

Virtual Memory Initialization and Management in the FreeBSD Operating System

(Inicjalizacja i zarządzanie
pamięcią wirtualną
w systemie operacyjnym FreeBSD)

Jakub Piecuch

Praca licencjacka

Promotor: dr Piotr Witkowski

Uniwersytet Wrocławski
Wydział Matematyki i Informatyki
Instytut Informatyki

3 września 2018

Abstract

Virtual memory is one of the most important abstractions provided to user space programs by modern operating system kernels. While it greatly simplifies application development, implementing it efficiently is a great challenge. The FreeBSD operating system kernel VM (Virtual Memory) subsystem succeeds at this goal. However, in doing so it employs complex data structures and algorithms, which make it difficult to understand for newcomers to the kernel. This thesis provides an overview of the architecture of the FreeBSD VM subsystem, along with fragments of real-world (albeit simplified) source code that implement its most important functions. Subsequently, it goes through the machine-dependent initialization of the VM subsystem on the x86 architecture. It also acts as a guide to reading the kernel's source code.

Pamięć wirtualna jest jedną z najważniejszych abstrakcji udostępnianych programom użytkownika przez jądra nowoczesnych systemów operacyjnych. Znacznie upraszcza ona budowanie aplikacji, jednak efektywna jej implementacja stanowi duże wyzwanie. Przykładem efektywnej implementacji jest podsystem pamięci wirtualnej jądra systemu operacyjnego FreeBSD. Niestety, osiągnięcie dobrej wydajności wymaga zastosowania złożonych struktur danych i algorytmów, co sprawia, że osobom niezaznajomionym z jądrem trudno jest zrozumieć działanie podsystemu pamięci wirtualnej. Poniższa praca stanowi przegląd architektury podsystemu pamięci wirtualnej jądra FreeBSD. Implementacja najważniejszych funkcji przedstawiona jest w postaci uproszczonych fragmentów prawdziwego kodu źródłowego wykorzystywanego w jądrze. Następnie przedstawiony jest proces inicjalizacji części podsystemu pamięci wirtualnej zależnej od architektury na przykładzie architektury x86. Praca ta służy też jako wprowadzenie do czytania kodu źródłowego jądra.

Contents

1	Introduction	7
1.1	What Is Virtual Memory?	7
1.2	The FreeBSD Operating System	9
1.2.1	Browsing the Source Code	9
1.2.2	Bibliographic Notes	10
2	An Overview of the FreeBSD VM Subsystem	11
2.1	vm_space	11
2.2	vm_map	11
2.3	vm_map_entry	13
2.4	vm_object	15
2.5	vm_page	17
2.6	pmap	19
2.6.1	x86 Page Table Layout	19
3	Virtual Memory Management	23
3.1	Address Space Creation – the <code>fork</code> System Call	24
3.2	Executing a New Program – the <code>execve</code> System Call	27
3.3	Handling Page Faults	31
3.4	Adding New Mappings - the <code>mmap</code> System Call	35
4	Virtual Memory Initialization on the x86 Architecture	41
4.1	Kernel Virtual Address Space Layout	41
4.2	Boot-time Virtual Memory Initialization	42

5 Conclusion	51
---------------------	-----------

Bibliography	53
---------------------	-----------

Chapter 1

Introduction

1.1 What Is Virtual Memory?

A system implementing virtual memory provides every process in the system with an address space of its own, its **Virtual Address Space** (VAS). The modifications that a process makes to its own VAS are not visible to other processes, unless they explicitly choose to share fragments of their address spaces. This protects running programs against unwanted interference.

The size of the VAS can (and typically is on 64 bit architectures) much larger than the amount of available **physical memory**, i.e. memory actually installed in the system. For instance, the VAS size on the AMD64 architecture is 2^{48} bytes, that is 256 TiB (tebibytes)¹.

The VAS of any process contains the whole **process image** (also called the **core image**), that is the instructions, data and the run-time stack. A process in execution accesses memory to fetch its instructions, as well as to read and write data. All instructions executed by the process access memory using **Virtual Addresses** (VAs), which are integers in the range from 0 to $N - 1$, where N is the VAS size.

Since the VAS size can be larger than the amount of available physical memory, there cannot be a 1-to-1 correspondence between VAs and **Physical Addresses** (PAs), which refer directly to locations in physical memory. It is necessary to allow either VAs which have no corresponding PAs, or multiple VAs with the same corresponding PA. To accomplish this, a scheme called **address translation** is used.

Modern architectures provide hardware support for address translation in the form of a **Memory Management Unit** (MMU). The MMU translates VAs generated by the program into PAs used to address the physical memory. As the details are architecture-specific, this thesis focuses on the x86 architecture.

¹ 1 TiB = 1024 GiB (gibibytes). A gibibyte is slightly larger than a gigabyte. See https://en.wikipedia.org/wiki/Binary_prefix.

On the x86 architecture, the translation from VAs to PAs can be done using two schemes: **segmentation** and **paging**. Since segmentation is rarely used in today's operating systems, the discussion will be focused on paging. In a paging scheme, the mapping between VAs and PAs is determined by an in-memory data structure called a **Page Table** (PT). The VAS is divided into units called **pages**. Their size is usually 4 KiB. Likewise, the physical memory is divided into page-sized units called **page frames**. Pages are commonly called virtual pages, and page frames are called physical pages. The page table describes the mapping between pages and page frames, instead of between individual virtual addresses and physical addresses. This significantly reduces the size of the page table.

In addition to specifying the mapping between pages and page frames, the page table also determines the protection attributes of each page. For instance, pages containing program code can be marked read-only, e.g. so that a malicious user can't exploit a bug in the program to overwrite the original instructions with malicious ones.

Not every page needs to be mapped to a page frame. Whenever the program references an address which is in a page that is not mapped, a **page fault** occurs. The MMU detects this and the CPU generates an exception. Execution then jumps to a routine provided by the OS, which handles the exception. The same is done when the programs tries to do something to a page that violates the page's protection attributes, e.g. write to a read-only page.

An important observation is that with virtual memory, processes can be partially **resident** in physical memory (i.e. only parts of their image have to be in physical memory). Unused pages can be unmapped, as an exception will generated only when a process references the page. Furthermore, even pages that are used by a process don't have to be resident at all times, as the OS can map them into the VAS in response to a page fault, and then restart the instruction that caused it. This strategy is called **demand paging** and is widely used. It enables the system to use memory resources efficiently by holding only as many mapped pages in physical memory as are necessary for any one process to run smoothly.

In summary, virtual memory has the following important characteristics:

- The address spaces in which processes live are isolated from one another
- A process may run even if the total amount of memory it requires is larger than the total amount of available physical memory
- A process can have its memory pages mapped into its address space on demand, and unmapped when the OS determines the process is unlikely to reference them in the near future

1.2 The FreeBSD Operating System

FreeBSD [2] is an open-source operating system, first released in 1993. It is the most popular OS in the BSD family of operating systems, which includes FreeBSD, NetBSD, OpenBSD and DragonFlyBSD. Unlike Linux, FreeBSD is a whole operating system, providing the kernel, drivers and a suite of utility programs.

It is released under a permissive BSD license, which makes it an attractive choice for commercial applications. FreeBSD has been used as the basis for operating systems such as Apple's MacOS, as well as the OS running on Sony's PlayStation 3 and PlayStation 4 consoles.

1.2.1 Browsing the Source Code

This thesis includes many source code listings from the FreeBSD repository. Because of its scope, the level of implementation detail covered is limited. Should the reader want to explore the source tree themselves (which is strongly encouraged), there are several websites providing the code of the whole FreeBSD kernel with support for identifier search:

(1) <http://fxr.watson.org/>

(2) <http://bxr.su/>

(1) allows the user to select any major version of the source tree, while (2) only presents the latest revision. On the other hand, (2) has a more aesthetically pleasing interface, which includes syntax highlighting. All source code discussion and listings apply specifically to the 12-CURRENT branch of the kernel, which is the latest revision and is still under development as of September 2018, and is available on both websites.

The kernel source code can be found in the `/sys` directory of the tree. Since other parts of the tree are not going to be discussed, all subsequent paths are relative to this directory.

The `kern` directory holds most of the kernel's portable code, like the scheduler, subsystem initialization and shutdown, and the Virtual File System.

The `vm` directory contains the machine-independent part of the virtual memory subsystem. This is where most of the code listings will come from. Machine-dependent parts of the subsystem can be found in the directories named after a particular architecture, e.g. the `i386` directory for the x86 architecture. This is quite confusing, since the kernel source tree contains a directory named `x86` – this directory contains architecture-specific code that is shared between the AMD64 and x86 architectures.

