

Lab10: mmap 实验报告

mmap是一种内存映射文件的方法，即将一个文件或者其它对象映射到进程的地址空间，实现文件磁盘地址和进程虚拟地址空间中一段虚拟地址的一一对映关系。

实现这样的映射关系后，进程就可以采用指针的方式读写操作这一段内存，而操作系统会自动回写脏页面到对应的文件磁盘上，即完成了对文件的操作而不必再调用read、write等系统调用函数。相应的，内核空间对这段区域的修改也直接反映用户空间，从而可以实现不同进程间的文件共享。

在本实验中，需要向 xv6 添加 `mmap` 和 `munmap` 系统调用。

- `mmap` 函数是将文件地址映射到虚拟内存中，返回映射后的地址，同时记录该进程映射到的文件信息。
- `munmap` 函数就是取消进程地址空间中，文件地址某一部分的映射。

1. 首先切换分支：

```
1 $ git fetch
2 $ git checkout mmap
3 $ make clean
```

mmap(hard)

1. 添加系统调用：

- 在 `user/user.h` 中添加系统调用声明：

```
26 void* mmap(void *addr, int length, int prot, int flags, int fd, uint offset);
27 int munmap(void *addr, int length);
```

- 在 `user/usys.pl` 中添加 entry：

```
39 entry("mmap");
40 entry("munmap");
```

- 在 Makefile 中添加编译声明:

```
178 $U/_zombie\
179 $U/_mmaptest\
```

- 在 kernel/syscall.h 中添加系统调用码:

```
23 #define SYS_mmap 22
24 #define SYS_munmap 23
```

- 在 kernel/syscall.c 中添加系统调用函数引用:

```
107 extern uint64 sys_mmap(void);
108 extern uint64 sys_munmap(void);
109
110 static uint64 (*syscalls[])(void) = {
111 [SYS_fork] sys_fork,
112 [SYS_exit] sys_exit,
113 [SYS_wait] sys_wait,
114 [SYS_pipe] sys_pipe,
115 [SYS_read] sys_read,
116 [SYS_kill] sys_kill,
117 [SYS_exec] sys_exec,
118 [SYS_fstat] sys_fstat,
119 [SYS_chdir] sys_chdir,
120 [SYS_dup] sys_dup,
121 [SYS_getpid] sys_getpid,
122 [SYS_sbrk] sys_sbrk,
123 [SYS_sleep] sys_sleep,
124 [SYS_uptime] sys_uptime,
125 [SYS_open] sys_open,
126 [SYS_write] sys_write,
127 [SYS_mknod] sys_mknod,
128 [SYS_unlink] sys_unlink,
129 [SYS_link] sys_link,
130 [SYS_mkdir] sys_mkdir,
131 [SYS_close] sys_close,
132 [SYS_mmap] sys_mmap,
133 [SYS_munmap] sys_munmap,
134 };
```

2. 在 kernel/proc.h 中构建一个 VMA 结构体数组, 长度为 16

```

83 #define VMASIZE 16
84 struct vma {
85     int used;
86     uint64 addr;
87     int length;
88     int prot;
89     int flags;
90     int fd;
91     int offset;
92     struct file *file;
93 };
94
95 enum procstate { UNUSED, SLEEPING, RUNNABLE, RUNNING, ZOMBIE };
96
97 // Per-process state
98 struct proc {
99     struct spinlock lock;
100
101     // p->lock must be held when using these:
102     enum procstate state; // Process state
103     struct proc *parent; // Parent process
104     void *chan; // If non-zero, sleeping on chan
105     int killed; // If non-zero, have been killed
106     int xstate; // Exit status to be returned to parent's wait
107     int pid; // Process ID
108
109     // these are private to the process, so p->lock need not be held.
110     uint64 kstack; // Virtual address of kernel stack
111     uint64 sz; // Size of process memory (bytes)
112     pagetable_t pagetable; // User page table
113     struct trapframe *trapframe; // data page for trampoline.S
114     struct context context; // swtch() here to run process
115     struct file *ofile[NOFILE]; // Open files
116     struct inode *cwd; // Current directory
117     char name[16]; // Process name (debugging)
118     struct vma vma[VMASIZE]; // vma of file system
119 };

```

3. 在 kernel/sysfile.c 中添加 sys_mmap 函数的声明和实现：

首先接收传来的参数，判断参数的合法性，然后遍历 VMA 数组，找到还没有使用的 vma，将参数信息添加进去。这里映射的虚拟地址，可以直接填写堆的最高地址，然后让堆继续生长。

```

488     uint64
489     sys_mmap(void)
490     {
491         uint64 addr;
492         int length, prot, flags, fd, offset;
493         struct file *file;
494         struct proc *p = myproc();
495         if(argaddr(0, &addr) || argint(1, &length) || argint(2, &prot) ||
496             argint(3, &flags) || argfd(4, &fd, &file) || argint(5, &offset)) {
497             return -1;
498         }
499         if(!file->writable && (prot & PROT_WRITE) && flags == MAP_SHARED)
500             return -1;
501         length = PGROUNDUP(length);
502         if(p->sz > MAXVA - length)
503             return -1;
504         for(int i = 0; i < VMASIZE; i++) {
505             if(p->vma[i].used == 0) {
506                 p->vma[i].used = 1;
507                 p->vma[i].addr = p->sz;
508                 p->vma[i].length = length;
509                 p->vma[i].prot = prot;
510                 p->vma[i].flags = flags;
511                 p->vma[i].fd = fd;
512                 p->vma[i].file = file;
513                 p->vma[i].offset = offset;
514                 filedup(file);
515                 p->sz += length;
516                 return p->vma[i].addr;
517             }
518         }
519         return -1;
520     }

