

部署

如果按照中做到了，希望已学到了一件关于 Elasticsearch 的事情并且准把的集群部署到生产环境。一章不是在生产中进行集群的尽指南，但是它涵盖了集群上之前需要考的事。

主要包括三个方面：

- 后勤方面的考，如硬件和部署策略的建
- 更合于生产环境的配置更改
- 部署后的考，例如安全，最大限度的索引性能和

硬件

按照正常的流程，可能已在自己的本或集群上使用了 Elasticsearch。但是当要部署 Elasticsearch 到生产环境，有一些建是需要考的。里没有什必要遵守的准，Elasticsearch 被用于在多的机器上理各任。基于我在生产环境使用 Elasticsearch 集群的，些建可以提供一个好的起点。

内存

如果有一源是最先被耗尽的，它可能是内存。排序和聚合都很耗内存，所以有足的堆空来付它是很重要的。即使堆空是比小的候，也能操作系文件存提供外的内存。因 Lucene 使用的多数据是基于磁的格式，Elasticsearch 利用操作系存能生很大效果。

64 GB 内存的机器是非常理想的，但是32 GB 和16 GB 机器也是很常的。少于8 GB 会得其反（最需要很多很多的小机器），大于64 GB 的机器也会有，我将在[堆内存:大小和交](#)中。

CPUs

大多数 Elasticsearch 部署往往 CPU 要求不高。因此，相其它源，具体配置多少个（CPU）不是那。具有多个内核的代埋器，常的集群使用到八个核的机器。

如果要在更快的 CPUs 和更多的核心之，更多的核心更好。多个内核提供的外并微快一点点的率。

硬

硬所有的集群都很重要，大量写入的集群更是加倍重要（例如那些存日志数据的）。硬是服器上最慢的子系，意味着那些写入量很大的集群很容易硬和，使得它成集群的瓶。

如果担得起 SSD，它将超出任何旋介（注：机械硬，磁等）。基于 SSD 的点，和索引性能都有提升。如果担得起，SSD 是一个好的。

的 I/O 调度程序

如果正在使用 SSDs, 保持的系 I/O 调度程序是配置正确的。当向硬盘写数据, I/O 调度程序决定何时把数据发送到硬盘。大多数 *nix 发行版下的调度程序都叫做 **cfq** (完全公平队列)。

调度程序分配 I/O 片到线程。并且优化一些到硬盘的多队列的。但它是旋转介质的: 机械硬盘的固有特性意味着它写入数据到基于物理布局的硬盘会更高效。

SSD 从来是低效的, 尽管里面没有涉及到机械硬盘。但是, **deadline** 或者 **noop** 被使用。**deadline** 调度程序基于写入等待行优化, **noop** 只是一个的 FIFO 队列。

一个小的更改可以带来显著的影响。是使用正确的调度程序, 我看到了500倍的写入能力提升。

如果使用旋转介质, 选取尽可能快的硬盘 (高性能服务器硬盘, 15k RPM 硬盘)。

使用 RAID 0 是提高硬盘速度的有效途径, 机械硬盘和 SSD 从来都是如此。没有必要使用像或其它 RAID 体系, 因为高可用已通常 replicas 内建于 Elasticsearch 之中。

最后, 避免使用附加存储 (NAS)。人常声称他的 NAS 解决方案比本地硬盘更快更可靠。除却一些声称, 我从没看到 NAS 能配得上它的大肆宣传。NAS 常常很慢, 露出更大的延迟和更大的平均延迟方差, 而且它是有点故障的。

快速可靠的网络虽然分布式系统的性能是很重要的。低延迟能帮助保证节点能容易的通信, 大带宽能帮助分片转移和恢复。现代数据中心 (1 GbE, 10 GbE) 大多数集群都是足够的。

即使数据中心近在咫尺, 也要避免集群跨越多个数据中心。要避免集群跨越大的地理距离。

Elasticsearch 假定所有节点都是平等的——并不会因为有一半的节点在150ms 以外的单一数据中心而有所不同。更大的延迟会加重分布式系统中的问题而且使得索引和排序更困难。

和 NAS 的争论类似, 个人都声称他的数据中心的路都是健壮和低延迟的。是真的——直到它不是 (故障终究是会发生的, 可以相信它)。从我的角度来看, 管理跨数据中心集群的麻烦是根本不得了的。

取真正的高配机器在今天是不可能的: 成百 GB 的 RAM 和几十个 CPU 核心。反之, 在云平台上串起成千的小虚拟机也是可能的, 例如 EC2。哪种方式是最好的?