1.2.2 Bibliographic Notes

A more detailed explanation of the concepts introduced in 1.1 can be found in [1, Section 3.3]. For a description of paging on the x86 architecture, see [3, Chapter 4].

Chapter 2

An Overview of the FreeBSD VM Subsystem

The goal of this chapter is to describe the data structures used for managing virtual memory in the FreeBSD Operating System. For each one, a brief high-level introduction is given, followed by a source code listing that shows the C language structure definition. Important member fields are then individually described below the listing.

2.1 vmSPACE

The `vmSPACE` structure is the highest-level structure describing the virtual address space of a process. It contains both the machine-independent (`vm_map`) and machine-dependent (`pmap`) structures used for describing a mapping. The other fields hold various statistics and parameters and are not relevant to the discussion.

```
struct vmSPACE {
    struct vm_map vm_map;           /* Machine-independent address map */
    caddr_t vm_maxsaddr;           /* User VA at max stack growth */
    volatile int vm_refcnt;         /* Number of references */
    struct pmap vm_pmap;            /* Machine-dependent 'physical' map */
};
```

Listing 1: `vm/vm_map.h`: Simplified definition of `struct vmSPACE`

2.2 vm_map

The `vm_map` structure represents the machine-independent part of a virtual address space. It is structured as a tree of `vm_map_entry` structures, each of which describes

a continuous fragment of the address space.

```

struct vm_map {
    struct vm_map_entry header; /* List of entries */
    struct sx lock;             /* Lock for map data */
    struct mtx system_mtx;
    int nentries;               /* Number of entries */
    vm_size_t size;             /* Virtual size */
    u_int timestamp;            /* Version number */
    u_char needs_wakeup;
    u_char system_map;          /* Am I a kernel map? */
    vm_flags_t flags;           /* Flags for this vm_map */
    vm_map_entry_t root;        /* Root of a binary search tree */
    pmap_t pmap;                /* Pointer to physical map */
};

```

Listing 2: vm/vm_map.h: Simplified definition of `struct vm_map`.

```

struct vm_map_entry header;

```

This `vm_map_entry` is used for holding the minimum and maximum virtual address available to the user. It also serves as the header of a linked list of all `vm_map_entry` structures in the `vm_map`, sorted by start address. It is used for quickly retrieving the immediate left and right neighbours of a `vm_map_entry`, as well as iterating over all entries in a `vm_map`.

```

struct sx lock;
struct mtx system_mtx;

```

These are synchronization tools used to manage concurrent access to the map. If the value of `system_map` is `TRUE`, `system_mtx` is used exclusively, otherwise only `lock` is used.

```

u_int timestamp;

```

The value of this field is incremented each time exclusive access is acquired to read or modify the map. This is so that algorithms that require relinquishing and reacquiring the lock after some time can detect if someone has possibly tampered with the data structure in the meantime.

```
u_char needs_wakeup;
```

This is a flag indicating that there is a thread (or threads) waiting for a large enough chunk of free address space to satisfy its allocation request. Whenever space is freed from the map and the flag is set, waiting threads are woken up to retry the allocation.

```
vm_map_entry_t root;
```

This is a pointer to the root of the binary tree used to look up entries in the map. The tree uses the self-balancing splay algorithm by Tarjan and Sleator [8]. The most recently looked up entry is at the root of the tree, which speeds up page fault handling by taking advantage of the spatial locality of page faults (i.e. when a page fault happens, the next one is likely to happen at an address close to the previous one).

2.3 vm_map_entry

The `vm_map_entry` structure describes a continuous, page-aligned segment of an address space. All virtual addresses within this segment have the same protection attributes. The `vm_map_entry` is not responsible for providing the contents of the segment, that's the job of the `vm_object` that is **backing** the entry.

```
struct vm_map_entry {
    struct vm_map_entry *prev; /* Previous entry */
    struct vm_map_entry *next; /* Next entry */
    struct vm_map_entry *left; /* Left child in binary search tree */
    struct vm_map_entry *right; /* Right child in binary search tree */
    vm_offset_t start; /* Start address */
    vm_offset_t end; /* End address */
    vm_offset_t next_read; /* Vaddr of the next sequential read */
    vm_size_t adj_free; /* Amount of adjacent free space */
    vm_size_t max_free; /* Max free space in subtree */
    union vm_map_object object; /* Object I point to */
    vm_offset_t offset; /* Offset into object */
    vm_eflags_t eflags; /* Map entry flags */
    vm_prot_t protection; /* Protection code */
    vm_prot_t max_protection; /* Maximum protection */
    vm_inherit_t inheritance; /* Inheritance */
    int wired_count; /* How many pages in the entry are wired? */
};
```

Listing 3: `vm/vm_map.h`: Definition of `struct vm_map_entry`.

```

struct vm_map_entry prev;
struct vm_map_entry next;

```

These pointers link the entry into the doubly linked list of all entries within a map, sorted by start address.

```
vm_offset_t next_read;
```

This field is used by the kernel to detect when a process accesses pages in the segment sequentially. When the kernel knows the accesses are sequential, it will bring into physical memory pages beyond the faulting one (since it is reasonable to assume that the pattern of accesses will continue) and therefore reduce the number of page faults generated by the process.

```

vm_size_t adj_free;
vm_size_t max_free;

```

`adj_free` is the amount of free space (in bytes) between this entry and the next entry to the right. It is used when searching the `vm_map` for a free segment of a certain size. `max_free` is simply the maximum of the values of `adj_free` of all entries inside the subtree rooted at this entry.

```
union vm_map_object object;
```

This field is either a pointer to a `struct vm_object`, or a pointer to another `struct vm_map` called a submap. Almost every time it is the former, as submaps are only used inside the kernel map to allocate address space in advance for certain data structures.

```
vm_ooffset_t offset;
```

The offset determines which part of the `vm_object` is accessible through the `vm_map_entry`. Let `ent` be a `vm_map_entry` backed by some object. Virtual addresses from `ent.start` to `ent.end - 1` (inclusive) map to offsets `ent.offset` to `ent.offset + (ent.end - ent.start) - 1` within the object.

```
vm_eflags_t eflags;
```

`eflags` contains various flags, some notable of which are:

- `MAP_ENTRY_COW`: indicates that the entry is a **copy-on-write** entry [1, Page 90], which means that initially all pages are marked read-only. As soon as the pro-

cess attempts to write to a page in the entry, the kernel will copy the faulting page to a new page and map the new page into the address space of the process, this time with the write bit set. This avoids unnecessary copying of pages when a process forks or requests a private mapping of a file (i.e. a mapping such that changes to it aren't reflected in the file).

- `MAP_BEHAV_{NORMAL, SEQUENTIAL, RANDOM}`: these flags specify the access pattern that is to be expected from the process. The default is `MAP_BEHAV_NORMAL`, which makes the kernel detect sequential access patterns and act accordingly.

```
vm_prot_t protection;
vm_prot_t max_protection;
```

These two fields specify the current protection attributes of the entry and the most permissive allowable protection attributes, respectively. Both are bitmasks composed of 3 fields: `VM_PROT_READ`, `VM_PROT_WRITE`, and `VM_PROT_EXECUTE`.

```
vm_inheritance_t inherit;
```

The `inherit` field determines what happens to the entry when the process forks. The possible behaviours are:

- `VM_INHERIT_SHARE`: The child gets an entry that shares the underlying object with the parent. Changes made by one process are visible to the other.
- `VM_INHERIT_COPY`: The child gets a copy-on-write entry with contents identical to the parent entry's contents. Changes made by one process are not visible to the other.
- `VM_INHERIT_ZERO`: The child gets an anonymous zero-filled entry with the same size and protection as the parent entry.
- `VM_INHERIT_NONE`: The entry won't appear in the child.

2.4 vm_object

A `vm_object` contains a `vm_map_entry`'s resident pages. Multiple entries belonging to different processes can have the same backing object, allowing for fast interprocess communication through shared memory.

Copy-on-write is implemented in BSD using special **shadow objects** [4, Page 304]. These objects have backing objects themselves, and hence can form chains ending in a non-shadow object. Shadow objects hold pages copied as a result of a copy-on-write fault.

