

# 苏宁ES平台实战

韩宝君

苏宁大数据平台

## 内容概要



平台介绍

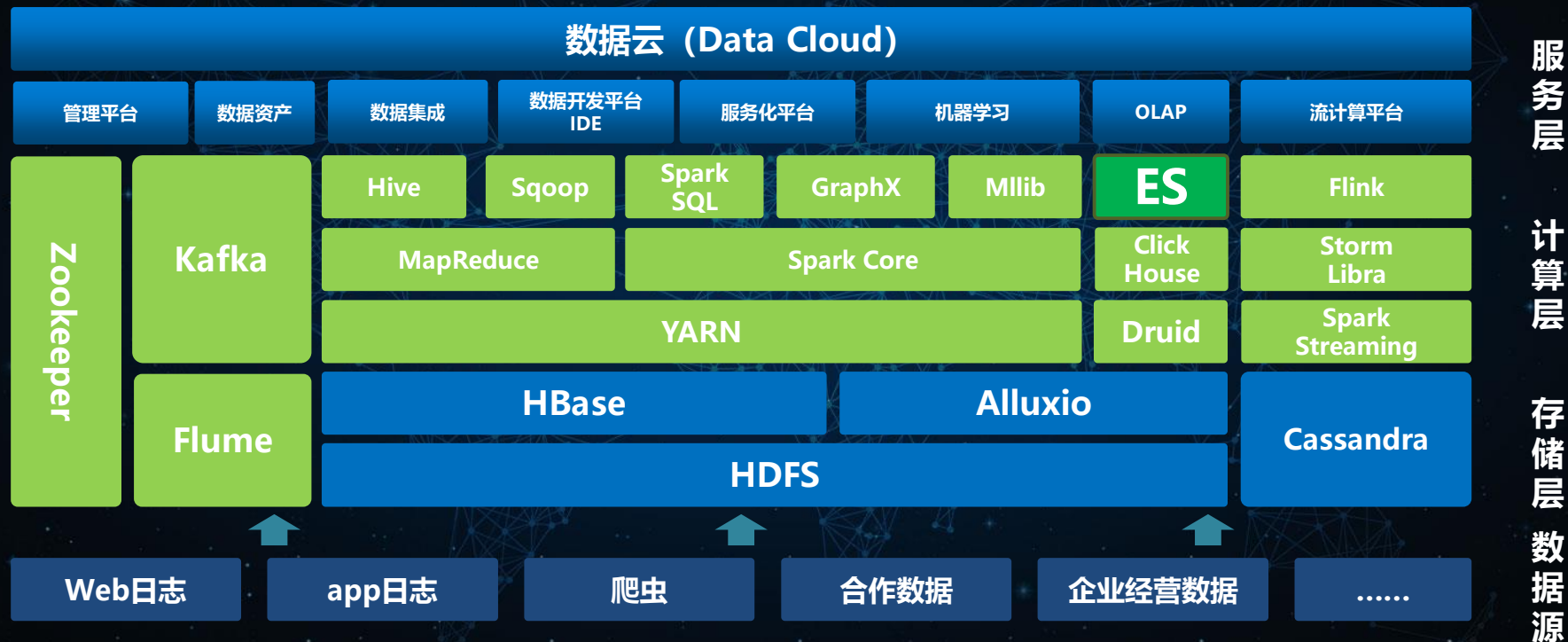


平台实战