通常, 中配或者高配机器更好。避免使用低配机器, 因为不会希望去管理有上千个节点的集群, 而且在一些低配机器上运行 Elasticsearch 也是困难的。

与此同时, 避免使用真正的高配机器。它通常会致资源使用不均衡 (例如, 所有的内存都被使用, 但 CPU 却没有) 而且在机器上运行多个节点, 会增加复杂度。

Java 虚拟机

始终运行最新版本的 Java 虚拟机 (JVM)，除非 Elasticsearch 站上有说明。Elasticsearch，特别是 Lucene，是一个高要求的组件。Lucene 的组件和集成经常暴露出 JVM 本身的 bug。这些 bug 的范围从微小的麻烦到严重故障，所以，最好尽可能的使用最新版本的 JVM。

Java 8 首先发布于 Java 7。不再支持 Java 6。Oracle 或者 OpenJDK 是可以接受的，它在性能和稳定性也差不多。

如果你的应用程序是用 Java 编写并正在使用客户端（注：Transport Client，下同）或节点客户端（注：Node Client，下同），保证应用程序的 JVM 和服务器的 JVM 是完全一致的。在 Elasticsearch 的几个地方，使用 Java 的本地序列化（IP 地址、常量等等）。不幸的是，Oracle 的 JVM 在几个小版本之间有修改序列化格式，从而导致奇怪的行为。这种情况很少，但最佳实践是客户端和服务端使用相同版本 JVM。

不要调整 JVM 配置

JVM 暴露出几十个（甚至数百）的配置、参数和配置。它允许进行微调，但 JVM 几乎是一个黑盒。当遇到一个性能问题，要怀疑它是人的本性。我们要求控制人的本性，而不要去调整 JVM 参数。Elasticsearch 是一个组件，并且我根据多年的使用情况调整了当前 JVM 配置。它很容易开始调整，并产生难以衡量的、未知的阴影，并最终使集群进入一个缓慢的、不稳定的混乱的效果。当调整集群时，第一往往是去除所有的自定义配置。多数情况下，这就可以恢复稳定性和性能。

Transport Client 与 Node Client

如果你使用的是 Java，可能想知道如何使用客户端（注：Transport Client，下同）与节点客户端（注：Node Client，下同）。在本文所述，客户端作为一个集群和应用程序之间的通信。它知道 API 并能自己在节点之间，嗅探集群等等。但它是集群外部的，和 REST 客户端类似。

一方面，节点客户端，它是一个集群中的节点（但不保存数据，不能成为主节点）。因为它是一个节点，它知道整个集群状态（所有节点保留，分片分布在哪些节点，等等）。这意味着它可以进行 API 调用但少了一个主节点。

这里有一个客户端案例的使用情况：

- 如果要将应用程序和 Elasticsearch 集群进行解耦，客户端是一个理想的选择。例如，如果你的应用程序需要快速的建立和连接到集群的连接，客户端比“节点客户端”更合适，因为它不是一个集群的一部分。

类似地，如果需要建成成千上万的连接，不想有成千上万的节点加入集群。客户端（Transport Client）将是一个更好的选择。

- 一方面，如果只需要有少数的、长期持久的连接接到集群，节点客户端可以更高效，因为它知道集群的布局。但是它会使得应用程序和集群耦合在一起，所以从防火的角度，它可能会成为问题。

配置管理

如果 已 使用配置管理（Puppet, Chef, Ansible）， 可以跳 此提示。

如果 没有使用配置管理工具，那 注意了！通 `parallel-ssh` 管理少量服 器
在可能正常工作，但伴随着集群的 它将成 一 梦。 在不犯 的情况下手 30
个配置文件几乎是不可能的。

