

xiaolincoding.com

5.7 线程崩溃了，进程也会崩溃吗？

小林coding

27-33 minutes

4.3 内存满了，会发生什么？

大家好，我是小林。

前几天有位读者留言说，面腾讯时，被问了两个内存管理的问题：



小林哥，腾讯面试官问了一个问题，为什么操作系统需要内存管理和虚拟内存，除了给进程分配内存和防止进程间相互影响，还有什么作用？



然后提到除了OOM，在内存满了之后会有什么处理？我当场就懵了😓

先来说说第一个问题：虚拟内存有什么作用？

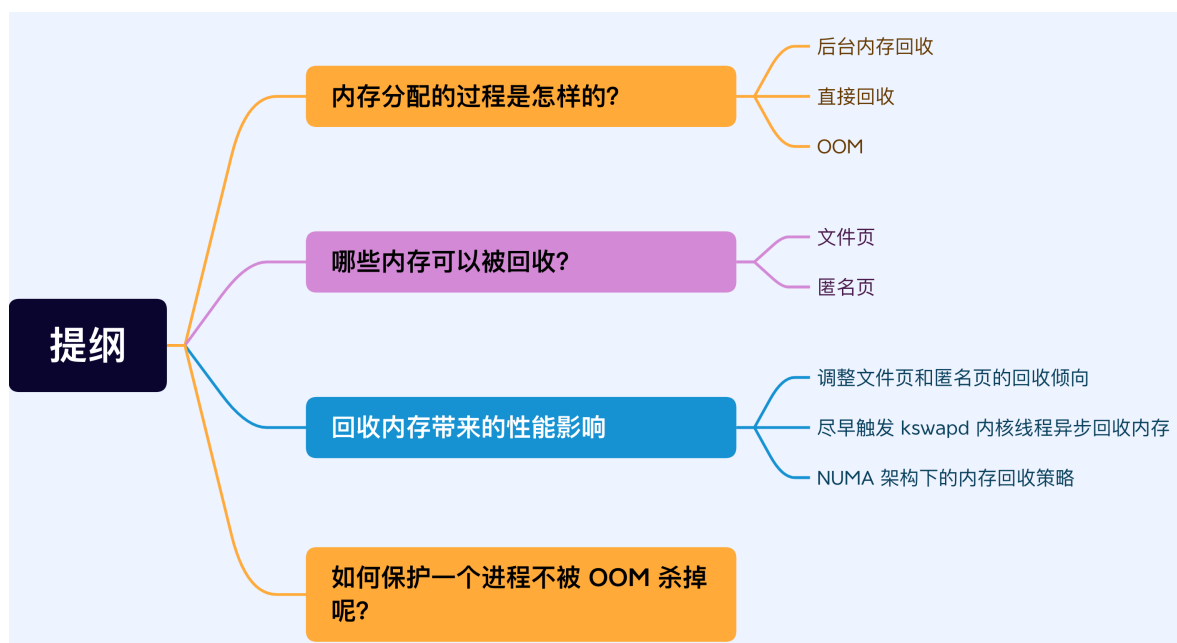
- 第一，虚拟内存可以使得进程对运行内存超过物理内存大小，因为程序运行符合局部性原理，CPU 访问内存会有很明显的重复访问的倾向性，对于那些没有被经常使

用到的内存，我们可以把它换出到物理内存之外，比如硬盘上的 swap 区域。

- 第二，由于每个进程都有自己的页表，所以每个进程的虚拟内存空间就是相互独立的。进程也没有办法访问其他进程的页表，所以这些页表是私有的，这就解决了多进程之间地址冲突的问题。
- 第三，页表里的页表项中除了物理地址之外，还有一些标记属性的比特，比如控制一个页的读写权限，标记该页是否存在等。在内存访问方面，操作系统提供了更好的安全性。

