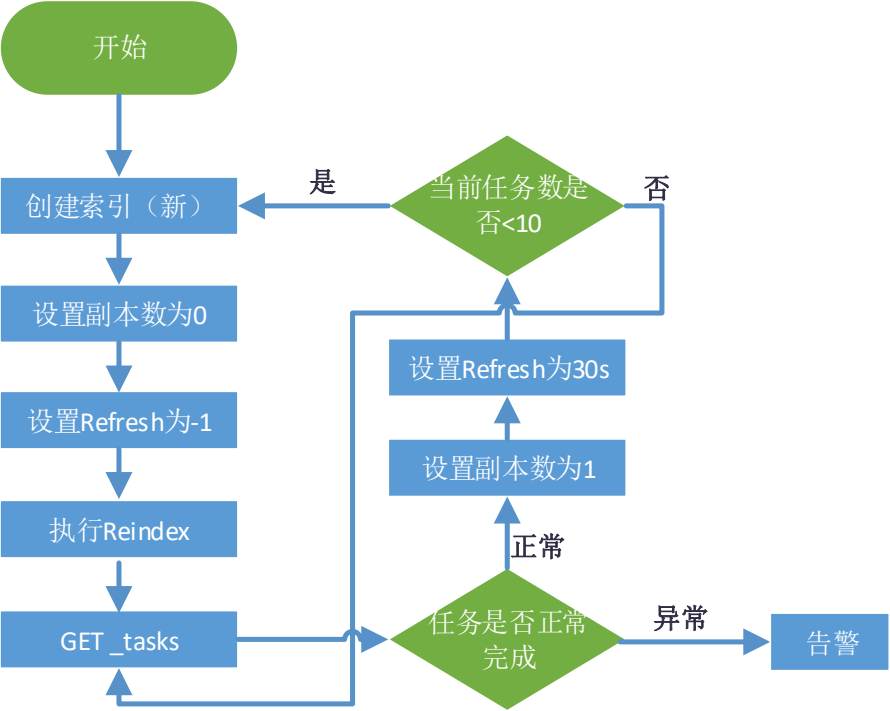


Reindex supports reindexing from a remote Elasticsearch cluster:

```
POST _reindex
{
  "source": {
    "remote": {
      "host": "http://otherhost:9200",
      "username": "user",
      "password": "pass"
    },
    "index": "source",
    "query": {
      "match": {
        "test": "data"
      }
    }
  },
  "dest": {
    "index": "dest"
  }
}
```