```

4. 在kernel/trap.c中修改usertrap函数，实现trap中断处理：

与lab5相似，由于是懒加载，在读取或写入相应的虚拟地址时，会存在地址未映射的情况。这时需要将物理地址上的数据读到虚拟地址中，然后重新进行读取或写入操作

```

72 } else if((which_dev = devintr()) != 0){
73     // ok
74 } else if(r_scause() == 13 || r_scause() == 15) {
75     uint64 va = r_stval();
76     if(va >= p->sz || va > MAXVA || PGROUNDUP(va) == PGROUNDDOWN(p->trapframe->sp)) p->killed = 1;
77     else {
78         struct vma *vma = 0;
79         for (int i = 0; i < VMASIZE; i++) {
80             if (p->vma[i].used == 1 && va >= p->vma[i].addr && va < p->vma[i].addr + p->vma[i].length) {
81                 vma = &p->vma[i];
82                 //printf("trap mmap: %d, addr: %d, used:%d\n", i, p->vma[i].addr, p->vma[i].used);
83                 break;
84             }
85         }
86         if(vma) {
87             va = PGROUNDDOWN(va);
88             uint64 offset = va - vma->addr;
89             uint64 mem = (uint64)kalloc();
90             if(mem == 0) {
91                 p->killed = 1;
92             } else {
93                 memset((void*)mem, 0, PGSIZE);
94                 ilock(vma->file->ip);
95                 readi(vma->file->ip, 0, mem, offset, PGSIZE);
96                 iunlock(vma->file->ip);
97                 int flag = PTE_U;
98                 if(vma->prot & PROT_READ) flag |= PTE_R;
99                 if(vma->prot & PROT_WRITE) flag |= PTE_W;
100                 if(vma->prot & PROT_EXEC) flag |= PTE_X;
101                 if(mappages(p->pagetable, va, PGSIZE, mem, flag) != 0) {
102                     kfree((void*)mem);
103                     p->killed = 1;
104                 }
105             }
106         }
107     }
108 } else {
109     printf("usertrap(): unexpected scause %p pid=%d\n", r_scause(), p->pid);

```

5. 在kernel/vm.c文件中修改uvmunmap和uvmcopy函数，由于一些地址并没有进行映射，因此在 walk 的时候，遇到这些地址直接跳过即可：

```

162 void
163 uvmunmap(pagetable_t pagetable, uint64 va, uint64 npages, int do_free)
164 {
165     uint64 a;
166     pte_t *pte;
167
168     if((va % PGSIZE) != 0)
169         panic("uvmunmap: not aligned");
170
171     for(a = va; a < va + npages*PGSIZE; a += PGSIZE){
172         if((pte = walk(pagetable, a, 0)) == 0)
173             panic("uvmunmap: walk");
174         if((*pte & PTE_V) == 0)
175             continue;
176         // panic("uvmunmap: not mapped");
177         if(PTE_FLAGS(*pte) == PTE_V)
178             panic("uvmunmap: not a leaf");
179         if(do_free){
180             uint64 pa = PTE2PA(*pte);
181             kfree((void*)pa);
182         }
183         *pte = 0;
184     }
185 }

```

```

298 int
299 uvmcopy(pagetable_t old, pagetable_t new, uint64 sz)
300 {
301     pte_t *pte;
302     uint64 pa, i;
303     uint flags;
304     char *mem;
305
306     for(i = 0; i < sz; i += PGSIZE){
307         if((pte = walk(old, i, 0)) == 0)
308             panic("uvmcopy: pte should exist");
309         if((*pte & PTE_V) == 0)
310             continue;
311         // panic("uvmcopy: page not present");
312         pa = PTE2PA(*pte);
313         flags = PTE_FLAGS(*pte);
314         if((mem = kalloc()) == 0)
315             goto err;
316         memmove(mem, (char*)pa, PGSIZE);
317         if(mappages(new, i, PGSIZE, (uint64)mem, flags) != 0){
318             kfree(mem);
319             goto err;
320         }
321     }
322     return 0;
323
324 err:
325     uvmunmap(new, 0, i / PGSIZE, 1);
326     return -1;
327 }

