
Assignment Report: Implementation of Memory Mapping in xv6

Filbert Hamijoyo - 122040012

1 Introduction

In this assignment, I implemented memory mapping functionalities in the xv6 operating system, specifically focusing on the `mmap` and `munmap` system calls. The primary objective was to enhance xv6's memory management capabilities by enabling processes to map files or devices directly into their address space. This facilitates efficient file I/O operations and inter-process communication by leveraging the operating system's memory management unit.

My work involved modifying the kernel to handle memory mapping requests, ensuring proper synchronization and protection mechanisms, and implementing robust error handling. Additionally, I used a given comprehensive test program, `mmaptest.c`, to validate the correctness and reliability of the implemented features. The project required a deep understanding of low-level memory management, system calls, and process synchronization within the xv6 environment.

2 Design

The design of the memory mapping functionality in xv6 encompasses several key components:

2.1 System Calls Implementation

2.1.1 `sys_mmap`

The `sys_mmap` system call facilitates the mapping of files or devices into a process's address space. The implementation involves:

- **Argument Validation:** Ensuring that the size and offset are valid and properly aligned.
- **Address Allocation:** Selecting an appropriate virtual address if none is specified, ensuring page alignment.
- **VMA Slot Allocation:** Allocating a Virtual Memory Area (VMA) slot to track the mapping details.
- **Page Mapping:** Allocating physical memory pages and mapping them into the process's page table with the specified protections and flags.
- **File Duplication:** Handling file reference counts appropriately.

Code Implementation:

```
1 uint64 sys_mmap(void) {  
2     uint64 addr, size, offset;  
3     int protection, flags, fd;  
4     struct file *f = NULL;
```

```
5     struct proc *p = myproc();
6
7     // Fetch arguments
8     argaddr(0, &addr);
9     argaddr(1, &size);
10    argint(2, &protection);
11    argint(3, &flags);
12    if (argfd(4, &fd, &f) < 0) {
13        return -1;
14    }
15    argaddr(5, &offset);
16
17    // Validate size and alignment
18    if (size <= 0 || offset % PGSIZE != 0) {
19        return -1;
20    }
21
22    // Choose address if unspecified
23    if (addr == 0) {
24        addr = PGROUNDUP(p->sz);
25    } else if (addr % PGSIZE != 0) {
26        return -1; // Ensure page alignment
27    }
28
29    // Map file offset and check for writable MAP_SHARED
30    if ((protection & PROT_WRITE) && (flags & MAP_SHARED) && !f->
        writable) {
31        return -1;
32    }
33
34    // Find a free VMA slot
35    struct vma *vma = NULL;
36    for (int i = 0; i < VMASIZE; i++) {
37        if (!p->vma[i].valid) {
38            vma = &p->vma[i];
39            break;
40        }
41    }
42    if (!vma) {
43        return -1;
44    }
45
46    // Set VMA fields
47    vma->start = addr;
48    vma->length = size;
49    vma->prot = protection;
50    vma->flags = flags;
51    vma->offset = offset;
52    vma->f = f;
53    vma->valid = 1;
54
55    if (f) filedup(f);
56
57    // Expand process size if necessary
58    if (addr + size > p->sz) p->sz = addr + size;
59
```

```

60     // Page permissions
61     int perm = PTE_U;
62     if (protection & PROT_READ) perm |= PTE_R;
63     if (protection & PROT_WRITE) perm |= PTE_W;
64     if (protection & PROT_EXEC) perm |= PTE_X;
65
66     // Map pages with file contents, handling MAP_PRIVATE and
        MAP_SHARED
67     begin_op();
68     for (uint64 va = addr; va < addr + size; va += PGSIZE) {
69         uint64 pa = (uint64)kalloc();
70         if (pa == 0) {
71             end_op();
72             return -1;
73         }
74         memset((void *)pa, 0, PGSIZE);
75
76         // Copy content from file starting at the specified offset
77         if (flags & (MAP_SHARED | MAP_PRIVATE)) {
78             ilock(f->ip);
79             int n = readi(f->ip, 0, pa, offset + (va - addr), PGSIZE);
80             iunlock(f->ip);
81
82             if (n < 0) {
83                 end_op();
84                 kfree((void *)pa);
85                 return -1;
86             }
87         }
88
89         if (mappages(p->pagetable, va, PGSIZE, pa, perm) < 0) {
90             end_op();
91             kfree((void *)pa);
92             return -1;
93         }
94     }
95     end_op();
96
97     return addr;
98 }

```

Listing 1: Implementation of `sys_mmap`**Explanation:**

The `sys_mmap` function begins by retrieving and validating the arguments provided by the user. It ensures that the requested size and offset are appropriate and aligned to page boundaries. If no address is specified, it automatically selects the next available address space. It then locates a free VMA slot to store the mapping details.

