

集群内的原理

充章

如前文所述，是充章中第一篇介绍Elasticsearch在分布式环境中的运行原理。在个章节中，我将会介绍cluster、node、shard等常用概念，Elasticsearch的容灾机制，以及如何理硬件故障的内容。

然个章节不是必须的——完全可以在不关注分片、副本和失效切换等内容的前提下使用Elasticsearch——但是将有助于了解Elasticsearch的内部工作流程。可以先快速浏览本章，将来有需要再次查看。

ElasticSearch的主旨是随需可用和按需扩容。而扩容可以通过性能更大（垂直扩容，或向扩容）或者数量更多的服务器（水平扩容，或横向扩容）来实现。

然Elasticsearch可以受益于更大的硬件，但是垂直扩容是有限的。真正的扩容能力是来自于水平扩容——集群添加更多的节点，并且将数据和索引力和索引分散到这些节点中。

于大多数的数据而言，通常需要对应用程序进行非常大的改造，才能利用上横向扩容的新来源。与之相反的是，ElasticSearch天生就是分布式的，它知道如何通过管理多个节点来提高容灾性和可用性。这意味着的应用无需关注节点。

本章将描述如何按需配置集群、节点和分片，并在硬件故障时保证数据安全。

空集群

如果我创建了一个空的节点，里面不包含任何的数据和索引，那么我的集群看起来就是一个包含空内容节点的集群。

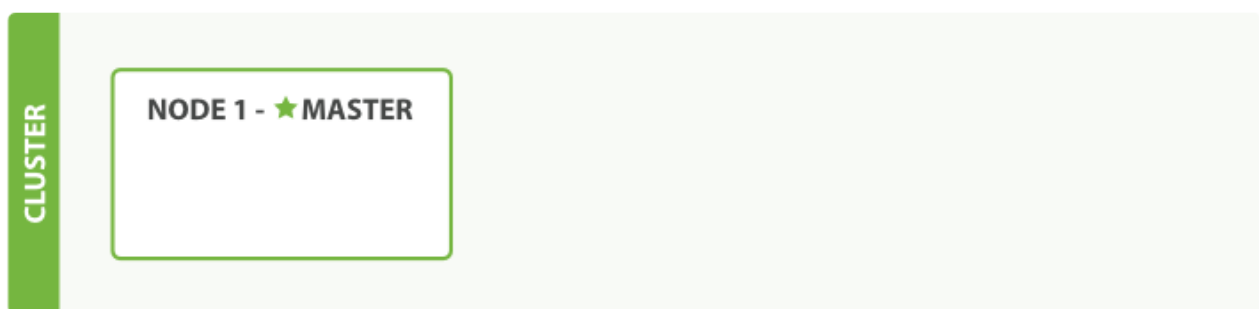


Figure 1. 包含空内容节点的集群

一个运行中的Elasticsearch实例称一个节点，而集群是由一个或者多个具有相同cluster.name配置的节点组成，它们共同承担数据和索引的负载。当有节点加入集群中或者从集群中移除节点，集群将会重新平均分布所有的数据。

当一个节点被选为主节点，它将管理集群内的所有更改，例如添加、删除索引，或者添加、删除节点等。而主节点并不需要涉及到文档的更改和搜索等操作，所以当集群只有一个主节点的情况下，即使流量的增加它也不会成为瓶颈。任何节点都可以成为主节点。我

的示例集群就只有一个节点，所以它同时也成了主节点。

作用，我可以将请求送到集群中的任何节点，包括主节点。每个节点都知道任意文档的位置，并且能将我的请求直接到存储我所需文档的节点。无论我将请求送到哪个节点，它都能从各个包含我所需文档的节点收集回数据，并将结果返回给客户端。Elasticsearch 一切的管理都是透明的。

集群健康

Elasticsearch 的集群控制信息中包含了多的数据，其中最重要的一就是集群健康，它在 `status` 字段中展示 `green`、`yellow` 或者 `red`。

```
GET /_cluster/health
```

在一个不包含任何索引的空集群中，它将会有一个类似于如下所示的返回内容：

```
{
  "cluster_name":      "elasticsearch",
  "status":            "green", ①
  "timed_out":         false,
  "number_of_nodes":   1,
  "number_of_data_nodes": 1,
  "active_primary_shards": 0,
  "active_shards":      0,
  "relocating_shards":  0,
  "initializing_shards": 0,
  "unassigned_shards":  0
}
```

