

```
1 #include "console.h"
2 #include "io.h"
3 #include "string.h"
4 #include "stdarg.h"
5
6 typedef enum
7 {
8     COLOUR_BLACK = 0,
9     COLOUR_BLUE = 1,
10    COLOUR_GREEN = 2,
11    COLOUR_CYAN = 3,
12    COLOUR_RED = 4,
13    COLOUR_PINK = 5,
14    COLOUR_BROWN = 6,
15    COLOUR_LIGHT_GREY = 7,
16    COLOUR_DARK_GREY = 8,
17    COLOUR_LIGHT_BLUE = 9,
18    COLOUR_LIGHT_GREEN = 10,
19    COLOUR_LIGHT_CYAN = 11,
20    COLOUR_LIGHT_RED = 12,
21    COLOUR_LIGHT_PINK = 13,
22    COLOUR_YELLOW = 14,
23    COLOUR_WHITE = 15,
24 } colour_t;
25
26 static uint8_t make_attr(colour_t background, colour_t foreground)
27 {
28     return (background << 4) | foreground;
29 }
30
31 typedef struct
32 {
33     uint8_t character;
34     uint8_t attr;
35 }
36
37 vchar_t;
38
39 static int x, y;
40
41 static const int width = 80, height = 25;
42
43 static vchar_t *const vram = (void *)0xb0000;
44
45 static uint16_t base_vga_port;
46
47 void console_init()
48 {
49     // read base vga port from bios data area
50     base_vga_port = *(uint16_t *) 0x463;
51
52     memset16(vram, make_attr(COLOUR_BLACK, COLOUR_LIGHT_GREY) << 8,
53              width * height);
54     memset16(vram + width * height,
55              (make_attr(COLOUR_RED, COLOUR_WHITE) << 8) | 'X', width);
56 }
57
58 static void scroll_up()
59 {
60     memcpy(vram, vram + width, width * (height - 1) * 2);
61
62     uint16_t empty_attr = make_attr(COLOUR_BLACK, COLOUR_LIGHT_GREY) << 8;
63     memset16(vram + width * (height - 1), empty_attr, width);
64 }
65
66 static void newline()
67 {
68     x = 0;
69     if (++y == height)
70     {
71         scroll_up();
72     }
```

```
73     y = height - 1;
74 }
75
76 static void update_cursor()
77 {
78     uint16_t pos = y * width + x;
79
80     outb(base_vga_port, 0x0e);
81     outb(base_vga_port + 1, (pos >> 8) & 0xff);
82
83     outb(base_vga_port, 0x0f);
84     outb(base_vga_port + 1, pos & 0xff);
85 }
86
87 static void putc(char c)
88 {
89     if (c == '\r')
90     {
91         x = 0;
92         return;
93     }
94
95     if (c == '\n')
96     {
97         newline();
98         return;
99     }
100
101     vram[y * width + x].attr = COLOUR_LIGHT_GREY;
102     vram[y * width + x].character = c;
103
104     if (++x == width)
105     {
106         newline();
107     }
108
109 }
110
111 void console_puts(const char *str, uint32_t len)
112 {
113     while (len--)
114     {
115         putc(*str++);
116     }
117
118     update_cursor();
119 }
120
121 static void itoa(int n, char *out)
122 {
123     int len = 0, negative = 0;
124
125     if (n == 0)
126     {
127         out[len++] = '0';
128         out[len++] = 0;
129         return;
130     }
131
132     if (n < 0)
133     {
134         negative = 1;
135         n *= -1;
136     }
137
138     while (n)
139     {
140         out[len++] = '0' + (n % 10);
141         n /= 10;
142     }
143
144     if (negative)
```

```
145 { out[len++] = '-';
146 }
147
148 out[len] = 0;
149
150 for (int i = 0; i < len / 2; i++)
151 {
152     int j = len - i - 1;
153     char c = out[j];
154     out[j] = out[i];
155     out[i] = c;
156 }
157
158 }
159
160 static char hex_dig(uint32_t dig)
161 {
162     if (dig <= 9)
163     {
164         return '0' + dig;
165     }
166     else if (dig >= 10 && dig <= 15)
167     {
168         return 'a' + dig - 10;
169     }
170     return '?';
171 }
172
173 static void utox(uint32_t u, char *out)
174 {
175     out[0] = hex_dig(u >> 28) & 0xf;
176     out[1] = hex_dig(u >> 24) & 0xf;
177     out[2] = hex_dig(u >> 20) & 0xf;
178     out[3] = hex_dig(u >> 16) & 0xf;
179     out[4] = hex_dig(u >> 12) & 0xf;
180     out[5] = hex_dig(u >> 8) & 0xf;
181     out[6] = hex_dig(u >> 4) & 0xf;
182     out[7] = hex_dig(u >> 0) & 0xf;
183     out[8] = 0;
184 }
185
186 void printf(const char *format, ...)