```

6. 在 kernel/sysfile.c 中添加sys_munmap函数的声明和实现:

```

522 uint64
523 sys_munmap(void)
524 {
525     uint64 addr;
526     int length;
527     struct proc *p = myproc();
528     struct vma *vma = 0;
529     if(argaddr(0, &addr) || argint(1, &length))
530         return -1;
531     addr = PGROUNDDOWN(addr);
532     length = PGROUNDUP(length);
533     for(int i = 0; i < VMASIZE; i++) {
534         if (addr >= p->vma[i].addr || addr < p->vma[i].addr + p->vma[i].length) {
535             vma = &p->vma[i];
536             break;
537         }
538     }
539     if(vma == 0) return 0;
540     if(vma->addr == addr) {
541         vma->addr += length;
542         vma->length -= length;
543         if(vma->flags & MAP_SHARED)
544             filewrite(vma->file, addr, length);
545         uvmunmap(p->pagetable, addr, length/PGSIZE, 1);
546         if(vma->length == 0) {
547             fileclose(vma->file);
548             vma->used = 0;
549         }
550     }
551     return 0;
552 }

```

7. 在kernel/proc.c中修改fork和exit函数，在进程创建和退出时，复制和清空相应的文件映射:

```

300     safestrcpy(np->name, p->name, sizeof(p->name));
301
302     pid = np->pid;
303
304     np->state = RUNNABLE;
305
306     for(int i = 0; i < VMASIZE; i++) {
307         if(p->vma[i].used){
308             memmove(&(np->vma[i]), &(p->vma[i]), sizeof(p->vma[i]));
309             filedup(p->vma[i].file);
310         }
311     }
312
313     release(&np->lock);
314
315     return pid;
316 }

```

```

355 // Close all open files.
356 for(int fd = 0; fd < NOFILE; fd++){
357     if(p->ofile[fd]){
358         struct file *f = p->ofile[fd];
359         fileclose(f);
360         p->ofile[fd] = 0;
361     }
362 }
363
364 for(int i = 0; i < VMASIZE; i++) {
365     if(p->vma[i].used) {
366         if(p->vma[i].flags & MAP_SHARED)
367             filewrite(p->vma[i].file, p->vma[i].addr, p->vma[i].length);
368         fileclose(p->vma[i].file);
369         uvmunmap(p->pagetable, p->vma[i].addr, p->vma[i].length/PGSIZE, 1);
370         p->vma[i].used = 0;
371     }
372 }
373
374 begin_op();
375 input(p->cwd);
376 end_op();
377 p->cwd = 0;
378
379 // we might re-parent a child to init. we can't be precise about
380 // waking up init, since we can't acquire its lock once we've
381 // acquired any other proc lock. so wake up init whether that's
382 // necessary or not. init may miss this wakeup, but that seems
383 // harmless.
384 acquire(&initproc->lock);
385 wakeup1(initproc);
386 release(&initproc->lock);

```

8. 执行测试命令: `mmapest`

```

220110512@comp4:~/xv6-labs-2020$ make qemu
qemu-system-riscv64 -machine virt -bios none -kernel kernel/kernel -m 128M -smp 3 -nographic -drive file=fs.img,if=none,format=raw,id=x0
-device virtio-blk-device,drive=x0,bus=virtio-mmio-bus.0

xv6 kernel is booting

hart 2 starting
hart 1 starting
init: starting sh
$ mmaptest
mmap_test starting
test mmap f
test mmap f: OK
test mmap private
test mmap private: OK
test mmap read-only
test mmap read-only: OK
test mmap read/write
test mmap read/write: OK
test mmap dirty
test mmap dirty: OK
test not-mapped unmap
test not-mapped unmap: OK
test mmap two files
test mmap two files: OK
mmap_test: ALL OK
fork_test starting
fork_test OK
mmaptest: all tests succeeded
$

```

执行测试命令： `usertests`

```

test sbrkbasic: OK
test sbrkmuch: OK
test kernmem: OK
test sbrkfail: OK
test sbrkarg: OK
test validatetest: OK
test stacktest: OK
test opentest: OK
test writetest: OK
test writebig: OK
test createtest: OK
test openinput: OK
test exitinput: OK
test input: OK
test mem: OK
test pipe1: OK
test preempt: kill... wait... OK
test exitwait: OK
test rmdot: OK
test fourteen: OK
test bigfile: OK
test dirfile: OK
test iref: OK
test forktest: OK
test bigdir: OK
ALL TESTS PASSED
$

```

结果截图

执行命令 `make grade`


```
make[1]: Leaving directory '/home/students/220110512/xv6-labs-2020'
== Test running mmaptest ==
$ make qemu-gdb
(5.4s)
== Test  mmaptest: mmap f ==
mmaptest: mmap f: OK
== Test  mmaptest: mmap private ==
mmaptest: mmap private: OK
== Test  mmaptest: mmap read-only ==
mmaptest: mmap read-only: OK
== Test  mmaptest: mmap read/write ==
mmaptest: mmap read/write: OK
== Test  mmaptest: mmap dirty ==
mmaptest: mmap dirty: OK
== Test  mmaptest: not-mapped unmap ==
mmaptest: not-mapped unmap: OK
== Test  mmaptest: two files ==
mmaptest: two files: OK
== Test  mmaptest: fork_test ==
mmaptest: fork_test: OK
== Test usertests ==
$ make qemu-gdb
usertests: OK (151.3s)
== Test time ==
time: OK
Score: 140/140
220110512@comp4:~/xv6-labs-2020$
```