Since `vm_object`s are only containers for resident pages, and not all pages are always resident, there must be some other component responsible for providing the contents of the pages. This is the job of a **pager** [4, Section 6.10] structure contained within the object. It provides the abstraction of a backing store from which pages can be filled with contents, and to which pages can be written back in case of memory shortage.

```
struct vm_object {
    struct rwlock lock;
    /* List entry used to link into the list of all vm_objects */
    TAILQ_ENTRY(vm_object) object_list;
    /* List of objects shadowing this object */
    LIST_HEAD(, vm_object) shadow_head;
    /* List entry used to link into shadow_head of shadowed object */
    LIST_ENTRY(vm_object) shadow_list;
    struct pglist memq;      /* List of resident pages */
    struct vm_radix rtree;   /* Root of resident page radix tree */
    vm_pindex_t size;        /* Object size */
    /* How many vm_map_entries/vm_objects reference this object? */
    int ref_count;
    int shadow_count;        /* Length of linked list at shadow_head */
    vm_memattr_t memattr;    /* Default memory attribute for pages */
    objtype_t type;          /* Type of pager */
    int resident_page_count; /* Number of resident pages */
    struct vm_object *backing_object; /* Object that I'm a shadow of */
    vm_ooffset_t backing_object_offset; /* Offset in backing object */
    /* List of all objects of this pager type */
    TAILQ_ENTRY(vm_object) pager_object_list;
    void *handle;           /* Opaque pointer used by the pager */
    union {
        /* Various pager structures */
    } un_pager;
};
```

Listing 4: `vm/vm_object.h`: Simplified definition of `struct vm_object`.

```
vm_memattr_t memattr;
```

`memattr` specifies the default cache behaviour [3, Section 11.3] of pages belonging to the object. For example, contents of pages with the `VM_MEMATTR_UNCACHEABLE` attribute cannot be cached and all read and write requests have to go directly to the physical memory. This attribute is used primarily for pages representing memory mapped devices, e.g. frame buffers.


```
objtype_t type;
```

There are several types of pagers that can provide the contents of pages to a `vm_object`. The `type` field determines the type of pager used by the object. The most widely used types of pagers are:

- **OBJT_SWAP**: The swap pager is used to provide contents of anonymous memory segments (i.e. segments that are not backed by a file and are initially filled with zeros). When asked to fill a page with contents, it first checks if the page had been swapped out before. If so, the contents are fetched from secondary storage – either a swap file or a dedicated swap partition. If not, the page is simply filled with zeros. In case of memory shortage, the swap pager may be asked to store the contents of a page in backing store.
- **OBJT_DEFAULT**: This is the pager type used in all newly created anonymous mappings (i.e. zero-filled mappings not backed by any file). It fills all new pages with zeros. This pager type is an optimization of the swap pager for the common case where the pages never need to be swapped out to backing store, as memory resources in today's systems are abundant. Not needing to keep track of swap space speeds up the pager's initialization procedures. When there is a need for an object with a default pager to write some of its pages to backing store, its pager is changed to the swap pager.
- **OBJT_VNODE**: The contents of pages backed by a file are supplied and written back by the vnode pager. An object with a vnode pager acts as a general-purpose cache of the backing file's pages, used not only by the virtual memory subsystem (i.e. when a process maps a file), but also when a process reads files via file descriptors.
- **OBJT_DEVICE**: The device pager manages pages belonging to objects which represent memory-mapped physical devices. The pages used by the device pager are different from ordinary pages in that they don't represent a frame of physical memory, or in FreeBSD terminology, they are **fictitious**.

2.5 vm_page

The `vm_page` structure represents a unit of physical address space that can be mapped into a virtual address space. Every `vm_page` has a physical address, although it does not need to be a physical address that can be used to address physical memory. Most pages, though, are pages representing frames of physical memory. Pages with physical addresses beyond the physical memory range are called **fictitious** pages and are primarily used for accessing memory-mapped devices.

```

struct vm_page {
    union {
        TAILQ_ENTRY(vm_page) q; /* Page queue or free list */
        /* Some union fields omitted */
    } plinks;
    TAILQ_ENTRY(vm_page) listq; /* List of pages in the same object */
    vm_object_t object;         /* Which object am I in */
    vm_pindex_t pindex;         /* Offset into object */
    vm_paddr_t phys_addr;       /* Physical address of page */
    struct md_page md;          /* Machine dependent part */
    u_int wire_count;           /* How many maps have this page wired? */
    uint16_t flags;             /* page PG_* flags */
    uint8_t aflags;             /* Flags with atomic access */
    uint8_t oflags;             /* page VPO_* flags */
    u_char act_count;           /* page usage count */
};

```

Listing 5: vm/vm_object.h: Simplified definition of `struct vm_object`.

```

struct md_page md;

```

This field contains per-page information maintained by the `pmap` module, i.e. the module responsible for the machine-dependent part of the mapping. Here is the definition of `struct md_page` for the x86 architecture:

```

struct md_page {
    TAILQ_HEAD(,pv_entry) pv_list;
    int pat_mode;
};

```

Listing 6: i386/include/pmap.h: Definition of `struct md_page`.

The `pv_list` field is a list of structures representing virtual addresses pointing to this page (there can be many of them, since a single `vm_object` can be shared by different `vm_map_entry` structures). The `pat_mode` field specifies caching attributes for the page.

```
uint16_t flags;
uint8_t aflags;
uint8_t oflags;
```

The flags for an object are split into several fields due to different characteristics (e.g. locks needed to read/modify the flag, whether operations need to be atomic). Flags contained in these fields include:

- **PGA_REFERENCED**: An access has recently been made to the page.
- **PG_FICTITIOUS**: The page doesn't represent a real physical page frame.

2.6 pmap

The **pmap** structure encapsulates the architecture-dependent aspects of the mapping. Most architectures provide some kind of hardware support for paging, for example the MMU on the x86 architecture. The MMU contains a hardware **page table walker**, which traverses the in-memory page table each time a memory access is made and the needed mapping isn't cached in the TLB. The TLB (Translation Look-aside Buffer) is simply a cache that contains mappings of recently accessed pages). The hardware page table walker expects the page table and page table entries to have a specific layout. This layout is exactly what is implemented in the x86 version of the **pmap** structure.

2.6.1 x86 Page Table Layout

For a detailed description of the layout, see [3, Section 4.2].

Logically, the page table is an array of 4-byte **Page Table Entries** (PTEs). The virtual address used to access memory is split into a **Virtual Page Number** (VPN) and an offset into the page. On x86, the VPN consists of the 20 most significant bits of the address, while the 12 least significant bits are the offset. The VPN is used to find the appropriate PTE.

Each entry contains a **valid bit**, which tells whether or not the corresponding page is mapped to a frame of physical memory. The bit is set if and only if the page is mapped. Every entry also contains the physical address of the frame to which the page is mapped, provided the valid bit is set. The entries also describe the page's protection attributes (whether write access is permitted), cacheability attributes (whether its contents can be cached at all, as well as whether the write-through or write-back policy should be applied), and minimum privilege level required to access the page (user or supervisor).

To help the OS determine which resident pages are actively used and which are not, the PTEs contain a **referenced bit** and a **dirty bit** which indicate that a page has been read from or written to, respectively, are set by the hardware and can be read and cleared by the OS e.g. during a page fault.

The last field contained in PTEs is the **global bit**. When it is set, the entry containing it won't get flushed from the TLB when the current page table is changed. Since the kernel used to (see section 4.1) be mapped into the same region of address space for every process, this bit could be used to avoid unnecessary page faults after performing a context switch by setting it in all entries mapping the kernel.

The size of the address space is 2^{32} bytes and each page is $2^{12} = 4096$ bytes large. Therefore, to describe the entire address space of a single process as a flat array of PTEs, we would need $4 * \frac{2^{32}}{2^{12}} = 2^{22}$ bytes, or 4 MiB. With hundreds of processes running concurrently on a single system not being uncommon, the amount of memory needed to store the page tables of resident processes would quickly become a bottleneck preventing the system from increasing its degree of multiprogramming (i.e. how many processes are resident in memory at the same time).

The solution to this problem is to give the page table a hierarchical structure. The virtual address, instead of being split into 2 parts, is split into 3 parts. The first and second parts (bits 22-31 and 12-21 respectively, where 0 is the least significant bit and 31 is the most significant bit) are used to find the PTE, and bits 0-11 are the offset into the page.

Bits 22-31 are used as an index into the **page directory**, an array of 4-byte **Page Directory Entries** (PDEs). Each PDE contains a pointer to a single page-sized chunk of the page table. However, if that chunk contains no valid entries, no physical memory is allocated for it and the PDE that points to it has its valid bit cleared. This allows for significant memory savings, since usually most of a process's virtual address space is unmapped. Once the physical address of the page table chunk is extracted from the PDE, bits 12-21 are used to index the chunk and retrieve the corresponding PTE.

```

/* NPGPTD - size of page directory in pages (1 in this case) */
struct pmap {
    pd_entry_t          *pm_pdir;          /* VA of page directory */
    TAILQ_HEAD(,pv_chunk) pm_pvchunk; /* List of chunks storing pv_entries */
    /* Set of CPUs using this page table right now */
    cpuset_t            pm_active;
    struct pmap_statistics pm_stats; /* pmap statistics */
    LIST_ENTRY(pmap)     pm_list;         /* List of all pmaps */
    struct vm_radix       pm_root;        /* Tree of spare page table pages */
    /* Table of pointers to page directory pages */
    vm_page_t            pm_ptdpg[NPGPTD];
};

```

Listing 7: i386/include/pmap.h: Simplified definition of `struct pmap`.

