



# OS Lab 4

page tables

王世烜

PB20151796

November 25, 2022

## Part 1: 实验要求

在本实验中，我们将探索页表并对其进行修改，以加快某些系统调用并检测哪些页面已经被访问。

## Part 2: Print a page table

### 2.1 问题简述

本部分要求我们可视化 RISC-V 页表，写一个函数输出页表的内容。

### 2.2 知识基础

一个 RISC-V 页表在逻辑上是一个  $2^{27}$  (134, 217, 728) **页表项** (Page Table Entry, PTE) 的数组。每个 PTE 包含一个 44 位的**物理页号** (Physical Page Number, PPN) 和一些**标志位**。分页硬件通过利用 39 位中的高 27 位索引到页表中找到一个 PTE 来转换一个虚拟地址，并计算出一个 56 位的物理地址，它的前 44 位来自于 PTE 中的 PPN，而它的后 12 位则是从原来的虚拟地址复制过来的，在逻辑上可以把页表看成是一个简单的 PTE 数组。

如 **图 1** 所示，实际转换分三步进行。一个页表以**三层树**的形式存储在物理内存中。树的根部是一个 4096 字节的页表页，它包含 512 个 PTE，这些 PTE 包含树的下一级页表页的物理地址。每一页都包含 512 个 PTE，用于指向下一个页表或物理地址。分页硬件用 27 位中的顶 9 位选择根页表页中的 PTE，用中间 9 位选择树中下一级页表页中的 PTE，用底 9 位选择最后的 PTE。

每个 PTE 包含**标志位**，告诉分页硬件如何允许使用相关的虚拟地址。**PTE\_V** 表示 **PTE** 是否存在：如果没有设置，对该页的引用会引起异常（即不允许）。**PTE\_R** 控制是否允许指令读取到页。**PTE\_W** 控制是否允许指令向写该页。**PTE\_X** 控制 CPU 是否可以将页面的内容解释为指令并执行。**PTE\_U** 控制是否允许用户模式下的指令访问页面；如果不设置 **PTE\_U**，PTE 只能在监督者模式下使用。

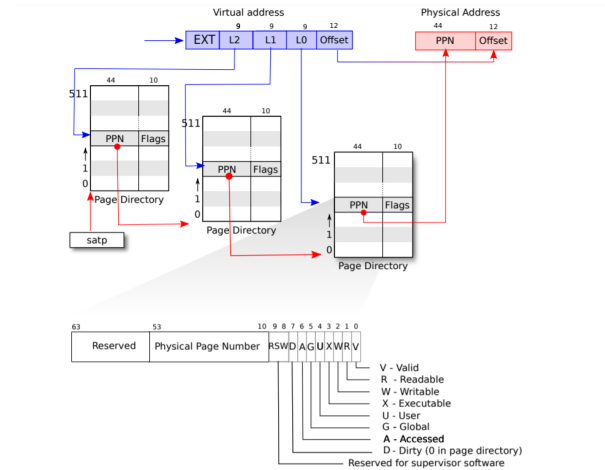


Figure 3.2: RISC-V address translation details.

图 1: RISC-V address translation details.

## 2.3 实验步骤

### (1) 分析 freewalk() 函数

根据所给提示，首先分析 kernel/vm.c 中的 freewalk() 函数：

```

1 // Recursively free page-table pages.
2 // All leaf mappings must already have been removed.
3 void
4 freewalk(pagetable_t pagetable)
5 {
6     // there are 2^9 = 512 PTEs in a page table.
7     for(int i = 0; i < 512; i++){
8         pte_t pte = pagetable[i];
9         if((pte & PTE_V) && (pte & (PTE_R|PTE_W|PTE_X)) == 0){
10             // this PTE points to a lower-level page table.
11             uint64 child = PTE2PA(pte);
12             freewalk((pagetable_t)child);
13             pagetable[i] = 0;
14         } else if(pte & PTE_V){
15             panic("freewalk: leaf");
16         }
17     }
18     kfree((void*)pagetable);
19 }