然后今天主要是聊聊第二个问题，「**系统内存紧张时，会发生什么？**」

发车！



内存分配的过程是怎样的？

应用程序通过 `malloc` 函数申请内存的时候，实际上申请的是虚拟内存，此时并不会分配物理内存。

当应用程序读写了这块虚拟内存，CPU 就会去访问这个虚拟内存，这时会发现这个虚拟内存没有映射到物理内存，CPU 就会产生**缺页中断**，进程会从用户态切换到内核态，并将缺页中断交给内核的 Page Fault Handler（缺页中断函数）处理。

缺页中断处理函数会看是否有空闲的物理内存，如果有，就直接分配物理内存，并建立虚拟内存与物理内存之间的映射关系。

如果没有空闲的物理内存，那么内核就会开始进行**回收内存**的工作，回收的方式主要是两种：直接内存回收和后台内存回收。

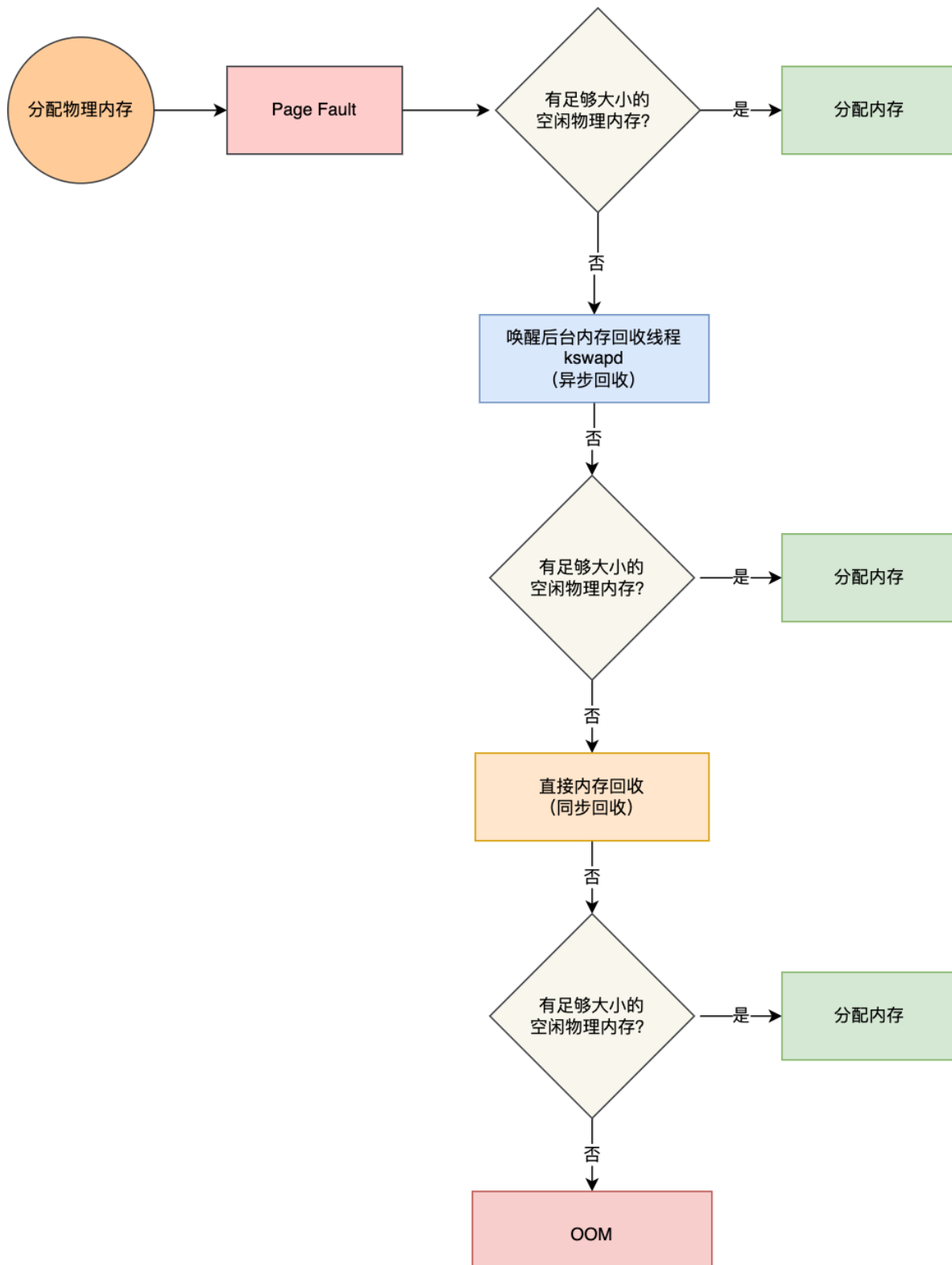
- **后台内存回收**（`kswapd`）：在物理内存紧张的时候，会唤醒 `kswapd` 内核线程来回收内存，这个回收内存的过程**异步**的，不会阻塞进程的执行。
- **直接内存回收**（`direct reclaim`）：如果后台异步回收跟不上进程内存申请的速度，就会开始直接回收，这个回收内存的过程是**同步**的，会阻塞进程的执行。