配置管理工具通 自 化更改配置的 程保持集群的一致性。 可能需要一点 来建立和学 ，但它本身
，随着 的推移会有 厚的回 。

重要配置的修改

Elasticsearch 已 有了 很好 的 ，特 是 及到性能相 的配置或者 。 如果 有疑
，最好就不要 它。我 已 目 了数十个因 的 置而 致 的集群， 因 它的管理者
改 一个配置或者 就可以 来 100 倍的提升。

NOTE

整 文章，所有的配置 都同等重要，和描述 序无 ， 所有的配置 ，并
用到 的集群中。

其它数据 可能需要 ，但 得来 ，Elasticsearch 不需要。 如果 遇到了性能
，解决方法通常是更好的数据布局或者更多的 点。 在 Elasticsearch 中很少有“神奇的配置 ”，
如果存在，我 也已 化了！

外，有些 上的 配置在生 境中是 整的。 些 整可能会 的工作更加 松，又或者因
没 法 定一个 （它取决于 的集群布局）。

指定名字

Elasticsearch 的集群名字叫 `elasticsearch` 。 最好 的生
境的集群改个名字，改名字的目的很 ， 就是防止某人的 本 加入了集群 意外。
修改成 `elasticsearch_production` 会很省心。

可以在 的 `elasticsearch.yml` 文件中修改：

```
cluster.name: elasticsearch_production
```

同 ，最好也修改 的点名字。就像 在可能 的那 ， Elasticsearch 会在 的点 的
候随机 它指定一个名字。 可能会 得 很有趣，但是当凌晨 3 点 的 候， 在 回
台物理机是 Tagak the Leopard Lord 的 候， 就不 得有趣了。

更重要的是， 些名字是在 的 候 生的， 次 点， 它都会得到一个新的名字。 会使日志
得很混乱，因 所有 点的名称都是不断 化的。

可能会 得 ，我 建 个 点 置一个有意 的、清楚的、描述性的名字，同 可以在
`elasticsearch.yml` 中配置：

```
node.name: elasticsearch_005_data
```

路径

情况下，Elasticsearch 会把 件、日志以及 最重要的数据放在安装目 下。 会 来不幸的事故，如果 重新安装 Elasticsearch 的 候不小心把安装目 覆 了。如果 不小心， 就可能把 的全部数据 掉了。

不要笑， 情况，我 很多次了。

最好的 就是把 的数据目 配置到安装目 以外的地方，同 也可以 移 的 件和日志目 。

可以更改如下：

```
path.data: /path/to/data1,/path/to/data2 ①

# Path to log files:
path.logs: /path/to/logs

# Path to where plugins are installed:
path.plugins: /path/to/plugins
```

① 注意： 可以通 逗号分隔指定多个目 。

数据可以保存到多个不同的目 ， 如果将 个目 分 挂 不同的硬 ， 可是一个 且高效 一个 磁 列（ RAID 0 ）的 法。Elasticsearch 会自 把条 化（注：RAID 0 又称 Stripe（条 化），在磁 列中，数据是以条 的方式 穿在磁 列所有硬 中的） 数据分隔到不同的目 ， 以便提高性能。

WARNING

多个数据路径的安全性和性能

如同任何磁 列（ RAID 0 ）的配置，只有 一的数据拷 保存到硬 器。如果 失去了一个硬 器， 肯定 会失去 计算机上的一部分数据。 气好的 的副本在集群的其他地方，可以用来恢 数据和最近的 。

Elasticsearch 将全部的条 化分片放到 个 器来保 最小程度的数据 失。 意味着 分片 0 将完全被放置在 个 器上。 Elasticsearch 没有一个条 化的分片跨越在多个 器，因 一个 器的 失会破坏整个分片。

性能 生的影 是：如果 添加多个 器来提高一个 独索引的性能，可能 助不大，因 大多数 点只有一个分片和 一个 的 器。多个数据路径只是 助如果 有 多索引/分片在 个 点上。

多个数据路径是一个非常方便的功能，但到 来，Elasticsearch 并不是 磁 列（ software RAID ）的 件。如果 需要更高 的、 健的、 活的配置， 我 建 使用 磁 列（ software RAID ）的 件，而不是多个数据路径的功能。