187 {
188     va_list va;
189     va_start(va, format);
190     vprintf(format, va);
191     va_end(va);
192 }
193
194 void vprintf(const char *format, va_list va)
195 {
196     while (1)
197     {
198         char c = *format++;
199
200         if (c == 0)
201         {
202             break;
203         }
204
205         if (c != '%')
206         {
207             putc(c);
208             continue;
209         }
210
211         switch (c = *format++)
212         {
213             case 'd':
214                 {
215                     char buff[16];
216                     itoa(va_arg(va, int), buff);
```

```
217 console_puts(buff, strlen(buff));
218 break;
219 }
220 case 'x':
221 {
222     char buff[16];
223     utox(va_arg(va, uint32_t), buff);
224     console_puts(buff, strlen(buff));
225     break;
226 }
227 case 's':
228 {
229     const char *str = va_arg(va, const char *);
230     console_puts(str, strlen(str));
231     break;
232 }
233 case 'c':
234 {
235     char c = va_arg(va, char);
236     putc(c);
237     break;
238 }
239 case 0:
240 {
241     goto ret;
242 }
243 default:
244 {
245     putc(c);
246     break;
247 }
248 }
249 ret:
250     update_cursor();
251 }
```

```
1 #ifndef CONSOLE_H
2 #define CONSOLE_H
3
4 #include "stdarg.h"
5
6 void console_init();
7
8 void console_puts(const char *str, uint32_t len);
9
10 void printf(const char *format, ...);
11
12 void vprintf(const char *format, va_list va);
13
14 #endif
```

```
1 use32
2
3 global __syscall0
4 global __syscall1
5 global __syscall2
6 global __syscall3
7
8 extern _main
9 extern _exit
10
11 section .crt0
12
13 %macro perform_syscall 0
14     push ecx
15     push edx
16     mov ecx, esp
17     mov edx, .ret
18     sysenter
19     .ret:
20     pop edx
21     pop ecx
22 %endmacro
23
24 start:
25     call _main
26     push eax
27     call _exit
28
29     __syscall0:
30     mov eax, [esp+0+4]
31     perform_syscall
32     ret
33
34     __syscall1:
35     push ebx
36     mov ebx, [esp+4+8]
37     mov eax, [esp+4+4]
38     perform_syscall
39     pop ebx
40     ret
41
42     __syscall2:
43     push ebx
44     push edi
45     mov edi, [esp+8+12]
46     mov ebx, [esp+8+8]
47     mov eax, [esp+8+4]
48     perform_syscall
49     pop edi
50     pop ebx
51     ret
52
53     __syscall3:
54     push ebx
55     push edi
56     push esi
57     mov esi, [esp+12+16]
58     mov edi, [esp+12+12]
59     mov ebx, [esp+12+8]
60     mov eax, [esp+12+4]
61     perform_syscall
62     pop esi
63     pop edi
64     pop ebx
65     ret
```

```
1 #include <radium.h>
2 #include "crt.h"
3
4 void regdump()
5 {
6     _syscall0(SYS_REGDUMP);
7 }
8
9 void exit(int status)
10 {
11     _syscall1(SYS_EXIT, status);
12 }
13
14 void yield()
15 {
16     _syscall0(SYS_YIELD);
17 }
18
19 uint32_t fork()
20 {
21     return _syscall0(SYS_FORK);
22 }
23
24 uint32_t wait(int *stat_loc)
25 {
26     return _syscall1(SYS_WAIT, (uint32_t) stat_loc);
27 }
28
29 void console_log(const char *str)
30 {
31     uint32_t len = 0;
32     for (const char *p = str; *p; p++)
33     {
34         len++;
35     }
36     _syscall12(SYS_CONSOLE_LOG, (uint32_t) str, len);
37 }
38
39 }
```

```
1 #ifndef CRT_H
2 #define CRT_H
3
4 typedef unsigned int uint32_t;
5
6 #define NULL ((void*)0)
7
8 uint32_t _syscall10(uint32_t number);
9
10 uint32_t _syscall11(uint32_t number, uint32_t arg1);
11
12 uint32_t _syscall12(uint32_t number, uint32_t arg1, uint32_t arg2);
13
14 uint32_t _syscall13(uint32_t number, uint32_t arg1, uint32_t arg2, uint32_t arg3);
15
16 void regdump();
17
18 void exit(int status);
19
20 void yield();
21
22 uint32_t fork();
23
24 uint32_t wait(int *stat_loc);
25
26 void console_log(const char *str);
27
28 int main();
29
30 #endif
```

```
1 #include "gdt.h"
2 #include "panic.h"
3 #include "types.h"
4
5 typedef struct
6 {
7     uint16_t limit_0_15;
8     uint16_t base_0_15;
9     uint8_t base_16_23;
10    uint8_t access;
11    uint8_t limit_16_19_and_flags;
12    uint8_t base_24_31;
13 }
14 gdt_entry_t;
15
16 static gdt_entry_t gdt[6];
17
18 volatile struct
19 {
20     uint16_t size;
21     gdt_entry_t *offset;
22 } __attribute__((packed)) gdt;
23
24 static void gdt_set_entry_raw(gdt_selector_t sel, uint32_t base, uint32_t limit,
25                             uint8_t access)
26 {
27     if (sel >= sizeof(gdt))
28     {
29         panic("GDT overflow");
30     }
31
32     uint8_t flags = 1 << 6;    // 32 bit segment
33
34     // 1 MiB is the maximum segment size that can be expressed with 1 byte
35     // granularity. If we need to express a size bigger than this, we need
36     // to divide the size by 4 KiB and use 4 KiB granularity.
37     if (limit >= (1 << 20))
38     {
39         limit /= 4096;
40         flags |= 1 << 7;    // 4 KiB granularity
41     }
42
43     gdt_entry_t ent;
44     ent.limit_0_15 = limit & 0xffff;
45     ent.base_0_15 = base & 0xffff;
46     ent.base_16_23 = (base >> 16) & 0xff;
47     ent.access = access;
48     ent.limit_16_19_and_flags = ((limit >> 16) & 0x0f) | flags;
49     ent.base_24_31 = (base >> 24) & 0xff;
50
51     gdt[sel / sizeof(gdt_entry_t)] = ent;
52 }
53
54 void gdt_set_entry(gdt_selector_t sel, uint32_t base, uint32_t limit,
55                  gdt_privilege_t priv, gdt_type_t type)
56 {
57     gdt_set_entry_raw(sel, base, limit, (1 << 7) | // present
58                          ((priv & 3) << 5) | // privilege
59                          (1 << 4) | // dumno lol
60                          ((type & 1) << 3) | // code/data
61                          (1 << 1) // data segments are always writable, code
62                                     // segments are always readable
63     );
64 }
65
66 void gdt_set_tss(gdt_selector_t sel, uint32_t base, uint32_t limit)
67 {
68     gdt_set_entry_raw(sel, base, limit, 0x89);
69 }
70
71 void gdt_reload();
72
```

```
73 void gdt_init()
74 {
75     gdt.size = sizeof(gdt) - 1;
76     gdt.offset = gdt;
77