For each page to be mapped, it allocates physical memory, initializes it, and copies the relevant file content if necessary. The pages are then mapped into the process's page table with the specified protections. Proper error handling ensures that any failures during the mapping process result in a clean exit.

2.1.2 sys_munmap

The `sys_munmap` system call handles the unmapping of previously mapped memory regions. The key steps include:

- **Argument Retrieval:** Obtaining the address and length of the region to unmap.
- **VMA Identification:** Locating the corresponding VMA slot that overlaps with the specified region.
- **Page Unmapping:** Removing the page mappings from the process's page table and freeing the associated physical memory.
- **VMA Slot Management:** Adjusting or splitting the VMA slot as necessary to reflect the unmapped region.
- **Dirty Page Handling:** For `MAP_SHARED` mappings, ensuring that any modified pages are written back to the file.

Code Implementation:

```

1  uint64
2  sys_munmap(void) {
3      struct proc *p = myproc();
4      uint64 addr, length;
5
6      // Retrieve arguments
7      argaddr(0, &addr);
8      argaddr(1, &length);
9
10     // Validate length
11     if (length <= 0) {
12         return -1;
13     }
14
15     // Round up length to nearest page boundary
16     uint64 unmap_limit = PGROUNDUP(addr + length);
17
18     // Iterate through VMAs to locate overlapping regions
19     for (int i = 0; i < VMASIZE; i++) {
20         struct vma *vma = &p->vma[i];
21         if (vma->valid && !(unmap_limit <= vma->start || addr >= vma->
22             start + vma->length)) {
23
24             uint64 start_unmap = addr > vma->start ? addr : vma->start
25                 ;
26             uint64 end_unmap = unmap_limit < (vma->start + vma->length
27                 ) ? unmap_limit : (vma->start + vma->length);
28
29             // If VMA is MAP_SHARED, write back modified pages to the
30             // file
31             if (vma->f && (vma->flags & MAP_SHARED)) {
32                 begin_op();
33                 for (uint64 va = start_unmap; va < end_unmap; va +=
34                     PGSIZE) {
35                     pte_t *pte = walk(p->pagetable, va, 0);
36                     if (pte && (*pte & PTE_V) && (*pte & PTE_R)) {
37                         uint64 pa = PTE2PA(*pte);
38                         char *kv = (char *)pa;

```

```

35         // Write back page and clear dirty bit
36         ilock(vma->f->ip);
37         int n = writei(vma->f->ip, 0, (uint64)kv, vma
38             ->offset + (va - vma->start), PGSIZE);
39         if (n != PGSIZE) {
40             printf("sys_munmap: writei failed for va=0
41                 x%x\n", va);
42         }
43         *pte &= ~PTE_R; // Clear dirty bit
44         iunlock(vma->f->ip);
45     }
46     end_op();
47 }
48
49 // Free pages within this region
50 int npages = (end_unmap - start_unmap) / PGSIZE;
51 uvmunmap(p->pagetable, start_unmap, npages, 1);
52
53 // Adjust or split the VMA as necessary
54 if (start_unmap == vma->start && end_unmap == vma->start +
55     vma->length) {
56     vma->valid = 0;
57     if (vma->f) {
58         fclose(vma->f);
59         vma->f = 0;
60     }
61 } else if (start_unmap == vma->start) {
62     vma->start = end_unmap;
63     vma->offset += end_unmap - vma->start;
64 } else if (end_unmap == vma->start + vma->length) {
65     vma->length = start_unmap - vma->start;
66 } else {
67     int j;
68     for (j = 0; j < VMASIZE; j++) {
69         if (!p->vma[j].valid) break;
70     }
71     if (j == VMASIZE) {
72         printf("sys_munmap: No free VMA slots available\n"
73             );
74         return -1;
75     }
76
77     p->vma[j] = *vma;
78     p->vma[j].start = end_unmap;
79     p->vma[j].length -= (end_unmap - vma->start);
80     p->vma[j].offset += end_unmap - vma->start;
81     vma->length = start_unmap - vma->start;
82 }
83 }
84
85 return 0;
86 }

```

Listing 2: Implementation of sys_munmap

Explanation:

The `sys_munmap` function begins by retrieving the address and length of the memory region to unmap. It then rounds up the length to ensure page alignment. The function iterates through the process's VMAs to find overlapping regions with the specified address range.

For `MAP_SHARED` mappings, it writes back any modified (dirty) pages to the underlying file before unmapping. It then proceeds to remove the page mappings from the process's page table and frees the associated physical memory. If the unmapping spans the entire VMA, the VMA slot is marked as invalid. Otherwise, the VMA is adjusted or split to accurately reflect the remaining mapped regions.