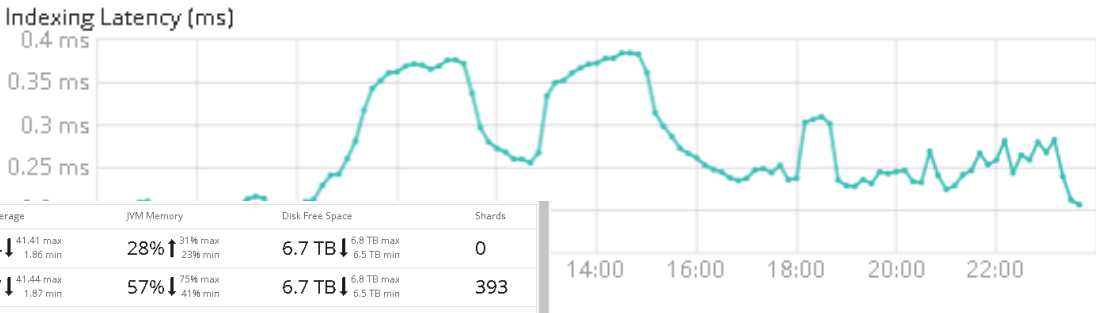
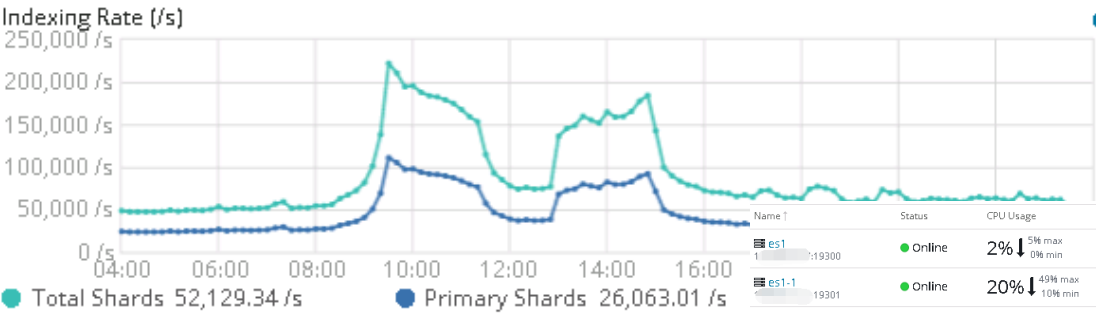
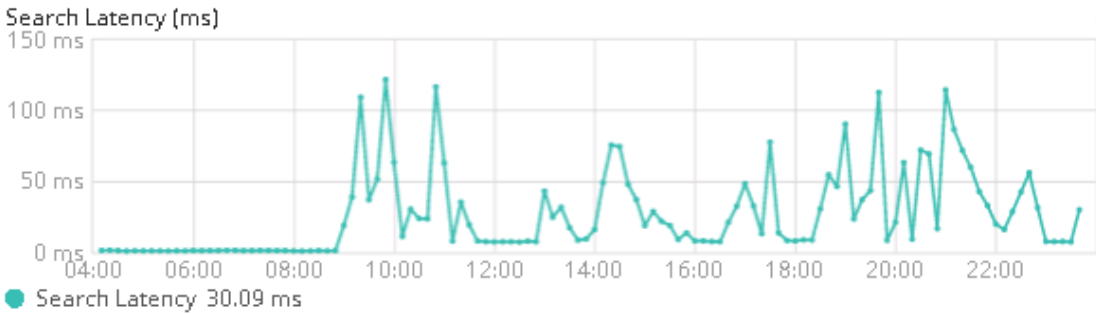
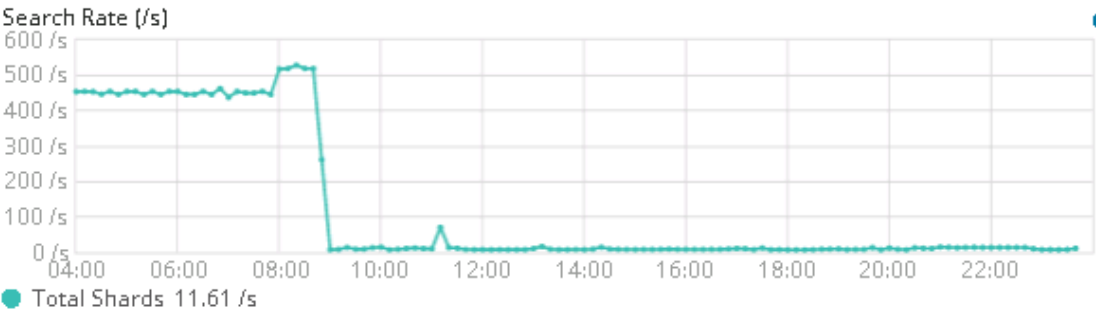
[COPY AS CURL](#) [VIEW IN CONSOLE](#)

跨版本升级 – 数据迁移



- 1. 创建索引
- 2. 设置副本数为0
`PUT xxxx_index/_settings { "number_of_replicas": 0 }`
- 3. 设置refresh时间间隔为不刷新
`PUTxxxx_index/_settings { "refresh_interval": -1 }`
- 4. 执行reindex操作
- 5. 定时通过task api 记录正在迁移的任务
`GET _tasks?detailed=true&actions=*reindex`
- 6. 判断已经完成的任务所对应的索引 日志条数是否一致。一致，设置副本数为1， refresh间隔为30s， 否则告警。
- 7.任务数小于10， 并且存在需要迁移的索引， 继续进行1操作

