


B站搜索自动化实践

管华新
BiliBili



战略级赞助商  HUAWEI

钻石级赞助商  普翔

白金级赞助商  华夏博格

 神州数码
Digital China

金牌级赞助商  iDataAPI

合作伙伴  开源中国
oschina.net

 掘金

 知乎

 IT大数据

 otpub

 Broadview
www.broadview.com.cn

 百格活动
bagevent.com

 MAXHUB
高效会议平台



B站搜索自动化实践



----管华新

目录

一 平台概况

二 痛点优化

三 SDK

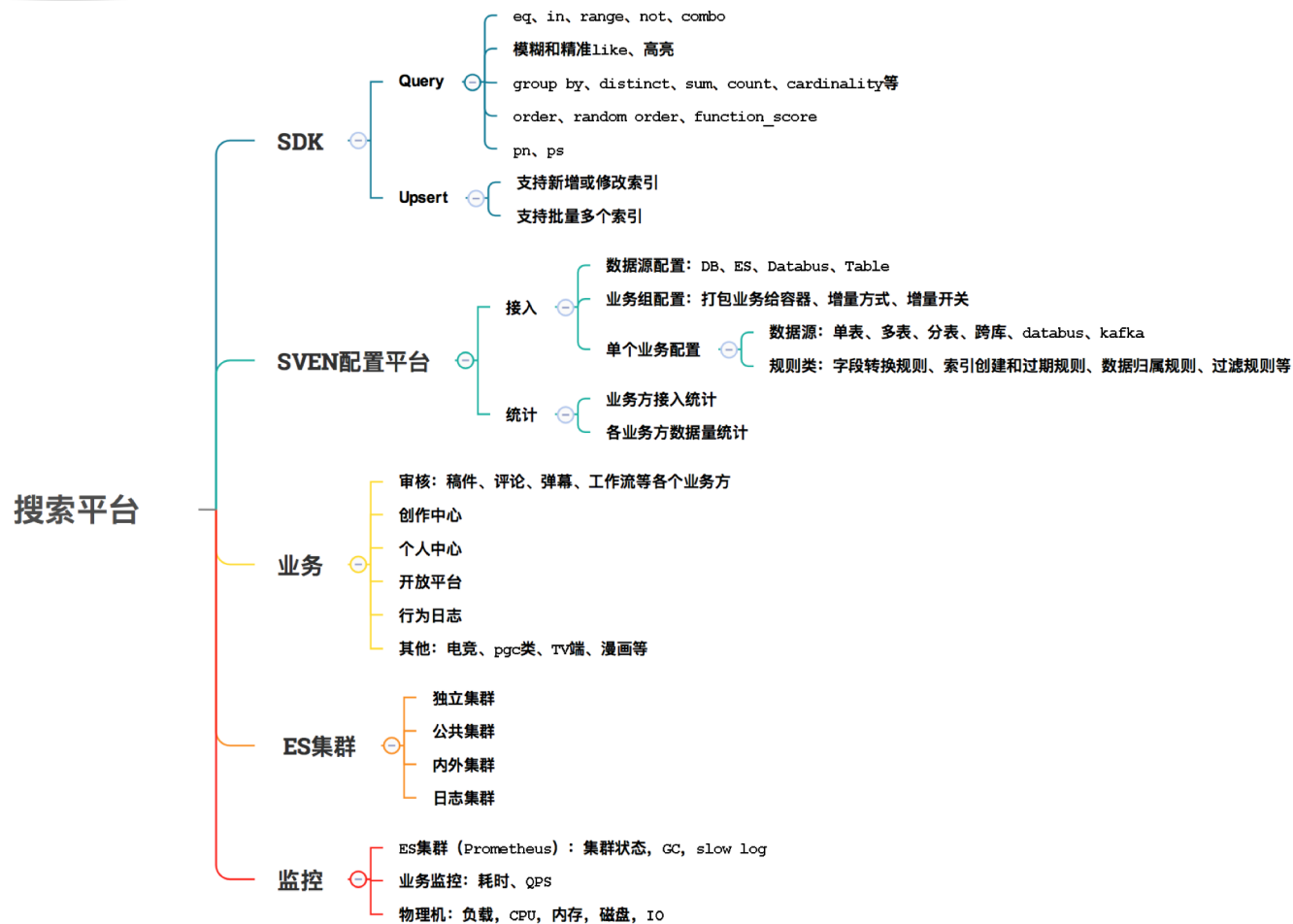
四 JOB

五 运维





平台概况





平台概况

- **数据：**
 1. 15个集群，45个节点
 2. 19条业务线，覆盖90个小业务
 3. 总量50+TB
 4. 5kw+条/天，1kw+次查询/天（对外请求占95%）
- **以前的问题：**
 1. 接入和迭代较慢，job和接口代码冗余
 2. 集群混用，相互耦合
 3. 单节点负载高，gc频繁