Chapter 3

Virtual Memory Management

In this chapter we take a look at the inner workings of the most important functions carried out by the operating system when it comes to managing a process's virtual address space. The order in which the functions will be discussed roughly corresponds to the order in which these functions are invoked during the lifetime of a process.

The first function we will discuss is the creation of a new address space. A new virtual address space is created whenever a new process enters the system. All processes in FreeBSD are created using a mechanism called **forking** [1, Section 1.6.1]. When a process forks, an almost identical copy of it is created by the kernel. The copy is called the **child**, and the forking process is called the **parent**. The layout of the child's virtual address space is identical to the one of its parent.

Usually, new processes are created in order to run some executable program stored in a file. To do that, a process invokes the **execve** [9] system call. The operating system then completely replaces the contents of the process's address space with the process image contained in the file.

Once the program begins to execute, it accesses instructions and data using virtual addresses. Some of these accesses will generate **page faults**, and the kernel has to handle them appropriately, fetching data from secondary storage or terminating the process if it tried to access an unmapped region.

A process can also modify the layout of its own address space. The **mmap** [10] system call allows a process to create a new segment of virtual memory that it can access, or overlay an existing one. The initial contents of these segments can be provided by a file from the file system or simply filled with zeros.

Each section of this chapter describes the steps taken by the kernel when carrying out each of these functions.

3.1 Address Space Creation – the fork System Call

The `fork` [11] system call creates an almost identical clone of the calling process, which includes an address space cloned from the address space of the calling process.

The `vm_space_fork` is used by the code implementing the system call to create a copy of an address space. The function takes a pointer to the `vm_space` structure of the calling process as an argument and returns a pointer to the copy which it has created.

The code listing below is a simplified version of the actual implementation of the `vm_space_fork` function. Fragments of code which deal with edge cases, synchronization and accounting have been removed for clarity.

```

/*
 * Create a vm_space for the child process given the vm_space of the
 * parent process. Entries inside parent's the vm_map are passed
 * down to the child according to the VM_INHERIT_* flags.
 */
struct vm_space *
vm_space_fork(struct vm_space *vm1)
{
    struct vm_space *vm2;
    vm_map_t new_map, old_map;
    vm_map_entry_t new_entry, old_entry;
    vm_object_t object;

    old_map = &vm1->vm_map;
    /*
     * Allocate a new vm_space structure with the same address range
     * as the parent's vm_space
     */
    vm2 = vm_space_alloc(vm_map_min(old_map), vm_map_max(old_map), NULL);
    /* Copy immutable fields of vm1 to vm2. */
    vm2->vm_taddr = vm1->vm_taddr;
    vm2->vm_daddr = vm1->vm_daddr;
    vm2->vm_maxsaddr = vm1->vm_maxsaddr;
    new_map = &vm2->vm_map;
    old_entry = old_map->header.next;

    /* Iterate over all entries in parent's vm_map */
    while (old_entry != &old_map->header) {

        /*
         * The inheritance field determines whether and how
         * the entry is inherited by the child vm_map.
         */
        switch (old_entry->inheritance) {

```



```

case VM_INHERIT_NONE:
    /* Don't clone the entry to new_map */
    break;

case VM_INHERIT_SHARE:
    /*
     * Clone the entry, but share the underlying vm_object between
     * both entries.
     */

    /*
     * Anonymous vm_objects are created lazily as an optimization,
     * so the object may not have been created yet.
     * Create it if that's the case.
     */
    object = old_entry->object.vm_object;
    if (object == NULL) {
        object = vm_object_allocate(OBJT_DEFAULT,
            atop(old_entry->end - old_entry->start));
        old_entry->object.vm_object = object;
        old_entry->offset = 0;
    }

    if (old_entry->eflags & MAP_ENTRY_NEEDS_COPY) {
        /*
         * old_entry is a copy-on-write entry and its shadow
         * object hasn't been created yet (which is possible
         * because shadow objects are created only when necessary).
         * We need to create the shadow object here.
         * If we didn't, the entries in the child and parent
         * would get separate shadow objects if they attempted
         * to write to the region of memory represented by this
         * entry, which is not what we want.
         * We want the changes made by the parent to be visible
         * to the child and vice versa.
         */
        vm_object_shadow(&old_entry->object.vm_object,
            &old_entry->offset,
            old_entry->end - old_entry->start);
        old_entry->eflags &= ~MAP_ENTRY_NEEDS_COPY;

        /*
         * Add a reference to the newly created shadow object,
         * bringing its reference count to 2 (parent + child).
         */
        vm_object_reference(old_entry->object.vm_object);

        object = old_entry->object.vm_object;
    }

```

```

/*
 * Clone the entry, referencing the shared object.
 */
new_entry = vm_map_entry_create(new_map);
*new_entry = *old_entry;

/*
 * Insert the entry into the new map -- we know we're
 * inserting at the end of the new map.
 */
vm_map_entry_link(new_map, new_map->header.prev, new_entry);

/*
 * Do some bookkeeping associated with the sizes of
 * the text, data and stack segments.
 */
vm_space_map_entry_forked(vm1, vm2, new_entry);

/*
 * Update the physical map
 */
pmap_copy(new_map->pmap, old_map->pmap,
          new_entry->start,
          (old_entry->end - old_entry->start),
          old_entry->start);
break;

case VM_INHERIT_COPY:
/*
 * Clone the entry and link into the map.
 */
new_entry = vm_map_entry_create(new_map);
*new_entry = *old_entry;

new_entry->object.vm_object = NULL;
vm_map_entry_link(new_map, new_map->header.prev,
                  new_entry);
vm_space_map_entry_forked(vm1, vm2, new_entry);

/*
 * Make old_entry and new_entry point to the same vm_object.
 * Both entries are marked copy-on-write so that the initial
 * contents of the memory segment are identical to the
 * parent and the child processes, but changes made by one
 * are not visible to the other.
 */
vm_map_copy_entry(old_map, new_map, old_entry, new_entry);
break;

```

```

case VM_INHERIT_ZERO:
    /*
     * Create a new anonymous mapping entry modelled from
     * the old one.
     */
    new_entry = vm_map_entry_create(new_map);
    memset(new_entry, 0, sizeof(*new_entry));

    new_entry->start = old_entry->start;
    new_entry->end = old_entry->end;
    new_entry->eflags = old_entry->eflags &
        ~(MAP_ENTRY_USER_WIRED | MAP_ENTRY_IN_TRANSITION |
          MAP_ENTRY_VN_WRITECNT);
    new_entry->protection = old_entry->protection;
    new_entry->max_protection = old_entry->max_protection;
    new_entry->inheritance = VM_INHERIT_ZERO;

    /* Add new entry to the map. */
    vm_map_entry_link(new_map, new_map->header.prev,
        new_entry);
    /* Bookkeeping */
    vm_space_map_entry_forked(vm1, vm2, new_entry);

    break;
}
old_entry = old_entry->next;
}

return (vm2);
}

```

Listing 8: `vm/vm_map.c`: Simplified implementation of the `vm_space_fork` function.

3.2 Executing a New Program – the `execve` System Call

A process may request the operating system to replace the program that is currently executing in the process with another program stored as a file in the file system. To do this, the process invokes the `execve` system call.

Among other things, the `execve` system call is responsible for wiping the address space of the calling process and populating it with the image contained in the file supplied by the caller.

The bulk of `execve` logic is contained in the `do_execve` function, defined in `kern/kern_exec.c`. However, most of the address space management logic is imple-

mented elsewhere. Hence, we will not show the code for `do_execve` here.

To understand the code responsible for populating an empty address space with a program contained in a file, it is helpful to understand the concept of **image activators** [5, Section 3.5].

The kernel supports multiple different file formats that contain an executable program. For instance, the file pointed to by the path passed to `execve` can be a plain-text script, with the first line containing the path to the interpreter program preceded by the characters `#!`. For example, if a script begins with the line `#!/bin/bash`, this is information to the kernel that the actual program to be executed is the interpreter, and the path to the script should be passed as an argument to the interpreter. Alternatively, an executable file can be an object file in the ELF (Executable and Linkable File) format. In that case, the kernel populates the address space according to the layout specified in the file.

Each executable file format has its own activator. At a certain point in the `do_execve` function, the kernel invokes every activator, passing the file to execute and other parameters to it. Each activator first checks the initial bytes of the file to see if they match a special constant, called the magic number, which is assigned to a specific file format. If the activator detects a match for its file format, it proceeds to load the program from the file. Otherwise, it returns an error code and the kernel tries the next activator.

