

一、页表是如何翻译地址的？(va->pa)

软件翻译：walkaddr(pagetable,va)

何处使用？copyin()或copyout()或copyinstr()

1-kernel/vm.c：‘copyin()’或‘copyout()’或‘copyinstr()’

```
1  int copyin(pagetable_t pagetable, char *dst, uint64 srcva, uint64 len);
2  int copyinstr(pagetable_t pagetable, char *dst, uint64 srcva, uint64 max);
3  int copyout(pagetable_t pagetable, uint64 dstva/* 目标虚拟地址*/, char *src,
4  uint64 len)
5  {
6      uint64 n, va0, pa0;
7      while(len > 0){
8          va0 = PGROUNDDOWN(dstva); // 将目标va向下取整，获得va所在的虚拟页的起始地址
9          pa0 = walkaddr(pagetable, va0); // 获取该虚拟页 对应的物理页的起始地址
10         ↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
11         // 下略....
12     }
13     return 0;
14 }
```

2-kernel/vm.c->kernel/vm.c：‘walkaddr()’

```
1  uint64
2  walkaddr(pagetable_t pagetable, uint64 va)
3  {
4      pte_t *pte;
5      uint64 pa;
6
7      if(va >= MAXVA) // 检查虚拟地址是否超出可用范围
8          return 0;
9
10     pte = walk(pagetable, va, 0); // 获取va所对应的虚拟页的pte的地址(回看PPT)
11     ↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
12     if(pte == 0)
13         return 0;
14     if((*pte & PTE_V) == 0)
15         return 0;
16     if((*pte & PTE_U) == 0)
17         return 0;
18     pa = PTE2PA(*pte); // *pte是pte条目的具体内容，根据pte的具体内容找到目标物理页面的
19     起始地址 ↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
20     return pa;
21 }
```

2-1kernel/vm.c->kernel/vm.c：‘walk()’

```

1 // The risc-v Sv39 scheme has three levels of page-table
2 // pages. A page-table page contains 512 64-bit PTEs.
3 // A 64-bit virtual address is split into five fields:
4 //   39..63 -- must be zero.
5 //   30..38 -- 9 bits of level-2 index. 高级页目录
6 //   21..29 -- 9 bits of level-1 index. 中级页目录
7 //   12..20 -- 9 bits of level-0 index. 低级页目录
8 //   0..11 -- 12 bits of byte offset within the page. 页内偏移量
9 pte_t * // 获取va所对应的虚拟页的pte的地址
10 walk.pagetable_t pagetable, uint64 va, int alloc)
11 {
12     if(va >= MAXVA) // 先检查va是否越界
13         panic("walk");
14
15     for(int level = 2; level > 0; level--) { // 再访问三级页表的前两级，从高级页目录
16         // 开始访问(level=2) (结合PPT)
17         pte_t *pte = &pagetable[PX(level, va)]; // 获取下一级页目录的映射信息(pte)所
18         // 在的地址 ↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
19         if(*pte & PTE_V) {
20             pagetable = (pagetable_t)PTE2PA(*pte); // 根据pte中的信息，获取下一级页目录
21             // 的所在地址，并将其作为新的页表↓↓↓↓↓↓↓
22         } else { // 若va还没被映射到内存中，则现场创建映射
23             if(!alloc || (pagetable = (pde_t*)kalloc()) == 0)
24                 return 0;
25             memset(pagetable, 0, PGSIZE);
26             *pte = PA2PTE(pagetable) | PTE_V;
27         }
28     }
29     return &pagetable[PX(0, va)]; // 这里的pagetable是最低级的页目录
30 }

```

2-1-1kernel/vm.c->kernel/riscv.h: 'PX()'和'PTE2PA()'

```

1 //看下图1
2 #define PGSHIFT 12 // 页内偏移量是12bit
3 #define PXMASK 0x1FF // 由于每级页目录中有512(2^9)个条目，所以每级页目录号
4 // 有9 bits,0x 1FF=1 1111 1111
5 #define PXSHIFT(level) (PGSHIFT+(9*(level))) // 计算取出第level级页目录号时，需
6 // 要右移多少位
7 #define PX(level, va) (((uint64)(va)) >> PXSHIFT(level)) & PXMASK // 先右
8 // 移，此时第level级页目录号就在最低9bit，再& 0x 1FF 操作相当于取出最低的9bit，其他bit全部
9 // 清零
10 // 看下图2
11 // pte中的最低10bit全是flag位 0x3FF= 11 1111 1111
12 #define PTE_FLAGS(pte) ((pte) & 0x3FF) // & 0x3FF 操作相当于取出最低的10bit，其
13 // 他bit全部清零
14 #define PA2PTE(pa) (((uint64)pa) >> 12) << 10 // 获取一个物理地址后，先右移
15 // 12bit(抹掉页内偏移量)，再左移10bit 给flag留空间

```

9bit 高级页目录号	9bit 中级页目录号	9bit 低级页目录号	12bit 页内偏移量
-------------	-------------	-------------	-------------

例：	公司编号：1	部门编号：3	小组编号：0	组内成员编号：2
----	--------	--------	--------	----------

27bit: 物理页号	10bit: flag
-------------	-------------

3-kernel/vm.c->kernel/riscv.h: 'PTE2PA()'

```
1 // 看上图
2 #define PTE2PA(pte) (((pte) >> 10) << 12) // 获取一个pte后，先右移10bit把flag抹掉，得到物理页框号，再左移12bit得到该物理页的起始物理地址
```

硬件翻译：MMU

何处使用？：正常访问数据(va)时，硬件就会隐式工作，比软件翻译快很多！

二、内存布局：kernel/memlayout.h(结合PPT)

