

start_kernel之物理内存管理

by zenhumay

2012-03-26——2012-04-26

目录

start_kernel之物理内存管理.....	1
目录	2
1. 概述	5
2. cpu初始化.....	5
3. page_address_init()初始化地址散列表	11
4. setup_arch处理结构相关的内容	12
4.1 setup_arch函数全貌	12
4.1.1 early_ioremap_init	24
4.1.1.1 内核地址空间简介.....	24
4.1.1.2 函数源码.....	25
4.1.2 setup_memory_map	31
4.1.3 初始化 0 号进程变量.....	34
4.1.4 x86_configure_nx() 设置页面不可执行标志	35
4.1.5 e820_end_of_ram_pfn获得物理页面中页面数（可能包含空洞）	35
4.1.6 find_low_pfn_range	37
4.1.7 init_memory_mapping建立永久内核页表	37
4.1.7.1 find_early_table_space.....	44
4.1.7.1.1 冲突检测.....	48
4.1.7.2 kernel_physical_mapping_init物理内存页表项设置	50
4.1.7.2.1 页表分配函数one_page_table_init.....	55
4.1.7.3 early_ioremap_page_table_range_init	57
4.1.7.4 load_cr3(swapper_pg_dir)	58
4.1.8 initmem_init(0, max_pfn, acpi, k8).....	58
4.1.8.1 e820_register_active_regions.....	60
4.1.8.1.1 e820_find_active_region.....	60
4.1.8.1.2 add_active_range.....	62
4.1.8.2 sparse_memory_present_with_active_regions.....	65
4.1.8.3 setup_bootmem_allocator().....	67
4.1.8.4 acpi_numa_init.....	71
4.1.8.5 get_memcfg_numa()	76
4.1.9 paging_init	80
4.1.9.1 永久内核页表分配.....	81
4.1.9.2 临时内核映射初始化.....	85
4.1.9.3 内存区域初始化.....	86
4.1.9.3.1 zone_size_init.....	86
4.1.9.3.2 free_area_init_ndoes	87
4.1.9.3.3 free_area_init_node.....	91
4.1.9.3.4 free_area_init_core.....	97
5. 节点备用列表初始化.....	104
5.1 build_all_zonelists	104

5.2	__build_all_zonelists	106
5.3	build_zonelists	107
5.4	build_zonelists_in_node_order	109
6.	利用early_res分配内存	110
6.1	__alloc_memory_core_early.....	115
7.	虚拟文件系统early初始化.....	117
8.	中断向量表初始化.....	119
9.	内存管理初始化.....	121
9.1	mm_init函数全貌	121
9.2	mem_init伙伴系统的建立	122
9.2.1	低端内存释放free_all_bootmem.....	125
9.2.1.1	free_all_memory_core_early	127
9.2.1.2	get_free_all_memory_range.....	127
9.2.1.3	add_from_early_node_map.....	130
9.2.1.4	subtract_range	131
9.2.1.5	__free_pages_memory(start, end)	134
9.2.2	高端内存释放set_highmem_pages_init	136
9.2.3	__free_pages释放页面	139
9.2.3.1	free_hot_cold_page	139
9.2.3.2	free_one_page	142
9.2.3.3	__free_one_page函数	143
9.2.4	save_pg_dir	147
9.2.5	zap_low_mappings.....	148
9.3	kmem_cache_init slab分配器初始化	149
9.3.1	slab分配器相关数据结构	149
9.3.1.1	kmem_cache数据结构.....	149
9.3.1.2	数据结构array_cache.....	152
9.3.1.3	kmem_list3 数据结构.....	153
9.3.2	kmem_cache_init 函数	154
9.4	vmalloc_init非连续内存区初始化.....	163
10.	其它重要的初始化.....	164
10.1	初始化调度程序.....	164
10.2	初始化中断处理系统.....	164
10.3	软中断初始化.....	164
10.4	定时器中断初始化.....	164
11.	slab后续初始化kmem_cache_init_late	164
11.1	enable_cpucache.....	165
11.2	do_tune_cpucache	167
11.3	alloc_kmemlist	169
12.	启动控制台输出.....	173
13.	fork_init.....	173
14.	proc_caches_init.....	174
15.	块缓存初始化buffer_init.....	175
16.	虚拟根文件系统安装vfs_caches_init.....	176

17. 安装PROC文件系统proc_root_init	177
18. 后start_kernel时代.....	178
18.1 Kthreadd.....	179
18.2 cpu_idle.....	181
18.3 Kernel_init	182
18.3.1 子系统初始化.....	184
18.3.2 启动SHELL	185
18.3.3 Initrd以及磁盘根文件系统的安	186
19. 参考书籍.....	189

1. 概述

本章中涉及的源代码全部来自 linux 内核 2.6.34

`start_kernel` 是内核初始化的最后一部分，包含了物理内存管理数据结构、中断、文件系统等各方面的初始化。由于涉及的内容较多，不可能一下子将清楚，在这里通过关注点（物理内存、文件系统、中断、进程调度）的不同从不同的角度来剖析该函数。

由于本篇是第一次剖析该函数，所以会将该函数的所有内容都做一个大致打概述，但关注的重点是物理内存管理，其他方面的部分只会大致的介绍器功能，留待已有分析该模块时在详细讲述。

2. cpu初始化

init/main.c

```
528 asmlinkage void __init start_kernel(void)
529 {
530     char * command_line;
531     extern struct kernel_param __start__param[],
__stop__param[];
532
533     smp_setup_processor_id();
534
535     /*
536      * Need to run as early as possible, to initialize the
537      * lockdep hash:
538      */
539     lockdep_init();
540     debug_objects_early_init();
```

```
541
542     /*
543      * Set up the the initial canary ASAP:
544      */
545     boot_init_stack_canary();
546
547     cgroup_init_early();
548
549     local_irq_disable();
550     early_boot_irqs_off();
551     early_init_irq_lock_class();
```

```
552
553 /*
554  * Interrupts are still disabled. Do necessary setups, then
555  * enable them
556  */
557     lock_kernel();
558     tick_init();
559     boot_cpu_init();
```

前面这些代码主要是与 SMP CPU、中断等有关的代码，在此不关心这些，后面的讲解将略过。

```
560     page_address_init(); //初始化散列表
561     printk(KERN_NOTICE "%s", linux_banner); //输出操作系统旗
标
562     setup_arch(&command_line); //处理体系结构相关的内容
563     mm_init_owner(&init_mm, &init_task);
564     setup_command_line(command_line);
565     setup_nr_cpu_ids();
566     setup_per_cpu_areas();
```

```
567     smp_prepare_boot_cpu(); /* arch-specific boot-cpu hooks */
568
569     build_all_zonelists(); //初始化内存管理区中备选列表
570     page_alloc_init(); //
571
572     printk(KERN_NOTICE "Kernel command line: %s\n",
boot_command_line);
573     parse_early_param();
574     parse_args("Booting kernel", static_command_line,
__start__param,
575             __stop__param - __start__param,
576             &unknown_bootoption);
577     /*
578      * These use large bootmem allocations and must precede
579      * kmem_cache_init()
580      */
581     pidhash_init();
582     vfs_caches_init_early();
583     sort_main_extable();
584     trap_init();
585     mm_init();
586     /*
587      * Set up the scheduler prior starting any interrupts (such as the
588      * timer interrupt). Full topology setup happens at smp_init()
589      * time - but meanwhile we still have a functioning scheduler.
590      */
591     sched_init();
592     /*
593      * Disable preemption - early bootup scheduling is extremely
594      * fragile until we cpu_idle() for the first time.
595      */
596     preempt_disable();
```

```
597     if (!irqs_disabled()) {
598         printk(KERN_WARNING "start_kernel(): bug: interrupts
were "
599             "enabled *very* early, fixing it\n");
600         local_irq_disable();
601     }
602     rcu_init();
603     radix_tree_init();
604     /* init some links before init_ISA_irqs() */
605     early_irq_init();
606     init_IRQ();
607     prio_tree_init();
608     init_timers();
609     hrtimers_init();
610     softirq_init();
611     timekeeping_init();
612     time_init();
613     profile_init();
614     if (!irqs_disabled())
615         printk(KERN_CRIT "start_kernel(): bug: interrupts were "
616             "enabled early\n");
617     early_boot_irqs_on();
618     local_irq_enable();
619
620     /* Interrupts are enabled now so all GFP allocations are safe. */
621     gfp_allowed_mask = __GFP_BITS_MASK;
622
623     kmem_cache_init_late();
624
625     /*
```



```

626      * HACK ALERT! This is early. We're enabling the console
before
627      * we've done PCI setups etc, and console_init() must be aware
of
628      * this. But we do want output early, in case something goes
wrong.
629      */
630      console_init();
631      if (panic_later)
632          panic(panic_later, panic_param);
633
634      lockdep_info();
635
636      /*
637      * Need to run this when irq's are enabled, because it wants
638      * to self-test [hard/soft]-irq's on/off lock inversion bugs
639      * too:
640      */
641      locking_selftest();
642
643      #ifdef CONFIG_BLK_DEV_INITRD
644          if (initrd_start && !initrd_below_start_ok &&
645              page_to_pfn(virt_to_page((void *)initrd_start)) <
min_low_pfn) {
646              printk(KERN_CRIT "initrd overwritten (0x%08lx < 0x%08lx)
- "
647                  "disabling it.\n",
648                  page_to_pfn(virt_to_page((void *)initrd_start)),
649                  min_low_pfn);
650              initrd_start = 0;
651          }
652      #endif

```

```
653     page_cgroup_init();
654     enable_debug_pagealloc();
655     kmemtrace_init();
656     kmemleak_init();
657     debug_objects_mem_init();
658     idr_init_cache();
659     setup_per_cpu_pageset();
660     numa_policy_init();
661     if (late_time_init)
662         late_time_init();
663     sched_clock_init();
664     calibrate_delay();
665     pidmap_init();
666     anon_vma_init();
667 #ifdef CONFIG_X86
668     if (efi_enabled)
669         efi_enter_virtual_mode();
670 #endif
671     thread_info_cache_init();
672     cred_init();
673     fork_init(totalram_pages);
674     proc_caches_init();
675     buffer_init();
676     key_init();
677     security_init();
678     vfs_caches_init(totalram_pages);
679     signals_init();
680     /* rootfs populating might need page-writeback */
681     page_writeback_init();
682 #ifdef CONFIG_PROC_FS
683     proc_root_init();
```

```
684 #endif
685     cgroup_init();
686     cpuset_init();
687     taskstats_init_early();
688     delayacct_init();
689
690     check_bugs();
691
692     acpi_early_init(); /* before LAPIC and SMP init */
693     sfi_init_late();
694
695     ftrace_init();
696
697     /* Do the rest non-__init'ed, we're now alive */
698     rest_init();
699 }
```

3. page_address_init()初始化地址散列表

mm/highmem.c

```
407 static struct page_address_map
page_address_maps[LAST_PKMAP];
408
409 void __init page_address_init(void)
410 {
411     int i;
412
413     INIT_LIST_HEAD(&page_address_pool);
414     for (i = 0; i < ARRAY_SIZE(page_address_maps); i++)
```

```

415         list_add(&page_address_maps[i].list,
&page_address_pool);
416     for (i = 0; i < ARRAY_SIZE(page_address_htable); i++) {
417         INIT_LIST_HEAD(&page_address_htable[i].lh);
418         spin_lock_init(&page_address_htable[i].lock);
419     }
420     spin_lock_init(&pool_lock);
421 }

```

初始化一个全局的 `page_address_pool` 列表。

`Page_address_map` 结构定义如下：

```

298  * Describes one page->virtual association
299  */
300 struct page_address_map {
301     struct page *page; //物理页面
302     void *virtual; //地址
303     struct list_head list;
304 };

```

4. setup_arch处理结构相关的内容

函数所在文件:arch/x86/kernel/setup.c

4.1 setup_arch函数全貌

该函数主要处理与 `cpu` 架构相关的内容，包括内存处理等。该函数总共 300 多行，在这里只分析我比较关心的部分，在下图中已经用蓝色标出。

```

713 void __init setup_arch(char **cmdline_p)
714 {
715     int acpi = 0;
716     int k8 = 0;
717

```

```

718 #ifdef CONFIG_X86_32
719     memcpy(&boot_cpu_data, &new_cpu_data,
sizeof(new_cpu_data));
720     visws_early_detect();
721 #else
722     printk(KERN_INFO "Command line: %s\n",
boot_command_line);
723 #endif
724
725     /* VMI may relocate the fixmap; do this before touching ioremap
area */
726     vmi_init();
727
728     early_cpu_init();
729     early_ioremap_init();
730
731     ROOT_DEV = old_decode_dev(boot_params.hdr.root_dev);
732     screen_info = boot_params.screen_info;
733     edid_info = boot_params.edid_info;
734 #ifdef CONFIG_X86_32
735     apm_info.bios = boot_params.apm_bios_info;
736     ist_info = boot_params.ist_info;
737     if (boot_params.sys_desc_table.length != 0) {
738         set_mca_bus(boot_params.sys_desc_table.table[3] &
0x2);
739         machine_id = boot_params.sys_desc_table.table[0];
740         machine_submodel_id =
boot_params.sys_desc_table.table[1];
741         BIOS_revision = boot_params.sys_desc_table.table[2];
742     }
743 #endif
744     saved_video_mode = boot_params.hdr.vid_mode;

```

```

745     bootloader_type = boot_params.hdr.type_of_loader;
746     if ((bootloader_type >> 4) == 0xe) {
747         bootloader_type &= 0xf;
748         bootloader_type |=
749 (boot_params.hdr.ext_loader_type+0x10) << 4;
749     }
750     bootloader_version = bootloader_type & 0xf;
751     bootloader_version |= boot_params.hdr.ext_loader_ver << 4;
752
753 #ifdef CONFIG_BLK_DEV_RAM
754     rd_image_start = boot_params.hdr.ram_size &
755 RAMDISK_IMAGE_START_MASK;
756     rd_prompt = ((boot_params.hdr.ram_size &
757 RAMDISK_PROMPT_FLAG) != 0);
758     rd_doload = ((boot_params.hdr.ram_size &
759 RAMDISK_LOAD_FLAG) != 0);
760 #endif
761 #ifdef CONFIG_EFI
762     if (!strncmp((char
763 *)&boot_params.efi_info.efi_loader_signature,
764 #ifdef CONFIG_X86_32
765         "EL32",
766 #else
767         "EL64",
768 #endif
769     4)) {
770         efi_enabled = 1;
771         efi_reserve_early();
772     }
773 #endif
774
775     x86_init.oem.arch_setup();

```

```

772
773     setup_memory_map();
774     parse_setup_data();
775     /* update the e820_saved too */
776     e820_reserve_setup_data();
777
778     copy_edd();
779
780     if (!boot_params.hdr.root_flags)
781         root_mountflags &= ~MS_RDONLY;
782     init_mm.start_code = (unsigned long) _text;
783     init_mm.end_code = (unsigned long) _etext;
784     init_mm.end_data = (unsigned long) _edata;
785     init_mm.brk = _brk_end;
786
787     code_resource.start = virt_to_phys(_text);
788     code_resource.end = virt_to_phys(_etext)-1;
789     data_resource.start = virt_to_phys(_etext);
790     data_resource.end = virt_to_phys(_edata)-1;
791     bss_resource.start = virt_to_phys(&__bss_start);
792     bss_resource.end = virt_to_phys(&__bss_stop)-1;
793
794 #ifdef CONFIG_CMDLINE_BOOL
795 #ifdef CONFIG_CMDLINE_OVERRIDE
796     strcpy(boot_command_line, builtin_cmdline,
COMMAND_LINE_SIZE);
797 #else
798     if (builtin_cmdline[0]) {
799         /* append boot loader cmdline to builtin */
800         strcat(builtin_cmdline, " ", COMMAND_LINE_SIZE);
801         strcat(builtin_cmdline, boot_command_line,
COMMAND_LINE_SIZE);

```

```
802      strcpy(boot_command_line, builtin_cmdline,
COMMAND_LINE_SIZE);
803  }
804 #endif
805 #endif
806
807      strcpy(command_line, boot_command_line,
COMMAND_LINE_SIZE);
808      *cmdline_p = command_line;
809
810      /*
811       * x86_configure_nx() is called before parse_early_param() to
detect
812       * whether hardware doesn't support NX (so that the early
EHCI debug
813       * console setup can safely call set_fixmap()). It may then be
called
814       * again from within noexec_setup() during parsing early
parameters
815       * to honor the respective command line option.
816       */
817      x86_configure_nx();
818
819      parse_early_param();
820
821      x86_report_nx();
822
823      /* Must be before kernel pagetables are setup */
824      vmi_activate();
825
826      /* after early param, so could get panic from serial */
827      reserve_early_setup_data();
```



```
828
829     if (acpi_mps_check()) {
830 #ifdef CONFIG_X86_LOCAL_APIC
831         disable_apic = 1;
832 #endif
833         setup_clear_cpu_cap(X86_FEATURE_APIC);
834     }
835
836 #ifdef CONFIG_PCI
837     if (pci_early_dump_regs)
838         early_dump_pci_devices();
839 #endif
840
841     finish_e820_parsing();
842
843     if (efi_enabled)
844         efi_init();
845
846     dmi_scan_machine();
847
848     dmi_check_system(bad_bios_dmi_table);
849
850     /*
851      * VMware detection requires dmi to be available, so this
852      * needs to be done after dmi_scan_machine, for the BP.
853      */
854     init_hypervisor_platform();
855
856     x86_init.resources.probe_oms();
857
858     /* after parse_early_param, so could debug it */
```

```
859     insert_resource(&iomem_resource, &code_resource);
860     insert_resource(&iomem_resource, &data_resource);
861     insert_resource(&iomem_resource, &bss_resource);
862
863     trim_bios_range();
864 #ifdef CONFIG_X86_32
865     if (ppro_with_ram_bug()) {
866         e820_update_range(0x70000000ULL, 0x40000ULL,
E820_RAM,
867             E820_RESERVED);
868         sanitize_e820_map(e820.map, ARRAY_SIZE(e820.map),
&e820.nr_map);
869         printk(KERN_INFO "fixed physical RAM map:\n");
870         e820_print_map("bad_ppro");
871     }
872 #else
873     early_gart_iommu_check();
874 #endif
875
876     /*
877      * partially used pages are not usable - thus
878      * we are rounding upwards:
879      */
880     max_pfn = e820_end_of_ram_pfn();
881
882     /* preallocate 4k for mptable mpc */
883     early_reserve_e820_mpc_new();
884     /* update e820 for memory not covered by WB MTRRs */
885     mtrr_bp_init();
886     if (mtrr_trim_uncached_memory(max_pfn))
887         max_pfn = e820_end_of_ram_pfn();
888
```

```
889 #ifdef CONFIG_X86_32
890     /* max_low_pfn get updated here */
891     find_low_pfn_range();
892 #else
893     num_physpages = max_pfn;
894
895     check_x2apic();
896
897     /* How many end-of-memory variables you have, grandma! */
898     /* need this before calling reserve_initrd */
899     if (max_pfn > (1UL<<(32 - PAGE_SHIFT)))
900         max_low_pfn = e820_end_of_low_ram_pfn();
901     else
902         max_low_pfn = max_pfn;
903
904     high_memory = (void *)__va(max_pfn * PAGE_SIZE - 1) + 1;
905     max_pfn_mapped = KERNEL_IMAGE_SIZE >> PAGE_SHIFT;
906 #endif
907
908 #ifdef CONFIG_X86_CHECK_BIOS_CORRUPTION
909     setup_bios_corruption_check();
910 #endif
911
912     printk(KERN_DEBUG "initial memory mapped : 0 - %08lx\n",
913            max_pfn_mapped<<PAGE_SHIFT);
914
915     reserve_brk();
916
917     /*
918      * Find and reserve possible boot-time SMP configuration:
919      */
```

```

920     find_smp_config();
921
922     reserve_ibft_region();
923
924     reserve_trampoline_memory();
925
926 #ifdef CONFIG_ACPI_SLEEP
927     /*
928      * Reserve low memory region for sleep support.
929      * even before init_memory_mapping
930      */
931     acpi_reserve_wakeup_memory();
932 #endif
933     init_gbpages();
934
935     /* max_pfn_mapped is updated here */
936     max_low_pfn_mapped = init_memory_mapping(0,
max_low_pfn<<PAGE_SHIFT);
937     max_pfn_mapped = max_low_pfn_mapped;
938
939 #ifdef CONFIG_X86_64
940     if (max_pfn > max_low_pfn) {
941         max_pfn_mapped = init_memory_mapping(1UL<<32,
942                                             max_pfn<<PAGE_SHIFT);
943         /* can we preseve max_low_pfn ?*/
944         max_low_pfn = max_pfn;
945     }
946 #endif
947
948     /*
949      * NOTE: On x86-32, only from this point on, fixmaps are ready
for use.

```

```
950     */
951
952 #ifdef CONFIG_PROVIDE_OHCI1394_DMA_INIT
953     if (init_ohci1394_dma_early)
954         init_ohci1394_dma_on_all_controllers();
955 #endif
956
957     reserve_initrd();
958
959     reserve_crashkernel();
960
961     vsmp_init();
962
963     io_delay_init();
964
965     /*
966      * Parse the ACPI tables for possible boot-time SMP
configuration.
967      */
968     acpi_boot_table_init();
969
970     early_acpi_boot_init();
971
972 #ifdef CONFIG_ACPI_NUMA
973     /*
974      * Parse SRAT to discover nodes.
975      */
976     acpi = acpi_numa_init();
977 #endif
978
979 #ifdef CONFIG_K8_NUMA
980     if (!acpi)
```

```
981         k8 = !k8_numa_init(0, max_pfn);
982 #endif
983
984     initmem_init(0, max_pfn, acpi, k8);
985 #ifndef CONFIG_NO_BOOTMEM
986     early_res_to_bootmem(0, max_low_pfn<<PAGE_SHIFT);
987 #endif
988
989     dma32_reserve_bootmem();
990
991 #ifdef CONFIG_KVM_CLOCK
992     kvmclock_init();
993 #endif
994
995     x86_init.paging.pagetable_setup_start(swapper_pg_dir);
996     paging_init();
997     x86_init.paging.pagetable_setup_done(swapper_pg_dir);
998
999     tboot_probe();
1000
1001 #ifdef CONFIG_X86_64
1002     map_vsyscall();
1003 #endif
1004
1005     generic_apic_probe();
1006
1007     early_quirks();
1008
1009     /*
1010      * Read APIC and some other early information from ACPI
tables.
1011      */
```

```
1012     acpi_boot_init();
1013
1014     sfi_init();
1015
1016     /*
1017      * get boot-time SMP configuration:
1018      */
1019     if (smp_found_config)
1020         get_smp_config();
1021
1022     prefill_possible_map();
1023
1024 #ifdef CONFIG_X86_64
1025     init_cpu_to_node();
1026 #endif
1027
1028     init_apic_mappings();
1029     ioapic_init_mappings();
1030
1031     /* need to wait for io_apic is mapped */
1032     probe_nr_irqs_gsi();
1033
1034     kvm_guest_init();
1035
1036     e820_reserve_resources();
1037     e820_mark_nosave_regions(max_low_pfn);
1038
1039     x86_init.resources.reserve_resources();
1040
1041     e820_setup_gap();
1042
```

```

1043 #ifdef CONFIG_VT
1044 #if defined(CONFIG_VGA_CONSOLE)
1045     if (!efi_enabled || (efi_mem_type(0xa0000) !=
EFI_CONVENTIONAL_MEMORY))
1046         conswitchp = &vga_con;
1047 #elif defined(CONFIG_DUMMY_CONSOLE)
1048     conswitchp = &dummy_con;
1049 #endif
1050 #endif
1051     x86_init.oem.banner();
1052
1053     mcheck_init();
1054 }

```

4.1.1 early_ioremap_init

所在文件 arch/x86/mm/ioremap.c

4.1.1.1 内核地址空间简介

在 x86-32 为系统的保护模式下，用户空间和内核空间的分配比例是 3:1
其中：

用户空间地址范围：0x00000000—0xbfffffff

内核空间地址范围：0xc0000000—0xffffffff

由于内核的线性地址和物理地址之间有线性的关系，对应的关系为：

内核物理地址=内核线性地址-0xc0000000

所以内核对应的物理地址空间为 0-1G。如果物理地址超过 1G，内核将无法直接访问，为了访问 1G 以上的内存，内核将

0-896M 的物理地址线性的映射到内核线性地址空间，896M 的物理地址成为高端内存，使用动态映射的方法，动态的映射到内核线性地址空间的后 128M。

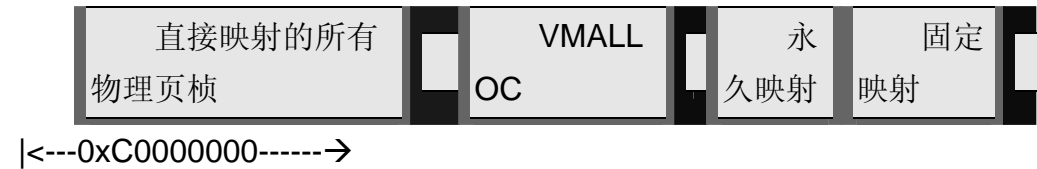
内核地址动态映射分为三种方式：

非连续内存分配：虚拟内存中连续、物理地址中不联系的内存区，使用 vmalloc 区域分配。

永久映射：将高端内存域中的非持久页映射到内核中。

固定映射：与物理地址空间中固定页关联的虚拟地址空间项，具体关联的页帧可以自由选择。可以自由选择。

内核示意图图如下：



4.1.1.2 函数源码

该函数的功能是将高端内存映射到内核后 128M 的线性地址空间内，使得内核可以访问高端内存。

```
369 void __init early_ioremap_init(void)
370 {
371     pmd_t *pmd;
372     int i;
373
374     if (early_ioremap_debug)
375         printk(KERN_INFO "early_ioremap_init()\n");
376     //FIX_BITMAPS_SLOTS=4 NR_FIX_BITMAPS=64
377     for (i = 0; i < FIX_BTMAPS_SLOTS; i++)
378         slot_virt[i] = __fix_to_virt(FIX_BTMAP_BEGIN -
NR_FIX_BTMAPS*i);
```

从后面__fix_to_virt 可以看出，一个固定的插槽索引，对应 4K 的地址空间，

```
379
380     pmd = early_ioremap_pmd(fix_to_virt(FIX_BTMAP_BEGIN));
```

380 行：获得固定索引地址 FIX_BITMAP_BEGIN 对应的虚拟地址的页中间目录项。

```
381     memset(bm_pte, 0, sizeof(bm_pte));
382     pmd_populate_kernel(&init_mm, pmd, bm_pte);
```

将页表 `bm_pte` 的地址赋值给 `pmd`。

```
383
384     /*
385      * The boot-ioremap range spans multiple pmds, for which
386      * we are not prepared:
387      */
388 #define __FIXADDR_TOP (-PAGE_SIZE)
389     BUILD_BUG_ON((__fix_to_virt(FIX_BTMAP_BEGIN) >>
PMD_SHIFT)
390                  != (__fix_to_virt(FIX_BTMAP_END) >>
PMD_SHIFT));
391 #undef __FIXADDR_TOP
392     if (pmd != early_ioremap_pmd(fix_to_virt(FIX_BTMAP_END))) {
393         WARN_ON(1);
394         printk(KERN_WARNING "pmd %p != %p\n",
395                pmd,
early_ioremap_pmd(fix_to_virt(FIX_BTMAP_END)));
396         printk(KERN_WARNING "fix_to_virt(FIX_BTMAP_BEGIN):
%08lx\n",
397                fix_to_virt(FIX_BTMAP_BEGIN));
398         printk(KERN_WARNING "fix_to_virt(FIX_BTMAP_END):
%08lx\n",
399                fix_to_virt(FIX_BTMAP_END));
400
401         printk(KERN_WARNING "FIX_BTMAP_END:          %d\n",
FIX_BTMAP_END);
402         printk(KERN_WARNING "FIX_BTMAP_BEGIN:
%d\n",
403                FIX_BTMAP_BEGIN);
404     }
405 }
```