78     gdt_set_entry(GDT_KERNEL_CODE, 0, 0xffffffff, GDT_KERNEL, GDT_CODE);
79     gdt_set_entry(GDT_KERNEL_DATA, 0, 0xffffffff, GDT_KERNEL, GDT_DATA);
80     gdt_set_entry(GDT_USER_CODE, 0, 0xffffffff, GDT_USER, GDT_CODE);
81     gdt_set_entry(GDT_USER_DATA, 0, 0xffffffff, GDT_USER, GDT_DATA);
82
83     gdt_reload();
84 }
```

```
1 extern gdt_r
2 global gdt_reload
3
4 section .text
5 gdt_reload:
6     jmp 0x00:flush_cs
7     .flush_cs:
8     mov eax, 0x10
9     mov ds, eax
10    mov es, eax
11    mov fs, eax
12    mov gs, eax
13    mov ss, eax
14    ret
15
```



```
1 #ifndef GDT_H
2 #define GDT_H
3
4 #include "types.h"
5
6 typedef enum
7 {
8     GDT_KERNEL_CODE = 0x08,
9     GDT_KERNEL_DATA = 0x10,
10    GDT_USER_CODE = 0x18,
11    GDT_USER_DATA = 0x20,
12    GDT_TSS = 0x28,
13 }
14 gdt_selector_t;
15
16 typedef enum
17 {
18     GDT_KERNEL = 0,
19     GDT_USER = 3,
20 }
21 gdt_privilege_t;
22
23 typedef enum
24 {
25     GDT_DATA = 0,
26     GDT_CODE = 1,
27 }
28 gdt_type_t;
29
30 void gdt_set_entry(gdt_selector_t sel, uint32_t base, uint32_t limit,
31                  gdt_privilege_t priv, gdt_type_t type);
32
33 void gdt_set_tss(gdt_selector_t sel, uint32_t base, uint32_t limit);
34
35 void gdt_reload();
36
37 void gdt_init();
38
39 #endif
```

```
1 #include "console.h"
2 #include "gdt.h"
3 #include "idt.h"
4 #include "io.h"
5 #include "types.h"
6
7 typedef struct
8 {
9     uint16_t offset_0_15;
10    uint16_t segment;
11    uint8_t zero;
12    uint8_t type_attr;
13    uint16_t offset_16_31;
14 }
15 idt_entry_t;
16
17 static idt_entry_t idt[256];
18
19 volatile struct
20 {
21     uint16_t size;
22     idt_entry_t *offset;
23 } __attribute__((packed)) idtr;
24
25 void interrupts_register_isr(uint8_t interrupt_no, uint32_t handler)
26 {
27     idt_entry_t ent;
28
29     ent.offset_0_15 = handler & 0xffff;
30     ent.segment = GDT_KERNEL_CODE;
31     ent.zero = 0;
32     ent.type_attr = (1 << 7) // present
33                   | 0xe     // interrupt gate
34                   ;
35     ent.offset_16_31 = (handler >> 16) & 0xffff;
36
37     idt[interrupt_no] = ent;
38 }
39
40 static void remap_irqs()
41 {
42     // remap IRQ table
43     outb(0x20, 0x11);
44     outb(0xa0, 0x11);
45     outb(0x21, 0x20);
46     outb(0xa1, 0x28);
47     outb(0x21, 0x04);
48     outb(0xa1, 0x02);
49     outb(0x21, 0x01);
50     outb(0xa1, 0x01);
51     outb(0x21, 0x00);
52     outb(0xa1, 0x00);
53 }
54
55 void idt_init_asm();
56
57 void idt_load();
58
59 void idt_init()
60 {
61     remap_irqs();
62     idt_init_asm();
63
64     idtr.size = sizeof(idt) - 1;
65     idtr.offset = idt;
66     idt_load();
67 }
```

```
1 #ifndef IDT_H
2 #define IDT_H
3
4 void idt_init();
5
6 #endif
```

```
1 #include "crt.h"
2
3 int main()
4 {
5     console_log("Hello world from init!\n");
6
7     if (fork())
8     {
9         wait(NULL);
10        return 123;
11    }
12    else
13    {
14        regdump();
15        yield();
16        regdump();
17        return 456;
18    }
19 }
```

```
1 #ifndef IO_H
2 #define IO_H
3
4 #include "types.h"
5
6 static void __attribute__((unused)) outb(uint16_t port, uint8_t value)
7 {
8     __asm__ volatile ("outb %1, %0"::"r" (value), "Nd"(port));
9 }
10
11 static uint8_t __attribute__((unused)) inb(uint16_t port)
12 {
13     volatile uint8_t value;
14     __asm__ volatile ("inb %0, %1"::"=r" (value):"Nd"(port));
15     return value;
16 }
17
18 static void __attribute__((unused)) outl(uint16_t port, uint32_t value)
19 {
20     __asm__ volatile ("outd %1, %0"::"r" (value), "Nd"(port));
21 }
22
23 static uint32_t __attribute__((unused)) inl(uint16_t port)
24 {
25     volatile uint32_t value;
26     __asm__ volatile ("ind %0, %1"::"=r" (value):"Nd"(port));
27     return value;
28 }
29
30 #endif
```

```
1 global idt_init_asm
2 global idt_load
3 extern interrupts_register_isr
4 extern idtr
5 extern panic
6 extern sched_switch
7
8 %macro register_isr 1
9     push isr_%1
10    push %1
11    call interrupts_register_isr
12    add esp, 8
13 %endmacro
14
15 %macro ack_irq 0
16     push ax
17     mov al, 0x20
18     out 0xa0, al
19     out 0x20, al
20     pop ax
21 %endmacro
22
23 %macro begin_isr 1
24     jmp isr_%1.end
25     isr_%1:
26 %endmacro
27
28 %macro end_isr 1
29     .end:
30     register_isr %1
31 %endmacro
32
33 %macro generic_exception 2
34     begin_isr %1
35     push .msg
36     call panic
37     .msg db %2, 0
38     end_isr %1
39 %endmacro
40
41 section .text
42 idt_load:
43     _lidt [idtr]
44     ret
45
46 idt_init_asm:
47
48     generic_exception 0, "divide by zero"
49     generic_exception 1, "debug"
50     generic_exception 2, "non-maskable interrupt"
51     generic_exception 3, "breakpoint"
52     generic_exception 4, "overflow"
53     generic_exception 5, "bound range exceeded"
54     generic_exception 6, "invalid opcode"
55     generic_exception 7, "device not available"
56     generic_exception 8, "double fault"
57     generic_exception 10, "invalid tss"
58     generic_exception 11, "segment not present"
59     generic_exception 12, "stack segment fault"
60     generic_exception 13, "general protection fault"
61
62     ; page fault
63     begin_isr 14
64         mov eax, cr2
65         push eax
66         push .msg
67         call panic
68         iret
69
70     .msg db "page fault at 0%x, error code: %x", 0
71     end_isr 14
72
```

```
73     ; PIT irq
74     begin_isr 32
75     ack_irq
76
77     push ebp
78     push dword 0
79     push dword 0
80     mov ebp, esp
81
82     call sched_switch
83
84     add esp, 8 ; fix up the fake stack frame we created
85     pop ebp
86
87     iret
88     end_isr 32
89
90     ; keyboard irq
91     begin_isr 33
92     ack_irq
93     iret
94     end_isr 33
95
96     ; spurious irq
97     begin_isr 39
98     iret
99     end_isr 39
100
101     ret
```

```
1 #include "kernel_page.h"
2 #include "paging.h"
3 #include "panic.h"
4 #include "string.h"
5
6 typedef union kernel_page
7 {
8     union kernel_page *next;
9     char mem[PAGE_SIZE];