In this thesis we will take a closer look at the activator for the ELF64 file format. The initial layout of the address space is described using **sections** and **segments**.

A section contains data that serves the same or similar purpose and has the same protection attributes. For instance, the `.data` section contains statically allocated, initialized, read-write data for use by the program, and the `.text` section contains read-only program code.

A segment consists of a number of sections and has a specified virtual address at which it is to be loaded. The sections comprising the segment are loaded into memory starting at the segment's initial address and continuing, one section after another.

The image activator for the ELF64 file format is a function called `exec_elf64_imgact`. We are only interested in the part which scans the table of program headers (which describe individual segments) and loads appropriate segments into memory, which is shown below.

```

/*
 * hdr - ELF header, contains general information about the file
 * phdr - array of program headers, which describe segments
 */
for (i = 0; i < hdr->e_phnum; i++) { /* For each segment... */
    switch (phdr[i].p_type) {
    case PT_LOAD: /* Loadable segment */
        if (phdr[i].p_memsz == 0)
            /* Nothing to load */
            break;

        /* Read the segment's protection attributes */
        prot = elf64_trans_prot(phdr[i].p_flags);
        /* Load the segment into the address space */
        elf64_load_section(imgp, phdr[i].p_offset,
            (caddr_t)(uintptr_t)phdr[i].p_vaddr,
            phdr[i].p_memsz, phdr[i].p_filesz, prot);

        break;
    /* Other cases omitted */
    default:
        break;
    }
}

```

Listing 9: kern/imgact_elf.c: Simplified fragment of the `exec_elf64_imgact` function mapping program segments into the address space of the calling process.

The confusingly named `elf64_load_section` function loads individual segments into the virtual address space. Most of its code handles the case when the size of a segment in memory is larger than the size of the segment in the executable file and the end of the mapping backed by the file is not on a page boundary. The mapping of the segment is done inside the `elf64_map_insert` function.

The `elf64_map_insert` function is responsible for mapping a range of offsets inside a `vm_object` to a range of virtual addresses in a `vm_map`. It also contains code to handle edge cases, but the essence has been distilled in the code listing below.

```

/*
 * Map a segment from an ELF executable into an address space.
 * map - vm_map into which the segment should be mapped
 * object - vm_object for the executable
 * offset - offset of the start of the segment within the executable
 * start/end - start/end virtual address of the segment
 * prot - protection attributes with which the segment should be mapped
 * cow - copy-on-write attributes
 */
static int
elf64_map_insert(vm_map_t map, vm_object_t object, vm_ooffset_t offset,
                vm_offset_t start, vm_offset_t end, vm_prot_t prot, int cow)
{
    /*
     * Handle edge cases where the start/end virtual addresses
     * are not page-aligned.
     */
    if (start != trunc_page(start)) {
        /* ... */
    }
    if (end != round_page(end)) {
        /* ... */
    }
    if ((offset & PAGE_MASK) != 0) {
        /*
         * Another edge case: (offset % PAGE_SIZE) != (start % PAGE_SIZE)
         * This means that we have to copy the data into a fresh set of
         * pages of anonymous memory instead of mapping the executable
         * vm_object.
         */
        /* ... */
    } else {
        /*
         * This is the most common code path: no edge cases,
         * simply map a range of pages in the vm_object to a range
         * of virtual pages in the address space
         */
        vm_object_reference(object);
        vm_map_fixed(map, object, offset, start, end - start,
                    prot, VM_PROT_ALL, cow | MAP_CHECK_EXCL);
    }
    return (KERN_SUCCESS);
}

```

Listing 10: kern/imgact_elf.c: Simplified fragment of the `elf64_map_insert` function mapping a program segment into the address space of a process.

3.3 Handling Page Faults

Once a process begins execution, it is almost inevitable that at some point it will attempt to access a page of virtual memory that isn't currently resident. When that happens, control is passed to the kernel, which attempts to bring the page into physical memory and resume the process which caused the page fault.

The main function responsible for handling page faults is `vm_fault_hold`. Its arguments include the `vm_map` of the faulting process, the virtual address that the process attempted to access, and what type of access was attempted (read or write).

It begins by looking up the `vm_map_entry` corresponding to the address. It then follows the chain of shadow objects starting at the entry's backing object, looking for an object that contains the needed page. If an object with the page is found and the fault is a write fault (the process attempted to write to the page), the page is copied to the first object in the chain. If it is just a read fault, the page can be simply mapped into the address space, without the need to copy or move it. The following code listing shows the algorithm in more detail.

```
/*
 * Handle a page fault by finding the needed page and mapping it
 * into the faulting process's address space.
 * map - vm_map of the faulting process
 * vaddr - virtual address that the process tried to access and
 * triggered the fault
 * fault_type - what kind of access was attempted (read/write), and
 * also what the cause for calling the function is (it's not only used
 * for handling page faults)
 */
int
vm_fault_hold(vm_map_t map, vm_offset_t vaddr, vm_prot_t fault_type)
{
    /*
     * The faultstate structure records the parts of the
     * current state of the handler which are most commonly
     * passed to helper functions. Its fields are:
     * map - vm_map of the faulting process
     * first_object, first_index, first_m and
     * object, index, m -
     *   object, offset into the object, and page corresponding
     *   to faulting address within the first and current
     *   vm_object in the chain
     * entry - vm_map_entry containing the faulting page
     */
    struct faultstate fs;
    vm_object_t next_object;
    vm_prot_t prot;
```

```
boolean_t wired;
int behind, ahead;
```

RetryFault:

```
/*
 * Find the backing store object and offset into it to begin the
 * search.
 */
fs.map = map;
vm_map_lookup(&fs.map, vaddr, fault_type |
    VM_PROT_FAULT_LOOKUP, &fs.entry, &fs.first_object,
    &fs.first_pindex, &prot, &wired);

fs.first_m = NULL;

/*
 * Set the initial (object, offset) pair that will be used
 * to query objects for the page we're looking for.
 */
fs.object = fs.first_object;
fs.pindex = fs.first_pindex;
while (TRUE) {

    /* See if page is resident. */
    fs.m = vm_page_lookup(fs.object, fs.pindex);
    if (fs.m != NULL) {
        /* Block and retry if the page is busy. */
        if (vm_page_bused(fs.m)) {
            if (fs.m == vm_page_lookup(fs.object, fs.pindex)) {
                vm_page_sleep_if_busy(fs.m);
            }
            goto RetryFault;
        }

        /* Mark page busy for other processes. */
        vm_page_xbused(fs.m);
        /* We found the page to be copied to fs.first_m */
        break; /* Break to PAGE HAS BEEN FOUND. */
    }

    /*
     * Page is not resident. If the pager might contain the page
     * or this is the beginning of the search, allocate a new
     * page. (Default objects are zero-fill, so there is no real
     * pager for them.)
     */
    if (fs.object->type != OBJT_DEFAULT ||
        fs.object == fs.first_object) {
```



```

    /* Allocate a new page inside the object. */
    fs.m = vm_page_alloc(fs.object, fs.pindex);
    if (fs.m == NULL) {
        /* No pages available. Block and retry. */
        vm_waitpfault();
        goto RetryFault;
    }
}

/*
 * Call the pager to retrieve the page if there is a chance
 * that the pager has it, and potentially retrieve additional
 * pages at the same time.
 */
if (fs.object->type != OBJT_DEFAULT) {

    /*
     * (omitted) Calculate how many surrounding pages we should
     * ask the pager for, determining the values of ahead and
     * behind.
     */

    /* Get the pages */
    rv = vm_pager_get_pages(fs.object, &fs.m, 1,
        &behind, &ahead);
    if (rv == VM_PAGER_OK) {
        /* Now we have the page to move to first_m */
        break; /* Break to PAGE HAS BEEN FOUND. */
    }

    /*
     * If an I/O error occurred or the requested page was
     * outside the range of the pager, clean up and return
     * an error.
     */
    if (rv == VM_PAGER_ERROR || rv == VM_PAGER_BAD) {
        vm_page_free(fs.m);
        fs.m = NULL;
        unlock_and_deallocate(&fs);
        return (rv == VM_PAGER_ERROR ? KERN_FAILURE :
            KERN_PROTECTION_FAILURE);
    }
}