痛点优化

- 集群优化
- 查询性能优化
- 一致性和实时性优化
- 使用通用接口
- 使用通用配置





集群优化：解决单节点负载高（优化内存管理）

perf性能分析页回收频繁，而且free cache 低

```
Samples: 284k of event 'cycles', Event count (approx.): 6155373846
49.51% [kernel] isolate_lru_pages.isra.46
1.44% [kernel] isolate_lru_page
1.86% [kernel] shrink_inactive_list
1.61% [kernel] __d_lookup_rcu
0.78% libjvm.so __bixtat64
0.71% libjvm.so __blockoffsetArrayNonContigSpace::block_start_unsafe(void const*) const
0.68% [kernel] compactibleFreeListSpace::block_size(Heapword const*) const
0.68% [kernel] _raw_spin_lock_irq
0.65% libjvm.so link_path_walk
0.56% [kernel] jwm_getStackAccessControlContext
0.41% libjvm.so strncpy_from_user
0.29% libjvm.so ClearNonCleanCardWrapper::do_MemRegion(MemRegion)
0.28% libjvm.so ParNewGeneration::real_forwardee_low(oopDesc*)
0.27% libjvm.so ObjectMonitor::trySpinVaryDuration(Thread*)
0.27% [kernel] codeheap::find_start(void*) const
0.27% libjvm.so copy_user_enhanced_fast_string
0.25% [kernel] ParNewGeneration::copy_to_survivor_space_avoiding_promotion_undo(ParScanThreadState*, oopDesc*, unsigned long, markoopDesc*)
0.25% libjvm.so vframeStreamCommon::fill_from_compiled_frame(int)
0.25% [kernel] lookup_fast
0.24% libjvm.so java_lang_Class::protection_domain(oopDesc*)
0.24% libjvm.so spin_pause
0.24% libjvm.so BacktraceBuilder::push(Method*, int, Thread*)
0.23% [kernel] _raw_spin_lock
0.22% libjvm.so vframeStreamCommon::next()
0.22% libjvm.so java_lang_Throwable::fill_in_stack_trace(Handle, methodHandle, Thread*)
0.22% perf-162246.map 0x00007f9565f9a9c4
0.21% [kernel] kmem_cache_alloc
0.20% [kernel] kmem_cache_free
0.20% [kernel] _raw_spin_lock_irqsave
0.19% [kernel] __sk_run_filter
0.19% libjvm.so shrink_lruvec
0.18% perf-162246.map typeArrayClass::allocate_common(int, bool, Thread*)
0.18% libjvm.so 0x00007f956923c086
0.18% [kernel] InstanceClass::oop_oop_iterate_nv(oopDesc*, ParScanWithoutBarrierClosure*)
0.18% [kernel] find_busiest_group
0.17% [kernel] generic_permission
0.17% libjvm.so FreeList<FreeChunk>::getFirstNChunksFromList(unsigned long, FreeList<FreeChunk>*)
0.17% perf-162246.map 0x00007f956923d0de
0.16% libpthread-2.19.so __l1l_lock_elision
0.15% [kernel] path_init
0.15% libjvm.so PromotionInfo::promoted_oops_iterate_nv(ParScanWithoutBarrierClosure*)
0.15% perf-162246.map 0x00007f956923bd33
0.14% [kernel] lockref_get_not_dead
0.13% libc-2.19.so realpath
0.13% libjvm.so CardTableModRefBS::process_chunk_boundaries(Space*, DirtyCardToOopClosure*, MemRegion, MemRegion, signed char**, unsigned long, unsigned long)
0.13% libjvm.so InstanceClass::oop_oop_iterate_nv(oopDesc*, ParScanWithoutBarrierClosure*)
0.12% libjvm.so ParScanThreadState::trim_queues(int)
0.12% perf-162246.map 0x00007f956534274f
0.11% libjvm.so javaCalls::call_helper(JavaValue*, methodHandle*, javaCallArguments*, Thread*)
0.11% libjvm.so ConcurrentMarkSweepGeneration::par_promote(int, oopDesc*, markoopDesc*, unsigned long)
0.11% [kernel] security_inode_getattr
0.11% libjvm.so ObjectMonitor::NotRunnable(Thread*, Thread*)
0.10% [kernel] __inode_permission
0.10% [kernel] system_call
0.10% [kernel] __schedule
0.10% [kernel] system_call_after_swaps
0.10% perf-162246.map 0x00007f9565f9a9af
0.10% perf-162246.map 0x00007f9567c353c1
0.10% [kernel] xfs_vn_getattr
0.09% libjvm.so Frame::sender_for_compiled_frame(RegisterMap*) const
0.09% perf-162246.map 0x00007f95653425bb
0.09% libjvm.so ParallelTaskTerminator::offer_termination(TerminatorTerminator*)
0.09% [kernel] lru_add_drain_cpu
0.09% perf-162246.map 0x00007f9565f9a9af
```