最小主 点数

`minimum_master_nodes` 定 的集群的 定 其 重要。 当 的集群中有 个 masters（注：主 点）的 候， 个配置有助于防止 裂，一 个主 点同 存在于一个集群的 象。

如果 的集群 生了 裂，那 的集群就会 在 失数据的危 中，因 主 点被 是 个集群的最高

治者，它决定了什么时候新的索引可以建，分片是如何移动的等等。如果有一个 masters 点，数据的完整性将得不到保证，因此有一个点他拥有集群的控制。

一个配置就是告诉 Elasticsearch 当没有足够 master 候选点的时候，就不要进行 master 点选举，等待 master 候选点足够了才行。

此配置一开始被配置 master 候选点的法定个数（大多数个）。法定个数就是 $(\text{master 候选点个数} / 2) + 1$ 。这里有几个例子：

- 如果有 10 个点（能保存数据，同时能成 master），法定数就是 6。
- 如果有 3 个候选 master 点，和 100 个 data 点，法定数就是 2，只要数数那些可以做 master 的点数就可以了。
- 如果有一个点，遇到了故障。法定数当然是 2，但是意味着如果有一个点挂掉，整个集群就不可用了。配置成 1 可以保证集群的功能，但是就无法保证集群不会裂了，像这种情况，最好至少保证有 3 个点。

可以在 `elasticsearch.yml` 文件中配置：

```
discovery.zen.minimum_master_nodes: 2
```

但是由于 Elasticsearch 是分布式的，可以很容易的添加和删除点，但是会改变法定个数。不得不修改一个索引点的配置并且重启整个集群只是配置生效，将是非常痛苦的一件事情。

基于这个原因，`minimum_master_nodes`（有一些其它配置）允许通过 API 用的方式进行配置。当集群在运行的时候，可以修改配置：

```
PUT /_cluster/settings
{
  "persistent": {
    "discovery.zen.minimum_master_nodes": 2
  }
}
```

将变成一个永久的配置，并且无论配置里配置的如何，都会生效。当添加和删除 master 点的时候，需要更改配置。

集群恢复方面的配置

当集群重启，几个配置影响的分片恢复的表。首先，我需要明白如果什么也没配置将会发生什么。

想象一下假设有 10 个点，每个点只保存一个分片，每个分片是一个主分片或者是一个副本分片，或者有一个有 5 个主分片 / 1 个副本分片的索引。有需要整个集群做（比如，安装了新的程序），当重启的集群，恰巧出了 5 个点已挂掉，有 5 个没挂掉的景象。

假其它 5 个点出故障，或者他们根本没有收到立即重启的命令。不管什么原因，有 5 个点在

上，五个点会相互通信，选出一个 master，从而形成一个集群。他注意到数据不再均匀分布，因有5个点在集群中挂了，所以他才会立即分片制。

最后，的其它5个点打加入了集群。些点会它的数据正在被制到其他点，所以他除本地数据（因数据要是多余的，要是的）。然后整个集群重新行平衡，因集群的大小已从5成了10。

在整个程中，的点会消耗磁和，来回移数据，因没有更好的法。于有TB数据的大集群，无用的数据需要很。如果等待所有的点重好了，整个集群再上，所有的本地的数据都不需要移。

在我知道的所在了，我可以修改一些置来解它。首先我要Elasticsearch一个格的限制：

```
gateway.recover_after_nodes: 8
```

将阻止 Elasticsearch 在存在至少8个点（数据点或者 master 点）之前行数据恢。个的定取决于个人喜好：整个集群提供服之前希望有多少个点在？情况下，我置8，意味着至少要有8个点，集群才可用。

在我要告 Elasticsearch 集群中有多少个点，以及我意些点等待多：

```
gateway.expected_nodes: 10
gateway.recover_after_time: 5m
```

