

Jos Lab 4: Preemptive Multitasking Lab Report

王勇杰 5120379068

2014.12.7

注：lab 简介以及 questions 都在此文档中，可通过目录直接跳转。

Contexts

1. MULTIPROCESSOR SUPPORT AND COOPERATIVE MULTITASKING	2
Exercise 1.....	2
Question 1	3
Exercise 2.....	3
Exercise 3.....	3
Exercise 4.....	4
Exercise 4.1	4
Question 2	5
Exercise 5.....	5
Question 3	6
Exercise 6.....	6
2. COPY-ON-WRITE FORK	9
Exercise 7.....	9
Exercise 8.....	9
Exercise 9.....	10
Exercise 10	10
Exercise 11	11
3. PREEMPTIVE MULTITASKING AND INTER-PROCESS COMMUNICATION (IPC).....	11
Exercise 12	11
Exercise 13	11
Exercise 14	11
4. CHALLENGE	12

1. Multiprocessor Support and Cooperative Multitasking

在这一部分介绍了如何实现对多处理器的支持。

Exercise 1

Read `boot_aps()` and `mp_main()` in `kern/init.c`, and the assembly code in `kern/mpentry.S`. Make sure you understand the control flow transfer during the bootstrap of APs. Then modify your implementation of `page_init()` in `kern/pmap.c` to avoid adding the page at `MPENTRY_PADDR` to the free list, so that we can safely copy and run AP bootstrap code at that physical address. Your code should pass the updated `check_page_free_list()` test (but might fail the updated `check_kern_pgdir()` test, which we will fix soon).

第一个练习主要是在初始化 `freepage` 链表的时候，如果物理地址在 `MPENTRY_PADDR` 就不能把它加入 `freepage` 链表中，具体只需在 `page_init` 时对插入的空页做一下判断，具体代码如下：

```
333     page_free_list = NULL;
334     for (i = 1; i < npages_basemem; i++)
335     {
336         pages[i].pp_ref = 0;
337         if(i == PGNUM(MPENTRY_PADDR)) continue;
338         pages[i].pp_link = page_free_list;
339         page_free_list = &pages[i];
340     }
341
342     for (i = IOPHYSMEM/PGSIZE; i < EXTPHYSMEM/PGSIZE; ++i)
343         pages[i].pp_ref = 1;
344
345     for (; i < PADDR(boot_alloc(0))/PGSIZE; ++i)
346         pages[i].pp_ref = 1;
347     for (; i < npages; ++i)
348     {
349         pages[i].pp_ref = 0;
350         if(i == PGNUM(MPENTRY_PADDR)) continue;
351         pages[i].pp_link = page_free_list;
352         page_free_list = &pages[i];
353     }
```

即在上面代码的 337 行与 350 行对要插入的页地址做一下判断即可。

Question 1

Compare `kern/mpentry.S` side by side with `boot/boot.S`. Bearing in mind that `kern/mpentry.S` is compiled and linked to run above `KERNBASE` just like everything else in the kernel, what is the purpose of macro `MPBOOTPHYS`? Why is it necessary in `kern/mpentry.S` but not in `boot/boot.S`? In other words, what could go wrong if it were omitted in `kern/mpentry.S`?

`MPBOOTPHYS` 这个宏的作用就是将内核的 Linear Address 转化为 Physical Address。而起必须在 `mpentry.S` 中使用是因为 BSP 启动 AP 时候，不同于 BSP 启动时候的 `KERNBASE` 上 4MB 空间映射，内存地址的计算方式并不相同。

Exercise 2

Modify `mem_init_mp()` (in `kern/pmap.c`) to map per-CPU stacks starting at `KSTACKTOP`, as shown in our revised `inc/memlayout.h`. The size of each stack is `KSTKSIZE` bytes plus `KSTKGAP` bytes of unmapped guard pages. Your code should pass the new check in `check_kern_pgdir()`

这个部分完成对每个 CPU 的映射环境，根据注释计算内存地址后调用 `boot_map_region()` 来进行映射，具体代码如下：

```
290 // LAB 4: Your code here:
291 size_t i;
292 for(i = 0; i < NCPU; ++i)
293 {
294     uint32_t kstacktop_i = KSTACKTOP - i * (KSTKSIZE + KSTKGAP);
295     boot_map_region(kern_pgdir, kstacktop_i - KSTKSIZE, KSTKSIZE,
296                   PADDR(percpu_kstacks[i]), PTE_W);
297 }
298 }
```