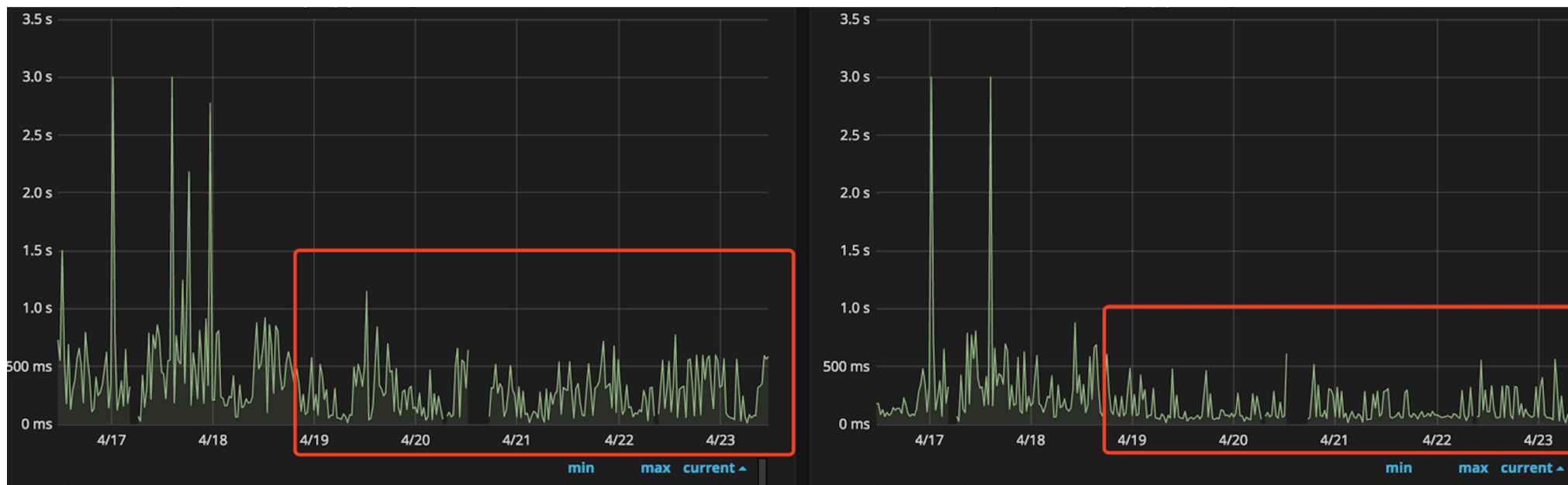
解决办法：

1. vm.swappiness = 1
2. vm.min_free_kbytes=655350
3. vm.zone_reclaim_mode = 0



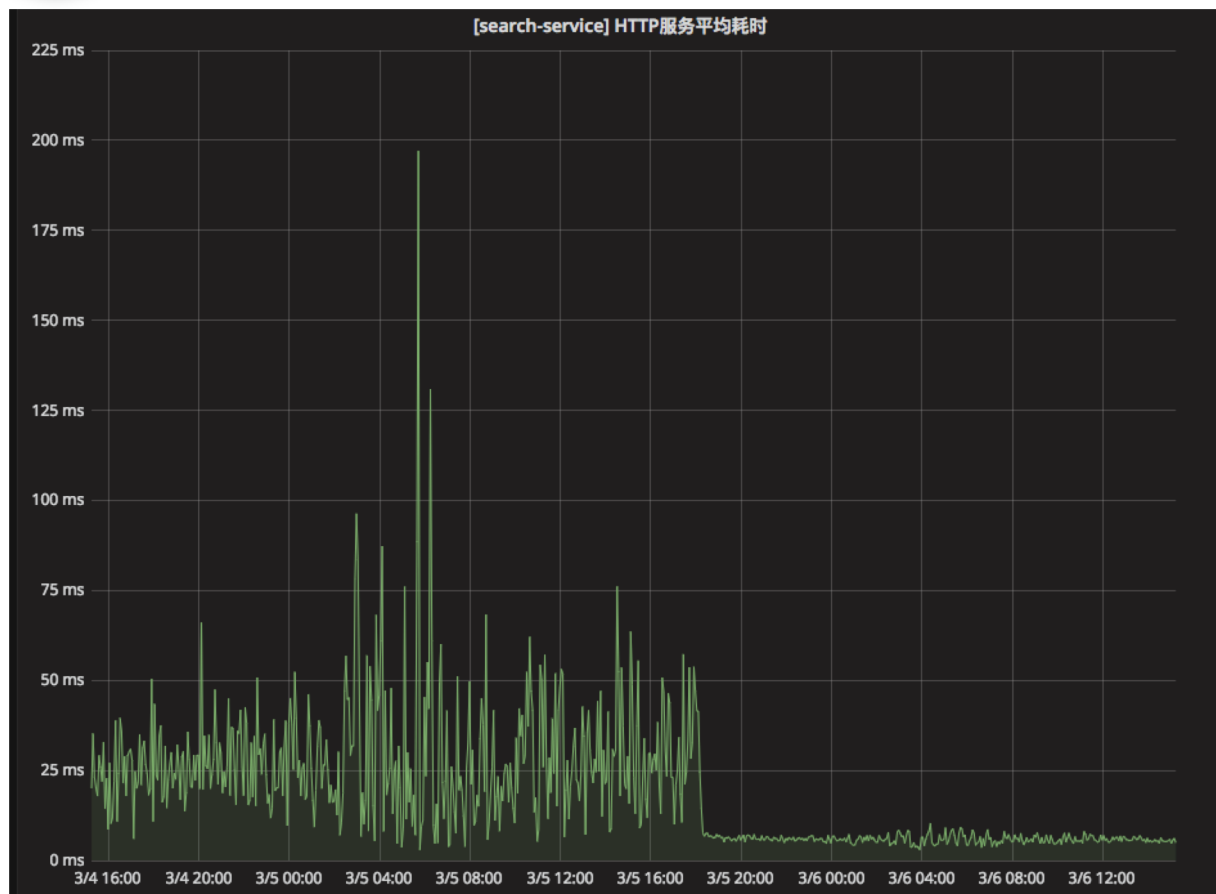


查询优化：Mapping优化 (int 转 keyword)