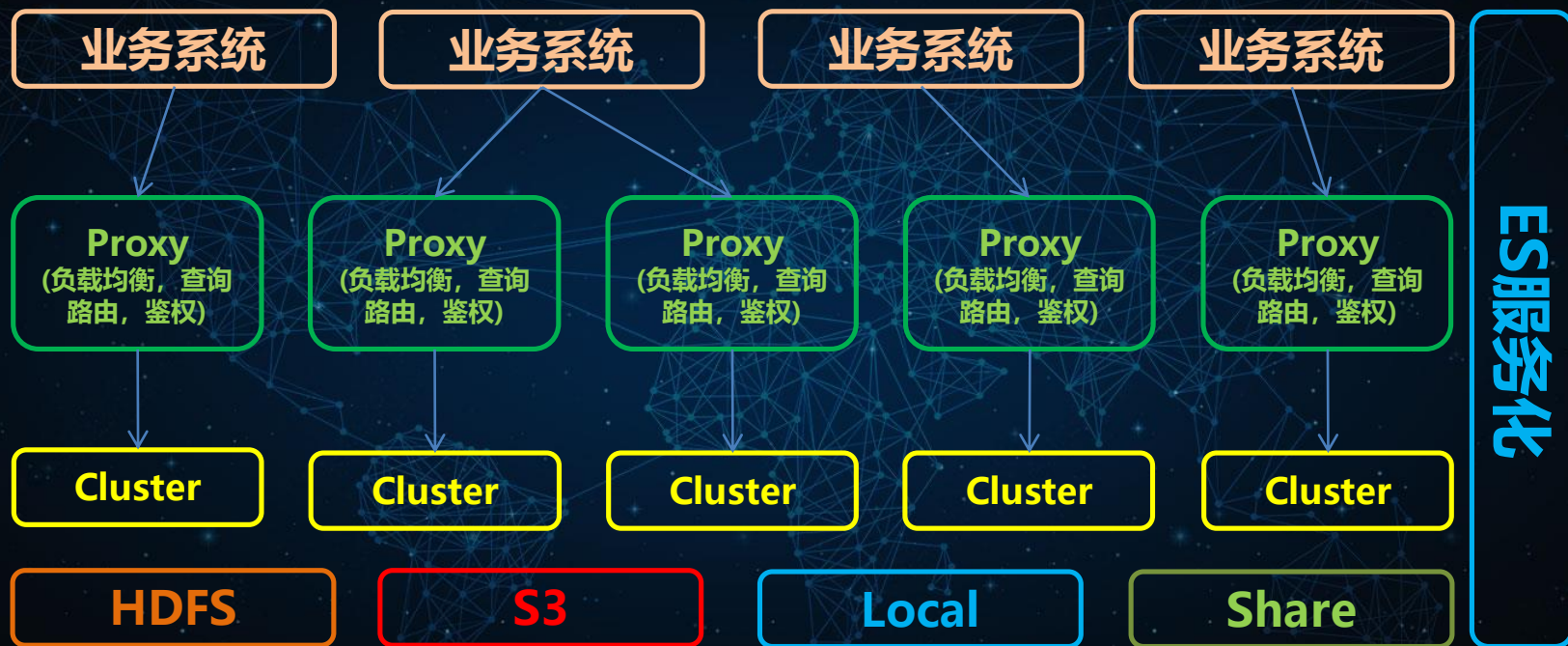


未来展望

# 苏宁大数据平台



## ES平台架构



# 平台规模

## 支持的业务

物流

(明细查询)

鹰眼

(流量分析)

金服

(订单搜索)