```

```

/*
 * The requested page does not exist at this object/
 * offset. Remove the invalid page from the object,
 * waking up anyone waiting for it, and continue on to
 * the next object. However, if this is the top-level
 * object, we must leave the busy page in place to
 * prevent another process from rushing past us, and
 * inserting the page in that object at the same time
 * that we are.
 */
if (fs.object != fs.first_object) {
    vm_page_free(fs.m);
    fs.m = NULL;
}
}

/*
 * We get here if the object has a default pager
 * or the pager doesn't have the page.
 */
if (fs.object == fs.first_object)
    fs.first_m = fs.m;

/* Move on to the next object. */
next_object = fs.object->backing_object;
if (next_object == NULL) {
    /*
     * If there's no object left, that means we accessed
     * a previously untouched page inside an anonymous mapping,
     * so we fill the page in the top object with zeros.
     */
    if (fs.object != fs.first_object) {
        fs.object = fs.first_object;
        fs.pindex = fs.first_pindex;
        fs.m = fs.first_m;
    }
    fs.first_m = NULL;

    /* Zero the page and mark it valid. */
    pmap_zero_page(fs.m);
    fs.m->valid = VM_PAGE_BITS_ALL;

    break; /* Break to PAGE HAS BEEN FOUND. */
}

/* Adjust object offset. */
fs.pindex += OFF_TO_IDX(fs.object->backing_object_offset);
fs.object = next_object;
}

```

```

    /* PAGE HAS BEEN FOUND. */

    /*
     * If the page is being written, but isn't already owned by the
     * top-level object, we have to copy it into a new page owned by the
     * top-level object.
    */
    if (fs.object != fs.first_object) {
        /*
         * We only really need to copy if we want to write it.
        */
        if ((fault_type & (VM_PROT_COPY | VM_PROT_WRITE)) != 0) {
            pmap_copy_page(fs.m, fs.first_m);
            fs.first_m->valid = VM_PAGE_BITS_ALL;
            /* We no longer need the old page or object. */
            release_page(&fs);
            /* Only use the new page below... */
            fs.object = fs.first_object;
            fs.pindex = fs.first_pindex;
            fs.m = fs.first_m;
        } else {
            prot &= ~VM_PROT_WRITE;
        }
    }

    /* Put this page into the physical map. */
    pmap_enter(fs.map->pmap, vaddr, fs.m, prot,
        fault_type | (wired ? PMAP_ENTER_WIRED : 0), 0);

    return (KERN_SUCCESS);
}

```

Listing 11: `vm/vm_fault.c`: Simplified implementation of the `vm_fault_hold` function.

3.4 Adding New Mappings - the `mmap` System Call

A process has the ability to introduce new mapped regions into its virtual address space using the `mmap` system call. It can create zero-filled anonymous mappings, as well as mappings whose contents are provided by a file. The `kern_mmap` function implements the whole system call.

```

/*
 * Add a new mapping to a thread's address space.
 * td - thread that called mmap
 * addr - virtual address where the mapping should begin
 * size - size of the mapping
 * prot - desired protection attributes of the mapping
 * flags - flags describing the kind of mapping
 * fd - file descriptor pointing to the file that will provide
 * the contents of the mapping
 * pos - position within the file where the contents of the
 * mapping begin
 */
int
kern_mmap(struct thread *td, uintptr_t addr, size_t size, int prot, int flags,
          int fd, off_t pos)
{
    struct vm_space *vms;
    struct file *fp;
    vm_offset_t vaddr;
    vm_size_t pageoff;
    vm_prot_t cap_maxprot;
    int align, error;
    cap_rights_t rights;

    vms = td->td_proc->p_vm_space;
    fp = NULL;

    if (flags & MAP_STACK) {
        /*
         * Stacks are not backed by any file, so we map them
         * anonymously.
         */
        flags |= MAP_ANON;
        pos = 0;
    }

    /*
     * Align the file position to a page boundary,
     * and save its page offset component.
     */
    pageoff = (pos & PAGE_MASK);
    pos -= pageoff;

    /* Adjust size for rounding (on both ends). */
    size += pageoff;          /* low end... */
    size = (vm_size_t) round_page(size); /* hi end */

```

```

/* Check for illegal addresses. */
if (flags & MAP_FIXED) {
    /*
     * MAP_FIXED means the mapping must begin at the
     * address supplied by the user.
     */

    /*
     * The specified address must have the same remainder
     * as the file offset taken modulo PAGE_SIZE, so it
     * should be aligned after adjustment by pageoff.
     */
    addr -= pageoff;
    if (addr & PAGE_MASK)
        return (EINVAL);

    /* Address range must be all in user VM space. */
    if (addr < vm_map_min(&vms->vm_map) ||
        addr + size > vm_map_max(&vms->vm_map))
        return (EINVAL);
    if (addr + size < addr)
        return (EINVAL);
} else {
    /*
     * For non-fixed mappings where no hint is provided or
     * the hint would fall in the potential heap space,
     * place it after the end of the largest possible heap.
     */
    if (addr == 0 ||
        (addr >= round_page((vm_offset_t)vms->vm_taddr) &&
         addr < round_page((vm_offset_t)vms->vm_daddr +
             lim_max(td, RLIMIT_DATA))))
        addr = round_page((vm_offset_t)vms->vm_daddr +
            lim_max(td, RLIMIT_DATA));
}

/*
 * Guard mappings are used only to reserve address space,
 * hence their protection is VM_PROT_NONE.
 */
if ((flags & MAP_GUARD) != 0) {
    error = vm_mmap_object(&vms->vm_map, &addr, size, VM_PROT_NONE,
        VM_PROT_NONE, flags, NULL, pos, FALSE, td);
} else if ((flags & MAP_ANON) != 0) {
    /* Mapping blank space is trivial. */
    error = vm_mmap_object(&vms->vm_map, &addr, size, prot,
        VM_PROT_ALL, flags, NULL, pos, FALSE, td);
} else {

```

```

    /*
     * Mapping file, get fp for validation and don't let the
     * descriptor disappear on us if we block. Check capability
     * rights, but also return the maximum rights to be combined
     * with maxprot later.
     */
    error = fget_mmap(td, fd, &rights, &cap_maxprot, &fp);
    if (error != 0)
        goto done;

    /* Actually map the file. */
    error = fo_mmap(fp, &vms->vm_map, &addr, size, prot,
        cap_maxprot, flags, pos, td);
}

done:
    if (fp)
        fdrop(fp, td);

    return (error);
}

```

Listing 12: vm/vm_mmap.c: Simplified implementation of the kern_mmap function.

The function that actually inserts the new mapping into the calling process's address space is the `vm_mmap_object` function. It does some basic sanity checks, after which it calls `vm_map_find_min` or `vm_map_fixed`. Both are functions that result in a new `vm_map_entry` being added to the map.

```

/*
 * Map a fragment of a vm_object into the given vm_map.
 * map - vm_map into which the object should be mapped
 * addr - pointer to the virtual address at which the
 * mapping should start
 * prot - protection attributes of the mapping
 * maxprot - maximum allowable protection attributes of the mapping
 * object - vm_object whose pages should be mapped
 * foff - offset into the object where the contents of the
 * mapping start
 * td - thread that called mmap
 */
int
vm_mmap_object(vm_map_t map, vm_offset_t *addr, vm_size_t size, vm_prot_t prot,
    vm_prot_t maxprot, int flags, vm_object_t object, vm_ooffset_t foff,
    struct thread *td)
{
    boolean_t fitit; /* FALSE if *addr is just a hint */
    vm_offset_t max_addr;

```

```

int docow, findspace, rv;

/*
 * MAP_FIXED means the mapping needs to begin exactly at
 * *addr and *addr needs to be page-aligned.
 */
if ((flags & MAP_FIXED) == 0) {
    fitit = TRUE;
    *addr = round_page(*addr);
} else {
    if (*addr != trunc_page(*addr))
        return (EINVAL);
    fitit = FALSE;
}

/*
 * (omitted) Perform various flag conversions, determining
 * the values of the docow and flags variables.
 */

if (fitit) {
    /*
     * Here *addr is only a hint, so the real start address
     * of the allocated entry may be different.
     */

    /*
     * (omitted) Convert alignment flags, determining the
     * value of findspace.
     */

    max_addr = 0; /* Means there's no maximum. */

    /* Find free space in the map and insert a new entry. */
    rv = vm_map_find_min(map, object, foff, addr, size,
        round_page((vm_offset_t)td->td_proc->p_vmspace->
            vm_daddr + lim_max(td, RLIMIT_DATA)), max_addr,
        findspace, prot, maxprot, docow);
} else {
    /* Try to insert an entry at exactly *addr. */
    rv = vm_map_fixed(map, object, foff, *addr, size,
        prot, maxprot, docow);
}

```

```
if (rv == KERN_SUCCESS) {  
    /*  
     * If the process has requested that all future mappings  
     * be wired, then heed this.  
     */  
    if (map->flags & MAP_WIREFUTURE) {  
        vm_map_wire(map, *addr, *addr + size,  
                     VM_MAP_WIRE_USER | ((flags & MAP_STACK) ?  
                     VM_MAP_WIRE_HOLESOK : VM_MAP_WIRE_NOHOLES));  
    }  
}  
return (vm_mmap_to_errno(rv));  
}
```

Listing 13: vm/vm_mmap.c: Simplified implementation of the vm_mmap_object function.