QEMU的物理内存布局

- 00001000 -- boot ROM：QEMU提供的启动ROM，用于启动过程。
- 02000000 -- CLINT：本地中断控制器（Core Local Interruptor），包含定时器等功能。
- 0C000000 -- PLIC：可编程中断控制器（Platform-Level Interrupt Controller），用于管理来自外设的中断。
- 10000000 -- uart0：UART 0设备的寄存器地址，用于串行通信。
- 10001000 -- virtio disk：VirtIO磁盘设备的内存映射I/O接口。
- 80000000 -- kernel entry/kernel base：QEMU的boot ROM跳转至此地址以启动内核，kernel 选项加载的内核镜像也放置在这里。
- unused RAM after 80000000：0x80000000 之后是未使用的RAM，供内核和用户程序使用，xv6只允许后续的128MB内存空间。

xv6的物理内存使用

- 80000000 -- entry.S, then kernel text and data：内核的入口点（entry.S）以及内核的代码和数据部分。
- end -- start of kernel page allocation area：end 符号之后是内核开始分配页面的区域。
- PHYSTOP -- end RAM used by the kernel：PHYSTOP 定义了内核使用的RAM的结束地址，通常设置为 KERNBASE + 128MB（假设有128MB的RAM供内核使用）。

三、xv6是如何跑起来的？

1-kernel/entry.s (结合PPT)

```
1  # qemu -kernel loads the kernel at 0x80000000 and causes each CPU to jump
   there.
2  # kernel.ld causes the following code to be placed at 0x80000000.
3  .section .text # 这行代码指示接下来的代码应该被放置在程序的文本（代码）段中。
4  _entry: # 这是程序的入口点标签，QEMU或其他引导程序会将CPU的PC设置到这个地址
   (0x80000000)，然后从这里开始执行。
5
6      la sp, stack0 # 将stack0的地址加载到堆栈指针（sp）寄存器中。stack0是在start.c中
   定义的，用于存放每个CPU的堆栈基地址。
7      li a0, 1024*4 # 将4096（即1024*4）加载到寄存器a0中，这里4096是每个CPU堆栈的大小
   （以字节为单位）。
8      csrr a1, mhartid # 从CSR（控制和状态寄存器）中读取mhartid的值到寄存器a1。
   mhartid是硬件线程ID，用于区分多核处理器中的不同CPU核心。
9      addi a1, a1, 1 # cpu的数量=mhartid+1
10     mul a0, a0, a1 # 将a0（堆栈大小）和a1相乘，结果存放在a0中。这一步计算了每个CPU核
   心堆栈的偏移量。
11     add sp, sp, a0 # 将sp（堆栈基地址）和a0（堆栈偏移量）相加，得到当前CPU核心的堆栈顶
   地址，并更新sp寄存器。
12     # jump to start() in start.c
13     call start # ↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
```

2-kernel/entry.s -> kernel/start.c: 'call start'

```
1  // entry.s jumps here in machine mode
2  __attribute__((aligned(16))) char stack0[4096 * NCPU];
3  void
4  start()
5  {
6      // 将特权模式（MPP, Machine Previous Privilege）设置为监督模式（Supervisor
   Mode）
7      unsigned long x = r_mstatus();
8      x &= ~MSTATUS_MPP_MASK;
9      x |= MSTATUS_MPP_S;
10     w_mstatus(x);
11
12     // 将mepc（机器模式异常程序计数器）寄存器设置为main函数的地址。如果执行mret指令，它将
   跳转到main函数
13     w_mepc((uint64)main); // ↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
14
15     // disable paging for now.
16     w_satp(0);
17
18     // delegate all interrupts and exceptions to supervisor mode.
19     w_medeleg(0xffff);
20     w_mideleg(0xffff);
21     w_sie(r_sie() | SIE_SEIE | SIE_STIE | SIE_SSIE);
22
23     // ask for clock interrupts.
24     timerinit();
```

```

25
26 // keep each CPU's hartid in its tp register, for cpuid().
27 int id = r_mhartid();
28 w_tp(id);
29
30 // switch to supervisor mode and jump to main().
31 asm volatile("mret");
32 }

```

3-kernel/start.c -> kernel/main.c: 'w_mepc((uint64)main)'

```

1 // start() jumps here in supervisor mode on all CPUs.
2 void
3 main()
4 {
5     if(cpuid() == 0){
6         consoleinit();
7         printfinit();
8         printf("\n");
9         printf("xv6 kernel is booting\n");
10        printf("\n");
11        kinit();           // 初始化空闲链表 ↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
12        kvmithart();       // 初始化内核页表, 完成直接映射 ↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
13        kvmithart();       // 初始化页表寄存器 ↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
14        procinit();        // 初始化每个进程的kstack ↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
15        trapinit();        // 陷阱初始化*
16        trapinithart();    // 陷阱初始化*
17        plicinit();        // 中断初始化*
18        plicinithart();    // 中断初始化*
19        binit();           // 文件系统相关初始化*
20        iinit();           // 文件系统相关初始化*
21        fileinit();        // 文件系统相关初始化*
22        virtio_disk_init(); // 模拟磁盘初始化*
23        userinit();        // 第一个用户进程
24        __sync_synchronize();
25        started = 1;
26    } else {
27        while(started == 0)
28            ;
29        __sync_synchronize();
30        printf("hart %d starting\n", cpuid());
31        kvmithart();       // turn on paging
32        trapinithart();    // install kernel trap vector
33        plicinithart();    // ask PLIC for device interrupts
34    }
35
36    scheduler();          // 开始调度 ↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
37 }
38