## 集群规模

500+物理机

30万+shard

80000+index

存储600+TB



## 内容概要



平台介绍



平台实战



未来展望

---

# 内核优化

---

## 问题概览

---

# 内核优化

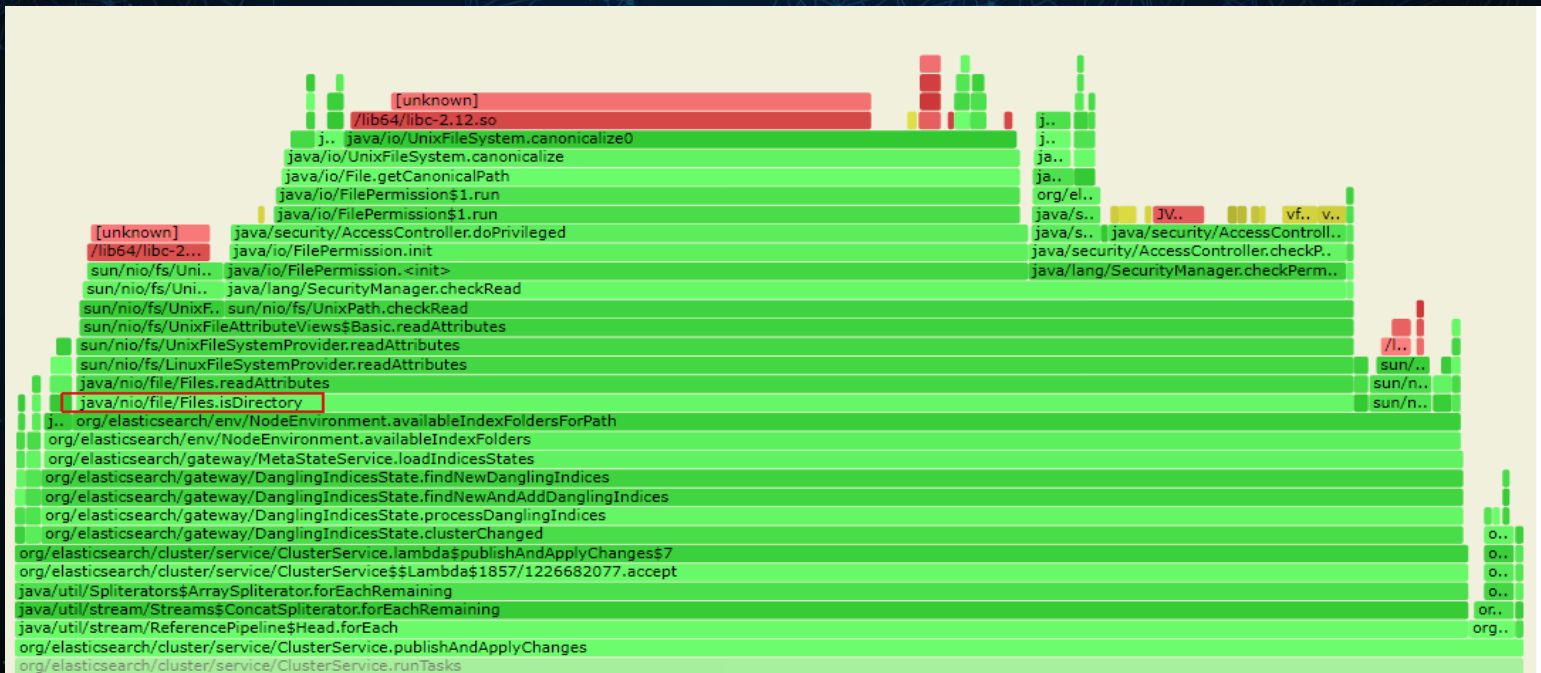
# 集群瓶颈

- 元数据操作任务堆积超时
- 遇到瓶颈的集群概览:
  - nodes: 37 (master: 3, data: 34)
  - indices:: 50000+
  - shards: 280000+
  - 创建索引、删除索引、增加删除别名等元数据操作超过30s, 元数据的task在队列中的等待时间甚至几分钟

insertOrder	timeInQueue	priority	source
96568	1.3m	URGENT	create-index [brmem_bc_m_star_1], cause [api]
96569	1.1m	URGENT	create-index [sts_bc_m_star_4], cause [api]
96571	50.3s	URGENT	create-index [brmem_bc_m_star_1], cause [api]
96570	1m	URGENT	create-index [sts_bc_m_11], cause [api]
96575	37s	URGENT	create-index [sts_bc_m_4], cause [api]
96572	38.1s	URGENT	shard-started
96573	38.1s	URGENT	shard-started
96574	38s	URGENT	shard-started
96577	32.5s	HIGH	add_listener
96576	33.6s	URGENT	create-index [sts_bc_m_star_11], cause [api]
96578	20.2s	URGENT	create-index [brmem_bc_1], cause [api]
96579	7s	URGENT	create-index [sts_bc_4], cause [api]
96580	3.5s	URGENT	create-index [sts_bc_11], cause [api]

# 火焰图统计

Files.isDirectory()这个方法占用百分之十几的时间



## 优化1: 僵尸索引

方案: 增加参数控制, 定时扫描  
dangling.indices.enabled: true  
dangling.indices.interval: "2h"