该函数涉及到的宏和变量定义如下：

```
367 static unsigned long slot_virt[FIX_BTMAPS_SLOTS] __initdata;
```

```
55  * Here we define all the compile-time 'special' virtual
56  * addresses. The point is to have a constant address at
57  * compile time, but to set the physical address only
58  * in the boot process.
59  * for x86_32: We allocate these special addresses
60  * from the end of virtual memory (0xffff000) backwards.
61  * Also this lets us do fail-safe vmalloc(), we
62  * can guarantee that these special addresses and
63  * vmalloc()-ed addresses never overlap.
64  *
65  * These 'compile-time allocated' memory buffers are
66  * fixed-size 4k pages (or larger if used with an increment
67  * higher than 1). Use set_fixmap(idx,phys) to associate
68  * physical memory with fixmap indices.
69  *
70  * TLB entries of such buffers will not be flushed across
71  * task switches.
72  */
73 enum fixed_addresses {
74 #ifdef CONFIG_X86_32
75     FIX_HOLE,
76     FIX_VDSO,
77 #else
78     VSYSCALL_LAST_PAGE,
79     VSYSCALL_FIRST_PAGE = VSYSCALL_LAST_PAGE
80                             + ((VSYSCALL_END-VSYSCALL_START) >>
PAGE_SHIFT) - 1,
81     VSYSCALL_HPET,
```

```

82 #endif
83     FIX_DBGP_BASE,
84     FIX_EARLYCON_MEM_BASE,
85 #ifdef CONFIG_PROVIDE_OHCI1394_DMA_INIT
86     FIX_OHCI1394_BASE,
87 #endif
88 #ifdef CONFIG_X86_LOCAL_APIC
89     FIX_APIC_BASE, /* local (CPU) APIC) -- required for SMP or
not */
90 #endif
91 #ifdef CONFIG_X86_IO_APIC
92     FIX_IO_APIC_BASE_0,
93     FIX_IO_APIC_BASE_END = FIX_IO_APIC_BASE_0 +
MAX_IO_APICS - 1,
94 #endif
95 #ifdef CONFIG_X86_VISWS_APIC
96     FIX_CO_CPU, /* Cobalt timer */
97     FIX_CO_APIC, /* Cobalt APIC Redirection Table */
98     FIX_LI_PCIA, /* Lithium PCI Bridge A */
99     FIX_LI_PCIB, /* Lithium PCI Bridge B */
100 #endif
101 #ifdef CONFIG_X86_F00F_BUG
102     FIX_F00F_IDT, /* Virtual mapping for IDT */
103 #endif
104 #ifdef CONFIG_X86_CYCLONE_TIMER
105     FIX_CYCLONE_TIMER, /*cyclone timer register*/
106 #endif
107 #ifdef CONFIG_X86_32
108     FIX\_KMAP\_BEGIN, /* reserved pte's for temporary kernel
mappings */
109     FIX_KMAP_END =
FIX_KMAP_BEGIN+(KM_TYPE_NR*NR_CPUS)-1,

```

```

110 #ifdef CONFIG_PCI_MMCONFIG
111     FIX_PCIE_MCFG,
112 #endif
113 #endif
114 #ifdef CONFIG_PARAVIRT
115     FIX_PARAVIRT_BOOTMAP,
116 #endif
117     FIX_TEXT_POKE1, /* reserve 2 pages for text_poke() */
118     FIX_TEXT_POKE0, /* first page is last, because allocation is
backward */
119     __end_of_permanent_fixed_addresses,
120     /*
121      * 256 temporary boot-time mappings, used by early_ioremap(),
122      * before ioremap() is functional.
123      *
124      * If necessary we round it up to the next 256 pages boundary so
125      * that we can have a single pgd entry and a single pte table:
126      */
127 #define NR_FIX_BTMAPS      64
128 #define FIX_BTMAPS_SLOTS   4
129 #define TOTAL_FIX_BTMAPS   (NR_FIX_BTMAPS *
FIX_BTMAPS_SLOTS)
130     FIX_BTMAP_END =
131         (__end_of_permanent_fixed_addresses ^
132         (__end_of_permanent_fixed_addresses +
TOTAL_FIX_BTMAPS - 1)) &
133         -PTRS_PER_PTE
134         ? __end_of_permanent_fixed_addresses +
TOTAL_FIX_BTMAPS -
135         (__end_of_permanent_fixed_addresses &
(TOTAL_FIX_BTMAPS - 1))
136         : __end_of_permanent_fixed_addresses,

```

```

137     FIX_BTMAP_BEGIN = FIX_BTMAP_END +
TOTAL_FIX_BTMAPS - 1,
138 #ifdef CONFIG_X86_32
139     FIX_WP_TEST,
140 #endif
141 #ifdef CONFIG_INTEL_TXT
142     FIX_TBOOT_BASE,
143 #endif
144     __end_of_fixed_addresses
145 };

```

可以看出 `fixed_addresses` 是一个枚举类型的变量，在其中定义的

`FIX_BTMAP_BEGIN`

`FIX_BTMAP_END`

`FIX_KMAP_BEGIN`

`FIX_KMAP_END`

`__end_of_permanent_fixed_addresses`

`__end_of_fixed_addresses`

都是枚举变量中的成员。

函数 `__fix_to_virt` 的定义如下

```

185 #define __fix_to_virt(x)    (FIXADDR_TOP - ((x) << PAGE_SHIFT))
186 #define __virt_to_fix(x)    ((FIXADDR_TOP - ((x)&PAGE_MASK))
>> PAGE_SHIFT)

```

其中 `FIXADDR_TOP` 定义如下

```

40 extern unsigned long __FIXADDR_TOP;
41 #define FIXADDR_TOP ((unsigned long)__FIXADDR_TOP)

```

其中 `__FIXADDR_TOP` 的定义在 `arch/x86/mm/pgtable32.c` 中如下：

```

99 unsigned long __FIXADDR_TOP = 0xfffff000;
100 EXPORT_SYMBOL(__FIXADDR_TOP);

```

__fix_to_virt 的含义为：从地址 0xffffffff(4G 后退一个页面的地址)回退固定地址(x)个页面后的虚拟地址。

4.1.2 setup_memory_map

进行内存管理，将前面通过 e820 调用获得的内存信息拷贝到 e820 变量中。该函数调用 x86_init.resources.memory_setup，获得内存信息，将 e820 中的信息赋值一份保存在 e820_saved 中，然后打印内存布局图。

```
1204 void __init setup_memory_map(void)
1205 {
1206     char *who;
1207
1208     who = x86_init.resources.memory_setup();
1209     memcpy(&e820_saved, &e820, sizeof(struct e820map));
1210     printk(KERN_INFO "BIOS-provided physical RAM map:\n");
1211     e820_print_map(who);
1212 }
```

X86_init 的定义在 arch/x86/kernel/x86_init.c 中：

```
32 struct x86_init_ops x86_init __initdata = {
33
34     .resources = {
35         .probe_roms      = x86_init_noop,
36         .reserve_resources = reserve_standard_io_resources,
37         .memory_setup     =
default_machine_specific_memory_setup,
38     },
39
40     .mpparse = {
41         .mpc_record      = x86_init_uint_noop,
42         .setup_ioapic_ids = x86_init_noop,
```

```
43     .mpc_apic_id      = default_mpc_apic_id,
44     .smp_read_mpc_oem  = default_smp_read_mpc_oem,
45     .mpc_oem_bus_info  = default_mpc_oem_bus_info,
46     .find_smp_config   = default_find_smp_config,
47     .get_smp_config     = default_get_smp_config,
48 },
49
50 .irqs = {
51     .pre_vector_init    = init_ISA_irqs,
52     .intr_init          = native_init_IRQ,
53     .trap_init          = x86_init_noop,
54 },
55
56 .oem = {
57     .arch_setup         = x86_init_noop,
58     .banner             = default_banner,
59 },
60
61 .paging = {
62     .pagetable_setup_start = native_pagetable_setup_start,
63     .pagetable_setup_done  = native_pagetable_setup_done,
64 },
65
66 .timers = {
67     .setup_percpu_clockev = setup_boot_APIC_clock,
68     .tsc_pre_init         = x86_init_noop,
69     .timer_init           = hpet_time_init,
70 },
71
72 .iommu = {
73     .iommu_init          = iommu_init_noop,
```



```

74     },
75
76     .pci = {
77         .init          = x86_default_pci_init,
78         .init_irq      = x86_default_pci_init_irq,
79         .fixup_irqs    = x86_default_pci_fixup_irqs,
80     },
81 };

```

对 `x86_init.resources.memory_setup()` 的调用将会调用函数，`default_machine_specific_memory_setup`，该函数定义在 `arch/x86/kernel/e820.c` 中：

该函数的将启动参数中内存分段的信息拷贝到全局变量 `e820` 中，过程比较简单，就不进一步分析了。

```

1166 char *__init default_machine_specific_memory_setup(void)
1167 {
1168     char *who = "BIOS-e820";
1169     u32 new_nr;
1170     /*
1171      * Try to copy the BIOS-supplied E820-map.
1172      *
1173      * Otherwise fake a memory map; one section from 0k->640k,
1174      * the next section from 1mb->appropriate_mem_k
1175      */
1176     new_nr = boot_params.e820_entries;
1177     sanitize_e820_map(boot_params.e820_map,
1178                     ARRAY_SIZE(boot_params.e820_map),
1179                     &new_nr);
1180     boot_params.e820_entries = new_nr;
1181     if (append_e820_map(boot_params.e820_map,
boot_params.e820_entries)
1182         < 0) {

```

```

1183         u64 mem_size;
1184
1185         /* compare results from other methods and take the
greater */
1186         if (boot_params.alt_mem_k
1187             < boot_params.screen_info.ext_mem_k) {
1188             mem_size = boot_params.screen_info.ext_mem_k;
1189             who = "BIOS-88";
1190         } else {
1191             mem_size = boot_params.alt_mem_k;
1192             who = "BIOS-e801";
1193         }
1194
1195         e820.nr_map = 0;
1196         e820_add_region(0, LOWMEMSIZE(), E820_RAM);
1197         e820_add_region(HIGH_MEMORY, mem_size << 10,
E820_RAM);
1198     }
1199
1200     /* In case someone cares... */
1201     return who;
1202 }

```

4.1.3 初始化 0 号进程变量

782 行到 792 行，初始化 0 号进程 `init_mm` 中的相关变量。

4.1.4 `x86_configure_nx()` 设置页面不可执行标志

函数定义在 `arch/x86/mm/setup_nx.c` 中：

```

32 void __cpuinit x86_configure_nx(void)
33 {

```

```
34     if (cpu_has_nx && !disable_nx)
35         __supported_pte_mask |= _PAGE_NX;
36     else
37         __supported_pte_mask &= ~_PAGE_NX;
38 }
```

由于之前没有区分代码段和堆、堆栈之间的区别，导致堆溢出和缓冲去攻击横行，后来 **cpu** 的硬件厂商支持 **nx** 技术，对非代码区页面设置 **nx** 标志，将不容许 **cpu** 在该页面上执行代码。

上述代码 34 行检测 **cpu** 是否支持 **nx** 技术，以及是否禁用 **nx** 技术，如果 **cpu** 支持且没有禁用 **nx**，则在 `__supported_pte_mask` 中职位 `_PAGE_NX` 位。

4.1.5 `e820_end_of_ram_pfn` 获得物理页面中页面数（可能包含空洞）

该函数所在文件 `arch/x86/kernel/e820.c` 中：

```
889 unsigned long __init e820_end_of_ram_pfn(void)
890 {
891     return e820_end_pfn(MAX_ARCH_PFN, E820_RAM);
892 }
```

891 行调用 `e820_end_pfn`。该函数检测 `e820` 变量中内存区段，将区段中所在内存的最大的内存页面号返回。过程比较简单，不进一步分析。

```
855 static unsigned long __init e820_end_pfn(unsigned long limit_pfn,
unsigned type)
856 {
857     int i;
858     unsigned long last_pfn = 0;
859     unsigned long max_arch_pfn = MAX_ARCH_PFN;
860
861     for (i = 0; i < e820.nr_map; i++) {
862         struct e820entry *ei = &e820.map[i];
```

```
863         unsigned long start_pfn;
864         unsigned long end_pfn;
865
866         if (ei->type != type)
867             continue;
868
869         start_pfn = ei->addr >> PAGE_SHIFT;
870         end_pfn = (ei->addr + ei->size) >> PAGE_SHIFT;
871
872         if (start_pfn >= limit_pfn)
873             continue;
874         if (end_pfn > limit_pfn) {
875             last_pfn = limit_pfn;
876             break;
877         }
878         if (end_pfn > last_pfn)
879             last_pfn = end_pfn;
880     }
881
882     if (last_pfn > max_arch_pfn)
883         last_pfn = max_arch_pfn;
884
885     printk(KERN_INFO "last_pfn = %#lx max_arch_pfn = %#lx\n",
886            last_pfn, max_arch_pfn);
887     return last_pfn;
888 }
```

4.1.6 find_low_pfn_range

物理内存被分为 DMA,NORMAL,HIGH 三个区。该函数获得 normal 内存区域的最大页帧号。

4.1.7 init_memory_mapping建立永久内核页表

在前面的分析过程中，已经看到初始化过程中为加载进内存的内核建立了页表，现在要为所有的物理内存建立页表。

函数所在文件：arch/x86/mm/init.c

```
120  * Setup the direct mapping of the physical memory at
PAGE_OFFSET.
121  * This runs before bootmem is initialized and gets pages directly
from
122  * the physical memory. To access them they are temporarily
mapped.
123  */
124 unsigned long __init_refok init_memory_mapping(unsigned long
start,
125                                         unsigned long end)
126 {
127     unsigned long page_size_mask = 0;
128     unsigned long start_pfn, end_pfn;
129     unsigned long ret = 0;
130     unsigned long pos;
131
132     struct map_range mr[NR_RANGE_MR];
133     int nr_range, i;
134     int use_pse, use_gbpages;
135
136     printk(KERN_INFO "init_memory_mapping: %016lx-%016lx\n",
start, end);
137
138     #if defined(CONFIG_DEBUG_PAGEALLOC) ||
defined(CONFIG_KMEMCHECK)
139         /*
140          * For CONFIG_DEBUG_PAGEALLOC, identity mapping will
use small pages.
```

```
141      * This will simplify cpa(), which otherwise needs to support
splitting
142      * large pages into small in interrupt context, etc.
143      */
144      use_pse = use_gbpages = 0;
145 #else
146      use_pse = cpu_has_pse;
147      use_gbpages = direct_gbpages;
148 #endif
149
150      /* Enable PSE if available */
151      if (cpu_has_pse)
152          set_in_cr4(X86_CR4_PSE);
153
154      /* Enable PGE if available */
155      if (cpu_has_pge) {
156          set_in_cr4(X86_CR4_PGE);
157          __supported_pte_mask |= _PAGE_GLOBAL;
158      }
159
160      if (use_gbpages)
161          page_size_mask |= 1 << PG_LEVEL_1G;
162      if (use_pse)
163          page_size_mask |= 1 << PG_LEVEL_2M;
164
165      memset(mr, 0, sizeof(mr));
166      nr_range = 0;
167
168      /* head if not big page alignment ? */
169      start_pfn = start >> PAGE_SHIFT;
170      pos = start_pfn << PAGE_SHIFT;
171 #ifdef CONFIG_X86_32
```

```

172     /*
173      * Don't use a large page for the first 2/4MB of memory
174      * because there are often fixed size MTRRs in there
175      * and overlapping MTRRs into large pages can cause
176      * slowdowns.
177      */
178     if (pos == 0)
179         end_pfn = 1<<(PMD_SHIFT - PAGE_SHIFT);
180     else
181         end_pfn = ((pos + (PMD_SIZE - 1))>>PMD_SHIFT)
182                 << (PMD_SHIFT - PAGE_SHIFT);
183 #else /* CONFIG_X86_64 */
184     end_pfn = ((pos + (PMD_SIZE - 1)) >> PMD_SHIFT)
185             << (PMD_SHIFT - PAGE_SHIFT);
186 #endif
187     if (end_pfn > (end >> PAGE_SHIFT))
188         end_pfn = end >> PAGE_SHIFT;
189     if (start_pfn < end_pfn) {
190         nr_range = save_mr(mr, nr_range, start_pfn, end_pfn, 0);
191         pos = end_pfn << PAGE_SHIFT;
192     }
193
194     /* big page (2M) range */
195     start_pfn = ((pos + (PMD_SIZE - 1))>>PMD_SHIFT)
196             << (PMD_SHIFT - PAGE_SHIFT);
197 #ifdef CONFIG_X86_32
198     end_pfn = (end>>PMD_SHIFT) << (PMD_SHIFT -
PAGE_SHIFT);
199 #else /* CONFIG_X86_64 */
200     end_pfn = ((pos + (PUD_SIZE - 1))>>PUD_SHIFT)
201             << (PUD_SHIFT - PAGE_SHIFT);

```

```

202     if (end_pfn > ((end >> PMD_SHIFT) << (PMD_SHIFT -
PAGE_SHIFT)))
203         end_pfn = ((end >> PMD_SHIFT) << (PMD_SHIFT -
PAGE_SHIFT));
204 #endif
205
206     if (start_pfn < end_pfn) {
207         nr_range = save_mr(mr, nr_range, start_pfn, end_pfn,
208             page_size_mask & (1 << PG_LEVEL_2M));
209         pos = end_pfn << PAGE_SHIFT;
210     }
211
212 #ifdef CONFIG_X86_64
213     /* big page (1G) range */
214     start_pfn = ((pos + (PUD_SIZE - 1)) >> PUD_SHIFT)
215         << (PUD_SHIFT - PAGE_SHIFT);
216     end_pfn = (end >> PUD_SHIFT) << (PUD_SHIFT -
PAGE_SHIFT);
217     if (start_pfn < end_pfn) {
218         nr_range = save_mr(mr, nr_range, start_pfn, end_pfn,
219             page_size_mask &
220             ((1 << PG_LEVEL_2M) | (1 << PG_LEVEL_1G)));
221         pos = end_pfn << PAGE_SHIFT;
222     }
223
224     /* tail is not big page (1G) alignment */
225     start_pfn = ((pos + (PMD_SIZE - 1)) >> PMD_SHIFT)
226         << (PMD_SHIFT - PAGE_SHIFT);
227     end_pfn = (end >> PMD_SHIFT) << (PMD_SHIFT -
PAGE_SHIFT);
228     if (start_pfn < end_pfn) {
229         nr_range = save_mr(mr, nr_range, start_pfn, end_pfn,

```



```

230             page_size_mask & (1<<PG_LEVEL_2M));
231     pos = end_pfn << PAGE_SHIFT;
232     }
233 #endif
234
235     /* tail is not big page (2M) alignment */
236     start_pfn = pos>>PAGE_SHIFT;
237     end_pfn = end>>PAGE_SHIFT;
238     nr_range = save_mr(mr, nr_range, start_pfn, end_pfn, 0);
239
240     /* try to merge same page size and continuous */
241     for (i = 0; nr_range > 1 && i < nr_range - 1; i++) {
242         unsigned long old_start;
243         if (mr[i].end != mr[i+1].start ||
244             mr[i].page_size_mask != mr[i+1].page_size_mask)
245             continue;
246         /* move it */
247         old_start = mr[i].start;
248         memmove(&mr[i], &mr[i+1],
249             (nr_range - 1 - i) * sizeof(struct map_range));
250         mr[i--].start = old_start;
251         nr_range--;
252     }
253
254     for (i = 0; i < nr_range; i++)
255         printk(KERN_DEBUG " %010lx - %010lx page %s\n",
256             mr[i].start, mr[i].end,
257             (mr[i].page_size_mask & (1<<PG_LEVEL_1G))?"1G":(
258             (mr[i].page_size_mask &
(1<<PG_LEVEL_2M))?"2M":"4k"));
259
260     /*

```

```

261      * Find space for the kernel direct mapping tables.
262      *
263      * Later we should allocate these tables in the local node of the
264      * memory mapped. Unfortunately this is done currently before
the
265      * nodes are discovered.
266      */

```

169-266 行是根据物理页面的不同大小（4K,2M），将物理内存分为不同的区，并将分区结果保存在变量 `mr` 中。

```

267      if (!after_bootmem)
268          find_early_table_space(end, use_pse, use_gbpages);
269
270      for (i = 0; i < nr_range; i++)
271          ret = kernel_physical_mapping_init(mr[i].start, mr[i].end,
272                                              mr[i].page_size_mask);
273
274 #ifdef CONFIG_X86_32
275      early_ioremap_page_table_range_init();
276
277      load_cr3(swapper_pg_dir);
278 #endif
279
280 #ifdef CONFIG_X86_64
280 #ifdef CONFIG_X86_64
281      if (!after_bootmem && !start) {
282          pud_t *pud;
283          pmd_t *pmd;
284
285          mmu_cr4_features = read_cr4();
286
287          /*

```

```

288      * _brk_end cannot change anymore, but it and _end may
be
289      * located on different 2M pages. cleanup_highmap(),
however,
290      * can only consider _end when it runs, so destroy any
291      * mappings beyond _brk_end here.
292      */
293      pud = pud_offset(pgd_offset_k(_brk_end), _brk_end);
294      pmd = pmd_offset(pud, _brk_end - 1);
295      while (++pmd <= pmd_offset(pud, (unsigned long)_end - 1))
296          pmd_clear(pmd);
297  }
298 #endif
299  __flush_tlb_all();
300
301  if (!after_bootmem && e820_table_end > e820_table_start)
302      reserve_early(e820_table_start << PAGE_SHIFT,
303                    e820_table_end << PAGE_SHIFT, "PGTABLE");
304
305  if (!after_bootmem)
306      early_memtest(start, end);
307
308  return ret >> PAGE_SHIFT;
309 }

```

该函数的调用过程是：

find_early_table_space: 查找可疑容纳所有物理内存页表项的连续内存区域。

kernel_physical_mapping_init: 为不同的分区建立物理内存映射。

early_ioremap_page_table_range_init()

load_cr3

4.1.7.1 find_early_table_space

该函数实现的功能是在 e820 内存区中，查找一块可以容纳所有物理内存页表项的连续物理内存块。实现过程比较简单，重点查看 find_e820_area 函数。

```
32 static void __init find_early_table_space(unsigned long end, int
use_pse,
33             int use_gbpages)
34 {
35     unsigned long puds, pmds, ptes, tables, start;
36
37     puds = (end + PUD_SIZE - 1) >> PUD_SHIFT;
38     tables = roundup(puds * sizeof(pud_t), PAGE_SIZE);
39
40     if (use_gbpages) {
41         unsigned long extra;
42
43         extra = end - ((end >> PUD_SHIFT) << PUD_SHIFT);
44         pmds = (extra + PMD_SIZE - 1) >> PMD_SHIFT;
45     } else
46         pmds = (end + PMD_SIZE - 1) >> PMD_SHIFT;
47
48     tables += roundup(pmds * sizeof(pmd_t), PAGE_SIZE);
49
50     if (use_pse) {
51         unsigned long extra;
52
53         extra = end - ((end >> PMD_SHIFT) << PMD_SHIFT);
54 #ifdef CONFIG_X86_32
55         extra += PMD_SIZE;
56 #endif
57         ptes = (extra + PAGE_SIZE - 1) >> PAGE_SHIFT;
58     } else
```

```

59         ptes = (end + PAGE_SIZE - 1) >> PAGE_SHIFT;
60
61         tables += roundup(ptes * sizeof(pte_t), PAGE_SIZE);
62
63 #ifdef CONFIG_X86_32
64         /* for fixmap */
65         tables += roundup(__end_of_fixed_addresses * sizeof(pte_t),
PAGE_SIZE);
66 #endif
67
68         /*
69          * RED-PEN putting page tables only on node 0 could
70          * cause a hotspot and fill up ZONE_DMA. The page tables
71          * need roughly 0.5KB per GB.
72          */
73 #ifdef CONFIG_X86_32
74         start = 0x7000;
75 #else
76         start = 0x8000;
77 #endif
78         e820_table_start = find_e820_area(start,
max_pfn_mapped<<PAGE_SHIFT,
79                                     tables, PAGE_SIZE);
80         if (e820_table_start == -1UL)
81             panic("Cannot find space for the kernel page tables");
82
83         e820_table_start >>= PAGE_SHIFT;
84         e820_table_end = e820_table_start;
85         e820_table_top = e820_table_start + (tables >> PAGE_SHIFT);

```

83-85 行：将获得的内存地址的起始页面号保存到 `e820_table_end` 中，终止页面号保存到 `e820_table_top` 中，以后分配页表时，`e820_table_end`，`e820_table_top` 将起到重要的作用。

```

86
87     printk(KERN_DEBUG "kernel direct mapping tables up to %lx @
88         end, e820_table_start << PAGE_SHIFT, e820_table_top <<
PAGE_SHIFT);
89 }

```

find_e820_area 所在文件 arch/x86/kernel/e820.c

```

743 u64 __init find_e820_area(u64 start, u64 end, u64 size, u64 align)
744 {
745     int i;
746
747     for (i = 0; i < e820.nr_map; i++) {
748         struct e820entry *ei = &e820.map[i];
749         u64 addr;
750         u64 ei_start, ei_last;
751         //必须是 RAM 区域
752         if (ei->type != E820_RAM)
753             continue;
754
755         ei_last = ei->addr + ei->size;
756         ei_start = ei->addr;
757         addr = find_early_area(ei_start, ei_last, start, end,
758             size, align);
759
760         if (addr != -1ULL)
761             return addr;
762     }
763     return -1ULL;
764 }

```

find_e820_area 转而调用在 kernel/early_res.c 中的函数 find_early_area:

该函数的六个参数如下:

ei_start:e820 内存区起始区域地址

ei_last:e820 内存区域终止地址

start:可供分配内存的起始地址

end:可供分配内存的终止地址

size:期望分配内存的大小

align:对齐方式

```
534 /*
535  * Find a free area with specified alignment in a specific range.
536  * only with the area.between start to end is active range from
early_node_map
537  * so they are good as RAM
538  */
539 u64 __init find_early_area(u64 ei_start, u64 ei_last, u64 start, u64
end,
540                          u64 size, u64 align)
541 {
542     u64 addr, last;
543
544     addr = round_up(ei_start, align);
545     if (addr < start)
546         addr = round_up(start, align);
```

addr 必须属于[start,end]

```
547     if (addr >= ei_last)
548         goto out;
```

addr 必须不超过 ei_last

```
549     while (bad_addr(&addr, size, align) && addr+size <= ei_last)
550         ;
551     last = addr + size;
552     if (last > ei_last)
```

553	goto out;
last 必须属于[ei_start,ei_last]	
554	if (last > end)
555	goto out;
556	
557	return addr;
558	
559	out:
560	return -1ULL;
561	}

上述函数的基本功能是：

在[ei_start,ei_last]区域分配一块大小为 **size** 的内存。如果 $ei_last - ei_start < size$ ，分配失败。在大于的情况下，还要保证分配内存的区间在 [start,end] 之间。

在上述函数的 549 行 `bad_addr(&addr, size, align)` 函数我们还没有介绍，该函数的作用是检测 [addr,addr+size] 是否落与某个已经分配出去的内存区域存在交集。如果是则继续查找合适的区域。

4.1.7.1.1 冲突检测

与冲突检测的相关数据结构

11	/*
12	* Early reserved memory areas.
13	*/
14	/*
15	* need to make sure this one is bigger enough before
16	* find_fw_memmap_area could be used
17	*/
18	#define MAX_EARLY_RES_X 32
19	
20	struct early_res {
21	u64 start, end;
22	char name[15];


```

23     char overlap_ok;
24 };
25 static struct early_res early_res_x[MAX_EARLY_RES_X] __initdata;
26
27 static int max_early_res __initdata = MAX_EARLY_RES_X;
28 static struct early_res *early_res __initdata = &early_res_x[0];
29 static int early_res_count __initdata;