```

freewalk() 函数通过递归调用来遍历所有的页表，并通过标志位来判断是否到达树的叶子结点。PTE2PA() 用于将逻辑地址转化为物理地址。

## (2) 编写 vmprint() 函数

仿照 freewalk() 函数可以编写出如下 vmprint() 函数：

```
1 void
2 _vmprint(pagetable_t pagetable, int level)
3 {
4     for (int i = 0; i < 512; i++)
5     {
6         pte_t pte = pagetable[i];
7         if (pte & PTE_V) // 页表项有效
8         {
9             uint64 child = PTE2PA(pte); // 转换到物理地址
10            for (int j = 0; j < level; j++)
11            {
12                if(j) printf(" ");
13                printf("..");
14            }
15            printf("%d: %pte %p pa %p\n", i, pte, child);
16
17            if((pte & (PTE_R|PTE_W|PTE_X)) == 0)
18            {
19                _vmprint((pagetable_t)child, level + 1);
20            }
21        }
22    }
23
24 }
25
26 void
27 vmprint(pagetable_t pagetable)
28 {
29     printf("page table %p\n", pagetable);
30     _vmprint(pagetable, 1);
31     return;
32 }
```

其中 level 用于记录到达第几层页表，以确定要输出几个..。

## (3) 添加声明

按照提示，在 kernel/exec.c 中的 return argc; 前添加 if(p->pid==1) vmprint(p->pagetable);:

```
1 if(p->pid==1) vmprint(p->pagetable);
2 return argc; // this ends up in a0, the first argument to main(argc, argv)
```

在 kernel/def.hs 中添加 vmprint() 函数的声明，使得 exec.c 可以调用 vmprint() 函数：

```

1 // vm.c
2 ...
3 void                vmprint(pagetable_t);

```

#### (4) 结果

上述操作完成后可以得到以下结果：

```

make[1]: 离开目录"/home/ubuntu/桌面/xv6-labs-202
== Test pte printout ==
$ make qemu-gdb
pte printout: OK (3.6s)

```

图 2: printout passed.

#### (5) 问题

Explain the output of vmprint in terms of Fig 3-4 from the text. What does page 0 contain? What is in page 2? When running in user mode, could the process read/write the memory mapped by page 1?

图 3-4 是哪个？是指 xv6book 中的图么？如果是，page0 又指什么？

## Part 3: A kernel page table per process

### 3.1 问题简述

本部分要求我们为每个进程新增一个内核态的页表，然后在该进程进入到内核态时，不使用公用的内核态页表，而是使用进程的内核态页表。

### 3.2 问题分析

### 3.3 实验步骤

#### (1) 添加内核页表字段：

在 proc 结构体中增加一个字段，用于进程的内核页表。

```

1 // Per-process state
2 struct proc {
3 ...
4     pagetable_t kpagetable;        // Kernel page table
5 ...
6 };

```

#### (2) 初始化内核页表：

仿照 kernel/vm.c 里 kvmmap() 与 kvminit() 函数, 重新写一个映射函数与一个初始化函数

```
1 void
2 ukvmmap(pagetable_t kpagetable, uint64 va, uint64 pa, uint64 sz, int perm)
3 {
4     if(mappages(kpagetable, va, sz, pa, perm) != 0)
5         panic("ukvmmap");
6 }
7
8 pagetable_t
9 ukvminit()
10 {
11     pagetable_t kpagetable = (pagetable_t) kalloc();
12     memset(kpagetable, 0, PGSIZE);
13     ukvmmap(kpagetable, UART0, UART0, PGSIZE, PTE_R | PTE_W);
14     ukvmmap(kpagetable, VIRTIO0, VIRTIO0, PGSIZE, PTE_R | PTE_W);
15     ukvmmap(kpagetable, CLINT, CLINT, 0x10000, PTE_R | PTE_W);
16     ukvmmap(kpagetable, PLIC, PLIC, 0x400000, PTE_R | PTE_W);
17     ukvmmap(kpagetable, KERNBASE, KERNBASE, (uint64)etext-KERNBASE,
18     PTE_R | PTE_X);
19     ukvmmap(kpagetable, (uint64)etext, (uint64)etext,
20     PHYSTOP-(uint64)etext, PTE_R | PTE_W);
21     ukvmmap(kpagetable, TRAMPOLINE, (uint64)trampoline, PGSIZE, PTE_R | PTE_X);
22     return kpagetable;
23 }
```