集群元数  
据变更

发布集群  
状态

通知  
监听器

僵尸索引  
扫描

延迟扫描

# 堆栈查看

线程长时间卡在awaitAllNodes方法处

```
java.lang.Thread.State: TIMED_WAITING (parking)
  at sun.misc.Unsafe.park(Native Method)
  - parking to wait for <0x00000000d2f231f8> (a java.util.concurrent.CountDownLatch$Sync)
  at java.util.concurrent.locks.LockSupport.parkNanos(LockSupport.java:215)
  at java.util.concurrent.locks.AbstractQueuedSynchronizer.doAcquireSharedNanos(AbstractQueuedSynchronizer.java:1064)
  at java.util.concurrent.locks.AbstractQueuedSynchronizer.tryAcquireSharedNanos(AbstractQueuedSynchronizer.java:1225)
  at java.util.concurrent.CountDownLatch.await(CountDownLatch.java:277)
  at org.elasticsearch.discovery.BlockingClusterStatePublishResponseHandler.awaitAllNodes(BlockingClusterStatePublishResponseHandler.java:100)
  at org.elasticsearch.discovery.zen.PublishClusterStateAction.innerPublish(PublishClusterStateAction.java:200)
  at org.elasticsearch.discovery.zen.PublishClusterStateAction.publish(PublishClusterStateAction.java:167)
  at org.elasticsearch.discovery.zen.ZenDiscovery.publish(ZenDiscovery.java:320)
  at org.elasticsearch.node.Node$$Lambda$1721/1471857648.accept(Unknown Source)
  at org.elasticsearch.cluster.service.ClusterService.publishAndApplyChanges(ClusterService.java:796)
  at org.elasticsearch.cluster.service.ClusterService.runTasks(ClusterService.java:640)
  at org.elasticsearch.cluster.service.ClusterService$ClusterServiceTaskBatcher.run(ClusterService.java:316)
  at org.elasticsearch.cluster.service.TaskBatcher.runIfNotProcessed(TaskBatcher.java:173)
  at org.elasticsearch.cluster.service.TaskBatcher$BatchedTask.run(TaskBatcher.java:211)
  at org.elasticsearch.common.util.concurrent.ThreadContext$ContextPreservingRunnable.run(ThreadContext.java:661)
  at org.elasticsearch.common.util.concurrent.PrioritizedThreadPoolExecutor$TieBreakingPrioritizedRunnable.run(PrioritizedThreadPoolExecutor.java:130)
  at org.elasticsearch.common.util.concurrent.PrioritizedThreadPoolExecutor$TieBreakingPrioritizedRunnable.run(PrioritizedThreadPoolExecutor.java:130)
  at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor.java:1142)
  at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor.java:617)
  at java.lang.Thread.run(Thread.java:748)
```

在日志中也可以看出等待的节点

```
2019-04-22T00:36:04.105 WARN [elasticsearch] [clusterService#updateTask][T#30] org.elasticsearch.discovery.zen.PublishClusterStateAction:innerPublish:212 - timed out waiting for all nodes to process published state [628929] (timeout [30s], pending nodes: [
  {eiUZR1pRXKblcdwtjcoCQ}{QBTL4_h3QRuSLAjcQI8g_A}{
    ml.enabled=true
  }])
```