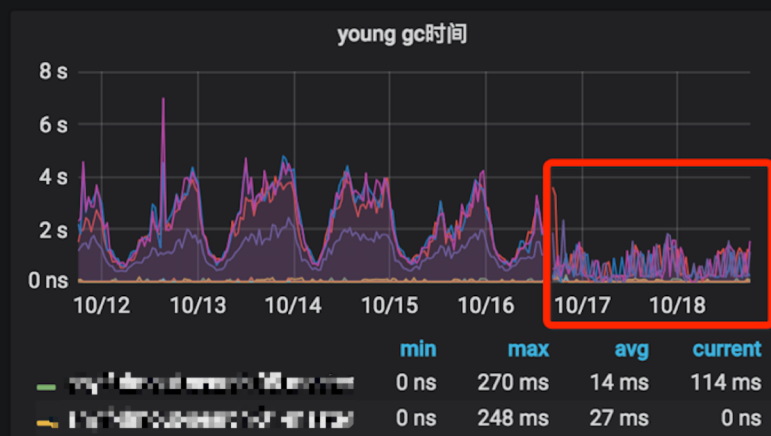
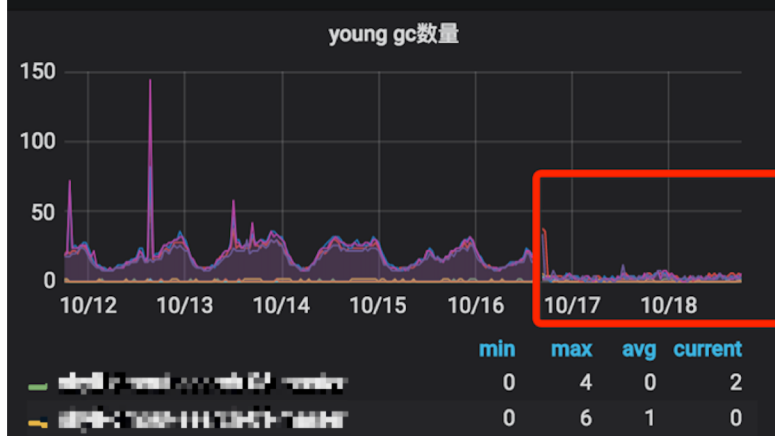
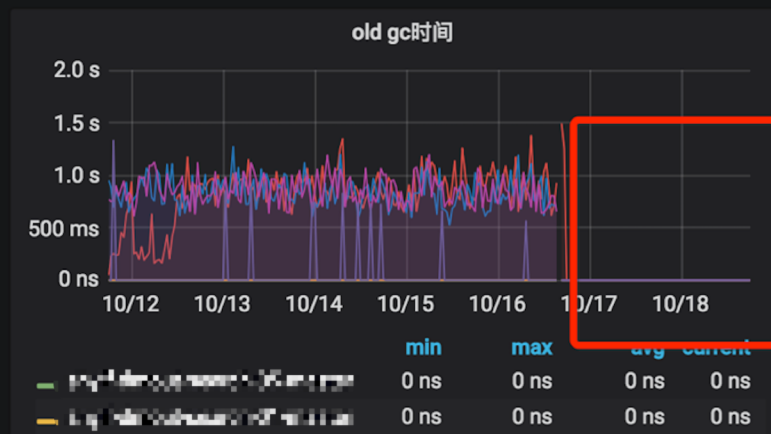
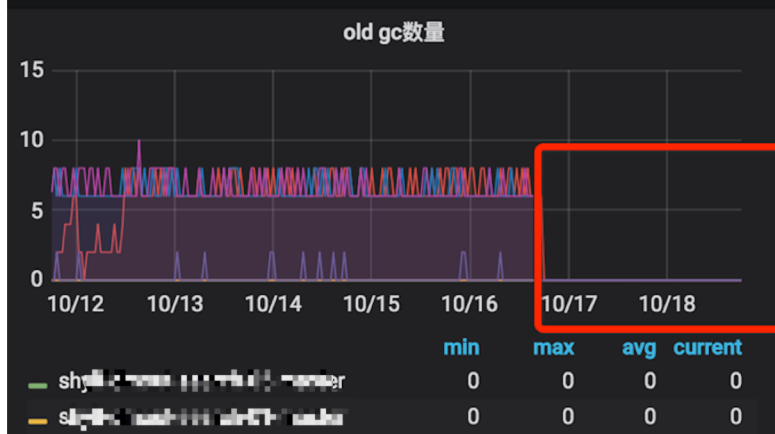
查询优化：SSD -> PCIE（接口响应平滑）





