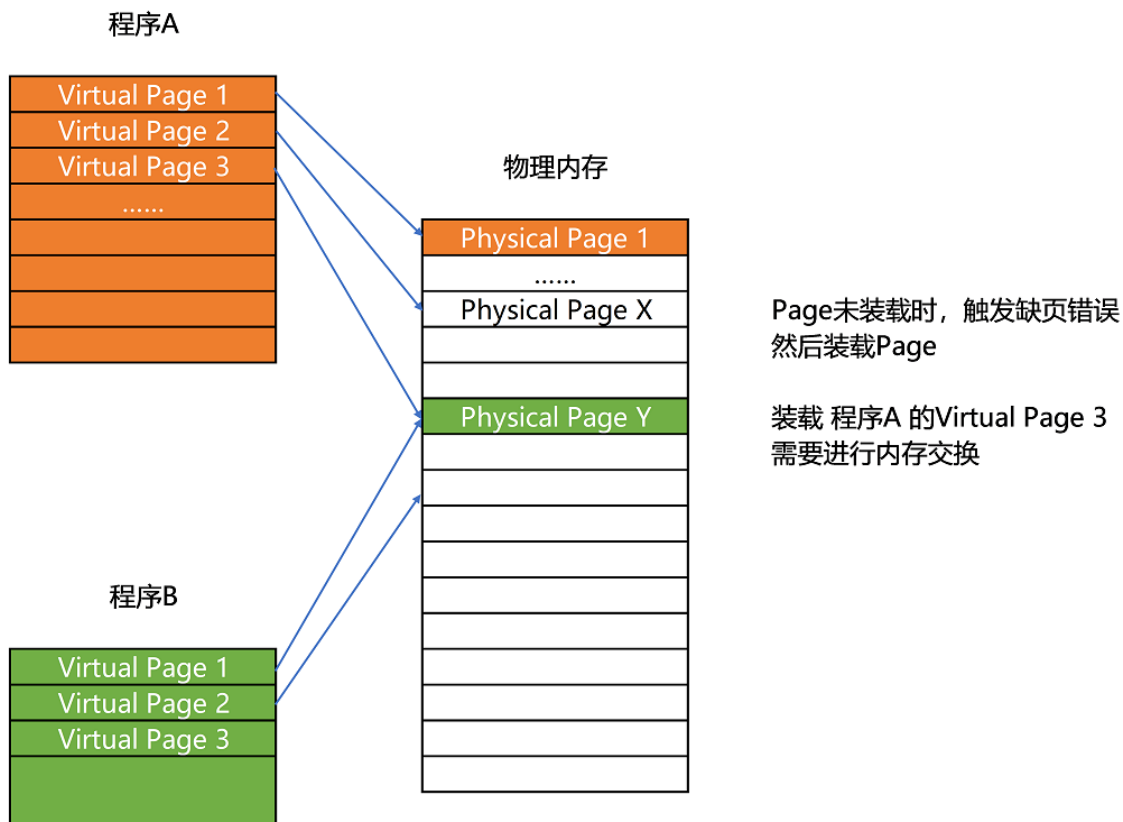


二

40 理解内存（上）：虚拟内存和内存保护是什么？

我们在专栏一开始说过，计算机有五大组成部分，分别是：运算器、控制器、存储器、输入设备和输出设备。如果说计算机最重要的组件，是承担了运算器和控制器作用的 CPU，那内存就是我们第二重要的组件了。内存是五大组成部分里面的存储器，我们的指令和数据，都需要先加载到内存里面，才会被 CPU 拿去执行。

专栏第 9 讲，我们讲了程序装载到内存的过程。可以知道，在我们日常使用的 Linux 或者 Windows 操作系统下，程序并不能直接访问物理内存。



我们的内存需要被分成固定大小的页（Page），然后再通过虚拟内存地址（Virtual Address）到物理内存地址（Physical Address）的地址转换（Address Translation），才能到达实际存放数据的物理内存位置。而我们的程序看到的内存地址，都是虚拟内存地址。