```

3-1kernel/main.c->kernel/kalloc.c: 'kinit()'

xv6的物理内存使用

- **80000000 -- entry.S, then kernel text and data**: 内核的入口点 (`entry.S`) 以及内核的代码和数据部分。
- **end -- start of kernel page allocation area**: `end` 符号之后是内核开始分配页面的区域。
- **PHYSTOP -- end RAM used by the kernel**: `PHYSTOP` 定义了内核使用的RAM的结束地址, 通常设置为 `KERNBASE + 128MB` (假设有128MB的RAM供内核使用)。

```
1 void
2 kinit() // 初始化物理空间
3 {
4     initlock(&kmem.lock, "kmem");
5     freerange(end, (void*)PHYSTOP); // 将能用的物理空间都free掉
6 }
7
8 void
9 freerange(void *pa_start, void *pa_end)
10 {
11     char *p;
12     p = (char*)PGROUNDUP((uint64)pa_start); // 将start地址转换为页的起始地址, 若
13     // start不是物理页的起始地址, 如: 0x100, 我们应该从它的下一页开始free
14     for(p + PGSIZE <= (char*)pa_end; p += PGSIZE)
15         kfree(p);
16 }
17
18 // Free the page of physical memory pointed at by v,
19 // which normally should have been returned by a
20 // call to kalloc(). (The exception is when
21 // initializing the allocator; see kinit above.)
22 void
23 kfree(void *pa)
24 {
25     struct run *r;
26
27     if(((uint64)pa % PGSIZE) != 0 || (char*)pa < end || (uint64)pa >= PHYSTOP)
28         panic("kfree");
29
30     // Fill with junk to catch dangling refs.
31     memset(pa, 1, PGSIZE);
32
33     r = (struct run*)pa; // 将闲置的物理(4KB)页, 作为struct run 结构体, 来存放空闲链
34     // 表的信息
35
36     acquire(&kmem.lock);
37     r->next = kmem.freelist;
38     kmem.freelist = r;
39     release(&kmem.lock);
40 }
```

3-2kernel/main.c-

>kernel/vm.c: 'kvminit()'+'kvminithart()'

QEMU的物理内存布局

- **00001000 -- boot ROM**: QEMU提供的启动ROM, 用于启动过程。
- **02000000 -- CLINT**: 本地中断控制器 (Core Local Interruptor), 包含定时器等功能。
- **0C000000 -- PLIC**: 可编程中断控制器 (Platform-Level Interrupt Controller), 用于管理来自外设的中断。
- **10000000 -- uart0**: UART 0设备的寄存器地址, 用于串行通信。
- **10001000 -- virtio disk**: VirtIO磁盘设备的内存映射I/O接口。
- **80000000 -- kernel entry**: QEMU的boot ROM跳转至此地址以启动内核, `-kernel` 选项加载的内核镜像也放置在这里。
- **unused RAM after 80000000**: `0x80000000` 之后是未使用的RAM, 供内核和用户程序使用。

xv6的物理内存使用

- **80000000 -- entry.S, then kernel text and data**: 内核的入口点 (`entry.S`) 以及内核的代码和数据部分。
- **end -- start of kernel page allocation area**: `end` 符号之后是内核开始分配页面的区域。
- **PHYSTOP -- end RAM used by the kernel**: `PHYSTOP` 定义了内核使用的RAM的结束地址, 通常设置为 `KERNBASE + 128MB` (假设有128MB的RAM供内核使用)。

在操作系统中, 内核页表中的空间分为两部分: 直接映射区和非直接映射区

- 直接映射区的虚拟空间布局完全与物理内存布局一样, `va=pa`包括:
 - qemu物理布局区
 - xv6物理区 (kernel base~phystop)
- 非直接映射区, 则是采用将指定`va`映射到随机申请/指定的物理空间的方式, 如:
 - trampoline、kstack

从低到高分别是:

- qemu物理布局区
- xv6物理区 (kernel base~phystop)
- kstack
- trampoline

```
1 void // 内核对于页面的映射是采用直接映射的方式: pa即va, 此处则是进行直接映射
2 kvm_init()
3 {
4     // kernel_pagetable变量的地址在内核数据区, 但它是个指针, 它指向的地址属于自由分配区
5     kernel_pagetable = (pagetable_t) kalloc();
6     memset(kernel_pagetable, 0, PGSIZE);
7
8     // 将qemu物理布局直接映射到内核页表 看ppt
9     kvmmap(UART0, UART0, PGSIZE, PTE_R | PTE_W);
10
11     kvmmap(VIRTIO0, VIRTIO0, PGSIZE, PTE_R | PTE_W);
12
13     kvmmap(CLINT, CLINT, 0x10000, PTE_R | PTE_W);
14
15     kvmmap(PLIC, PLIC, 0x400000, PTE_R | PTE_W);
```

```

16     kvmmap(KERNBASE, KERNBASE, (uint64)etext-KERNBASE, PTE_R | PTE_X);
17
18
19     // etext就是end, 所有的内核初始代码和静态变量都放在kernbase-end之间
20     kvmmap((uint64)etext, (uint64)etext, PHYSTOP-(uint64)etext, PTE_R |
PTE_W);
21
22     // trampoline.S的代码, 属于非直接映射区 看ppt
23     kvmmap(TRAMPOLINE, (uint64)trampoline, PGSIZE, PTE_R | PTE_X);
24 }
25
26 void // 页表的地址存放在satp寄存器中, 所以切换页表的过程就是更改satp寄存器中数据的过程
27     // 一开始satp=0, 不会开启分页, 也就是说不会通过页表进行地址翻译工作, 此处给satp赋值,
开启分页了
28 kvmminithart()
29 {
30     w_satp(MAKE_SATP(kernel_pagetable)); // 内存管理magic时刻!!!! 此前, 访问地址
都是物理地址, 此后访问地址是虚拟地址->物理地址
31     sfence_vma(); // 刷新tlb
32 }