如果直接内存回收后，空闲的物理内存仍然无法满足此次物理内存的申请，那么内核就会放最后的大招了——**触发 OOM（Out of Memory）机制**。

OOM Killer 机制会根据算法选择一个占用物理内存较高

的进程，然后将其杀死，以便释放内存资源，如果物理内存依然不足，OOM Killer 会继续杀死占用物理内存较高的进程，直到释放足够的内存位置。

申请物理内存的过程如下图：



哪些内存可以被回收？

系统内存紧张的时候，就会进行回收内存的工作，那具体哪些内存是可以被回收的呢？

主要有两类内存可以被回收，而且它们的回收方式也不同。

- **文件页**（File-backed Page）：内核缓存的磁盘数据（Buffer）和内核缓存的文件数据（Cache）都叫作文件页。大部分文件页，都可以直接释放内存，以后有需要时，再从磁盘重新读取就可以了。而那些被应用程序修改过，并且暂时还没写入磁盘的数据（也就是脏页），就得先写入磁盘，然后才能进行内存释放。所以，**回收干净页的方式是直接释放内存，回收脏页的方式是先写回磁盘后再释放内存。**
- **匿名页**（Anonymous Page）：这部分内存没有实际载体，不像文件缓存有硬盘文件这样一个载体，比如堆、栈数据等。这部分内存很可能还要再次被访问，所以不能直接释放内存，它们**回收的方式是通过 Linux 的 Swap 机制**，Swap 会把不常访问的内存先写到磁盘中，然后释放这些内存，给其他更需要的进程使用。再次访问这些内存时，重新从磁盘读入内存就可以了。

文件页和匿名页的回收都是基于 LRU 算法，也就是优先回收不常访问的内存。LRU 回收算法，实际上维护着

active 和 inactive 两个双向链表，其中：

- **active_list** 活跃内存页链表，这里存放的是最近被访问过（活跃）的内存页；
- **inactive_list** 不活跃内存页链表，这里存放的是很少被访问（非活跃）的内存页；

越接近链表尾部，就表示内存页越不常访问。这样，在回收内存时，系统就可以根据活跃程度，优先回收不活跃的内存。

活跃和非活跃的内存页，按照类型的不同，又分别分为文件页和匿名页。可以从 /proc/meminfo 中，查询它们的大小，比如：

回收内存带来的性能影响

在前面我们知道了回收内存有两种方式。

- 一种是后台内存回收，也就是唤醒 kswapd 内核线程，这种方式是异步回收的，不会阻塞进程。
- 一种是直接内存回收，这种方式是同步回收的，会阻塞进程，这样就会造成很长时间的延迟，以及系统的 CPU 利用率会升高，最终引起系统负荷飙高。