查询优化：CMS -> G1 GC (减少old gc和yong gc)

内存状态





查询优化：CMS -> G1 GC (减少old gc和young gc)



danielmitterdorfer Elastic Team Member

27d

Hi,

we will be supporting G1 with the upcoming version of Elasticsearch 6.5.0 but only if you run Elasticsearch on JDK 10 or later. For more context: the support has been added in [#33685](#) 17.

Daniel

1. ES6.5正式支持G1
2. 官方建议jdk10或以上版本使用G1



Add commented out JVM options for G1GC #33685

Merged dakrone merged 3 commits into elastic:master from dakrone:add-g1gc-for-jdk-10 on 22 Sep

Conversation 10 Commits 3 Checks 0 Files changed 1

Changes from all commits Jump to... +8 -0

Diff settings

Review changes

8 distribution/src/config/jvm.options		View
37	37	-XX:CMSInitiatingOccupancyFraction=75
38	38	-XX:+UseCMSInitiatingOccupancyOnly
39	39	
40	40	+ ## G1GC Configuration
41	41	+ # NOTE: G1GC is only supported on JDK version 10 or later.
42	42	+ # To use G1GC uncomment the lines below.
43	43	+ # 10-:-XX:-UseConcMarkSweepGC
44	44	+ # 10-:-XX:-UseCMSInitiatingOccupancyOnly
45	45	+ # 10-:-XX:+UseG1GC
46	46	+ # 10-:-XX:InitiatingHeapOccupancyPercent=75
47	47	+
48	48	## optimizations
49	49	
50	50	# pre-touch memory pages used by the JVM during initialization





优化：一致性和实时性优化

- 消息队列替代扫表增量
- 统一使用Bulk和doc_as_upsert，对版本conflict做retry
- 当消费慢的时候，业务方对实时性又比较敏感，建议业务方使用update接口阻塞式更新，或是业务方查缓存和db对数据剔除
- 增量和全量交给业务方，一致性交给他们，我们保障集群的稳定性