跨版本升级 – 现有集群



Name	Status	CPU Usage	Load Average	JVM Memory	Disk Free Space	Shards
es1	Online	2% ↓ 5% max 0% min	4.64 ↓ 41.41 max 1.86 min	28% ↓ 31% max 23% min	6.7 TB ↓ 6.8 TB max 6.5 TB min	0
es1-1	Online	20% ↓ 49% max 10% min	4.47 ↓ 41.44 max 1.87 min	57% ↓ 75% max 41% min	6.7 TB ↓ 6.8 TB max 6.5 TB min	393
es2	Online	1% ↓ 3% max 0% min	3.93 ↓ 34.66 max 1.79 min	74% ↓ 77% max 65% min	6.8 TB ↓ 6.8 TB max 6.5 TB min	0
es2-1	Online	14% ↓ 37% max 7% min	3.77 ↓ 34.66 max 1.79 min	74% ↓ 76% max 45% min	6.8 TB ↓ 6.8 TB max 6.5 TB min	392
es3	Online	2% ↓ 5% max 0% min	6.41 ↓ 18.56 max 1.76 min	61% ↓ 62% max 53% min	9.3 TB ↓ 9.3 TB max 9.0 TB min	0
es3-1	Online	18% ↓ 48% max 9% min	6.79 ↓ 18.72 max 1.69 min	60% ↓ 75% max 38% min	9.3 TB ↓ 9.3 TB max 9.0 TB min	393
es4	Online	2% ↓ 5% max 0% min	5.64 ↓ 9.39 max 1.79 min	56% ↓ 60% max 52% min	9.3 TB ↓ 9.3 TB max 9.0 TB min	0
es4-1	Online	21% ↓ 56% max 13% min	5.64 ↓ 9.39 max 1.79 min	54% ↓ 76% max 44% min	9.3 TB ↓ 9.3 TB max 9.0 TB min	393
es5	Online	2% ↓ 5% max 0% min	3.05 ↓ 11.56 max 1.68 min	45% ↓ 48% max 39% min	9.1 TB ↓ 9.1 TB max 8.9 TB min	0
es5-1	Online	20% ↓ 52% max 10% min	2.88 ↓ 12.05 max 1.74 min	50% ↓ 76% max 43% min	9.1 TB ↓ 9.1 TB max 8.9 TB min	393
es6-1	Online	5% ↑ 10% max 2% min	1.75 ↓ 29.05 max 0.4 min	52% ↓ 75% max 49% min	2.5 TB ↑ 2.7 TB max 2.5 TB min	484
es6-2	Online	7% ↑ 18% max 2% min	1.75 ↓ 29.14 max 0.4 min	59% ↓ 59% max 51% min	2.5 TB ↑ 2.8 TB max 2.5 TB min	484
es7-1	Online	5% ↑ 12% max 2% min	3.42 ↓ 27.31 max 0.5 min	54% ↓ 56% max 47% min	2.5 TB ↑ 2.8 TB max 2.5 TB min	485
es7-2	Online	4% ↓ 13% max 1% min	3.8 ↓ 26.2 max 0.52 min	49% ↓ 75% max 46% min	2.5 TB ↑ 2.7 TB max 2.5 TB min	484
es8-1	Online	6% ↑ 10% max 2% min	2.57 ↓ 27.97 max 0.42 min	53% ↓ 54% max 46% min	2.3 TB ↑ 2.5 TB max 2.3 TB min	483
es8-2	Online	6% ↓ 12% max 1% min	2.57 ↓ 27.97 max 0.42 min	64% ↓ 65% max 57% min	2.5 TB ↑ 2.7 TB max 2.5 TB min	484