可被回收的内存类型有文件页和匿名页：

- 文件页的回收：对于干净页是直接释放内存，这个操作不会影响性能，而对于脏页会先写回到磁盘再释放内

存，这个操作会发生磁盘 I/O 的，这个操作是会影响系统性能的。

- 匿名页的回收：如果开启了 Swap 机制，那么 Swap 机制会将不常访问的匿名页换出到磁盘中，下次访问时，再从磁盘换入到内存中，这个操作是会影响系统性能的。

可以看到，回收内存的操作基本都会发生磁盘 I/O 的，如果回收内存的操作很频繁，意味着磁盘 I/O 次数会很多，这个过程势必会影响系统的性能，整个系统给人的感觉就是很卡。

下面针对回收内存导致的性能影响，说说常见的解决方式。

调整文件页和匿名页的回收倾向

从文件页和匿名页的回收操作来看，文件页的回收操作对系统的影响相比匿名页的回收操作会少一点，因为文件页对于干净页回收是不会发生磁盘 I/O 的，而匿名页的 Swap 换入换出这两个操作都会发生磁盘 I/O。

Linux 提供了一个 `/proc/sys/vm/swappiness` 选项，用来调整文件页和匿名页的回收倾向。

`swappiness` 的范围是 0-100，数值越大，越积极使用 Swap，也就是更倾向于回收匿名页；数值越小，越消极使用 Swap，也就是更倾向于回收文件页。

一般建议 `swappiness` 设置为 0（默认值是 60），这样在回收内存的时候，会更倾向于文件页的回收，但是并不代表不会回收匿名页。

尽早触发 `kswapd` 内核线程异步回收内存

如何查看系统的直接内存回收和后台内存回收的指标？

我们可以使用 `sar -B 1` 命令来观察：

```
[root@xiaolin ~]# sar -B 1
Linux 3.10.0-514.26.2.el7.x86_64 (xiaolin)      05/23/2022      _x86_64_      (1 CPU)

11:53:10 AM pgpgin/s pgpgout/s  fault/s  majflt/s  pgfree/s  pgscank/s pgscand/s pgsteal/s  %vmeff
11:53:11 AM      0.00      0.00    84.69      0.00    35.71      0.00      0.00      0.00      0.00
11:53:12 AM      0.00      0.00    29.00      0.00    21.00      0.00      0.00      0.00      0.00
11:53:13 AM      0.00      0.00    19.19      0.00    21.21      0.00      0.00      0.00      0.00
11:53:14 AM      0.00      0.00    19.19      0.00    21.21      0.00      0.00      0.00      0.00
```

图中红色框住的就是后台内存回收和直接内存回收的指标，它们分别表示：

- `pgscank/s` : `kswapd`(后台回收线程) 每秒扫描的 page 个数。
- `pgscand/s`: 应用程序在内存申请过程中每秒直接扫描的 page 个数。
- `pgsteal/s`: 扫描的 page 中每秒被回收的个数 (`pgscank+pgscand`) 。