```

内核定义了全局变量 `early_res_x`，该变量是一个包含了 32 个 `early_res` 结构体元素的变量，用来保存已经分配出去的内存区域。

`early_res` 变量指向该数组的起始地址。

`max_early_res` 表示数组的大小。

下面看看 `early_res_x` 中的元素是如何被设置的。

当调用 `__alloc_memory_core_early` 函数分配内存后，会调用 `reserve_early_checkout_check` 函数将已分配的内存保存到 `early_res_x` 中。

```

298 void __init reserve_early_without_check(u64 start, u64 end, char
*name)
299 {
300     struct early_res *r;
301
302     if (start >= end)
303         return;
304
305     __check_and_double_early_res(start, end);

```

如果 `early_res` 内存空间不足，则扩充内存。

```

306
307     r = &early_res[early_res_count];
308
309     r->start = start;
310     r->end = end;
311     r->overlap_ok = 0;

```

309-311 设置已分配内存区。

```
312     if (name)
313         strncpy(r->name, name, sizeof(r->name) - 1);
314     early_res_count++;
315 }
```

4.1.7.2 kernel_physical_mapping_init物理内存页表项设置

函数所在文件 arch/x86/mm/init_32.c

函数参数含义：

start:起始虚拟地址

end: 结束虚拟地址

page_size_mask:页表大小标志

```
238 unsigned long __init
239 kernel_physical_mapping_init(unsigned long start,
240                             unsigned long end,
241                             unsigned long page_size_mask)
242 {
243     int use_pse = page_size_mask == (1<<PG_LEVEL_2M);
244     unsigned long last_map_addr = end;
245     unsigned long start_pfn, end_pfn;
246     pgd_t *pgd_base = swapper_pg_dir;
```

246 行: pgd_base(页目录起始地址)页目录第一项

```
247     int pgd_idx, pmd_idx, pte_ofs;
248     unsigned long pfn;
249     pgd_t *pgd;
250     pmd_t *pmd;
251     pte_t *pte;
252     unsigned pages_2m, pages_4k;
253     int mapping_iter;
```

254	
255	start_pfn = start >> PAGE_SHIFT;
256	end_pfn = end >> PAGE_SHIFT;

255-256 行：将虚拟地址装换为页面号

257	
258	/*
259	* First iteration will setup identity mapping using large/small pages
260	* based on use_pse, with other attributes same as set by
261	* the early code in head_32.S
262	*
263	* Second iteration will setup the appropriate attributes (NX, GLOBAL..)
264	* as desired for the kernel identity mapping.
265	*
266	* This two pass mechanism conforms to the TLB app note which says:
267	*
268	* "Software should not write to a paging-structure entry in a way
269	* that would change, for any linear address, both the page size
270	* and either the page frame or attributes."
271	*/
272	mapping_iter = 1;
273	
274	if (!cpu_has_pse)
275	use_pse = 0;

检测 cpu 是否设置物理页面扩充位。

276	
277	repeat:
278	pages_2m = pages_4k = 0;

```

279     pfn = start_pfn;
280     pgd_idx = pgd_index((pfn<<PAGE_SHIFT) +
PAGE_OFFSET);
281     pgd = pgd_base + pgd_idx;

```

pgd: 保存虚拟地址对应的页目录项

```

282     for (; pgd_idx < PTRS_PER_PGD; pgd++, pgd_idx++) {
283         pmd = one_md_table_init(pgd);

```

在二级页表映射中，pmd=pgd

```

284
285         if (pfn >= end_pfn)
286             continue;
287 #ifdef CONFIG_X86_PAE
288         pmd_idx = pmd_index((pfn<<PAGE_SHIFT) +
PAGE_OFFSET);
289         pmd += pmd_idx;
290 #else
291         pmd_idx = 0;
292 #endif
293         for (; pmd_idx < PTRS_PER_PMD && pfn < end_pfn;
294             pmd++, pmd_idx++) { //在二级页表映射中，for 循环只
执行一次
295             unsigned int addr = pfn * PAGE_SIZE +
PAGE_OFFSET;
296
297             /*
298              * Map with big pages if possible, otherwise
299              * create normal page tables:
300              */
301             if (use_pse) {
302                 unsigned int addr2;
303                 pgprot_t prot = PAGE_KERNEL_LARGE;
304                 /*

```

```

305          * first pass will use the same initial
306          * identity mapping attribute + _PAGE_PSE.
307          */
308          pgprot_t init_prot =
309              __pgprot(PTE_IDENT_ATTR |
310                      _PAGE_PSE);
311
312          addr2 = (pfn + PTRS_PER_PTE-1) * PAGE_SIZE
+
313                  PAGE_OFFSET + PAGE_SIZE-1;
314
315          if (is_kernel_text(addr) ||
316              is_kernel_text(addr2))
317              prot = PAGE_KERNEL_LARGE_EXEC;
318
319          pages_2m++;
320          if (mapping_iter == 1)
321              set_pmd(pmd, pfn_pmd(pfn, init_prot));
322          else
323              set_pmd(pmd, pfn_pmd(pfn, prot));
324
325          pfn += PTRS_PER_PTE;
326          continue;
327      }
328      pte = one_page_table_init(pmd);

```

328 行：为目录项分配页表

```

329
330          pte_ofs = pte_index((pfn<<PAGE_SHIFT) +
PAGE_OFFSET);
331          pte += pte_ofs;

```

330 行：求出虚拟地址在页表 `pte` 中的偏移

331 行: pte 中保存对应的页表项

332	for (; pte_ofs < PTRS_PER_PTE && pfn < end_pfn;
333	pte++, pfn++, pte_ofs++, addr += PAGE_SIZE) {

332-333 行: for 循环用来设置页表项, 指向相应的物理页面的地址。

334	pgprot_t prot = PAGE_KERNEL;
335	/*
336	* first pass will use the same initial
337	* identity mapping attribute.
338	*/
339	pgprot_t init_prot =
340	__pgprot(PTE_IDENT_ATTR);
341	if (is_kernel_text(addr))
342	prot = PAGE_KERNEL_EXEC;
343	
344	pages_4k++;
345	if (mapping_iter == 1) {
346	set_pte(pte, pfn_pte(pfn, init_prot));
347	last_map_addr = (pfn << PAGE_SHIFT) +
348	PAGE_SIZE;
349	} else
350	set_pte(pte, pfn_pte(pfn, prot));

设置页表项。

350	}
351	}
352	}
353	if (mapping_iter == 1) {
354	/*
355	* update direct mapping page count only in the first
356	* iteration.
357	*/

```

358     update_page_count(PG_LEVEL_2M, pages_2m);
359     update_page_count(PG_LEVEL_4K, pages_4k);
360
361     /*
362      * local global flush tlb, which will flush the previous
363      * mappings present in both small and large page TLB's.
364      */
365     __flush_tlb_all();
366
367     /*
368      * Second iteration will set the actual desired PTE
attributes.
369      */
370     mapping_iter = 2;
371     goto repeat;
372 }
373 return last_map_addr;
374 }

```

4.1.7.2.1 页表分配函数one_page_table_init

```

105  * Create a page table and place a pointer to it in a middle page
106  * directory entry:
107  */
108 static pte_t * __init one_page_table_init(pmd_t *pmd)
109 {
110     if (!(pmd_val(*pmd) & _PAGE_PRESENT)) {
111         pte_t *page_table = NULL;
112
113         if (after_bootmem) {

```

113 行：现在还未设置 after_bootmem

```

114 #if defined(CONFIG_DEBUG_PAGEALLOC) ||
defined(CONFIG_KMEMCHECK)
115         page_table = (pte_t *)
alloc_bootmem_pages(PAGE_SIZE);
116 #endif
117         if (!page_table)
118             page_table =
119             (pte_t *)alloc_bootmem_pages(PAGE_SIZE);
120     } else
121         page_table = (pte_t *)alloc_low_page();
122
123     paravirt_alloc_pte(&init_mm, __pa(page_table) >>
PAGE_SHIFT);
124     set_pmd(pmd, __pmd(__pa(page_table) |
_PAGE_TABLE));

```

124 行: pmd 保存刚分配页表的物理地址

```

125         BUG_ON(page_table != pte_offset_kernel(pmd, 0));
126     }
127
128     return pte_offset_kernel(pmd, 0);
129 }

```

在 bootmem 分配机制还没建立起来之前, 该函数调用 alloc_low_page 分配页表

```

61 static __init void *alloc_low_page(void)
62 {
63     unsigned long pfn = e820_table_end++;

```

在 find_early_table_space 中分配的页表项在这里派上了用场。

```

64     void *adr;
65
66     if (pfn >= e820_table_top)
67         panic("alloc_low_page: ran out of memory");
68

```


69	adr = __va(pfn * PAGE_SIZE);
70	memset(adr, 0, PAGE_SIZE);
71	return adr;

69-71 行：将页面号转换为虚拟地址，将页表内容设置为零，返回页表的虚拟地址。

72 }

4.1.7.3 early_ioremap_page_table_range_init

固定内核映射区页表设置，不设置页表项。

该函数所在文件：

Arch/x86/mm/init_32.c

529	void __init early_ioremap_page_table_range_init(void)
530	{
531	pgd_t *pgd_base = swapper_pg_dir;
532	unsigned long vaddr, end;
533	
534	/*
535	* Fixed mappings, only the page table structure has to be
536	* created - mappings will be set by set_fixmap():
537	*/
538	vaddr = __fix_to_virt(__end_of_fixed_addresses - 1) &
	PMD_MASK;
539	end = (FIXADDR_TOP + PMD_SIZE - 1) & PMD_MASK;

538-539 行：固定内存区域起始地址与终止地址对齐设置

540	page_table_range_init(vaddr, end, pgd_base);
541	early_ioremap_reset();
542	}

4.1.7.4 load_cr3(swapper_pg_dir)

将控制 swapper_pg_dir 送入控制寄存器 cr3。每当重新设置 cr3 时，CPU 就会将页面映射目录所在的页面装入 CPU 内部高速缓存中的 TLB 部分。现在内存中(实际上是高速缓存中)的映射目录变了，就要再让 CPU 装入一次。由于页面映射机制本来就是开启着的，所以从这条指令以后就扩大了系统空间中有映射区域的大小，使整个映射覆盖到整个物理内存(高端内存)除外。实际上此时 swapper_pg_dir 中已经改变的目录项很可能还在高速缓存中，所以还要通过 __flush_tlb_all() 将高速缓存中的内容冲刷到内存中，这样才能保证内存中映射目录内容的一致性。最后，init_memory_mapping 返回了所映射页的数量值，并将值保存在 max_low_pfn_mapped 中。而 32 位 x86 体系中 max_pfn_mapped 的值也为 max_low_pfn_mapped。

4.1.8 initmem_init(0, max_pfn, acpi, k8)

该函数启用初始化期间的内存管理器 early。根据是否设置编译选项 CONFIG_NEED_MULTIPLE_NODES，该函数分为 UMA 版本和 NUMA 版本。尽管这两个版本的代码有较大差异，但是最后都殊途同归：建立 early_node_map 数据结构，存储内核可用内存（RAM），为后面的内存管理做好准备。

下面分别讨论这里个版本的函数。

一、UMA 版本

该版本的函数在 arch/x86/mm/init_32.c 中：

```
707 #ifndef CONFIG_NEED_MULTIPLE_NODES
708 void __init initmem_init(unsigned long start_pfn, unsigned long
end_pfn,
709                          int acpi, int k8)
710 {
711     #ifdef CONFIG_HIGHMEM
712         highstart_pfn = highend_pfn = max_pfn;
713         if (max_pfn > max_low_pfn)
714             highstart_pfn = max_low_pfn;
715         e820_register_active_regions(0, 0, highend_pfn);
716         sparse_memory_present_with_active_regions(0);
717         printk(KERN_NOTICE "%ldMB HIGHMEM available.\n",
```

```

718         pages_to_mb(highend_pfn - highstart_pfn));
719     num_physpages = highend_pfn;
720     high_memory = (void *) __va(highstart_pfn * PAGE_SIZE - 1) +
1;
721 #else
722     e820_register_active_regions(0, 0, max_low_pfn);
723     sparse_memory_present_with_active_regions(0);
724     num_physpages = max_low_pfn;
725     high_memory = (void *) __va(max_low_pfn * PAGE_SIZE - 1) +
1;
726 #endif
727 #ifdef CONFIG_FLATMEM
728     max_mapnr = num_physpages;
729 #endif
730     __vmalloc_start_set = true;
731
732     printk(KERN_NOTICE "%ldMB LOWMEM available.\n",
733           pages_to_mb(max_low_pfn));
734
735     setup_bootmem_allocator();
736 }
737 #endif /* !CONFIG_NEED_MULTIPLE_NODES */

```

该函数调用 `e820_register_active_regions`,
`sparse_memory_present_with_active_regions` 完成相关任务。

4.1.8.1 e820_register_active_regions

函数参数:

`nid`: cpu 节点号

`start_pfn`:起始页面号

`last_pfn`:终止页面号


```
905             unsigned long *ei_startpfn,  
906             unsigned long *ei_endpfn)  
907 {  
908     u64 align = PAGE_SIZE;  
909  
910     *ei_startpfn = round_up(ei->addr, align) >> PAGE_SHIFT;  
911     *ei_endpfn = round_down(ei->addr + ei->size, align) >>  
PAGE_SHIFT;
```

E820 区间内存起始和终止节点页面对齐

```
912  
913     /* Skip map entries smaller than a page */  
914     if (*ei_startpfn >= *ei_endpfn)  
915         return 0;
```

914-915 行：略过小于一页的区间

```
916  
917     /* Skip if map is outside the node */  
918     if (ei->type != E820_RAM || *ei_endpfn <= start_pfn ||  
919         *ei_startpfn >= last_pfn)  
920         return 0;
```

918-920 行：跳过不属于节点的内存区域

```
921  
922     /* Check for overlaps */  
923     if (*ei_startpfn < start_pfn)  
924         *ei_startpfn = start_pfn;
```

923-924 行：e820 中的起始页面号小于节点（node）的起始页面号，则 ei_startpfn 赋值为 node 起始页面号。

```
925     if (*ei_endpfn > last_pfn)  
926         *ei_endpfn = last_pfn;
```

925-926 行：e820 中的终止页面号大于节点（node）的终止页面号，则 ei_endpfn 赋值为终止页面号。

```
927
```

```
928     return 1;
929 }
```

4.1.8.1.2 add_active_range

该函数所在文件为 mm/page_alloc.c 中：

从注释中可以看出：

nid:CPU 节点号

start_pfn: 可用物理内存的起始页面号

end_pfn: 终止节点的起始页面号

该函数的功能：将[start_pfn,end_pfn]对应的范围保存到 early_node_map 数组中，在后面的过程中该数组将会被 free_area_init_nodes 调用，用于计算各个区域的大小以及是否存在空洞。

```
3972 /**
3973  * add_active_range - Register a range of PFNs backed by physical
memory
3974  * @nid: The node ID the range resides on
3975  * @start_pfn: The start PFN of the available physical memory
3976  * @end_pfn: The end PFN of the available physical memory
3977  *
3978  * These ranges are stored in an early_node_map[] and later used
by
3979  * free_area_init_nodes() to calculate zone sizes and holes. If the
3980  * range spans a memory hole, it is up to the architecture to ensure
3981  * the memory is not freed by the bootmem allocator. If possible
3982  * the range being registered will be merged with existing ranges.
3983  */
3984 void __init add_active_range(unsigned int nid, unsigned long
start_pfn,
3985                               unsigned long end_pfn)
3986 {
3987     int i;
```

```

3988
3989     mminit_dprintk(MMINIT_TRACE, "memory_register",
3990                     "Entering add_active_range(%d, %#lx, %#lx) "
3991                     "%d entries of %d used\n",
3992                     nid, start_pfn, end_pfn,
3993                     nr_nodemap_entries, MAX_ACTIVE_REGIONS);
3994
3995     mminit_validate_memmodel_limits(&start_pfn, &end_pfn);
3996
3997     /* Merge with existing active regions if possible */
3998     for (i = 0; i < nr_nodemap_entries; i++) {
3999         if (early_node_map[i].nid != nid)
4000             continue;

```

3999-4000: 节点号不匹配, 跳入下一次循环

```

4001
4002         /* Skip if an existing region covers this new one */
4003         if (start_pfn >= early_node_map[i].start_pfn &&
4004             end_pfn <= early_node_map[i].end_pfn)
4005             return;

```

4003-4005: 已经存在的区域, 直接返回

```

4006
4007         /* Merge forward if suitable */
4008         if (start_pfn <= early_node_map[i].end_pfn &&
4009             end_pfn > early_node_map[i].end_pfn) {
4010             early_node_map[i].end_pfn = end_pfn;
4011             return;
4012         }

```

4008-4012: 前向合并

```

4013
4014         /* Merge backward if suitable */
4015         if (start_pfn < early_node_map[i].start_pfn &&

```

```
4016             end_pfn >= early_node_map[i].start_pfn) {
4017             early_node_map[i].start_pfn = start_pfn;
4018             return;
4019         }
```

4014-4019: 后向合并

```
4020     }
4021
4022     /* Check that early_node_map is large enough */
4023     if (i >= MAX_ACTIVE_REGIONS) {
4024         printk(KERN_CRIT "More than %d memory regions,
truncating\n",
4025                 MAX_ACTIVE_REGIONS);
4026         return;
4027     }
4028
4029     early_node_map[i].nid = nid;
4030     early_node_map[i].start_pfn = start_pfn;
4031     early_node_map[i].end_pfn = end_pfn;
4032     nr_nodemap_entries = i + 1;
```

4029-4032: 如果无法合并, 这建立一个新的项, 同时将 nr_nodemap_entries 加一

```
4033 }
```

有一个小小的疑问: 为啥不考虑

start_pfn <= early_node_map[i].start_pfn &&

end_pfn > early_node_map[i].end_pfn 的情况?

4.1.8.2 sparse_memory_present_with_active_regions

该函数所在文件为 mm/page_alloc.c 中:

```
3463 /**
3464  * sparse_memory_present_with_active_regions - Call
memory_present for each active range
```



```

3465  * @nid: The node to call memory_present for. If
MAX_NUMNODES, all nodes will be used.
3466  *
3467  * If an architecture guarantees that all ranges registered with
3468  * add_active_ranges() contain no holes and may be freed, this
3469  * function may be used instead of calling memory_present()
manually.
3470  */
3471 void __init sparse_memory_present_with_active_regions(int nid)
3472 {
3473     int i;
3474
3475     for_each_active_range_index_in_nid(i, nid)
3476         memory_present(early_node_map[i].nid,
3477                       early_node_map[i].start_pfn,
3478                       early_node_map[i].end_pfn);
3479 }

```

该函数使用 `for_each_active_range_index_in_nid(i, nid)` 宏（具体定义见下），探测 `early_node_map` 中节点号为 `nid` 的所有内存区域。`memory_present` 函数就不细看了。

```

3360 /* Basic iterator support to walk early_node_map[] */
3361 #define for_each_active_range_index_in_nid(i, nid) \
3362     for (i = first_active_region_index_in_nid(nid); i != -1; \
3363          i = next_active_region_index_in_nid(i, nid))
3364

```

```

3287  * Basic iterator support. Return the first range of PFNs for a node
3288  * Note: nid == MAX_NUMNODES returns first region regardless of
node
3289  */
3290 static int __meminit first_active_region_index_in_nid(int nid)

```

```

3291 {
3292     int i;
3293
3294     for (i = 0; i < nr_nodemap_entries; i++)
3295         if (nid == MAX_NUMNODES || early_node_map[i].nid ==
nid)
3296             return i;

```

3295: 如果 `nid` 不是用来指定 `cpu` 节点号，则直接返回 0，如果 `early_node_map` 中有 `nid` 的内存区域，则返回第一个匹配的内存区域号。

```

3297
3298     return -1;
3299 }

```

```

3302  * Basic iterator support. Return the next active range of PFNs for a
node
3303  * Note: nid == MAX_NUMNODES returns next region regardless of
node
3304  */
3305 static int __meminit next_active_region_index_in_nid(int index, int
nid)
3306 {
3307     for (index = index + 1; index < nr_nodemap_entries; index++)
3308         if (nid == MAX_NUMNODES || early_node_map[index].nid
== nid)
3309             return index;

```

该函数的基本功能同上，只不过是查找的起始号从传入的参数 `index` 开始。

```

3310
3311     return -1;
3312 }

```

4.1.8.3 setup_bootmem_allocator()

```
775 void __init setup_bootmem_allocator(void)
776 {
777     #ifndef CONFIG_NO_BOOTMEM
778         int nodeid;
779         unsigned long bootmap_size, bootmap;
780         /*
781          * Initialize the boot-time allocator (with low memory only):
782          */
783         bootmap_size =
bootmem_bootmap_pages(max_low_pfn)<<PAGE_SHIFT;
784         bootmap = find_e820_area(0,
max_pfn_mapped<<PAGE_SHIFT, bootmap_size,
785                                 PAGE_SIZE);
786         if (bootmap == -1L)
787             panic("Cannot find bootmem map of size %ld\n",
bootmap_size);
788         reserve_early(bootmap, bootmap + bootmap_size,
"BOOTMAP");
789     #endif
790
791     printk(KERN_INFO "    mapped low ram: 0 - %08lx\n",
792            max_pfn_mapped<<PAGE_SHIFT);
793     printk(KERN_INFO "    low ram: 0 - %08lx\n",
max_low_pfn<<PAGE_SHIFT);
794
795     #ifndef CONFIG_NO_BOOTMEM
796         for_each_online_node(nodeid) {
797             unsigned long start_pfn, end_pfn;
798
799             #ifdef CONFIG_NEED_MULTIPLE_NODES
```

```

800      start_pfn = node_start_pfn[nodeid];
801      end_pfn = node_end_pfn[nodeid];
802      if (start_pfn > max_low_pfn)
803          continue;
804      if (end_pfn > max_low_pfn)
805          end_pfn = max_low_pfn;
806 #else
807      start_pfn = 0;
808      end_pfn = max_low_pfn;
809 #endif
810      bootmap = setup_node_bootmem(nodeid, start_pfn,
end_pfn,
811                                bootmap);
812  }
813 #endif
814
815  after_bootmem = 1;
816 }

```

该函数在设置了 **CONFIG_NO_BOOTMEM** 后，基本上啥事也不敢，只是将 **after_bootmem** 设置为 1。现在一般弃用 **BOOTMEM**，所以没有设置 **CONFIG_NO_BOOTMEM** 的过程这里不讨论了。

二 、 NUMA 版

NUMA 版的函数在 `arch/x86/mm/numa_32.c` 中：

```

350 void __init initmem_init(unsigned long start_pfn, unsigned long
end_pfn,
351                          int acpi, int k8)
352 {
353     int nid;
354     long kva_target_pfn;
355
356     /*

```

```

357      * When mapping a NUMA machine we allocate the
node_mem_map arrays
358      * from node local memory.  They are then mapped directly into
KVA
359      * between zone normal and vmalloc space.  Calculate the size
of
360      * this space and use it to adjust the boundary between
ZONE_NORMAL
361      * and ZONE_HIGHMEM.
362      */
363
364      get_memcfg_numa();
365
366      kva_pages = roundup(calculate_numa_remap_pages(),
PTRS_PER_PTE);
367
368      kva_target_pfn = round_down(max_low_pfn - kva_pages,
PTRS_PER_PTE);
369      do {
370          kva_start_pfn =
find_e820_area(kva_target_pfn<<PAGE_SHIFT,
371                max_low_pfn<<PAGE_SHIFT,
372                kva_pages<<PAGE_SHIFT,
373                PTRS_PER_PTE<<PAGE_SHIFT) >>
PAGE_SHIFT;
374          kva_target_pfn -= PTRS_PER_PTE;
375      } while (kva_start_pfn == -1UL && kva_target_pfn >
min_low_pfn);
376
377      if (kva_start_pfn == -1UL)
378          panic("Can not get kva space\n");
379

```

```

380     printk(KERN_INFO "kva_start_pfn ~ %lx max_low_pfn ~ %lx\n",
381             kva_start_pfn, max_low_pfn);
382     printk(KERN_INFO "max_pfn = %lx\n", max_pfn);
383
384     /* avoid clash with initrd */
385     reserve_early(kva_start_pfn<<PAGE_SHIFT,
386                 (kva_start_pfn + kva_pages)<<PAGE_SHIFT,
387                 "KVA PG");
388 #ifdef CONFIG_HIGHMEM
389     highstart_pfn = highend_pfn = max_pfn;
390     if (max_pfn > max_low_pfn)
391         highstart_pfn = max_low_pfn;
392     printk(KERN_NOTICE "%ldMB HIGHMEM available.\n",
393           pages_to_mb(highend_pfn - highstart_pfn));
394     num_physpages = highend_pfn;
395     high_memory = (void *) __va(highstart_pfn * PAGE_SIZE - 1) +
1;
396 #else
397     num_physpages = max_low_pfn;
398     high_memory = (void *) __va(max_low_pfn * PAGE_SIZE - 1) +
1;
399 #endif
400     printk(KERN_NOTICE "%ldMB LOWMEM available.\n",
401           pages_to_mb(max_low_pfn));
402     printk(KERN_DEBUG "max_low_pfn = %lx, highstart_pfn =
%lx\n",
403           max_low_pfn, highstart_pfn);
404
405     printk(KERN_DEBUG "Low memory ends at vaddr %08lx\n",
406           (ulong) pfn_to_kaddr(max_low_pfn));
407     for_each_online_node(nid) {
408         init_remap_allocator(nid);

```

```

409
410     allocate_pgdat(nid);
411 }
412     remap_numa_kva();
413
414     printk(KERN_DEBUG "High memory starts at vaddr %08lx\n",
415            (ulong) pfn_to_kaddr(highstart_pfn));
416     for_each_online_node(nid)
417         propagate_e820_map_node(nid);
418
419     for_each_online_node(nid) {
420         memset(NODE_DATA(nid), 0, sizeof(struct pglist_data));
421         NODE_DATA(nid)->node_id = nid;
422 #ifndef CONFIG_NO_BOOTMEM
423         NODE_DATA(nid)->bdata = &bootmem_node_data[nid];
424 #endif
425     }
426
427     setup_bootmem_allocator();
428 }

```

4.1.8.4 acpi_numa_init

要查看该函数，需要回去看看 `Setup_arch` 函数中的第 976 行的 `acpi_numa_init` 函数。

```

972 #ifdef CONFIG_ACPI_NUMA
973     /*
974      * Parse SRAT to discover nodes.
975      */
976     acpi = acpi_numa_init();
977 #endif

```

该函数所在的文件为 drivers/acpi/numa.c 中：

```
277 int __init acpi_numa_init(void)
278 {
279     int ret = 0;
280
281     /* SRAT: Static Resource Affinity Table */
282     if (!acpi_table_parse(ACPI_SIG_SRAT, acpi_parse_srat)) {
283
acpi_table_parse_srat(ACPI_SRAT_TYPE_X2APIC_CPU_AFFINITY,
284                     acpi_parse_x2apic_affinity, nr_cpu_ids);
285
acpi_table_parse_srat(ACPI_SRAT_TYPE_CPU_AFFINITY,
286                     acpi_parse_processor_affinity, nr_cpu_ids);
287     ret =
acpi_table_parse_srat(ACPI_SRAT_TYPE_MEMORY_AFFINITY,
288                     acpi_parse_memory_affinity,
289                     NR_NODE_MEMBLKS);
290     }
```

该函数使用 ACPI 的 SRAT 来处理与内存与处理器相关的函数，这里只关心 [acpi_parse_memory_affinity](#) 函数。

```
238 static int __init
239 acpi_parse_memory_affinity(struct acpi_subtable_header * header,
240                           const unsigned long end)
241 {
242     struct acpi_srat_mem_affinity *memory_affinity;
243
244     memory_affinity = (struct acpi_srat_mem_affinity *)header;
245     if (!memory_affinity)
246         return -EINVAL;
247
248     acpi_table_print_srat_entry(header);
249 }
```



```

250     /* let architecture-dependent part to do it */
251     acpi_numa_memory_affinity_init(memory_affinity);
252
253     return 0;
254 }