## 优化2：元数据提交

方案：只等待全部master和大部分data节点



## 继续查看堆栈

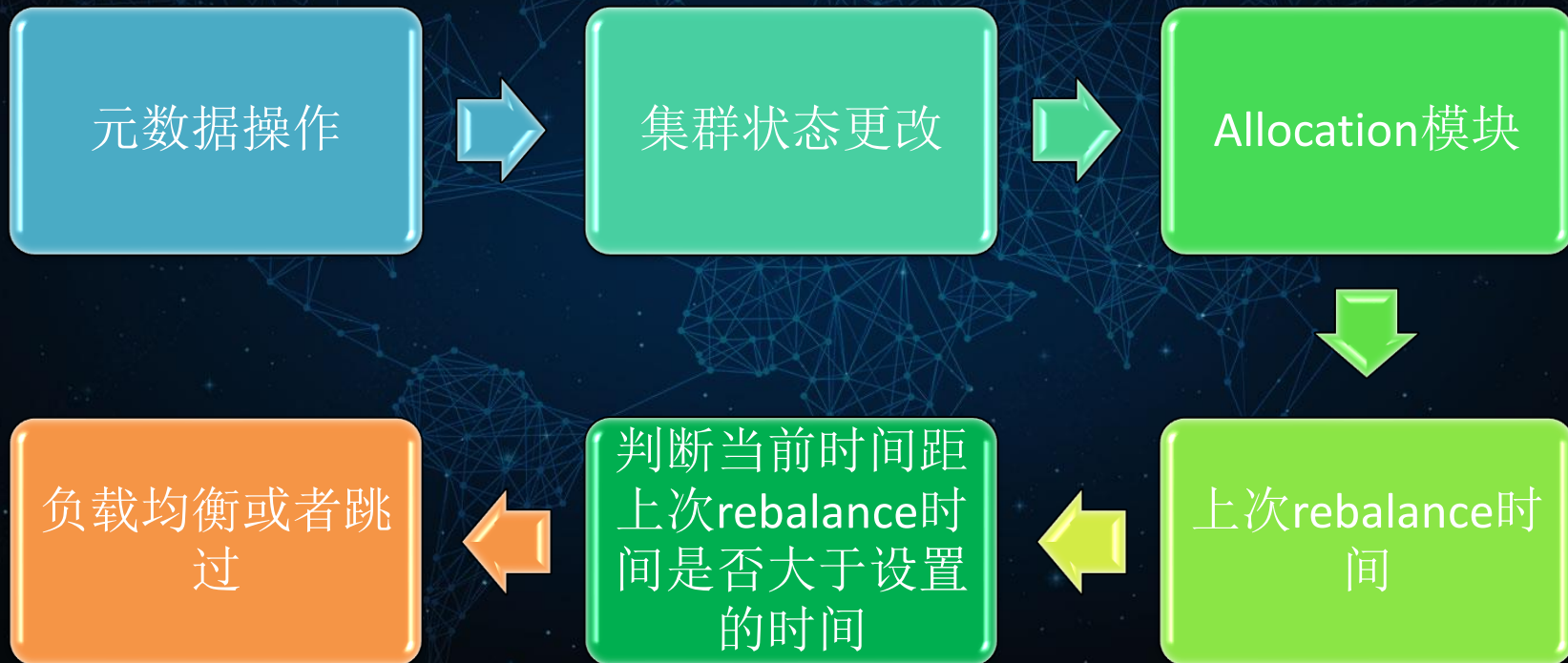
线程长时间卡在balanceByWeights方法处

```
java.lang.Thread.State: RUNNABLE
    at java.util.Collections$UnmodifiableCollection$1.<init>(Collections.java:1039)
    at java.util.Collections$UnmodifiableCollection.iterator(Collections.java:1038)
    at org.elasticsearch.cluster.routing.allocation.decider.AllocationDeciders.canAllocate(AllocationDeciders.java:72)
    at org.elasticsearch.cluster.routing.allocation allocator.BalancedShardsAllocator$Balancer.tryRelocateShard(BalancedShardsAllocator$Balancer.java:100)
    at org.elasticsearch.cluster.routing.allocation allocator.BalancedShardsAllocator$Balancer.balanceByWeights(BalancedShardsAllocator$Balancer.java:100)
    at org.elasticsearch.cluster.routing.allocation allocator.BalancedShardsAllocator$Balancer.balance(BalancedShardsAllocator$Balancer.java:100)
    at org.elasticsearch.cluster.routing.allocation allocator.BalancedShardsAllocator$Balancer.access$100(BalancedShardsAllocator$Balancer.java:100)
    at org.elasticsearch.cluster.routing.allocation allocator.BalancedShardsAllocator.allocate(BalancedShardsAllocator.java:100)
    at org.elasticsearch.cluster.routing.allocation.AllocationService.reroute(AllocationService.java:396)
    at org.elasticsearch.cluster.routing.allocation.AllocationService.reroute(AllocationService.java:363)
    at org.elasticsearch.cluster.routing.allocation.AllocationService.reroute(AllocationService.java:334)
    at org.elasticsearch.cluster.routing.allocation.AllocationService$1.execute(AllocationService.java:112)
    at org.elasticsearch.cluster.ClusterStateUpdateTask.execute(ClusterStateUpdateTask.java:45)
    at org.elasticsearch.cluster.service.ClusterService.executeTasks(ClusterService.java:687)
    at org.elasticsearch.cluster.service.ClusterService.calculateTaskOutputs(ClusterService.java:665)
    at org.elasticsearch.cluster.service.ClusterService.runTasks(ClusterService.java:624)
    at org.elasticsearch.cluster.service.ClusterService$ClusterServiceTaskBatcher.run(ClusterService.java:316)
    at org.elasticsearch.cluster.service.TaskBatcher.runIfNotProcessed(TaskBatcher.java:173)
    at org.elasticsearch.cluster.service.TaskBatchers$BatchedTask.run(TaskBatcher.java:211)
    at org.elasticsearch.common.util.concurrent.ThreadContext$ContextPreservingRunnable.run(ThreadContext.java:569)
    at org.elasticsearch.common.util.concurrent.PrioritizedEsThreadPoolExecutor$TieBreakingPrioritizedRunnable.runAndClear(ThreadContext.java:569)
    at org.elasticsearch.common.util.concurrent.PrioritizedEsThreadPoolExecutor$TieBreakingPrioritizedRunnable.run(PrioritizedEsThreadPoolExecutor.java:569)
    at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor.java:1142)
    at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor.java:617)
    at java.lang.Thread.run(Thread.java:748)
```

