

Lab 5: Process Management and System Call

习题解答

1. 进程管理

1.1 关于 PCB 设计

在 proc（即 PCB）中仅存储了进程的 trapframe 与 context 的指针，请说明 trapframe 与 context 的实例存在何处，为什么要这样设计？

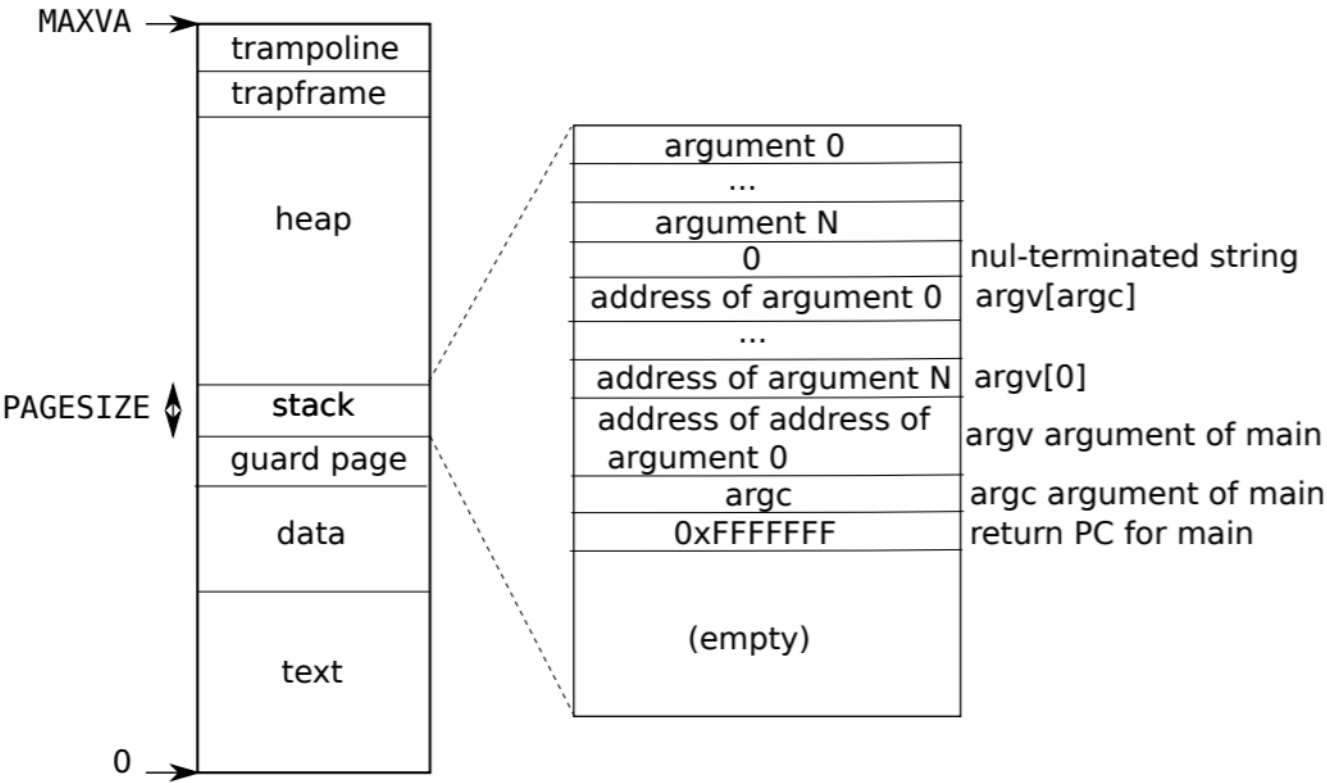


Figure 3.4: A process's user address space, with its initial stack.

本图引自 *xv6: a simple, Unix-like teaching operating system* [1]。

如图所示，trapframe 的实例存在进程的用户地址空间（user address space）。为什么 PCB 中仅存储 trapframe 的指针？因为 trapframe 是一个保存了所有通用寄存器和一些特殊寄存器的结构。一方面，trap 时需要用到 trapframe 中的数据，所以它应当以某种形式存储在 PCB 中；另一方面，trap 时我们只能传入 trapframe 指针，因为硬件没有提供足够多的寄存器来在 trap 时传入整个 trapframe 结构。因此，在 PCB 中仅存储 trapframe 指针是一个节省空间的方案，trap 时只需一个寄存器用于传入 trapframe 指针即可 [1:1]。

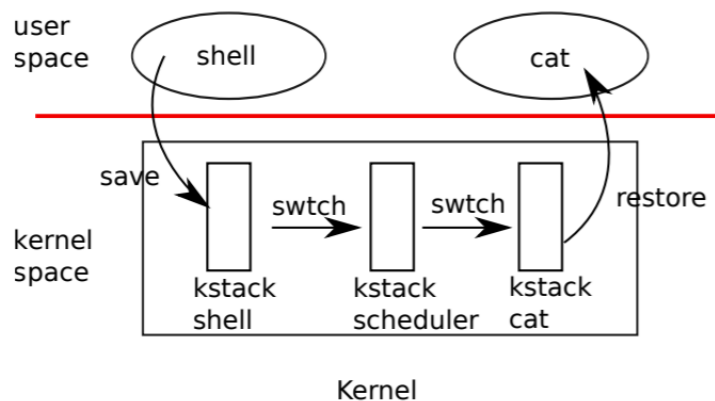


Figure 7.1: Switching from one user process to another. In this example, xv6 runs with one CPU (and thus one scheduler thread).