坑和运维

04

A

GC问题

内存的设置和ES2不通，
导致采用的默认内存大小

B

Bulk操作报错

ES写入压力过大时，ES
会报错写入操作

C

Jvm使用率

随着分片数的增多，jvm
的使用率偏高，GC的越来越频繁

D

磁盘占用

裸盘的节点还是会出现使用率差距很大的问题导致告警

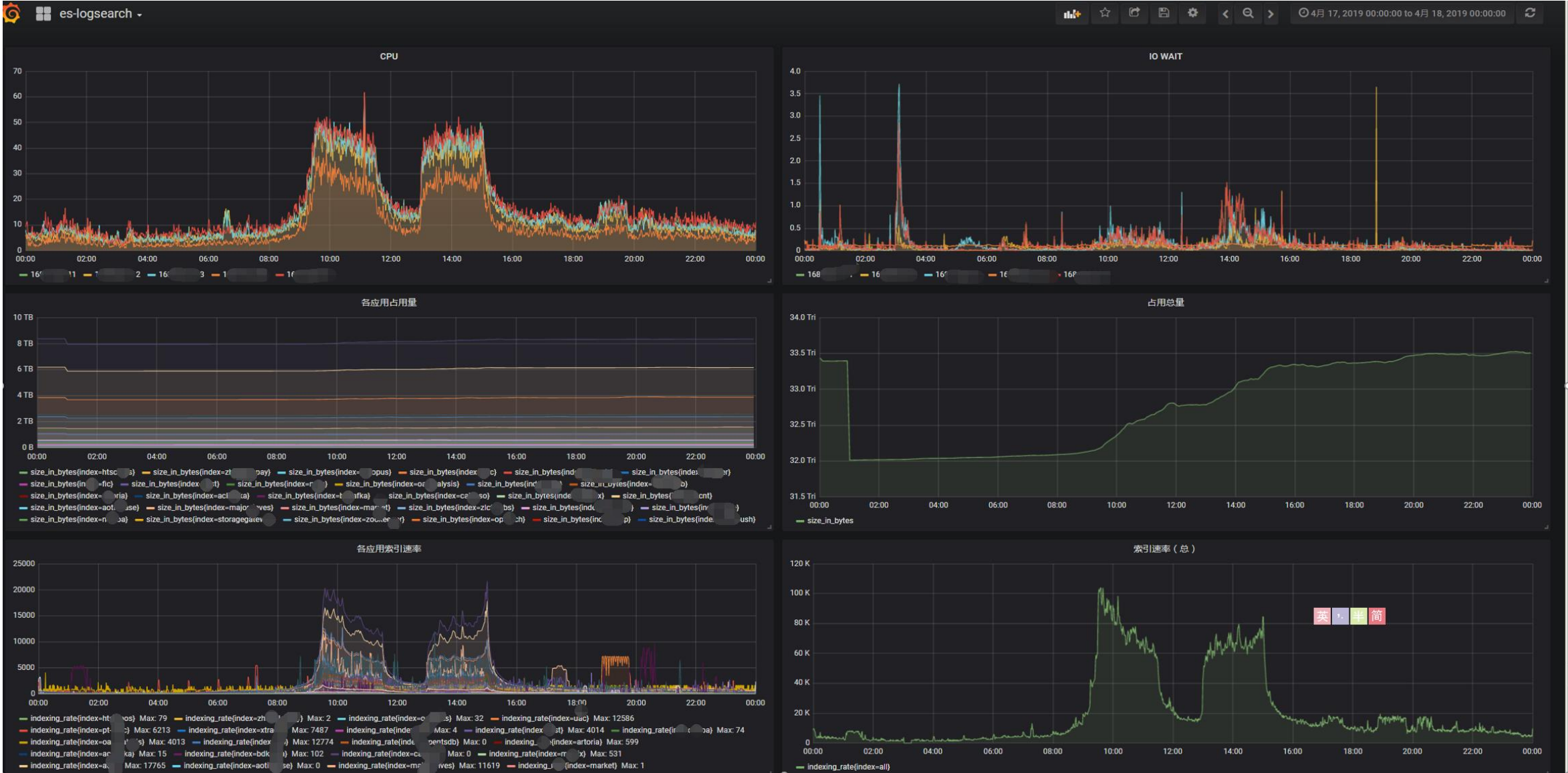
```