10 }
11 kernel_page_t;
12
13 static kernel_page_t *allocated_to;
14
15 static kernel_page_t *end;
16
17 static kernel_page_t *next_free;
18
19 void kernel_page_init(virt_t begin, virt_t end_)
20 {
21     allocated_to = (void *)begin;
22     end = (void *)end_;
23 }
24
25 void *kernel_page_alloc()
26 {
27     if (next_free)
28     {
29         kernel_page_t *page = next_free;
30         next_free = next_free->next;
31         return page;
32     }
33
34     if (allocated_to >= end)
35     {
36         return NULL;
37     }
38
39     page_map(virt_t) allocated_to, page_alloc(), PE_PRESENT | PE_READ_WRITE);
40     void *retn = allocated_to++;
41
42     return retn;
43 }
44
45 void *kernel_page_alloc_zeroed()
46 {
47     void *page = kernel_page_alloc();
48     memset32(page, 0, 1024);
49     return page;
50 }
51
52 void kernel_page_free(void *ptr)
53 {
54     kernel_page_t *page = ptr;
55     page->next = next_free;
56     next_free = page;
57 }
58 }
```

```
1 #ifndef KALLOC_H
2 #define KALLOC_H
3
4 #include "types.h"
5
6 void kernel_page_init(virt_t begin, virt_t end);
7
8 void *kernel_page_alloc();
9
10 void *kernel_page_alloc_zeroed();
11
12 void kernel_page_free(void *);
13
14 #endif
```



```
1 global loader
2 global end_of_image
3
4 extern kmain
5
6 section .multiboot_header
7 align 4
8 multiboot_header:
9     dd 0x1badb002      ; magic
10     dd 3               ; flags
11     dd -(0x1badb002 + 3) ; checksum = -(flags + magic)
12
13 section .text
14 align 4
15 loader:
16     mov esp, stack
17     push dword 0
18     push dword 0
19     mov ebp, esp
20
21     push eax ; multiboot magic number
22     push ebx ; pointer to multiboot struct
23
24     fninit
25     mov eax, cr0
26     or eax, 1 << 5 ; FPU NE bit
27     mov cr0, eax
28
29     call kmain
30
31 section .bss
32 align 4
33 resb 65536
34 stack:
35
36 section .end_of_image
37 end_of_image:
```

```
1 #include "console.h"
2 #include "gdt.h"
3 #include "idt.h"
4 #include "kernel_page.h"
5 #include "multiboot.h"
6 #include "paging.h"
7 #include "panic.h"
8 #include "pit.h"
9 #include "sched.h"
10 #include "string.h"
11 #include "syscall.h"
12 #include "task.h"
13 #include "types.h"
14
15 static multiboot_info_t *mb;
16
17 static multiboot_module_t *find_module(const char *name)
18 {
19     multiboot_module_t *mods = (void *)mb->mods_addr;
20
21     for (size_t i = 0; i < mb->mods_count; i++)
22     {
23         if (strcmp((const char *)mods[i].cmdline, name))
24         {
25             return &mods[i];
26         }
27     }
28
29     return NULL;
30 }
31
32 void kmain(multiboot_info_t *mb_, uint32_t magic)
33 {
34     (void)magic;
35     mb = mb_;
36
37     console_init();
38
39     printf("Radium booting from %s.\n", (const char *)mb->boot_loader_name);
40
41     for (size_t i = 0; i < mb->mods_count; i++)
42     {
43         multiboot_module_t *mods = (void *)mb->mods_addr;
44         paging_set_allocatable_start(mods[i].mod_end);
45     }
46
47     gdt_init();
48     idt_init();
49     pit_set_frequency(100);
50     paging_init(mb);
51     task_init();
52     syscall_init();
53
54     multiboot_module_t *mod = find_module("/init.bin");
55
56     task_boot_init((const char *)mod->mod_start,
57                   mod->mod_end - mod->mod_start);
58 }
```

```
1 /* multiboot.h - Multiboot header file. */
2 /* Copyright (c) 1999,2003,2007,2008,2009 Free Software Foundation, Inc. *
3 Permission is hereby granted, free of charge, to any person obtaining a
4 copy * of this software and associated documentation files (the
5 "Software"), to * deal in the Software without restriction, including
6 without limitation the * rights to use, copy, modify, merge, publish,
7 distribute, sublicense, and/or * sell copies of the Software, and to permit
8 persons to whom the Software is * furnished to do so, subject to the
9 following conditions: * * The above copyright notice and this permission
10 notice shall be included in * all copies or substantial portions of the
11 Software. * * THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY
12 KIND, EXPRESS OR * IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF
13 MERCHANTABILITY, * FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT.
14 IN NO EVENT SHALL ANY * DEVELOPER OR DISTRIBUTOR BE LIABLE FOR ANY CLAIM,
15 DAMAGES OR OTHER LIABILITY, * WHETHER IN AN ACTION OF CONTRACT, TORT OR
16 OTHERWISE, ARISING FROM, OUT OF OR * IN CONNECTION WITH THE SOFTWARE OR THE
17 USE OR OTHER DEALINGS IN THE SOFTWARE. */
18
19 #ifndef MULTIBOOT_HEADER
20 #define MULTIBOOT_HEADER 1
21
22 /* * How many bytes from the start of the file we search for the header. */
23 #define MULTIBOOT_SEARCH 8192
24
25 /* * The magic field should contain this. */
26 #define MULTIBOOT_HEADER_MAGIC 0x1BADB002
27
28 /* * This should be in %eax. */
29 #define MULTIBOOT_BOOTLOADER_MAGIC 0x2BADB002
30
31 /* * The bits in the required part of flags field we don't support. */
32 #define MULTIBOOT_UNSUPPORTED 0x0000ffff
33
34 /* * Alignment of multiboot modules. */
35 #define MULTIBOOT_MOD_ALIGN 0x00001000
36
37 /* * Alignment of the multiboot info structure. */
38 #define MULTIBOOT_INFO_ALIGN 0x00000004
39
40 /* * Flags set in the 'flags' member of the multiboot header. */
41
42 /* * Align all boot modules on i386 page (4KB) boundaries. */
43 #define MULTIBOOT_PAGE_ALIGN 0x00000001
44
45 /* * Must pass memory information to OS. */
46 #define MULTIBOOT_MEMORY_INFO 0x00000002
47
48 /* * Must pass video information to OS. */
49 #define MULTIBOOT_VIDEO_MODE 0x00000004
50
51 /* * This flag indicates the use of the address fields in the header. */
52 #define MULTIBOOT_AOUT_KLUDGE 0x00010000
53
54 /* * Flags to be set in the 'flags' member of the multiboot info structure. */
55
56 /* * is there basic lower/upper memory information? */
57 #define MULTIBOOT_INFO_MEMORY 0x00000001
58
59 /* * is there a boot device set? */
60 #define MULTIBOOT_INFO_BOOTDEV 0x00000002
61