```

3-3kernel/main.c ->kernel/proc.c: 'procinit()'

```

1 void // 初始化每个进程的kstack 看ppt
2 procinit(void)
3 {
4     struct proc *p;
5
6     initlock(&pid_lock, "nextpid");
7     for(p = proc; p < &proc[NPROC]; p++) { // 遍历每个进程的pcb
8         initlock(&p->lock, "proc");
9
10        char *pa = kalloc(); // kstack采用的是非直接映射, 所以pa由kalloc()申请而来
11        if(pa == 0)
12            panic("kalloc");
13        uint64 va = KSTACK((int) (p - proc)); // (p-proc)代表着目前是第几个进程,
此处获取第(p - proc)个进程的kstack的虚拟地址
14        kvmmap(va, (uint64)pa, PGSIZE, PTE_R | PTE_W); // 映射kstack
15        p->kstack = va;
16    }
17    kvmminithart();
18 }

```

3-4kernel/main.c->kernel/proc.c: 'userinit()'

trapframe:

这段代码定义了一个名为 `trapframe` 的结构体, 它用于在 RISC-V 架构的操作系统中保存和恢复进程在执行过程中被陷阱 (trap) 中断时的上下文信息。这个结构体特别用于 `trampoline.s` 汇编文件中的陷阱处理代码, 以及用户态和内核态之间的上下文切换。

`trapframe` 结构体包含了多个字段，每个字段都对应着进程在执行过程中需要保存或恢复的寄存器或状态信息。这些字段包括：

- `kernel_satp`：内核页表基址寄存器（SATP）的值，用于在用户态和内核态之间切换页表。
- `kernel_sp`：内核栈顶指针，指向进程内核栈的顶部。
- `kernel_trap`：陷阱处理函数的地址，当从用户态陷入内核态时，会跳转到这个函数执行。
- `epc`：异常程序计数器（EPC），保存了发生陷阱时的用户程序计数器（PC）值。
- `kernel_hartid`：内核的硬件线程ID（hartid），通常保存在线程指针（TP）寄存器中，但在用户态和内核态之间切换时可能需要保存和恢复。

接下来是用户态寄存器的保存区域，包括：

- `ra`：返回地址寄存器。
- `sp`：堆栈指针。
- `gp`：全局指针。
- `tp`：线程指针（在用户态下可能不使用，但在某些情况下需要保存）。
- `t0` 到 `t6`：临时寄存器。
- `s0` 到 `s11`：保存寄存器，用于保存跨函数调用的寄存器值。
- `a0` 到 `a7`：函数参数/返回值寄存器。

这些寄存器在用户态被陷阱中断时，其值会被保存到 `trapframe` 结构体中相应的字段里。当陷阱处理完成后，如果需要返回到用户态继续执行，**这些值会被从 `trapframe` 中恢复回相应的寄存器中**。

需要注意的是，`trapframe` 结构体位于用户页表的某个页面中，但在内核页表中并不特别映射。

```
1 void // 设置第一个应用进程的代码内容、指针等
2 userinit(void)
3 {
4     struct proc *p;
5
6     p = allocproc(); // 第一个应用进程 ↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
7     initproc = p;
8
9     uvminit(p->pagetable, initcode, sizeof(initcode)); // 将initcode.S的内容映射
    到页表中(从地址0开始↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
10    p->sz = PGSIZE;
11
12    p->trapframe->epc = 0; // 在返回用户态时，epc的值会被复制到pc中，pc就会从该
    处开始执行。将epc设置为0，这样该程序在可执行后，就会执行上面的initcode代码了
13    p->trapframe->sp = PGSIZE; // user stack pointer
14
15    safestrcpy(p->name, "initcode", sizeof(p->name));
16    p->cwd = namei("/");
17
18    p->state = RUNNABLE; // 允许第一个应用进程执行
19
20    release(&p->lock);
21 }
22
23
24 // ↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
25 // 向内核申请一个新进程，并返回该进程的pcb
26 static struct proc*
27 allocproc(void)
28 {
```

```

29 struct proc *p;
30
31 for(p = proc; p < &proc[NPROC]; p++) { // 遍历进程队列，如果找到了unused的空闲
pcb, 则执行found
32     acquire(&p->lock); // 获取相关的锁
33     if(p->state == UNUSED) {
34         goto found;
35     } else {
36         release(&p->lock); // 释放相关的锁
37     }
38 }
39 return 0;
40
41 found:
42 p->pid = allocpid(); // 获取pid 进程号
43
44 if((p->trapframe = (struct trapframe *)kalloc()) == 0){ // 向内核申请页面给该
进程的trapframe, 但不映射到内核页表中
45     release(&p->lock); // 释放相关的锁
46     return 0;
47 }
48
49 p->pagetable = proc_pagetable(p); // 为该进程申请页表
↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
50 if(p->pagetable == 0){
51     freeproc(p); // 如果页表创建失败，则释放该进程
52     release(&p->lock); // 释放相关的锁
53     return 0;
54 }
55 // 如果进程从用户态进入内核态，则会把数据存到trapframe里
56 // 但如果是切换进程(进程A切换到进程B)，则A要把自己的寄存器数据都放在p->context里
57 memset(&p->context, 0, sizeof(p->context)); // 清空新进程的上下文数据
58 p->context.ra = (uint64)forkret; // 进程被调度时，会回到上一次还没执行的代码处继续
执行(ra)，而刚刚被创建的进程没有上一次的记录，所以需要给第一次执行的进程一个初始化的函数—
forkret
59 p->context.sp = p->kstack + PGSIZE; // 初始化内核栈顶
60
61 return p;
62 }
63 //
↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
↓↓↓↓↓↓↓↓
64 pagetable_t
65 proc_pagetable(struct proc *p)
66 {
67     pagetable_t pagetable;
68
69     // An empty page table.
70     pagetable = uvmcreate(); // 申请一块物理空间存放页表
71     if(pagetable == 0)
72         return 0;
73
74     if(mappages(pagetable, TRAMPOLINE, PGSIZE, // 用户页表和内核页表中，trampoline
都被映射在同样的va（TRAMPOLINE处）看ppt
75             (uint64)trampoline, PTE_R | PTE_X) < 0){

