Gabriele Proietti Mattia

Advanced Operating Systems and Virtualization

[2] Step 4: Kernel Boot



Department of Computer, Control and Management Engineering "A. Ruberti", Sapienza University of Rome

gpm.name · proiettimattia@diag.uniroma1.it

A.Y. 2020/2021 · v6

Outline

[2] Step 4: Kernel Boot

- 1. Initial Life of the Linux Kernel
- startup_32()
- 3. start_kernel()
 - 1. A Primer on Memory Organization
 - 2. Bootmem and Memblock Allocators
 - 3. Paging Introduction
 - 4. Paging Initialization
 - 5. TLB
 - 6. Final operations and recap

2.1

2. Step 4: Kernel Boot

Initial Life of the Linux Kernel



Advanced Operating Systems and Virtualization

Outline

- 1. **BIOS/UEFI** Actual hardware setup
- 2. **Bootloader Stage 1** Executes the stage 2 bootloader (skipped for UEFI)
- 3. **Bootloader Stage 2** Loads and starts the kernel

4. Kernel

Takes control and initializes the machine (machine-dependent operations)

5. Init (or systemd)

First process: basic environment initialization

6. Runlevels/Targets

Initializes the user environment

Initial Life

The Stage 2 Bootloader (or UEFI) loads in RAM the image of the kernel. This image is really different from the one that we have at steady state. We remind that when CPU starts it is in Real Mode with 1MB of addressable memory.

Where is the **entry point**? When the kernel switches to **Protected Mode**?



https://manybutfinite.com/post/kernel-boot-process/ http://lxr.linux.no/#linux+v2.6.25.6/Documentation/i386/boot.txt

From Real to Protected Mode

Real Mode

v2.6



https://manybutfinite.com/post/kernel-boot-process/

First Instruction

v2.6

Real Mode

The first executed instruction is a 2-byte jump to start_of_setup directly written in machine code.

110	.globl	_start
111	_start:	
112		# Explicitly enter this as bytes, or the assembler
113		<pre># tries to generate a 3-byte jump here, which causes</pre>
114		<pre># everything else to push off to the wrong offset.</pre>
115		.byte 0xeb # short (2-byte) jump
116		.byte start_of_setup-1f

https://elixir.bootlin.com/linux/v2.6.25.6/source/arch/x86/boot/header.S#L110

start_of_setup()

This is a short routine that makes some initial setup:

- it sets up a **stack**
- zeroes the **bss** section
- jumps to main() in arch/x86/boot/main.c

In this portion of code the kernel still runs in real mode, and the function implements part of the Kernel Boot Protocol, for example it loads the boot options in memory.

V2.6

main()

v2.6

Real Mode

The goal of the <u>main()</u> function is to prepare the machine to enter protected and then to the switch. Therefore it:

- enables the A20 line
- sets up the Interrupt Descriptor Table (IDT) and the Global Descriptor Table (GDT)
- sets up memory, asks BIOS which is the available memory for creating a physical address map. As a general rule, the kernel is installed in RAM starting from the physical address oxoo100000, i.e. from the second megabyte. For kernel 2.6, a typical amount of required RAM is 3MB.

In the end the function calls go_to_protected_mode() in arch/x86/boot/pm.c

go_to_protected_mode()

```
99
       1*
100
        * Actual invocation sequence
        */
101
       void go_to_protected_mode(void)
102
103
       {
               /* Hook before leaving real mode, also disables interrupts */
104
               realmode switch hook();
105
106
107
               /* Enable the A20 gate */
               if (enable_a20()) {
108
                       puts("A20 gate not responding, unable to boot...\n");
109
110
                       die();
               }
111
112
               reset_coprocessor();
               /* Mask all interrupts in the PIC */
116
               mask_all_interrupts();
117
118
119
               /* Actual transition to protected mode... */
120
               setup_idt();
               setup_gdt();
121
               protected_mode_jump(boot_params.hdr.code32_start,
122
123
                                   (u32)\&boot params + (ds() << 4));
124
```

https://elixir.bootlin.com/linux/v5.11/source/arch/x86/boot/pm.c#L99

V5.11

Real Mode

go_to_protected_mode()

Interrupt Descriptor Table

V5.11 Real Mode

In real mode the Interrupt Vector Table is always at address o. The IDTR register is set up in the following way:

90	/*
91	* Set up the IDT
92	*/
93	<pre>static void setup_idt(void)</pre>
94	{
95	<pre>static const struct gdt_ptr null_idt = {0, 0};</pre>
96	asm volatile ("lidtl %0" : : "m" (null_idt));
97	}
0.0	

https://elixir.bootlin.com/linux/v5.11/source/arch/x86/boot/pm.c#L93

go_to_protected_mode()

Global Descriptor Table (GDT)

64	<pre>static void setup_gdt(void)</pre>
65	{
66	/* There are machines which are known to not boot with the GDT
67	being 8-byte unaligned. Intel recommends 16 byte alignment. *
68	<pre>static const u64 boot_gdt[]attribute((aligned(16))) = {</pre>
69	/* CS: code, read/execute, 4 GB, base 0 */
70	<pre>[GDT_ENTRY_BOOT_CS] = GDT_ENTRY(0xc09b, 0, 0xfffff),</pre>
71	/* DS: data, read/write, 4 GB, base 0 */
72	<pre>[GDT_ENTRY_BOOT_DS] = GDT_ENTRY(0xc093, 0, 0xfffff),</pre>
73	/* TSS: 32-bit tss, 104 bytes, base 4096 */
74	<pre>/* We only have a TSS here to keep Intel VT happy;</pre>
75	we don't actually use it for anything. */
76	[GDT_ENTRY_BOOT_TSS] = GDT_ENTRY(0x0089, 4096, 103),
77	};
78	/* Xen HVM incorrectly stores a pointer to the gdt_ptr, instead
79	of the gdt_ptr contents. Thus, make it static so it will
80	stay in memory, at least long enough that we switch to the
81	proper kernel GDT. */
82	static struct gdt_ptr gdt ;
83	
84	<pre>gdt.len = sizeof(boot_gdt)-1;</pre>
85	<pre>gdt.ptr = (u32)&boot_gdt + (ds() << 4);</pre>
86	
87	<pre>asm volatile("lgdtl %0" : : "m" (gdt));</pre>
88	}

https://elixir.bootlin.com/linux/v5.11/source/arch/x86/boot/pm.c#L64

protected_mode_jump()