62 /* * is the command-line defined? */
63 #define MULTIBOOT_INFO_CMDLINE 0x00000004
64
65 /* * are there modules to do something with? */
66 #define MULTIBOOT_INFO_MODULES 0x00000008
67
68 /* * These next two are mutually exclusive */
69
70 /* * is there a symbol table loaded? */
71 #define MULTIBOOT_INFO_AOUT_SYMS 0x00000010
72
73 /* * is there an ELF section header table? */
74 #define MULTIBOOT_INFO_ELF_SHDR 0x00000020
75
76 /* * is there a full memory map? */
```

```
73 #define MULTIBOOT_INFO_MEM_MAP 0x00000040
74
75 /* * Is there drive info? */
76 #define MULTIBOOT_INFO_DRIVE_INFO 0x00000080
77
78 /* * Is there a config table? */
79 #define MULTIBOOT_INFO_CONFIG_TABLE 0x00000100
80
81 /* * Is there a boot loader name? */
82 #define MULTIBOOT_INFO_BOOT_LOADER_NAME 0x00000200
83
84 /* * Is there a APM table? */
85 #define MULTIBOOT_INFO_APM_TABLE 0x00000400
86
87 /* * Is there video information? */
88 #define MULTIBOOT_INFO_VIDEO_INFO 0x00000800
89
90 #ifndef ASM_FILE
91
92 typedef unsigned short multiboot_uint16_t;
93 typedef unsigned int multiboot_uint32_t;
94 typedef unsigned long long multiboot_uint64_t;
95
96 struct multiboot_header
97 {
98     /* * Must be MULTIBOOT_MAGIC - see above. */
99     multiboot_uint32_t magic;
100
101     /* * Feature flags. */
102     multiboot_uint32_t flags;
103
104     /* * The above fields plus this one must equal 0 mod 2^32. */
105     multiboot_uint32_t checksum;
106
107     /* * These are only valid if MULTIBOOT_AOUT_KLUDGE is set. */
108     multiboot_uint32_t header_addr;
109     multiboot_uint32_t load_addr;
110     multiboot_uint32_t load_end_addr;
111     multiboot_uint32_t bss_end_addr;
112     multiboot_uint32_t entry_addr;
113
114     /* * These are only valid if MULTIBOOT_VIDEO_MODE is set. */
115     multiboot_uint32_t mode_type;
116     multiboot_uint32_t width;
117     multiboot_uint32_t height;
118     multiboot_uint32_t depth;
119 };
120
121 /* * The symbol table for a.out. */
122 struct multiboot_aout_symbol_table
123 {
124     multiboot_uint32_t tabsize;
125     multiboot_uint32_t strsize;
126     multiboot_uint32_t addr;
127     multiboot_uint32_t reserved;
128 };
129
130 typedef struct multiboot_aout_symbol_table multiboot_aout_symbol_table_t;
131
132 /* * The section header table for ELF. */
133 struct multiboot_elf_section_header_table
134 {
135     multiboot_uint32_t num;
136     multiboot_uint32_t size;
137     multiboot_uint32_t addr;
138     multiboot_uint32_t shndx;
139 };
140
141 typedef struct multiboot_elf_section_header_table multiboot_elf_section_header_table_t;
142
143 struct multiboot_info
144 {
145     /* * Multiboot info version number */
```

```
145 multiboot_uint32_t flags;
146
147 /* Available memory from BIOS */
148 multiboot_uint32_t mem_lower;
149 multiboot_uint32_t mem_upper;
150
151 /* "root" partition */
152 multiboot_uint32_t boot_device;
153
154 /* Kernel command line */
155 multiboot_uint32_t cmdline;
156
157 /* Boot-Module List */
158 multiboot_uint32_t mods_count;
159 multiboot_uint32_t mods_addr;
160
161 union
162 {
163     multiboot_aout_symbol_table_t aout_sym;
164     multiboot_elf_section_header_table_t elf_sec;
165 } u;
166
167 /* Memory Mapping buffer */
168 multiboot_uint32_t mmap_length;
169 multiboot_uint32_t mmap_addr;
170
171 /* Drive Info buffer */
172 multiboot_uint32_t drives_length;
173 multiboot_uint32_t drives_addr;
174
175 /* ROM configuration table */
176 multiboot_uint32_t config_table;
177
178 /* Boot Loader Name */
179 multiboot_uint32_t boot_loader_name;
180
181 /* APM table */
182 multiboot_uint32_t apm_table;
183
184 /* Video */
185 multiboot_uint32_t vbe_control_info;
186 multiboot_uint32_t vbe_mode_info;
187 multiboot_uint16_t vbe_mode;
188 multiboot_uint16_t vbe_interface_seg;
189 multiboot_uint16_t vbe_interface_off;
190 multiboot_uint16_t vbe_interface_len;
191 };
192 typedef struct multiboot_info multiboot_info_t;
193
194 struct multiboot_mmap_entry
195 {
196     multiboot_uint32_t size;
197     multiboot_uint64_t addr;
198     multiboot_uint64_t len;
199     #define MULTIBOOT_MEMORY_AVAILABLE 1
200     #define MULTIBOOT_MEMORY_RESERVED 2
201     multiboot_uint32_t type;
202     __attribute__((packed))
203 } __attribute__((packed));
204 typedef struct multiboot_mmap_entry multiboot_memory_map_t;
205
206 struct multiboot_mod_list
207 {
208     /* the memory used goes from bytes 'mod_start' to 'mod_end-1' inclusive */
209     multiboot_uint32_t mod_start;
210     multiboot_uint32_t mod_end;
211
212     /* Module command line */
213     multiboot_uint32_t cmdline;
214
215     /* padding to take it to 16 bytes (must be zero) */
216     multiboot_uint32_t pad;
217 };
```

```
217 typedef struct multiboot_mod_list multiboot_module_t;
218
219 #endif /* ! ASM_FILE */
220
221 #endif /* ! MULTIBOOT_HEADER */
```

```

1 #include "console.h"
2 #include "paging.h"
3 #include "panic.h"
4 #include "string.h"
5
6 static phys_t next_free_page;
7
8 #define FL_PAGING_ENABLED (1 <= 31)
9
10 static bool paging_enabled()
11 {
12     uint32_t cr0;
13     __asm__ ("mov %0, cr0" : "=r"(cr0));
14     return !(cr0 & FL_PAGING_ENABLED);
15 }
16
17 void set_page_directory(phys_t page_directory)
18 {
19     uint32_t cr0;
20
21     __asm__ ("mov cr3, %0" : "=r"(page_directory));
22
23     __asm__ ("mov %0, cr0" : "=r"(cr0));
24     cr0 |= FL_PAGING_ENABLED;
25     __asm__ ("mov cr0, %0" : "=r"(cr0); "memory");
26 }
27
28 static void invlpg(virt_t virt)
29 {
30     __asm__ volatile ("invlpg [%0]" : : "r" (virt) : "memory");
31 }
32
33 phys_t page_alloc()
34 {
35     phys_t page = next_free_page;
36
37     if (paging_enabled())
38     {
39         phys_t *temp_mapping = page_temp_map(page);
40         next_free_page = *temp_mapping;
41         page_temp_unmap();
42     }
43     else
44     {
45         next_free_page = *(phys_t *) page;
46     }
47
48     return page;
49 }
50
51 void page_free(phys_t addr)
52 {
53     if (paging_enabled())
54     {