然后, 在 allocproc 中调用该函数:

```
1 static struct proc*
2 allocproc(void)
3 {
4     ...
5     // An empty user page table.
6     p->pagetable = proc_pagetable(p);
7     if(p->pagetable == 0){
8         freeproc(p);
9         release(&p->lock);
10        return 0;
11    }
12    // 下面是新添加的内容
13    // An empty user kernel page table.
14    p->kpagetable = ukvminit(p);
15    if(p->kpagetable == 0){
```

```

16     freeproc(p);
17     release(&p->lock);
18     return 0;
19 }
20 ...
21 }

```

并在 kernel/defs.h 添加函数声明: `pagetable_t ukvminit(void)`; 以及 `void ukvmmmap(pagetable_t, uint64, uint64, uint64, int)`;

### (3) 初始化内核栈:

根据提示, 接下来要将 `procinit()` 函数中部分或全部功能转移到 `allocproc()` 函数。首先观察 `procinit()` 函数:

```

1 // initialize the proc table at boot time.
2 void
3 procinit(void)
4 {
5     struct proc *p;
6
7     initlock(&pid_lock, "nextpid");
8     for(p = proc; p < &proc[NPROC]; p++) {
9         initlock(&p->lock, "proc");
10
11         // Allocate a page for the process's kernel stack.
12         // Map it high in memory, followed by an invalid
13         // guard page.
14         char *pa = kalloc();
15         if(pa == 0)
16             panic("kalloc");
17         uint64 va = KSTACK((int) (p - proc));
18         kvmmmap(va, (uint64)pa, PGSIZE, PTE_R | PTE_W);
19         p->kstack = va;
20     }
21     kvminithart();
22 }

```

该函数是用于将每个进程的内核页表映射到一个针对该进程的内核堆栈。

于是转移到 `allocproc()` 函数, 用上之前写好的 `ukvmmmap()` 函数以初始化内核堆栈:

```

1 // An empty user kernel page table.
2 p->kpagetable = ukvminit(p);
3 if(p->pagetable == 0){
4     freeproc(p);
5     release(&p->lock);

```

```

6     return 0;
7 }
8
9 // 初始化内核栈
10 char *pa = kalloc();
11 if(pa == 0)
12     panic("kalloc");
13 uint64 va = KSTACK((int) (p - proc));
14 ukvmmap(p->kpagetable, va, (uint64)pa, PGSIZE, PTE_R | PTE_W);
15 p->kstack = va;

```

#### (4) 修改 scheduler() 函数，在进程调度时，切换内核页：

按照提示，仿照 kvminithart() 函数，修改 scheduler() 函数将进程的内核页表加载到内核的 satp，并且寄存器中在没有进程运行时使用 kernel\_pagetable

```

1 void
2 scheduler(void)
3 {
4     ...
5     for(;;){
6         // Avoid deadlock by ensuring that devices can interrupt.
7         intr_on();
8
9         int found = 0;
10        for(p = proc; p < &proc[NPROC]; p++) {
11            acquire(&p->lock);
12            if(p->state == RUNNABLE) {
13                // Switch to chosen process. It is the process's job
14                // to release its lock and then reacquire it
15                // before jumping back to us.
16                p->state = RUNNING;
17                c->proc = p;
18
19                // load the process's kernel page table into the core's satp
20                // register
21                w_satp(MAKE_SATP(p->kpagetable));
22                sfence_vma();
23
24                swtch(&c->context, &p->context);
25
26                // Process is done running for now.
27                // It should have changed its p->state before coming back.
28                c->proc = 0;