本图引自 *xv6: a simple, Unix-like teaching operating system* [1:2]。

context 的实例存在执行 context switch 的内核所对应的 kernel stack 处。与 trapframe 类似，context 也是一个保存了一组通用寄存器的结构。scheduler 在 context switch 时需要交换进程的 context，然而我们并没有这么多寄存器来在函数调用时传入整个 context 结构。因此，我们在 PCB 中仅存储 context 的指针，这样在调度时 scheduler 就只需用到两个寄存器，分别存放了将被调入和调出的新旧进程的 context 指针 [1:3]。

1.2 Context switch

请完成 inc/proc.h 中 struct context 的定义以及 kern/swtch.S 中 context switch 的实现。

1.2.1 Context 设计：struct context

context 中需要保存所有的 callee-saved 寄存器 [1:4]，根据 ARM 开发文档 [2]，即通用寄存器 X19 ~ X28。此外，我们额外保存寄存器 X29 (Frame Pointer) 和 X30 (Procedure Link Register)，其中 X30 用于指定用户进程初次运行的地址。

```

1 // inc/proc.h
2
3 struct context {
4     // Callee-saved Registers
5     uint64_t x19;
6     uint64_t x20;
7     uint64_t x21;
8     uint64_t x22;
9     uint64_t x23;
10    uint64_t x24;
11    uint64_t x25;
12    uint64_t x26;
13    uint64_t x27;
14    uint64_t x28;
15
16    uint64_t x29; // Frame Pointer
17    uint64_t x30; // Procedure Link Register
18 };

```

1.2.2 Context switch 实现

context switch 主要做了以下几件事情 [3]:

1. 将当前（将被调出的）旧进程的 callee-saved 寄存器压栈
2. 将当前栈指针的地址保存在 `*old`（其中 `old` 是函数调用传入的第一个参数，位于寄存器 X0），构建旧进程的 context
3. 将 `new`（即函数调用传入的第二个参数，位于寄存器 X1）的值覆盖当前栈指针的地址，切换到（将被调入的）新进程的 context
4. 将新进程的 callee-saved 寄存器弹栈，函数 `swtch` 返回（`ret`，等价于 `mov pc, x30`）

```
1  # kern/swtch.S
2
3  /*
4   * Context switch
5   *
6   * void swtch(struct context **old, struct context *new);
7   *
8   * Save current register context on the stack,
9   * creating a struct context, and save its address in *old.
10  * Switch stacks to new and pop previously-saved registers.
11  */
12 .global swtch
13
14 swtch:
15     # Save old callee-saved registers
16     stp x29, x30, [sp, #-16]!
17     stp x27, x28, [sp, #-16]!
18     stp x25, x26, [sp, #-16]!
19     stp x23, x24, [sp, #-16]!
20     stp x21, x22, [sp, #-16]!
21     stp x19, x20, [sp, #-16]!
22
23     # Switch stacks
24     mov x19, sp
25     str x19, [x0]
26     mov sp, x1
27
28     # Load new callee-saved registers
29     ldp x19, x20, [sp], #16
30     ldp x21, x22, [sp], #16
31     ldp x23, x24, [sp], #16
32     ldp x25, x26, [sp], #16
33     ldp x27, x28, [sp], #16
34     ldp x29, x30, [sp], #16
35
36     ret
```

1.3 关于 Context switch 设计

1.3.1

在 `kern/proc.c` 中将 `swtch` 声明为 `void swtch(struct context**, struct context*)`，请说明为什么要这样设计？

因为如果第一个参数传的是 `struct context*`，那么在函数 `swtch` 中对第一个参数值的修改（也就是将栈指针的地址保存在寄存器 `X0`）将无法反映到函数外部。即在函数返回后，这个局部变量就会失效，这样也就无法保存旧进程的 `context` 指针。而传入 `struct context**`，就可以通过修改这个指针所指向的地址，来将旧进程的 `context` 地址传给函数外部。