V5.11

After setting the initial IDT and GDT, the kernel jumps to protected via protected_mode_jump() in arch/x86/boot/pmjump.S. This routine:

- sets PE in CRo
- issues a ljmp to its very next instruction to load in CS the boot CS sector
- sets up a data segment for flat 32-bit mode
- sets up a temporary stack

37		movl	%сг0, %edx		
38		огь	\$X86_CR0_PE, %dl	#	Protected mode
39		movl	%edx, %cr0		
40					
41		# Trans	ition to 32- bit mode		
42		.byte	0x66, 0xea	#	ljmpl opcode
43	2:	.long	.Lin_pm32	#	offset
44		.word	BOOT_CS	#	segment

https://elixir.bootlin.com/linux/v5.11/source/arch/x86/boot/pmjump.S#L24

startup_32() #primary

protected_mode_jump() jumps into <u>startup 32()</u> arch/x86/boot/compressed/head_32.S and this routine does the following:

- sets the segments to known values (__BOOT_DS)
- loads a new stack
- clears again the BSS section
- determines the actual position in memory via a call/pop (image below)
- calls <u>extract kernel()</u> (previously named decompress_kernel())

50	/*					
51	* Calculate	e the delta between where we were compiled to run				
52	* at and wh	* at and where we were actually loaded at. This can only be done				
53	* with a sh	nort local call on x86. Nothing else will tell us wha				
54	* address v	we are running at. The reserved chunk of the real-mode				
55	* data at 🤅	<pre>Dx1e4 (defined as a scratch field) are used as the stack</pre>				
56	<pre>* for this</pre>	calculation. Only 4 bytes are needed.				
57	*/					
58	leal	(BP_scratch+4)(%esi), %esp				
59	call	l 1f				
60	1: pop	%edx				
61	addl	l \$_GLOBAL_OFFSET_TABLE_+(1b), % edx				

https://elixir.bootlin.com/linux/v5.11/source/arch/x86/boot/compressed/head_32.S#L50

V5.11

Protected

KASLR

Kernel Address Space Layout Randomization

In order to prevent that an attacker patches the kernel memory image, at the boot time the kernel **randomly decides** where to decompress itself in memory relying on the most accurate source of entropy available. However, since the kernel in mapped using 2MB aligned pages, the number of valid slots is limited.

The current layout of the kernel's virtual address space only leaves **512M for the kernel code**—and 1.5G for modules. Since there is no need for that much module space, his patches reduce that to 1G, leaving 1G for the kernel, thus 512 possible slots (as it needs to be 2M aligned). The number of slots may increase when the modules' location is added to KASLR.

-- https://lwn.net/Articles/569635/

2.2

2. Step 4: Kernel Boot

startup_32()



Advanced Operating Systems and Virtualization

startup_32() #secondary

After the decompression the true image of kernel can run, and this is done by a jump to <u>startup 32()</u> at arch/x86/kernel/head_32.S. This routine sets up the environment for the first Linux process (process 0):

- initializes the segmentation registers with their final values
- clears again the bss
- builds the page table
- enables paging
- creates the final IDT
- jumps to the architecture-dependent kernel entry point (i.e. start_kernel() at init/kernel.c)

Memory

During the initialization the steady-state kernel must take control of the available physical memory. This because it will have to manage it with respect to the virtual address spaces of all processes, in particular it needs to be able to:

- allocate and deallocate memory
- swap

For this reason, upon starting, the kernel must have an early organization setup out of the box. For this reason the kernel use a set of statically generated page tables.

Process Page Tables

On 32bit architecture, the process address space is divided in two parts:

User/Kernel

- linear addresses from 0x00000000 to 0xbfffffff (about 3GB) can be addressed when a process runs in User or Kernel Mode
- linear addresses from 0xc0000000 to 0xfffffff (about 1GB) can be addressed when a process runs in Kernel Mode

When the process runs in User Mode it issues linear addresses < 0xc0000000, when in Kernel Mode >= 0xc00000000.

What should be kept in mind is that addresses lower than 0xc0000000 (value often referred as PAGE_0FFSET) depend on the specific process, the others **are the same for every process** and **equal to the corresponding entries of the Master Kernel Page General Directory**.

Only Kernel

The kernel maintains a set of page tables of its own use rooted at a so-called "**Master Kernel Page Global Directory**". After the system initialization this set of pages tables is never used by any process or kernel thread, but the highest entries will be the reference model for the corresponding entries of the Page Global Directories of every regular process in the system (we will see that every process has a PGD).

The setup of these tables is a two step activity:

- the kernel first creates a limited address space, including code, data, the initial Page Tables and a dynamic area (of 128KB) -- this structure is known at compile time
- 2. the kernel takes advantage of all of the existing RAM and sets up the page table properly

Provisional Kernel Page Tables

A provisional **Page Global Directory** (PGD) is initialized statically during the kernel compilation, while the provisional **Page Tables** are initialized by startup_32().