# xms represents the initial size of total heap space
# xmx represents the maximum size of total heap space

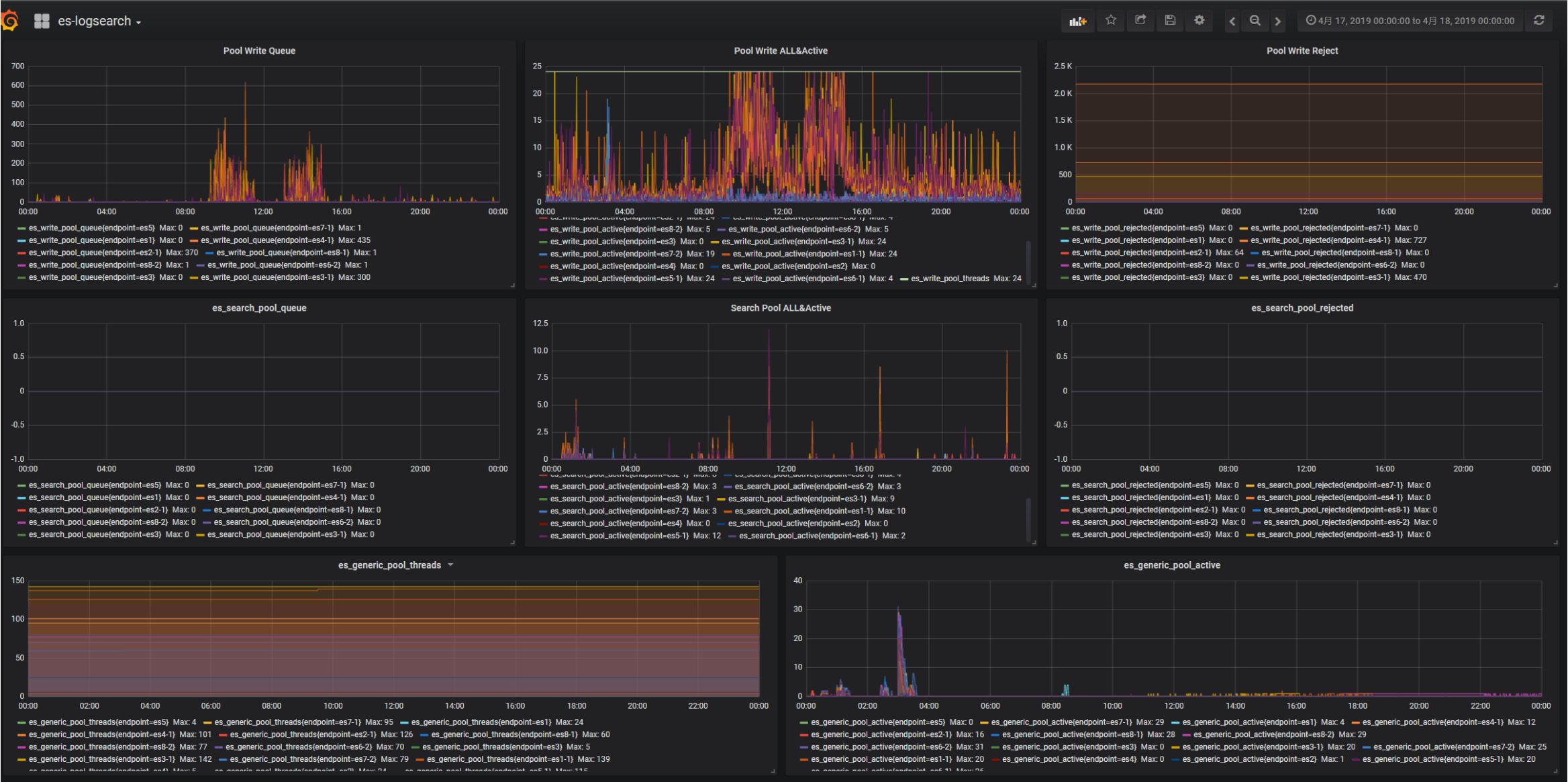
-Xms30g
-Xmx30g
```

```
indices.fielddata.cache.size: 1gb
indices.queries.cache.size: 1gb
```