Chapter 4

Virtual Memory Initialization on the x86 Architecture

Every architecture has a lot of aspects whose specifics differ between different architectures, and sometimes even between revisions of the same architecture. One example is boot-time initialization of the virtual memory subsystem, which includes setting up the initial kernel page table (if any) and reserving memory used by low-level components. Of course, we would like the operating system to have as much portable code as possible, therefore it is crucial to properly identify machine-dependent aspects of system operation and abstract them away behind a common interface.

When it comes to virtual memory management, the FreeBSD kernel maximizes portability by hiding architecture-specific implementation details inside the `pmmap` [12] module, which exposes an abstract interface for low-level VM initialization and page table management. This chapter focuses on the boot-time machine-dependent initialization process as implemented in the 32-bit x86 version of the kernel.

4.1 Kernel Virtual Address Space Layout

Until recently, the kernel and user address spaces used to occupy the same virtual address space. That is, the address space of every user process had the kernel mapped into the upper 1GiB, so that switching between executing user and kernel code was less costly in terms of performance.

However, recently discovered security vulnerabilities affecting (among others) x86 processors, dubbed “Meltdown” [6] and “Spectre” [7], allow an unprivileged process to read the contents of its entire address space, including possibly sensitive information stored in the kernel portion of the address space. To mitigate these vulnerabilities, the kernel no longer shares its address space with user processes.

Any user process has at its disposal almost the entire 4GiB of virtual address space, and so does the kernel. This comes at a cost of a greater performance penalty when switching between executing user and kernel code. Previously, it wasn't necessary to flush the TLB (i.e. invalidate all entries) when jumping from user to kernel code and vice versa, but now, since it involves changing the current page table, it has become necessary to do so, which results in a higher rate of page faults.

There remains a little bit of address space inside a user process that is still inaccessible to it. That part of the address space is called the **trampoline**, since usually executions jumps into it only to quickly jump out of it. For instance, code that handles exceptions and interrupts and switches the page table to the kernel page table can be found in the trampoline area. After the kernel page table is installed, execution jumps to lower addresses, where the remainder of kernel code resides.

The exact layout of the kernel virtual address space is as follows:

- 0x00000000 - 0x003fffff (0-4 MiB): Transient identity mapping of low memory, normally disabled to catch NULL dereference attempts.
- 0x00400000 - 0x007fffff (4-8 MiB): Mapping of the same low memory, but this mapping isn't transient
- 0x00800000 - 0xffbfffff: General purpose area containing kernel text and data. Managed by the machine-independent part of the VM subsystem.
- 0xffc00000 - 0xffdfffff: Recursive¹ mapping of the kernel page table.
- 0xffe00000 - 0xffffffff: Trampoline area shared between the kernel and user address spaces.

4.2 Boot-time Virtual Memory Initialization

When the system starts up, the MMU is disabled, so the CPU addresses physical memory directly. The first code that is executed by the CPU is the firmware found in the motherboard's ROM (Read-Only Memory). That firmware (also called the BIOS - Basic Input/Output System) [4, Section 15.1] initializes the hardware at the lowest level of abstraction and provides the rest of the software with a more abstract and uniform interface to the hardware.

The BIOS then finds a device to boot from and loads into memory a program called the **bootloader** [4, Section 15.2]. The bootloader's job is to find the file containing the kernel image, load it into memory, and transfer control to it.

Once the kernel code starts running, the CPU is still addressing physical memory directly. The kernel's entry point on the x86 architecture can be found in the

¹ https://wiki.osdev.org/Page_Tables#Recursive_mapping (accessed 03/09/2018)

file `i386/i386/locore.s`. The code in the file identifies the CPU and converts boot information passed to the kernel into a common format. Other than that, it calls C procedures that initialize the system.

The first stage of VM subsystem initialization sets up a minimal kernel page table that maps the range of virtual addresses containing the kernel image to identical physical addresses, and then enables the MMU. All this is done in the `pmap_cold` procedure presented below.

```

/*
 * Called from locore.s before paging is enabled. Sets up the first
 * kernel page table.
 */
void
pmap_cold(void)
{
    u_long a;
    u_int cr3;

    /* physfree - physical address of first free physical page */
    physfree = (u_long)&_end;
    /*
     * NBPDR - size in bytes of the range mapped by a single
     * PDE (Page Directory Entry). In our case it's 4MiB.
     */
    physfree = roundup2(physfree, NBPDR);
    KERNend = physfree;

    /*
     * Allocate kernel page table pages.
     * The allocpages function simply returns the value of
     * physfree and increments it by the requested number of pages.
     * NKPT - initial number of kernel page table pages. Must be
     * sufficient to map the array of vm_page structures representing
     * physical memory.
     */
    KPTphys = allocpages(NKPT, &physfree);
    /*
     * KPTmap points to a linear mapping of kernel page table pages.
     * Since NKPT is not enough pages to contain the entire page table,
     * later we'll allocate a larger number of pages of virtual address space
     * (not physical pages) to be able to provide a view of the entire
     * page table.
     */
    KPTmap = (pt_entry_t *)KPTphys;

```

```

/*
 * Allocate page table directory.
 * NPGPTD - size of the Page Table Directory (PTD) in pages.
 */
IdlePTD = (pd_entry_t *)allocpages(NPGPTD, &physfree);

/*
 * Install page tables into PTD. Page table page 1 is wasted.
 * The ptoa function converts a page number to an address simply
 * by multiplying the page number by 4096 (the page size).
 */
for (a = 0; a < NKPT; a++)
    IdlePTD[a] = (KPTphys + ptoa(a)) | PG_V | PG_RW | PG_A | PG_M;

/*
 * Install recursive mapping for kernel page tables into
 * itself.
 * PTDPTDI - PTD index of the recursive mapping
 * (i.e. the page directory entry at index PTDPTDI points to
 * the page directory itself).
 */
for (a = 0; a < NPGPTD; a++)
    IdlePTD[PTDPTDI + a] = ((u_int)IdlePTD + ptoa(a)) | PG_V |
        PG_RW;

/*
 * Initialize page table pages mapping physical address zero
 * through the (physical) end of the kernel. Many of these
 * pages must be reserved, and we reserve them all and map
 * them linearly for convenience.
 * This and all other page table entries allow read and write
 * access for various reasons. Kernel mappings never have any
 * access restrictions.
 *
 * The pmap_cold_mapident(pa, cnt) function maps cnt pages
 * starting at address pa to virtual pages starting at the same
 * address.
 * The pmap_cold(pa, va, cnt) function does the same thing,
 * but now the virtual address can be specified.
 */

/*
 * Map first 4 MiB of VAS to first 4 MiB of PAS
 * The atop function converts an address to the number of
 * the page containing that address.
 */
pmap_cold_mapident(0, atop(NBPDR));

```

```

/* Map second 4 MiB of VAS to first 4 MiB of PAS */
pmap_cold_map(0, NBPDR, atop(NBPDR));

/*
 * Identity map the kernel image.
 * KERNBASE - physical load address of the kernel.
 */
pmap_cold_mapident(KERNBASE, atop(KERNend - KERNBASE));

/* Map page table directory. */
pmap_cold_mapident((u_long)IdlePTD, NPGPTD);

/* Map early KPTmap. It is really pmap_cold_mapident. */
pmap_cold_map(KPTphys, (u_long)KPTmap, NKPT);

/*
 * The cr3 register contains the physical address of the
 * page directory
 */
cr3 = (u_int)IdlePTD;
load_cr3(cr3);
/* Set flag in cr0 that enables paging. */
load_cr0(rcr0() | CR0_PG);

/* We're now using virtual addresses. */

/*
 * Remove the lowest part of the double mapping of low memory
 * to get some null pointer checks.
 */
IdlePTD[0] = 0;
load_cr3(cr3);      /* Invalidate TLB */
}

```

Listing 14: i386/i386/pmap.c: Simplified implementation of the `pmap_cold` function.