Figure 2-13. The first 768 page frames (3 MB) in Linux 2.6

Bovet, Daniel P., and Marco Cesati. Understanding the Linux Kernel: from I/O ports to process management. "O'Reilly Media, Inc.", 2005.

Provisional Kernel Page Tables

The **Page Global Directory** (PGD) is stored in the swapper_page_dir variable. Now suppose that all the kernel segments, the provisional page tables and the dynamic area fits 8MB of RAM. In the early paging, with pages of size 4MB we needed 2 entries in the Page Table.

Now, the objective of this phase of paging is to allow these 8MB of RAM to be easily addressed both in real mode and protected mode. Therefore the kernel must create a mapping from both the linear address 0x00000000 through 0x007fffff and the linear addresses 0xc0000000 through 0xc07fffff into the physical 0x00000000 through 0x007fffff.

This mapping will be explained in Slide 61.

Provisional Kernel Page Tables



Enabling Paging

https://elixir.bootlin.com/linux/v5.11/source/arch/x86/kernel/head_32.S#L31



https://elixir.bootlin.com/linux/v5.11/source/arch/x86/kernel/head 32.S#L250

Kernel-Level MM Data Structures

The main data structures for memory management in the kernel are:

- **Kernel Page Tables**, that keeps the memory mapping for kernel level code and data, it will pointed by swapper_pg_dir
- **Core Map**, that keeps the status information for any frame (or page) of the physical memory and the free memory frames for any NUMA node



2.3

2. Step 4: Kernel Boot

start_kernel()



Advanced Operating Systems and Virtualization

Kernel Initialization



Cores Initialization

start_kernel() executes on a single core (the master). All of the other cores keep waiting
that the master has finished.

The kernel uses the function smp_processor_id() for obtaining the ID of the current core. The function is architecture-dependent a written in assembly code by using a specific hardware identification protocol. In modern version it uses the APIC. The function can be used both at kernel startup and at steady state.

Kernel initialization signature

V5.11

The <u>start kernel()</u> function is declared as

asmlinkage __visible void __init __no_sanitize_address start_kernel(void)

Where:

- **asmlinkage** tells the compiler that the calling convention is such that parameters are passed on stack
- __visible prevent Link-Time Optimization (since gcc 4.5)
- __init tells the kernel that the function is only used at initialization phase so memory can be freed' after
- __no_sanitize_address prevent address sanitizing (since gcc 4.8)

Main operations

The main operations carried out by <u>start kernel()</u> (init/main.c) are:

- 1. setup_arch() that initializes the architecture
- 2. build_all_zonelists() builds the memory zones
- 3. page_alloc_init() / mem_init() the steady state allocator (Buddy System) is initialized and the boot one removed
- 4. sched_init() initializes the scheduler
- 5. trap_init() the final IDT is built
- 6. time_init() the system time is initialized
- 7. kmem_cache_init() the slab allocator is initialized
- 8. arch_call_rest_init() / <u>rest_init()</u> prepares the environment
 - a. kernel_thread(kernel init) starts the kernel thread for process 1 is created
 - i. <u>kernel init freeable()</u> -> <u>prepare namespace()</u> -> initrd_load() mounts the initramfs, a temporary filesystem used to start the init process
 - ii. <u>run init process()</u> -> kernel_execve() Execute /bin/init
 - b. <u>cpu startup entry()</u> -> <u>do idle()</u> starts the idle process

setup_arch()

The main operations carried out by <u>setup arch()</u> (/arch/x86/kernel/setup.c) are:

- 1. load_cr3() initializes kernel page tables
- 2. __flush_tlb_all() flush the TLB
- 3. init_bootmem() initializes the bootmem allocator (v < 5)
- 4. e820__memory_setup() / e820__reserve_resources() initializes the available memory (also for memblock allocator)
- 5. <u>x86 init.paging.pagetable init()</u> -> paging_init() initializes paging
 - -> <u>native pagetable init()</u> ->

2.3.1

2. Step 4: Kernel Boot
3. start_kernel()

A Primer on Memory Organization



Advanced Operating Systems and Virtualization

NUMA

Linux is available for a great number of architectures, for this reason a machine independent way of describing memory is needed.

Large scale machines memory may be arranged into banks that incur a different access delay depending on the distance from the CPU. For this reason a memory bank can be assigned to each CPU, or a bank can be suitable for Direct Memory Access (DMA) near the devices.

Each of these banks is called in linux a **node** and the concept of accessing the memory in nodes is called *Non-Uniform Memory Access* (**NUMA**) (with a single node we have a UMA architecture). Each node is represented by the struct pg_data_t and all nodes are kept in linked list.



Possible node0 for CPU0



Zones

Each node is divided in a number of blocks called zones, which represents ranges within the memory. On x86 there are three kinds of zone:

- ZONE_DMA is **directly mapped** by the kernel in the lower part of memory and it is destined to ISA (Industry Standard Architecture) devices, in x86 first 16 MB
- ZONE_NORMAL is **directly mapped** by the kernel into the upper region of the linear address space, in x86 from 16MB to 896MB
- ZONE_HIGHMEM is the remaining available memory and it is **not directly mapped** by the kernel, in x86 from 896MB to end of memory.

The Page table is usually located at the top beginning of ZONE_NORMAL. To access memory between 1GB and 4GB the kernel temporarily maps pages from high memory to ZONE_NORMAL.

ZONE_NORMAL is fixed in size, addressing 16GiB can require 176MB of data structures!



Process Virtual Address Space



Physical Memory on x86



ZONE_DMA
2.3.2

2. Step 4: Kernel Boot
3. start_kernel()

Bootmem and Memblock Allocators



Advanced Operating Systems and Virtualization

Bootmem allocator

In the previous section we concluded that the memory map of the initial kernel image is known at **compile** time. It's impractical to initialize all of the core kernel memory structures at compile time. The Linux kernel have a link-time memory manager, embedded into the kernel image, called **bootmem allocator** (linux/bootmem.h).