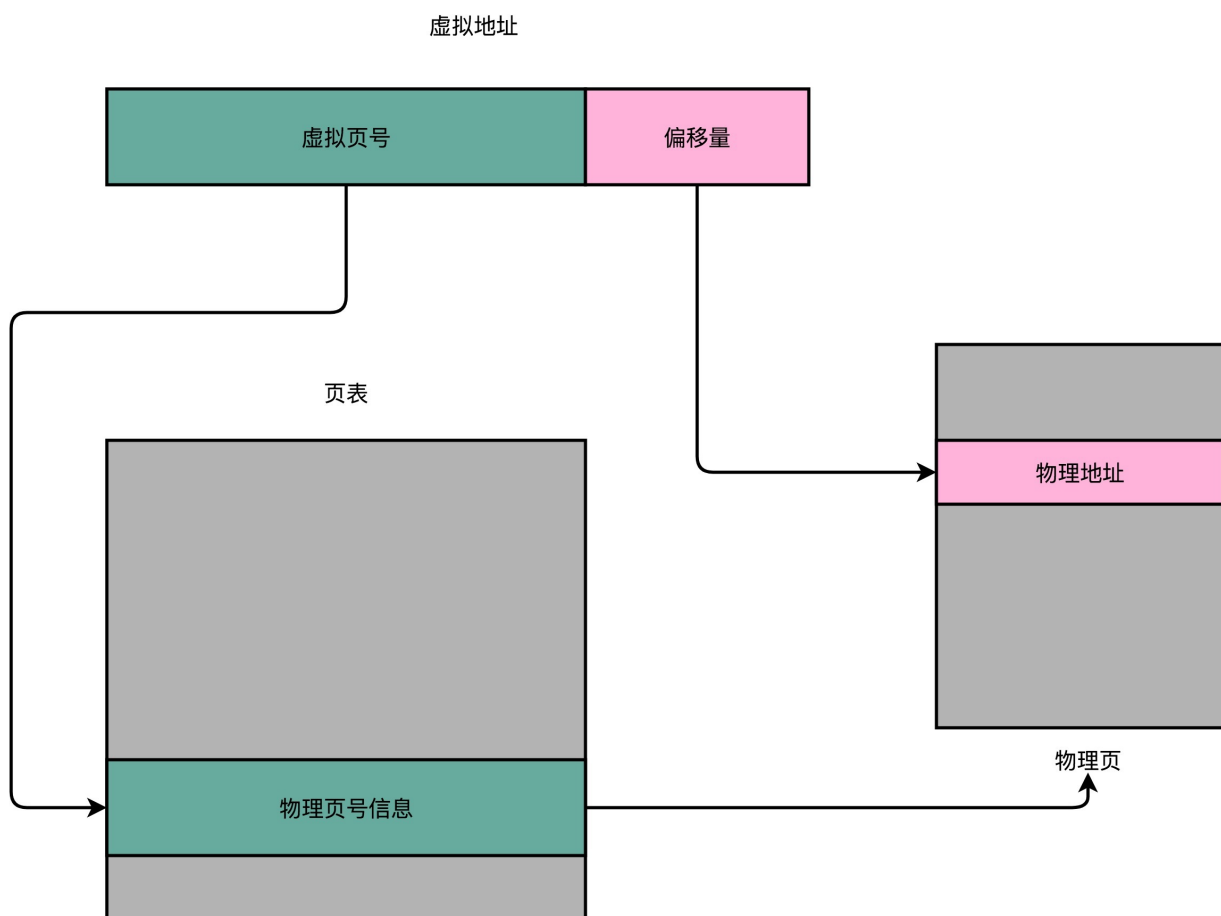
既然如此，这些虚拟内存地址究竟是怎么转换成物理内存地址的呢？这一讲里，我们就来看一看。