如果系统时不时发生抖动，并且在抖动的时间段里如果通过 `sar -B` 观察到 `pgscand` 数值很大，那大概率是因为「直接内存回收」导致的。

针对这个问题，解决的办法就是，可以通过尽早的触发

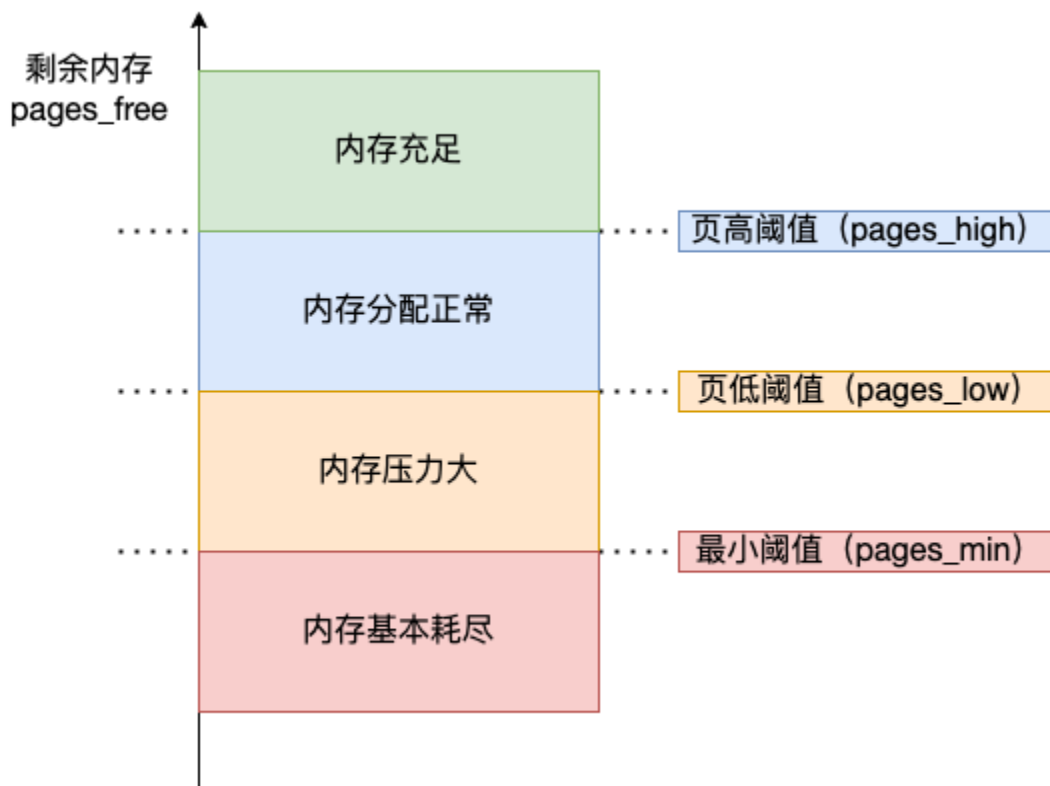
「后台内存回收」来避免应用程序进行直接内存回收。

什么条件下才能触发 kswapd 内核线程回收内存呢？

内核定义了三个内存阈值（watermark，也称为水位），用来衡量当前剩余内存（pages_free）是否充裕或者紧张，分别是：

- 页最小阈值（pages_min）；
- 页低阈值（pages_low）；
- 页高阈值（pages_high）；

这三个内存阈值会划分为四种内存使用情况，如下图：



kswapd 会定期扫描内存的使用情况，根据剩余内存（pages_free）的情况来进行内存回收的工作。

- 图中绿色部分：如果剩余内存（pages_free）大于 页高

阈值 (`pages_high`)，说明剩余内存是充足的；

- 图中蓝色部分：如果剩余内存 (`pages_free`) 在页高阈值 (`pages_high`) 和页低阈值 (`pages_low`) 之间，说明内存有一定压力，但还可以满足应用程序申请内存的请求；
- 图中橙色部分：如果剩余内存 (`pages_free`) 在页低阈值 (`pages_low`) 和页最小阈值 (`pages_min`) 之间，说明内存压力比较大，剩余内存不多了。**这时 `kswapd0` 会执行内存回收，直到剩余内存大于高阈值 (`pages_high`) 为止。**虽然会触发内存回收，但是不会阻塞应用程序，因为两者关系是异步的。
- 图中红色部分：如果剩余内存 (`pages_free`) 小于页最小阈值 (`pages_min`)，说明用户可用内存都耗尽了，此时就会**触发直接内存回收**，这时应用程序就会被阻塞，因为两者关系是同步的。