Exercise 3

The code in `trap_init_percpu()` (`kern/trap.c`) initializes the TSS and TSS descriptor for the BSP. It worked in Lab 3, but is incorrect when running on other CPUs. Change the code so that it can work on all CPUs. (Note: your new code should not use the global `ts` variable any more.)

```
196 // Get CPU ID
197 int cpuID = thiscpu->cpu_id;
198
199 // Setup a TSS so that we get the right stack
200 // when we trap to the kernel.
201 thiscpu->cpu_ts.ts_esp0 = KSTACKTOP - (KSTKSIZE + KSTKGAP)*cpuID;
202 thiscpu->cpu_ts.ts_ss0 = GD_KD;
203
204 // Initialize the TSS slot of the gdt.
205 int gdtOffset = GD_TSS0 + (cpuID<<3);
206 int gdtIndex = gdtOffset>>3;
207 gdt[gdtIndex] = SEG16(STS_T32A, (uint32_t) &thiscpu->cpu_ts, sizeof
208 gdt[gdtIndex].sd_s = 0;
209
210 // Load the TSS selector (like other segment selectors, the
211 // bottom three bits are special; we leave them 0)
212 ltr(gdtOffset);
213
214 // Load the IDT
215 lidt(&idt_pd);
216
```

这个练习是为了每个 CPU 初始化 TSS 与 TSS 描述符，具体代码如下。

Exercise 4

Apply the big kernel lock as described above, by calling `lock_kernel()` and `unlock_kernel()` at the proper locations.

这个练习要求在一些位置应用全局锁，具体需加锁解锁的位置在 lab 介绍中也说明得很明确。

Exercise 4.1

Implement ticket spinlock in `kern/spinlock.c`. You can define a macro `USE_TICKET_SPIN_LOCK` at the beginning of `kern/spinlock.h` to make it work. After you correctly implement ticket spinlock and define the macro, your code should pass `spinlock_test()`. Don't use ticket spinlock before you implement all the rest exercises. Ticket spinlock may cause stresssched and primes fail because of its poor performance. (Optional) Consider why ticket spinlock is slow here.

这个练习需要完成排队自旋锁，具体代码可见 `spinlock.c`，运行测试输出如下：

```
Pass check_page_free_list
check_page_installed_pgdir() succeeded!
SMP: CPU 0 found 1 CPU(s)
enabled interrupts: 1 2
spinlock_test() succeeded on CPU 0!
000000001 000000001 00001000
```

Question 2

It seems that using the big kernel lock guarantees that only one CPU can run the kernel code at a time. Why do we still need separate kernel stacks for each CPU? Describe a scenario in which using a shared kernel stack will go wrong, even with the protection of the big kernel lock.

在中断发生后拿 kernel 锁之前，原来 environment 的数据会被 push 到栈上，这时候如果共享 kernel stack，就会出现保存现场状态出错。因此需要为每个 CPU 分配不同的内存栈。

Exercise 5

Implement round-robin scheduling in `sched_yield()` as described above. Don't forget to modify `syscall()` to dispatch `sys_yield()`.

Modify `kern/init.c` to create three (or more!) environments that all run the program `user/yield.c`.

After the `yield` programs exit, when only idle environments are runnable, the scheduler should invoke the JOS kernel monitor. If any of this does not happen, then fix your code before proceeding.

这个练习是实现 round-robin 的进程调度，具体代码如下：

```

32     if(curenv == NULL) i = 0;
33     else i = ENVX(curenv->env_id) + 1;
34
35     int times = 0;
36     while(times < NENV){
37         int tem = (i + times) % NENV;
38         if(envs[tem].env_type != ENV_TYPE_IDLE &&
39             envs[tem].env_status == ENV_RUNNABLE){
40             env_run(envs + tem);
41         }
42         times++;
43     }
44
45     if(curenv && curenv->env_status == ENV_RUNNING)
46         env_run(curenv);

```

首先在 32-33 行获取当前正在运行的线程 id，而后在 36-43 行做 NENV 次循环，找到接下来需要运行的线程。

Question 3

In your implementation of `env_run()` you should have called `lcr3()`. Before and after the call to `lcr3()`, your code makes references (at least it should) to the variable `e`, the argument to `env_run`. Upon loading the `%cr3` register, the addressing context used by the MMU is instantly changed. But a virtual address (namely `e`) has meaning relative to a given address context--the address context specifies the physical address to which the virtual address maps. Why can the pointer `e` be dereferenced both before and after the addressing switch?

`e` 指向的是内核地址空间，页表切换时不切换 kernel 的寻址空间的地址。

Exercise 6

Implement the system calls described above in `kern/syscall.c`. You will need to use various functions in `kern/pmap.c` and `kern/env.c`, particularly `envid2env()`. For now, whenever you call `envid2env()`, pass 1 in the `checkperm` parameter. Be sure you check for any invalid system call arguments, returning `-E_INVALID` in that case. Test your JOS kernel with `user/dumbfork` and make sure it works before proceeding.