```

该函数调用 `acpi_numa_memory_affinity_init` 函数，函数所在文件 `arch/x86/mm/srat_32.c`

该函数是供 ACPI 驱动调用的。用来设置 `node_memory_chunk` 全局变量。该变量的定义如下：

```

47 #define MAX_CHUNKS_PER_NODE 3
48 #define MAXCHUNKS      (MAX_CHUNKS_PER_NODE *
MAX_NUMNODES)
49 struct node_memory_chunk_s {
50     unsigned long    start_pfn; //起始页面号
51     unsigned long    end_pfn; //终止页面号
52     u8    pxm;        // proximity domain of node 邻近域的节点
53     u8    nid;        // which cnode contains this chunk?
54     u8    bank;       // which mem bank on this node
55 };
56 static struct node_memory_chunk_s __initdata
node_memory_chunk[MAXCHUNKS];

```

56 行：全局变量 `node_memory_chunk`，记录每个内存 chunk 的数组

```

57
58 static int __initdata num_memory_chunks; /* total number of memory
chunks */

```

58 行：全局变量 `num_memory_chunks`，记录总的内存 chunk 的数目。

下面看函数的具体代码：

```

105 acpi_numa_memory_affinity_init(struct acpi_srat_mem_affinity
*memory_affinity)
106 {
107     unsigned long long paddr, size;

```

```

108     unsigned long start_pfn, end_pfn;
109     u8 pxm;
110     struct node_memory_chunk_s *p, *q, *pend;
111
112     if (srat_disabled())
113         return;
114     if (memory_affinity->header.length !=
115         sizeof(struct acpi_srat_mem_affinity)) {
116         bad_srat();
117         return;
118     }
119
120     if ((memory_affinity->flags & ACPI_SRAT_MEM_ENABLED) ==
0)
121         return;    /* empty entry */
122
123     pxm = memory_affinity->proximity_domain & 0xff;
124
125     /* mark this node as "seen" in node bitmap */
126     BMAP_SET(pxm_bitmap, pxm);
127
128     /* calculate info for memory chunk structure */
129     paddr = memory_affinity->base_address;
130     size = memory_affinity->length;
131
132     start_pfn = paddr >> PAGE_SHIFT;
133     end_pfn = (paddr + size) >> PAGE_SHIFT;
134
135
136     if (num_memory_chunks >= MAXCHUNKS) {
137         printk(KERN_WARNING "Too many mem chunks in SRAT."
138             " Ignoring %lld MBytes at %llx\n",

```

```
139         size/(1024*1024), paddr);
140     return;
141 }
```

136-141: 如果 num_memory_chunks 超过内核设置上限, 则退出函数

```
142
143     /* Insertion sort based on base address */
144     pend = &node_memory_chunk[num_memory_chunks];
145     for (p = &node_memory_chunk[0]; p < pend; p++) {
146         if (start_pfn < p->start_pfn)
147             break;
148     }
149     if (p < pend) {
150         for (q = pend; q >= p; q--)
151             *(q + 1) = *q;
152     }
153     p->start_pfn = start_pfn;
154     p->end_pfn = end_pfn;
155     p->pxm = pxm;
156
157     num_memory_chunks++;
```

153-157 行: 设置 chunk 的各个变量。

```
158
159     printk(KERN_DEBUG "Memory range %08lx to %08lx"
160            " in proximity domain %02x %s\n",
161            start_pfn, end_pfn,
162            pxm,
163            ((memory_affinity->flags &
ACPI_SRAT_MEM_HOT_PLUGGABLE) ?
164            "enabled and removable" : "enabled" ) );
165 }
```

4.1.8.5 get_memcfg_numa()

在 `acpi_numa_init` 中设置了 `node_memory_chunk` 后，`get_memcfg_numa` 将使用改变了来初始化 `node_early_map` 变量。

该函数所在的文件为： `arch/x86/include/asm/mmzone_32.h`

```
21  * This allows any one NUMA architecture to be compiled
22  * for, and still fall back to the flat function if it
23  * fails.
24  */
25 static inline void get_memcfg_numa(void)
26 {
27
28     if (get_memcfg_numaq()) //IBM NUMA-Q
29         return;
30     if (get_memcfg_from_srat()) //intel srat
31         return;
32     get_memcfg_numa_flat();
33 }
```

INTEL 处理器通过 ACPI 规范对硬件资源进行管理。其中 SRAT 给出了全系统中的 CPU 和 MEM 的分布情况。上面 `get_memcfg_from_srat()` 就是该过程的实现。

下面只分析 `get_memcfg_from_srat()` 函数的实现：

该函数所在的文件为 `arch/x86/mm/srat_32.c`

```
209 int __init get_memcfg_from_srat(void)
210 {
211     int i, j, nid;
212
213
214     if (srat_disabled())
215         goto out_fail;
```

214-215: 检测是否支持 srat。

```
216
```

```

217     if (num_memory_chunks == 0) {
218         printk(KERN_DEBUG
219             "could not find any ACPI SRAT memory areas.\n");
220         goto out_fail;
221     }

```

217-221: 如果 srat 内存区域为 0, 则跳转到 out_fail

```

222
223     /* Calculate total number of nodes in system from PXM bitmap
and create
224     * a set of sequential node IDs starting at zero.  (ACPI doesn't
seem
225     * to specify the range of _PXM values.)
226     */
227     /*
228     * MCD - we no longer HAVE to number nodes sequentially.
PXM domain
229     * numbers could go as high as 256, and MAX_NUMNODES for
i386 is typically
230     * 32, so we will continue numbering them in this manner until
MAX_NUMNODES
231     * approaches MAX_PXM_DOMAINS for i386.
232     */
233     nodes_clear(node_online_map);
234     for (i = 0; i < MAX_PXM_DOMAINS; i++) {
235         if (BMAP_TEST(pxm_bitmap, i)) {
236             int nid = acpi_map_pxm_to_node(i);
237             node_set_online(nid);
238         }
239     }
240     BUG_ON(num_online_nodes() == 0);
241
242     /* set cnode id in memory chunk structure */

```

```

243     for (i = 0; i < num_memory_chunks; i++)
244         node_memory_chunk[i].nid =
pxm_to_node(node_memory_chunk[i].pxm);
245
246     printk(KERN_DEBUG "pxm bitmap: ");
247     for (i = 0; i < sizeof(pxm_bitmap); i++) {
248         printk(KERN_CONT "%02x ", pxm_bitmap[i]);
249     }
250     printk(KERN_CONT "\n");
251     printk(KERN_DEBUG "Number of logical nodes in system =
%d\n",
252             num_online_nodes());
253     printk(KERN_DEBUG "Number of memory chunks in system =
%d\n",
254             num_memory_chunks);
255
256     for (i = 0; i < MAX_APICID; i++)
257         apicid_2_node[i] = pxm_to_node(apicid_to_pxm[i]);

```

233-257 行：初始化一系列变量后，来到 259 行

```

258
259     for (j = 0; j < num_memory_chunks; j++){
260         struct node_memory_chunk_s * chunk =
&node_memory_chunk[j];
261         printk(KERN_DEBUG
262             "chunk %d nid %d start_pfn %08lx end_pfn %08lx\n",
263             j, chunk->nid, chunk->start_pfn, chunk->end_pfn);
264         if (node_read_chunk(chunk->nid, chunk))
265             continue;

```

264-265 行：验证 chunk 的有效性，有需要的话修改节点 node 的 node_start_pfn 变量。

```

266

```

```

267         e820_register_active_regions(chunk->nid,
chunk->start_pfn,
268         min(chunk->end_pfn, max_pfn));
269     }

```

267-269 行：根据 chunk 初始化 early_node_map 全局变量，该函数在前面已经分析过。

```

270     /* for out of order entries in SRAT */
271     sort_node_map();
272
273     for_each_online_node(nid) {
274         unsigned long start = node_start_pfn[nid];
275         unsigned long end = min(node_end_pfn[nid], max_pfn);
276
277         memory_present(nid, start, end);
278         node_remap_size[nid] = node_memmap_size_bytes(nid,
start, end);
279     }
280     return 1;
281 out_fail:
282     printk(KERN_DEBUG "failed to get NUMA memory information
from SRAT"
283           " table\n");
284     return 0;
285 }

```

NUMA 函数版的 initmem_init 深入的函数在这就不在分析了。

对比一下 UMA 和 NUMA 的 initmem_init 函数：

UMA 的使用 E820 建立的数据结构，来初始化 early_node_map 全局变量。

NUMA 使用 ACPI 的 SRAT 建立的数据结构，来初始化 early_node_map 全局变量。

当记录内存区域的变量 early_node_map 建立好后，在其基础上调用 paging_init 建立物理内存管理的相关数据结构。

4.1.9 paging_init

该函数所在的文件为 arch/x86/mm/init_32.c

```
818 /*
819  * paging_init() sets up the page tables - note that the first 8MB are
820  * already mapped by head.S.
821  *
822  * This routines also unmaps the page at virtual kernel address 0,
so
823  * that we can trap those pesky NULL-reference errors in the kernel.
824  */
825 void __init paging_init(void)
826 {
827     pagetable_init();
828
829     __flush_tlb_all();
830
831     kmap_init();
832
833     /*
834      * NOTE: at this point the bootmem allocator is fully available.
835      */
836     sparse_init();
837     zone_sizes_init();
838 }
```

该函数分别调用：

Pageable_init: 初始化永久内核页表

Kmap_init: 临时内核映射

Sparse_init 初始化

zone_sizes_init: 内存区域数据结构初始化。

4.1.9.1 永久内核页表分配

`pagetable_init`: 初始化永久内核页表。分配一个页框存储页表，但页表中的页表项并没有指向相应的页框。

```
544 static void __init pagetable_init(void)
545 {
546     pgd_t *pgd_base = swapper_pg_dir; //pgd_base 指向页目录
547
548     permanent_kmaps_init(pgd_base);
549 }
```

546 行: `pgd_base` 指向全局页目录表

548 行: 调用 `permanent_kmaps_init` 执行永久内核区域的初始化

```
398 #ifdef CONFIG_HIGHMEM
399 static void __init permanent_kmaps_init(pgd_t *pgd_base)
400 {
401     unsigned long vaddr;
402     pgd_t *pgd;
403     pud_t *pud;
404     pmd_t *pmd;
405     pte_t *pte;
406
407     #define PKMAP_BASE ((FIXADDR_BOOT_START -
PAGE_SIZE * (LAST_PKMAP + 1)) & PMD_MASK)
408     vaddr = PKMAP_BASE; //永久内核映射虚拟地址开始区
//LAST_PKMAP=1024 PAGE_SIZE=4K
409     page_table_range_init(vaddr, vaddr +
PAGE_SIZE*LAST_PKMAP, pgd_base);
410
411     pgd = swapper_pg_dir + pgd_index(vaddr);
412     pud = pud_offset(pgd, vaddr);
413     pmd = pmd_offset(pud, vaddr);
```

```

413     pte = pte_offset_kernel(pmd, vaddr);
414     pkmap_page_table = pte;
415 }

```

407 行: `vaddr` 是永久内核映射的起始地址。被赋值为 `PKMAP_BASE`。
`PKMAP_BASE` 的定义如下:

```

#define PKMAP_BASE ((FIXADDR_BOOT_START - PAGE_SIZE *
(LAST_PKMAP + 1)) & PMD_MASK)

```

其中, `FIXADDR_BOOT_START` 为固定映射的起始虚拟地址, `LASK_PKMAP` 为 1024, 整个表达式的含义是 永久内核映射的起始地址 为: 固定映射的起始地址 -4M-4K, 并且地址为 `PMD_MASK` 对齐。之所以在永久内核映射与固定映射之间有 4K 的空隙, 适用于内核捕获内存访问异常使用的。

408 行: 调用 `page_table_range_init` 函数, 为虚拟地址[`vaddr,vaddr+4M`]在全局页目录表中设置页目录项, 并且将页目录项指向分配的页表。

410-414 行: 在全局页目录表中设置好页目录项后, 通过 `pud_offset`、`pmd_offset`、`pte_offset_kernel` 函数的到分配的页表的虚拟地址, 并且将虚拟地址赋值给全局变量 `pkmap_page_table`。

[这里需要重点讲解一下 pte_offset_kernel 函数。](#)

我们知道, 在页目录中的页目录项、页表中的页表项存放的都是对应页框的物理地址, `pte_offset_kernel` 函数如下:

```

384 static inline pte_t *pte_offset_kernel(pmd_t *pmd, unsigned long
address)
385 {
386     return (pte_t *)pmd_page_vaddr(*pmd) + pte_index(address);
387 }

```

`(pte_t *)pmd_page_vaddr(*pmd)`将也中间目录项的物理地址转换为虚拟地址, `pte_index(address)`的到虚拟地址 `address` 在 `pmd` 指向的页表中的偏移。

整个函数的作用: 虚拟地址 `address` 在 `pmd` (物理地址) 指向的页表中有一个对应的表项, 有一个指向这个表项的指针 `p_pte`, `pte_offset_kernel` 的作用就是将 `p_pte`(物理地址)转换为虚拟地址。

[page_table_range_init](#) 为虚拟地址[`vaddr,vaddr+4M`]在全局页目录表中设置页目录项, 并且将页目录项指向分配的页表。

```

189 /*
190  * This function initializes a certain range of kernel virtual memory

```

```

191  * with new bootmem page tables, everywhere page tables are
missing in
192  * the given range.
193  *
194  * NOTE: The pagetables are allocated contiguous on the physical
space
195  * so we can cache the place of the first one and move around
without
196  * checking the pgd every time.
197  */
198 static void __init
199 page_table_range_init(unsigned long start, unsigned long end,
pgd_t *pgd_base)
200 {
201     int pgd_idx, pmd_idx;
202     unsigned long vaddr;
203     pgd_t *pgd;
204     pmd_t *pmd;
205     pte_t *pte = NULL;
206
207     vaddr = start;
208     pgd_idx = pgd_index(vaddr); //页目录索引
209     pmd_idx = pmd_index(vaddr); //中间页表索引
210     pgd = pgd_base + pgd_idx;
211 //内外 for 循环都指执行一次
212     for ( ; (pgd_idx < PTRS_PER_PGD) && (vaddr != end); pgd++,
pgd_idx++) {
213         pmd = one_md_table_init(pgd); //pmd 指向中间页表
214         pmd = pmd + pmd_index(vaddr); //中间页表项
215         for ( ; (pmd_idx < PTRS_PER_PMD) && (vaddr != end);
216             pmd++, pmd_idx++) {
217             pte=page_table_kmap_check(one_page_table_init(pmd),

```

```

218                                     pmd, vaddr, pte);
219
220             vaddr += PMD_SIZE;
221         }
222         pmd_idx = 0;
223     }
224 }

```

212 215 内外 for 循环都只是执行一次。

213 行: `one_md_table_init` 函数在没有启用物理地址扩展选项，并且内核采用二级页目录时，只是将 `pgd` 赋值给 `pmd`。

217 行: `one_page_table_init(pmd)`: 分配一个页框，并将页框的物理地址赋值给 `pmd` 指向的页目录项（二级页面寻址）。

```

105  * Create a page table and place a pointer to it in a middle page
106  * directory entry:
107  */
108 static pte_t * __init one_page_table_init(pmd_t *pmd)
109 {
110     if (!(pmd_val(*pmd) & _PAGE_PRESENT)) {
111         pte_t *page_table = NULL;
112
113         if (after_bootmem) {
114             #if defined(CONFIG_DEBUG_PAGEALLOC) ||
defined(CONFIG_KMEMCHECK)
115                 page_table = (pte_t *)
alloc_bootmem_pages(PAGE_SIZE);
116             #endif
117             if (!page_table)
118                 page_table =
119                 (pte_t *)alloc_bootmem_pages(PAGE_SIZE);
120         } else

```

```

121          page_table = (pte_t *)alloc_low_page();//分配一个页
框，用来作为页表
122
123          paravirt_alloc_pte(&init_mm, __pa(page_table) >>
PAGE_SHIFT);
124          set_pmd(pmd, __pmd(__pa(page_table) |
_PAGE_TABLE));
125          BUG_ON(page_table != pte_offset_kernel(pmd, 0));
126      }
127
128      return pte_offset_kernel(pmd, 0);
129 }

```

4.1.9.2 临时内核映射初始化

kmap_init()初始化临时内核映射。

```

385 static void __init kmap_init(void)
386 {
387     unsigned long kmap_vstart;
388
389     /*
390      * Cache the first kmap pte:
391      */
392     kmap_vstart = __fix_to_virt(FIX_KMAP_BEGIN);
393     kmap_pte = kmap_get_fixmap_pte(kmap_vstart);
394
395     kmap_prot = PAGE_KERNEL;
396 }

```

392 行：获得临时内核映射虚拟地址的起始地址。

393 行：将临时内核映射的第一个页表的虚拟地址保存在全局变量 kmap_pte 中。

4.1.9.3 内存区域初始化

4.1.9.3.1 zone_size_init

Linux 内核采用一个三级的结构来管理内存，从节点->区域->页面。

Zone_size_init 用来初始化节点、区域中的部分数据结构。

```
#define MAX_DMA_ADDRESS      (PAGE_OFFSET + 0x1000000)
739 static void __init zone_sizes_init(void)
740 {
741     unsigned long max_zone_pfns[MAX_NR_ZONES];
742     memset(max_zone_pfns, 0, sizeof(max_zone_pfns));
743     max_zone_pfns[ZONE_DMA] =
744         virt_to_phys((char *)MAX_DMA_ADDRESS) >>
PAGE_SHIFT;
745     max_zone_pfns[ZONE_NORMAL] = max_low_pfn;
746 #ifdef CONFIG_HIGHMEM
747     max_zone_pfns[ZONE_HIGHMEM] = highend_pfn;
748 #endif
749
750     free_area_init_nodes(max_zone_pfns);
751 }
```

Zone_sizes_init 的主要作用是初始化 max_zone_pfns 数组，然后将 max_zone_pfns 作为参数调用 free_area_init_nodes 函数。

743、745、747 行分别设置 DMA 区、NORMAL 区、HIGHMEM 的边界。

DMA 区的边界为：0-16M。

NORMAL 的边界为：16M-max_low_pfn*4k，max_low_pfn 是在前面内核初始化过程中计算出来的，如果有高端内存，则 max_low_pfn*4k 为 896M。

HIGHMEM 的边界为：896M-高端内存的结尾（小于 4G）。

4.1.9.3.2 free_area_init_nodes

free_area_init_nodes 函数设置必要的变量后，初始化每一个节点。

Mm/page_alloc. C

```

4346 * free_area_init_nodes - Initialise all pg_data_t and zone data
4347 * @max_zone_pfn: an array of max PFNs for each zone
4348 *
4349 * This will call free_area_init_node() for each active node in the system.
4350 * Using the page ranges provided by add_active_range(), the size of each
4351 * zone in each node and their holes is calculated. If the maximum PFN
4352 * between two adjacent zones match, it is assumed that the zone is empty.
4353 * For example, if arch_max_dma_pfn == arch_max_dma32_pfn, it is assumed
4354 * that arch_max_dma32_pfn has no pages. It is also assumed that a zone
4355 * starts where the previous one ended. For example, ZONE_DMA32 starts
4356 * at arch_max_dma_pfn.
4357 */
4358 void __init free_area_init_nodes(unsigned long *max_zone_pfn)
4359 {
4360     unsigned long nid;
4361     int i;
4362
4363     /* Sort early_node_map as initialisation assumes it is sorted */
4364     sort_node_map();
4365
4366     /* Record where the zone boundaries are */
4367     memset(arch_zone_lowest_possible_pfn, 0,
4368            sizeof(arch_zone_lowest_possible_pfn));
4369     memset(arch_zone_highest_possible_pfn, 0,
4370            sizeof(arch_zone_highest_possible_pfn));

```

4367、4369 初始化静态数组 `arch_zone_lowest_possible_pfn`，`arch_zone_highest_possible_pfn`，两个数组同一索引出的元素分别记录某一个内存区（zone）的起始页面号和终止页面号。

```

180     static unsigned long __meminitdata
arch_zone_lowest_possible_pfn[MAX_NR_ZONES];

```

```
181 static unsigned long __meminitdata
arch_zone_highest_possible_pfn[MAX_NR_ZONES];
```

```
4371 arch_zone_lowest_possible_pfn[0] = find_min_pfn_with_active_regions();
```

find_min_pfn_with_active_regions 函数调用 find_min_pfn_for_node 函数查找指定（该处为所有节点）节点中最小的页面号。

```
4136 /* Find the lowest pfn for a node */
4137 static unsigned long __init find_min_pfn_for_node(int nid)
4138 {
4139     int i;
4140     unsigned long min_pfn = ULONG_MAX;
4141
4142     /* Assuming a sorted map, the first range found has the starting pfn */
4143     for_each_active_range_index_in_nid(i, nid)
4144         min_pfn = min(min_pfn, early_node_map[i].start_pfn);
4145
4146     if (min_pfn == ULONG_MAX) {
4147         printk(KERN_WARNING
4148             "Could not find start_pfn for node %d\n", nid);
4149         return 0;
4150     }
4151
4152     return min_pfn;
4153 }
```

4143-4144 行：使用内核在之前初始化的数据结构 early_node_map，查找所有节点中页面号最小的页面。

```
4372 arch_zone_highest_possible_pfn[0] = max_zone_pfn[0];
4373 for (i = 1; i < MAX_NR_ZONES; i++) {
4374     if (i == ZONE_MOVABLE)
4375         continue;
```



```
4376     arch_zone_lowest_possible_pfn[i] =
4377         arch_zone_highest_possible_pfn[i-1];
4378     arch_zone_highest_possible_pfn[i] =
4379         max(max_zone_pfn[i], arch_zone_lowest_possible_pfn[i]);
4380 }
```

4372-4380: 根据 max_zone_pfn 来计算各初始化数组。

```
4381     arch_zone_lowest_possible_pfn[ZONE_MOVABLE] = 0;
4382     arch_zone_highest_possible_pfn[ZONE_MOVABLE] = 0;
4383
4384     /* Find the PFNs that ZONE_MOVABLE begins at in each node */
4385     memset(zone_movable_pfn, 0, sizeof(zone_movable_pfn));
4386     find_zone_movable_pfns_for_nodes(zone_movable_pfn);
4387
4388     /* Print out the zone ranges */
4389     printk("Zone PFN ranges:\n");
4390     for (i = 0; i < MAX_NR_ZONES; i++) {
4391         if (i == ZONE_MOVABLE)
4392             continue;
4393         printk("  %-8s ", zone_names[i]);
4394         if (arch_zone_lowest_possible_pfn[i] ==
4395             arch_zone_highest_possible_pfn[i])
4396             printk("empty\n");
4397         else
4398             printk("%0#10lx -> %0#10lx\n",
4399                 arch_zone_lowest_possible_pfn[i],
4400                 arch_zone_highest_possible_pfn[i]);
4401     }
4402
4403     /* Print out the PFNs ZONE_MOVABLE begins at in each node */
4404     printk("Movable zone start PFN for each node\n");
4405     for (i = 0; i < MAX_NUMNODES; i++) {
```

```

4406         if (zone_movable_pfn[i])
4407             printk("  Node %d: %lu\n", i, zone_movable_pfn[i]);
4408     }
4409
4410     /* Print out the early_node_map[] */
4411     printk("early_node_map[%d] active PFN ranges\n", nr_nodemap_entries);
4412     for (i = 0; i < nr_nodemap_entries; i++)
4413         printk("  %3d: %0#10lx -> %0#10lx\n", early_node_map[i].nid,
4414             early_node_map[i].start_pfn,
4415             early_node_map[i].end_pfn);
4416
4417     /* Initialise every node */
4418     mminit_verify_pageflags_layout();
4419     setup_nr_node_ids();
4420     for_each_online_node(nid) {
4421         pg_data_t *pgdat = NODE_DATA(nid); //节点数据
4422         free_area_init_node(nid, NULL,
4423             find_min_pfn_for_node(nid), NULL);
4424
4425         /* Any memory on that node */
4426         if (pgdat->node_present_pages)
4427             node_set_state(nid, N_HIGH_MEMORY);
4428         check_for_regular_memory(pgdat);
4429     }
4430 }

```

4422 行调用 `free_area_init_node` 来初始化一个节点。

4.1.9.3.3 free_area_init_node

//初始化指定节点的相关数据结构

```

3932 void __paginginit free_area_init_node(int nid, unsigned long *zones_size,

```

```

3933         unsigned long node_start_pfn, unsigned long *zholes_size)
3934 {
3935     pg_data_t *pgdat = NODE_DATA(nid);
3936
3937     pgdat->node_id = nid;
3938     pgdat->node_start_pfn = node_start_pfn;
3939     calculate_node_totalpages(pgdat, zones_size, zholes_size);
3940
3941     alloc_node_mem_map(pgdat);
3942 #ifdef CONFIG_FLAT_NODE_MEM_MAP
3943     printk(KERN_DEBUG "free_area_init_node: node %d, pgdat %08lx,
node_mem_map %08lx\n",
3944         nid, (unsigned long)pgdat,
3945         (unsigned long)pgdat->node_mem_map);
3946 #endif
3947
3948     free_area_init_core(pgdat, zones_size, zholes_size);
3949 }

```

3937-3938 行：初始话节点的相关变量。

3939 行：计算节点中所有内存区总的节点数，空洞数。

3941 行：为节点分配表示页框的数据结构 **page** 所需要的连续内存。

3949 行：调用 **free_area_init_core** 初始化节点。

//计算节点中总的页面数以及空洞数

```

3699 static void __meminit calculate_node_totalpages(struct pglist_data *pgdat,
3700         unsigned long *zones_size, unsigned long *zholes_size)
3701 {
3702     unsigned long realtotalpages, totalpages = 0;
3703     enum zone_type i;
3704
3705     for (i = 0; i < MAX_NR_ZONES; i++)
3706         totalpages += zone_spanned_pages_in_node(pgdat->node_id, i,

```

```

3707                                     zones_size);
3708     pgdat->node_spanned_pages = totalpages;
3709
3710     realtotalpages = totalpages;
3711     for (i = 0; i < MAX_NR_ZONES; i++)
3712         realtotalpages -=
3713             zone_absent_pages_in_node(pgdat->node_id, i,
3714                                     zholes_size);
3715     pgdat->node_present_pages = realtotalpages;
3716     printk(KERN_DEBUG "On node %d totalpages: %lu\n", pgdat->node_id,
3717            realtotalpages);
3718 }

```

3705-3708 行：计算制定节点中所有区域的中的页面数，并赋值给节点的 `node_spanned_pages` 变量。

3711-3714 行：计算节点中各个区域的空洞页数，用 `totalpages` 减去空洞页数，得到实际的页面数，赋值给节点的 `node_present_pages` 变量。

```

3565 /*
3566  * Return the number of pages a zone spans in a node, including holes
3567  * present_pages = zone_spanned_pages_in_node() - zone_absent_pages_in_node()
3568  */
3569 static unsigned long __meminit zone_spanned_pages_in_node(int nid,
3570                  unsigned long zone_type,
3571                  unsigned long *ignored)
3572 {
3573     unsigned long node_start_pfn, node_end_pfn;
3574     unsigned long zone_start_pfn, zone_end_pfn;
3575
3576     /* Get the start and end of the node and zone */
3577     get_pfn_range_for_nid(nid, &node_start_pfn, &node_end_pfn);
3578     zone_start_pfn = arch_zone_lowest_possible_pfn[zone_type];

```

```

3579     zone_end_pfn = arch_zone_highest_possible_pfn[zone_type];
3580     adjust_zone_range_for_zone_movable(nid, zone_type,
3581         node_start_pfn, node_end_pfn,
3582         &zone_start_pfn, &zone_end_pfn);
3583
3584     /* Check that this node has pages within the zone's required range */
3585     if (zone_end_pfn < node_start_pfn || zone_start_pfn > node_end_pfn)
3586         return 0;
3587
3588     /* Move the zone boundaries inside the node if necessary */
3589     zone_end_pfn = min(zone_end_pfn, node_end_pfn);
3590     zone_start_pfn = max(zone_start_pfn, node_start_pfn);
3591
3592     /* Return the spanned pages */
3593     return zone_end_pfn - zone_start_pfn;
3594 }