1.3.2

`context` 中仅需要存储 callee-saved registers，请结合 PCS 说明为什么？

因为根据 PCS (Procedure Call Standard) [2:1]，函数调用时，callee 只需要确保约定的 callee-saved 寄存器中的数据不被损坏（corrupt），而其他寄存器中的数据是可以损坏的。因此在 `context switch` 中，`context` 不需要存储 callee-saved 寄存器以外的其他寄存器，因为即使这些数据在 `context switch` 的过程中被损坏了也没有关系。`context` 只需保护 callee-saved 寄存器中的数据不受 `context switch` 影响即可。

1.3.3

与 `trapframe` 对比，请说明为什么 `trapframe` 需要存储这么多信息？

因为 `trap` 过程不是函数调用，没有 caller 和 callee 的说法，不遵循也无法遵循 PCS 规范。例如系统中断时，内核可以直接中断用户程序，用户程序并不会有机会提前保存所谓的 caller-saved 寄存器，但这些数据同样是不应在 `trap` 后被内核程序损坏的。因此 `trapframe` 需要存储所有通用寄存器，才能保证之后回到用户态时可以正确还原用户程序的数据。

1.3.4

`trapframe` 似乎已经包含了 `context` 中的内容，为什么上下文切换时还需要先 `trap` 再 `switch`？

因为 `trap` 过程是从用户态切换到内核态的过程，`switch` 过程是内核态中的过程。上下文切换需要在内核态中进行，因此还是要先 `trap` 再 `switch`。虽然 `trapframe` 似乎包含了 `context` 中的内容，但它们完全是两个不同的东西，有着不同的用途，保存在不同的位置，因此也无法复用其中的数据。

1.4 内核进程管理模块

请根据 `kern/proc.c` 中相应代码的注释完成内核进程管理模块以支持调度第一个用户进程 `user/initcode.S`。

1.4.1 PCB 设计： `struct proc`

每个用户进程的 PCB 中保存了以下数据，具体作用参见注释 [4]：

```

1 // inc/proc.h
2
3 struct proc {
4     struct spinlock lock;
5
6     // p->lock must be held when using these:
7     enum procstate state; // Process state
8     void* chan;           // If non-zero, sleeping on chan
9     int killed;           // If non-zero, have been killed
10    int xstate;            // Exit status to be returned to parent's wait
11    int pid;               // Process ID
12
13    // wait_lock must be held when using these:
14    struct proc* parent; // Parent process
15
16    // no lock needs to be held when using these:
17    char* kstack;          // Bottom of kernel stack for this process
18    uint64_t sz;           // Size of process memory (bytes)
19    uint64_t* pgdir;       // Page table
20    struct trapframe* tf;  // Trapframe for current syscall
21    struct context* context; // switch() here to run process
22    char name[16];         // Process name (debugging)
23 };

```

1.4.2 锁的初始化: `proc_init`

函数 `proc_init` 的主要工作是完成 `ptable` 锁的初始化，以处理多核的并发问题。这里我们不是在整个 `struct ptable` 中，而是选择在每个 `struct proc` 中新增一个自旋锁 `proc_lock`。这样做的目的是为了锁的控制粒度更细，实际上这也是 Xv6 for RISC-V [\[4:1\]](#) 的实现方法。

在获取进程 `pid` 时，同样也存在并发问题，因此这里也为 `pid` 新增了一个自旋锁 `pid_lock`。自旋锁 `wait_lock` 则是为了确保父进程在子进程返回前维持等待状态，防止子进程 `p` 在访问父进程 `p->parent` 时发现父进程已丢失。

```

1 // kern/proc.c
2
3 /*
4  * Initialize the spinlock for ptable to serialize the access to ptable
5  */
6 void
7 proc_init()
8 {
9     initlock(&wait_lock, "wait_lock");
10    initlock(&pid_lock, "pid_lock");
11    for (struct proc* p = ptable.proc; p < &ptable.proc[NPROC]; ++p) {
12        initlock(&p->lock, "proc_lock");
13    }
14    cprintf("proc_init: success.\n");
15 }

```

1.4.3 创建新进程: `proc_alloc`

函数 `proc_alloc` 的主要工作是遍历进程表 `ptable`，找到一个 UNUSED 进程，进行内核部分的初始化工作，最后返回进程的 `proc` 指针。具体来说：

1. 利用函数 `pid_next (kern/proc.c)` 分配 PID
2. 利用函数 `kalloc (kern/kalloc.c)` 分配内核栈 `kstack`
3. 在 `kstack` 的栈顶分配一块空间作为 `trapframe`
4. 在 `trapframe` 下面再分配一块空间作为 `context`，并进行初始化；其中寄存器 X30 保存函数 `forkret` 的地址，作为进程初次从函数 `swtch` 返回时的返回地址；这里函数 `forkret` 只需在进程第一次被 `scheduler` 调度时进入一次，之后就不再需要进入了，关于调度时所有函数调用的过程还会在 1.4.5 节细讲
5. 设置进程状态为 EMBRYO