```


3-5-1kernel/proc.c->kernel/swtch.S: 'swtch(&c->context, &p->context)'

```
1  # 将现在的上下文存到old context处，再从new context处读取下一个进程的上下文 读到寄存器
   # 中，这样就可以执行新的进程了
2  # 由于allocproc()中 将新进程的ra设置为了forkret()，而执行最后一行的ret后会将ra的值赋
   # 给PC，PC就能从forkret()处继续执行了
3  .globl swtch
4  swtch:
5      sd ra, 0(a0)
6      sd sp, 8(a0)
7      sd s0, 16(a0)
8      sd s1, 24(a0)
9      sd s2, 32(a0)
10     sd s3, 40(a0)
11     sd s4, 48(a0)
12     sd s5, 56(a0)
13     sd s6, 64(a0)
14     sd s7, 72(a0)
15     sd s8, 80(a0)
16     sd s9, 88(a0)
17     sd s10, 96(a0)
18     sd s11, 104(a0)
19
20     ld ra, 0(a1)
21     ld sp, 8(a1)
22     ld s0, 16(a1)
23     ld s1, 24(a1)
24     ld s2, 32(a1)
25     ld s3, 40(a1)
26     ld s4, 48(a1)
27     ld s5, 56(a1)
28     ld s6, 64(a1)
29     ld s7, 72(a1)
30     ld s8, 80(a1)
31     ld s9, 88(a1)
32     ld s10, 96(a1)
33     ld s11, 104(a1)
34
35     ret # ↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
36
37
38
```

4-kernel/swtch.S->kernel/proc.c: 'ret' (userinit())创造的进程——系统的第一个用户进程)

```
1  // 子进程在调度器（scheduler）首次调度它时的一些初始化工作。
2  void
3  forkret(void)
4  {
5      static int first = 1;
6
```

```

7 // Still holding p->lock from scheduler.
8 release(&myproc()->lock);
9
10 if (first) { //
11     first = 0;
12     fsinit(ROOTDEV);
13 }
14 // 所有程序被首次调度时都必须做这件事
15 usertrapret(); // 从内核态返回到用户态，并且恢复用户程序执行前的状态（如寄存器值、堆栈等）。↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
16 }

```

5-kernel/pro.c->kernel/trap.c: 'usertrapret()' (userinit())创造的进程——系统的第一个用户进程

```

1 void
2 usertrapret(void)
3 {
4     struct proc *p = myproc();
5
6     intr_off();
7
8     // 在发生陷阱时,CPU会将STVEC寄存器的值复制到PC (Program Counter) 中，此处是将
    trampoline.S中的uservec的地址写入
9     // 如此，之后发生陷阱时，就会直接跳转到trampoline.S的uservec开始执行
10    w_stvec(TRAMPOLINE + (uservec - trampoline));
11
12    p->trapframe->kernel_satp = r_satp(); // 即将从内核态返回用户态，保存内核页表
13    p->trapframe->kernel_sp = p->kstack + PGSIZE; // kstack为内核栈的底部地址，由于栈是向下增长的，该底部地址实际为guard page的起始地址，+ PGSIZE后才是栈的顶部（起始）地址
14    p->trapframe->kernel_trap = (uint64)usertrap; // 设置陷阱处理函数，在trampoline.S中使用到
15    p->trapframe->kernel_hartid = r_tp(); // 保存CPU号
16
17
18
19    unsigned long x = r_sstatus();
20    x &= ~SSTATUS_SPP; // 将ssp-bit设置为0，准备进入用户态
21    x |= SSTATUS_SPIE;
22    w_sstatus(x);
23
24
25    w_sepc(p->trapframe->epc); // 将进入内核态之前的用户PC写入sepc，返回用户态后，系统会将sepc的值复制到pc里，用户就会返回之前的位置继续执行
26
27    // 将用户页表放在satp变量中
28    uint64 satp = MAKE_SATP(p->pagetable);
29
30    uint64 fn = TRAMPOLINE + (userret - trampoline); // fn被设置为trampoline.S中userret的函数地址，也就是一个函数指针

```

```

31 ((void (*)(uint64,uint64))fn)(TRAPFRAME, satp); // 等效于执行
    userret(TRAPFRAME, satp)↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
32 }

