

超越物理内存.机制

如何超越物理内存

到目前为止,我们一直假定地址空间非常小,能放入物理内存(页表存于内存中).事实上,我们假设每个正在运行的进程的地址空间都能放入内存.

现在假设我们需要支持许多同时运行的巨大地址空间.为了达到这个目的,需要在内存层级上再加一层(---> 硬盘)操作系统需要把当前没有在用的那部分地址空间找个地方存储起来

有了巨大的地址空间,不必担心程序的数据结构是否有足够空间存储,只需自然地编写程序,根据需要分配内存.这是操作系统提供的一个强大的假象.

一个反面例子是,一些早期系统使用"内存覆盖(memory overlays)",它需要程序员根据需要手动移入或移出内存中的代码或数据

为什么我们要为进程支持巨大的地址空间? 答案是方便和易用性

不仅是一个进程,增加交换空间让操作系统为多个并发运行的进程都提供巨大地址空间的假象.多道程序的出现,强烈要求能够换出一些页,因为早期的机器显然不能将所有进程需要的所有页都放在内存中.

多道程序和易用性都需要操作系统支持比物理内存更大的地址空间.这是所有现代虚拟内存系统都会做的事情,也是现在我们要进一步学习的内容

交换空间

要做的第一件事情就是在硬盘上开辟一部分空间用于物理页的移入和移出,这样的空间称为交换空间(swap space)

假设操作系统能够以页大小为单元读取或者写入交换空间,因为我们将内存中的页交换到其中,并在需要的时候又交换回去.为了达到这个目的,操作系统需要记住给定页的硬盘地址(disk address)

例子中,有一个4页的物理内存和一个8页的交换空间.进程 0,1,2 共享物理内存,并且3个都只有一部分有效页在内存中,剩下的在硬盘的交换空间中.进程 3的所有页都被交换到硬盘上,因此很清楚它目前没有运行.

需要注意的是,交换空间不是唯一的硬盘交换目的地.例如运行一个二进制程序(如ls或者你自己编译的 main程序),这个二进制程序的代码页最开始是在硬盘上.但程序运行的时候,它们被加载到内存中(要么在程序开始运行时全部加载,要么在现代操作系统中按需要一页一页加载).但是,如果系统需要在物理内存中腾出空间以满足其他需求,则可以安全地重新使用这些代码页的内存空间,因为稍后它又可以重新从硬盘上的二进制文件加载.

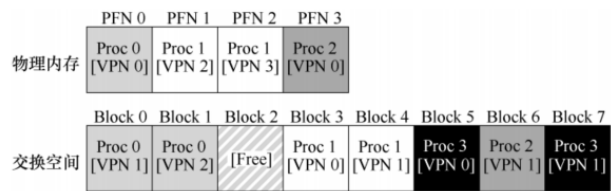


图 21.1 物理内存和交换空间

存在位

现在我们在硬盘上有一些空间,需要在系统中增加一些更高级的机制,来支持从硬盘交换页.简单起见,假设有一个硬件管理 TLB 的系统

内存访问时,如果TLB 未命中,则硬件在内存中查找页表.如果页有效且存在于物理内存中,则硬件从 PTE中获得 PFN,将其插入TLB,并重试该指令,这次产生 TLB 命中

如果希望允许页交换到硬盘,必须添加更多的机制.具体来说,当硬件在 PTE中查找时,可能发现页不在物理内存中.硬件操作系统判断是否在内存中的方法,是通过页表项中的一条新信息,即存在位(present bit)

如果存在位设置为1,则表示该页存在于物理内存中.如果存在位设置为0,则页不在内存中,而在硬盘上.访问不在物理内存中的页,这种行为通常被称为 page fault

page fault 时,操作系统被唤醒来处理 page fault.一段称为 page-fault handler 的代码会执行来处理

页错误 page fault

在 TLB 未命中的时,两种类型的系统:硬件管理 TLB和软件管理 TLB.不论在哪种系统中,如果页不存在,都由操作系统负责处理页错误.操作系统的 page-fault handler 确定要做什么.几乎所有的系统都在软件中处理页错误

如果一个页已被交换到硬盘,在处理 page fault 的时候,操作系统需要将该页交换到内存中.问题来了: 操作系统如何知道所需的页在哪儿? 在许多系统中,页表是存储这些信息最自然的地方.因此,操作系统可以用 PTE中的某些位来存储硬盘地址.这些位通常用来存储像页的 PFN这样的数据.当操作系统接收到页错误时,它会在 PTE中查找地址,并将请求发送到硬盘,将页读取到内存中.

从 TLB的经验中得知,硬件设计者不愿意信任操作系统做所有事情.那为什么他们相信操作系统来处理页错误呢?