The Bootmem allocator relies on **bitmaps** (instead of linked list of free blocks) that tells if any 4KB page in the currently reachable memory is busy or free. It also offers **API** (only at boot time) to get free buffers, i.e. sets of contiguous page-aligned areas.

Bootmem allocator is a *First Fit* allocator. To satisfy the allocations that are less than a page, the allocator records the last allocated Page Frame Number (PFN) and the offset. Subsequent allocations are stored in the same page.

Bootmem organization



Initializing bootmem

The information used by the bootmem allocator is represented by struct bootmem_data. An array to hold up to MAX_NUMNODES such structures is statically allocated and then it is discarded when the system initialization completes. Each entry in this array corresponds to a node with memory. For UMA systems only entry o is used.

Initializing bootmem

The bootmem allocator is initialized during early architecture specific setup. Each architecture is required to supply a setup_arch() (called by start_kernel()) function which, among other tasks, is responsible for acquiring the necessary parameters to initialise the boot memory allocator. These parameters define limits of usable physical memory:

- **min_low_pfn** the lowest PFN that is available in the system
- max_low_pfn the highest PFN that may be addressed by low memory (ZONE_NORMAL)
- **max_pfn** the last PFN available to the system.

After those limits are determined, the init_bootmem() or init_bootmem_node() function should be called to initialize the bootmem allocator. The UMA case should use the init_bootmem function. It will initialize contig_page_data structure that represents the only memory node in the system. In the NUMA case the init_bootmem_node function should be called to initialize the bootmem allocator for each node.

Bootmem APIs

#include <linux/bootmem.h>

void *alloc_bootmem(unsigned long size);

Allocate size number of bytes from ZONE_NORMAL. The allocation will be aligned to the L1 hardware cache to get the maximum benefit from the hardware cache

void *alloc_bootmem_low(unsigned long size);

Allocate size number of bytes from ZONE_DMA. The allocation will be aligned to the L1 hardware cache

void *alloc_bootmem_pages(unsigned long size);

Allocate size number of bytes from ZONE_NORMAL aligned on a page size so that full pages will be returned to the caller

void*alloc_bootmem_low_pages(unsignedlongsize);Allocate size number of bytes from ZONE_NORMAL aligned on a page size so that full pages will be returned to the callervoid free_bootmem(unsigned long addr, unsigned long size);

Bootmem API is only available for code linked in the kernel image.

https://www.kernel.org/doc/gorman/html/understand/understandoo8.html

Memblock is a method of managing memory regions during the early boot period when the usual kernel memory allocators are not up and running. The memblock allocator, differently from the bootmem does not use bitmaps for keeping track of allocated regions, but collections of regions.

Memblock views the system memory as collections of contiguous regions. There are several types of these collections:

- memory describes the physical memory available to the kernel; this may differ from the actual physical memory installed in the system, for instance when the memory is restricted with mem= command line parameter
- reserved describes the regions that were allocated
- physmem describes the actual physical memory available during boot regardless of the possible restrictions and memory hot(un)plug; the physmem type is only available on some architectures.

https://www.kernel.org/doc/html/latest/core-api/boot-time-mm.html



V5.11

Architecture

Each region is represented by struct memblock_region that defines the region extents, its attributes and NUMA node id on NUMA systems. Every memory type is described by the struct memblock_type which contains an array of memory regions along with the allocator metadata.

The "memory" and "reserved" types are nicely wrapped with struct memblock. This structure is statically initialized at build time. The region arrays are initially sized to INIT_MEMBLOCK_REGIONS for "memory" and INIT_MEMBLOCK_RESERVED_REGIONS for "reserved". The region array for "physmem" is initially sized to INIT_PHYSMEM_REGIONS.

The memblock_allow_resize() enables automatic resizing of the region arrays during addition of new regions. This feature should be used with care so that memory allocated for the region array will not overlap with areas that should be reserved, for example initrd.

API

After the initialization of memory regions done by setup_arch() with functions memblock_add() or memblock_add_node() functions. We can use the following APIs:

- memblock_phys_alloc*() these functions return the physical address of the allocated memory
 - memblock_phys_alloc_range(phys_addr_t size, phys_addr_t align, phys_addr_t start, phys_addr_t end);
 - o memblock_phys_alloc_try_nid(phys_addr_t size, phys_addr_t align, int nid);
 - o memblock_phys_alloc(phys_addr_t size, phys_addr_t align)
- memblock_alloc*() these functions return the virtual address of the allocated memory:
 - o memblock_alloc(phys_addr_t size, phys_addr_t align)
 - memblock_alloc_raw(phys_addr_t size, phys_addr_t align)
 - memblock_alloc_from(phys_addr_t size, phys_addr_t align, phys_addr_t min_addr)
 - o memblock_alloc_low(phys_addr_t size, phys_addr_t align)
 - o memblock_alloc_node(phys_addr_t size, phys_addr_t align, int nid)

After boot

As the system boot progresses, the architecture specific mem_init() function frees all the memory to the buddy page allocator.

Unless an architecture enables CONFIG_ARCH_KEEP_MEMBLOCK, the memblock data structures (except "physmem") will be discarded after the system initialization completes.

