

40 | Kubernetes的资源模型与资源管理

2018-11-23 张磊

深入剖析Kubernetes

[进入课程 >](#)



讲述：张磊

时长 10:25 大小 4.78M



你好，我是张磊。今天我和你分享的主题是：Kubernetes 的资源模型与资源管理。

作为一个容器集群编排与管理项目，Kubernetes 为用户提供的基础设施能力，不仅包括了我前面为你讲述的应用定义和描述的部分，还包括了对应用的资源管理和调度的处理。那么，从今天这篇文章开始，我就来为你详细讲解一下后面这部分内容。

而作为 Kubernetes 的资源管理与调度部分的基础，我们要从它的资源模型开始说起。

我在前面的文章中已经提到过，在 Kubernetes 里，Pod 是最小的原子调度单位。这也就意味着，所有跟调度和资源管理相关的属性都应该是属于 Pod 对象的字段。而这其中最重要的部分，就是 Pod 的 CPU 和内存配置，如下所示：

```
1 apiVersion: v1
2 kind: Pod
3 metadata:
4   name: frontend
5 spec:
6   containers:
7   - name: db
8     image: mysql
9     env:
10    - name: MYSQL_ROOT_PASSWORD
11      value: "password"
12    resources:
13      requests:
14        memory: "64Mi"
15        cpu: "250m"
16      limits:
17        memory: "128Mi"
18        cpu: "500m"
19   - name: wp
20     image: wordpress
21     resources:
22       requests:
23         memory: "64Mi"
24         cpu: "250m"
25       limits:
26         memory: "128Mi"
27         cpu: "500m"
```

备注：关于哪些属性属于 Pod 对象，而哪些属性属于 Container，你可以在回顾一下第 14 篇文章 [《深入解析 Pod 对象（一）：基本概念》](#) 中的相关内容。

在 Kubernetes 中，像 CPU 这样的资源被称作“可压缩资源”（compressible resources）。它的典型特点是，当可压缩资源不足时，Pod 只会“饥饿”，但不会退出。

而像内存这样的资源，则被称作“不可压缩资源（incompressible resources）”。当不可压缩资源不足时，Pod 就会因为 OOM（Out-Of-Memory）被内核杀掉。

而由于 Pod 可以由多个 Container 组成，所以 CPU 和内存资源的限额，是要配置在每个 Container 的定义上的。这样，Pod 整体的资源配置，就由这些 Container 的配置值累加得到。

其中，Kubernetes 里为 CPU 设置的单位是“CPU 的个数”。比如，`cpu=1` 指的就是，这个 Pod 的 CPU 限额是 1 个 CPU。当然，具体“1 个 CPU”在宿主机上如何解释，是 1 个 CPU 核心，还是 1 个 vCPU，还是 1 个 CPU 的超线程（Hyperthread），完全取决于宿主机的 CPU 实现方式。Kubernetes 只负责保证 Pod 能够使用到“1 个 CPU”的计算能力。

此外，Kubernetes 允许你将 CPU 限额设置为分数，比如在我们的例子里，CPU limits 的值就是 500m。所谓 500m，指的就是 500 millicpu，也就是 0.5 个 CPU 的意思。这样，这个 Pod 就会被分配到 1 个 CPU 一半的计算能力。

当然，你也可以直接把这个配置写成 `cpu=0.5`。但在实际使用时，我还是推荐你使用 500m 的写法，毕竟这才是 Kubernetes 内部通用的 CPU 表示方式。

而对于内存资源来说，它的单位自然就是 bytes。Kubernetes 支持你使用 Ei、Pi、Ti、Gi、Mi、Ki（或者 E、P、T、G、M、K）的方式来作为 bytes 的值。比如，在我们的例子里，Memory requests 的值就是 64MiB（2 的 26 次方 bytes）。这里要注意区分 MiB（mebibyte）和 MB（megabyte）的区别。