```

6-kernel/trap.c-

>kernel/trampoline.S: 'userret(TRAPFRAME, satp)'(userinit())创造的进程——系统的第一个用户进程)

```

1  .globl userret
2  userret:
3      # userret(TRAPFRAME, pagetable)
4      # switch from kernel to user.
5      # usertrapret() calls here.
6      # a0: TRAPFRAME, in user page table.
7      # a1: user page table, for satp.
8
9      # switch to the user page table.
10     csrw satp, a1 # a1是用户页表，此处是将用户页表写入satp中，实现了从内核页表->
        用户页表的切换
11     sfence.vma zero, zero # 刷新TLB
12
13     ld t0, 112(a0) # 112(a0)是trapframe->a0，保存着系统调用的返回值，此时
        t0=trapframe->a0
14     csrw sscratch, t0 # 将trapframe->a0先写入sscratch中
15
16     # 将trapframe中保存的数据(进入内核态之前的用户态数据)重新读回对应的寄存器中
17     ld ra, 40(a0)
18     ld sp, 48(a0)
19     ld gp, 56(a0)
20     ld tp, 64(a0)
21     ld t0, 72(a0)
22     ld t1, 80(a0)
23     ld t2, 88(a0)
24     ld s0, 96(a0)
25     ld s1, 104(a0)
26     ld a1, 120(a0)
27     ld a2, 128(a0)
28     ld a3, 136(a0)
29     ld a4, 144(a0)
30     ld a5, 152(a0)
31     ld a6, 160(a0)
32     ld a7, 168(a0)
33     ld s2, 176(a0)
34     ld s3, 184(a0)
35     ld s4, 192(a0)
36     ld s5, 200(a0)
37     ld s6, 208(a0)
38     ld s7, 216(a0)
39     ld s8, 224(a0)
40     ld s9, 232(a0)
41     ld s10, 240(a0)
42     ld s11, 248(a0)
43     ld t3, 256(a0)

```

```

44         ld t4, 264(a0)
45         ld t5, 272(a0)
46         ld t6, 280(a0)
47
48         # 将a0与sscratch的值互换, 此时a0=trapframe->a0(即系统调研返回
值), sscratch=TRAPFRAME首地址
49         csrrw a0, sscratch, a0
50
51         # sret会将sepc的值写回pc, 这样用户进程就可以回到之前的工作的地方
52         # 记性好的同学可能还记得, 对于系统的第一个进程来说, epc在userinit()中被设定为0
了, 而地址0处装载着initcode.S的代码, 所以第一个进程在返回用户态后就会直接执行initcode.S
53         sret # ↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
54

```

7-kernel/proc.c->user/initcode.S: 'sret' (userinit())创造的进程——系统的第一个用户进程)

```

1  # exec(init, argv)
2  .globl start
3  start:
4      la a0, init
5      la a1, argv
6      li a7, SYS_exec
7      ecall # ↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
8
9  # char init[] = "/init\0";
10 init:
11     .string "/init\0"
12
13 # char *argv[] = { init, 0 };
14 .p2align 2
15 argv:
16     .long init
17     .long 0
18

```

8-initcode.S->kernel/trampoline.S: 'ecall' (userinit())创造的进程——系统的第一个用户进程)

```

1  # ecall执行后, 会将stvec寄存器的值复制到pc中
2  # 记性好的同学或许记得, 刚刚在usertrapret()函数中, 我们已经将stvec的值设定为
trampoline.S的uservec, 所以开始执行uservec
3  uservec:
4      # 记性好的同学应该也记得, 刚刚在userret处, 我们将sscratch的值改成了TRAPFRAME
5      # 所以sscratch记录着TRAPFRAME的虚拟地址, 此操作交换sscratch与a0, 于是a0就记
录着trapframe的地址
6      # 想了解40(a0)到底是什么, 可以去proc.h中查询trapframe结构体的定义, 注释中有标
注
7      csrrw a0, sscratch, a0
8
9      # 将之前说的寄存器都存起来
10     sd ra, 40(a0)

```

```

11      sd sp, 48(a0)
12      sd gp, 56(a0)
13      sd tp, 64(a0)
14      sd t0, 72(a0)
15      sd t1, 80(a0)
16      sd t2, 88(a0)
17      sd s0, 96(a0)
18      sd s1, 104(a0)
19      sd a1, 120(a0)
20      sd a2, 128(a0)
21      sd a3, 136(a0)
22      sd a4, 144(a0)
23      sd a5, 152(a0)
24      sd a6, 160(a0)
25      sd a7, 168(a0)
26      sd s2, 176(a0)
27      sd s3, 184(a0)
28      sd s4, 192(a0)
29      sd s5, 200(a0)
30      sd s6, 208(a0)
31      sd s7, 216(a0)
32      sd s8, 224(a0)
33      sd s9, 232(a0)
34      sd s10, 240(a0)
35      sd s11, 248(a0)
36      sd t3, 256(a0)
37      sd t4, 264(a0)
38      sd t5, 272(a0)
39      sd t6, 280(a0)
40
41      # a0的初始值存在sscratch里, 现在交换sscratch和t0, 再将t0(a0的初始值)存到
      trapframe->a0 即:112(a0)里
42      csrr t0, sscratch
43      sd t0, 112(a0)
44
45      # 8(a0)是kernel_stack_pointer, 由于现在要转入内核模式, 所以要将sp换成 ksp
46      ld sp, 8(a0)
47
48      # 32(a0)是cpu的id, 将该id存到tp寄存器中
49      ld tp, 32(a0)
50
51      # 16(a0)是kernel_trap
52      # 再次, 记性好的同学应该还记得, 在trap.c的usertrapret()中, 我们将
      kernel_trap的值改成了usertrap()函数地址
53      # 所以此时t0即为usertrap()函数的地址
54      ld t0, 16(a0)
55
56      # 0(a0)是内核页表, 由于现在要转入内核模式, 所以要将页表换成内核页表
57      ld t1, 0(a0)
58      csrw satp, t1
59      sfence.vma zero, zero
60
61      # 之前t0中已经存储了usertrap()函数地址, 现在跳转到该地址去执行usertrap()
62      jr t0 # ↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
63