2.2 Page Fault Handling

To support memory-mapped regions, the page fault handler in `trap.c` was extended to manage faults arising from these mappings. The enhancements include:

- **VMA Lookup:** Determining if the faulting address lies within a valid VMA.
- **Copy-On-Write (COW):** Implementing COW semantics for `MAP_PRIVATE` mappings to ensure that modifications do not affect the underlying file.
- **Page Allocation and Mapping:** Allocating new physical pages as needed and updating the page tables with appropriate permissions.
- **Error Handling:** Ensuring that any invalid accesses result in process termination to maintain system stability.

Code Implementation:

```

1  else if (r_scause() == 13 || r_scause() == 15) { // Page fault
2      handling
3          uint64 va = r_stval(); // faulting address
4          uint64 fault_page_start = PGROUNDDOWN(va);
5
6          struct vma *vma = NULL;
7          for (int i = 0; i < VMASIZE; i++) {
8              if (p->vma[i].valid && va >= p->vma[i].start && va < p->
9                  vma[i].start + p->vma[i].length) {
10                 vma = &p->vma[i];
11                 break;
12             }
13         }
14
15         if (vma == NULL) {
16             printf("usertrap: no VMA found for address 0x%lx, pid=%d\n",
17                 va, p->pid);
18             setkilled(p);
19             goto err;
20         }
21
22         bool is_write = (r_scause() == 15); // Write page fault
23
24         // Check if page is mapped
25         pte_t *pte = walk(p->pagetable, fault_page_start, 0);
26         if (pte == 0 || (*pte & PTE_V) == 0) {
27             // Page not yet mapped, proceed with allocation
28             goto allocate_page;
29         }

```

```

28     // Handle Copy-On-Write for MAP_PRIVATE
29     if (is_write && (vma->flags & MAP_PRIVATE)) {
30         uint64 pa = (*pte) >> 10;
31         pa = pa << 12;
32
33         char *old_mem = (char*)pa;
34
35         char *new_mem = kalloc();
36         if (new_mem == 0) {
37             printf("usertrap: kalloc failed for address 0x%lx, pid
38                 =%d\n", va, p->pid);
39             setkilled(p);
40             goto err;
41         }
42
43         // Copy the contents to the new page
44         memmove(new_mem, old_mem, PGSIZE);
45
46         uint64 new_pa = (uint64)new_mem;
47
48         // Update the PTE to point to the new physical page with
49         // write permissions
50         *pte = (new_pa >> 12) << 10 | ((*pte) & 0x3FF) | PTE_W;
51         sfence_vma(); // Flush TLB
52     } else {
53         allocate_page: {
54             char *mem_kva = kalloc();
55             if (mem_kva == 0) {
56                 printf("usertrap: kalloc failed for address 0x%lx, pid
57                     =%d\n", va, p->pid);
58                 setkilled(p);
59                 goto err;
60             }
61
62             memset(mem_kva, 0, PGSIZE);
63
64             if (vma->f != NULL) {
65                 uint64 file_offset = vma->offset + (fault_page_start -
66                     vma->start);
67                 mapfile(vma->f, mem_kva, file_offset);
68             }
69
70             // Determine page permissions based on VMA protection
71             // flags
72             int perm = PTE_U | PTE_V;
73             if (vma->prot & PROT_READ)
74                 perm |= PTE_R;
75             if (vma->prot & PROT_WRITE && !(vma->flags & MAP_PRIVATE))
76                 perm |= PTE_W; // Allow write only if not MAP_PRIVATE
77             if (vma->prot & PROT_EXEC)
78                 perm |= PTE_X;
79
80             uint64 pa = (uint64)mem_kva;

```

```

78         if (mappages(p->pagetable, fault_page_start, PGSIZE, pa,
79             perm) < 0) {
80             printf("usertrap: mappages failed for address 0x%lx,
81                 pid=%d\n", va, p->pid);
82             kfree(mem_kva);
83             setkilled(p);
84             goto err;
85         }
86         sfence_vma();
87     } else {
88         printf("usertrap(): unexpected scause 0x%lx pid=%d\n",
89             r_scause(), p->pid);
90         printf("             sepc=0x%lx stval=0x%lx\n", r_sepc(),
91             r_stval());
92         setkilled(p);
93     }
94     if (killed(p))
95         exit(-1);
96     if (which_dev == 2)
97         yield();
98     usertrapret();
99 err:
100     if (killed(p))
101         exit(-1);
102     usertrapret();
103 }

```

Listing 3: Page Fault Handling in trap.c

Explanation:

The enhanced page fault handler begins by identifying the faulting virtual address and locating the corresponding VMA. If no VMA is found for the address, the process is terminated to maintain system integrity. For write faults (MAP_PRIVATE), the handler implements Copy-On-Write (COW) by allocating a new physical page, copying the contents from the original page, and updating the page table entry to point to the new page with write permissions.