意味着 Elasticsearch 会采取如下操作：

- 等待集群至少存在8个点
- 等待5分，或者10个点上后，才行数据恢，取决于个条件先到。

三个置可以在集群重的候避免多的分片交。可能会数据恢从数个小短几秒。

注意：些配置只能置在 `config/elasticsearch.yml` 文件中或者是在命令行里（它不能更新）它只在整个集群重的候有性作用。

最好使用播代替播

Elasticsearch 被配置使用播，以防止点无意中加入集群。只有在同一台机器上行的点才会自成集群。

然播然作件提供，但它永不被使用在生境了，否在得到的果就是一个点意外的加入到了的生境，是因他收到了一个的播信号。于播本身并没有，播会致一些愚蠢的，并且致集群的脆弱（比如，一个工程正在鼓，而没有告，会所有的点突然不了方了）。

使用播，可以Elasticsearch 提供一些它去接的点列表。当一个点系到播列表中的成，它就会得到整个集群所有点的状，然后它会系 master 点，并加入集群。

意味着 的 播列表不需要包含 的集群中的所有 点, 它只是需要足 的 点, 当一个新 点 系上其中一个并且 上 就可以了。如果 使用 master 候 点作 播列表, 只要列出三个就可以了。 个配置在 `elasticsearch.yml` 文件中:

```
discovery.zen.ping.unicast.hosts: ["host1", "host2:port"]
```

于 Elasticsearch 点 的 信息, 参 [Zen Discovery](#) Elasticsearch 文献。

不要触 些配置!

在 Elasticsearch 中有一些 点, 人 可能不可避免的会 到。 我 理解的, 所有的 整就是 了 化, 但是 些 整, 真的不需要理会它。因 它 常会被乱用, 从而造成系 的不 定或者糟 的性 能, 甚至 者都有可能。

回收器

里已 要介 了 [\[garbage_collector_primer\]](#), JVM 使用一个 回收器来 放不再使用的内存。 篇内容的 是上一篇的一个延 , 但是因 重要, 所以 得 独拿出来作 一 。

不要更改 的 回收器!

Elasticsearch 的 回收器 (GC) 是 CMS。 个 回收器可以和 用并行 理, 以便它可以最小化停 。然而, 它有 个 stop-the-world 段, 理大内存也有点吃力。

尽管有 些 点, 它 是目前 于像 Elasticsearch 低延 需求 件的最佳 回收器。官方建 使用 CMS。

在有一款新的 回收器, 叫 G1 回收器 (G1GC)。 款新的 GC 被 , 旨在比 CMS 更小的 停 , 以及 大内存的 理能力。 它的原理是把内存分成 多区域, 并且 些区域最有可能需要回收内存。通 先收集 些区域 (*garbage first*), 生更小的 停 , 从而能 更大的内存。

听起来很棒! 憾的是, G1GC 是太新了, 常 新的 bugs。 些 通常是段 (*segfault*) 型, 便造成硬 的崩 。 Lucene 的 套件 回收算法要求 格, 看起来 些 陷 G1GC 并没有很好地解决。

我 很希望在将来某一天推 使用 G1GC, 但是 于 在, 它 不能足 定的 足 Elasticsearch 和 Lucene 的要求。

程池

多人 喜 整 程池。 无 什 原因, 人 都 加 程数无法抵抗。索引太多了? 加 程! 搜索太多了? 加 程! 点空 率低于 95%? 加 程!

Elasticsearch 的 程 置已 是很合理的了。 于所有的 程池 (除了 搜索), 程个数是根据 CPU 核心数 置的。 如果 有 8 个核, 可以同 行的只有 8 个 程, 只分配 8 个 程 任何特定的 程池是有道理的。

搜索 程池 置的大一点, 配置 `int ((核心数 * 3) / 2) + 1`。

可能会 某些 程可能会阻塞（如磁 上的 I/O 操作），所以 才想加大 程的。 于 Elasticsearch 来 并不是一个 ：因 大多数 I/O 的操作是由 Lucene 程管理的，而不是 Elasticsearch。

此外， 程池通 彼此之 的工作配合。 不必再因 它正在等待磁 写操作而担心 程阻塞，因 程早已把 个工作交 外的 程池，并且 行了 。

最后， 的 理器的 算能力是有限的， 有更多的 程会 致 的 理器 繁切 程上下文。 一个 理器同 只能 行一个 程。所以当它需要切 到其它不同的 程的 候，它会存 当前的状 （寄存器 等等），然后加 外一个 程。 如果幸 的 ， 个切 生在同一个核心，如果不幸的 ， 个切 可能 生在不同的核心， 就需要在内核 上 行 。