Memblock use case

e820

In the x86 architecture, the e820 is the shorthand for obtaining the memory map of the system. At boot time, by looking at the dmesg you will find:

+0.000000]	e820: BIOS-pr	ovided physical RAM map:
+0.000000]	BIOS-e820: [m	em 0x00000000000000000000000000000000000
+0.000000]	BIOS-e820: [m	em 0x000000000009f800-0x000000000009ffff] reserved
+0.000000]	BIOS-e820: [m	em 0x00000000000f0000-0x00000000000fffff] reserved
+0.000000]	BIOS-e820: [m	em 0x00000000000100000-0x00000000bfd9ffff] usable
+0.000000]	BIOS-e820: [m	em 0x00000000bfda0000-0x0000000bfdd0fff] ACPI NVS
+0.000000]	BIOS-e820: [m	em 0x00000000bfdd1000-0x00000000bfdfffff] ACPI data
+0.000000]	BIOS-e820: [m	em 0x00000000bfe00000-0x0000000bfefffff] reserved
+0.000000]	BIOS-e820: [m	em 0x00000000000000000000000000000000000
+0.000000]	BIOS-e820: [m	em 0x00000000fec00000-0x00000000ffffffff] reserved
+0.000000]	BIOS-e820: [m	em 0x0000000100000000-0x000000043effffff] usable

The e820 subsystem must mark some specific areas to be reserved, for example the can be reserved for the kernel, for BIOS facilities or drivers. The code on the right just does this by using memblock.

```
1129
1130
         * Mark E820 reserved areas as busy for the resource manager
         */
                                                                   V5.11
        static struct resource __initdata *e820_res;
1134
        void __init e820__reserve_resources(void)
1136
                int i:
1138
                struct resource *res:
1139
                u64 end:
1140
1141
                res = memblock_alloc(sizeof(*res) * e820_table->nr_entries,
1142
                                     SMP_CACHE_BYTES);
1143
                if (!res)
                        panic("%s: Failed to allocate %zu bytes\n", __func__,
1144
                              sizeof(*res) * e820 table->nr entries);
1145
1146
                e820 res = res:
1147
1148
                for (i = 0; i < e820_table->nr_entries; i++) {
                        struct e820_entry *entry = e820_table->entries + i;
1149
                        end = entry->addr + entry->size - 1;
                        if (end != (resource_size_t)end) {
                                res++:
                                continue:
1156
                        res->start = entry->addr;
                        res->end = end:
1158
                        res->name = e820_type_to_string(entry);
                        res->flags = e820_type_to_iomem_type(entry);
                        res->desc = e820_type_to_iores_desc(entry);
1162
                         * Don't register the region that could be conflicted with
1164
                         * PCI device BAR resources and insert them later in
                         * pcibios resource survey():
1166
                         */
1167
                        if (do_mark_busy(entry->type, res)) {
1168
                                res->flags |= IORESOURCE BUSY:
1169
                                insert_resource(&iomem_resource, res);
1170
                        res++:
          https://elixir.bootlin.com/linux/v5.11.2/source/arch/x86/kernel/e820.c#L1129
1172
```

1*

Superseding

In recent versions of the kernel (5+), the bootmem allocator has been removed in favour of the memblock allocator on almost all architectures. See this patch <u>https://lwn.net/Articles/764807/</u>

mm: remove bootmem allocator

From:Mike Rapoport <rppt-AT-linux.vnet.ibm.com>To:linux-mm-AT-kvack.orgSubject:[PATCH 00/30] mm: remove bootmem allocatorDate:Fri, 14 Sep 2018 15:10:15 +0300Message-ID:<1536927045-23536-1-git-send-email-rppt@linux.vnet.ibm.com>

Hi,

These patches switch early memory management to use memblock directly without any bootmem compatibility wrappers. As the result both bootmem and nobootmem are removed.

The patchset survived allyesconfig builds on arm, arm64, i386, mips, nds32, parisc, powerpc, riscv, s390 and x86 and most of the *_defconfig builds for all architectures except unicore32.

The patchset is based on v4.19-rc3-mmotm-2018-09-12-16-40, so I needed a small PSI fix from [1] for some of the builds.

I did my best to verify that the failures are not caused by my changes, but I may have missed something. Most defconfig build failures I've seen were caused by assembler being unhappy about unsupported opcode, wrong encoding 2.3.3

2. Step 4: Kernel Boot
3. start_kernel()

Paging Introduction



Advanced Operating Systems and Virtualization

Pages handling

Prior to version 2.6.11 the Linux paging model consisted of 3 indirection levels, next versions introduced another level of indirections for a total of 4.





Splitting the address

For splitting the linear address there are three kinds of macros that can be used:

- SHIFT macros specify the length in bits mapped to each PT level
- MASK macros AND'd with an address mask out all the upper bits and they are often used for understanding if an address is aligned to a given level within the page table
- SIZE macros reveal how many bytes are addressed by each entry at each level



Figure 3.2. Linear Address Bit Size Macros

Figure 3.3. Linear Address Size and Mask Macros

Gorman, Mel. Understanding the Linux virtual memory manager. Upper Saddle River: Prentice Hall, 2004.

Configuring the PT

Those macros are declared in the following source files:

- arch/x86/include/asm/pgtable-2level types.h
- arch/x86/include/asm/pgtable-3level types.h
- arch/x86/include/asm/pgtable 64 types.h

With other kinds of macros like PTRS_PER_* which describe the number of entries in PTs

```
27
       * PGDIR SHIFT determines what a top-level page table entry can map
29
       */
     #define PGDIR SHIFT
30
                             30
     #define PTRS PER PGD
31
                           4
32
33
     /*
34
      * PMD SHIFT determines the size of the area a middle-level
35
      * page table can map
36
       */
     #define PMD SHIFT
37
                             21
     #define PTRS PER PMD 512
38
39
40
     /*
41
      * entries per page directory level
42
       */
      #define PTRS PER PTE
43
                             512
```

https://elixir.bootlin.com/linux/v5.11.2/source/arch/x86/include/asm/pgtable-3level_types.h#L27

Page Table data structures

As already introduced earlier, swapper_pg_dir keeps the virtual memory address of the PGD (PDE) portion of the kernel page table. The data structure is initialized at compile time, depending on the memory layout defined for the kernel bootable image.