备注：1Mi=1024*1024；1M=1000*1000

此外，不难看到，Kubernetes 里 Pod 的 CPU 和内存资源，实际上还要分为 limits 和 requests 两种情况，如下所示：

 复制代码

```
1 spec.containers[].resources.limits.cpu
2 spec.containers[].resources.limits.memory
3 spec.containers[].resources.requests.cpu
4 spec.containers[].resources.requests.memory
```

这两者的区别其实非常简单：在调度的时候，kube-scheduler 只会按照 requests 的值进行计算。而在真正设置 Cgroups 限制的时候，kubelet 则会按照 limits 的值来进行设置。

更确切地说，当你指定了 `requests.cpu=250m` 之后，相当于将 Cgroups 的 `cpu.shares` 的值设置为 $(250/1000)*1024$ 。而当你没有设置 `requests.cpu` 的时候，`cpu.shares` 默认则是 1024。这样，Kubernetes 就通过 `cpu.shares` 完成了对 CPU 时间的按比例分配。

而如果你指定了 `limits.cpu=500m` 之后，则相当于将 Cgroups 的 `cpu.cfs_quota_us` 的值设置为 $(500/1000)*100ms$ ，而 `cpu.cfs_period_us` 的值始终是 `100ms`。这样，Kubernetes 就为你设置了这个容器只能用到 CPU 的 50%。

而对于内存来说，当你指定了 `limits.memory=128Mi` 之后，相当于将 Cgroups 的 `memory.limit_in_bytes` 设置为 $128 * 1024 * 1024$ 。而需要注意的是，在调度的时候，调度器只会使用 `requests.memory=64Mi` 来进行判断。


Kubernetes 这种对 CPU 和内存资源限额的设计，实际上参考了 Borg 论文中对“动态资源边界”的定义，既：容器化作业在提交时所设置的资源边界，并不一定是调度系统所必须严格遵守的，这是因为在实际场景中，大多数作业使用到的资源其实远小于它所请求的资源限额。

基于这种假设，Borg 在作业被提交后，会主动减小它的资源限额配置，以便容纳更多的作业、提升资源利用率。而当作业资源使用量增加到一定阈值时，Borg 会通过“快速恢复”过程，还原作业原始的资源限额，防止出现异常情况。

而 Kubernetes 的 `requests+limits` 的做法，其实就是上述思路的一个简化版：用户在提交 Pod 时，可以声明一个相对较小的 `requests` 值供调度器使用，而 Kubernetes 真正设置给容器 Cgroups 的，则是相对较大的 `limits` 值。不难看到，这跟 Borg 的思路相通的。

在理解了 Kubernetes 资源模型的设计之后，我再来和你谈谈 Kubernetes 里的 QoS 模型。在 Kubernetes 中，不同的 `requests` 和 `limits` 的设置方式，其实会将这个 Pod 划分到不同的 QoS 级别当中。

当 Pod 里的每一个 Container 都同时设置了 `requests` 和 `limits`，并且 `requests` 和 `limits` 值相等的时候，这个 Pod 就属于 Guaranteed 类别，如下所示：


 复制代码

```
1 apiVersion: v1
2 kind: Pod
3 metadata:
4   name: qos-demo
5   namespace: qos-example
6 spec:
7   containers:
8     - name: qos-demo-ctr
9       image: nginx
```

```
10     resources:
11         limits:
12             memory: "200Mi"
13             cpu: "700m"
14     requests:
15         memory: "200Mi"
16         cpu: "700m"
```


当这个 Pod 创建之后，它的 qosClass 字段就会被 Kubernetes 自动设置为 Guaranteed。需要注意的是，当 Pod 仅设置了 limits 没有设置 requests 的时候，Kubernetes 会自动为它设置与 limits 相同的 requests 值，所以，这也属于 Guaranteed 情况。

而当 Pod 不满足 Guaranteed 的条件，但至少有一个 Container 设置了 requests。那么这个 Pod 就会被划分到 Burstable 类别。比如下面这个例子：

 复制代码

```
1  apiVersion: v1
2  kind: Pod
3  metadata:
4    name: qos-demo-2
5    namespace: qos-example
6  spec:
7    containers:
8      - name: qos-demo-2-ctr
9        image: nginx
10     resources:
11         limits
12             memory: "200Mi"
13         requests:
14             memory: "100Mi"
```