## 优化3: rebalance

创建和删除索引需要立即负载均衡吗?

方案: 禁止rebalance, 后台定时rebalance



## 优化4: Master节点

方案: 不接收业务请求, 专注于处理集群级别请求



## 优化5: 磁盘决策者

方案: 在分配分片或者rebalance的时候**不考虑RELOCATING和INITIALIZING的分片**



## 优化6：去除遍历所有分片

方案：在获取未分配的分片个数和最小时间时，从当前集群状态直接获取，不用调用该方法遍历所有分片  
`routingTable.shardsWithState(ShardRoutingState.UNASSIGNED)`



# 问题概览

## 问题概览



1、任务堆积太多影响集群稳定性，想取消怎么办？



2、怎么提高任务优先级，让任务优先或者延后执行？



3、到达内存阈值后熔断了元数据请求导致任务失败？



4、怎么更方便的部署集群，集群扩缩容？



5、业务怎么对集群设置权限？



6、大量导出导致CPU负载高，影响别的业务查询



7、模糊查询、Range查询负载高，可能导致节点down机



8、索引实时入库超时怎么解决？

# 问题1

任务堆积太多影响集群稳定性，想取消怎么办？

方案：对ES功能进行扩展，增加如下参数

```
PUT _cluster/settings
```

```
{  
  "transient": {  
    "cluster.service.update.tasks.cancel.source": ["shard-started", "shard-  
failed", "update-settings"]  
  }  
}
```

```
PUT _cluster/settings
```

```
{  
  "transient": {  
    "cluster.service.update.tasks.cancel.order": [96568,96569,96570,96571]  
  }  
}
```