Any entry in the PGD is accessed with displacement, but the main types for defining page table entries are explicitly defined, even if they are just unsigned integers:

```
typedef struct { unsigned long pte_low; } pte_t;
typedef struct { unsigned long pmd; } pmd_t;
typedef struct { unsigned long pgd; } pgd_t;
```

This is done essentially for enforcing type protection and for supporting PAE (where additional 4 bits are used for addressing more than 4GB of RAM).

An additional structure is used for storing page protection bits: pgprot_t.

(C and weak typing)

Remember that C language is weak typed, therefore the following code will compile and execute with nor error or warning:

```
typedef unsigned long pgd_t;
typedef unsigned long pte_t;
pgd_t x; pte_t y;
x = y;
y = x;
```

Bit fields

12	#define _PAGE_BIT_PR	RESENT)
13	#define _PAGE_BIT_RW	I 1	L
14	#define _PAGE_BIT_US	SER 2	2
15	#define _PAGE_BIT_PM	IT 3	3
16	#define _PAGE_BIT_PC	CD 4	1
17	#define _PAGE_BIT_AC	CESSED	5
18	#define _PAGE_BIT_DI	ERTY 6	5
	#define _PAGE_BIT_PS	E 7	7
	#define _PAGE_BIT_PA	T 7	7
	#define _PAGE_BIT_GL	.OBAL 8	
	#define _PAGE_BIT_SO	DFTW1 9	
	#define _PAGE_BIT_SO)FTW2	
	#define _PAGE_BIT_SO	FTW3	
	#define _PAGE_BIT_PA	T_LARGE 1	12
	#define _PAGE_BIT_SO)FTW4 5	
	#define _PAGE_BIT_PK	EY_BITO	
	#define _PAGE_BIT_PK	KEY_BIT1 e	
	#define _PAGE_BIT_PK	KEY_BIT2 e	
	#define _PAGE_BIT_PK	KEY_BIT3 e	
	#define _PAGE_BIT_NX	(6	

/* is present */ /* writeable */ /* userspace addressable */ /* page write through */ /* page cache disabled */ /* was accessed (raised by CPU) */ /* was written to (raised by CPU) */

https://elixir.bootlin.com/linux/v5.11.2/source/arch/x86/include/asm/pgtable_types.h#L12

Bit fields and Masks

Type casting macros are defined in asm/page.h which takes the previous types and returns the relevant part of the struct. They are pte_val(), pmd_val(), pgd_val() and pgprot_val(). For reverse type casting we have __pte(), __pmd(), __pgd(), __pgprot().

In the following example code, a check if a page is present is carried out:

```
93
       * Create a page table and place a pointer to it in a middle page
94
        * directory entry:
95
96
        */
97
      static pte t * init one page table init(pmd t *pmd)
98
99
              if (!(pmd_val(*pmd) & _PAGE_PRESENT)) {
                       pte_t *page_table = (pte_t *)alloc_low_page();
101
102
                       paravirt_alloc_pte(&init_mm, __pa(page_table) >> PAGE_SHIFT);
                       set_pmd(pmd, __pmd(__pa(page_table) | _PAGE_TABLE));
103
104
                       BUG ON(page table != pte offset kernel(pmd, 0));
105
106
107
               return pte offset kernel(pmd, 0);
108
```

Different PD Entries

Different kind of page entries are again described with macros in /arch/x86/include/asm/pgtable_types.h

```
#define _PAGE_TABLE \
(_PAGE_PRESENT | _PAGE_RW | \
_PAGE_USER | _PAGE_ACCESSED | \
_PAGE_DIRTY)
```

```
#define _KERNPG_TABLE \
(_PAGE_PRESENT | _PAGE_RW | \
_PAGE_ACCESSED | _PAGE_DIRTY)
```

2.3.4

2. Step 4: Kernel Boot
3. start_kernel()

Paging Initialization



Advanced Operating Systems and Virtualization

As already introduced, startup_32() enables the Paging unit. While the kernel code is compiled with base address at PAGE_OFFSET + 1MB, the kernel is actually loaded at the beginning of physical memory. The initialization of kernel page tables begins at compile time, statically defining an array called swapper_pg_dir (at 0x00101000), that establishes page table entries for **2 pages of 4MB each**, pg0 and pg1. These two pointers covers the addresses from 1MB to 9MB but they are placed at PAGE_OFFSET + 1MB.



The 8MB of addressable memory must be addressed both in real mode than in protected mode. For this reason, for that memory area the physical address must be equal to the virtual one. This strategy is realized by declaring statically **4 entries** in the swapper_pg_dir:

- Entry 0 and 0x300 (768) point to pg0
- Entry 1 and 0x301 (769) point to pg1

These entries have set bits P,R/W,U/S and cleared A,D,PCD,PWD and Page Size.



The **Provisional Page Table** with only two pages is set with the following assembly instructions

The rest of kernel page tables are initialized by paging_init() called by setup_arch().