个上下文的切 ， 会 CPU 周期 来管理 度的 ；在 代的 CPUs 上， 估 高 30 μ s。也就是 程会被堵塞超 30 μ s，如果 个 用于 程的 行， 有可能早就 束了。

人 常稀里糊 的 置 程池的 。8 个核的 CPU，我 遇到 有人配了 60、100 甚至 1000 个 程。 些 置只会 CPU 工作效率更低。

所以，下次 不要 整 程池的 程数。如果 真 想 整 ， 一定要 注 的 CPU 核心数，最多 置成核心数的 倍，再多了都是浪 。

堆内存:大小和交

Elasticsearch 安装后 置的堆内存是 1 GB。 于任何一个 部署来 ， 个 置都太小了。如果 正在使用 些 堆内存配置， 的集群可能会出 。

里有 方式修改 Elasticsearch 的堆内存。最 的一个方法就是指定 `ES_HEAP_SIZE` 境 量。服 程在 候会 取 个 量，并相 的 置堆的大小。比如， 可以用下面的命令 置它：

```
export ES_HEAP_SIZE=10g
```

此外， 也可以通 命令行参数的形式，在程序 的 候把内存大小 它，如果 得 更 的：

```
./bin/elasticsearch -Xmx10g -Xms10g ①
```

① 保堆内存最小 （ `Xms` ）与最大 （ `Xmx` ）的大小是相同的，防止程序在 行 改 堆内存大小， 是一个很耗系 源的 程。

通常来 ， 置 `ES_HEAP_SIZE` 境 量，比直接写 `-Xmx -Xms` 更好一点。

把 的内存的（少于）一半 **Lucene**

一个常 的 是 Elasticsearch 分配的内存 太大了。假 有一个 64 GB 内存的机器，天 ，我要把 64 GB 内存全都 Elasticsearch。因 越多越好 ！

当然，内存 于 Elasticsearch 来 是重要的，它可以被 多内存数据 使用来提供更快的操作。但是 到 里， 有 外一个内存消耗大 非堆内存 （off-heap）：Lucene。

Lucene 被 可以利用操作系 底 机制来 存内存数据 。 Lucene 的段是分 存 到 个文件中的。因 段是不可 的， 些文件也都不会 化， 是 存友好的，同 操作系 也会把 些 段文件 存起来，以便更快的 。

Lucene 的性能取决于和操作系 的相互作用。如果 把所有的内存都分配 Elasticsearch 的堆内存，那将不会有剩余的内存交 Lucene。 将 重地影 全文 索的性能。

准的建 是把 50% 的可用内存作 Elasticsearch 的堆内存，保留剩下的 50%。当然它也不会被浪 ， Lucene 会很 意利用起余下的内存。

如果 不需要 分 字符串做聚合 算（例如，不需要 `fielddata` ）可以考 降低堆内存。堆内存越小，Elasticsearch（更快的 GC）和 Lucene（更多的内存用于 存）的性能越好。

不要超 32 GB！

里有 外一个原因不分配大内存 Elasticsearch。事 上， JVM 在内存小于 32 GB 的 候会采用一个内存 象指 技 。

在 Java 中，所有的 象都分配在堆上，并通 一个指 行引用。 普通 象指 （OOP）指向 些 象，通常 CPU 字 的大小：32 位或 64 位，取决于 的 理器。指 引用的就是 个 OOP 的字 位置。

于 32 位的系 ，意味着堆内存大小最大 4 GB。 于 64 位的系 ， 可以使用更大的内存，但是 64 位的指 意味着更大的浪 ，因 的指 本身大了。更糟 的是， 更大的指 在主内存和各 存（例如 LLC，L1 等）之 移 数据的 候，会占用更多的 。

Java 使用一个叫作 内存指 （compressed oops）的技 来解决 个 。 它的指 不再表示 象在内存中的精 位置，而是表示 偏移量 。 意味着 32 位的指 可以引用 40 个 象， 而不是 40 个字 。最 ，也就是 堆内存 到 32 GB 的物理内存，也可以用 32 位的指 表示。