SDK：通用查询体特点

1. 加快接入和迭代

业务方按照自己的查询场景，自助写查询语句，现在90%业务做到0发版。

2. 简化复杂度

官方的DSL对业务方比较绕，而且容易采坑。简化nested，group by，cardinality，混合逻辑查询等的使用。

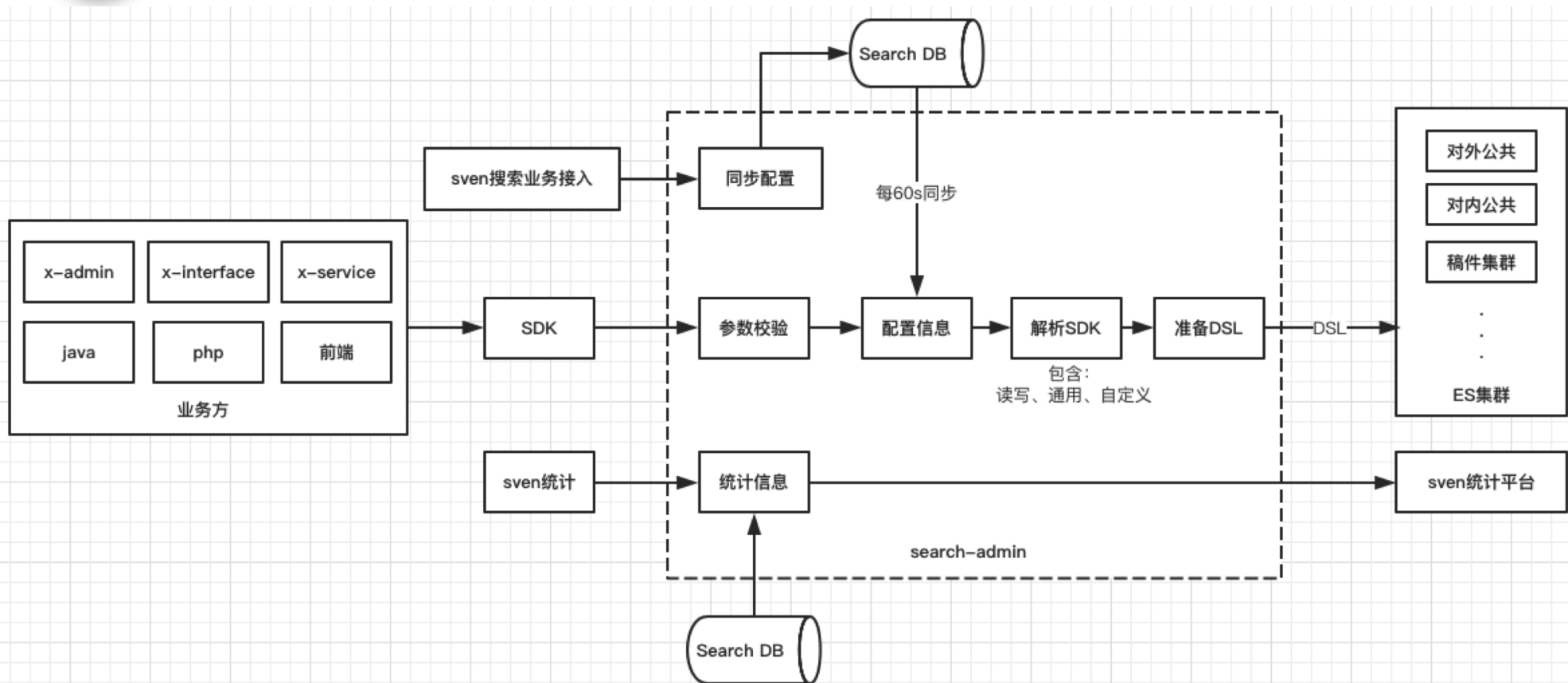
3. 安全、规范、效率

业务方直接使用DSL查询容易对集群有隐患，给不同的业务分配：最多查询索引个数、max_result_window、scroll、是否集群嗅探等权限点。





SDK: 接口服务设计





SDK : 查询功能

1. 过滤 : Eq、in、range、not、combo
2. 关键词 : 模糊、精准、高亮
3. 统计 : group by、distinct、sum、count、cardinality等
4. 排序 : order、random order、function_score
5. 分页 : pn、ps、scroll
6. 定制 : 继续使用通用SDK传参 , 只需要很少量的硬编码append到对应业务的通用查询中