Figure 3.4. Call Graph: paging_init()

Gorman, Mel. Understanding the Linux virtual memory manager. Upper Saddle River: Prentice Hall, 2004.

pagetable_init()

The initialization of kernel page tables starts with function paging_init() that initializes the necessary pages for addressing ZONE_DMA and ZONE_NORMAL from PAGE_OFFSET.

```
for (; i < PTRS PER PGD; pgd++, i++) {</pre>
    vaddr = i*PGDIR SIZE; /* i is set to map from 3 GB */
    if (end && (vaddr >= end)) break;
    pmd = (pmd_t *) pgd;/* pgd initialized to (swapper_pg_dir+i) */
    . . . . . . . . .
    for (j = 0; j < PTRS PER PMD; pmd++, j++) {
        pte_base = pte = (pte_t *) alloc_bootmem_low_pages(PAGE_SIZE);
        for (k = 0; k < PTRS_PER_PTE; pte++, k++) {</pre>
            vaddr = i*PGDIR SIZE + j*PMD SIZE + k*PAGE SIZE;
            if (end && (vaddr >= end)) break;
            *pte = mk_pte_phys(__pa(vaddr), PAGE_KERNEL);
        }
        set pmd(pmd, pmd( KERNPG TABLE + pa(pte base)));
         . . . . . . . . .
```

The starting address vaddr is set to start from 3GB because we are mapping the virtual addresses of kernel pages.

The kernel prefers to use 4MB pages if the CPU support them, this for reducing the TLB miss rate and speeding up the address translation.

https://elixir.bootlin.com/linux/2.4.31/source/arch/i386/mm/init.c#L205

After pagetable_init() execution, all the ZONE_NORMAL and ZONE_DMA is directly **mapped** in the kernel. Memory is not allocated, is just <u>mapped</u>. Memory is allocated only for the Page Tables.



set_pmd() and __pa()/__va()

#define set_pmd(pmdptr, pmdval) (*(pmdptr) = pmdval)

Parameters are:

- pmdptr, pointing to an entry of the PMD, of type pmd_t. The value to assign, of pmd_t type is computed by using the macro

#define __pa(x)((unsigned long)(x)-PAGE_OFFSET)

Linux sets up a direct mapping from the physical address 0 to the virtual address PAGE_OFFSET at 3GB on x86. The opposite can be done using the __va(x) macro.

mk_pte_sys()

The function creates a page table entry given the physical address and the protection metadata.

311 /* This takes a physical page address that is used by the remapping functions */
312 #define mk_pte_phys(physpage, pgprot) __mk_pte((physpage) >> PAGE_SHIFT, pgprot)

https://elixir.bootlin.com/linux/2.4.22/source/include/asm-i386/pgtable.h#L312

The input parameters are:

- A frame physical address physpage, of type unsigned long
- A bit string pgprot for a PTE, of type pgprot_t

The macro builds a complete PTE entry, which includes the physical address of the target frame. The return type is pte_t and it can be then assigned to one PTE entry.

Loading the page table



https://elixir.bootlin.com/linux/2.4.22/source/arch/i386/mm/init.c#L351

V2.4





The load cr₃ instruction is directly mapped to the assembly code which loads the address of the PGD into CR₃.

v5.11

Latest versions of the kernel uses the **paravirtualization** scheme to map all of the basic functions that for example regards the mmu, like writing CR₃, creating PGD, PTE and so on. So in modern kernel you will find load_cr₃ mapped to

138 static inline void write_cr3(unsigned long x)
139 {
140 PVOP_VCALL1(mmu.write_cr3, x);
141 }

https://elixir.bootlin.com/linux/v5.12-rc2/source/arch/x86/include/asm/paravirt.h#L138

Final <u>virtual</u> kernel memory map

This is the final virtual **memory** map of the kernel, it has nothing to do with the <u>physical</u> counterpart.



Figure 3-15: Division of the kernel address space on IA-32 systems.

Mauerer, Wolfgang. Professional Linux kernel architecture. John Wiley & Sons, 2010.

As you can see the kernel may only address only 896MB because the last 128MB are reserved:

- **VMALLOC** virtual contiguous memory areas that are not contiguous in physical memory, especially for user processes.
- **persistent mappings** used for mapping highmem pages
- fixmaps virtual addresses customly mapped to selectable physical frames

2.3.5

2. Step 4: Kernel Boot
3. start_kernel()

TLB



Advanced Operating Systems and Virtualization

Implicit and explicit operations

The degree of automation in the management process of TLB entries depends on the hardware architecture. Kernel hooks exist for explicit management of TLB operations (mapped at compile time to nops in case of fully-automated TLB management)

On x86, automation is only **partial**: automatic TLB flushes occur upon updates of the CR₃ register (e.g. page table changes) but Changes inside the current page table are not automatically reflected into the TLB.
TLB Relevant events

Scale classification

- **global**: dealing with virtual addresses accessible by every CPU/core in real-time-concurrency
- **local**: dealing with virtual addresses accessible in timesharing concurrency

Typology classification

- Virtual to physical address remapping
- Virtual address access rule modification (read only vs write access)

The typical management is TLB implicit renewal via flush operations

TLB Flush Costs

Direct costs

- the latency of the firmware level protocol for TLB entries invalidation (selective vs non-selective)
- the latency for cross-CPU coordination in case of global TLB flushes

Indirect costs

- TLB renewal latency by the MMU firmware upon misses in the translation process of virtual to physical addresses and this cost depends on the amount of entries to be refilled
- Tradeoff vs TLB API and software complexity inside the kernel (selective vs non-selective flush/renewal)

Linux full TLB flush



This flushes the **entire TLB** on **all** processors running in the system (most expensive TLB flush operation). After it completes, all modifications to the page tables are globally visible. This is required after the kernel page tables, which are global in nature, have been modified.