简单页表

想要把虚拟内存地址，映射到物理内存地址，最直观的办法，就是来建一张映射表。这个映射表，能够实现虚拟内存里面的页，到物理内存里面的页的一一映射。这个映射表，在计算机里面，就叫作**页表**（Page Table）。

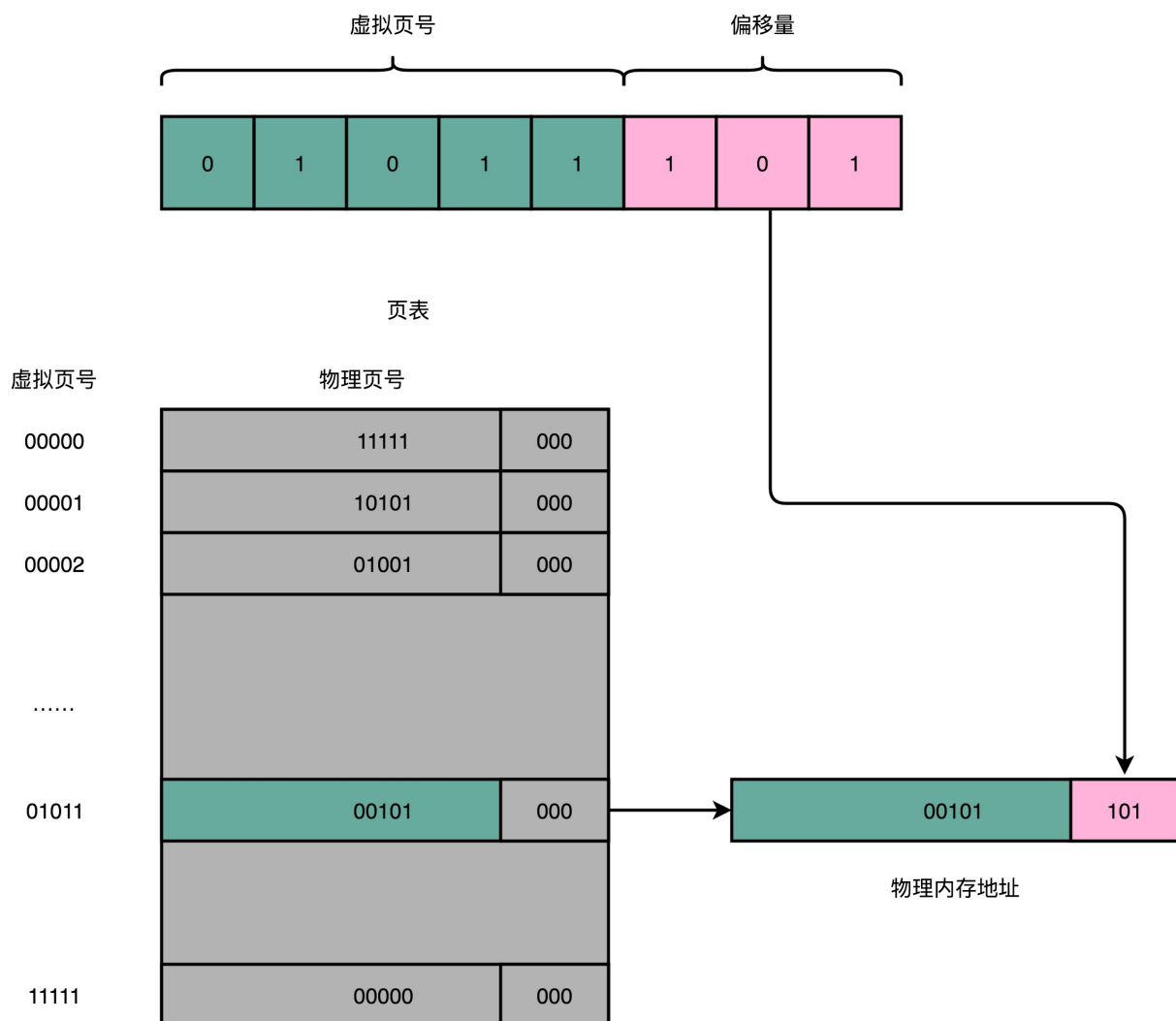
页表这个地址转换的办法，会把一个内存地址分成**页号**（Directory）和**偏移量**（Offset）两个部分。这么说太理论了，我以一个 32 位的内存地址为例，帮你理解这个概念。