可以看到，当剩余内存页 (`pages_free`) 小于页低阈值 (`pages_low`)，就会触发 `kswapd` 进行后台回收，然后 `kswapd` 会一直回收到剩余内存页 (`pages_free`) 大于页高阈值 (`pages_high`)。

也就是说 `kswapd` 的活动空间只有 `pages_low` 与 `pages_min` 之间的这段区域，如果剩余内存低于了 `pages_min` 会触发直接内存回收，高于了 `pages_high` 又不会唤醒 `kswapd`。

页低阈值 (`pages_low`) 可以通过内核选项 `/proc/sys/vm/min_free_kbytes` (该参数代表系统所保留空闲内存的最低限) 来间接设置。

`min_free_kbytes` 虽然设置的是页最小阈值 (`pages_min`)，但是页高阈值 (`pages_high`) 和页低阈值 (`pages_low`) 都是根据页最小阈值 (`pages_min`) 计算生成的，它们之间的计算关系如下：

如果系统时不时发生抖动，并且通过 `sar -B` 观察到 `pgscand` 数值很大，那大概率是因为直接内存回收导致的，这时可以增大 `min_free_kbytes` 这个配置选项来及早地触发后台回收，然后继续观察 `pgscand` 是否会降为 0。

增大了 `min_free_kbytes` 配置后，这会使得系统预留过多的空闲内存，从而在一定程度上降低了应用程序可使用的内存量，这在一定程度上浪费了内存。极端情况下设置 `min_free_kbytes` 接近实际物理内存大小时，留给应用程序的内存就会太少而可能会频繁地导致 OOM 的发生。

所以在调整 `min_free_kbytes` 之前，需要先思考一下，应用程序更加关注什么，如果关注延迟那就适当地增大 `min_free_kbytes`，如果关注内存的使用量那就适当地调小 `min_free_kbytes`。