V5.12

Linux TLB Flush

67

```
38
     #define flush tlb()
             do {
39
40
                     unsigned int tmpreg;
41
42
                     asm volatile (
43
                             "movl %%cr3, %0; # flush TLB \n"
44
                             "movl %0, %%cr3;
                                                           \n"
45
                             : "=r" (tmprea)
46
                             :: "memory");
47
             } while (0)
48
49
     /*
50
      * Global pages have to be flushed a bit differently. Not a real
51
      * performance problem because this does not happen often.
52
      */
53
     #define __flush_tlb_global()
54
             do {
55
                     unsigned int tmpreq;
56
57
                     asm volatile (
58
                             "movl %1, %%cr4; # turn off PGE
                                                                  n''
59
                             "movl %%cr3, %0; # flush TLB
                                                                  n''
60
                             "movl %0, %%cr3;
                                                                  n''
61
                             "movl %2, %%cr4; # turn PGE back on \n"
                             : "=&r" (tmpreg)
62
                             : "r" (mmu cr4 features & ~X86 CR4 PGE),
63
64
                               "r" (mmu cr4 features)
                             : "memory");
65
66
             \} while (0)
```

Linux partial TLB flush

void flush_tlb_mm(struct mm_struct *mm)

This flushes all TLB entries related to a **portion** of the userspace memory context. On some architectures (e.g. MIPS), this is required for all cores (usually it is confined to the local processor).

This is called only after an operation affecting the entire address space:

- when cloning a process with a fork()
- when, in general, there is an interaction with the Copy-On-Write protection

void flush_tlb_page(struct vm_area_struct *vma, unsigned long a);

This API flushes a single page from the TLB. The two most common uses of it are to flush the TLB after a page has been faulted in or has been paged out.

Linux partial TLB flush

void flush_tlb_range(struct mm_struct *mm, unsigned long start, unsigned long end);

This flushes all entries within the requested user space range for the mm context. This is used after a region has been moved (mremap()) or when changing permissions (mprotect()). This API is provided for architectures that can remove ranges of TLB entries quicker than iterating with flush_tlb_page().

Used when the page tables are being torn down and free'd. Some platforms cache the lowest level of the page table, which needs to be flushed when the pages are being deleted (e.g. Sparc64). This is called when a region is being unmapped and the page directory entries are being reclaimed.

Linux partial TLB flush

void update_mmu_cache(struct vm_area_struct *vma, unsigned long addr, pte_t *ptep);

Only called after a page fault completes. It tells that a new translation now exists at pte for the virtual address addr. Each architecture decides how this information should be used.

For example, Sparc64 uses the information to decide if the local CPU needs to flush its data cache. In some cases it is also used for preloading TLB entries

2.3.6

2. Step 4: Kernel Boot
3. start_kernel()

Final Operations and Recap



Advanced Operating Systems and Virtualization

Kernel Boot Flow



Main operations

The main operations carried out by start kernel() (init/main.c) are:

- 1. setup_arch() that initializes the architecture
- 2. build_all_zonelists() builds the memory zones
- 3. page_alloc_init() / mem_init() the steady state allocator (Buddy System) is initialized and the boot one removed
- 4. sched_init() initializes the scheduler
- 5. trap_init() the final IDT is built
- 6. time_init() the system time is initialized
- 7. kmem_cache_init() the slab allocator is initialized
- 8. arch_call_rest_init() / <u>rest_init()</u> prepares the environment
 - a. kernel_thread(kernel init) starts the kernel thread for process 1 is created
 - i. <u>kernel init freeable()</u> -> <u>prepare namespace()</u> -> initrd_load() mounts the initramfs, a temporary filesystem used to start the init process
 - ii. <u>run init process()</u> -> kernel_execve() Execute /bin/init
 - b. <u>cpu startup entry()</u> -> <u>do idle()</u> starts the idle process

Final GDT

Linux's GDT	Segment Selectors	
null	0x0	
reserved		
reserved		PNF
reserved		PNF
not used		PNF
not used		PNF
TLS #1	0x33	PNF
TLS #2	0x3b	APN
TLS #3	0x43	APN
reserved		
reserved		
reserved		
kernel code	Ox60 (KERNEL_CS)	
kernel data	Ox68 (KERNEL_DS)	
user code	0x73 (USER_CS)	
user data	Ox7b (USER_DS)	d
	Linux's GDT null reserved reserved reserved not used not used TLS #1 TLS #2 TLS #3 reserved reserved kernel code kernel data user code user data	Linux's GDTSegment SelectorsnullOxOreservedreservednot usednot usedOx33TLS #1Ox3bOx43Ox43reservedreservedreservedkernel codeOx60 (KERNEL_CS)user codeOx73 (USER_CS)user dataOx7b (USER_DS)

Linux's GDT	Segment Selectors
TSS	0x80
LDT	0x88
PNPBIOS 32-bit code	0x90
PNPBIOS 16-bit code	0x98
PNPBIOS 16-bit data	0xa0
PNPBIOS 16-bit data	0xa8
PNPBIOS 16-bit data	0xb0
APMBIOS 32-bit code	0xb8
APMBIOS 16-bit code	0xc0
APMBIOS data	0xc8
not used	
double fault TSS	0xf8

Bovet, Daniel P., and Marco Cesati. Understanding the Linux Kernel: from I/O ports to process management. "O'Reilly Media, Inc.", 2005.



End of the Kernel Boot

The idle loop is the **ending** of the kernel booting process.

Since the very first long jump ljmp \$0xf000,\$0xe05b at the reset vector at F000:FFF0 which activated the BIOS, we have worked hard to setup a system which is spinning forever.

This is the end of the "romantic" Kernel boot procedure: we infinitely loop into a hlt instruction or ...

Advanced Operating Systems and Virtualization

[2] Step 4: Kernel Boot

LECTURER Gabriele **Proietti Mattia**

BASED ON WORK BY <u>http://www.ce.uniroma2.it/~pellegrini/</u>



gpm.name · proiettimattia@diag.uniroma1.it