其实，前面的高位，就是内存地址的页号。后面的低位，就是内存地址里面的偏移量。做地址转换的页表，只需要保留虚拟内存地址的页号和物理内存地址的页号之间的映射关系就可以了。同一个页里面的内存，在物理层面是连续的。以一个页的大小是 4K 比特（4KiB）为例，我们需要 20 位的高位，12 位的低位。



总结一下，对于一个内存地址转换，其实就是这样三个步骤：

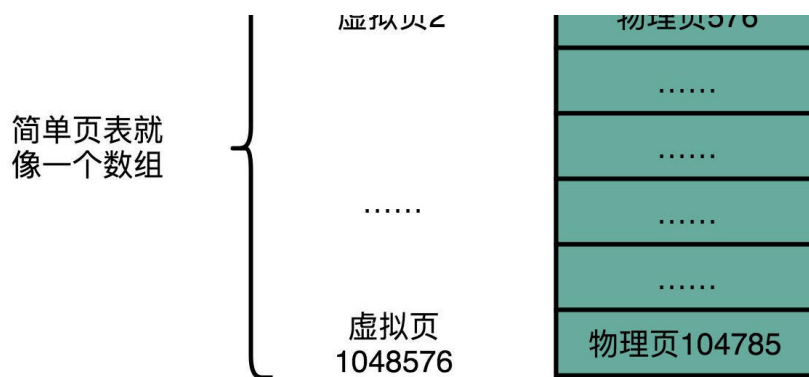
1. 把虚拟内存地址，切分成页号和偏移量的组合；
2. 从页表里面，查询出虚拟页号，对应的物理页号；
3. 直接拿物理页号，加上前面的偏移量，就得到了物理内存地址。



看起来这个逻辑似乎很简单，很容易理解，不过问题马上就来了。你能算一算，这样一个页表需要多大的空间吗？我们以 32 位的内存地址空间为例，你可以暂停一下，拿出纸笔算一算。

不知道你算出的数字是多少？32 位的内存地址空间，页表一共需要记录 2^{20} 个到物理页号的映射关系。这个存储关系，就好比一个 2^{20} 大小的数组。一个页号是完整的 32 位的 4 字节 (Byte)，这样一个页表就需要 4MB 的空间。听起来 4MB 的空间好像还不小啊，毕竟我们现在的内存至少也有 4GB，服务器上有个几十 GB 的内存是很正常。





不过，这个空间可不是只占用一份哦。我们每一个进程，都有属于自己独立的虚拟内存地址空间。这也就意味着，每一个进程都需要这样一个页表。不管我们这个进程，是个本身只有几 KB 大小的程序，还是需要几 GB 的内存空间，都需要这样一个页表。如果你用的是 Windows，你可以打开你自己电脑上的任务管理器看看，现在你的计算机里同时在跑多少个进程，用这样的方式，页表需要占用多大的内存。

这还只是 32 位的内存地址空间，现在大家用的内存，多半已经超过了 4GB，也已经用上了 64 位的计算机和操作系统。这样的话，用上面这个数组的数据结构来保存页面，内存占用就更大了。那么，我们有没有什么更好的解决办法呢？你可以先仔细思考一下。

多级页表

仔细想一想，我们其实没有必要存下这 2^{20} 个物理页表啊。大部分进程所占用的内存是有限的，需要的页也自然是有限的。我们只需要去存那些用到的页之间的映射关系就好了。如果你对数据结构比较熟悉，你可能要说了，那我们是不是应该用哈希表（Hash Map）这样的数据结构呢？