```

```

29
30     kvminithart(); // use kernel_pagetable when no process is running.
31     ...
32 }

```

#### (5) 在 freeproc() 函数中释放一个进程的内核页表:

释放一个进程的内核页表首先要释放页表内的内核栈, 因为页表内存储的内核栈地址本身就是一个虚拟地址, 需要先将这个地址指向的物理地址进行释放:

```

1  if (p->kstack)
2  {
3  pte_t * pte = walk(p->pagetable, p->kstack, 0);
4  kfree((void*)PA2PTE(*pte));
5  }
6  p->kstack = 0;

```

其中根据 xv6 book 介绍, walk() 函数模仿 RISC-V 分页硬件查找虚拟地址的 PTE。walk 每次降低 3 级页表的 9 位。它使用每一级的 9 位虚拟地址来查找下一级页表或最后一级的 PTE。

需要注意的是, 这里需要在 kernel/defs.h 中添加 walk 函数的声明, 否则无法直接引用。然后是释放页表, 直接遍历所有的页表, 释放所有有效的页表项即可。仿照 freewalk 函数。由于 freewalk 函数将对应的物理地址也直接释放了, 我们这里释放的内核页表仅仅只是用户进程的一个备份, 释放时仅释放页表的映射关系即可, 不能将真实的物理地址也释放了。因此不能直接调用 freewalk 函数, 而是需要进行更改:

```

1  void
2  proc_freewalk(pagetable_t pagetable)
3  {
4      for (int i = 0; i < 512; i++) {
5          pte_t pte = pagetable[i];
6          if (pte & PTE_V) {
7              pagetable[i] = 0;
8              if ((pte & (PTE_R | PTE_W | PTE_X)) == 0) {
9                  uint64 child = PTE2PA(pte);
10                 proc_freewalk((pagetable_t)child);
11             }
12         }
13     }
14     kfree((void*)pagetable);
15 }

```

再在 freeproc 中进行调用:



```

1 // delete kernel pagetable
2 if(p->kpagetable) {
3   proc_freewalk(p->kpagetable);
4 }
5 p->kpagetable = 0;

```

#### (6) 切换进程的内核页表:

最后, 在 vm.c 中添加头文件:

```

1 #include "spinlock.h"
2 #include "proc.h"

```

然后更改 kvmpa 函数 (否则会报错: "panic:kvmpa"):

```

1 pte = walk(myproc()->kpagetable, va, 0);

```

```

      qemu output saved to xv6.out.log
== Test usertests ==
$ make qemu-gdb
(285.6s)
== Test   usertests: copyin ==
   usertests: copyin: OK
== Test   usertests: copyinstr1 ==
   usertests: copyinstr1: OK
== Test   usertests: copyinstr2 ==
   usertests: copyinstr2: OK
== Test   usertests: copyinstr3 ==
   usertests: copyinstr3: OK
== Test   usertests: sbrkmuch ==
   usertests: sbrkmuch: OK
== Test   usertests: all tests ==
   usertests: all tests: OK

```

图 3: usertests passed

## Part 4: Simplify copyin/copyinstr

### 4.1 问题简述

该部分主要目的是将用户进程页表的所有内容都复制到内核页表中.