如果创建进程失败，则返回 `NULL`。

```

1 // kern/proc.c
2
3 /*
4  * Look through the process table for an UNUSED proc.
5  * If found, change state to EMBRYO and initialize
6  * state required to run in the kernel.
7  * Otherwise return 0.
8  */
9 static struct proc*
10 proc_alloc()
11 {
12     for (struct proc* p = ptable.proc; p < &ptable.proc[NPROC]; ++p) {
13         acquire(&p->lock);
14         if (p->state != UNUSED) {
15             release(&p->lock);
16             continue;
17         }
18
19         p->pid = pid_next();
20
21         // Allocate kernel stack.
22         if (!(p->kstack = kalloc())) {
23             proc_free(p);
24             release(&p->lock);
25             return NULL;
26         }
27         char* sp = p->kstack + KSTACKSIZE;
28
29         // Leave room for trapframe.
30         sp -= sizeof(*p->tf);
31         p->tf = (struct trapframe*)sp;
32
33         // Set up new context to start executing at forkret.
34         sp -= sizeof(*p->context);
35         p->context = (struct context*)sp;
36         memset(p->context, 0, sizeof(*p->context));
37         p->context->x30 = (uint64_t)forkret;
38
39         p->state = EMBRYO;
40         cprintf("proc_alloc: proc %d success.\n", p->pid);
41         return p;
42     }
43     return NULL;
44 }

```

其中，函数 `proc_free` 的作用是清空进程的 PCB，并利用函数 `kfree` 释放申请的内存。

```

1 // kern/proc.c
2
3 /*
4  * Free a proc structure and the data hanging from it,
5  * including user pages.
6  * p->lock must be held.
7  */
8 static void
9 proc_free(struct proc* p)
10 {
11     p->chan = NULL;
12     p->killed = 0;
13     p->xstate = 0;
14     p->pid = 0;
15     p->parent = NULL;
16     if (p->kstack) kfree(p->kstack);
17     p->kstack = NULL;
18     p->sz = 0;
19     if (p->pgdir) kfree((char*)p->pgdir);
20     p->pgdir = NULL;
21     p->tf = NULL;
22     p->name[0] = '\0';
23     p->state = UNUSED;
24 }

```

1.4.4 初始化用户进程：user_init

函数 `user_init` 的主要工作是初始化第一个用户进程。具体来说：

1. 利用函数 `proc_alloc` (`kern/proc.c`) 进行内核部分的初始化
2. 利用函数 `pgdir_init` (`kern/vm.c`) 分配一个用户页表，并指定进程的内存空间为一个页表的大小 `PGSIZE`
3. 利用函数 `uvm_init` (`kern/vm.c`) 将初始化二进制码 `initcode` 加载到页表的起始位置
4. 清空 trapframe，并进行初始化；其中寄存器 `SP_ELO` 设置为 `PGSIZE`，其余寄存器设置为 `0`；部分寄存器会在函数 `trapret` 返回 (`eret`) 时用到，之后会在 1.4.5 节细讲
5. 设置进程名为 `initproc`
6. 设置进程状态为 `RUNNABLE`

需要注意的是，函数 `user_init` 只需在 CPU0 中执行一次，而无需在其他 CPU 上执行（在这个坑上我花了 3 个多小时调试 orz）。


```

1 // kern/proc.c
2
3 /*
4  * Set up first user process (only used once).
5  * Set trapframe for the new process to run
6  * from the beginning of the user process determined
7  * by uvm_init
8  */
9 void
10 user_init()
11 {
12     extern char _binary_obj_user_initcode_start[];
13     extern char _binary_obj_user_initcode_size[];
14
15     struct proc* p = proc_alloc();
16     if (!p) panic("\tuser_init: process failed to allocate.\n");
17     initproc = p;
18
19     // Allocate a user page table.
20     if (!(p->pgdir = pgdir_init()))
21         panic("\tuser_init: page table failed to allocate.\n");
22     p->sz = PGSIZE;
23
24     // Copy initcode into the page table.
25     uvm_init(
26         p->pgdir, _binary_obj_user_initcode_start,
27         (uint64_t)_binary_obj_user_initcode_size);
28
29     // Set up trapframe to prepare for the first "return" from kernel to user.
30     memset(p->tf, 0, sizeof(*p->tf));
31     p->tf->x30 = 0; // initcode start address
32     p->tf->sp_el0 = PGSIZE; // user stack pointer
33     p->tf->spsr_el1 = 0; // program status register
34     p->tf->elr_el1 = 0; // exception link register
35
36     strncpy(p->name, "initproc", sizeof(p->name));
37     p->state = RUNNABLE;
38     release(&p->lock);
39
40     cprintf("user_init: proc %d (%s) success.\n", p->pid, p->name, cpuid());
41 }