SDK : 关键词查询

1. wildcard
2. ngram(2,2)
3. match





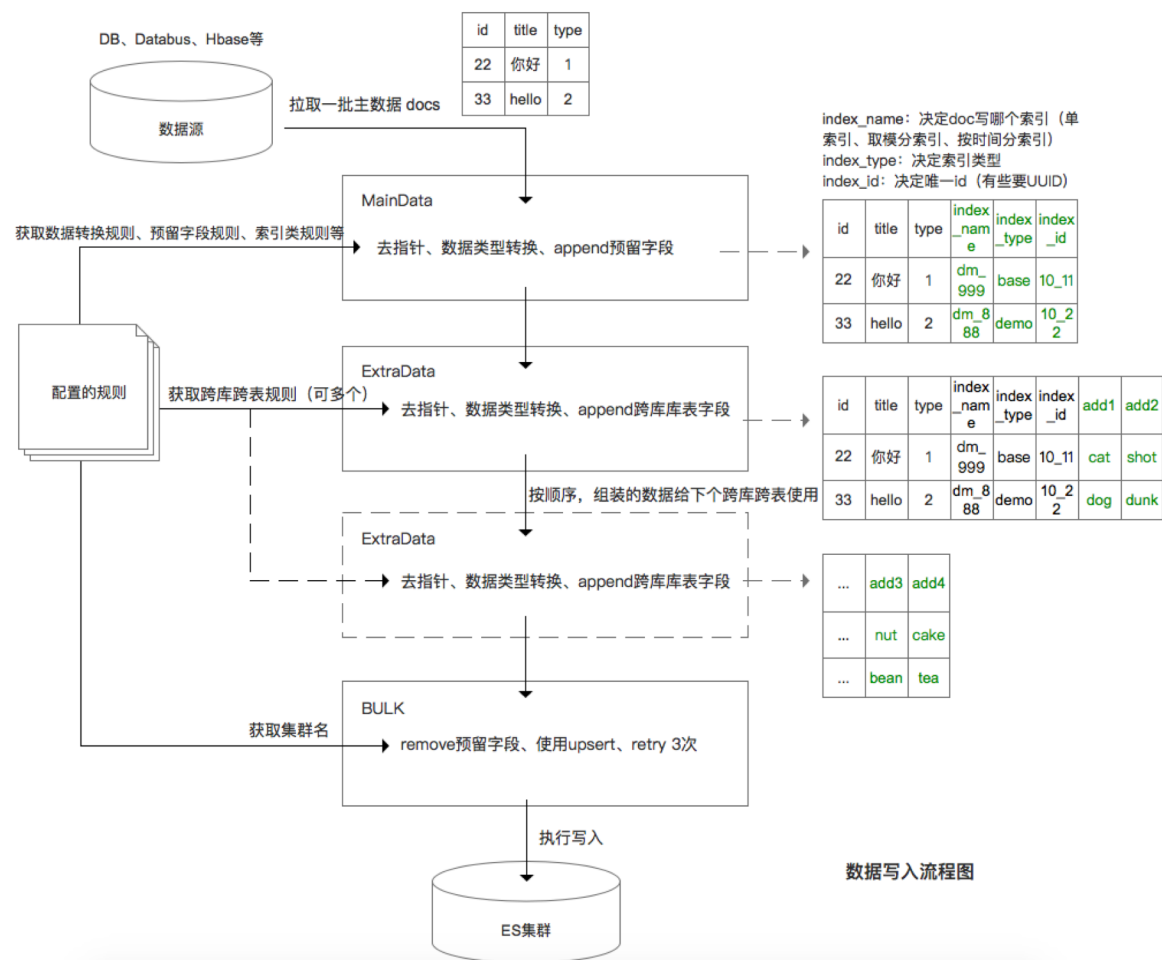
SDK: 通用 Update & Insert

简化BULK，面向多个索引批量更新、新增数据



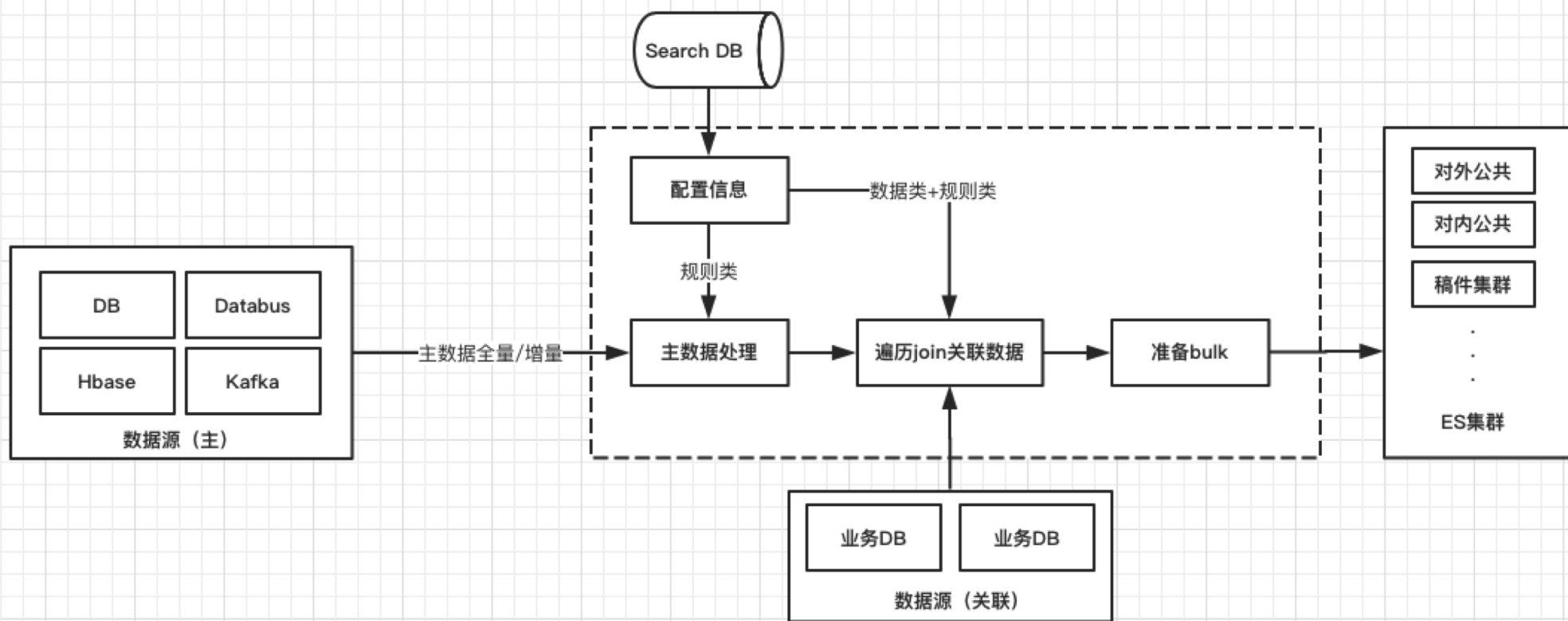
四

JOB: 数据聚合设计



四

JOB: 通用job服务设计



四

JOB: SEVN配置平台

首页 / 平台管理 / 搜索 / 业务组 / 应用配置

academy_archive academy_archive

基本信息 Mapping / 字段 跨表跨库

基础模块

AppId academy_archive

描述 创作中心 - up主学院

数据源模块

DB 库 bilibili_creative

ES 集群 pcie_pub_out01

Table bilibili_creative:academy_archive

Databus 请选择 Databus

索引模块

索引前缀 academy_archive

索引版本 _v1

索引格式 single,0,0,0,fixed

索引_type base

索引_id %d_%d,oid,business

基本信息 Mapping / 字段 跨表跨库

Mapping

```
1 {
2   "mappings": {
3     "base": {
4       "properties": {
5         "business": {
6           "type": "keyword"
7         },
8         "click": {
9           "type": "integer"
10        },
11        "copyright": {
12          "type": "keyword"
13        },
14        "fav": {
15          "type": "integer"
16        },
17        "oid": {
18          "type": "keyword"
19        },
20        "pubtime": {
21          "type": "date",
22          "format": "yyyy-MM-dd HH:mm:ss"
23        },
24        "state": {
25          "type": "keyword"
26        },
27        "arc_state": {
28          "type": "keyword"
29        },
30        "tid": {
31          "type": "keyword"
32        },
33        "tid_name": {
34          "type": "keyword"
35        },
36      }
37    }
38  }
39 }
```

字段

```
1 {
2   "es": "_id",
3   "field": "id",
4   "sql": "id as _id",
5   "expect": "int64",
6   "stored": "n",
7   "in_dtb": "n"
8 },
9 {
10  "es": "_mtime",