```

该函数通过节点的起始页面和终止页面与内存区域的起始页面和终止页面，计算内存区域包含的页面数。

计算的方法见 **3589-3590** 行：用内存区域与节点终止页面较小的一个减去 内存区域与节点起始页面较大的一个 得到内存区域总的页面数。

```

//计算一个节点中 hole 的数目

3659 /* Return the number of page frames in holes in a zone on a node */
3660 static unsigned long __meminit zone_absent_pages_in_node(int nid,
3661     unsigned long zone_type,
3662     unsigned long *ignored)
3663 {
3664     unsigned long node_start_pfn, node_end_pfn;
3665     unsigned long zone_start_pfn, zone_end_pfn;
3666
3667     get_pfn_range_for_nid(nid, &node_start_pfn, &node_end_pfn);
3668     zone_start_pfn = max(arch_zone_lowest_possible_pfn[zone_type],

```

```

3669             node_start_pfn);
3670     zone_end_pfn = min(arch_zone_highest_possible_pfn[zone_type],
3671             node_end_pfn);
3672
3673     adjust_zone_range_for_zone_movable(nid, zone_type,
3674             node_start_pfn, node_end_pfn,
3675             &zone_start_pfn, &zone_end_pfn);
3676     return __absent_pages_in_range(nid, zone_start_pfn, zone_end_pfn);
3677 }

```

3668-3671 行：得到区域起始和终止页面号。

3676:行：调用__absent_pages_in_range 计算起始和终止页面中是空洞的页面数目。

```

3596 /*
3597  * Return the number of holes in a range on a node. If nid is MAX_NUMNODES,
3598  * then all holes in the requested range will be accounted for.
3599  */
3600 unsigned long __meminit __absent_pages_in_range(int nid,
3601         unsigned long range_start_pfn,
3602         unsigned long range_end_pfn)
3603 {
3604     int i = 0;
3605     unsigned long prev_end_pfn = 0, hole_pages = 0;
3606     unsigned long start_pfn;
3607
3608     /* Find the end_pfn of the first active range of pfns in the node */
3609     i = first_active_region_index_in_nid(nid);
3610     if (i == -1)
3611         return 0;
3612
3613     prev_end_pfn = min(early_node_map[i].start_pfn, range_end_pfn);
3614
3615     /* Account for ranges before physical memory on this node */

```

```

3616     if (early_node_map[i].start_pfn > range_start_pfn)
3617         hole_pages = prev_end_pfn - range_start_pfn;
3618
3619     /* Find all holes for the zone within the node */
3620     for (; i != -1; i = next_active_region_index_in_nid(i, nid)) {
3621
3622         /* No need to continue if prev_end_pfn is outside the zone */
3623         if (prev_end_pfn >= range_end_pfn)
3624             break;
3625
3626         /* Make sure the end of the zone is not within the hole */
3627         start_pfn = min(early_node_map[i].start_pfn, range_end_pfn);
3628         prev_end_pfn = max(prev_end_pfn, range_start_pfn);
3629
3630         /* Update the hole size count and move on */
3631         if (start_pfn > range_start_pfn) {
3632             BUG_ON(prev_end_pfn > start_pfn);
3633             hole_pages += start_pfn - prev_end_pfn;
3634         }
3635         prev_end_pfn = early_node_map[i].end_pfn;
3636     }
3637
3638     /* Account for ranges past physical memory on this node */
3639     if (range_end_pfn > prev_end_pfn)
3640         hole_pages += range_end_pfn -
3641             max(range_start_pfn, prev_end_pfn);
3642
3643     return hole_pages;
3644 }

```

该函数通过计算[range_start_pfn, range_end_pfn]中包含的 early_node_map 内存区数目，如果这些内存区中存在间隙，则[range_start_pfn, range_end_pfn]存在空洞，并且可以计算空洞的数目。

//为节点中的叶框分配 page 数据结构所需的 连续内存。

```
3891 static void __init_refok alloc_node_mem_map(struct pglist_data *pgdat)
3892 {
3893     /* Skip empty nodes */
3894     if (!pgdat->node_spanned_pages)
3895         return;
3896
3897 #ifdef CONFIG_FLAT_NODE_MEM_MAP
3898     /* ia64 gets its own node_mem_map, before this, without bootmem */
3899     if (!pgdat->node_mem_map) {
3900         unsigned long size, start, end;
3901         struct page *map;
3902
3903         /*
3904          * The zone's endpoints aren't required to be MAX_ORDER
3905          * aligned but the node_mem_map endpoints must be in order
3906          * for the buddy allocator to function correctly.
3907          */
3908         start = pgdat->node_start_pfn & ~(MAX_ORDER_NR_PAGES - 1);
3909         end = pgdat->node_start_pfn + pgdat->node_spanned_pages;
3910         end = ALIGN(end, MAX_ORDER_NR_PAGES);
3911         size = (end - start) * sizeof(struct page);
3912         map = alloc_remap(pgdat->node_id, size);
3913         if (!map)
3914             map = alloc_bootmem_node(pgdat, size);
3915         pgdat->node_mem_map = map + (pgdat->node_start_pfn - start);
3916     }
3917 #endif CONFIG_NEED_MULTIPLE_NODES
```



```

3918    /*
3919    * With no DISCONTIG, the global mem_map is just set as node 0's
3920    */
3921    if (pgdat == NODE_DATA(0)) {
3922        mem_map = NODE_DATA(0)->node_mem_map;
3923    #ifdef CONFIG_ARCH_POPULATES_NODE_MAP
3924        if (page_to_pfn(mem_map) != pgdat->node_start_pfn)
3925            mem_map -= (pgdat->node_start_pfn - ARCH_PFN_OFFSET);
3926    #endif /* CONFIG_ARCH_POPULATES_NODE_MAP */
3927    }
3928 #endif
3929 #endif /* CONFIG_FLAT_NODE_MEM_MAP */
3930 }

```

该函数通过节点中页面数（`node_spanned_pages`，包含空洞页面）计算所需的 `page` 的数目，然后分配一块连续的物理内存来存放这些 `page`，并将 `page` 的起始地址保存在 `pgdat->node_mem_map` 中。

4.1.9.3.4 free_area_init_core

该函数初始化节点中的每一个内存区域,初始化的内容包括:

- 1、计算内存区域中包含的页面数，空洞页面数
- 2、设置内存区域中所有页面为保留
- 3、设置伙伴系统中所有的队列为空
- 4、清空内存位图

```

3793 /*
3794  * Set up the zone data structures:
3795  *   - mark all pages reserved
3796  *   - mark all memory queues empty
3797  *   - clear the memory bitmaps
3798  */
3799 static void __paginginit free_area_init_core(struct pglist_data *pgdat,
3800        unsigned long *zones_size, unsigned long *zhole_size)

```

```
3801 {
3802     enum zone_type j;
3803     int nid = pgdat->node_id;
3804     unsigned long zone_start_pfn = pgdat->node_start_pfn;
3805     int ret;
3806
3807     pgdat_resize_init(pgdat);
3808     pgdat->nr_zones = 0;
3809     init_waitqueue_head(&pgdat->kswapd_wait);
3810     pgdat->kswapd_max_order = 0;
3811     pgdat_page_cgroup_init(pgdat);
3812     //初始化节点这中的每个区
3813     for (j = 0; j < MAX_NR_ZONES; j++) {
3814         struct zone *zone = pgdat->node_zones + j;
3815         unsigned long size, realsize, memmap_pages;
3816         enum lru_list l;
3817
3818         size = zone_spanned_pages_in_node(nid, j, zones_size);
3819         realsize = size - zone_absent_pages_in_node(nid, j,
3820             zholes_size);
3821
3822         /*
3823          * Adjust realsize so that it accounts for how much memory
3824          * is used by this zone for memmap. This affects the watermark
3825          * and per-cpu initialisations
3826          */
3827         memmap_pages =
3828             PAGE_ALIGN(size * sizeof(struct page)) >> PAGE_SHIFT;
3829         if (realsize >= memmap_pages) {
3830             realsize -= memmap_pages;
3831             if (memmap_pages)
```

```

3832         printk(KERN_DEBUG
3833                 "   %s zone: %lu pages used for memmap\n",
3834                 zone_names[j], memmap_pages);
3835     } else
3836         printk(KERN_WARNING
3837                 "   %s zone: %lu pages exceeds realsize %lu\n",
3838                 zone_names[j], memmap_pages, realsize);
3839
3840     /* Account for reserved pages */
3841     if (j == 0 && realsize > dma_reserve) {
3842         realsize -= dma_reserve;
3843         printk(KERN_DEBUG "   %s zone: %lu pages reserved\n",
3844                 zone_names[0], dma_reserve);
3845     }
3846
3847     if (!is_highmem_idx(j))
3848         nr_kernel_pages += realsize;
3849     nr_all_pages += realsize;
3850
3851     zone->spanned_pages = size;
3852     zone->present_pages = realsize;

```

3818-3852 行： 计算内存区域中页面数目，空洞页面数目。

```

3853 #ifdef CONFIG_NUMA
3854     zone->node = nid;
3855     zone->min_unmapped_pages = (realsize*sysctl_min_unmapped_ratio)
3856                               / 100;
3857     zone->min_slab_pages = (realsize * sysctl_min_slab_ratio) / 100;
3858 #endif
3859     zone->name = zone_names[j];
3860     spin_lock_init(&zone->lock);
3861     spin_lock_init(&zone->lru_lock);

```

```

3862     zone_seqlock_init(zone);
3863     zone->zone_pgdat = pgdat;
3864
3865     zone->prev_priority = DEF_PRIORITY;
3866
3867     zone_pcp_init(zone);
3868     for_each_lru(l) {
3869         INIT_LIST_HEAD(&zone->lru[l].list);
3870         zone->reclaim_stat.nr_saved_scan[l] = 0;
3871     }
3872     zone->reclaim_stat.recent_rotated[0] = 0;
3873     zone->reclaim_stat.recent_rotated[1] = 0;
3874     zone->reclaim_stat.recent_scanned[0] = 0;
3875     zone->reclaim_stat.recent_scanned[1] = 0;
3876     zap_zone_vm_stats(zone);
3877     zone->flags = 0;
3878     if (!size)
3879         continue;
3880
3881     set_pageblock_order(pageblock_default_order());

```

3853-3881 行： 初始化内存区域 **zone** 中的各个变量。

```

3882     setup_usemap(pgdat, zone, size);

```

3882 行： 清空内存位图

```

3883     ret = init_currently_empty_zone(zone, zone_start_pfn,
3884                                     size, MEMMAP_EARLY);
3885     BUG_ON(ret);
3886     memmap_init(size, nid, j, zone_start_pfn);
3887     zone_start_pfn += size;
3888 }
3889 }

```

3883 行： 初始化伙伴系统

3886 行：设置页面为保留

```
//初始化和伙伴系统有关的数据结构
3260 __meminit int init_currently_empty_zone(struct zone *zone,
3261                                     unsigned long zone_start_pfn,
3262                                     unsigned long size,
3263                                     enum memmap_context context)
3264 {
3265     struct pglist_data *pgdat = zone->zone_pgdat;
3266     int ret;
3267     ret = zone_wait_table_init(zone, size);
3268     if (ret)
3269         return ret;
3270     pgdat->nr_zones = zone_idx(zone) + 1;
3271
3272     zone->zone_start_pfn = zone_start_pfn;
3273
3274     mminit_dprintk(MMINIT_TRACE, "memmap_init",
3275                   "Initialising map node %d zone %lu pfns %lu ->
%lu\n",
3276                   pgdat->node_id,
3277                   (unsigned long)zone_idx(zone),
3278                   zone_start_pfn, (zone_start_pfn + size));
3279
3280     zone_init_free_lists(zone);
3281
3282     return 0;
3283 }
```

3280 行：调用 zone_init_free_lists 函数初始化伙伴系统数据结构 free_list。

```
3053 static void __meminit zone_init_free_lists(struct zone *zone)
3054 {
```

```
3055     int order, t;
3056     for_each_migratetype_order(order, t) {
3057         INIT_LIST_HEAD(&zone->free_area[order].free_list[t]);
3058         zone->free_area[order].nr_free = 0;
3059     }
3060 }
```

```
3062 #ifndef __HAVE_ARCH_MEMMAP_INIT
3063 #define memmap_init(size, nid, zone, start_pfn) \
3064     memmap_init_zone((size), (nid), (zone), (start_pfn), MEMMAP_EARLY)
3065 #endif
```

```
2990 /*
2991  * Initially all pages are reserved - free ones are freed
2992  * up by free_all_bootmem() once the early boot process is
2993  * done. Non-atomic initialization, single-pass.
2994  */
2995 void __meminit memmap_init_zone(unsigned long size, int nid, unsigned long zone,
2996     unsigned long start_pfn, enum memmap_context context)
2997 {
2998     struct page *page;
2999     unsigned long end_pfn = start_pfn + size;
3000     unsigned long pfn;
3001     struct zone *z;
3002
3003     if (highest_memmap_pfn < end_pfn - 1)
3004         highest_memmap_pfn = end_pfn - 1;
3005
3006     z = &NODE_DATA(nid)->node_zones[zone];
3007     for (pfn = start_pfn; pfn < end_pfn; pfn++) {
3008         /*
```

```

3009      * There can be holes in boot-time mem_map[]s
3010      * handed to this function.  They do not
3011      * exist on hotplugged memory.
3012      */
3013      if (context == MEMMAP_EARLY) {
3014          if (!early_pfn_valid(pfn))
3015              continue;
3016          if (!early_pfn_in_nid(pfn, nid))
3017              continue;
3018      }
3019      page = pfn_to_page(pfn);
3020      set_page_links(page, zone, nid, pfn);
3021      mminit_verify_page_links(page, zone, nid, pfn);
3022      init_page_count(page);
3023      reset_page_mapcount(page);
3024      SetPageReserved(page);

```

设置页面为保留。

```

3025      /*
3026      * Mark the block movable so that blocks are reserved for
3027      * movable at startup. This will force kernel allocations
3028      * to reserve their blocks rather than leaking throughout
3029      * the address space during boot when many long-lived
3030      * kernel allocations are made. Later some blocks near
3031      * the start are marked MIGRATE_RESERVE by
3032      * setup_zone_migrate_reserve()
3033      *
3034      * bitmap is created for zone's valid pfn range. but memmap
3035      * can be created for invalid pages (for alignment)
3036      * check here not to call set_pageblock_migratetype() against
3037      * pfn out of zone.
3038      */

```

```

3039         if ((z->zone_start_pfn <= pfn)
3040             && (pfn < z->zone_start_pfn + z->spanned_pages)
3041             && !(pfn & (pageblock_nr_pages - 1)))
3042             set_pageblock_migratetype(page, MIGRATE_MOVABLE);
3043
3044         INIT_LIST_HEAD(&page->lru);
3045 #ifdef WANT_PAGE_VIRTUAL
3046         /* The shift won't overflow because ZONE_NORMAL is below 4G. */
3047         if (!lis_highmem_idx(zone))
3048             set_page_address(page, __va(pfn << PAGE_SHIFT));
3049 #endif
3050     }
3051 }

```

回到 `start_kernel` 中，调用 `build_all_zonelists` 建立节点备用列表。

5. 节点备用列表初始化

5.1 `build_all_zonelists`

初始化节点备用分配内存区域分配顺序

```

2815 void build_all_zonelists(void)
2816 {
2817     set_zonelist_order();
2818
2819     if (system_state == SYSTEM_BOOTING) {
2820         __build_all_zonelists(NULL);
2821         mminit_verify_zonelist();
2822         cpuset_init_current_mems_allowed();
2823     } else {
2824         /* we have to stop all cpus to guarantee there is no user

```



```

2825         of zonelist */
2826         stop_machine(__build_all_zonelists, NULL, NULL);
2827         /* cpuset refresh routine should be here */
2828     }
2829     vm_total_pages = nr_free_pagecache_pages();
2830     /*
2831      * Disable grouping by mobility if the number of pages in the
2832      * system is too low to allow the mechanism to work. It would be
2833      * more accurate, but expensive to check per-zone. This check is
2834      * made on memory-hotadd so a system can start with mobility
2835      * disabled and enable it later
2836      */
2837     if (vm_total_pages < (pageblock_nr_pages * MIGRATE_TYPES))
2838         page_group_by_mobility_disabled = 1;
2839     else
2840         page_group_by_mobility_disabled = 0;
2841
2842     printk("Built %i zonelists in %s order, mobility grouping %s.  "
2843           "Total pages: %ld\n",
2844           nr_online_nodes,
2845           zonelist_order_name[current_zonelist_order],
2846           page_group_by_mobility_disabled ? "off" : "on",
2847           vm_total_pages);
2848 #ifdef CONFIG_NUMA
2849     printk("Policy zone: %s\n", zone_names[policy_zone]);
2850 #endif
2851 }

```

该函数调用__build_all_zonelists 来完成大部分的工作。

5.2 __build_all_zonelists

```
2780 /* return values int ....just for stop_machine() */
2781 static int __build_all_zonelists(void *dummy)
2782 {
2783     int nid;
2784     int cpu;
2785
2786 #ifdef CONFIG_NUMA
2787     memset(node_load, 0, sizeof(node_load));
2788 #endif
2789     for_each_online_node(nid) {
2790         pg_data_t *pgdat = NODE_DATA(nid);
2791
2792         build_zonelists(pgdat);
2793         build_zonelist_cache(pgdat);
2794     }
2795
2796     /*
2797      * Initialize the boot_pagesets that are going to be used
2798      * for bootstrapping processors. The real pagesets for
2799      * each zone will be allocated later when the per cpu
2800      * allocator is available.
2801      *
2802      * boot_pagesets are used also for bootstrapping offline
2803      * cpus if the system is already booted because the pagesets
2804      * are needed to initialize allocators on a specific cpu too.
2805      * F.e. the percpu allocator needs the page allocator which
2806      * needs the percpu allocator in order to allocate its pagesets
2807      * (a chicken-egg dilemma).
2808      */
```

```
2809     for_each_possible_cpu(cpu)
2810         setup_pageset(&per_cpu(boot_pageset, cpu), 0);
2811
2812     return 0;
2813 }
```

2789-2794 行：对每个节点，调用 `build_zonelists` 函数建立备用列表。

5.3 build_zonelists

```
2637 static void build_zonelists(pg_data_t *pgdat)
2638 {
2639     int j, node, load;
2640     enum zone_type i;
2641     nodemask_t used_mask;
2642     int local_node, prev_node;
2643     struct zonelist *zonelist;
2644     int order = current_zonelist_order;
2645
2646     /* initialize zonelists */
2647     for (i = 0; i < MAX_ZONELISTS; i++) {
2648         zonelist = pgdat->node_zonelists + i;
2649         zonelist->_zonerefs[0].zone = NULL;
2650         zonelist->_zonerefs[0].zone_idx = 0;
2651     }
2652
2653     /* NUMA-aware ordering of nodes */
2654     local_node = pgdat->node_id;
2655     load = nr_online_nodes;
2656     prev_node = local_node;
2657     nodes_clear(used_mask);
2658 }
```

```
2659     memset(node_order, 0, sizeof(node_order));
2660     j = 0;
2661
2662     while ((node = find_next_best_node(local_node, &used_mask)) >= 0) {
2663         int distance = node_distance(local_node, node);
2664
2665         /*
2666          * If another node is sufficiently far away then it is better
2667          * to reclaim pages in a zone before going off node.
2668          */
2669         if (distance > RECLAIM_DISTANCE)
2670             zone_reclaim_mode = 1;
2671
2672         /*
2673          * We don't want to pressure a particular node.
2674          * So adding penalty to the first node in same
2675          * distance group to make it round-robin.
2676          */
2677         if (distance != node_distance(local_node, prev_node))
2678             node_load[node] = load;
2679
2680         prev_node = node;
2681         load--;
2682         if (order == ZONELIST_ORDER_NODE)
2683             build_zonelists_in_node_order(pgdat, node);
2684         else
2685             node_order[j++] = node; /* remember order */
2686     }
2687
2688     if (order == ZONELIST_ORDER_ZONE) {
2689         /* calculate node order -- i.e., DMA last! */
```

```
2690         build_zonelists_in_zone_order(pgdat, j);
2691     }
2692
2693     build_thisnode_zonelists(pgdat);
2694 }
```

5.4 build_zonelists_in_node_order

```
2510 /*
2511  * Build zonelists ordered by node and zones within node.
2512  * This results in maximum locality--normal zone overflows into
local
2513  * DMA zone, if any--but risks exhausting DMA zone.
2514  */
2515 static void build_zonelists_in_node_order(pg_data_t *pgdat, int
node)
2516 {
2517     int j;
2518     struct zonelist *zonelist;
2519
2520     zonelist = &pgdat->node_zonelists[0];
2521     for (j = 0; zonelist->_zonerefs[j].zone != NULL; j++)
2522         ;
2523     j = build_zonelists_node(NODE_DATA(node), zonelist, j,
2524                             MAX_NR_ZONES - 1);
2525     zonelist->_zonerefs[j].zone = NULL;
2526     zonelist->_zonerefs[j].zone_idx = 0;
2527 }
```

6. 利用early_res分配内存

在 slab 分配器还不存在的时候，只有前面刚刚建立好的初始化阶段的内存管理体系，如何分配一个内存空间来存放 pid 散列表。而 pid 散列表这么一个东西，它跟物理内存大小有关，也就是 1GB 的内存空间就可能有 16~4096 个散列项。不管怎样，调用 alloc_large_system_hash。

```
501 void __init pidhash_init(void)
502 {
503     int i, pidhash_size;
504
505     pid_hash = alloc_large_system_hash("PID", sizeof(*pid_hash),
0, 18,
506                                     HASH_EARLY | HASH_SMALL,
507                                     &pidhash_shift, NULL, 4096);
508     pidhash_size = 1 << pidhash_shift;
509
510     for (i = 0; i < pidhash_size; i++)
511         INIT_HLIST_HEAD(&pid_hash[i]);
512 }
```

在这里重点分析 alloc_large_system_hash 函数。

```
4857 /*
4858  * allocate a large system hash table from bootmem
4859  * - it is assumed that the hash table must contain an exact power-of-2
4860  *   quantity of entries
4861  * - limit is the number of hash buckets, not the total allocation size
4862  */
4863 void *__init alloc_large_system_hash(const char *tablename,
4864                                     unsigned long bucketsize,
4865                                     unsigned long numentries,
4866                                     int scale,
4867                                     int flags,
```

```
4868             unsigned int *_hash_shift,
4869             unsigned int *_hash_mask,
4870             unsigned long limit)
4871 {
4872     unsigned long long max = limit;
4873     unsigned long log2qty, size;
4874     void *table = NULL;
4875
4876     /* allow the kernel cmdline to have a say */
4877     if (!numentries) {
4878         /* round applicable memory size up to nearest megabyte */
4879         numentries = nr_kernel_pages;
4880         numentries += (1UL << (20 - PAGE_SHIFT)) - 1;
4881         numentries >>= 20 - PAGE_SHIFT;
4882         numentries <= 20 - PAGE_SHIFT;
4883
4884         /* limit to 1 bucket per 2^scale bytes of low memory */
4885         if (scale > PAGE_SHIFT)
4886             numentries >>= (scale - PAGE_SHIFT);
4887         else
4888             numentries <= (PAGE_SHIFT - scale);
4889
4890         /* Make sure we've got at least a 0-order allocation.. */
4891         if (unlikely(flags & HASH_SMALL)) {
4892             /* Makes no sense without HASH_EARLY */
4893             WARN_ON(!(flags & HASH_EARLY));
4894             if (!(numentries >> *_hash_shift)) {
4895                 numentries = 1UL << *_hash_shift;
4896                 BUG_ON(!numentries);
4897             }
4898         } else if (unlikely((numentries * bucketsize) < PAGE_SIZE))
```

```
4899         numentries = PAGE_SIZE / bucketsize;
4900     }
4901     numentries = roundup_pow_of_two(numentries);
4902
4903     /* limit allocation size to 1/16 total memory by default */
4904     if (max == 0) {
4905         max = ((unsigned long long)nr_all_pages << PAGE_SHIFT) >> 4;
4906         do_div(max, bucketsize);
4907     }
4908
4909     if (numentries > max)
4910         numentries = max;
4911
4912     log2qty = ilog2(numentries);
4913
4914     do {
4915         size = bucketsize << log2qty;
4916         if (flags & HASH_EARLY)
4917             table = alloc_bootmem_nopanic(size);
4918         else if (hashdist)
4919             table = __vmalloc(size, GFP_ATOMIC, PAGE_KERNEL);
4920         else {
4921             /*
4922              * If bucketsize is not a power-of-two, we may free
4923              * some pages at the end of hash table which
4924              * alloc_pages_exact() automatically does
4925              */
4926             if (get_order(size) < MAX_ORDER) {
4927                 table = alloc_pages_exact(size, GFP_ATOMIC);
4928                 kmemleak_alloc(table, size, 1, GFP_ATOMIC);
4929             }
```



```

4930     }
4931 } while (!table && size > PAGE_SIZE && --log2qty);
4932
4933 if (!table)
4934     panic("Failed to allocate %s hash table\n", tablename);
4935
4936 printk(KERN_INFO "%s hash table entries: %d (order: %d, %lu bytes)\n",
4937         tablename,
4938         (1U << log2qty),
4939         ilog2(size) - PAGE_SHIFT,
4940         size);
4941
4942 if (_hash_shift)
4943     *_hash_shift = log2qty;
4944 if (_hash_mask)
4945     *_hash_mask = (1 << log2qty) - 1;
4946
4947 return table;
4948}

```

该函数具体内容不细看了，主要查看分配内存的函数。

`alloc_bootmem_nopanic` 函数最终会调用 `__alloc_bootmem_nopanic` 函数分配内存。

```

681 static void * __init __alloc_bootmem_nopanic(unsigned long size,
682         unsigned long align,
683         unsigned long goal,
684         unsigned long limit)
685 {
686 #ifdef CONFIG_NO_BOOTMEM
687     void *ptr;
688
689     if (WARN_ON_ONCE(slab_is_available()))
690         return kzalloc(size, GFP_NOWAIT);

```

```
691
692 restart:
693
694     ptr = __alloc_memory_core_early(MAX_NUMNODES, size,
align, goal, limit);
695         ⌀] 990L, 24400C
696     if (ptr)
697         return ptr;
698
699     if (goal != 0) {
700         goal = 0;
701         goto restart;
702     }
703
704     return NULL;
705 #else
706     bootmem_data_t *bdata;
707     void *region;
708
709 restart:
710     region = alloc_arch_preferred_bootmem(NULL, size, align, goal,
limit);
711     if (region)
712         return region;
713
714     list_for_each_entry(bdata, &bdata_list, list) {
715         if (goal && bdata->node_low_pfn <= PFN_DOWN(goal))
716             continue;
717         if (limit && bdata->node_min_pfn >= PFN_DOWN(limit))
718             break;
719
720         region = alloc_bootmem_core(bdata, size, align, goal, limit);
```

```

721         if (region)
722             return region;
723     }
724
725     if (goal) {
726         goal = 0;
727         goto restart;
728     }
729
730     return NULL;
731 #endif
732 }
733

```

6.1 __alloc_memory_core_early

```

3411 #ifdef CONFIG_NO_BOOTMEM
3412 void * __init __alloc_memory_core_early(int nid, u64 size, u64 align,
3413                                         u64 goal, u64 limit)
3414 {
3415     int i;
3416     void *ptr;
3417
3418     /* need to go over early_node_map to find out good range for node */
3419     for_each_active_range_index_in_nid(i, nid) {
3420         u64 addr;
3421         u64 ei_start, ei_last;
3422
3423         ei_last = early_node_map[i].end_pfn;
3424         ei_last <=<= PAGE_SHIFT;
3425         ei_start = early_node_map[i].start_pfn;