```

其中，函数 `pgdir_init` 用于获取一个新页表。

```

1 // kern/vm.c
2
3 /*
4  * Get a new page table.
5  */
6 uint64_t*
7 pgdir_init()
8 {
9     uint64_t* pgdir;
10    if (!(pgdir = (uint64_t*)kalloc())) return NULL;
11    memset(pgdir, 0, PGSIZE);
12    return pgdir;
13 }

```

函数 `uvm_init` 则用于将二进制码加载到页表地址 `0x0` 的位置。

```

1 // kern/vm.c
2
3 /*
4  * Load binary code into address 0 of pgdir.
5  * sz must be less than a page.
6  * The page table entry should be set with
7  * additional PTE_USER|PTE_RW|PTE_PAGE permission
8  */
9 void
10 uvm_init(uint64_t* pgdir, char* binary, uint64_t sz)
11 {
12     char* mem;
13     if (sz ≥ PGSIZE) panic("\tuvvm_init: sz must be less than a page.\n");
14     if (!(mem = kalloc())) panic("\tuvvm_init: not enough memory.\n");
15     memset(mem, 0, PGSIZE);
16     map_region(
17         pgdir, (void*)0, PGSIZE, (uint64_t)mem, PTE_USER | PTE_RW | PTE_PAGE);
18     memmove((void*)mem, (const void*)binary, sz);
19 }

```

1.4.5 内核调度：scheduler

函数 `scheduler` 的主要工作是调度进程，这就是个大工程了。以下我们按函数调用顺序慢慢展开。

函数 `scheduler` 的主体部分如下：

```

1 // kern/proc.c
2
3 /*
4  * Per-CPU process scheduler
5  * Each CPU calls scheduler() after setting itself up.
6  * Scheduler never returns. It loops, doing:
7  * - choose a process to run
8  * - swtch to start running that process
9  * - eventually that process transfers control
10  *   via swtch back to the scheduler.
11  */
12 void
13 scheduler()
14 {
15     struct cpu* c = thiscpu;
16     c->proc = NULL;
17
18     while (1) {
19         // Loop over process table looking for process to run.
20         for (struct proc* p = ptable.proc; p < &ptable.proc[NPROC]; ++p) {
21             acquire(&p->lock);
22             if (p->state != RUNNABLE) {
23                 release(&p->lock);
24                 continue;
25             }
26
27             // Switch to chosen process. It is the process's job
28             // to release its lock and then reacquire it
29             // before jumping back to us.
30             c->proc = p;
31             uvm_switch(p);
32             p->state = RUNNING;
33             cprintf(
34                 "scheduler: switch to proc %d at CPU %d.\n", p->pid, cpuid());
35
36             swtch(&c->scheduler, p->context);
37
38             // Process is done running for now.
39             // It should have changed its p->state before coming back.
40             c->proc = NULL;
41             release(&p->lock);
42         }
43     }
44 }

```

首先，`scheduler` 遍历进程表 `ptable`，找到一个 `RUNNABLE` 进程 `p`，然后利用函数 `uvm_switch` 切换到该进程的用户页表 `p->ptable`，并设置进程状态为 `RUNNING`。

```

1 // kern/vm.c
2
3 /*
4  * Switch to the process's own page table for execution of it.
5  */
6 void
7 uvm_switch(struct proc* p)
8 {
9     if (!p) panic("\tuvm_switch: no process.\n");
10    if (!p->kstack) panic("\tuvm_switch: no kstack.\n");
11    if (!p->pgdir) panic("\tuvm_switch: no pgdir.\n");
12
13    lttbr0(V2P(p->pgdir)); // Switch to process's address space
14 }

```

随后调用函数 `swtch`，切换到该进程的 context，参见 1.2.2 节。`swtch` 的返回地址由 `p->context` 保存的寄存器 X30 决定。1.4.3 节中我们提到，X30 保存的是函数 `forkret` 的地址，因此进程初次被 `scheduler` 调度，从函数 `swtch` 返回时，将返回到 `forkret`。此后，进程就按照每次切换 context 时 X30 保存的地址，返回到用户地址空间的相应位置即可。

```

1 // kern/proc.c
2
3 /*
4  * A fork child's very first scheduling by scheduler()
5  * will switch to forkret. "Return" to user space.
6  */
7 void
8 forkret()
9 {
10    struct proc* p = thiscpu->proc;
11
12    // Still holding p->lock from scheduler.
13    release(&p->lock);
14
15    // Pass trapframe pointer as an argument when calling trapret.
16    usertrapret(p->tf);
17 }

```

函数 `forkret` 的作用是在进程初次被调度时，释放 `scheduler` 持有的进程锁 `p->lock`，并进行一些必须在用户进程中才能进行的初始化工作，例如文件系统的初始化（因为需要调用 `sleep` 休眠当前进程，故不能在函数 `main` 中执行）。由于目前我们还没有实现文件系统，因此目前 `forkret` 只是为这些初始化工作预留一下位置。