这个练习需要完成一个 fork 的功能，具体实现如下：

1. `sys_exofork()` :

```
100     struct Env* env;
101     if(env_alloc(&env, thiscpu->cpu_env->env_id))
102         return -E_NO_FREE_ENV;
103
104     env->env_status = ENV_NOT_RUNNABLE;
105     env->env_tf = thiscpu->cpu_env->env_tf;
106     env->env_tf.tf_regs.reg_eax = 0;
107
108     return env->env_id;
```

这个函数的作用就是复制当前的环境，返回新的环境 id。

2. sys_env_set_status() :

```
130     struct Env* env;
131
132     if(status != ENV_RUNNABLE &&
133        status != ENV_NOT_RUNNABLE)
134         return -E_INVAL;
135
136     if(envid2env(envid, &env, 1))
137         return -E_BAD_ENV;
138
139     env->env_status = status;
140     return 0;
```

这个函数用户设置环境状态，根据目前的 status 进行设置。

3. sys_page_alloc() :

```
186     // LAB 4: Your code here.
187     if((uint32_t)va>=UTOP || ROUNDUP(va, PGSIZE)!=va)
188         return -E_INVAL;
189
190     if(!(perm&PTE_U) || !(perm&PTE_P) || (perm&~PTE_SYSCALL))
191         return -E_INVAL;
192
193     struct Env* env;
194     if(envid2env(envid, &env, 1)<0) return -E_BAD_ENV;
195
196     struct Page* page;
197     page = page_alloc(ALLOC_ZERO);
198     if(!page) return -E_NO_MEM;
199
200     if(page_insert(env->env_pgdir, page, va, perm)){
201         page_free(page);
202         return -E_NO_MEM;
203     }
204
205     return 0;
```

这个函数是为 env_id 的环境申请新的物理页，并做虚拟地址的映射。

4. sys_page_map() :

```
237     if((uint32_t)srcva>=UTOP ||
238         ROUNDUP(srcva, PGSIZE)!=srcva ||
239         (uint32_t)dstva>=UTOP ||
240         ROUNDUP(dstva, PGSIZE)!=dstva){
241         return -E_INVAL;
242     }
243
244     if(!(perm&PTE_U) || !(perm&PTE_P) || (perm&~PTE_SYSCALL))
245         return -E_INVAL;
246
247     struct Env* env_s, * env_d;
248     if(envid2env(srcenvid, &env_s, 1)||
249         envid2env(dstenvid, &env_d, 1))
250         return -E_BAD_ENV;
251
252     pte_t* pte;
253     struct Page* page = page_lookup(env_s->env_pgdir, srcva, &pte);
254     if(!page)
255         return -E_INVAL;
256
257     if((perm & PTE_W) && !(*pte & PTE_W))
258         return -E_INVAL;
259
260     if(page_insert(env_d->env_pgdir, page, dstva, perm))
261         return -E_NO_MEM;
262
263     return 0;
```

5. sys_page_unmap() :

```
279     // LAB 4: Your code here.
280     if((uint32_t)va>=UTOP ||
281         ROUNDUP(va, PGSIZE)!=va)
282         return -E_INVAL;
283
284     struct Env* env;
285     if(envid2env(envid, &env, 1))
286         return -E_BAD_ENV;
287
288     page_remove(env->env_pgdir, va);
289     return 0;
```


2. Copy-on-Write Fork

Exercise 7

Implement the `sys_env_set_pgfault_upcall` system call. Be sure to enable permission checking when looking up the environment ID of the target environment, since this is a "dangerous" system call.

```
153 static int
154 sys_env_set_pgfault_upcall(envid_t envid, void *func)
155 {
156     // LAB 4: Your code here.
157     struct Env* env;
158     if(envid2env(envid, &env, 1))
159         return -E_BAD_ENV;
160
161     env->env_pgfault_upcall = func;
162     return 0;
163
164     //panic("sys_env_set_pgfault_upcall not implemented");
165 }
```