55         phys_t *temp_mapping = page_temp_map(addr);
56         *temp_mapping = next_free_page;
57         page_temp_unmap();
58     }
59     else
60     {
61         *(phys_t *) addr = next_free_page;
62     }
63
64     next_free_page = addr;
65 }
66
67 void page_map(virt_t virt_page, phys_t phys_page, page_flags_t flags)
68 {
69     // page directories are always recursively mapped into themselves:
70     uint32_t *current_page_directory = (uint32_t *) CURRENT_PAGE_DIRECTORY;
71
72     size_t page_dir_i = (virt_page / 4096) / 1024;

```

```

73     size_t page_tab_i = (virt_page / 4096) % 1024;
74
75     uint32_t *page_table =
76         (uint32_t *) (CURRENT_PAGE_TABLE_BASE + page_dir_i * PAGE_SIZE);
77
78     uint32_t pd_entry = current_page_directory[page_dir_i];
79     if (!(pd_entry & PE_PRESENT))
80     {
81         current_page_directory[page_dir_i] =
82             page_alloc() | (flags & PE_FLAG_MASK);
83         invlpg((virt_t) page_table);
84         memset(page_table, 0, 4096);
85     }
86
87     page_table[page_tab_i] =
88         (phys_page & PE_ADDR_MASK) | (flags & PE_FLAG_MASK);
89     invlpg((virt_t) virt_page);
90 }
91
92 void page_unmap(virt_t virt_page)
93 {
94     page_map(virt_page, 0, 0);
95 }
96
97 phys_t virt_to_phys(virt_t virt)
98 {
99     uint32_t *current_page_directory = (uint32_t *) CURRENT_PAGE_DIRECTORY;
100
101     size_t page_dir_i = (virt / 4096) / 1024;
102     size_t page_tab_i = (virt / 4096) % 1024;
103     size_t page_offset = virt % 4096;
104
105     if (!(current_page_directory[page_dir_i] & PE_PRESENT))
106     {
107         return 0;
108     }
109
110     uint32_t *page_table =
111         (uint32_t *) (CURRENT_PAGE_TABLE_BASE + page_dir_i * PAGE_SIZE);
112
113     uint32_t page_tab_entry = page_table[page_tab_i];
114
115     if (!(page_tab_entry & PE_PRESENT))
116     {
117         return 0;
118     }
119
120     return (page_tab_entry & PE_ADDR_MASK) + page_offset;
121 }
122
123 static uint32_t *const temp_page_entry = (void *)CURRENT_PAGE_TABLE_BASE;
124
125 static uint32_t old_null_page;
126
127 void *page_temp_map(phys_t phys_page)
128 {
129     old_null_page = *temp_page_entry;
130     *temp_page_entry = phys_page | PE_PRESENT | PE_READ_WRITE;
131     invlpg(0);
132     return NULL;
133     // we map in the temp page at NULL - todo
134     // fix...
135 }
136
137 void page_temp_unmap()
138 {
139     *temp_page_entry = old_null_page;
140     invlpg(0);
141 }
142
143 bool page_is_user_mapped(virt_t virt)
144 {
145     uint32_t dir_i = virt / (1024 * PAGE_SIZE);

```

```
145  uint32_t base_i = (virt / PAGE_SIZE);
146
147  uint32_t *page_directory = (uint32_t *) CURRENT_PAGE_DIRECTORY;
148
149  if (!page_directory[dir_i] & (PE_PRESENT | PE_USER)))
150  {
151      return false;
152  }
153
154  uint32_t *page_table_base = (uint32_t *) CURRENT_PAGE_TABLE_BASE;
155
156  if (!page_table_base[base_i] & (PE_PRESENT | PE_USER)))
157  {
158      return false;
159  }
160
161  return true;
162 }
```

```
1 #ifndef PAGING_H
2 #define PAGING_H
3
4 #include "types.h"
5 #include "multiboot.h"
6
7 #define PAGE_SIZE 4096
8
9 #define PE_FLAG_MASK (PAGE_SIZE - 1)
10 #define PE_ADDR_MASK (~PE_FLAG_MASK)
11
12 #define KERNEL_STACK_BEGIN 0x0fc00000ul
13 #define KERNEL_STACK_END 0x0fffffc0ul
14 #define USER_STACK_BEGIN 0x100000000ul
15 #define USER_STACK_END 0xffc00000ul
16 #define CURRENT_PAGE_DIRECTORY 0xfffff000ul
17 #define CURRENT_PAGE_TABLE_BASE 0xffc00000ul
18
19 typedef enum
20 {
21     PE_PRESENT = 1 << 0,
22     PE_READ_WRITE = 1 << 1,
23     PE_USER = 1 << 2,
24 }
25 page_flags_t;
26
27 void paging_set_allocatable_start(phys_t addr);
28
29 void paging_init(struct multiboot_info *mb);
30
31 void set_page_directory(phys_t page_directory);
32
33 phys_t page_alloc();
34
35 void page_free(phys_t addr);
36
37 void page_map(phys_t virt_page, phys_t phys_page, page_flags_t flags);
38
39 void page_unmap(phys_t virt_page);
40
41 phys_t virt_to_phys(phys_t virt);
42
43 void *page_temp_map(phys_t phys_page);
44
45 void page_temp_unmap();
46
47 bool page_is_user_mapped(phys_t virt);
48
49 #endif
```

```
1 #include "console.h"
2 #include "kernel_page.h"
3 #include "multiboot.h"
4 #include "paging.h"
5 #include "string.h"
6 #include "util.h"
7
8 static phys_t kernel_end;
9
10 static phys_t alloc_zeroed_page()
11 {
12     phys_t page = page_alloc();
13     memset32((void *)page, 0, 1024);
14     return page;
15 }
16
17 static size_t register_available_memory_region(multiboot_memory_map_t * region)
18 {
19     size_t pages_registered = 0;
20
21     for (size_t offset = 0; offset + PAGE_SIZE <= region->len;
22          offset += PAGE_SIZE)
23     {
24         phys_t addr = region->addr + offset;
25
26         if (addr < kernel_end)
27         {
28             // don't put memory before the end of the kernel in the free list
29             continue;
30         }
31
32         page_free(addr);
33         pages_registered++;
34     }
35
36     return pages_registered;
37 }
38
39 static size_t register_available_memory(multiboot_info_t * mb)
40 {
41     multiboot_memory_map_t *mmap = (void *)mb->mmap_addr;
42     size_t pages_registered = 0;
43
44     for (size_t i = 0; i < mb->mmap_length / sizeof(multiboot_memory_map_t);
45          i++)
46     {
47         if (mmap[i].type == MULTIBOOT_MEMORY_AVAILABLE)
48         {
49             pages_registered += register_available_memory_region(mmap + i);
50         }
51     }
52
53     return pages_registered;
54 }
55
56 static void create_page_tables_for_kernel_space(uint32_t * page_directory)
57 {
58     for (size_t i = 0; i < KERNEL_STACK_END / (4 * 1024 * 1024); i++)
59     {
60         page_directory[i] = alloc_zeroed_page() | PE_PRESENT | PE_READ_WRITE;
61     }
62 }
63
64 static void identity_map_kernel(uint32_t * page_directory)
65 {
66     // we start looping from PAGE_SIZE in order to leave the null page
67     // unmapped
68     // accessing it will cause a page fault.