① `status` 字段是我感兴趣的。

`status` 字段指示着当前集群在整体上是否工作正常。它的三色含义如下：

`green`

所有的主分片和副本分片都正常运行。

`yellow`

所有的主分片都正常运行，但不是所有的副本分片都正常运行。

`red`

有主分片没能正常运行。

在本章剩余的部分，我将解释什么是主分片和副本分片，以及上面提到的一些色的意义。

添加索引

我往 Elasticsearch 添加数据需要用到索引——保存相关数据的地方。索引上是指向一个或者多个物理分片的命名空间。

一个分片是一个底的工作元，它保存了全部数据中的一部分。在分片内部机制中，我将介绍分片是如何工作的，而在我只需知道一个分片是一个 Lucene 的例，以及它本身就是一个完整的搜索引擎。我的文档被存和索引到分片内，但是用程序是直接和索引而不是与分片行交互。

Elasticsearch 是利用分片将数据分到集群内各的。分片是数据的容器，文档保存在分片内，分片又被分配到集群内的各个节点里。当的集群模大或者小，Elasticsearch 会自的在各节点中移分片，使得数据然均分布在集群里。

一个分片可以是主分片或者副本分片。索引内任意一个文档都属于一个主分片，所以主分片的数目决定着索引能保存的最大数据量。

NOTE 技术上讲，一个主分片最大能存 Integer.MAX_VALUE - 128 个文档，但是最大需要参考的使用场景：包括使用的硬件，文档的大小和程度，索引和文档的方式以及期望的。

一个副本分片只是一个主分片的拷贝。副本分片作硬件故障保数据不丢失的冗余，并搜索和返回文档等操作提供服务。

在索引建立的时候就已经定了主分片数，但是副本分片数可以随修改。

我在包含一个空节点的集群内建名为 **blogs** 的索引。索引在一般情况下会被分配5个主分片，但是为了演示目的，我将分配3个主分片和一个副本（一个主分片有一个副本分片）：

```
PUT /blogs
{
  "settings" : {
    "number_of_shards" : 3,
    "number_of_replicas" : 1
  }
}
```

我的集群是在有一个索引的节点集群。所有3个主分片都被分配在 Node 1。

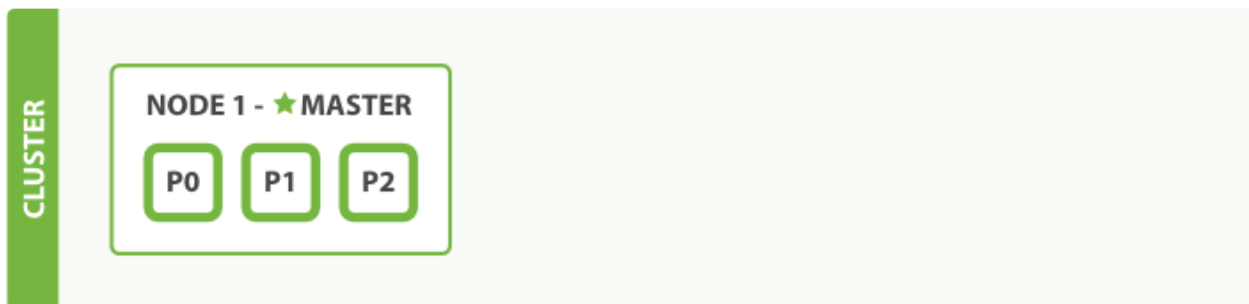


Figure 2. 有一个索引的节点集群

如果我正在看集群健康，我将看到如下内容：

```
{
  "cluster_name": "elasticsearch",
  "status": "yellow", ①
  "timed_out": false,
  "number_of_nodes": 1,
  "number_of_data_nodes": 1,
  "active_primary_shards": 3,
  "active_shards": 3,
  "relocating_shards": 0,
  "initializing_shards": 0,
  "unassigned_shards": 3, ②
  "delayed_unassigned_shards": 0,
  "number_of_pending_tasks": 0,
  "number_of_in_flight_fetch": 0,
  "task_max_waiting_in_queue_millis": 0,
  "active_shards_percent_as_number": 50
}
```

① 集群 `status` 为 `yellow`。

② 没有被分配到任何节点的副本数。

集群的健康状况 `yellow` 表示全部主分片都正常运行（集群可以正常服务所有请求），但是副本分片没有全部在正常状态。当前，所有3个副本分片都是 `unassigned` —— 它们都没有被分配到任何节点。在同一个节点上既保存原始数据又保存副本是没有意义的，因为一旦失去了那个节点，我也将失去该节点上的所有副本数据。

当前我的集群是正常运行的，但是在硬件故障有丢失数据的风险。

添加故障转移

当集群中只有一个节点在运行，意味着会有一个节点故障——没有冗余。幸运的是，我只需再添加一个节点即可防止数据丢失。

第二个节点

了解了第二个节点后的情况，可以在同一个目录内，完全依照第一个节点的方式来添加一个新节点（参考[\[running-elasticsearch\]](#)）。多个节点可以共享同一个目录。