```


9-trampoline.S->kernel/trap.c: 'jr t0:usertrap()'

```
1 // handle an interrupt, exception, or system call from user space.
2 // called from trampoline.S
3 //
4 void
5 usertrap(void)
6 {
7     int which_dev = 0;
8     // status中有一个bit叫做SPP, SPP位=1代表是内核模式, =0代表是用户模式, 此处要求SPP=0
    因为usertrap只能接受来自用户的请求
9     if((r_sstatus() & SSTATUS_SPP) != 0)
10         panic("usertrap: not from user mode");
11
12     w_stvec((uint64)kernelvec);
13
14     // 由于现在进入内核态了, 内核不知道之前是哪个进程调用这个函数
15     // myproc()函数获取当前正在执行的进程指针, 以便能够访问该进程的上下文信息。
16     struct proc *p = myproc();
17
18     // 将当前程序计数器(PC)的值(即sepc寄存器的值)保存到当前进程的trapframe中的epc字
    段。这对于之后恢复用户程序执行是必要的
19     p->trapframe->epc = r_sepc();
20
21     if(r_scause() == 8){
22         // system call
23
24         if(p->killed)
25             exit(-1);
26
27         // 之前我们将PC的值存在epc中, 该值为ecall所在地址, 之后返回要返回到ecall的下一条指
    令中, 所以+4
28         p->trapframe->epc += 4;
29
30         intr_on();
31
32         syscall(); // ↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
33     } else if((which_dev = devintr()) != 0){
34         // ok
35     } else {
36         printf("usertrap(): unexpected scause %p pid=%d\n", r_scause(), p->pid);
37         printf("             sepc=%p stval=%p\n", r_sepc(), r_stval());
38         p->killed = 1;
39     }
40
41     if(p->killed)
42         exit(-1);
43
44     // give up the CPU if this is a timer interrupt.
45     if(which_dev == 2)
46         yield();
47
48     usertrapret(); // ↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
49 }
```

10-kernel/trap.c->kernel/syscall.c: 'exec init'

```
1 void
2 syscall(void)
3 {
4     int num;
5     struct proc *p = myproc();
6
7     num = p->trapframe->a7; // a7处存了系统调用编号
8     if(num > 0 && num < NELEM(syscalls) && syscalls[num]) { //
9         ↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓判断系统调用编号是否合法
10        p->trapframe->a0 = syscalls[num](); // 取出对应的系统调用函数指针,执行系统调用
11        指针
12    } else { // 返回错误信息
13        printf("%d %s: unknown sys call %d\n",
14               p->pid, p->name, num);
15        p->trapframe->a0 = -1;
16    }
17 }
```

11-kernel/syscall.c:exec(init)->user/init.c: 'syscalls[sys_init]()'

```
1 char *argv[] = { "sh", 0 };
2
3 int
4 main(void)
5 {
6     int pid, wpid;
7
8     if(open("console", O_RDWR) < 0){
9         mknod("console", CONSOLE, 0);
10        open("console", O_RDWR);
11    }
12    dup(0);
13    dup(0);
14
15    for(;;){
16        printf("init: starting sh\n");
17        pid = fork(); // ↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
18        if(pid < 0){
19            printf("init: fork failed\n");
20            exit(1);
21        }
22        if(pid == 0){ // 子进程
23            exec("sh", argv); // ↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓子进程exec运行sh.c
24            printf("init: exec sh failed\n");
25            exit(1);
26        }
27
28        for(;;){ // 父进程等待子进程结束
29
30            wpid = wait((int *) 0);
```

```

31     if(wpid == pid){
32         // the shell exited; restart it.
33         break;
34     } else if(wpid < 0){
35         printf("init: wait returned an error\n");
36         exit(1);
37     } else {
38         // it was a parentless process; do nothing.
39     }
40 }
41 }
42 }
43

```

11-1 user/init.c->kernel/proc.c: 'fork()'

```

1  // 创建一个子进程，父进程的返回值为子进程的pid，子进程返回值为0
2  int
3  fork(void)
4  {
5      int i, pid;
6      struct proc *np;
7      struct proc *p = myproc(); // 获取当前进程(父进程)的pcb
8
9      // Allocate process.
10     if((np = allocproc()) == 0){ // 为子进程申请新的pcb
11         return -1;
12     }
13
14     // Copy user memory from parent to child.
15     if(uvmcopy(p->pagetable, np->pagetable, p->sz) < 0){ // 将父进程的数据copy到
        子进程的内存空间中↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓↓
16         freeproc(np);
17         release(&np->lock);
18         return -1;
19     }
20     np->sz = p->sz; // 复制父进程的size数据
21
22     np->parent = p; // 确定parent关系
23
24     // copy saved user registers.
25     *(np->trapframe) = *(p->trapframe); // 将父进程的trapframe内容复制给子进程
26
27     // Cause fork to return 0 in the child.
28     np->trapframe->a0 = 0; // a0一般代表着系统调用的返回值，将子进程的a0设置为0 即可让
        子进程的系统调用返回0，用来区分 父/子
29
30     // increment reference counts on open file descriptors.
31     for(i = 0; i < NOFILE; i++) // 子进程也要打开父进程使用的文件
32         if(p->ofile[i])
33             np->ofile[i] = filedup(p->ofile[i]); // 文件被进程打开时会有reference记
        录，子进程也打开这些文件所以ref值要增加
34     np->cwd = idup(p->cwd);
35

```

```

36     safestrcpy(np->name, p->name, sizeof(p->name)); // 将父进程的名字复制给子进程
37
38     pid = np->pid; // 获取子进程的pid以返回
39
40     np->state = RUNNABLE; // 将子进程设定为可执行状态，调度后则会开始运行
41
42     release(&np->lock);
43
44     return pid; // 父进程返回子进程的pid，子进程返回值为0
45 }