一旦 越 那个神奇的 ~32 GB 的 界，指 就会切回普通 象的指 。 个 象的指 都 了，就会使用更多的 CPU 内存 ，也就是 上失去了更多的内存。事 上，当内存到 40–50 GB 的 候，有效内存才相当于使用内存 象指 技 候的 32 GB 内存。

段描述的意思就是 ：即便 有足 的内存，也尽量不要 超 32 GB。因 它浪 了内存，降低了 CPU 的性能， 要 GC 大内存。

到底需要低于 32 GB多少，来 置我的 JVM？

憾的是， 需要看情况。 切的 分要根据 JVMs 和操作系 而定。 如果 想保 其安全可 ， 置堆内存 31 GB 是一个安全的 。 外， 可以在 的 JVM 置里添加 `-XX:+PrintFlagsFinal` 用来 JVM 的 界 ， 并且 UseCompressedOops 的 是否 true。 于 自己使用的 JVM 和操作系 ， 将 到最合 的堆内存 界 。

例如，我 在一台安装 Java 1.7 的 MacOSX 上 ， 可以看到指 在被禁用之前，最大堆内存大 是在 32600 mb (~31.83 gb)：

```
$ JAVA_HOME='/usr/libexec/java_home -v 1.7' java -Xmx32600m -XX:+PrintFlagsFinal 2>
/dev/null | grep UseCompressedOops
    bool UseCompressedOops    := true
$ JAVA_HOME='/usr/libexec/java_home -v 1.7' java -Xmx32766m -XX:+PrintFlagsFinal 2>
/dev/null | grep UseCompressedOops
    bool UseCompressedOops    = false
```

相比之下，同一台机器安装 Java 1.8，可以看到指 在禁用之前，最大堆内存大 是在 32766 mb (~31.99 gb)：

```
$ JAVA_HOME='/usr/libexec/java_home -v 1.8' java -Xmx32766m -XX:+PrintFlagsFinal 2>
/dev/null | grep UseCompressedOops
    bool UseCompressedOops    := true
$ JAVA_HOME='/usr/libexec/java_home -v 1.8' java -Xmx32767m -XX:+PrintFlagsFinal 2>
/dev/null | grep UseCompressedOops
    bool UseCompressedOops    = false
```

个例子告 我，影 内存指 使用的 界， 是会根据 JVM 的不同而 化的。所以从其他地方 取的例子，需要 慎使用，要 操作系 配置和 JVM。

如果使用的是 Elasticsearch v2.2.0， 日志其 会告 JVM 是否正在使用内存指 。 会看到像 的日志消息：

```
[2015-12-16 13:53:33,417][INFO ][env] [Illyana Rasputin] heap size [989.8mb],
compressed ordinary object pointers [true]
```

表明内存指 正在被使用。如果没有，日志消息会 示 `[false]`。

我有一个 1 TB 内存的机器！

一个 32 GB 的分割是很重要的。那如果 的机器有很大的内存 ？ 一台有着 512GB 内存的服务器 常。

首先，我建议避免使用 的高配机器（参考[硬件](#)）。

但是如果已经有了 的机器，有三个可：

- 主要做全文索引？考虑 Elasticsearch 4 - 32 GB 的内存，Lucene 通过操作系统文件存来利用余下的内存。那些内存都会用来存 segments，来索引全文索引。
- 需要更多的排序和聚合？而且大部分的聚合 算是在数字、日期、地理点和 非分字符串上？很幸， 的聚合 算将在内存友好的 doc values 上完成！ Elasticsearch 4 到 32 GB 的内存，其余部分 操作系统 存内存中的 doc values。
- 在分字符串做大量的排序和聚合（例如， 或者 SigTerms，等等）不幸的是，意味着需要 fielddata，意味着需要堆空。考在 个机器上行 个或多个点，而不是有大量 RAM 的一个点。然要持 50% 原。

假 有个机器有 128 GB 的内存，可以建 个点， 个点内存分配不超 32 GB。也就是不超 64 GB 内存 ES 的堆内存，剩下的超 64 GB 的内存 Lucene。