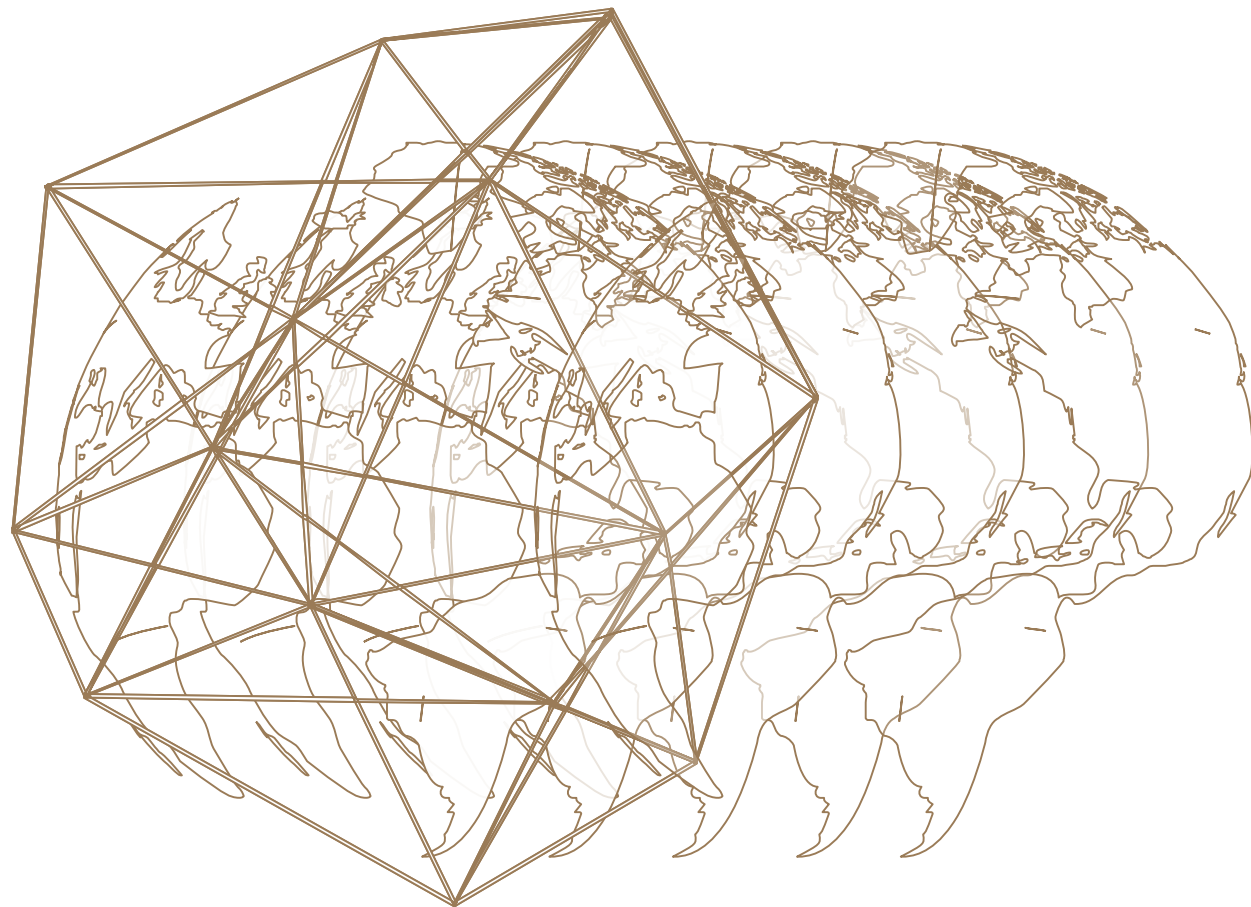
```
thread_pool:
  write:
    size: 24
    queue_size: 2000
```

```
"cluster": {
  "routing": {
    "allocation": {
      "disk": {
        "watermark": {
          "low": "75%",
          "high": "80%"
        }
      }
    }
  }
}
```





感谢聆听





专业、垂直、纯粹的 Elastic 开源技术交流社区
<https://elasticsearch.cn/>