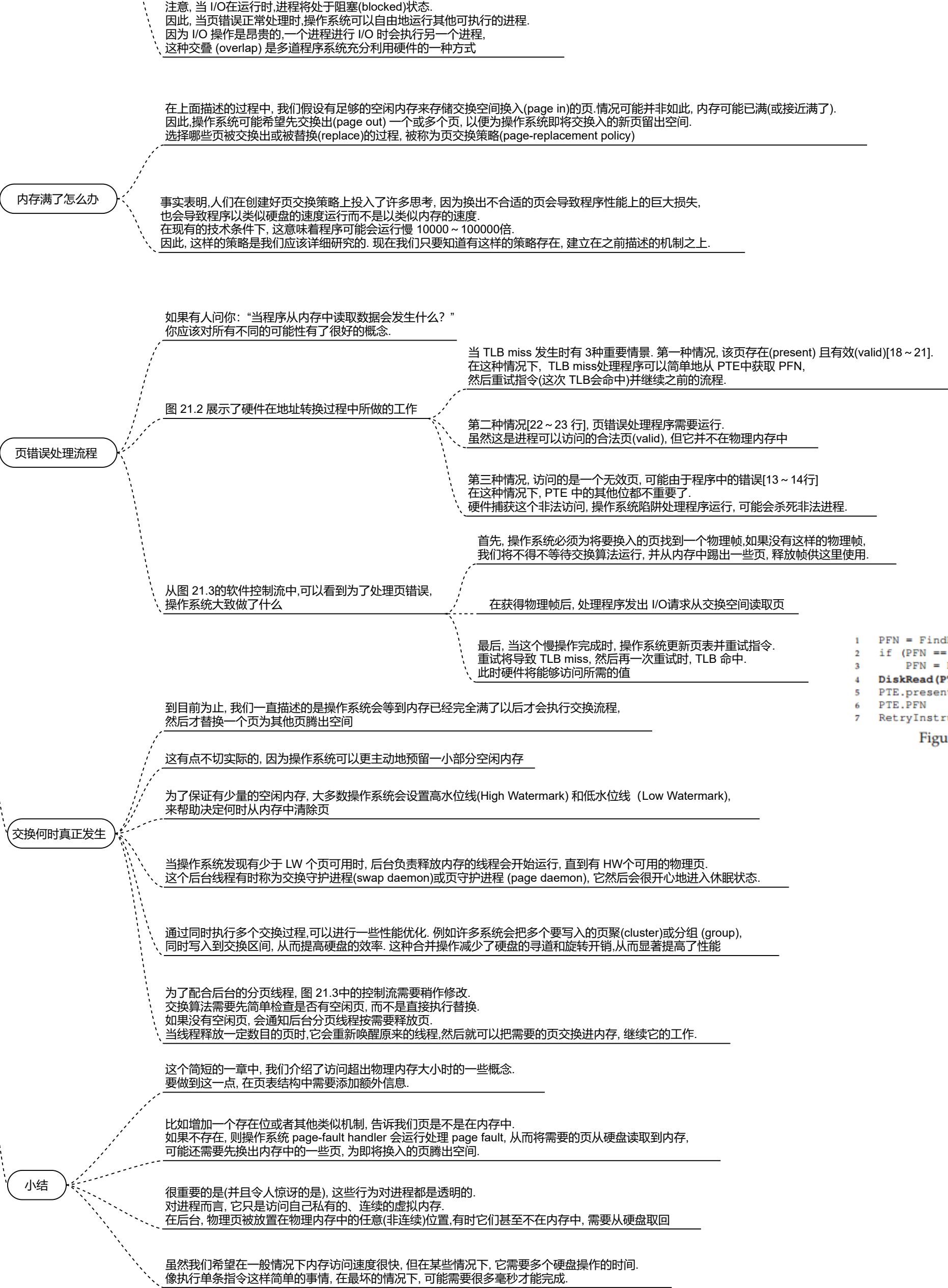
为什么硬件不能处理页错误

首先,页错误导致的硬盘操作很慢.即使操作系统需要很长时间来处理故障,执行大量的指令,但相比于硬盘操作,这些额外开销是很小的.

其次,为了能够处理 page fault,硬件必须了解交换空间,如何向硬盘发起 I/O 操作,以及很多它当前所不知道的细节.

因此,由于性能和简单的原因,操作系统来处理页错误,即使硬件人员也很开心.

当硬盘 I/O 完成时,操作系统会更新页表,将此页标记为存在,更新页表项的 PFN字段以记录新获取页的内存位置,并重试指令.下一次重新访问 TLB还是未命中,然而这次因为页在内存中,因此会将页表中的地址更新到 TLB中(也可以在处理页错误时更新 TLB以避免此步骤).最后的重试操作会在 TLB中找到转换映射,从已转换的内存物理地址,获取所需的数据或指令



```
1  VPN = (VirtualAddress & VPN_MASK) >> SHIFT
2  (Success, TlbEntry) = TLB_Lookup(VPN)
3  if (Success == True)    // TLB Hit
4      if (CanAccess(TlbEntry.ProtectBits) == True)
5          Offset  = VirtualAddress & OFFSET_MASK
6          PhysAddr = (TlbEntry.PFN << SHIFT) | Offset
7          Register = AccessMemory(PhysAddr)
8      else
9          RaiseException(PROTECTION_FAULT)
10 else    // TLB Miss
11     PTEAddr = PTBR + (VPN * sizeof(PTE))
12     PTE = AccessMemory(PTEAddr)
13     if (PTE.Valid == False)
14         RaiseException(SEGMENTATION_FAULT)
15     else
16         if (CanAccess(PTE.ProtectBits) == False)
17             RaiseException(PROTECTION_FAULT)
18         else if (PTE.Present == True)
19             // assuming hardware-managed TLB
20             TLB_Insert(VPN, PTE.PFN, PTE.ProtectBits)
21             RetryInstruction()
22         else if (PTE.Present == False)
23             RaiseException(PAGE_FAULT)
```

Figure 21.2: Page-Fault Control Flow Algorithm (Hardware)

```
1  PFN = FindFreePhysicalPage()
2  if (PFN == -1)    // no free page found
3      PFN = EvictPage()    // run replacement algorithm
4  DiskRead(PTE.DiskAddr, PFN)    // sleep (waiting for I/O)
5  PTE.present = True    // update page table with present
6  PTE.PFN     = PFN     // bit and translation (PFN)
7  RetryInstruction()    // retry instruction
```

Figure 21.3: Page-Fault Control Flow Algorithm (Software)