接下来函数 `forkret` 应该返回到函数 `trapret`。这里一个非常 tricky 的点在于，如何返回？关于这点我研究了 7 个多小时，阅读了大量手册和源码。这项工作的难点在于，如果直接返回，那么由于 1.2.2 节我们设置的 context 中寄存器 X30 的值为函数 `forkret` 的地址，而且后续没有地方修改过，因此这里 `forkret` 还是会返回到 `forkret`，导致死循环。那如果直接调用 `trapret` 呢？由于当前栈指针 SP 保存的地址指向函数 `forkret` 目前栈帧的栈顶，显然不是进程 trapframe 的地址 `p->tf`。然而 `trapret` 在还原寄存器时需要用到 SP 的值，且该值应当为 `p->tf`，错误的 SP 值将导致 `trapret` 无法正常工作。

```

1  # kern/trapasm.S
2
3  /* Return falls through to trapret. */
4  .global trapret
5  trapret:
6      /* Restore registers. */
7      ldp x1, x2, [sp], #16
8      ldp x3, x0, [sp], #16
9      msr sp_el0, x1
10     msr spsr_el1, x2
11     msr elr_el1, x3
12
13     ldp x1, x2, [sp], #16
14     ldp x3, x4, [sp], #16
15     ldp x5, x6, [sp], #16
16     ldp x7, x8, [sp], #16
17     ldp x9, x10, [sp], #16
18     ldp x11, x12, [sp], #16
19     ldp x13, x14, [sp], #16
20     ldp x15, x16, [sp], #16
21     ldp x17, x18, [sp], #16
22     ldp x19, x20, [sp], #16
23     ldp x21, x22, [sp], #16
24     ldp x23, x24, [sp], #16
25     ldp x25, x26, [sp], #16
26     ldp x27, x28, [sp], #16
27     ldp x29, x30, [sp], #16
28
29     eret

```

由于 X86 下函数返回时要先弹栈，而此前栈底预先保存的值即为函数返回地址。因此 Xv6 [\[3:1\]](#) 的解决方案是，在 kstack 中 context 部分的上方保存 trapret 的地址。由于之前函数是直接跳转到 forkret 的，而没有进行正常函数调用前必要的将函数返回地址压栈的操作。因此利用这个 trick，可以使得 forkret 返回时，弹栈得到的返回地址为 trapret 的地址，实现跳转。同时这样的好处是，弹栈后，位于返回地址上方的地址正好就是 trapframe 的地址，因此栈指针 SP 的值也是正确的，恰好指向 p->tf。

```

1  // https://github.com/mit-pdos/xv6-public/blob/master/proc.c
2
3  sp = p->kstack + KSTACKSIZE;
4
5  // Leave room for trap frame.
6  sp -= sizeof *p->tf;
7  p->tf = (struct trapframe*)sp;
8
9  // Set up new context to start executing at forkret,
10 // which returns to trapret.
11 sp -= 4;
12 *(uint*)sp = (uint)trapret;
13
14 sp -= sizeof *p->context;
15 p->context = (struct context*)sp;
16 memset(p->context, 0, sizeof *p->context);
17 p->context->eip = (uint)forkret;

```

然而，ARM 下的函数返回机制与 X86 不同，是将寄存器 X30 的值作为返回地址，而没有这步弹栈操作。因此这个方案在 ARM 下不可行。

Xv6 for RISC-V [4:2] 的解决方案是，不采用 X86 下直接返回的方式，而是调用函数 `usertrapret`。在 `usertrapret` 的最后，这里其实是调用了函数 `userret`，有点类似于我们的 `trapret`。

```
1 // https://github.com/mit-pdos/xv6-riscv/blob/riscv/kernel/trap.c
2
3 // jump to trampoline.S at the top of memory, which
4 // switches to the user page table, restores user registers,
5 // and switches to user mode with sret.
6 uint64 fn = TRAMPOLINE + (userret - trampoline);
7 ((void (*)(uint64,uint64))fn)(TRAPFRAME, satp);
```

```
1 # https://github.com/mit-pdos/xv6-riscv/blob/riscv/kernel/trampoline.S
2
3 .globl userret
4 userret:
5     # userret(TRAPFRAME, pagetable)
6     # switch from kernel to user.
7     # usertrapret() calls here.
8     # a0: TRAPFRAME, in user page table.
9     # a1: user page table, for satp.
```

但需要注意的是，Xv6 for RISC-V 的实现中，函数 `userret` 是可以传入参数的，而参数中正包含了 `trapframe` 的地址，这其实也可以解决我们的问题。然而不幸的是，我们的 `trapret` 是不接受参数的，这意味着我们并不能将 `trapframe` 的地址传入 `trapret`，从而覆盖 SP 的值。