```

```

3426         ei_start <=<= PAGE_SHIFT;
3427         addr = find_early_area(ei_start, ei_last,
3428                                goal, limit, size, align);
3429
3430         if (addr == -1ULL)
3431             continue;
3432
3433 #if 0
3434         printk(KERN_DEBUG "alloc (nid=%d %llx - %llx) (%llx - %llx) %llx
%llx => %llx\n",
3435                nid,
3436                ei_start, ei_last, goal, limit, size,
3437                align, addr);
3438 #endif
3439
3440         ptr = phys_to_virt(addr);
3441         memset(ptr, 0, size);
3442         reserve_early_without_check(addr, addr + size, "BOOTMEM");
3443         return ptr;
3444     }
3445
3446     return NULL;
3447 }
3448 #endif

```

该函数调用 `find_early_area` 来分配函数，`find_early_area` 在之前的章节中已经分析过。

至此，针对初始化期间的内存管理就全介绍完了，不过我还有一点要说。初始化期间为各个数据结构分配物理页面都是通过 `_alloc_memory_core_early` 函数。内核发展到现在，已经抛弃了以前的 **BOOTMEM** 分配体系，而只是通过


```
2299         dhash_entries,
2300         13,
2301         HASH_EARLY,
2302         &d_hash_shift,
2303         &d_hash_mask,
2304         0);
2305
2306     for (loop = 0; loop < (1 << d_hash_shift); loop++)
2307         INIT_HLIST_HEAD(&dentry_hashtable[loop]);
2308 }
```

```
1536 /*
1537  * Initialize the waitqueues and inode hash table.
1538  */
1539 void __init inode_init_early(void)
1540 {
1541     int loop;
1542
1543     /* If hashes are distributed across NUMA nodes, defer
1544      * hash allocation until vmalloc space is available.
1545      */
1546     if (hashdist)
1547         return;
1548
1549     inode_hashtable =
1550         alloc_large_system_hash("Inode-cache",
1551                                 sizeof(struct hlist_head),
1552                                 ihash_entries,
1553                                 14,
1554                                 HASH_EARLY,
1555                                 &i_hash_shift,
```

```
1556             &i_hash_mask,
1557             0);
1558
1559     for (loop = 0; loop < (1 << i_hash_shift); loop++)
1560         INIT_HLIST_HEAD(&inode_hashtable[loop]);
1561 }
```

8. 中断向量表初始化

中断向量表初始化，该过程就不详细分析了，在分析中断时再讨论。

```
882 void __init trap_init(void)
883 {
884     int i;
885
886 #ifdef CONFIG_EISA
887     void __iomem *p = early_ioremap(0xFFFFD9, 4);
888
889     if (readl(p) == 'E' + ('I' << 8) + ('S' << 16) + ('A' << 24))
890         EISA_bus = 1;
891     early_iounmap(p, 4);
892 #endif
893
894     set_intr_gate(0, &divide_error);
895     set_intr_gate_ist(1, &debug, DEBUG_STACK);
896     set_intr_gate_ist(2, &nmi, NMI_STACK);
897     /* int3 can be called from all */
898     set_system_intr_gate_ist(3, &int3, DEBUG_STACK);
899     /* int4 can be called from all */
900     set_system_intr_gate(4, &overflow);
```

```
901     set_intr_gate(5, &bounds);
902     set_intr_gate(6, &invalid_op);
903     set_intr_gate(7, &device_not_available);
904 #ifdef CONFIG_X86_32
905     set_task_gate(8, GDT_ENTRY_DOUBLEFAULT_TSS);
906 #else
907     set_intr_gate_ist(8, &double_fault, DOUBLEFAULT_STACK);
908 #endif
909     set_intr_gate(9, &coprocessor_segment_overrun);
910     set_intr_gate(10, &invalid_TSS);
911     set_intr_gate(11, &segment_not_present);
912     set_intr_gate_ist(12, &stack_segment, STACKFAULT_STACK);
913     set_intr_gate(13, &general_protection);
914     set_intr_gate(14, &page_fault);
915     set_intr_gate(15, &spurious_interrupt_bug);
916     set_intr_gate(16, &coprocessor_error);
917     set_intr_gate(17, &alignment_check);
918 #ifdef CONFIG_X86_MCE
919     set_intr_gate_ist(18, &machine_check, MCE_STACK);
920 #endif
921     set_intr_gate(19, &simd_coprocessor_error);
922
923     /* Reserve all the builtin and the syscall vector: */
924     for (i = 0; i < FIRST_EXTERNAL_VECTOR; i++)
925         set_bit(i, used_vectors);
926
927 #ifdef CONFIG_IA32_EMULATION
928     set_system_intr_gate(IA32_SYSCALL_VECTOR, ia32_syscall);
929     set_bit(IA32_SYSCALL_VECTOR, used_vectors);
930 #endif
931
```



```

932 #ifdef CONFIG_X86_32
933     if (cpu_has_fxsr) {
934         printk(KERN_INFO "Enabling fast FPU save and restore... ");
935         set_in_cr4(X86_CR4_OSFCSR);
936         printk("done.\n");
937     }
938     if (cpu_has_xmm) {
939         printk(KERN_INFO
940             "Enabling unmasked SIMD FPU exception support... ");
941         set_in_cr4(X86_CR4_OSXMMEXCPT);
942         printk("done.\n");
943     }
944 }
945 set_system_trap_gate(SYSCALL_VECTOR, &system_call);
946 set_bit(SYSCALL_VECTOR, used_vectors);
947 #endif
948
949 /*
950  * Should be a barrier for any external CPU state:
951  */
952 cpu_init();
953
954 x86_init.irqs.trap_init();
955 }

```

9. 内存管理初始化

9.1 mm_init函数全貌

该函数是内核初始化中最后和内存管理相关的函数。

- 1、mem_init 初始化伙伴系统。
- 2、kmem_cache_init 初始化 slab 管理数据结构
- 3、vmalloc_init:初始化高端内存管理中的非连续内存分配结构

```
512 /*
513  * Set up kernel memory allocators
514  */
515 static void __init mm_init(void)
516 {
517     /*
518      * page_cgroup requires countinuous pages as memmap
519      * and it's bigger than MAX_ORDER unless SPARSEMEM.
520      */
521     page_cgroup_init_flatmem();
522     mem_init();
523     kmem_cache_init();
524     pgtable_cache_init();
525     vmalloc_init();
526 }
```

9.2 mem_init伙伴系统的建立

```
867 void __init mem_init(void)
868 {
869     int codesize, reservedpages, datasize, initsize;
870     int tmp;
871
872     pci_iommu_alloc();
873
874 #ifdef CONFIG_FLATMEM
875     BUG_ON(!mem_map);
876 #endif
877     /* this will put all low memory onto the freelists */
878     totalram_pages += free_all_bootmem();
```

878 行：释放低端内存到伙伴系统。

```
879
880     reservedpages = 0;
```

```

881     for (tmp = 0; tmp < max_low_pfn; tmp++)
882         /*
883          * Only count reserved RAM pages:
884          */
885         if (page_is_ram(tmp) && PageReserved(pfn_to_page(tmp)))
886             reservedpages++;

```

881-886: 计算低端内存中保留页面的数目

```

887
888     set_highmem_pages_init();

```

888 行: 初始化高端内存到伙伴系统

```

889
890     codesize = (unsigned long) &_etext - (unsigned long) &_text;
891     datasize = (unsigned long) &_edata - (unsigned long) &_etext;
892     initsize = (unsigned long) &__init_end - (unsigned long) &__init_begin;
893
894     printk(KERN_INFO "Memory: %luk/%luk available (%dk kernel code, "
895                "%dk reserved, %dk data, %dk init, %ldk highmem)\n",
896            nr_free_pages() << (PAGE_SHIFT-10),
897            num_physpages << (PAGE_SHIFT-10),
898            codesize >> 10,
899            reservedpages << (PAGE_SHIFT-10),
900            datasize >> 10,
901            initsize >> 10,
902            totalhigh_pages << (PAGE_SHIFT-10));
903
904     printk(KERN_INFO "virtual kernel memory layout:\n"
905            "    fixmap   : 0x%08lx - 0x%08lx   (%4ld kB)\n"
906#ifdef CONFIG_HIGHMEM
907            "    pkmap    : 0x%08lx - 0x%08lx   (%4ld kB)\n"
908#endif
909            "    vmalloc : 0x%08lx - 0x%08lx   (%4ld MB)\n"

```

```

910     "    lowmem    : 0x%08lx - 0x%08lx    (%4ld MB)\n"
911     "    .init : 0x%08lx - 0x%08lx    (%4ld kB)\n"
912     "    .data : 0x%08lx - 0x%08lx    (%4ld kB)\n"
913     "    .text : 0x%08lx - 0x%08lx    (%4ld kB)\n",
914     FIXADDR_START, FIXADDR_TOP,
915     (FIXADDR_TOP - FIXADDR_START) >> 10,
916
917 #ifdef CONFIG_HIGHMEM
918     PKMAP_BASE, PKMAP_BASE+LAST_PKMAP*PAGE_SIZE,
919     (LAST_PKMAP*PAGE_SIZE) >> 10,
920 #endif
921
922     VMALLOC_START, VMALLOC_END,
923     (VMALLOC_END - VMALLOC_START) >> 20,
924
925     (unsigned long)__va(0), (unsigned long)high_memory,
926     ((unsigned long)high_memory - (unsigned long)__va(0)) >> 20,
927
928     (unsigned long)&__init_begin, (unsigned long)&__init_end,
929     ((unsigned long)&__init_end -
930      (unsigned long)&__init_begin) >> 10,
931
932     (unsigned long)&_etext, (unsigned long)&_edata,
933     ((unsigned long)&_edata - (unsigned long)&_etext) >> 10,
934
935     (unsigned long)&_text, (unsigned long)&_etext,
936     ((unsigned long)&_etext - (unsigned long)&_text) >> 10);
937
938     /*
939     * Check boundaries twice: Some fundamental inconsistencies can
940     * be detected at build time already.

```

```

941      */

942 #define __FIXADDR_TOP (-PAGE_SIZE)

943 #ifdef CONFIG_HIGHMEM

944     BUILD_BUG_ON(PKMAP_BASE + LAST_PKMAP*PAGE_SIZE > FIXADDR_START);

945     BUILD_BUG_ON(VMALLOC_END > PKMAP_BASE);

946 #endif

947 #define high_memory (-128UL << 20)

948     BUILD_BUG_ON(VMALLOC_START >= VMALLOC_END);

949 #undef high_memory

950 #undef __FIXADDR_TOP

951

952 #ifdef CONFIG_HIGHMEM

953     BUG_ON(PKMAP_BASE + LAST_PKMAP*PAGE_SIZE > FIXADDR_START);

954     BUG_ON(VMALLOC_END > PKMAP_BASE);

955 #endif

956     BUG_ON(VMALLOC_START >= VMALLOC_END);

957     BUG_ON((unsigned long)high_memory > VMALLOC_START);

958

959     if (boot_cpu_data.wp_works_ok < 0)

960         test_wp_bit();

```

889-960: 输出内存布局相关信息

```

961

962     save_pg_dir();

```

962:备份页全局变量 swapper_pg_dir

```

963     zap_low_mappings(true);

```

963 行:

```

964 }

```

9.2.1 低端内存释放free_all_bootmem

释放空闲页面到伙伴系统分配器中

```

Mm/bootmem.c
299 /**
300  * free_all_bootmem - release free pages to the buddy allocator
301  *
302  * Returns the number of pages actually released.
303  */
304 unsigned long __init free_all_bootmem(void)
305 {
306 #ifdef CONFIG_NO_BOOTMEM
307     /*
308      * We need to use MAX_NUMNODES instead of
309      * because in some case like Node0 doesnt have RAM
310      * low ram will be on Node1
311      * Use MAX_NUMNODES will make sure all ranges in
312      * will be used instead of only Node0 related
313      */
314     return free_all_memory_core_early(MAX_NUMNODES);
315 #else
316     unsigned long total_pages = 0;
317     bootmem_data_t *bdata;
318
319     list_for_each_entry(bdata, &bdata_list, list)
320         total_pages += free_all_bootmem_core(bdata);
321
322     return total_pages;
323 #endif
324 }

```

因为内核倾向不使用 **BOOTMEM**，所以在这分析
`free_all_memory_core_early(MAX_NUMNODES)`函数。

9.2.1.1 free_all_memory_core_early

该函数释放由 `early_res` 分配的内存。

```
200 unsigned long __init free_all_memory_core_early(int nodeid)
201 {
202     int i;
203     u64 start, end;
204     unsigned long count = 0;
205     struct range *range = NULL;
206     int nr_range;
207
208     nr_range = get_free_all_memory_range(&range, nodeid);
209
210     for (i = 0; i < nr_range; i++) {
211         start = range[i].start;
212         end = range[i].end;
213         count += end - start;
214         __free_pages_memory(start, end);
215     }
216
217     return count;
218 }
```

208 行：获得空闲 `page` 的区域。

210-215 行：根据 208 行获得的 `page` 区域，释放 `page` 到伙伴系统分配器。

9.2.1.2 get_free_all_memory_range

```
393 int __init get_free_all_memory_range(struct range **rangep, int
nodeid)
394 {
395     int i, count;
```

```

396     u64 start = 0, end;
397     u64 size;
398     u64 mem;
399     struct range *range;
400     int nr_range;
401
402     count = 0;
403     for (i = 0; i < max_early_res && early_res[i].end; i++)
404         count++;
405
406     count *= 2;
407
408     size = sizeof(struct range) * count;
409     end = get_max_mapped();//获得最大的物理页面号
410 #ifdef MAX_DMA32_PFN
411     if (end > (MAX_DMA32_PFN << PAGE_SHIFT))
412         start = MAX_DMA32_PFN << PAGE_SHIFT;
413 #endif
414     mem = find_fw_memmap_area(start, end, size, sizeof(struct
range));

```

从[start,end]分配一块连续的内存，存放 **range** 结构，用于存放空闲页面区间的数据结构。

Range 的类型如下：

```

struct range {
    u64    start;
    u64    end;
};

```

```

766 u64 __init find_fw_memmap_area(u64 start, u64 end, u64 size, u64
align)
767 {
768     return find_e820_area(start, end, size, align);
769 }

```



```

415     if (mem == -1ULL)
416         panic("can not find more space for range free");
417
418     range = __va(mem);
419     /* use early_node_map[] and early_res to get range array at first
*/
420     memset(range, 0, size);
421     nr_range = 0;

```

初始化 range 数组

```

422
423     /* need to go over early_node_map to find out good range for
node */
424     nr_range = add_from_early_node_map(range, count, nr_range,
nodeid);
425 #ifdef CONFIG_X86_32
426     subtract_range(range, count, max_low_pfn, -1ULL);

```

426 行：从 range 中去掉高端内存范围内的元素

```

427 #endif
428     subtract_early_res(range, count);

```

428 行：从 range 中去掉 early_res 中占用的页面区域

```

429     nr_range = clean_sort_range(range, count);
430
431     /* need to clear it ? */
432     if (nodeid == MAX_NUMNODES) {
433         memset(&early_res[0], 0,
434             sizeof(struct early_res) * max_early_res);
435         early_res = NULL;
436         max_early_res = 0;
437     }

```

431-437 行：清空 early_res 数组。

```
438
439     *rangep = range;
440     return nr_range;
441 }
```

9.2.1.3 add_from_early_node_map

```
3396 int __init add_from_early_node_map(struct range *range, int az,
3397                                     int nr_range, int nid)
3398 {
3399     int i;
3400     u64 start, end;
3401
3402     /* need to go over early_node_map to find out good range for
node */
3403     for_each_active_range_index_in_nid(i, nid) {
3404         start = early_node_map[i].start_pfn;
3405         end = early_node_map[i].end_pfn;
3406         nr_range = add_range(range, az, nr_range, start, end);
3407     }
3408     return nr_range;
3409 }
```

该函数将 `early_node_map` 中的页面添加到 `range` 数组中，不管页面是否保留（也就是非空闲的页面也会包含到 `range` 数组中。）

3404-3405 行：获得 `early_node_map` 索引 `i` 的起始和终止页面号

3406 行：将 `start,end` 添加到 `range` 数组中。

```
14 int add_range(struct range *range, int az, int nr_range, u64 start, u64
end)
15 {
16     if (start >= end)
17         return nr_range; //范围检查
```

```
18
19     /* Out of slots: */
20     if (nr_range >= az)
21         return nr_range;//超出 range 数组的范围。
22
23     range[nr_range].start = start;
24     range[nr_range].end = end;
25
26     nr_range++;
27
28     return nr_range;
29 }
```

9.2.1.4 subtract_range

```
359 static void __init subtract_early_res(struct range *range, int az)
360 {
361     int i, count;
362     u64 final_start, final_end;
363     int idx = 0;
364
365     count = 0;
366     for (i = 0; i < max_early_res && early_res[i].end; i++)
367         count++;
368
369     /* need to skip first one ?*/
370     if (early_res != early_res_x)
371         idx = 1;
372
373 #define DEBUG_PRINT_EARLY_RES 1
374
```

```

375 #if DEBUG_PRINT_EARLY_RES
376     printk(KERN_INFO "Subtract (%d early reservations)\n", count);
377 #endif
378     for (i = idx; i < count; i++) {
379         struct early_res *r = &early_res[i];
380 #if DEBUG_PRINT_EARLY_RES
381         printk(KERN_INFO "    #%d [%010llx - %010llx] %15s\n", i,
382             r->start, r->end, r->name);
383 #endif
384         final_start = PFN_DOWN(r->start);
385         final_end = PFN_UP(r->end);
386         if (final_start >= final_end)
387             continue;
388         subtract_range(range, az, final_start, final_end);

```

378-379 行：去掉 early_res 中占用的 early_res。

```

389     }

```

```

64 void subtract_range(struct range *range, int az, u64 start, u64 end)
65 {
66     int i, j;
67
68     if (start >= end)
69         return;
70
71     for (j = 0; j < az; j++) {
72         if (!range[j].end)
73             continue;
74
75         if (start <= range[j].start && end >= range[j].end) {
76             range[j].start = 0;
77             range[j].end = 0;

```

```
78         continue;
79     }
```

75-79 行: 如果 `range[i]` 落于 `[start,end]` 之间, 则去掉。

```
80
81     if (start <= range[j].start && end < range[j].end &&
82         range[j].start < end) {
83         range[j].start = end;
84         continue;
85     }
```

81-85 行: `range[i]` 与 `[start,end]` 有交集, 去掉交集部分

```
86
87
88     if (start > range[j].start && end >= range[j].end &&
89         range[j].end > start) {
90         range[j].end = start;
91         continue;
92     }
```

88-92 行: `range[i]` 与 `[start,end]` 有交集, 去掉交集部分

```
93
94     if (start > range[j].start && end < range[j].end) {
95         /* Find the new spare: */
96         for (i = 0; i < az; i++) {
97             if (range[i].end == 0)
98                 break;
99         }
100        if (i < az) {
101            range[i].end = range[j].end;
102            range[i].start = end;
103        } else {
104            printk(KERN_ERR "run of slot in ranges\n");
105        }
```

```
106         range[j].end = start;
107         continue;
108     }
```

94-108 行: range[i]包含[start,end], 则拆分 range。

```
109     }
110 }
```

9.2.1.5 __free_pages_memory(start, end)

获得空闲的 range 后, 在 free_all_memory_core_early 中调用 __free_pages_memory 释放内存到伙伴系统。

```
174 static void __init __free_pages_memory(unsigned long start,
unsigned long end)
175 {
176     int i;
177     unsigned long start_aligned, end_aligned;
178     int order = ilog2(BITS_PER_LONG);
179
180     start_aligned = (start + (BITS_PER_LONG - 1)) &
~(BITS_PER_LONG - 1);
181     end_aligned = end & ~(BITS_PER_LONG - 1);
182
183     if (end_aligned <= start_aligned) {
184         for (i = start; i < end; i++)
185             __free_pages_bootmem(pfn_to_page(i), 0);
186
187         return;
188     }
189
190     for (i = start; i < start_aligned; i++)
191         __free_pages_bootmem(pfn_to_page(i), 0);
```

```

192
193     for (i = start_aligned; i < end_aligned; i += BITS_PER_LONG)
194         __free_pages_bootmem(pfn_to_page(i), order);
195
196     for (i = end_aligned; i < end; i++)
197         __free_pages_bootmem(pfn_to_page(i), 0);
198 }

```

该函数调用__free_pages_bootmem 函数释放内存。

```

637 void __meminit __free_pages_bootmem(struct page *page, unsigned
int order)
638 {
639     if (order == 0) {
640         __ClearPageReserved(page);
641         set_page_count(page, 0);
642         set_page_refcounted(page);
643         __free_page(page);
644     } else {
645         int loop;
646
647         prefetchw(page);
648         for (loop = 0; loop < BITS_PER_LONG; loop++) {
649             struct page *p = &page[loop];
650
651             if (loop + 1 < BITS_PER_LONG)
652                 prefetchw(p + 1);
653             __ClearPageReserved(p);
654             set_page_count(p, 0);
655         }
656
657         set_page_refcounted(page);
658         __free_pages(page, order);

```

```
659     }
```

该函数调用__free_page 释放内存。__free_page 始释放页面的最终接口，待会在分析该函数。

9.2.2 高端内存释放set_highmem_pages_init

arch/x86/mm/ highmem_32.c

```
09 void __init set_highmem_pages_init(void)
110 {
111     struct zone *zone;
112     int nid;
113
114     for_each_zone(zone) {
115         unsigned long zone_start_pfn, zone_end_pfn;
116
117         if (!is_highmem(zone))
118             continue;
119
120         zone_start_pfn = zone->zone_start_pfn;
121         zone_end_pfn = zone_start_pfn + zone->spanned_pages;
122
123         nid = zone_to_nid(zone);
124         printk(KERN_INFO "Initializing %s for node %d
(%08lx:%08lx)\n",
125                zone->name, nid, zone_start_pfn, zone_end_pfn);
126
127         add_highpages_with_active_regions(nid, zone_start_pfn,
128                zone_end_pfn);
129     }
130     totalram_pages += totalhigh_pages;
131 }
```


117-118 行：检查区域是否属于高端内存。

127 行：调用 `add_highpages_with_active_regions` 函数进一步处理。

```
457 void __init add_highpages_with_active_regions(int nid, unsigned long
start_pfn,
458             unsigned long end_pfn)
459 {
460     struct add_highpages_data data;
461
462     data.start_pfn = start_pfn;
463     data.end_pfn = end_pfn;
464
465     work_with_active_regions(nid, add_highpages_work_fn,
&data);
466 }
```

`add_highpages_with_active_regions` 调用 `work_with_active_regions` 函数进一步处理，并传入函数 `add_highpages_work_fn` 作为参数。

```
451 void __init work_with_active_regions(int nid, work_fn_t work_fn, void
*data)
3452 {
3453     int i;
3454     int ret;
3455
3456     for_each_active_range_index_in_nid(i, nid) {
3457         ret = work_fn(early_node_map[i].start_pfn,
3458             early_node_map[i].end_pfn, data);
3459         if (ret)
3460             break;
3461     }
3462 }
```

3457-3458 行：调用 `add_highpages_work_fn` 进一步处理。

```

430 static int __init add_highpages_work_fn(unsigned long start_pfn,
431                                         unsigned long end_pfn, void *datax)
432 {
433     int node_pfn;
434     struct page *page;
435     unsigned long final_start_pfn, final_end_pfn;
436     struct add_highpages_data *data;
437
438     data = (struct add_highpages_data *)datax;
439
440     final_start_pfn = max(start_pfn, data->start_pfn);
441     final_end_pfn = min(end_pfn, data->end_pfn);
442     if (final_start_pfn >= final_end_pfn)
443         return 0;
444
445     for (node_pfn = final_start_pfn; node_pfn < final_end_pfn;
446         node_pfn++) {
447         if (!pfn_valid(node_pfn))
448             continue;
449         page = pfn_to_page(node_pfn);
450         add_one_highpage_init(page);
451     }
452
453     return 0;
454
455 }

```

调用 `add_one_highpage_init` 进一步处理。

```

417 static void __init add_one_highpage_init(struct page *page)

```

```
418 {  
419     ClearPageReserved(page);  
420     init_page_count(page);  
421     __free_page(page);  
422     totalhigh_pages++;  
423 }
```

419 行：去掉页面保留标志

420 行：初始化引用计数

421 行：调用__free_page 释放页面到伙伴系统。

9.2.3 __free_pages释放页面

```
2028 void __free_pages(struct page *page, unsigned int order)  
2029 {  
2030     if (put_page_testzero(page)) {  
2031         if (order == 0)  
2032             free_hot_cold_page(page, 0);  
2033         else  
2034             __free_pages_ok(page, order);  
2035     }  
2036 }
```

order=0，调用 free_hot_cold_page(page, 0);

9.2.3.1 free_hot_cold_page

```
1099  * Free a 0-order page  
1100  * cold == 1 ? free a cold page : free a hot page  
1101  */  
1102 void free_hot_cold_page(struct page *page, int cold)  
1103 {  
1104     struct zone *zone = page_zone(page);
```

```

1105     struct per_cpu_pages *pcp;
1106     unsigned long flags;
1107     int migratetype;
1108     int wasMlocked = __TestClearPageMlocked(page);
1109
1110     trace_mm_page_free_direct(page, 0);
1111     kmemcheck_free_shadow(page, 0);
1112
1113     if (PageAnon(page))
1114         page->mapping = NULL;
1115     if (free_pages_check(page))
1116         return;

```

1115-1116 行：检测 page 是否可以被释放。

- 1、page mapcount 必须等于-1
- 2、page _count 等于 0
- 3、page 标志位不能含有如下的标志：

```

/*
 * Flags checked when a page is freed.  Pages being freed should not
have
 * these flags set.  If they are, there is a problem.
 */
#define PAGE_FLAGS_CHECK_AT_FREE \
    (1 << PG_lru    | 1 << PG_locked    | \
    1 << PG_private | 1 << PG_private_2 | \
    1 << PG_buddy    | 1 << PG_writeback | 1 << PG_reserved | \
    1 << PG_slab    | 1 << PG_swapcache | 1 << PG_active | \
    1 << PG_unevictable | __PG_MLOCKED | __PG_HWPOISON)

```

```

524 static inline int free_pages_check(struct page *page)
525 {
526     if (unlikely(page_mapcount(page) |

```

```
527         (page->mapping != NULL) |
528         (atomic_read(&page->_count) != 0) |
529         (page->flags & PAGE_FLAGS_CHECK_AT_FREE))) {
530         bad_page(page);
531         return 1;
532     }
533     if (page->flags & PAGE_FLAGS_CHECK_AT_PREP)
534         page->flags &= ~PAGE_FLAGS_CHECK_AT_PREP;
535     return 0;
536 }
```

```
1117
1118     if (!PageHighMem(page)) {
1119         debug_check_no_locks_freed(page_address(page), PAGE_SIZE);
1120         debug_check_no_obj_freed(page_address(page), PAGE_SIZE);
1121     }
1122     arch_free_page(page, 0);
1123     kernel_map_pages(page, 1, 0);
1124
1125     migratetype = get_pageblock_migratetype(page);
1126     set_page_private(page, migratetype);
1127     local_irq_save(flags);
1128     if (unlikely(wasMlocked))
1129         free_page_mlock(page);
1130     __count_vm_event(PGFREE);
1131
1132     /*
1133     * We only track unmovable, reclaimable and movable on pcp lists.
1134     * Free ISOLATE pages back to the allocator because they are being
1135     * offlined but treat RESERVE as movable pages so we can get those
1136     * areas back if necessary. Otherwise, we may have to free
```

```

1137      * excessively into the page allocator
1138      */
1139      if (migratetype >= MIGRATE_PCPTYPES) {
1140          if (unlikely(migratetype == MIGRATE_ISOLATE)) {
1141              free_one_page(zone, page, 0, migratetype);
1142              goto out;
1143          }
1144          migratetype = MIGRATE_MOVABLE;
1145      }
1146
1147      pcp = &this_cpu_ptr(zone->pageset)->pcp;
1148      if (cold)
1149          list_add_tail(&page->lru, &pcp->lists[migratetype]);
1150      else
1151          list_add(&page->lru, &pcp->lists[migratetype]);
1152      pcp->count++;
1153      if (pcp->count >= pcp->high) {
1154          free_pcppages_bulk(zone, pcp->batch, pcp);
1155          pcp->count -= pcp->batch;
1156      }
1157
1158 out:
1159      local_irq_restore(flags);
1160 }