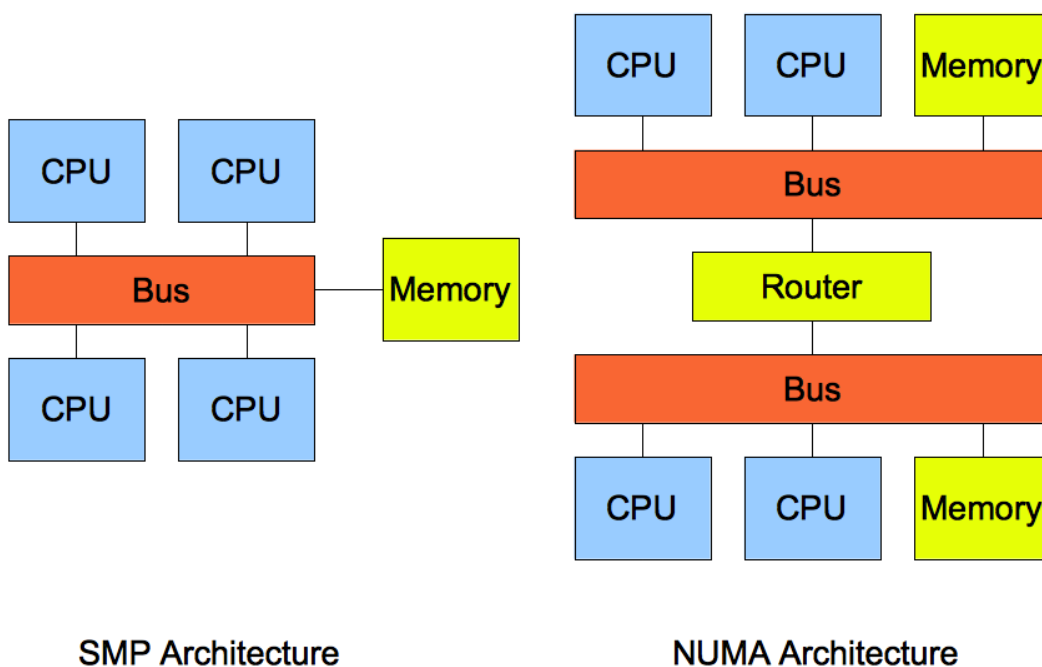
NUMA 架构下的内存回收策略

什么是 NUMA 架构？

再说 NUMA 架构前，先给大家说说 SMP 架构，这两个架构都是针对 CPU 的。

SMP 指的是一种**多个 CPU 处理器共享资源的电脑硬件架构**，也就是说每个 CPU 地位平等，它们共享相同的物理资源，包括总线、内存、IO、操作系统等。每个 CPU 访问内存所用时间都是相同的，因此，这种系统也被称为一致存储访问结构（UMA, Uniform Memory Access）。

随着 CPU 处理器核数的增多，多个 CPU 都通过一个总线访问内存，这样总线的带宽压力会越来越大，同时每个 CPU 可用带宽会减少，这也就是 SMP 架构的问题。



为了解决 SMP 架构的问题，就研制出了 NUMA 结构，即非一致存储访问结构（Non-uniform memory access, NUMA）。

NUMA 架构将每个 CPU 进行了分组，每一组 CPU 用 Node 来表示，一个 Node 可能包含多个 CPU。

每个 Node 有自己独立的资源，包括内存、IO 等，每个 Node 之间可以通过互联模块总线（QPI）进行通信，所以，也就意味着每个 Node 上的 CPU 都可以访问到整个系统中的所有内存。但是，访问远端 Node 的内存比访问本地内存要耗时很多。

〔NUMA 架构跟回收内存有什么关系？

在 NUMA 架构下，当某个 Node 内存不足时，系统可以从其他 Node 寻找空闲内存，也可以从本地内存中回收内存。

具体选哪种模式，可以通过 `/proc/sys/vm/zone_reclaim_mode` 来控制。它支持以下几个选项：

- 0（默认值）：在回收本地内存之前，在其他 Node 寻找空闲内存；
- 1：只回收本地内存；
- 2：只回收本地内存，在本地回收内存时，可以将文件页中的脏页写回硬盘，以回收内存。

- 4: 只回收本地内存，在本地回收内存时，可以用 swap 方式回收内存。

在使用 NUMA 架构的服务器，如果系统出现还有一半内存的时候，却发现系统频繁触发「直接内存回收」，导致了影响了系统性能，那么大概率是因为 `zone_reclaim_mode` 没有设置为 0，导致当本地内存不足的时候，只选择回收本地内存的方式，而不去使用其他 Node 的空闲内存。

虽然说访问远端 Node 的内存比访问本地内存要耗时很多，但是相比内存回收的危害而言，访问远端 Node 的内存带来的性能影响还是比较小的。因此，`zone_reclaim_mode` 一般建议设置为 0。

如何保护一个进程不被 OOM 杀掉呢？

在系统空闲内存不足的情况，进程申请了一个很大的内存，如果直接内存回收都无法回收出足够大的空闲内存，那么就会触发 OOM 机制，内核就会根据算法选择一个进程杀掉。

Linux 到底是根据什么标准来选择被杀的进程呢？这就要提到一个在 Linux 内核里有一个 `oom_badness()` 函数，它会把系统中可以被杀掉的进程扫描一遍，并对每个进程打分，得分最高的进程就会被首先杀掉。

进程得分的结果受下面这两个方面影响：

- 第一，进程已经使用的物理内存页面数。
- 第二，每个进程的 OOM 校准值 `oom_score_adj`。它是可以通过 `/proc/[pid]/oom_score_adj` 来配置的。我们可以在设置 -1000 到 1000 之间的任意一个数值，调整进程被 OOM Kill 的几率。