如果 一，需要配置 `cluster.routing.allocation.same_shard.host: true`。会防止同一个分片（shard）的主副本存在同一个物理机上（因 如果存在一个机器上，副本的高可用性就没有了）。

Swapping 是性能的 墓

是而易的，但是 是有必要 的更清楚一点：内存交 到磁 服务器性能来 是 致命的。想想看：一个内存操作必 能被快速行。

如果内存交 到磁 上，一个 100 微秒的操作可能 成 10 秒。再想想那多 10 微秒的操作 延累加起来。不 看出 swapping 于性能是多 可怕。

最好的 法就是在 的操作系统 中完全禁用 swap。 可以 禁用：

```
sudo swapoff -a
```

如果需要永久禁用， 可能需要修改 `/etc/fstab` 文件， 要参考 的操作系统 相 文。

如果 并不打算完全禁用 swap，也可以 降低 `swappiness` 的。 个 决定操作系统 交 内存的率。 可以 防正常情况下 生交，但 允 操作系统 在 急情况下 生交。

于大部分Linux操作系统，可以在 `sysctl` 中 配置：

```
vm.swappiness = 1 ①
```

① `swappiness` 置 1 比 置 0 要好，因 在一些内核版本 `swappiness` 置 0 会触 系 OOM-

killer（注：Linux 内核的 Out of Memory（OOM）killer 机制）。

最后，如果上面的方法都不合用，需要打配置文件中的 `mlockall`。它的作用就是允许 JVM 住内存，禁止操作系统交换出去。在 `elasticsearch.yml` 文件中，配置如下：

```
bootstrap.mlockall: true
```

文件描述符和 MMap

Lucene 使用了大量的文件。同样，Elasticsearch 在点和 HTTP 客户端之间通信也使用了大量的套接字（注：sockets）。所有这一切都需要足够的文件描述符。

可悲的是，多代的 Linux 发行版本，一个进程允许一个微不足道的 1024 文件描述符。一个小的 Elasticsearch 节点实在是太低了，更不用说一个数以百索引的节点。

增加的文件描述符，设置一个很大的值，如 64,000。这个进程困得人火，它高度依赖于特定操作系统和分布。参考操作系统文档来定如何最好地修改允许的文件描述符数量。

一旦已经修改了它，Elasticsearch，以确保它的真的起作用并且有足够的文件描述符：

```
GET /_nodes/process

{
  "cluster_name": "elasticsearch__zach",
  "nodes": {
    "TGn9i02_QQKb0kavcLbnDw": {
      "name": "Zach",
      "transport_address": "inet[/192.168.1.131:9300]",
      "host": "zacharys-air",
      "ip": "192.168.1.131",
      "version": "2.0.0-SNAPSHOT",
      "build": "612f461",
      "http_address": "inet[/192.168.1.131:9200]",
      "process": {
        "refresh_interval_in_millis": 1000,
        "id": 19808,
        "max_file_descriptors": 64000, ①
        "mlockall": true
      }
    }
  }
}
```

① `max_file_descriptors` 字段 示 Elasticsearch 进程可以使用的可用文件描述符数量。

Elasticsearch 各文件混合使用了 NioFs（注：非阻塞文件系统）和 MMapFs（注：内存映射文件系统）。确保配置的最大映射数量，以便有足够的虚拟内存可用于 mmaped 文件。可以设置：

```
sysctl -w vm.max_map_count=262144
```

或者 可以在 `/etc/sysctl.conf` 通 修改 `vm.max_map_count` 永久 置它。

在生 之前，重温 个列表

在 入生 之前， 可能 了本 。本章中 及的 非常好，一般是可以知道的，但 是，正 部署到生 境之前需要重温 个列表。

一些 会 地阻止（如：可用的文件描述符太少）。因 他 很快 出来， 些都是容易 的。其他的一些 ，如 裂和内存 置，只有在糟 的事情 生之后才可 。在 一点上，解决 法往往是凌乱和繁 的。

在 生 之前 ，通 当配置集群来主 阻止 些情况 生，是更好的 。所以如果 想要从整本 的一个部分折角（或保存 ），本章将是一个很好的 。在部署到生 境的前一周， 地 里 出的列表，并 所有的建 。