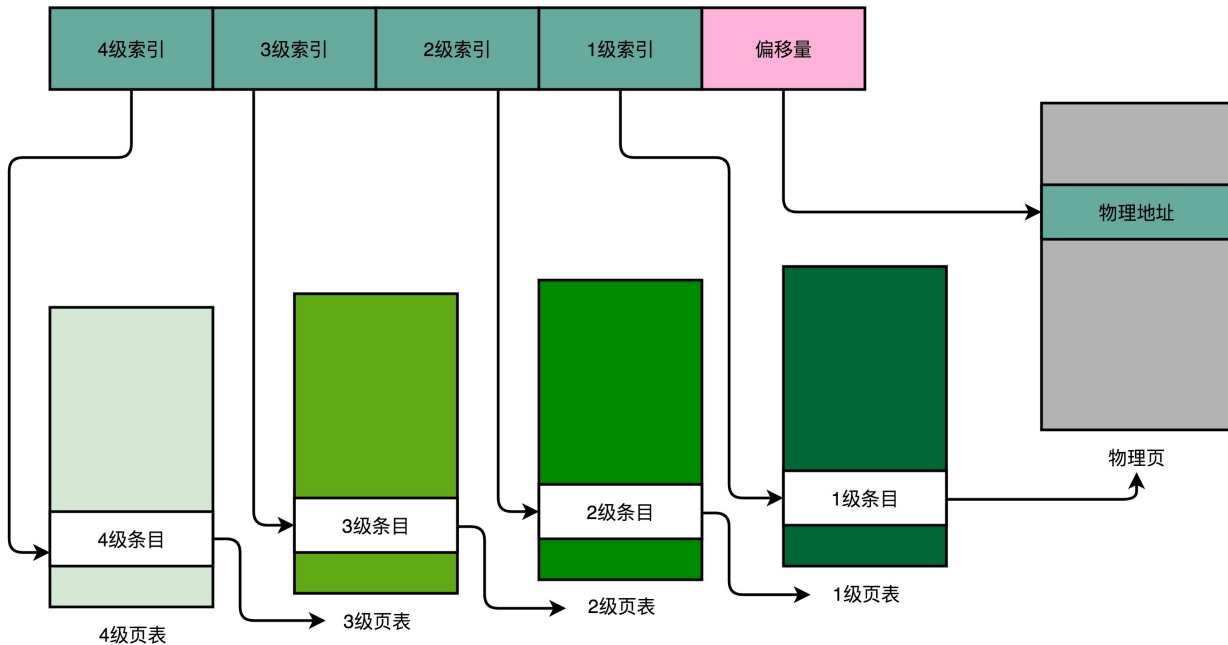
很可惜你猜错了：）。在实践中，我们其实采用的是一种叫作多级页表（Multi-Level Page Table）的解决方案。这是为什么呢？为什么我们不用哈希表而用多级页表呢？别着急，听我慢慢跟你讲。

我们先来看一看，一个进程的内存地址空间是怎么分配的。在整个进程的内存地址空间，通常是“两头实、中间空”。在程序运行的时候，内存地址从顶部往下，不断分配占用的栈的空间。而堆的空间，内存地址则是从底部往上，是不断分配占用的。

所以，在一个实际的程序进程里面，虚拟内存占用的地址空间，通常是两段连续的空间。而不是完全散落的随机的内存地址。而多级页表，就特别适合这样的内存地址分布。

我们以一个 4 级的多级页表为例，来看一下。同样一个虚拟内存地址，偏移量的部分和上面简单页表一样不变，但是原先的页号部分，我们把它拆成四段，从高到低，分成 4 级到 1