69     for (phys_t addr = PAGE_SIZE; addr <= kernel_end; addr += PAGE_SIZE)
70     {
71         size_t page_dir_i = addr / 4096 / 1024;
72         size_t page_tab_i = addr / 4096 % 1024;
```

```
73 uint32_t page_dir_ent = page_directory[page_dir_i];
74 if (!page_dir_ent)
75 {
76     page_dir_ent = page_directory[page_dir_i] =
77         alloc_zeroed_page() | PE_PRESENT | PE_READ_WRITE;
78 }
79 uint32_t *page_tab = (uint32_t *) (page_dir_ent & PE_ADDR_MASK);
80 page_tab[page_tab_i] = addr | PE_PRESENT | PE_READ_WRITE;
81 }
82 }
83
84 static void recursively_map_page_directory(uint32_t * page_directory)
85 {
86     page_directory[1023] =
87         (phys_t) page_directory | PE_PRESENT | PE_READ_WRITE;
88 }
89
90 void paging_set_allocatable_start(phys_t addr)
91 {
92     if (addr > kernel_end)
93     {
94         kernel_end = round_up(addr, PAGE_SIZE);
95     }
96 }
97
98 void paging_init(multiboot_info_t * mb)
99 {
100     extern int end_of_image;
101     paging_set_allocatable_start((phys_t) &end_of_image);
102
103     size_t pages_registered = register_available_memory(mb);
104     printf("%d MiB available useful memory.\n",
105            pages_registered * PAGE_SIZE / 1024 / 1024);
106
107     uint32_t *page_directory = (uint32_t *) alloc_zeroed_page();
108
109     create_page_tables_for_kernel_space(page_directory);
110     identity_map_kernel(page_directory);
111     recursively_map_page_directory(page_directory);
112
113     set_page_directory((phys_t) page_directory);
114
115     kernel_page_init(kernel_end, KERNEL_STACK_BEGIN);
116 }
```



```
1 #include "console.h"
2 #include "panic.h"
3 #include "stdarg.h"
4 #include "types.h"
5
6 #define RECORD_SIZE 32
7
8 typedef struct
9 {
10     uint32_t addr;
11     char name[RECORD_SIZE - sizeof(uint32_t)];
12 }
13 symbol_record_t;
14
15 #define SYMBOL_ENTRIES (65536 / RECORD_SIZE)
16
17 extern symbol_record_t panic_symbols[SYMBOL_ENTRIES];
18
19 void panic_print_backtrace() __attribute__((noreturn));
20
21 void panic(const char *format, ...)
22 {
23     __asm__ volatile ("cli");
24
25     va_list va;
26     va_start(va, format);
27     printf("\npanic: ");
28     vprintf(format, va);
29     printf("\n");
30     va_end(va);
31
32     panic_print_backtrace();
33 }
34
35 void panic_print_backtrace_item(uint32_t addr)
36 {
37     uint32_t base = 0;
38     const char *name = "?";
39
40     for (size_t i = 0; i < SYMBOL_ENTRIES; i++)
41     {
42         if (panic_symbols[i].addr == 0)
43         {
44             break;
45         }
46         if (panic_symbols[i].addr < addr && panic_symbols[i].addr > base)
47         {
48             base = panic_symbols[i].addr;
49             name = panic_symbols[i].name;
50         }
51     }
52
53     printf(" [<0x%x>] %s +%d\n", addr, name, addr - base);
54 }
```

```
1 global panic_print_backtrace
2 global panic_symbols
3
4 extern printf
5 extern vprintf
6 extern panic_print_backtrace_item
7
8 section .rodata
9
10 panic_symbols:
11     db '@@@PANIC SYMBOL TABLE @@@'
12     times 65536 - ($ - panic_symbols) db 0
13
14 section .text
15 panic_print_backtrace:
16     mov esp, ebp
17     pop ebp
18     cmp dword [esp], 0
19     je .end_loop
20     call panic_print_backtrace_item
21     jmp panic_print_backtrace
22 .end_loop:
23     hlt
```

```
1 #ifndef PANIC_H
2 #define PANIC_H
3
4 void panic(const char *message, ...) __attribute__((noreturn));
5
6 #endif
```

```
1 #include "io.h"
2 #include "panic.h"
3 #include "pit.h"
4
5 #define PIT_CLOCK_FREQUENCY 1193180
6
7 void pit_set_frequency(uint32_t hz)
8 {
9     uint32_t divisor = PIT_CLOCK_FREQUENCY / hz;
10
11     if (divisor > 0xffff)
12     {
13         panic("frequency too low");
14     }
15
16     outb(0x43, 0x36);
17     outb(0x40, divisor & 0xff);
18     outb(0x40, (divisor >> 8) & 0xff);
19 }
```

```
1 #ifndef PIT_H
2 #define PIT_H
3
4 #include "types.h"
5
6 void pit_set_frequency(uint32_t hz);
7
8 #endif
```

```
1 #ifndef RADIUM_H
2 #define RADIUM_H
3
4 #define ENOSYS 1 // no such syscall
5 #define EFAULT 2 // bad address
6
7 // debug syscalls:
8 #define SYS_REGDUMP 0
9 #define SYS_CONSOLE_LOG 1
10
11 // process management syscalls:
12 #define SYS_EXIT 2
13 #define SYS_YIELD 3
14 #define SYS_FORK 4
15 #define SYS_WAIT 5
16
17 #endif
```

1	# Radium		
2			
3	Hobby OS.		
4			
5	## Memory Map		
6			
7	Begin	End (incl.)	Purpose
8	-----	-----	-----
9	`0000_0000`	`0000_0fff`	Null page. Not mapped in.
10	`0000_1000`	`000f_ffff`	Low memory. Untouched.
11	`0010_0000`	`*???*`	The kernel is loaded here by GRUB.
12	`*???*`	`0fbf_ffff`	Kernel dynamic allocation area.