For read faults or non-COW scenarios, the handler allocates a new physical page, initializes it, and maps it into the process's address space. Proper synchronization and error handling ensure that any issues during page allocation or mapping result in safe termination of the offending process.

2.3 Virtual Memory Areas (VMA)

A VMA structure was introduced to keep track of memory-mapped regions:

Code Implementation:

```

1 struct vma {
2     uint64 start;           // Starting virtual address of the VMA
3     uint64 length;         // Length of the memory region in bytes
4     int prot;               // Protection flags (e.g., PROT_READ,
5                             PROT_WRITE)
6     int flags;              // Mapping flags (e.g., MAP_SHARED,
7                             MAP_PRIVATE)

```



```
6     uint64 offset;           // Offset within the file
7     struct file *f;         // Pointer to the mapped file
8     int valid;              // Indicates if the VMA slot is in use
9 };
```

Listing 4: Definition of VMA Structure

Explanation:

The `struct vma` defines the properties of a virtual memory area, including its start address, length, protection and mapping flags, file offset, and associated file pointer. The `valid` flag indicates whether the VMA slot is currently in use. Each process maintains an array of such VMAs (`struct vma vma[VMASIZE];`), allowing multiple memory-mapped regions to coexist within a single process's address space.

3 Environment and Execution

3.1 Running Environment

The implementation was developed and tested using the following setup:

- **Virtual Machine:** UTM
- **Image:** CSC3150_a3_xv6.qcow2
- **Network Mode:** Emulated VLAN
- **Port Forwarding:** Host port 2200 mapped to guest port 22 for SSH access

3.2 Setup and Access

To set up and access the virtual machine, follow these steps:

1. **Login Credentials:**
 - **Username:** csc3150
 - **Password:** csc3150
2. **Assign IP Address:**
 - Execute `sudo dhclient` to obtain an IP address.
 - Verify the assignment using `ip a`.
3. **SSH Access:**
 - Connect via SSH using the command:

```
1 ssh -p 2200 csc3150@127.0.0.1
```

Listing 5: SSH Command

3.3 Compilation and Execution

The project was compiled and executed as follows:

1. Navigate to the `xv6-labs-2022` directory:

```
1 cd xv6-labs-2022
```

Listing 6: Navigation Command

2. Compile and run `xv6` using `QEMU`:

```
1 make clean
2 make qemu
```

Listing 7: Compilation Commands

3. Within the xv6 environment, execute the test program:

```
1 mmaptest
```

Listing 8: Execution Command

3.4 Execution Output

The successful execution of the test program is demonstrated in the screenshot below.

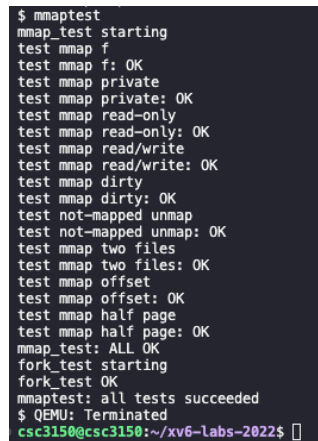
A terminal window showing the output of the 'mmaptest' program. The output consists of a series of test cases and their results, all of which are 'OK'. The tests include: mmap starting, mmap f, mmap private, mmap private: OK, mmap read-only, mmap read-only: OK, mmap read/write, mmap read/write: OK, mmap dirty, mmap dirty: OK, not-mapped unmap, not-mapped unmap: OK, mmap two files, mmap two files: OK, mmap offset, mmap offset: OK, mmap half page, mmap half page: OK, mmaptest: ALL OK, fork test starting, fork test OK, mmaptest: all tests succeeded, and QEMU: Terminated. The prompt at the bottom is 'csc3150@csc3150:~/xv6-labs-2022\$'.

Figure 1: Output of mmaptest Execution

As shown in Figure 1, all memory mapping tests passed successfully, indicating the correct implementation of the `mmap` and `munmap` functionalities. The output confirms that memory regions were correctly mapped, accessed, and unmapped, with appropriate handling of shared and private mappings, as well as proper synchronization between parent and child processes during fork operations.

4 Conclusion

This assignment provided an in-depth exploration of memory management within operating systems, specifically focusing on the implementation of memory mapping in xv6. By developing and integrating the `mmap` and `munmap` system calls, I enhanced xv6's ability to handle file and device mappings, which are fundamental for efficient I/O operations and memory utilization.

Through this project, I gained valuable insights into low-level memory management, process synchronization, and kernel-user space interactions. The successful execution of the test program validated the robustness of the implemented features and reinforced my understanding of operating system principles. Additionally, addressing challenges such as Copy-On-Write and ensuring proper synchronization mechanisms deepened my proficiency in systems programming and kernel development. This experience has significantly strengthened my skills and prepared me for more complex tasks in operating system design and implementation.