11  "field": "mtime",
12  "sql": "mtime as _mtime",
13  "expect": "time",
14  "stored": "n",
15  "in_dtb": "n"
16 },
17 {
18  "es": "id",
19  "field": "id",
20  "sql": "id",
21  "expect": "string",
22  "stored": "y",
23  "in_dtb": "y"
24 },
25 {
26  "es": "oid",
27  "field": "oid",
28  "sql": "oid",
29  "expect": "string",
30  "stored": "y",
31  "in_dtb": "y"
32 },
33 {
34  "es": "business",
35 }
```



四

JOB : 跨库跨表 (实时join)

基本信息

Mapping / 字段

跨表跨库

跨表跨库

```
1 {
2   {
3     "condition": {
4       "in_field": "oid",
5       "include": "business=1"
6     },
7     "dbname": "*****",
8     "in_field": "id",
9     "fields_str": "id:id:int64:y,mid:mid:int64:y,title:title:string:y,pubtime:pubtime:string:y,copyright:copyright:string:y,arc_state:S",
10    "remove_fields": [
11      "id"
12    ],
13    "fields": [
14      "id",
15      "mid",
16      "title",
17      "pubtime",
18      "copyright",
19      "arc_state"
20    ],
21    "sql": "select id,mid, title, pubtime,copyright,state as arc_state from ***** where id in (%s)"
22  },
23  }
```



四

JOB: 索引平滑升级

1. job初始化索引时，实体索引统一加版本号，加别名。默认使用别名读写。
2. 修改索引mapping时，新建带有新版本号的实体索引，双写两个索引，通过控制切换别名的方式，达到切换mapping的目的



四

JOB: 支持索引管理

1. 单个索引、按表名一致方式建索引（单表或分表都可以）、取模拆分、按时间拆分（年月日周）
2. 不同业务方，配置单次请求最多索引个数，默认为1



五

运维：部署和监控

- Ansible统一部署，es版本一致
- ES-Exporter接入Prometheus，集群监控、GC、线程池等增加阈值告警
- 物理机监控，内存、CPU、磁盘等阈值告警
- 业务方接口监控，耗时告警
- 消息队列监控，延迟告警



五

运维：索引管理

- 基于官方curator
- 面向按时间拆分的索引，提前创建、别名、template（年月日周）
- 定期挪动冷索引到冷节点
- 定期删除过期索引



谢谢





elastic 中文社区

专业、垂直、纯粹的 Elastic 开源技术交流社区

<https://elasticsearch.cn/>