```

9.2.3.2 free_one_page

```

590 static void free_one_page(struct zone *zone, struct page *page, int
order,
591                          int migratetype)
592 {

```

```

593     spin_lock(&zone->lock);
594     zone->all_unreclaimable = 0;
595     zone->pages_scanned = 0;
596
597     __mod_zone_page_state(zone, NR_FREE_PAGES, 1 <<
order);
598     __free_one_page(page, zone, order, migratetype);
599     spin_unlock(&zone->lock);
600 }

```

9.2.3.3 __free_one_page 函数

```

449 /*
450  * Freeing function for a buddy system allocator.
451  *
452  * The concept of a buddy system is to maintain direct-mapped table
453  * (containing bit values) for memory blocks of various "orders".
454  * The bottom level table contains the map for the smallest allocatable
455  * units of memory (here, pages), and each level above it describes
456  * pairs of units from the levels below, hence, "buddies".
457  * At a high level, all that happens here is marking the table entry
458  * at the bottom level available, and propagating the changes upward
459  * as necessary, plus some accounting needed to play nicely with other
460  * parts of the VM system.
461  * At each level, we keep a list of pages, which are heads of continuous
462  * free pages of length of (1 << order) and marked with PG_buddy. Page's
463  * order is recorded in page_private(page) field.
464  * So when we are allocating or freeing one, we can derive the state of the
465  * other. That is, if we allocate a small block, and both were
466  * free, the remainder of the region must be split into blocks.
467  * If a block is freed, and its buddy is also free, then this

```

```

468  * triggers coalescing into a block of larger size.
469  *
470  * -- wli
471  */
472
473 static inline void __free_one_page(struct page *page,
474     struct zone *zone, unsigned int order,
475     int migratetype)
476 {
477     unsigned long page_idx;
478
479     if (unlikely(PageCompound(page)))
480         if (unlikely(destroy_compound_page(page, order)))
481             return;
482
483     VM_BUG_ON(migratetype == -1);
484
485     page_idx = page_to_pfn(page) & ((1 << MAX_ORDER) - 1);

```

485 行: `page_idx` 表示 `page` 在 `mem_map` 数组中的下标模 1023 的值。

```

486
487     VM_BUG_ON(page_idx & ((1 << order) - 1));
488     VM_BUG_ON(bad_range(zone, page));
489
490     while (order < MAX_ORDER-1) {
491         unsigned long combined_idx;
492         struct page *buddy;
493
494         buddy = __page_find_buddy(page, page_idx, order);

```

494 行: 寻找 `page` 块的伙伴。

406 static inline struct page *


```

407 __page_find_buddy(struct page *page, unsigned long page_idx,
unsigned int order)
408 {
409     unsigned long buddy_idx = page_idx ^ (1 << order);
410
411     return page + (buddy_idx - page_idx);
412 }

```

要看懂上述函数，必须清楚如下几点：

1、伙伴算法的块的分配的方法，大小为 2^x 个页框的块，他的起始地址一定是 2^x 的倍数

2、在把两个大小为 $2^{(x-1)}$ 的块合并为大小为 2^x 的块的时候，需要考虑目前回收的块是 **buddy** 里面的第一个还是第二个，判断的依据就是 **page_idx** 的 2^{order} 这一位，如果为 1，则必须向前寻找 **buddy**，如果为 0，则必须向后寻找 **buddy**，这也是异或的精妙之处，如果写成 **if else** 可能更好理解，但是有了奇技淫巧才显得 **linux** 源代码的高深。

3、当 **order** 固定时，一个 **page_idx** 在指数为 **order** 的 **free_list** 中只可能有一个伙伴。

```

495     if (!page_is_buddy(page, buddy, order))
496         break;

```

495-496 行：判断 **buddy** 是否为 **page** 的伙伴。

420-427 行的注释很清楚的说明了伙伴的条件。

- 1、**buddy** 不能再空洞中。
- 2、伙伴在伙伴系统中。
- 3、伙伴和 **page** 需要有相同的阶数
- 4、**page** 和 **buddy** 在相同的内存区域中。

```

420 /*
421  * This function checks whether a page is free && is the buddy
422  * we can do coalesce a page and its buddy if
423  * (a) the buddy is not in a hole &&
424  * (b) the buddy is in the buddy system &&
425  * (c) a page and its buddy have the same order &&
426  * (d) a page and its buddy are in the same zone.

```

```

427  *
428  * For recording whether a page is in the buddy system, we use
PG_buddy.
429  * Setting, clearing, and testing PG_buddy is serialized by
zone->lock.
430  *
431  * For recording page's order, we use page_private(page).
432  */
433 static inline int page_is_buddy(struct page *page, struct page
*buddy,
434                                int order)
435 {
436     if (!pfn_valid_within(page_to_pfn(buddy)))
437         return 0;
438
439     if (page_zone_id(page) != page_zone_id(buddy))
440         return 0;
441
442     if (PageBuddy(buddy) && page_order(buddy) == order) {
443         VM_BUG_ON(page_count(buddy) != 0);
444         return 1;
445     }
446     return 0;
447 }

```

```

497
498     /* Our buddy is free, merge with it and move up one order. */
499     list_del(&buddy->lru);
500     zone->free_area[order].nr_free--;

```

499-500 行：将伙伴冲阶 **order** 中移除，同时减少阶 **order** 中 **nr_free**。

```

501     rmv_page_order(buddy);

```

502	<code>combined_idx = __find_combined_index(page_idx, order);</code>
502 行：计算伙伴合并后新的 idx	
503	<code>page = page + (combined_idx - page_idx);</code>
504	<code>page_idx = combined_idx;</code>
505	<code>order++;</code>
503-505 行：计算伙伴的头 page，得到伙伴的 idx，order++，进行下一轮循环。	
506	<code>}</code>
507	<code>set_page_order(page, order);</code>
507 行：设置 page 的阶数	
508	<code>list_add(&page->lru,</code>
509	<code>&zone->free_area[order].free_list[migratetype]);</code>
510	<code>zone->free_area[order].nr_free++;</code>
508-510 行：将伙伴添加到对应的阶中，且该阶的 nr_free++。	
511	<code>}</code>

9.2.4 save_pg_dir

备份页全局目录，原因注释中讲的很清楚。

551	<code>#ifdef CONFIG_ACPI_SLEEP</code>
552	<code>/*</code>
553	<code>* ACPI suspend needs this for resume, because things like the</code>
intel-agp	
554	<code>* driver might have split up a kernel 4MB mapping.</code>
555	<code>*/</code>
556	<code>char swsusp_pg_dir[PAGE_SIZE]</code>
557	<code>__attribute__((aligned(PAGE_SIZE)));</code>
558	
559	<code>static inline void save_pg_dir(void)</code>
560	<code>{</code>
561	<code>memcpy(swsusp_pg_dir, swapper_pg_dir, PAGE_SIZE);</code>

9.2.5 zap_low_mappings

这个函数很简单，就是把前面我们在 `arch/x86/kernel/head_32.S` 中设置的页全局目录的前若干项清零。这若干项到底是多少项呢？我们看看 `KERNEL_PGD_BOUNDARY` 是什么东西：

```
#define KERNEL_PGD_BOUNDARY pgd_index(PAGE_OFFSET)
#define pgd_index(address) (((address) >> PGDIR_SHIFT) &
(PTRS_PER_PGD - 1))
#define PGDIR_SHIFT 22
#define PTRS_PER_PGD 1024
```

不错， $0xc0000000 \gg 22 \& 1023 = 768$ ，这些也全局目录项代表虚拟地址前 3G 的页面，也就是所谓的用户区，我们在这里把它全清零了。

```
569 void zap_low_mappings(bool early)
570 {
571     int i;
572
573     /*
574      * Zap initial low-memory mappings.
575      *
576      * Note that "pgd_clear()" doesn't do it for
577      * us, because pgd_clear() is a no-op on i386.
578      */
579     for (i = 0; i < KERNEL_PGD_BOUNDARY; i++) {
580 #ifdef CONFIG_X86_PAE
581         set_pgd(swapper_pg_dir+i, __pgd(1 +
__pa(empty_zero_page)));
582 #else
583         set_pgd(swapper_pg_dir+i, __pgd(0));
584 #endif
585     }
```

```
586
587     if (early)
588         __flush_tlb();
589     else
590         flush_tlb_all();
591 }
```

9.3 kmem_cache_init slab分配器初始化

9.3.1 slab分配器相关数据结构

在 `kmem_cache_init` 函数之前，内核使用 `early_res` 机制分配内存，但 `early_res` 毕竟是一个临时的过渡方案，只考虑分配的效率没有考虑内存管理中内存碎片的问题。上节的伙伴系统按页来分配内存，这个分配单位太大。现在是时候初始化 `slab` 分配器聊。

`slab` 分配器中的一个重要概念是缓存。一个缓存包含一系列的 `slab`，每一个 `slab` 由几个连续的物理页面组成，这些页面包含相同类型的对象。

在内核中，缓存的结构如下：

Include/linux/slab_def.h

9.3.1.1 kmem_cache数据结构

```
19 /*
20  * struct kmem_cache
21  *
22  * manages a cache.
23  */
24
25 struct kmem_cache {
26     /* 1) per-cpu data, touched during every alloc/free */
27     struct array_cache *array[NR_CPUS];
```

27:每 cpu 数据，每次分配/释放时访问，加快分配速度。

`array` 是一个指针数组，数组的长度是系统中 CPU 数目的个数。每个数组项指向一个 `array_cache` 示例，表示系统中的一个 CPU。

```
28 /* 2) Cache tunables. Protected by cache_chain_mutex */
29     unsigned int batchcount;
30     unsigned int limit;
31     unsigned int shared;
32
33     unsigned int buffer_size;
34     u32 reciprocal_buffer_size;
```

28-34 行：可调整的缓存参数

batchcount：在 per-CPU 列表为空时，从缓存的 **slab** 中获取对象的数目。它还表示在缓存增长时分配的对象数目。

limit：指定了 per-CPU 列表中保存的对象的最大数目。

```
35 /* 3) touched by every alloc & free from the backend */
36
37     unsigned int flags;        /* constant flags */
38     unsigned int num;         /* # of objs per slab */
39
```

37-38 行：常数标志。

```
40 /* 4) cache_grow/shrink */
41     /* order of pgs per slab (2^n) */
42     unsigned int gfporder;
43
44     /* force GFP flags, e.g. GFP_DMA */
45     gfp_t gfpflags;
46
47     size_t colour;            /* cache colouring range */
48     unsigned int colour_off;  /* colour offset */
49     struct kmem_cache *slabp_cache;
50     unsigned int slab_size;
51     unsigned int dflags;      /* dynamic flags */
```

```
52                                ⌘] 223L, 5300C
53    /* constructor func */
54    void (*ctor)(void *obj);
55
```

41-55 行: 缓存的增长/缩减变量

```
56 /* 5) cache creation/removal */
57    const char *name;
58    struct list_head next;
59
```

57-58 行: 缓存创建/删去相关变量

```
60 /* 6) statistics */

61 #ifdef CONFIG_DEBUG_SLAB
62    unsigned long num_active;
63    unsigned long num_allocations;
64    unsigned long high_mark;
65    unsigned long grown;
66    unsigned long reaped;
67    unsigned long errors;
68    unsigned long max_freeable;
69    unsigned long node_allocs;
70    unsigned long node_frees;
71    unsigned long node_overflow; ⌘] 223L, 5300C
72    atomic_t allochit;
73    atomic_t allocmiss;
74    atomic_t freehit;
75    atomic_t freemiss;
76
77    /*
78     * If debugging is enabled, then the allocator can add additional
79     * fields and/or padding to every object. buffer_size contains the total
```

```
80      * object size including these internal fields, the following two
81      * variables contain the offset to the user object and its size.
82      */
83      int obj_offset;
84      int obj_size;
85 #endif /* CONFIG_DEBUG_SLAB */
```

61-85 行：统计量

```
86
87      /*
88      * We put nodelists[] at the end of kmem_cache, because we want to size
89      * this array to nr_node_ids slots instead of MAX_NUMNODES
90      * (see kmem_cache_init())
91      * We still use [MAX_NUMNODES] and not [1] or [0] because cache_cache
92      * is statically defined, so we reserve the max number of nodes.
93      */
94      struct kmem_list3 *nodelists[MAX_NUMNODES];
```

94 行：slab 链表表头数据结构。Nodelists 是一个数组，每个数组项对应系统中一个可能的节点。

```
95      /*
96      * Do not add fields after nodelists[]
97      */
98 };
```

9.3.1.2 数据结构array_cache

```
255 /*
256  * struct array_cache
257  *
258  * Purpose:
259  * - LIFO ordering, to hand out cache-warm objects from _alloc
260  * - reduce the number of linked list operations
```



```

261  * - reduce spinlock operations
262  *
263  * The limit is stored in the per-cpu structure to reduce the data cache
264  * footprint.
265  *
266  */
267 struct array_cache {
268     unsigned int avail;
269     unsigned int limit;
270     unsigned int batchcount;
271     unsigned int touched;
272     spinlock_t lock;
273     void *entry[]; /*
274                     * Must have this definition in here for the proper
275                     * alignment of array_cache. Also simplifies accessing
276                     * the entries.
277                     */
278 };

```

avail:保存当前可用对象的数目。

entry: 一个伪数组，方便访问 array_cache 实例之后缓存中的各个对象。

9.3.1.3 kmem_list3 数据结构

```

290 /*
291  * The slab lists for all objects.
292  */
293 struct kmem_list3 {
294     struct list_head slabs_partial; /* partial list first, better asm code
*/
295     struct list_head slabs_full;
296     struct list_head slabs_free;

```

部分、满、空闲 slab 的链表头

```

297     unsigned long free_objects;
298     unsigned int free_limit;
299     unsigned int colour_next;    /* Per-node cache coloring */
300     spinlock_t list_lock;
301     struct array_cache *shared; /* shared per node */
302     struct array_cache **alien; /* on other nodes */
303     unsigned long next_reap;     /* updated without locking */
304     int free_touched;           /* updated without locking */
305 };

```

297-305 行：和页面回收的相关数据结构

9.3.2 kmem_cache_init 函数

为初始化 slab 数据结构，内核若干小于 1 页的内存，这些内存最适合使用 kmalloc 调用，但在 slab 分配器建立之前，是不能使用 kmalloc 分配内存的。

内核使用编译时创建的静态数据，为 slab 的初始化提供内存。

Kmem_cache_init 的主要过程是：

- 1、创建系统中的第一个 slab 缓存，为 kmem_cache 实例提供内存。第一个 slab 缓存的实例使用编译时的静态数据结构 initarray_cache 用作 per-CPU 数组，缓存名称为 cache_cache。
- 2、初始化一般性的的缓存，为使用用 kmalloc 的做好准备。
- 3、

```

1372  * Initialisation.  Called after the page allocator have been initialised and
1373  * before smp_init().
1374  */
1375 void __init kmem_cache_init(void)
1376 {
1377     size_t left_over;
1378     struct cache_sizes *sizes;
1379     struct cache_names *names;
1380     int i;
1381     int order;

```

```
1382     int node;
1383
1384     if (num_possible_nodes() == 1)
1385         use_alien_caches = 0;
1386
```

```
1387     for (i = 0; i < NUM_INIT_LISTS; i++) {
1388         kmem_list3_init(&initkmem_list3[i]);
1389         if (i < MAX_NUMNODES)
1390             cache_cache.nodelists[i] = NULL;
1391     }
```

初始化全局变量 initkmem_list3

```
307 /*
308  * Need this for bootstrapping a per node allocator.
309  */
310 #define NUM_INIT_LISTS (3 * MAX_NUMNODES)
311 struct kmem_list3 __initdata initkmem_list3[NUM_INIT_LISTS];
312 #define CACHE_CACHE 0
313 #define SIZE_AC MAX_NUMNODES
314 #define SIZE_L3 (2 * MAX_NUMNODES)
```

initkmem_list3 是一个数组，数组的每个元素为 kmem_list3,数组的个数为最大节点数的 3 倍。

```
1392     set_up_list3s(&cache_cache, CACHE_CACHE);
```

初始化全局变量 cache_cache.

```
1393
1394     /*
1395      * Fragmentation resistance on low memory - only use bigger
1396      * page orders on machines with more than 32MB of memory.
1397      */
1398     if (totalram_pages > (32 << 20) >> PAGE_SHIFT)
1399         slab_break_gfp_order = BREAK_GFP_ORDER_HI;
```

```

1400
1401     /* Bootstrap is tricky, because several objects are allocated
1402     * from caches that do not exist yet:
1403     * 1) initialize the cache_cache cache: it contains the struct
1404     *    kmem_cache structures of all caches, except cache_cache itself:
1405     *    cache_cache is statically allocated.
1406     *    Initially an __init data area is used for the head array and the
1407     *    kmem_list3 structures, it's replaced with a kmalloc allocated
1408     *    array at the end of the bootstrap.
1409     * 2) Create the first kmalloc cache.
1410     *    The struct kmem_cache for the new cache is allocated normally.
1411     *    An __init data area is used for the head array.
1412     * 3) Create the remaining kmalloc caches, with minimally sized
1413     *    head arrays.
1414     * 4) Replace the __init data head arrays for cache_cache and the first
1415     *    kmalloc cache with kmalloc allocated arrays.
1416     * 5) Replace the __init data for kmem_list3 for cache_cache and
1417     *    the other cache's with kmalloc allocated memory.
1418     * 6) Resize the head arrays of the kmalloc caches to their final sizes.
1419     */
1420
1421     node = numa_node_id();
1422
1423     /* 1) create the cache_cache */
1424     INIT_LIST_HEAD(&cache_chain);
1425     list_add(&cache_cache.next, &cache_chain);

```

将 `cache_cache` 连接到链表头 `cache_chain` 中。

```

1426     cache_cache.colour_off = cache_line_size();
1427     cache_cache.array[smp_processor_id()] = &initarray_cache.cache;

```

使用静态全局变量作为 `cache_cache` 缓存的 per-CPU 实例。

```

591 static struct arraycache_init initarray_cache __initdata =

```

```

592     { {0, BOOT_CPUCACHE_ENTRIES, 1, 0} };
593 static struct arraycache_init initarray_generic =
594     { {0, BOOT_CPUCACHE_ENTRIES, 1, 0} };
595
596 /* internal cache of cache description objs */
597 static struct kmem_cache cache_cache = {
598     .batchcount = 1,
599     .limit = BOOT_CPUCACHE_ENTRIES,
600     .shared = 1,
601     .buffer_size = sizeof(struct kmem_cache),
602     .name = "kmem_cache",
603 };

```

initarray_cache 是内型为 arraycache_init 结构体，它的元素包含一个 array_cache 的示例，以及一个指针数组。

```

284 #define BOOT_CPUCACHE_ENTRIES    1
285 struct arraycache_init {
286     struct array_cache cache;
287     void *entries[BOOT_CPUCACHE_ENTRIES];
288 };

```

```

1428     cache_cache.nodelists[node] = &initkmem_list3[CACHE_CACHE + node];
1429
1430     /*
1431      * struct kmem_cache size depends on nr_node_ids, which
1432      * can be less than MAX_NUMNODES.
1433      */
1434     cache_cache.buffer_size = offsetof(struct kmem_cache, nodelists) +
1435                               nr_node_ids * sizeof(struct kmem_list3 *);
1436 #if DEBUG
1437     cache_cache.obj_size = cache_cache.buffer_size;
1438 #endif

```

```

1439     cache_cache.buffer_size = ALIGN(cache_cache.buffer_size,
1440                                     cache_line_size());
1441     cache_cache.reciprocal_buffer_size =
1442         reciprocal_value(cache_cache.buffer_size);
1443
1444     for (order = 0; order < MAX_ORDER; order++) {
1445         cache_estimate(order, cache_cache.buffer_size,
1446                       cache_line_size(), 0, &left_over, &cache_cache.num);
1447         if (cache_cache.num)
1448             break;
1449     }
1450     BUG_ON(!cache_cache.num);
1451     cache_cache.gfporder = order;
1452     cache_cache.colour = left_over / cache_cache.colour_off;
1453     cache_cache.slab_size = ALIGN(cache_cache.num *
sizeof(kmem_bufctl_t) +
1454                                   sizeof(struct slab), cache_line_size());
1455

```

1426-1455 行：初始化 `cache_cache` 字段的相关变量

```

1456     /* 2+3) create the kmalloc caches */
1457     sizes = malloc_sizes;
1458     names = cache_names;

```

`Malloc_sizes`, `cache_names` 是两个指针数组，其中分别存放一般性缓存的大小和名称。

```

1459
1460     /*
1461      * Initialize the caches that provide memory for the array cache and the
1462      * kmem_list3 structures first.  Without this, further allocations will
1463      * bug.
1464      */
1465

```

```

1466     sizes[INDEX_AC].cs_cachep =
kmem_cache_create(names[INDEX_AC].name,
1467                 sizes[INDEX_AC].cs_size,
1468                 ARCH_KMALLOC_MINALIGN,
1469                 ARCH_KMALLOC_FLAGS|SLAB_PANIC,
1470                 NULL);

```

创建大小为(sizeof(struct arraycache_init)的普通缓存。

```

349 #define INDEX_AC index_of(sizeof(struct arraycache_init))
350 #define INDEX_L3 index_of(sizeof(struct kmem_list3))

```

```

1471
1472     if (INDEX_AC != INDEX_L3) {
1473         sizes[INDEX_L3].cs_cachep =
1474             kmem_cache_create(names[INDEX_L3].name,
1475                             sizes[INDEX_L3].cs_size,
1476                             ARCH_KMALLOC_MINALIGN,
1477                             ARCH_KMALLOC_FLAGS|SLAB_PANIC,
1478                             NULL);
1479     }

```

如果 `arraycache_init` 和 `kmem_list3` 大小相等，则共用同一普通缓存，否则创建大小为 `sizeof(struct kmem_list3)` 的缓存。

这两个普通缓存非常重要，因为它是下面其它缓存的基础。因为每个 `kmem_cache` 缓存都需要分配 `arraycache_init` 和 `kmem_list3` 大小的对象，而这些对象需要通过 `kmalloc` 来分配，而 `kmalloc` 的分配有是建立的上述两个普通缓存的基础上的。

这里有一个问题大家必须明白，即为分配 `arraycache_init` 和 `kmem_list3` 建立的缓存也是一个 `kmem_cache` 实例，它们也需要 `arraycache_init` 和 `kmem_list3`，所以也需要 `kmalloc` 的分配，但是 `kmalloc` 的分配在 `arraycache_init` 和 `kmem_list3` 的缓存还没建立起来是不可使用的，这里就存在一个鸡与蛋的问题，起始从上面的代码可以看出，内核的解决办法是：为分配 `arraycache_init` 和 `kmem_list3` 的而建立的缓存的 `arraycache_init` 和 `kmem_list3` 使用静态的全局变量。具体的实现就不深入进去了，等到讲解虚拟内存系统是在详细讨论。

```

1480
1481     slab_early_init = 0;
1482
1483     while (sizes->cs_size != ULONG_MAX) {
1484         /*
1485          * For performance, all the general caches are L1 aligned.
1486          * This should be particularly beneficial on SMP boxes, as it
1487          * eliminates "false sharing".
1488          * Note for systems short on memory removing the alignment will
1489          * allow tighter packing of the smaller caches.
1490          */
1491         if (!sizes->cs_cachep) {
1492             sizes->cs_cachep = kmem_cache_create(names->name,
1493                 sizes->cs_size,
1494                 ARCH_KMALLOC_MINALIGN,
1495                 ARCH_KMALLOC_FLAGS|SLAB_PANIC,
1496                 NULL);
1497         }
1498 #ifdef CONFIG_ZONE_DMA
1499         sizes->cs_dmacachep = kmem_cache_create(
1500             names->name_dma,
1501             sizes->cs_size,
1502             ARCH_KMALLOC_MINALIGN,
1503             ARCH_KMALLOC_FLAGS|SLAB_CACHE_DMA|
1504                 SLAB_PANIC,
1505             NULL);
1506 #endif
1507         sizes++;
1508         names++;
1509     }

```

1456-1509 行：建立剩下的一般性缓存。


```

1510     /* 4) Replace the bootstrap head arrays */
1511     {
1512         struct array_cache *ptr;
1513
1514         ptr = kmalloc(sizeof(struct arraycache_init), GFP_NOWAIT);

```

前面的一般性缓存已经建立，现在可以使用 **kmalloc** 分配内存了。

```

1515
1516     BUG_ON(cpu_cache_get(&cache_cache) !=
&initarray_cache.cache);

```

1516 行：首先判断 **cache_cache** 的 per-CPU 变量是否是从静态全局变量 **initarray_cache** 中来的，如果不是这内核初始化肯定是出问题了。

```

1517         memcpy(ptr, cpu_cache_get(&cache_cache),
1518             sizeof(struct arraycache_init));
1519     /*
1520     * Do not assume that spinlocks can be initialized via memcpy:
1521     */
1522     spin_lock_init(&ptr->lock);
1523
1524     cache_cache.array[smp_processor_id()] = ptr;

```

替换掉之前的静态 per-CPU 变量。

```

1525
1526     ptr = kmalloc(sizeof(struct arraycache_init), GFP_NOWAIT);
1527
1528     BUG_ON(cpu_cache_get(malloc_sizes[INDEX_AC].cs_cache)
1529         != &initarray_generic.cache);
1530     memcpy(ptr, cpu_cache_get(malloc_sizes[INDEX_AC].cs_cache),
1531         sizeof(struct arraycache_init));
1532     /*
1533     * Do not assume that spinlocks can be initialized via memcpy:
1534     */
1535     spin_lock_init(&ptr->lock);

```

```
1536
1537     malloc_sizes[INDEX_AC].cs_cachep->array[smp_processor_id()] =
1538         ptr;
1539 }
```

1510-1539 行：替换 bootstrap

```
1540     /* 5) Replace the bootstrap kmem_list3's */
1541     {
1542         int nid;
1543
1544         for_each_online_node(nid) {
1545             init_list(&cache_cache, &initkmem_list3[CACHE_CACHE + nid],
nid);
1546
1547             init_list(malloc_sizes[INDEX_AC].cs_cachep,
1548                 &initkmem_list3[SIZE_AC + nid], nid);
1549
1550             if (INDEX_AC != INDEX_L3) {
1551                 init_list(malloc_sizes[INDEX_L3].cs_cachep,
1552                     &initkmem_list3[SIZE_L3 + nid], nid);
1553             }
```

1545-1553: 替换掉缓存 cache_cache, arraycache_init, kmem_list3 的 kmem_list3 变量。

```
1554     }
1555 }
1556
1557     g_cpucache_up = EARLY;
```

设置 g_cpucache_up 为 EARLY。g_cpucache_up 的状态控制着建立缓存时 per-cache 的分配机制。

```
1558 }
```

9.4 vmalloc_init非连续内存区初始化

```
1088 void __init vmalloc_init(void)
1089 {
1090     struct vmmap_area *va;
1091     struct vm_struct *tmp;
1092     int i;
1093
1094     for_each_possible_cpu(i) {
1095         struct vmmap_block_queue *vbq;
1096
1097         vbq = &per_cpu(vmmap_block_queue, i);
1098         spin_lock_init(&vbq->lock);
1099         INIT_LIST_HEAD(&vbq->free);
1100     }
1101
1102     /* Import existing vmmap entries. */
1103     for (tmp = vmmap; tmp; tmp = tmp->next) {
1104         va = kzalloc(sizeof(struct vmmap_area), GFP_NOWAIT);
1105         va->flags = tmp->flags | VM_VM_AREA;
1106         va->va_start = (unsigned long)tmp->addr;
1107         va->va_end = va->va_start + tmp->size;
1108         __insert_vmmap_area(va);
```