## 问题2

怎么提高任务优先级，让任务优先或者延后执行？

**方案：**对ES功能进行扩展，增加如下参数

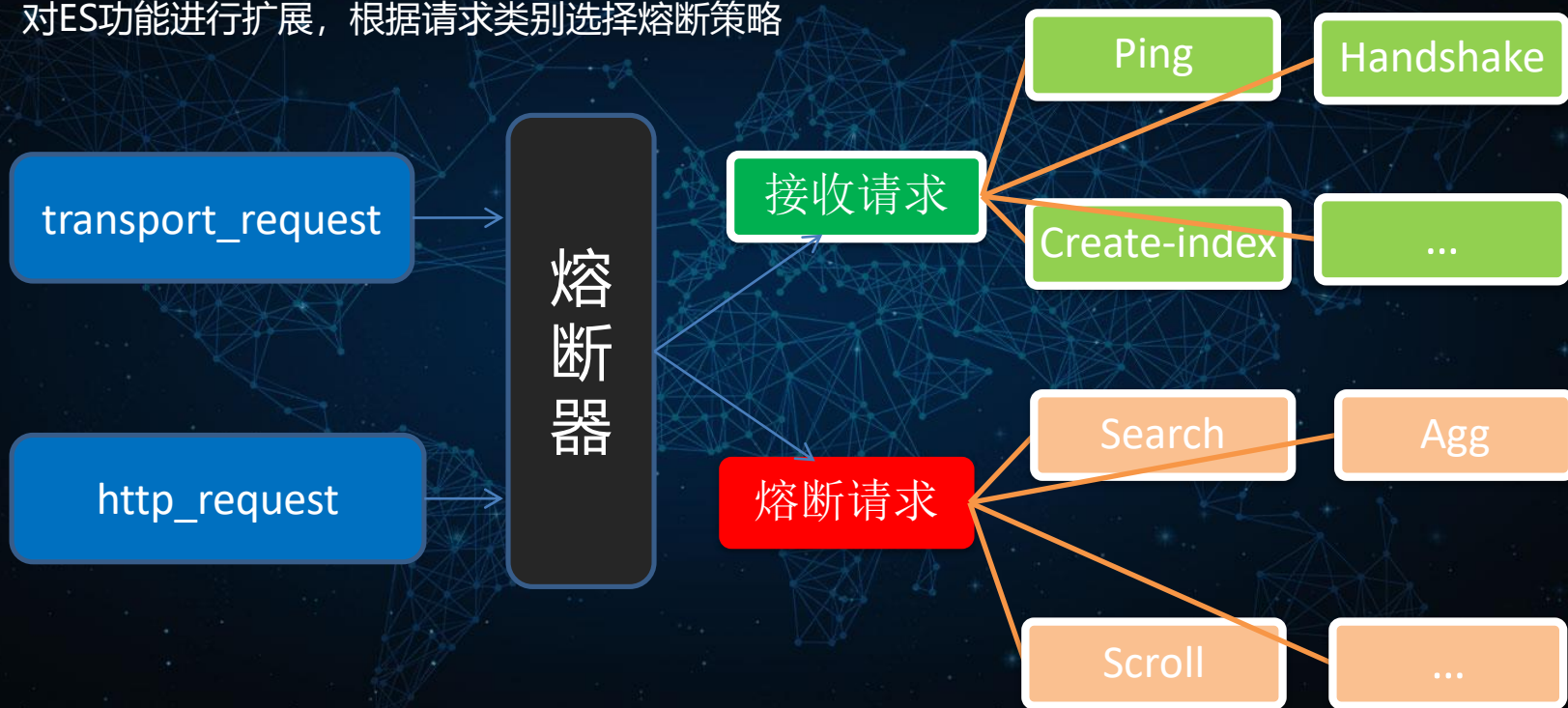
```
PUT _cluster/settings
{
  "transient": {
    "cluster.service.update.tasks. IMMEDIATE.source": ["shard-started","shard-
failed","update-settings"]
  }
}
```

```
PUT _cluster/settings
{
  "transient": {
    "cluster.service.update.tasks. HIGH.order": [96568,96569,96570,96571]
  }
}
```

# 问题3

到达内存阈值后熔断了元数据请求导致任务失败？

方案：对ES功能进行扩展，根据请求类别选择熔断策略



## 问题4

怎么更方便的部署集群，集群扩缩容？

方案：修改ES加载配置逻辑，读取自定义的配置文件



## 问题5

业务怎么对自己的集群设置权限？

方案：

对ES功能进行扩展，增加黑白名单（支持IPv4、IPv6和通配符）  
参数如下：

1. `http.filter.enabled`
2. `transport.filter.enabled`
3. `http.filter.allow`
4. `http.filter.deny`
5. `transport.filter.allow`
6. `transport.filter.deny`

## 问题6

大量导出导致CPU负载高，影响别的业务查询

方案：

对ES功能进行扩展，添加相应的控制参数，可以控制导出的并发量和数据量。

参数：

1. `scroll.enabled`
2. `scroll.interval`
3. `scroll.concurrent.indices`
4. `scroll.limit`

## 问题7

模糊查询、Range查询负载高，可能导致节点down机

### 方案：

- 1、模糊查询时如果字段 长度太长，直接熔断
- 2、修改Lucene源码，添加相应的控制参数，控制栈的大小最大为500，防止栈溢出

## 问题8

索引实时入库超时怎么解决？

方案：

Linux设置以下参数：

```
vm.dirty_background_ratio = 5
```

```
vm.dirty_ratio = 10
```

```
vm.dirty_writeback_centisecs = 500
```

```
vm.dirty_expire_centisecs = 3000
```

大索引（分片大小超过30G）后台手动merge：

```
PUT index_name/_settings
```

```
{
```

```
  "index.merge.policy.max_merged_segment": "500mb"
```

```
}
```

## 内容概要



平台介绍



平台实战



未来展望

# 未来展望



SQL



bulkload



proxy



多活

# THANKS

公众号



个人号





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