而如果一个 Pod 既没有设置 requests，也没有设置 limits，那么它的 QoS 类别就是 BestEffort。比如下面这个例子：

 复制代码

```
1  apiVersion: v1
2  kind: Pod
3  metadata:
4    name: qos-demo-3
```

```
5 namespace: qos-example
6 spec:
7   containers:
8   - name: qos-demo-3-ctr
9     image: nginx
```

那么，Kubernetes 为 Pod 设置这样三种 QoS 类别，具体有什么作用呢？

实际上，**QoS 划分的主要应用场景，是当宿主机资源紧张的时候，kubelet 对 Pod 进行 Eviction（即资源回收）时需要用到的。**

具体地说，当 Kubernetes 所管理的宿主机上不可压缩资源短缺时，就有可能触发 Eviction。比如，可用内存（memory.available）、可用的宿主机磁盘空间（nodefs.available），以及容器运行时镜像存储空间（imagefs.available）等等。

目前，Kubernetes 为你设置的 Eviction 的默认阈值如下所示：

 复制代码

```
1 memory.available<100Mi
2 nodefs.available<10%
3 nodefs.inodesFree<5%
4 imagefs.available<15%
```

当然，上述各个触发条件在 kubelet 里都是可配置的。比如下面这个例子：

 复制代码

```
1 kubelet --eviction-hard=imagefs.available<10%,memory.available<500Mi,nodefs.available<5%
```

在这个配置中，你可以看到**Eviction 在 Kubernetes 里其实分为 Soft 和 Hard 两种模式。**

其中，Soft Eviction 允许你为 Eviction 过程设置一段“优雅时间”，比如上面例子里的 imagefs.available=2m，就意味着当 imagefs 不足的阈值达到 2 分钟之后，kubelet 才会

开始 Eviction 的过程。

而 Hard Eviction 模式下，Eviction 过程就会在阈值达到之后立刻开始。

Kubernetes 计算 Eviction 阈值的数据来源，主要依赖于从 Cgroups 读取到的值，以及使用 cAdvisor 监控到的数据。

当宿主机 Eviction 阈值达到后，就会进入 MemoryPressure 或者 DiskPressure 状态，从而避免新的 Pod 被调度到这台宿主机上。

而当 Eviction 发生的时候，kubelet 具体会挑选哪些 Pod 进行删除操作，就需要参考这些 Pod 的 QoS 类别了。

首当其冲的，自然是 BestEffort 类别的 Pod。

其次，是属于 Burstable 类别、并且发生“饥饿”的资源使用量已经超出了 requests 的 Pod。

最后，才是 Guaranteed 类别。并且，Kubernetes 会保证只有当 Guaranteed 类别的 Pod 的资源使用量超过了其 limits 的限制，或者宿主机本身正处于 Memory Pressure 状态时，Guaranteed 的 Pod 才可能被选中进行 Eviction 操作。

当然，对于同 QoS 类别的 Pod 来说，Kubernetes 还会根据 Pod 的优先级来进行进一步地排序和选择。

在理解了 Kubernetes 里的 QoS 类别的设计之后，我再来为你讲解一下 **Kubernetes 里一个非常有用的特性：cpuset 的设置**。

我们知道，在使用容器的时候，你可以通过设置 cpuset 把容器绑定到某个 CPU 的核上，而不是像 cpushare 那样共享 CPU 的计算能力。

这种情况下，由于操作系统在 CPU 之间进行上下文切换的次数大大减少，容器里应用的性能会得到大幅提升。事实上，**cpuset 方式，是生产环境里部署在线应用类型的 Pod 时，非常常用的一种方式。**


可是，这样的需求在 Kubernetes 里又该如何实现呢？

其实非常简单。

首先，你的 Pod 必须是 Guaranteed 的 QoS 类型；

然后，你只需要将 Pod 的 CPU 资源的 requests 和 limits 设置为同一个相等的整数值即可。

比如下面这个例子：

 复制代码

```
1 spec:
2   containers:
3   - name: nginx
4     image: nginx
5     resources:
6       limits:
7         memory: "200Mi"
8         cpu: "2"
9       requests:
10        memory: "200Mi"
11        cpu: "2"
```