```

11-1-1 kernel/proc.c->kernel/vm.c: 'uvmcopy()'

```

1  // 将old页表中的数据copy到新页表的空间中
2  int
3  uvmcopy(pagetable_t old, pagetable_t new, uint64 sz)
4  {
5      pte_t *pte;
6      uint64 pa, i;
7      uint flags;
8      char *mem;
9
10     for(i = 0; i < sz; i += PGSIZE){ // 每次复制一个page
11         if((pte = walk(old, i, 0)) == 0) // 寻找旧页表中 va=i所对应的页表项PTE
12             panic("uvmcopy: pte should exist");
13         if((*pte & PTE_V) == 0) // PTE应该要是有效的
14             panic("uvmcopy: page not present");
15         pa = PTE2PA(*pte); // 根据PTE中的信息获取old页面的起始物理地址
16         flags = PTE_FLAGS(*pte); // 根据PTE中的信息获取该页面的flags
17         if((mem = kalloc()) == 0) // 申请新页面
18             goto err;
19         memmove(mem, (char*)pa, PGSIZE); // 将old页面的数据copy到新页面
20         if(mappages(new, i, PGSIZE, (uint64)mem, flags) != 0){ // 将新页面映射到
new页表中
21             kfree(mem);
22             goto err;
23         }
24     }
25     return 0;
26
27     err: // 如果这个过程出现失败，则要回退之前做过的所有copy
28         uvmunmap(new, 0, i / PGSIZE, 1);
29         return -1;
30 }

```

12-user/init->user/sh: 'exec("sh",argv)'

```

1  int // 这就是我们看到的shell
2  main(void)
3  {
4      static char buf[100];
5      int fd;
6

```

```

7 // Ensure that three file descriptors are open.
8 while((fd = open("console", O_RDWR)) >= 0){
9     if(fd >= 3){
10         close(fd);
11         break;
12     }
13 }
14
15
16 while(getcmd(buf, sizeof(buf)) >= 0){ // shell 会不停地接收指令，执行指令
17     if(buf[0] == 'c' && buf[1] == 'd' && buf[2] == ' '){
18         // chdir must be called by the parent, not the child.
19         buf[strlen(buf)-1] = 0; // chop \n
20         if(chdir(buf+3) < 0)
21             fprintf(2, "cannot cd %s\n", buf+3);
22         continue;
23     }
24     if(fork1() == 0) // shell 每次收到指令就fork一个子进程去运行指令，而自己则
        wait(0)等待子进程直接完毕才接收下一条指令
25         runcmd(parsecmd(buf));
26     wait(0);
27 }
28 exit(0);
29 }

```

三、任务一：

- 任务名称：打印页表
- 任务要求：
 - 定义一个vmprint (pagetable_t pagetable) 函数（推荐放在vm.c文件中），可以打印出该页表中所有pte（三级页目录都要访问过去哦）
 - 若pte无效，则不打印
 - 在exec.c的exec()函数返回前增加一个条件：if(p->pid==1)，则调用vmprint (p->pagetable)
 - 输出要求如下：
 - 第一行为"page table "+<页表地址>后续行输出 ".."* ++": "+"pte "+"pa"+
 - level=1—最高级页目录，2—中级页目录，3—最低级页目录
 - index=页目录号，参照前面的ppt
 - “..”间要有空格
 - 使用深度优先遍历

```

page table 0x0000000087f6e000
..0: pte 0x0000000021fda801 pa 0x0000000087f6a000
.. ..0: pte 0x0000000021fda401 pa 0x0000000087f69000
.. .. .0: pte 0x0000000021fdac1f pa 0x0000000087f6b000
.. .. .1: pte 0x0000000021fda00f pa 0x0000000087f68000
.. .. .2: pte 0x0000000021fd9c1f pa 0x0000000087f67000
..255: pte 0x0000000021fdb401 pa 0x0000000087f6d000
.. ..511: pte 0x0000000021fdb001 pa 0x0000000087f6c000
.. .. .510: pte 0x0000000021fdd807 pa 0x0000000087f76000
.. .. .511: pte 0x0000000020001c0b pa 0x0000000080007000

```

- 提示：
 - 如何遍历页表？可以参考walk()函数
 - 如何判断pte无效？与‘PTE_V’有关

四、任务二、三：

- 任务名称：为每个用户进程新增一个独享的内核页表并应用
- 任务背景说明：
 - 由于内核页表为所有进程共享的公共资源，所以不可能保存用户进程的用户页面映射信息，这就导致在内核态使用copyin()、copyout()、copyinstr()等函数访问用户页面时，需要通过walkaddr()软件翻译用户页面的地址，速度很慢。
 - 若是为每个用户进程创建一个独享的内核页表，这个内核页表中记录了用户页面的映射信息，并在进入内核态时将satp(页表寄存器)中的页表换为独享内核页表(instead of 共享内核页表)，则内核在访问用户页面的时，可以直接访问用户页面地址（自动触发硬件翻译，很快）
- 任务要求：
 - 在proc.h中为每个pcb增加一个独享内核页表变量，也可以叫kernel_pagetable
 - 将所有的copyin()和copyinstr()替换为vmcopyin.c中的copyin_new()、copyinstr_new()
 - 以合适的方式维护该独享内核页表，如：
 - 在alloc_proc()时创建、初始化它
 - 只映射自己的kstack，别把别人的kstack映射到你的独享内核页表中（参考procinit()函数）
 - 在修改用户页表时，同步该改动，如：释放用户页表（freeproc()）时、父进程复制数据给子进程时(fork())，exec()修改页表时，growproc()时，userinit()时 等等，大家自己搜索相关的操作
 - 在内核态时都要使用内核页表，但内核态分为有程序运行状态和无程序状态，若无程序在运行，则用共享的内核页表，若有程序在运行，则使用用户程序独享的内核页表（scheduler()函数处见分晓）
 - 虚拟地址不能无限增长，不然QEMU物理内存所占用的条目就被刷新掉了（此处VA限制为PLIC），在growproc()函数中要判断VA是否越界