函数 `oom_badness()` 里的最终计算方法是这样的：

用「系统总的可用页面数」乘以「OOM 校准值 `oom_score_adj`」再除以 1000，最后再加上进程已经使用的物理页面数，计算出来的值越大，那么这个进程被 OOM Kill 的几率也就越大。

每个进程的 `oom_score_adj` 默认值都为 0，所以最终得分跟进程自身消耗的内存有关，消耗的内存越大越容易被杀掉。我们可以通过调整 `oom_score_adj` 的数值，来改成进程的得分结果：

- 如果你不想某个进程被首先杀掉，那你可以调整该进程的 `oom_score_adj`，从而改变这个进程的得分结果，降低该进程被 OOM 杀死的概率。
- 如果你想某个进程无论如何都不能被杀掉，那你可以将 `oom_score_adj` 配置为 -1000。

我们最好将一些很重要的系统服务的 `oom_score_adj` 配置为 -1000，比如 `sshd`，因为这些系统服务一旦被杀掉，我们就很难再登陆进系统了。

但是，不建议将我们自己的业务程序的 `oom_score_adj` 设置为 -1000，因为业务程序一旦发生了内存泄漏，而它又不能杀掉，这就会导致随着它的内存开销变大，OOM killer 不停地被唤醒，从而把其他进程一个个给杀掉。

参考资料：

- <https://time.geekbang.org/column/article/277358>
- <https://time.geekbang.org/column/article/75797>
- <https://www.jianshu.com/p/e40e8813842f>

总结

内核在给应用程序分配物理内存的时候，如果空闲物理内存不够，那么就会进行内存回收的工作，主要有两种方式：

- 后台内存回收：在物理内存紧张的时候，会唤醒 `kswapd` 内核线程来回收内存，这个回收内存的过程异步的，不会阻塞进程的执行。
- 直接内存回收：如果后台异步回收跟不上进程内存申请的速度，就会开始直接回收，这个回收内存的过程是同步的，会阻塞进程的执行。

可被回收的内存类型有文件页和匿名页：

- 文件页的回收：对于干净页是直接释放内存，这个操作

不会影响性能，而对于脏页会先写回到磁盘再释放内存，这个操作会发生磁盘 I/O 的，这个操作是会影响系统性能的。

- 匿名页的回收：如果开启了 Swap 机制，那么 Swap 机制会将不常访问的匿名页换出到磁盘中，下次访问时，再从磁盘换入到内存中，这个操作是会影响系统性能的。

文件页和匿名页的回收都是基于 LRU 算法，也就是优先回收不常访问的内存。回收内存的操作基本都会发生磁盘 I/O 的，如果回收内存的操作很频繁，意味着磁盘 I/O 次数会很多，这个过程势必会影响系统的性能。

针对回收内存导致的性能影响，常见的解决方式。

- 设置 `/proc/sys/vm/swappiness`，调整文件页和匿名页的回收倾向，尽量倾向于回收文件页；
- 设置 `/proc/sys/vm/min_free_kbytes`，调整 kswapd 内核线程异步回收内存的时机；
- 设置 `/proc/sys/vm/zone_reclaim_mode`，调整 NUMA 架构下内存回收策略，建议设置为 0，这样在回收本地内存之前，会在其他 Node 寻找空闲内存，从而避免在系统还有很多空闲内存的情况下，因本地 Node 的本地内存不足，发生频繁直接内存回收导致性能下降的问题；

在经历完直接内存回收后，空闲的物理内存大小依然不

够，那么就会触发 OOM 机制，OOM killer 就会根据每个进程的内存占用情况和 `oom_score_adj` 的值进行打分，得分最高的进程就会被首先杀掉。

我们可以通过调整进程的 `/proc/[pid]/oom_score_adj` 值，来降低被 OOM killer 杀掉的概率。

完！

新的图解文章都在公众号首发，别忘记关注了哦！如果你想加入百人技术交流群，扫码下方二维码回复「加群」。