这时候，该 Pod 就会被绑定在 2 个独占的 CPU 核上。当然，具体是哪两个 CPU 核，是由 kubelet 为你分配的。

以上，就是 Kubernetes 的资源模型和 QoS 类别相关的主要内容。

总结

在本篇文章中，我先为你详细讲解了 Kubernetes 里对资源的定义方式和资源模型的设计。然后，我为你讲述了 Kubernetes 里对 Pod 进行 Eviction 的具体策略和实践方式。

正是基于上述讲述，在实际的使用中，我强烈建议你将 DaemonSet 的 Pod 都设置为 Guaranteed 的 QoS 类型。否则，一旦 DaemonSet 的 Pod 被回收，它又会立即在原宿主机上被重建出来，这就使得前面资源回收的动作，完全没有意义了。

思考题

为什么宿主机进入 MemoryPressure 或者 DiskPressure 状态后，新的 Pod 就不会被调度到这台宿主机上呢？

感谢你的收听，欢迎你给我留言，也欢迎分享给更多的朋友一起阅读。



深入剖析 Kubernetes

Kubernetes 原来可以如此简单

张磊
Kubernetes 社区
资深成员与项目维护者



新版升级：点击「 请朋友读」，10位好友免费读，邀请订阅更有**现金**奖励。

© 版权归极客邦科技所有，未经许可不得传播售卖。页面已增加防盗追踪，如有侵权极客邦将依法追究其法律责任。

上一篇 39 | 谈谈Service与Ingress

下一篇 41 | 十字路口上的Kubernetes默认调度器

精选留言 (20)

写留言



DJH

2018-11-23

9

“为什么宿主机进入 MemoryPressure 或者 Dis...”

这是因为给宿主机打了污点标记吗？

展开

作者回复: 对



wilder

2018-11-23

👍 4

极客时间里面最爱的课程，没有之一，哈哈哈哈哈哈

展开 ▾



虎虎

2018-11-23

👍 3

能否分享一下给namespace 设置quota的经验呢？

如果设置的太小，会造成资源的浪费。如果设置太大，又怕起不到限制的作用。一个namespace使用资源太多可能会影响其他namespace用户的使用。

...

展开 ▾



gogo

2018-11-23

👍 3

老师您好，cpu设置limit之后，容器的cpu使用率永远不会超过这个限制对吗？而mem设置limit之后，容器mem使用率有可能超过这个限制而被kill掉，也就是说设置了cpu limit之后，容器永远不会因为cpu超过限制而被kill对吗

展开 ▾

作者回复: 是



zfei

2018-11-30

👍 2

请问老师，cpuset为2，这个Pod就独占两个cpu核上，假如宿主机总共只有10个cpu核，那么这台机就只能运行5个cpuset=2的Pod吗

展开 ▾



刘岚乔月

2018-11-25

👍 2

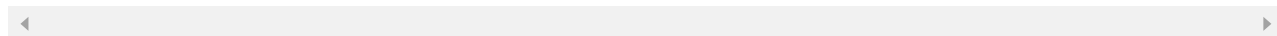
1、3、5都在追文章，一直都有一个疑问，请作者能解惑下。

对于主java是其他语言（非运维）的同学来说，我们是否需要深入了解k8s和docker（还是停留在使用层面）我想一直跟着学的同学大部门还是冲着能找到更好的工作去的（有情怀的同学请忽略）

目前更大公司的招聘对于要求掌握k8s和docker的基本上都是运维岗位。...

展开 ▾

作者回复: kubernetes 是云时代的开发者工具。重要的事情只说一遍。



阿鹏

2018-11-27

👍 1

老师，关于资源隔离我有三个问题想请教一下。

第一，正如您所说，/proc是不能被隔离的，但是我们可以通过lxcfs或者高版本的jdk版本来让容器里的服务知道自己的资源限制，或者还有方式，老师有推荐的吗？

第二，我使用lxcfs隔离后，容器内/proc/meminfo文件确实是限制后的内存大小，但是容器内/proc/cpuinfo的信息跟宿主机是一样的，那么容器内的应该要怎么知道自己正确的...