当在同一台机器上添加第二个节点，只要它和第一个节点有相同的 `cluster.name` 配置，它就会自动加入集群并加入到其中。但是在不同机器上添加节点的时候，需要加入同一集群，需要配置一个可接收的广播主机列表。更多信息看[\[unicast\]](#)

如果添加第二个节点，我的集群将会如[有 2 个节点的集群——所有主分片和副本分片都被分配](#)所示。

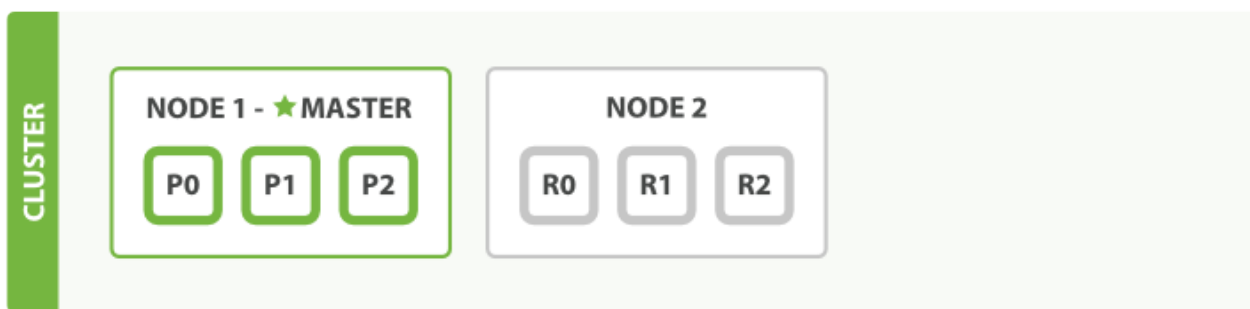


Figure 3. 有两个节点的集群——所有主分片和副本分片都已被分配

当第二个节点加入到集群后，3个副本分片将会分配到2个节点上——一个主分片，一个副本分片。这意味着当集群内任何一个节点出现故障，我的数据都完好无损。

所有新近被索引的文档都会保存在主分片上，然后被并行的复制到副本分片上。这就保证了我既可以从主分片又可以从副本分片上得到文档。

`cluster-health` 在展示的状态是 `green`，表示所有6个分片（包括3个主分片和3个副本分片）都在正常运行。

```
{
  "cluster_name": "elasticsearch",
  "status": "green", ①
  "timed_out": false,
  "number_of_nodes": 2,
  "number_of_data_nodes": 2,
  "active_primary_shards": 3,
  "active_shards": 6,
  "relocating_shards": 0,
  "initializing_shards": 0,
  "unassigned_shards": 0,
  "delayed_unassigned_shards": 0,
  "number_of_pending_tasks": 0,
  "number_of_in_flight_fetch": 0,
  "task_max_waiting_in_queue_millis": 0,
  "active_shards_percent_as_number": 100
}
```

① 集群 `status` 是 `green`。

我的集群现在不是正常运行的，并且处于始用状态。

水平扩容

我正在使用的应用程序按需扩容？当添加了第三个节点，我的集群将会看起来如**有三个节点的集群——为了分散而分片行重新分配**所示。

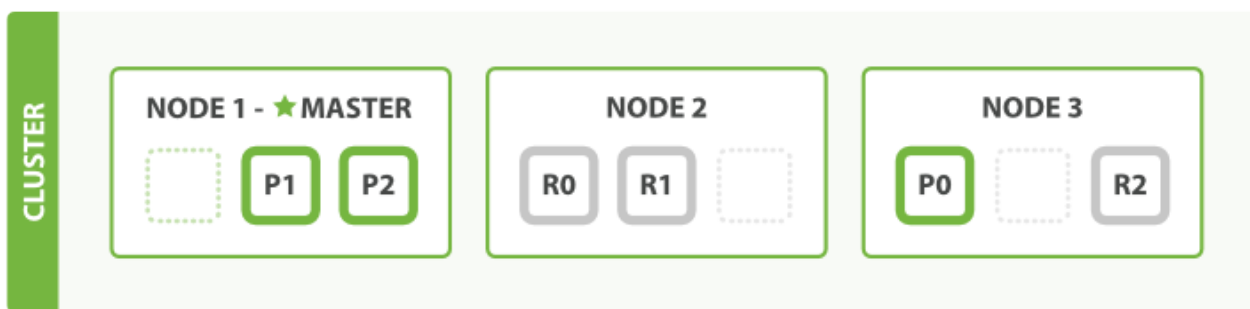


Figure 4. 有三个节点的集群——为了分散负载而分片进行重新分配

Node 1 和 Node 2 上各有一个分片被移到了新的 Node 3 节点，在这个节点上都有 2 个分片，而不是之前的 3 个。这表示一个节点的硬件资源（CPU，RAM，I/O）将被更少的分片所共享，一个分片的性能将会得到提升。