级这样 4 个页表索引。



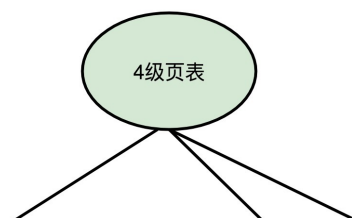
对应的，一个进程会有一个 4 级页表。我们先通过 4 级页表索引，找到 4 级页表里面对应的条目（Entry）。这个条目里存放的是一张 3 级页表所在的位置。4 级页面里面的每一个条目，都对应着一张 3 级页表，所以我们可能有多张 3 级页表。

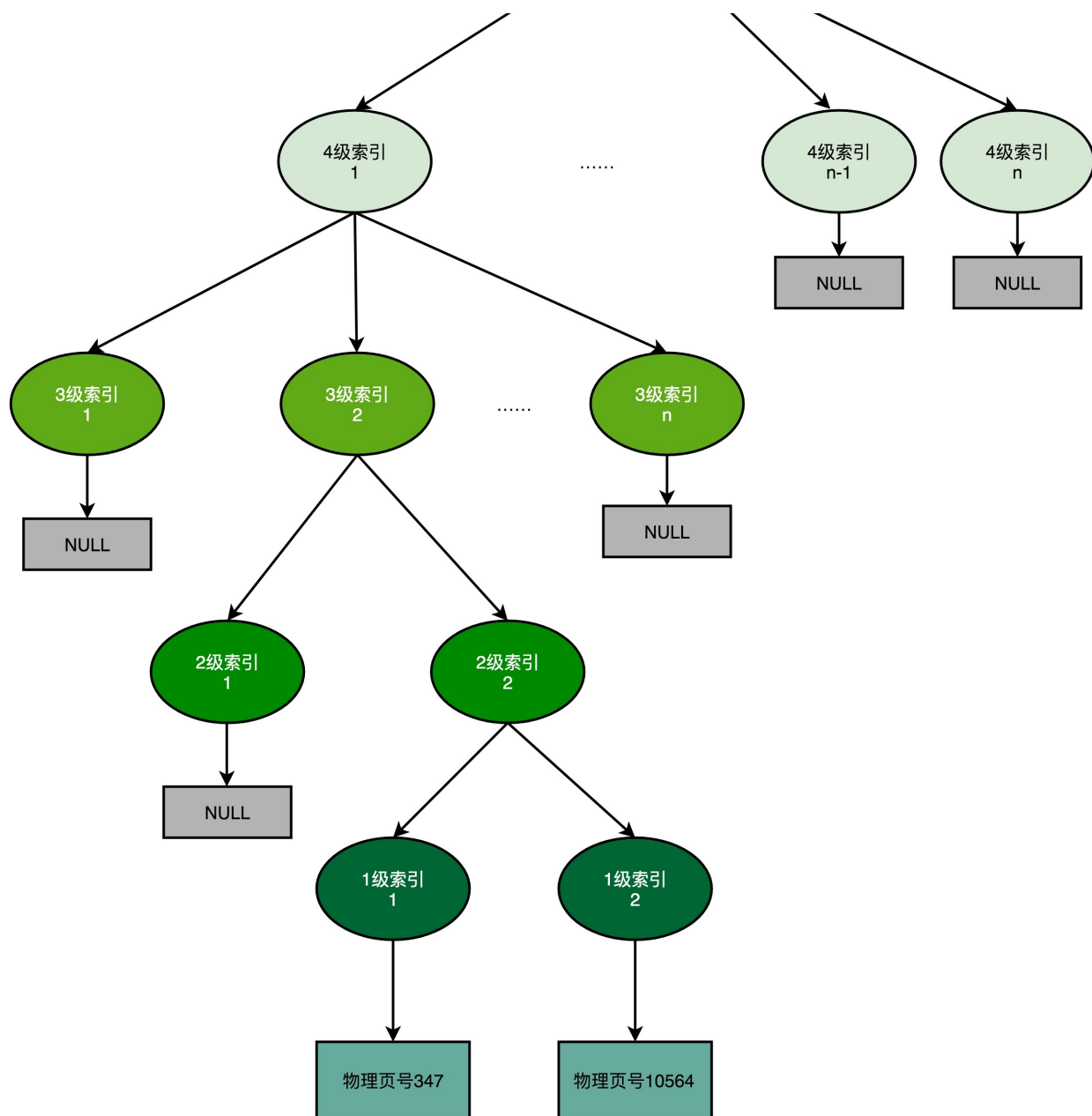
找到对应这张 3 级页表之后，我们用 3 级索引去找到对应的 3 级索引的条目。3 级索引的条目再会指向一个 2 级页表。同样的，2 级页表里我们可以用 2 级索引指向一个 1 级页表。