将已经存在的 vmalloc 区域添加到 vmmap 中。

```
1109     }
1110
1111     vmmap_area_pcpu_hole = VMALLOC_END;
1112
1113     vmmap_initialized = true;
```

设置 vmmap_initialized 全局变量，表示已经初始化。

```
1114 }
```

10. 其它重要的初始化

下面这些内容的初始化非常重要，但是和内存管理关系不大，等到分析相关主题时再详细讨论，下面只是列出这些目录。

10.1 初始化调度程序

`sched_init`

10.2 初始化中断处理系统

`init_IRQ`

10.3 软中断初始化

`softirq_init`

10.4 定时器中断初始化

`time_init`

11. slab后续初始化

`kmem_cache_init_late`

```
1560 void __init kmem_cache_init_late(void)
1561 {
1562     struct kmem_cache *cachep;
1563
1564     /* 6) resize the head arrays to their final sizes */
```

```
1565     mutex_lock(&cache_chain_mutex);
1566     list_for_each_entry(cachep, &cache_chain, next)
1567         if (enable_cpucache(cachep, GFP_NOWAIT))
```

1567 行：初始化阶段 local cache 的大小是固定的，现在要根据对象大小重新计算

```
1568         BUG();
1569     mutex_unlock(&cache_chain_mutex);
1570
1571     /* Done! */
1572     g_cpucache_up = FULL;
```

到 1572 行：slab 分配器初始化全部完成。

```
1573
1574     /* Annotate slab for lockdep -- annotate the malloc caches */
1575     init_lock_keys();
1576
1577     /*
1578      * Register a cpu startup notifier callback that initializes
1579      * cpu_cache_get for all new cpus
1580      */
1581     register_cpu_notifier(&cpucache_notifier);
1582
1583     /*
1584      * The reap timers are started later, with a module init call: That part
1585      * of the kernel is not yet operational.
1586      */
1587 }
```

11.1 enable_cpucache

```
3947 static int enable_cpucache(struct kmem_cache *cachep, gfp_t gfp)
3948 {
```

```

3949     int err;
3950     int limit, shared;
3951
3952     /*
3953      * The head array serves three purposes:
3954      * - create a LIFO ordering, i.e. return objects that are cache-warm
3955      * - reduce the number of spinlock operations.
3956      * - reduce the number of linked list operations on the slab and
3957      *   bufctl chains: array operations are cheaper.
3958      * The numbers are guessed, we should auto-tune as described by
3959      * Bonwick.
3960      */
3961     if (cachep->buffer_size > 131072)
3962         limit = 1;
3963     else if (cachep->buffer_size > PAGE_SIZE)
3964         limit = 8;
3965     else if (cachep->buffer_size > 1024)
3966         limit = 24;
3967     else if (cachep->buffer_size > 256)
3968         limit = 54;
3969     else
3970         limit = 120;

```

根据对象大小，计算本地 **cpu** 缓存中对象的数目。

```

3971
3972     /*
3973      * CPU bound tasks (e.g. network routing) can exhibit cpu bound
3974      * allocation behaviour: Most allocs on one cpu, most free operations
3975      * on another cpu. For these cases, an efficient object passing between
3976      * cpus is necessary. This is provided by a shared array. The array
3977      * replaces Bonwick's magazine layer.
3978      * On uniprocessor, it's functionally equivalent (but less efficient)

```

```

3979      * to a larger limit. Thus disabled by default.
3980      */
3981      shared = 0;
3982      if (cachep->buffer_size <= PAGE_SIZE && num_possible_cpus() > 1)

3983          shared = 8;
3984
3985 #if DEBUG
3986      /*
3987      * With debugging enabled, large batchcount lead to excessively long
3988      * periods with disabled local interrupts. Limit the batchcount
3989      */
3990      if (limit > 32)
3991          limit = 32;
3992 #endif
3993      err = do_tune_cpucache(cachep, limit, (limit + 1) / 2, shared, gfp);

```

3993 行：配置 local cache

```

3994      if (err)
3995          printk(KERN_ERR "enable_cpucache failed for %s, error %d.\n",
3996                  cachep->name, -err);
3997      return err;
3998 }

```

11.2 do_tune_cpucache

```

3904 static int do_tune_cpucache(struct kmem_cache *cachep, int limit,
3905                             int batchcount, int shared, gfp_t gfp)
3906 {
3907     struct ccupdate_struct *new;
3908     int i;
3909

```

```

3910     new = kzalloc(sizeof(*new), gfp);
3911     if (!new)
3912         return -ENOMEM;
3913
3914     for_each_online_cpu(i) {
3915         new->new[i] = alloc_arraycache(cpu_to_node(i), limit,
3916                                     batchcount, gfp);
3917         if (!new->new[i]) {
3918             for (i--; i >= 0; i--)
3919                 kfree(new->new[i]);
3920             kfree(new);
3921             return -ENOMEM;
3922         }
3923     }

```

为每个 cpu 分配新的 struct array_cache 对象

```

3924     new->cachep = cachep;
3925
3926     on_each_cpu(do_ccupdate_local, (void *)new, 1);

```

用新的 struct array_cache 对象替换旧的 struct array_cache 对象

```

3927
3928     check_irq_on();
3929     cachep->batchcount = batchcount;
3930     cachep->limit = limit;
3931     cachep->shared = shared;
3932
3933     for_each_online_cpu(i) {
3934         struct array_cache *ccold = new->new[i];
3935         if (!ccold)
3936             continue;
3937         spin_lock_irq(&cachep->nodelists[cpu_to_node(i)]->list_lock);
3938         free_block(cachep, ccold->entry, ccold->avail, cpu_to_node(i));

```



```

3939         spin_unlock_irq(&cachep->nodelists[cpu_to_node(i)]->list_lock);
3940         kfree(ccold);
3941     }

```

释放旧的 struct array_cache 对象

```

3942     kfree(new);
3943     return alloc_kmemlist(cachep, gfp);

```

初始化 shared local cache 和 slab 三链

```

3944 }

```

```

3891 static void do_ccupdate_local(void *info)
3892 {
3893     struct ccupdate_struct *new = info;
3894     struct array_cache *old;
3895
3896     check_irq_off();
3897     old = cpu_cache_get(new->cachep);
3898
3899     new->cachep->array[smp_processor_id()] =
new->new[smp_processor_id()];

```

指向新的 struct array_cache 对象

```

3900     new->new[smp_processor_id()] = old;

```

保存旧的 struct array_cache 对象*/

```

3901 }

```

11.3 alloc_kmemlist

初始化 shared local cache 和 slab 三链，初始化完成后，slab 三链中没有任何 slab

```

3799  * This initializes kmem_list3 or resizes various caches for all nodes.
3800  */

```

```

3801 static int alloc_kmemlist(struct kmem_cache *cachep, gfp_t gfp)
3802 {
3803     int node;
3804     struct kmem_list3 *l3;
3805     struct array_cache *new_shared;
3806     struct array_cache **new_alien = NULL;
3807
3808     for_each_online_node(node) {
3809
3810         if (use_alien_caches) {
3811             new_alien = alloc_alien_cache(node, cachep->limit,
gfp);
3812             if (!new_alien)
3813                 goto fail;
3814         }
3815
3816         new_shared = NULL;
3817         if (cachep->shared) {
3818             new_shared = alloc_arraycache(node,
3819             cachep->shared*cachep->batchcount,
3820             0xbaadf00d, gfp);
3821             /* 分配 shared local cache */
3822             if (!new_shared) {
3823                 free_alien_cache(new_alien);
3824                 goto fail;
3825             }
3826
3827             获得旧的 slab 三链
3828             l3 = cachep->nodelists[node];
3829             if (l3) {
3830                 struct array_cache *shared = l3->shared;

```

```
3830
3831         spin_lock_irq(&l3->list_lock);
3832
3833         if (shared)
3834             //旧 slab 三链指针不为空，需要先释放旧的资源
3835             free_block(cachep, shared->entry,
3836                       shared->avail, node);
3837
3838             //指向新的 shared local cache
3839             l3->shared = new_shared;
3840             if (!l3->alien) {
3841                 l3->alien = new_alien;
3842                 new_alien = NULL;
3843             }
3844
3845             //计算 cache 中空闲对象的上限
3846             l3->free_limit = (1 + nr_cpus_node(node)) *
3847                             cachep->batchcount + cachep->num;
3848             spin_unlock_irq(&l3->list_lock);
3849             kfree(shared);
3850             free_alien_cache(new_alien);
3851
3852             //释放旧 shared local cache 的 struct array_cache 对象
3853             continue;
3854         }
3855
3856         //如果没有旧的 l3，分配新的 slab 三链
3857         l3 = kmalloc_node(sizeof(struct kmem_list3), gfp, node);
3858         if (!l3) {
3859             free_alien_cache(new_alien);
3860             kfree(new_shared);
3861             goto fail;
3862         }
3863
3864         l3->shared = new_shared;
3865         l3->alien = new_alien;
3866         new_alien = NULL;
3867         l3->free_limit = (1 + nr_cpus_node(node)) *
3868                         cachep->batchcount + cachep->num;
3869         spin_unlock_irq(&l3->list_lock);
3870         kfree(shared);
3871         free_alien_cache(new_alien);
3872
3873         //释放旧 shared local cache 的 struct array_cache 对象
3874         continue;
3875     }
3876
3877     //如果没有旧的 l3，分配新的 slab 三链
3878     l3 = kmalloc_node(sizeof(struct kmem_list3), gfp, node);
3879     if (!l3)
3880         goto fail;
3881
3882     l3->shared = new_shared;
3883     l3->alien = new_alien;
3884     new_alien = NULL;
3885     l3->free_limit = (1 + nr_cpus_node(node)) *
3886                     cachep->batchcount + cachep->num;
3887     spin_unlock_irq(&l3->list_lock);
3888     kfree(shared);
3889     free_alien_cache(new_alien);
3890
3891     //释放旧 shared local cache 的 struct array_cache 对象
3892     continue;
3893 }
```

//初始化 slab 三链

```
3856         kmem_list3_init(l3);
3857         l3->next_reap = jiffies + REAPTIMEOUT_LIST3 +
3858             ((unsigned long)cachep) % REAPTIMEOUT_LIST3;
3859         l3->shared = new_shared;
3860         l3->alien = new_alien;
3861         l3->free_limit = (1 + nr_cpus_node(node)) *
3862             cachep->batchcount + cachep->num;
3863         cachep->nodelists[node] = l3;
3864     }
3865     return 0;
3866
3867 fail:
3868     if (!cachep->next.next) {
3869         /* Cache is not active yet. Roll back what we did */
3870         node--;
3871         while (node >= 0) {
3872             if (cachep->nodelists[node]) {
3873                 l3 = cachep->nodelists[node];
3874
3875                 kfree(l3->shared);
3876                 free_alien_cache(l3->alien);
3877                 kfree(l3);
3878                 cachep->nodelists[node] = NULL;
3879             }
3880             node--;
3881         }
3882     }
3883     return -ENOMEM;
3884 }
```

12. 启动控制台输出

console_init

将保存在缓存中的内容输出到屏幕上。

13. fork_init

根据物理内存大小计算允许创建进程/线程的数量

```
189 void __init fork_init(unsigned long mempages)
190 {
191     #ifndef __HAVE_ARCH_TASK_STRUCT_ALLOCATOR
192     #ifndef ARCH_MIN_TASKALIGN
193     #define ARCH_MIN_TASKALIGN    L1_CACHE_BYTES
194     #endif
195     /* create a slab on which task_structs can be allocated */
196     task_struct_cachep =
197         kmem_cache_create("task_struct", sizeof(struct task_struct),
198         ARCH_MIN_TASKALIGN, SLAB_PANIC | SLAB_NOTRACK,
199 NULL);
200
201     #endif
202
203     /* do the arch specific task caches init */
204     arch_task_cache_init();
205
206     /*
207     * The default maximum number of threads is set to a safe
208     * value: the thread structures can take up at most half
209     * of memory.
210     */
```

```

209     max_threads = mempages / (8 * THREAD_SIZE / PAGE_SIZE);
210
211     /*
212      * we need to allow at least 20 threads to boot a system
213      */
214     if(max_threads < 20)
215         max_threads = 20;
216
217     init_task.signal->rlim[RLIMIT_NPROC].rlim_cur = max_threads/2;
218     init_task.signal->rlim[RLIMIT_NPROC].rlim_max = max_threads/2;
219     init_task.signal->rlim[RLIMIT_SIGPENDING] =
220         init_task.signal->rlim[RLIMIT_NPROC];
221 }

```

14. proc_caches_init

与进程相关数据结构的缓存创建

```

1470 void __init proc_caches_init(void)
1471 {
1472     sighand_cachep = kmem_cache_create("sighand_cache",
1473         sizeof(struct sighand_struct), 0,
1474         SLAB_HWCACHE_ALIGN|SLAB_PANIC|SLAB_DESTROY_BY_RCU|
1475         SLAB_NOTRACK, sighand_ctor);
1476     signal_cachep = kmem_cache_create("signal_cache",
1477         sizeof(struct signal_struct), 0,
1478         SLAB_HWCACHE_ALIGN|SLAB_PANIC|SLAB_NOTRACK,
1479         NULL);
1479     files_cachep = kmem_cache_create("files_cache",
1480         sizeof(struct files_struct), 0,

```

```

1481          SLAB_HWCACHE_ALIGN|SLAB_PANIC|SLAB_NOTRACK,
NULL);
1482      fs_cachep = kmem_cache_create("fs_cache",
1483          sizeof(struct fs_struct), 0,
1484          SLAB_HWCACHE_ALIGN|SLAB_PANIC|SLAB_NOTRACK,
NULL);
1485      mm_cachep = kmem_cache_create("mm_struct",
1486          sizeof(struct mm_struct), ARCH_MIN_MMSTRUCT_ALIGN,
1487          SLAB_HWCACHE_ALIGN|SLAB_PANIC|SLAB_NOTRACK,
NULL);
1488      vm_area_cachep = KMEM_CACHE(vm_area_struct, SLAB_PANIC);
1489      mmap_init();
1490 }

```

15. 块缓存初始化buffer_init

```

3355 void __init buffer_init(void)
3356 {
3357     int nrpages;
3358
3359     bh_cachep = kmem_cache_create("buffer_head",
3360         sizeof(struct buffer_head), 0,
3361         (SLAB_RECLAIM_ACCOUNT|SLAB_PANIC|
3362         SLAB_MEM_SPREAD),
3363         NULL);
3364
3365     /*
3366      * Limit the bh occupancy to 10% of ZONE_NORMAL
3367      */
3368     nrpages = (nr_free_buffer_pages() * 10) / 100;

```

```
3369     max_buffer_heads = nrpages * (PAGE_SIZE / sizeof(struct buffer_head));
3370     hotcpu_notifier(buffer_cpu_notify, 0);
3371 }
```

16. 虚拟根文件系统安装vfs_caches_init

该函数对文件系统来说非常重要，等讲解文件系统时在详细分析。

```
2354 void __init vfs_caches_init(unsigned long mempages)
2355 {
2356     unsigned long reserve;
2357
2358     /* Base hash sizes on available memory, with a reserve equal to
2359        150% of current kernel size */
2360
2361     reserve = min((mempages - nr_free_pages()) * 3/2, mempages - 1);
2362     mempages -= reserve;
2363
2364     names_cache = kmem_cache_create("names_cache", PATH_MAX, 0,
2365                                     SLAB_HWCACHE_ALIGN|SLAB_PANIC, NULL);
2366
2367     dcache_init();
2368     inode_init();
2369     files_init(mempages);
2370     mnt_init();
2371     bdev_cache_init();
2372     chrdev_init();
2373 }
```


17. 安装PROC文件系统proc_root_init

```
104 void __init proc_root_init(void)
105 {
106     int err;
107
108     proc_init_inodecache();
109     err = register_filesystem(&proc_fs_type);
110     if (err)
111         return;
112     proc_mnt = kern_mount_data(&proc_fs_type, &init_pid_ns);
113     err = PTR_ERR(proc_mnt);
114     if (IS_ERR(proc_mnt)) {
115         unregister_filesystem(&proc_fs_type);
116         return;
117     }
118
119     proc_symlink("mounts", NULL, "self/mounts");
120
121     proc_net_init();
122
123 #ifdef CONFIG_SYSVIPC
124     proc_mkdir("sysvipc", NULL);
125 #endif
126     proc_mkdir("fs", NULL);
127     proc_mkdir("driver", NULL);
128     proc_mkdir("fs/nfsd", NULL); /* somewhere for the nfsd filesystem to be
mounted */
129 #if defined(CONFIG_SUN_OPENPROMFS) ||
defined(CONFIG_SUN_OPENPROMFS_MODULE)
130     /* just give it a mountpoint */
```

```
131     proc_mkdir("openprom", NULL);
132 #endif
133     proc_tty_init();
134 #ifdef CONFIG_PROC_DEVICETREE
135     proc_device_tree_init();
136 #endif
137     proc_mkdir("bus", NULL);
138     proc_sys_init();
139 }
```

18. 后start_kernel时代

start_kernel 最后调用 rest_init 函数进行初始化。

```
424 static noinline void __init_refok rest_init(void)
425     __releases(kernel_lock)
426 {
427     int pid;
428
429     rcu_scheduler_starting();
430     kernel_thread(kernel_init, NULL, CLONE_FS | CLONE_SIGHAND);
431     numa_default_policy();
432     pid = kernel_thread(kthreadd, NULL, CLONE_FS | CLONE_FILES);
433     rcu_read_lock();
434     kthreadd_task = find_task_by_pid_ns(pid, &init_pid_ns);
435     rcu_read_unlock();
436     unlock_kernel();
437
438     /*
439      * The boot idle thread must execute schedule()
```

```
440     * at least once to get things moving:
441     */
442     init_idle_bootup_task(current);
443     preempt_enable_no_resched();
444     schedule();
445     preempt_disable();
446
447     /* Call into cpu_idle with preempt disabled */
448     cpu_idle();
449 }
```

`rest_init` 使用 `kernel_thread` 创建内核线程 `init`(1 号进程)。

`rest_init` 使用 `kernel_thread` 创建内核线程 `kthreadd`。`kthreadd` 根据内核需要，动态的创建其他内核线程。

0 号进程进入 `cpu_idle` 中，在系统中没有其他就绪进程时，就进入 0 号进程运行。

18.1 Kthreadd

```
214 int kthreadd(void *unused)
215 {
216     struct task_struct *tsk = current;
217
218     /* Setup a clean context for our children to inherit. */
219     set_task_comm(tsk, "kthreadd");
220     ignore_signals(tsk);
221     set_cpus_allowed_ptr(tsk, cpu_all_mask);
222     set_mems_allowed(node_states[N_HIGH_MEMORY]);
223
224     current->flags |= PF_NOFREEZE | PF_FREEZER_NOSIG;
225
226     for (;;) {
```

```
227         set_current_state(TASK_INTERRUPTIBLE);
```

将本进程设置为可中断睡眠模式。

```
228         if (list_empty(&kthread_create_list))
```

如果没有需要创建的内核线程，则让出 CPU。

```
229             schedule();
230             __set_current_state(TASK_RUNNING);
231
232             spin_lock(&kthread_create_lock);
233             while (!list_empty(&kthread_create_list)) {
234                 struct kthread_create_info *create;
235
236                 create = list_entry(kthread_create_list.next,
237                                     struct kthread_create_info, list);
238                 list_del_init(&create->list);
239                 spin_unlock(&kthread_create_lock);
240
241                 create_kthread(create);
```

创建内核线程。

kthread_create_list 在内核运行的过程中会变化，接受创建线程的请求。

```
242
243             spin_lock(&kthread_create_lock);
244         }
245         spin_unlock(&kthread_create_lock);
246     }
247
248     return 0;
249 }
```

18.2 cpu_idle

空转函数。

```
78 /*
79  * The idle thread. There's no useful work to be
80  * done, so just try to conserve power and have a
81  * low exit latency (ie sit in a loop waiting for
82  * somebody to say that they'd like to reschedule)
83  */
84 void cpu_idle(void)
85 {
86     int cpu = smp_processor_id();
87
88     /*
89      * If we're the non-boot CPU, nothing set the stack canary up
90      * for us. CPU0 already has it initialized but no harm in
91      * doing it again. This is a good place for updating it, as
92      * we wont ever return from this function (so the invalid
93      * canaries already on the stack wont ever trigger).
94      */
95     boot_init_stack_canary();
96
97     current_thread_info()->status |= TS_POLLING;
98
99     /* endless idle loop with no priority at all */
100    while (1) {
101        tick_nohz_stop_sched_tick(1);
102        while (!need_resched()) {
103
104            check_pgt_cache();
105            rmb();
106
107            if (cpu_is_offline(cpu))
108                play_dead();
```

```
109
110     local_irq_disable();
111     /* Don't trace irqs off for idle */
112     stop_critical_timings();
113     pm_idle();
114     start_critical_timings();
115 }
116 tick_nohz_restart_sched_tick();
117 preempt_enable_no_resched();
118 schedule();
119 preempt_disable();
120 }
121 }
```

18.3 Kernel_init

```
854 static int __init kernel_init(void * unused)
855 {
856     lock_kernel();
857
858     /*
859      * init can allocate pages on any node
860      */
861     set_mems_allowed(node_states[N_HIGH_MEMORY]);
```

init 可以在所有节点上分配内存

```
862     /*
863      * init can run on any cpu.
864      */
865     set_cpus_allowed_ptr(current, cpu_all_mask);
```

init 可以在所有节点上运行。

```
866     /*
```

```
867      * Tell the world that we're going to be the grim
868      * reaper of innocent orphaned children.
869      *
870      * We don't want people to have to make incorrect
871      * assumptions about where in the task array this
872      * can be found.
873      */
874      init_pid_ns.child_reaper = current;//设置当前进程为接受孤儿进程的进程
875
876      cad_pid = task_pid(current);
877
878      smp_prepare_cpus(setup_max_cpus);
879
880      do_pre_smp_initcalls();
881      start_boot_trace();
882
883      smp_init();
884      sched_init_smp();
885
886      do_basic_setup();
887
888      /* Open the /dev/console on the rootfs, this should never fail */
889      if (sys_open((const char __user *) "/dev/console", O_RDWR, 0) < 0)
890          printk(KERN_WARNING "Warning: unable to open an initial
console.\n");
891
892      (void) sys_dup(0);
893      (void) sys_dup(0);
894      /*
895       * check if there is an early userspace init.  If yes, let it do all
896       * the work
897       */
```

```
898
899     if (!ramdisk_execute_command)
900         ramdisk_execute_command = "/init";
901
902     if (sys_access((const char __user *) ramdisk_execute_command, 0) != 0) {
903         ramdisk_execute_command = NULL;
904         prepare_namespace();
905     }
906
907     /*
908      * Ok, we have completed the initial bootup, and
909      * we're essentially up and running. Get rid of the
910      * initmem segments and start the user-mode stuff..
911      */
912
913     init_post();
914     return 0;
915 }
```

18.3.1 子系统初始化

```
785 static void __init do_basic_setup(void)
786 {
787     init_workqueues();//初始化工作队列 events
788     cpuset_init_smp();//空函数
789     usermodehelper_init();//初始化工作队列 khelper
790     init_tmpfs();//初始化临时文件系统
791     driver_init();//建立设备驱动模型
792     init_irq_proc();//向/proc 文件系统中增加子目录 irq
793     do_ctors();//初始化内核构造函数
```



```
794     do_initcalls();//首先说一下，内核中有一个专门的节，用来存放初
初始化末尾要被调用的函数，改函数就是用来调用该节中的函数的。
795 }
```

18.3.2 启动SHELL

```
812  * makes it inline to init() and it becomes part of init.text section
813  */
814 static noline int init_post(void)
815     __releases(kernel_lock)
816 {
817     /* need to finish all async __init code before freeing the memory */
818     async_synchronize_full();
819     free_initmem();
820     unlock_kernel();
821     mark_rodata_ro();
822     system_state = SYSTEM_RUNNING;
823     numa_default_policy();
824
825
826     current->signal->flags |= SIGNAL_UNKILLABLE;
827
828     if (ramdisk_execute_command) {
829         run_init_process(ramdisk_execute_command);
830         printk(KERN_WARNING "Failed to execute %s\n",
831             ramdisk_execute_command);
832     }
833
834     /*
835      * We try each of these until one succeeds.
836      *
```

```

837      * The Bourne shell can be used instead of init if we are
838      * trying to recover a really broken machine.
839      */
840      if (execute_command) {
841          run_init_process(execute_command);
842          printk(KERN_WARNING "Failed to execute %s.  Attempting "
843                      "defaults...\n", execute_command);
844      }
845      run_init_process("/sbin/init");
846      run_init_process("/etc/init");
847      run_init_process("/bin/init");
848      run_init_process("/bin/sh");
849
850      panic("No init found.  Try passing init= option to kernel. "
851            "See Linux Documentation/init.txt for guidance.");
852 }

```

```

805 static void run_init_process(char *init_filename)
806 {
807     argv_init[0] = init_filename;
808     kernel_execve(init_filename, argv_init, envp_init);
809 }

```

808 行，从内核调用 `execve` 函数，进入 `shell` 进程，内核第一次进入用户空间（从内核空间进入用户空间）。

18.3.3 Initrd以及磁盘根文件系统的安

在 `Kernel_init` 的 904 `prepare_namespace()` 函数，会进行根文件系统的安装。

```

363 /*
364  * Prepare the namespace - decide what/where to mount, load ramdisks, etc.
365  */

```

```
366 void __init prepare_namespace(void)
367 {
368     int is_floppy;
369
370     if (root_delay) {
371         printk(KERN_INFO "Waiting %dsec before mounting root device...\n",
372             root_delay);
373         ssleep(root_delay);
374     }
375
376     /*
377      * wait for the known devices to complete their probing
378      *
379      * Note: this is a potential source of long boot delays.
380      * For example, it is not atypical to wait 5 seconds here
381      * for the touchpad of a laptop to initialize.
382      */
383     wait_for_device_probe();
384
385     md_run_setup();
386
387     if (saved_root_name[0]) {
388         root_device_name = saved_root_name;
389         if (!strncmp(root_device_name, "mtd", 3) ||
390             !strncmp(root_device_name, "ubi", 3)) {
391             mount_block_root(root_device_name, root_mountflags);
392             goto out;
393         }
394         ROOT_DEV = name_to_dev_t(root_device_name);
395         if (strncmp(root_device_name, "/dev/", 5) == 0)
```

```
396         root_device_name += 5;
397     }
398
399     if (initrd_load())
400         goto out;
```

根据内核启动参数的配置，决定是否加载 `initrd`（`initrd` 不是必须的）。

```
401
402     /* wait for any asynchronous scanning to complete */
403     if ((ROOT_DEV == 0) && root_wait) {
404         printk(KERN_INFO "Waiting for root device %s...\n",
405             saved_root_name);
406         while (driver_probe_done() != 0 ||
407             (ROOT_DEV = name_to_dev_t(saved_root_name)) == 0)
408             msleep(100);
409         async_synchronize_full();
410     }
411
412     is_floppy = MAJOR(ROOT_DEV) == FLOPPY_MAJOR;
413
414     if (is_floppy && rd_doload && rd_load_disk(0))
415         ROOT_DEV = Root_RAM0;
416
417     mount_root();
```

加载磁盘根文件系统

```
418 out:
419     devtmpfs_mount("dev");
420     sys_mount(".", "/", NULL, MS_MOVE, NULL);
421     sys_chroot(".");
422 }
```

19. 参考书籍

- [1] linux 内核 2.6.34 源码
- [2] 深入了解 linux 内核 （第三版）
- [3] linux 内核源代码情景分析 毛德超
- [4] 深入 linux 内核架构
- [5] <http://blog.csdn.net/yunsongice>(网络资源)