思来想去，最后还是没有想到什么优雅的解决方案（能用 C 语言解决都已经算优雅了）。因此我就暴力地用汇编写了个函数 `void usertrapret(struct trapframe*)`，本质是对 `trapret` 的重载，区别在于可以接受一个 `trapframe` 指针作为参数了。

```
1 # kern/trapasm.S
2
3 /* Help forkret to call trapret in an expected way. */
4 .global usertrapret
5 usertrapret:
6     /* Overlay stack pointer in trapret. */
7     mov sp, x0
8     b trapret
```

虽然暴力，但简单明了。

终于，我们跳转到了函数 `trapret`，其作用主要是载入 `trapframe`，初始化所有寄存器。1.4.4 节中我们提到，寄存器 X30 和 ELR_EL1 设置为 0，其实指的是 `initcode` 在页表中的起始地址；寄存器 SP_EL0 设置为 `PGSIZE`，指的是用户栈的栈底地址，作为栈指针 SP 的初始值；寄存器 SPSR_EL1 设置为 0，表示用户态（EL0）。于是，`trapret` 在异常返回（`eret`）时，将返回到用户态下 `initcode` 的起始地址。至此，用户程序 `initcode` 开始执行。

1.4.6 进程切换：yield

每当时间片耗尽，程序就要被强制暂停执行。这时我们通过 trap 调用函数 `yield` 来切换当前使用 CPU 的程序。trap 的部分我们放在 2.1 节再讲，这里我们只关注进程管理的部分。

函数 `yield` 的工作很简单，就是设置进程状态为 `RUNNABLE`，然后调用函数 `sched`。

```
1 // kern/proc.c
2
3 /*
4  * Give up the CPU for one scheduling round.
5  */
6 void
7 yield()
8 {
9     struct proc* p = thiscpu->proc;
10    acquire(&p->lock);
11    p->state = RUNNABLE;
12    cprintf("yield: proc %d gives up CPU %d.\n", p->pid, cpuid());
13    sched();
14    release(&p->lock);
15 }
```

函数 `sched` 的工作也很简单，就是调用函数 `swtch` 切换 context，执行权回到函数 `scheduler`，然后由 `scheduler` 来决定下一个执行的程序，如此循环。

```
1 // kern/proc.c
2
3 /*
4  * Enter scheduler. Must hold only p->lock
5  * and have changed p->state.
6  */
7 void
8 sched()
9 {
10    struct cpu* c = thiscpu;
11    struct proc* p = c->proc;
12
13    if (!holding(&p->lock)) panic("\tsched: process not locked.\n");
14    if (p->state == RUNNING) panic("\tsched: process running.\n");
15
16    swtch(&p->context, c->scheduler);
17 }
```

1.4.7 进程退出：exit

假设程序已经执行完了，最后系统函数 `sys_exit` 会调用函数 `exit` 退出。函数 `exit` 的主要工作是将当前进程的子进程托管给第一个用户进程 `initproc`，设置进程状态为 `ZOMBIE`，最后调用函数 `sched` 回到 `scheduler`。由于目前我们没有实现函数 `wait`，因此其实进程申请的资源暂时还无法释放，父子进程的概念暂时也没有什么用处，可以先不管。

需要注意的是，暂时我们为了省事，是直接在 `initproc` 里执行的用户程序代码。但显然我们并不应该这么做，因为这样在程序退出时，也就把这第一个用户程序 `initproc` 给退出了，这样其实之后就无法利用 `initproc` 启动新程序了。未来我们实现了文件系统和其他一些必要的函数（如 `fork`，`sleep`，`wait`），真正开始执行用户程序后，这里需要将第 14 行取消注释。

```
1 // kern/proc.c
2
3 /*
4  * Exit the current process. Does not return.
5  * An exited process remains in the zombie state
6  * until its parent calls wait() to find out it exited.
7  */
8 void
9 exit(int status)
10 {
11     struct proc* p = thiscpu->proc;
12
13     // Temporarily disabled before user processes are implemented.
14     // if (p == initproc) panic("\textit: initproc exiting.\n");
15
16     acquire(&wait_lock);
17
18     // Give any children to init.
19     reparent(p);
20
21     acquire(&p->lock);
22     p->xstate = status;
23     p->state = ZOMBIE;
24
25     release(&wait_lock);
26
27     // Jump into the scheduler, never return.
28     sched();
29     panic("\textit: zombie returned!\n");
30 }
```