而最后一层的 1 级页表里面的条目，对应的数据内容就是物理页号了。在拿到了物理页号之后，我们同样可以用“页号 + 偏移量”的方式，来获取最终的物理内存地址。

我们可能有很多张 1 级页表、2 级页表，乃至 3 级页表。但是，因为实际的虚拟内存空间通常是连续的，我们很可能只需要很少的 2 级页表，甚至只需要 1 张 3 级页表就够了。

事实上，多级页表就像一个多叉树的数据结构，所以我们常常称它为**页表树**（Page Table Tree）。因为虚拟内存地址分布的连续性，树的第一层节点的指针，很多就是空的，也就不需要有对应的子树了。所谓不需要子树，其实就是不需要对应的 2 级、3 级的页表。找到最终的物理页号，就好像通过一个特定的访问路径，走到树最底层的叶子节点。





以这样的分成4级的多级页表来看，每一级如果都用5个比特表示。那么每一张某1级的页表，只需要 $2^5=32$ 个条目。如果每个条目还是4个字节，那么一共需要128个字节。而一个1级索引表，对应32个4KiB的也就是16KB的大小。一个填满的2级索引表，对应的就是32个1级索引表，也就是512KB的大小。

我们可以一起来测算一下，一个进程如果占用了1MB的内存空间，分成了2个512KB的连续空间。那么，它一共需要2个独立的、填满的2级索引表，也就意味着64个1级索引表，2个独立的3级索引表，1个4级索引表。一共需要69个索引表，每个128字节，大概就是9KB的空间。比起4MB来说，只有差不多1/500。

不过，多级页表虽然节约了我们的存储空间，却带来了时间上的开销，所以它其实是一个“以时间换空间”的策略。原本我们进行一次地址转换，只需要访问一次内存就能找到物理页

号，算出物理内存地址。但是，用了 4 级页表，我们就需要访问 4 次内存，才能找到物理页号了。

我们在前面两讲讲过，内存访问其实比 Cache 要慢很多。我们本来只是要做一个简单的地址转换，反而是一下子要多访问好多次内存。对于这个时间层面的性能损失，我们有没有什么更好的解决办法呢？那请你一定要关注下一讲的内容哦！

总结延伸

好了，这一讲的内容差不多了，我们来总结一下。

我们从最简单的进行虚拟页号——映射的简单页表说起，仔细讲解了现在实际应用的多级页表。多级页表就像是一颗树。因为一个进程的内存地址相对集中和连续，所以采用这种页表树的方式，可以大大节省页表所需要的空间。而因为每个进程都需要一个独立的页表，这个空间的节省是非常可观的。

在优化页表的过程中，我们可以观察到，数组这样的紧凑的数据结构，以及树这样稀疏的数据结构，在时间复杂度和空间复杂度的差异。另外，纯粹理论软件的数据结构和硬件的设计也是高度相关的。

推荐阅读

对于虚拟内存的知识点，你可以再深入读一读《计算机组成与设计：硬件 / 软件接口》的第 5.7 章节。如果你觉得还不过瘾，可以进一步去读一读《[What Every Programmer Should Know About Memory](#)》的第 4 部分，也就是 Virtual Memory。

[上一页](#)

[下一页](#)