这个函数用于设置当发生 page fault 时用哪个函数来处理。

Exercise 8

Implement the code in `page_fault_handler` in `kern/trap.c` required to dispatch page faults to the user-mode handler. Be sure to take appropriate precautions when writing into the exception stack. (What happens if the user environment runs out of space on the exception stack?)

```

440     if(curenv->env_pgfault_upcall){
441
442         uint32_t stktop = (tf->tf_esp >= UXSTACKTOP - PGSIZE &&
443             tf->tf_esp < UXSTACKTOP) ?
444             tf->tf_esp - sizeof(struct UTrapframe) - 4 :
445             UXSTACKTOP - sizeof(struct UTrapframe);
446
447         user_mem_assert(curenv, (void*) stktop, sizeof(struct UTrapframe), PTE_U|PTE_W);
448
449         struct UTrapframe* utf = (struct UTrapframe*) stktop;
450         utf->utf_eflags = tf->tf_eflags;
451         utf->utf_eip = tf->tf_eip;
452         utf->utf_err = tf->tf_err;
453         utf->utf_esp = tf->tf_esp;
454         utf->utf_fault_va = fault_va;
455         utf->utf_regs = tf->tf_regs;
456
457         tf->tf_eip = (uint32_t) curenv->env_pgfault_upcall;
458         tf->tf_esp = stktop;
459
460         env_run(curenv);

```

这个函数为当发生用户态下的 page fault 时，如果该进程的环境设置了 upcall，那么就将 eip 指向 upcall，并在调用页处理函数前将当前状态压入用户异常栈。

Exercise 9

Implement the `_pgfault_upcall` routine in `lib/pfentry.S`. The interesting part is returning to the original point in the user code that caused the page fault. You'll return directly there, without going back through the kernel. The hard part is simultaneously switching stacks and re-loading the EIP.

这个练习需要实现页错误处理函数的入口，具体实现见 `lib/pfentry.S`。

Exercise 10

Finish `set_pgfault_handler()` in `lib/pgfault.c`

```

// LAB 4: Your code here.
if(sys_page_alloc(0, (void*) (UXSTACKTOP - PGSIZE), PTE_W|PTE_U|PTE_P)){
    cprintf("lib/pgfault:set_pgfault_handler: out of memory.");
    return;
}

sys_env_set_pgfault_upcall(0, _pgfault_upcall);

```

Exercise 11

Implement `fork`, `duppage` and `pgfault` in `lib/fork.c`

这里需要实现 `fork` 的一些主要函数，其中 `duppage` 主要将页 `pn` 映射到 `envid` 中，`pgfault` 为用户页错误处理函数，`fork` 函数则在 `exofork` 的基础上做拓展，当发生 `page fault` 时可以由用户定义函数来处理。具体代码实现见 `lib/fork.c`。

3. Preemptive Multitasking and Inter-Process communication (IPC)

Exercise 12

Modify `kern/trapentry.S` and `kern/trap.c` to initialize the appropriate entries in the IDT and provide handlers for IRQs 0 through 15. Then modify the code in `env_alloc()` in `kern/env.c` to ensure that user environments are always run with interrupts enabled.

这个练习需要初始化 IRQ，详见具体代码。

Exercise 13

Modify the kernel's `trap_dispatch()` function so that it calls `sched_yield()` to find and run a different environment whenever a clock interrupt takes place.

这个练习要求在发生时钟中断时调用 `sched_yield`，具体代码如下

```
if(tf->tf_trapno == IRQ_OFFSET + IRQ_TIMER){
    lapic_eoi();
    sched_yield();
    return;
}
```

Exercise 14

Implement `sys_ipc_recv` and `sys_ipc_try_send` in `kern/syscall.c`.

这个练习需要实现进程之间的通信，具体实现见代码部分。

4. Challenge

Add a less trivial scheduling policy to the kernel, such as a fixed-priority scheduler that allows each environment to be assigned a priority and ensures that higher-priority environments are always chosen in preference to lower-priority environments. If you're feeling really adventurous, try implementing a Unix-style adjustable-priority scheduler or even a lottery or stride scheduler. (Look up "lottery scheduling" and "stride scheduling" in Google.)

加入了调整程序自己优先级的 system call，每次 sched_yield() 中，会优先搜索最高优先级的进程运行，然后默认，然后低优先级。测试程序通过修改 user/hello.c 的代码实现。测试结果如下

```
[00001008] High Priority
[00001008] High Priority
[00001008] High Priority
[00001008] High Priority
[00001008] exiting gracefully
[00001008] free env 00001008
[00001009] new env 00002008
[0000100a] Default Priority
[0000100a] Default Priority
[0000100a] Default Priority
[0000100a] Default Priority
[0000100a] exiting gracefully
[0000100a] free env 0000100a
[00001009] Default Priority
[00001009] Default Priority
[00001009] Default Priority
[00001009] Default Priority
[00001009] Default Priority
[00001009] exiting gracefully
[00001009] free env 00001009
[00002008] Low Priority
[00002008] Low Priority
```