分片是一个功能完整的搜索引擎，它有使用一个节点上的所有资源的能力。我们一共有 6 个分片（3 个主分片和 3 个副本分片）的索引可以最大容纳到 6 个节点，一个节点上存在一个分片，并且一个分片有该节点的全部资源。

更多的容量

但是如果我想要容量超过 6 个节点呢？

主分片的数目在索引建立就已经定了下来。因此，一个数目定义了一个索引能存储的最大数据量。（大小取决于数据、硬件和使用场景。）但是，复制操作——搜索和返回数据——可以同被主分片或副本分片所处理，所以当有越多的副本分片，也将有越高的吞吐量。

在索引中的集群上是可以调整副本分片数目的，我可以按需伸缩集群。我把副本数从原来的 1 加到 2：

```
PUT /blogs/_settings
{
  "number_of_replicas" : 2
}
```

如将参数 `number_of_replicas` 大到 2 所示，`blogs` 索引现在有 9 个分片：3 个主分片和 6 个副本分片。这意味着我可以将集群容量扩大到 9 个节点，一个节点上一个分片。相比原来 3 个节点，集群搜索性能可以提升 3 倍。

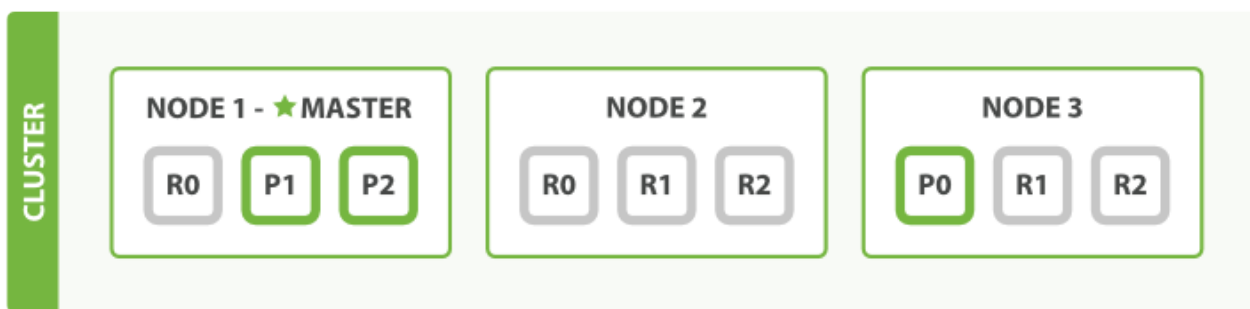


Figure 5. 将参数 `number_of_replicas` 大到 2

NOTE

当然，如果只是在相同节点数目的集群上加更多的副本分片并不能提高性能，因为一个分片从节点上得到的资源会少。需要加更多的硬件资源来提升吞吐量。

但是更多的副本分片数提高了数据冗余量：按照上面的节点配置，我可以在失去2个节点的情况下不丢失任何数据。

故障

我之前 Elasticsearch 可以一点故障，接下来我下一个功能。如果我第一个节点，集群的状态 [变成了一个点后的集群](#)



Figure 6. 变成了一个点后的集群

我的节点是一个主节点。而集群必须有一个主节点来保证正常工作，所以产生的第一件事情就是一个新的主节点：**Node 2**。

在我 **Node 1** 的同时也失去了主分片 **1** 和 **2**，并且在失去主分片的时候索引也不能正常工作。如果此来集群的状况，我看到的状态将会 **red**：不是所有主分片都在正常工作。

幸运的是，在其它节点上存在着一个主分片的完整副本，所以新的主节点立即将一些分片在 **Node 2** 和 **Node 3** 上的副本分片提升为主分片，此集群的状态将会 **yellow**。一个提升主分片的过程是瞬时的，如同按下一个按钮一般。

为什么集群状态是 **yellow** 而不是 **green**？虽然我所有的三个主分片，但是同时置了一个主分片需要 2 个副本分片，而此只存在一副本分片。所以集群不能 **green** 的状态，我不必担心：如果我同时 **Node 2**，我的程序依然可以保持在不丢失任何数据的情况下运行，因为 **Node 3** 一个分片都保留着一副本。

如果我重新 `Node 1`，集群可以将 失的副本分片再次 行分配，那 集群的状 也将如[将参数 number_of_replicas 大到 2](#)所示。如果 `Node 1` 依然 有着之前的分片，它将 去重用它，同 从主分片 制 生了修改的数据文件。

到目前 止， 分片如何使得 Elasticsearch 行水平 容以及数据保障等知 有了一定了解。接下来我 将 述 于分片生命周期的更多 。