### 4.2 实验步骤

(1) 替换 copyin(),copyinstr() 函数:

根据提示, 将 vm.c 中 copyin(),copyinstr() 函数的主体内容替换为已经写好的 vm-copyin.c 中的 copyin\_new(),copyinstr\_new() 函数:

```

1  int
2  copyin(pagetable_t pagetable, char *dst, uint64 srcva, uint64 len)
3  {
4      return copyin_new(pagetable, dst, srcva, len);
5  }
6
7  int
8  copyinstr(pagetable_t pagetable, char *dst, uint64 srcva, uint64 max)
9  {
10     return copyinstr_new(pagetable, dst, srcva, max);
11 }

```

并在 defs.h 中添加声明:

```

1  // vmcopyin.c
2  int          copyin_new(pagetable_t, char *, uint64, uint64);
3  int          copyinstr_new(pagetable_t, char *, uint64, uint64);

```

(2) 修改 fork(),exec(),sbrk() 函数:

在 fork() 函数中, 需要将父进程的 pagetable 复制给子进程的 kpagetable。而 fork() 函数中使用了 vm.c 中的 uvmcopy() 函数来将父进程的 pagetable 复制给子进程的 pagetable, 所以我们要仿照此函数, 写出一个函数 u2kvmcopy() 用来将父进程的 pagetable 复制给子进程的 kpagetable。

```

1  void
2  u2kvmcopy(pagetable_t pagetable, pagetable_t kpagetable, uint64 oldsz,
3  uint64 newsz)
4  {
5      pte_t *pte_from, *pte_to;
6      uint64 a, pa;
7      uint flags;
8
9      if (newsz < oldsz)
10         return;
11     oldsz = PGROUNDUP(oldsz);
12     for (a = oldsz; a < newsz; a += PGSIZE)
13     {
14         if ((pte_from = walk(pagetable, a, 0)) == 0)
15             panic("u2kvmcopy: pte should exist");
16         if ((pte_to = walk(kpagetable, a, 1)) == 0)
17             // copy the pte to the same address at kernel page table
18             panic("u2kvmcopy: walk fails");
19         pa = PTE2PA(*pte_from);

```

```

20     flags = (PTE_FLAGS(*pte_from) & (~PTE_U));
21     *pte_to = PA2PTE(pa) | flags;
22 }
23 }

```

这个函数与 `uvmcopy()` 函数不同的地方在于 `PTE_U`，由第四条提示，一个设置了 `PTE_U` 的页不能在内核模式下被访问，所以将在 line14 中将这个标志位去掉。

然后修改 `fork()` 函数，实现上述目标：

```

1  int
2  fork(void)
3  {
4  ...
5      u2kvmcopy(np->pagetable, np->kpagetable, 0, np->sz);
6  ...
7  }

```

`exec()` 函数在 `kernel/exec.c` 中，在执行新的程序前，初始化之后，进行页表拷贝：

```

1  int
2  exec(char *path, char **argv)
3  {
4  ...
5      // Load program into memory.
6      for(i=0, off=elf.phoff; i<elf.phnum; i++, off+=sizeof(ph)){
7  ...
8          if((sz1 = uvmmalloc(pagetable, sz, ph.vaddr + ph.memsz)) == 0)
9              goto bad;
10         // 添加检测，防止程序大小超过 PLIC
11         if(sz1 >= PLIC)
12             goto bad;
13         sz = sz1;
14     ...
15     }
16     iunlockput(ip);
17     ...
18     stackbase = sp - PGSIZE;
19
20     u2kvmcopy(pagetable, p->kpagetable, 0, sz);
21     // Push argument strings, prepare rest of stack in ustack.
22     for(argc = 0; argv[argc]; argc++) {
23     ...
24     }

```

sbrk() 函数是系统调用，位于 sysproc.c 中，具体内容如下：

```
1 uint64
2 sys__sbrk(void)
3 {
4     int addr;
5     int n;
6
7     if(argint(0, &n) < 0)
8         return -1;
9     addr = myproc()->sz;
10    if(growproc(n) < 0)
11        return -1;
12    return addr;
13 }
```