展开 ▾



汪浩

2018-11-25

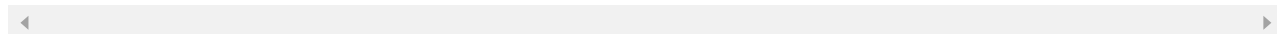
👍 1

被称作“不可压缩资源（compressible resources）

应该是 uncompressible

展开 ▾

作者回复: 对，得然后编辑改一下



unique

2018-11-23

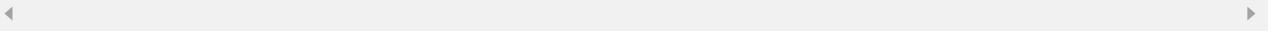
👍 1

这时候，该 Pod 就会被绑定在 2 个独占的 CPU 核上。

独占的意思就是其它pod 不能使用这两个CPU了么？

展开 ▾

作者回复: 对



1520145770...

2019-03-18



DiskPressure nodefs.available, nodefs.inodesFree, imagefs.available, or imagefs.inodesFree Available disk space and inodes on either the node' s root filesystem or image filesystem has satisfied an eviction threshold 老师这个有点懵，请问是什么意思？

展开 ▾



Lucius

2019-03-05



"将 DaemonSet 的 Pod 都设置为 Guarant..." 不太懂, Guaranteed和重新创建有什么关系



(!0+)[...

2018-12-20



思考题

因为controller看到node上有污点？



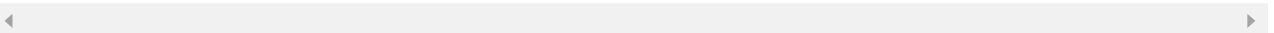
初学者

2018-12-01



老师，能不能讲一下kubernetes是如何划分和管理宿主机上的cgroups结构的？

作者回复: 这个非常简单，就是按qos划分成三种



大彬

2018-11-30



我在实际操作中遇到一个问题，我用 aws kops 部署了一个 3 master，3 node 的集群，

服务器选的是 t2.medium (2核4G)。用 helm 安装 stable/elasticsearch 无法启动，node 会出现 Not Ready，System OOM Killed 这样的错误，也无法 ssh 登录了。我想知道有没有办法避免这样的故障，能给 node 系统至少预留一些资源，可以允许 pod 无法启动，也不至于节点崩溃没有响应了？

展开 ▾



虎虎

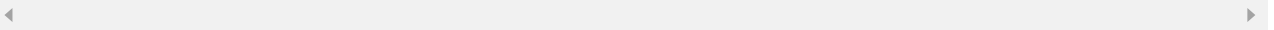
2018-11-27



在namespace limitRange 里面设置了default request 和 default limit之后，创建出来的pod即使不显式指定limit和request，是不是也是guaranteed？

展开 ▾

作者回复: 是的



leo

2018-11-26



老师，请教您个问题，当我们使用istio的时候，业务服务的鉴权如何做？在spring cloud 里可以在api-gateway层完成。可istio的gateway应该没有这样的职责。



每日都想上...

2018-11-23



每天看每天都有吸收知识的课程

展开 ▾



蜗牛

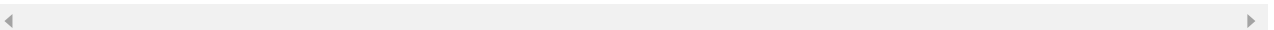
2018-11-23



如果某个Guaranteed Pod 的 Mem 设定为了256，在宿主机资源不紧张但该Pod 的的 mem 使用量达到了256以后会出现什么情况？会被oom杀掉吗？

展开 ▾

作者回复: 见QoS部分





jaxzhai

2018-11-23



因为调度器在预选的时候会判断节点资源是否充足，如果资源不够将不会被调度。



维苏威的血

2018-11-23



不可压缩资源 typo

展开 ∨