After calling `pmap_cold`, the `init386` procedure is called to perform most of the machine-dependent initialization of the whole system. Only a small part of the procedure is concerned with initializing the VM subsystem. The `getmemsize` function is called to detect available segments of physical memory, as well as test it by writing and reading back bit patterns from every page of physical memory (other than the ones containing the kernel). In order to be able to do that, pieces of kernel address space have to be reserved, for instance to temporarily map the physical page being tested. This is done in the `pmap_bootstrap` procedure, which is called from `getmemsize`. Its simplified code is presented below.

```

/*
 * Bootstrap the system enough to run with virtual memory.
 *
 * On the i386 this is called after mapping has already been enabled
 * in locore.s with the page table created in pmap_cold(),
 * and just syncs the pmap module with what has already been done.
 */

void
pmap_bootstrap(vm_paddr_t firstaddr)
{
    vm_offset_t va;
    pt_entry_t *pte;
    struct pcpu *pc;
    int i;

    /*
     * Add a physical memory segment (vm_phys_seg) corresponding to the
     * preallocated kernel page table pages so that vm_page structures
     * representing these pages will be created. The vm_page structures
     * are required for promotion of the corresponding kernel virtual
     * addresses to superpage mappings.
     */
    vm_phys_add_seg(KPTphys, KPTphys + ptoa(nkpt));

    /* Initialize the first available kernel virtual address. */
    virtual_avail = (vm_offset_t)firstaddr;

    virtual_end = VM_MAX_KERNEL_ADDRESS;

    /*
     * Initialize the kernel pmap (which is statically allocated).
     */
    kernel_pmap->pm_pdir = IdlePTD;

    /*
     * Reserve some special page table entries/VA space for temporary
     * mapping of pages.
     */
#define SYSMAP(c, p, v, n) \
    v = (c)va; va += ((n)*PAGE_SIZE); p = pte; pte += (n);

    va = virtual_avail;
    pte = vtopte(va);

```

```

/*
 * Initialize temporary map objects on the current CPU for use
 * during early boot.
 * CMAP1/CMAP2 are used for zeroing and copying pages.
 * CMAP3 is used for the boot-time memory test.
 */
pc = get_pcpu();
mtx_init(&pc->pc_cmap_lock, "SYSMAPS", NULL, MTX_DEF);
SYSMAP(caddr_t, pc->pc_cmap_pte1, pc->pc_cmap_addr1, 1)
SYSMAP(caddr_t, pc->pc_cmap_pte2, pc->pc_cmap_addr2, 1)
SYSMAP(vm_offset_t, pte, pc->pc_qmap_addr, 1)

SYSMAP(caddr_t, CMAP3, CADDR3, 1);

/*
 * KPTmap is used by pmap_kextract().
 *
 * KPTmap is first initialized by pmap_cold. However, that initial
 * KPTmap can only support NKPT page table pages. Here, a larger
 * KPTmap is created that can support KVA_PAGES page table pages.
 */
SYSMAP(pt_entry_t *, KPTD, KPTmap, KVA_PAGES)

for (i = 0; i < NKPT; i++)
    KPTD[i] = (KPTphys + ptoa(i)) | PG_RW | PG_V;

virtual_avail = va;
}

```

Listing 15: i386/i386/pmap.c: Simplified implementation of the `pmap_bootstrap` function.

The last step in initializing the `pmap` module comes after kernel memory allocators have been bootstrapped and are ready to use. The `pmap_init` procedure allocates global data structures used by the module and initializes the `vm_page` structures used for the kernel page table. Simplified code of the procedure is listed below.

```

/*
 * Initialize the pmap module.
 * Called by vm_init, to initialize any structures that the pmap
 * system needs to map virtual memory.
 */
void
pmap_init(void)
{
    vm_page_t mpte;

```

```

vm_size_t s;
int i, pv_npg;

/*
 * Initialize the vm_page array entries for the kernel pmap's
 * page table pages.
 */
for (i = 0; i < NKPT; i++) {
    mpte = PHYS_TO_VM_PAGE(KPTphys + ptoa(i));
    mpte->pindex = i + KPTDI;
    mpte->phys_addr = KPTphys + ptoa(i);
    mpte->wire_count = 1;
}
vm_wire_add(NKPT);

/*
 * Initialize the address space for the pv_entries. Set a
 * high water mark so that the system can recover from excessive
 * numbers of pv entries.
 */

/*
 * shpgperproc - maximum number of pages that a process can
 * share with other processes
 * pv_entry_max - maximum number of pv_entries in the system
 */
pv_entry_max = shpgperproc * maxproc + vm_cnt.v_page_count;
pv_entry_max = roundup(pv_entry_max, _NPCPV);
pv_entry_high_water = 9 * (pv_entry_max / 10);

/*
 * Calculate the size of the pv head table for superpages.
 */
pv_npg = trunc_4mpage(vm_phys_segs[vm_phys_nsegs - 1].end -
    PAGE_SIZE) / NBPDR + 1;

/*
 * Every struct vm_page in the global vm_page array represents
 * a physical page of size PAGE_SIZE has a pv_entry list header
 * embedded into it. We also need a pv_entry list header for each
 * physical superpage. Here we allocate an array containing a
 * pv_entry list header for each physical superpage.
 */
s = (vm_size_t)(pv_npg * sizeof(struct md_page));
s = round_page(s);
pv_table = (struct md_page *)kmem_malloc(s, M_WAITOK | M_ZERO);
for (i = 0; i < pv_npg; i++)
    TAILQ_INIT(&pv_table[i].pv_list);

```



```
/*
 * A pv_chunk is a page-sized structures holding pv_entries associated
 * with the same pmap. Here we determine how much address space we need
 * to reserve to be able to hold as many pv_chunks as required.
 */
pv_maxchunks = MAX(pv_entry_max / _NPCPV, maxproc);
/* Allocate kernel address space for the chunks. */
pv_chunkbase = (struct pv_chunk *)kva_alloc(PAGE_SIZE * pv_maxchunks);
pmap_ptelist_init(&pv_vafree, pv_chunkbase, pv_maxchunks);

pmap_initialized = 1;
/* Initialize the trampoline area. */
pmap_init_trm();
}
```

Listing 16: i386/i386/pmap.c: Simplified implementation of the `pmap_init` function.

Chapter 5

Conclusion

There remain many aspects of the FreeBSD virtual memory subsystem which have not been discussed in this thesis. However, the goal of this thesis has been not to provide a complete treatment, but to explain at a high level its foundations, which are the data structures and the role they play in performing the system's basic functions, while at the same time showing that real kernel source code is based on really simple algorithms and ideas.

Unfortunately, real kernel code is less readable due to performance and concurrency concerns. Nevertheless, this text should make it easier for the reader to see the bigger picture and not be overwhelmed by the amount of noise while reading the sources.

Bibliography

- [1] Andrew S. Tanenbaum and Herbert Bos. 2014. *Modern Operating Systems* (4th ed.). Prentice Hall Press, Upper Saddle River, NJ, USA.
- [2] Main Page of the FreeBSD project: <https://www.freebsd.org/> (accessed 30/08/2018)
- [3] *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3: System Programming Guide*, available at <https://software.intel.com/en-us/articles/intel-sdm> (accessed 30/08/2018)
- [4] Marshall Kirk McKusick, George Neville-Neil, and Robert N.M. Watson. 2014. *The Design and Implementation of the FreeBSD Operating System* (2nd ed.). Addison-Wesley Professional.
- [5] Marshall Kirk McKusick and George V. Neville-Neil. 2004. *The Design and Implementation of the FreeBSD Operating System* (1st ed.). Pearson Education.
- [6] Moritz Lipp, Michael Schwarz, Daniel Gruss, Thomas Prescher, Werner Haas, Stefan Mangard, Paul Kocher, Daniel Genkin, Yuval Yarom and Michael Hamburg. 2018. *Meltdown*. <https://meltdownattack.com> (accessed 30/08/2018)
- [7] Paul Kocher, Daniel Genkin, Daniel Gruss, Werner Haas, Michael Hamburg, Moritz Lipp, Stefan Mangard, Thomas Prescher, Michael Schwarz and Yuval Yarom. 2018. *Spectre Attacks: Exploiting Speculative Execution*. <https://meltdownattack.com> (accessed 30/08/2018)
- [8] Daniel Dominic Sleator and Robert Endre Tarjan. 1985. *Self-adjusting binary search trees*. J. ACM 32, 3 (July 1985), 652-686.
- [9] *FreeBSD Manual Pages*, `execve(2)`, <https://www.freebsd.org/cgi/man.cgi?query=execve> (accessed 30/08/2018)
- [10] *FreeBSD Manual Pages*, `mmap(2)`, <https://www.freebsd.org/cgi/man.cgi?query=mmap> (accessed 30/08/2018)
- [11] *FreeBSD Manual Pages*, `fork(2)`, <https://www.freebsd.org/cgi/man.cgi?query=fork> (accessed 30/08/2018)

- [12] *FreeBSD Manual Pages*, `pmap(9)`, <https://www.freebsd.org/cgi/man.cgi?query=pmap> (accessed 30/08/2018)