可见其调用了 growproc() 函数，下面对于该函数进行修改。

```
1 int
2 growproc(int n)
3 {
4     uint sz;
5     struct proc *p = myproc();
6
7     sz = p->sz;
8     if(n > 0){
9         if((sz = uvmmalloc(p->pagetable, sz, sz + n)) == 0) {
10             return -1;
11         }
12         // 添加
13         u2kvmcopy(p->pagetable, p->kpagetable, sz-n, sz);
14     } else if(n < 0){
15         sz = uvmmdealloc(p->pagetable, sz, sz + n);
16     }
17     p->sz = sz;
18     return 0;
19 }
```

最后，在 userinit 中把第一个进程的用户页表包含在它的内核页表中。

```
1 void
2 userinit(void)
3 {
4     struct proc *p;
5 }
```

```

6  p = allocproc();
7  initproc = p;
8
9  // allocate one user page and copy init's instructions
10 // and data into it.
11 uvminit(p->pagetable, initcode, sizeof(initcode));
12 p->sz = PGSIZE;
13 // 添加
14 u2kvmcopy(p->pagetable, p->kpagetable, 0, p->sz);
15
16 // prepare for the very first "return" from kernel to user.
17 p->trapframe->epc = 0;      // user program counter
18 p->trapframe->sp = PGSIZE;  // user stack pointer
19
20 safestrcpy(p->name, "initcode", sizeof(p->name));
21 p->cwd = namei("/");
22
23 p->state = RUNNABLE;
24
25 release(&p->lock);
26 }

```

## Part5 Question

1、Explain why the third test  $srcva + len < srcva$  is necessary in `copyin_new()`: give values for `srcva` and `len` for which the first two test fail (i.e., they will not cause to return -1) but for which the third one is true (resulting in returning -1).

实际上这样操作就是为了防止溢出。xv6book 中有提到

For example if  $(ph.vaddr + ph.memsz < ph.vaddr)$  checks for whether the sum overflows a 64-bit integer.

例如：

```

1  #include <iostream>
2  using namespace std;
3  int main()
4  {
5      uint64_t srcva = UINT64_MAX / 2;
6      uint64_t len = UINT64_MAX - srcva + 1;
7      cout << srcva + len << "□" << srcva << "□" << len << endl;
8      return 0;
9  }

```

输出结果

```
PS C:\wsd\vscode\code> g++ 1.cpp -o 1.exe
$?) { &'./1.exe' }
0 9223372036854775807 9223372036854775809
PS C:\wsd\vscode\code>
```

图 4: Result

## Part 6: 实验结果

运行 make grade, 检验结果:

```
riscv64-linux-gnu-objdump -t kernel/kernel | sed '1,/SYMBOL TABLE/d; s/ .* / /; /^$/d' > 1
make[1]: 离开目录"/home/ubuntu/桌面/xv6-labs-2020-1"
== Test pte printout ==
$ make qemu-gdb
pte printout: OK (4.2s)
(Old xv6.out.pteprint failure log removed)
== Test answers-pgtbl.txt == answers-pgtbl.txt: OK
== Test count copyin ==
$ make qemu-gdb
count copyin: OK (2.0s)
(Old xv6.out.count failure log removed)
== Test userests ==
$ make qemu-gdb
(298.1s)
== Test userests: copyin ==
userests: copyin: OK
== Test userests: copyinstr1 ==
userests: copyinstr1: OK
== Test userests: copyinstr2 ==
userests: copyinstr2: OK
== Test userests: copyinstr3 ==
userests: copyinstr3: OK
== Test userests: sbrkmuch ==
userests: sbrkmuch: OK
== Test userests: all tests ==
userests: all tests: OK
== Test time ==
time: OK
Score: 66/66
ubuntu@VM5153-OSVM:~/桌面/xv6-labs-2020-1$
```

图 5: Result

本次实验完成了 Page Table 有关功能的实现和实际应用。在实验过程中得到了以下收获:

- 对于多级页表这个知识点有了深刻的认识
- 对于虚拟地址和物理地址的转换以及内存结构有了切身的感受
- 认识到了操作系统内存设计合理性以及正确性