2. 系统调用

2.1 系统调用模块

目前内核已经支持基本的异常处理，在本实验中还需要进一步完善内核的系统调用模块。

从函数 `trap` 开始说起。在 `trap` 后，如果判断当前为 timer 中断，则调用函数 `yield`，触发进程切换，具体参见 1.4.6 节。如果判断当前为系统调用，则清空寄存器 ESR (Exception Syndrome Register)，设置 `trapframe`，并调用函数 `syscall`。


```

1 // kern/trap.c
2
3 void
4 trap(struct trapframe* tf)
5 {
6     struct proc* p = thiscpu->proc;
7     int src = get32(IRQ_SRC_CORE(cpuid()));
8     int bad = 0;
9
10    if (src & IRQ_CNPNSIRQ) {
11        timer(), timer_reset(), yield();
12    } else if (src & IRQ_TIMER) {
13        clock(), clock_reset();
14    } else if (src & IRQ_GPU) {
15        if (get32(IRQ_PENDING_1) & AUX_INT)
16            uart_intr();
17        else
18            bad = 1;
19    } else {
20        switch (resr() >> EC_SHIFT) {
21            case EC_SVC64:
22                lesr(0); /* Clear esr. */
23                /* Jump to syscall to handle the system call from user process */
24                if (p->killed) exit(1);
25                p->tf = tf;
26                syscall();
27                if (p->killed) exit(1);
28                break;
29            default: bad = 1;
30        }
31    }
32    if (bad) panic("\ttrap: unexpected irq.\n");
33 }

```

函数 `syscall` 根据传入的第一个参数（即寄存器 X0 的值），也就是 system call number 决定跳转到哪个系统函数。这里我们仿照 Xv6 [\[3:2\]](#)，采用了一个函数指针数组 `syscalls` 作为路由。

```

1 // inc/syscallno.h
2
3 #define SYS_exec 0
4 #define SYS_exit 1

```

```

1 // kern/syscall.c
2
3 extern int sys_exec();
4 extern int sys_exit();
5
6 static int (*syscalls[])() = {
7     [SYS_exec] sys_exec,
8     [SYS_exit] sys_exit,
9 };
10
11 int
12 syscall()
13 {
14     struct proc* p = thiscpu->proc;
15     int num = p->tf->x0;
16     if (num ≥ 0 && num < NELEM(syscalls) && syscalls[num]) {
17         cprintf("syscall: proc %d calls syscall %d.\n", p->pid, num);
18         return syscalls[num]();
19     } else {
20         cprintf(
21             "syscall: unknown syscall %d from proc %d (%s).\n", num, p->pid,
22             p->name);
23         return -1;
24     }
25     return 0;
26 }

```

3. 调整主循环

由于 1.4.4 节中踩到的坑，我仔细检查了一遍哪些初始化函数是只能在 CPU0 上被调用一次的。目前的函数

`main` 如下所示：

```

1 // kern/main.c
2
3 volatile static int started = 0;
4
5 void
6 main()
7 {
8     extern char edata[], end[], vectors[];
9
10    if (cpuid() == 0) {
11        memset(edata, 0, end - edata);
12        console_init();
13        cprintf("main: [CPU 0] init started.\n");
14        alloc_init();
15        proc_init();
16        lvbar(vectors);
17        irq_init();
18        timer_init();
19        user_init();
20        started = 1; // allow APs to run
21    } else {
22        while (!started) {}
23        cprintf("main: [CPU %d] init started.\n", cpuid());
24        lvbar(vectors);
25        timer_init();
26    }
27
28    cprintf("main: [CPU %d] init success.\n", cpuid());
29
30    scheduler();
31 }

```

运行结果

```
> make qemu
```

```
console_init: success.
main: [CPU 0] init started.
alloc_init: success.
proc_init: success.
irq_init: success.
timer_init: success at CPU 0.
proc_alloc: proc 1 success.
user_init: proc 1 (initproc) success.
main: [CPU 0] init success.
main: [CPU 3] init started.
main: [CPU 2] init started.
timer_init: success at CPU 3.
main: [CPU 1] init started.
main: [CPU 3] init success.
timer_init: success at CPU 2.
main: [CPU 2] init success.
scheduler: switch to proc 1 at CPU 0.
timer_init: success at CPU 1.
main: [CPU 1] init success.
syscall: proc 1 calls syscall 0.
sys_exec: executing /init with parameters: /init
syscall: proc 1 calls syscall 1.
sys_exit: in exit.
```

测试环境

- OS: Ubuntu 18.04.5 LTS (WSL2 4.19.128-microsoft-standard)
- Compiler: gcc version 8.4.0 (Ubuntu/Linaro 8.4.0-1ubuntu1~18.04)
 - Target: aarch64-linux-gnu
- Debugger: GNU gdb 8.2 (Ubuntu 8.2-0ubuntu1~18.04)
 - Target: aarch64-linux-gnu
- Emulator: QEMU emulator version 5.0.50
- Using GNU Make 4.1

-
1. [xv6: a simple, Unix-like teaching operating system - MIT](#) ↩ ↩ ↩ ↩ ↩
 2. [AArch64 Instruction Set Architecture | Procedure Call Standard – Arm Developer](#) ↩ ↩
 3. [mit-pdos/xv6-public: xv6 OS - GitHub](#) ↩ ↩ ↩
 4. [mit-pdos/xv6-riscv: Xv6 for RISC-V - GitHub](#) ↩ ↩ ↩