13	`0fc0_0000`	`0fc0_0fff`	Kernel stack guard page. Not mapped in.
14	`0fc0_1000`	`0fff_ffff`	Kernel stack.
15	`1000_0000`	`ffbf_ffff`	User address space.
16	`ffc0_0000`	`ffff_efff`	Recursively mapped page tables
17	`ffff_f000`	`ffff_ffff`	Recursively mapped page directory

```
1 #ifndef SCHED_H
2 #define SCHED_H
3
4 void sched_begin() __attribute__((noreturn));
5
6 void sched_switch();
7
8 #endif
```



```
1 use32
2
3 global sched_begin
4 global sched_switch
5 global task_fork
6
7 extern current_task
8 extern sched_next
9 extern task_fork_inner
10
11 %define USER_CODE (0x18 | 3)
12 %define USER_DATA (0x20 | 3)
13
14 %define task_fpu_state(task) [(task) + 0]
15 %define task_esp(task) [(task) + 512]
16 %define task_eip(task) [(task) + 516]
17 %define task_kernel_stack(task) [(task) + 520]
18 %define task_page_dir_phys(task) [(task) + 528]
19
20 sched_begin:
21 mov ax, USER_DATA
22 mov ds, ax
23 mov es, ax
24 mov fs, ax
25 mov gs, ax
26
27 mov ecx, 0xffc00000 ; userland stack end
28 mov edx, 0x10000000 ; task entry point
29 sti
30 sysexit
31
32 sched_switch:
33 ; save old task state
34 pusha
35 mov eax, [current_task]
36 fxsave task_fpu_state(eax)
37 mov task_esp(eax), esp
38 mov task_eip(eax), dword .return
39
40 call sched_next
41 mov [current_task], eax
42
43 mov ebx, task_page_dir_phys(eax)
44 mov cr3, ebx
45
46 fxrstor task_fpu_state(eax)
47 mov esp, task_esp(eax)
48 jmp task_eip(eax)
49 .return:
50 popa
51 ret
52
53 task_fork:
54 xor eax, eax
55
56 pusha
57 call task_fork_inner
58
59 ; we need to save the return value of task_fork_inner so that the parent
60 ; task that calls this function gets the right task_t* for the newly
61 ; created child task:
62 mov [esp + 7*4], eax ; EAX from pusha
63
64 mov task_esp(eax), esp
65 mov task_eip(eax), dword .return
66 .return:
67 popa
68 ret
```

```

1  #include <radium.h>
2
3  #include "console.h"
4  #include "paging.h"
5  #include "panic.h"
6  #include "sched.h"
7  #include "syscall.h"
8  #include "task.h"
9  #include "util.h"
10
11 #define REG_VECTOR(regs) ((regs)->eax)
12 #define REG_RETURN(regs) ((regs)->eax)
13
14 #define REG_ARG1(regs) ((regs)->ebx)
15 #define REG_ARG2(regs) ((regs)->edx)
16 #define REG_ARG3(regs) ((regs)->esi)
17
18 static bool valid_user_buffer(virt_t ptr, size_t len)
19 {
20     if ((0xffffffff - len) < ptr)
21     {
22         return false;
23     }
24
25     size_t page_offset = ptr % PAGE_SIZE;
26
27     ptr -= page_offset;
28     len += page_offset;
29
30     virt_t end = ptr + len;
31
32     while (ptr < end)
33     {
34         if (!page_is_user_mapped(ptr))
35         {
36             return false;
37         }
38
39         ptr += PAGE_SIZE;
40     }
41
42     return true;
43 }
44
45 static uint32_t syscall_regdump(registers_t * regs)
46 {
47     printf("process %d register dump:\n", current_task->pid);
48     printf("    eax=%x, ebx=%x, ecx=%x, edx=%x\n", regs->eax, regs->ebx,
49           regs->ecx, regs->edx);
50     printf("    esp=%x, ebp=%x, esi=%x, edi=%x\n", regs->esp, regs->ebp,
51           regs->esi, regs->edi);
52
53     return 0;
54 }
55
56 static uint32_t syscall_exit(registers_t * regs)
57 {
58     uint8_t status = REG_ARG1(regs) & 0xff;
59
60     if (current_task->pid == 1)
61     {
62         panic("init exited with status: %d", status);
63     }
64
65     task_kill(current_task, status);
66
67     // goodbye!
68     sched_switch();
69
70     return 0;
71 }
72

```

```

73 static uint32_t syscall_yield(registers_t * regs)
74 {
75     (void)regs;
76     sched_switch();
77     return 0;
78 }
79
80 static uint32_t syscall_fork(registers_t * regs)
81 {
82     (void)regs;
83
84     task_t *new_task = task_fork();
85
86     if (new_task)
87     {
88         // parent
89         return new_task->pid;
90     }
91     else
92     {
93         // child
94         return 0;
95     }
96 }
97
98 static uint32_t syscall_wait(registers_t * regs)
99 {
100     if (REG_ARG1(regs) != 0 && !valid_user_buffer(REG_ARG1(regs), sizeof(int)))
101     {
102         return -EFAULT;
103     }
104
105     int *stat_loc = (int *)REG_ARG1(regs);
106
107     again:
108     if (current_task->wait_queue.live.head)
109     {
110         task_t *child = current_task->wait_queue.live.head;
111         current_task->wait_queue.live.head = child->wait_queue.dead.next;
112         if (!current_task->wait_queue.live.head)
113         {
114             current_task->wait_queue.live.tail = NULL;
115         }
116
117         uint32_t child_pid = child->pid;
118         uint8_t child_status = child->exit_status;
119
120         task_destroy(child);
121
122         if (stat_loc)
123         {
124             *stat_loc = child_status;
125         }
126
127         return child_pid;
128     }
129
130     current_task->state = TASK_BLOCK_WAIT;
131     sched_switch();
132     goto again;
133 }
134
135 static uint32_t syscall_console_log(registers_t * regs)
136 {
137     if (valid_user_buffer(REG_ARG1(regs), REG_ARG2(regs)))
138     {
139         console_puts((const char *)REG_ARG1(regs), REG_ARG2(regs));
140         return 0;
141     }
142     else
143     {
144         return -EFAULT;
145     }
146 }
147

```

```
145     }
146 }
147
148 typedef uint32_t(syscall_t) (registers_t *);
149
150 static syscall_t *syscall_table[] = {
151     [SYS_REGDUMP] = syscall_regdump,
152     [SYS_EXIT] = syscall_exit,
153     [SYS_YIELD] = syscall_yield,
154     [SYS_FORK] = syscall_fork,
155     [SYS_WAIT] = syscall_wait,
156     [SYS_CONSOLE_LOG] = syscall_console_log,
157 };
158
159 void syscall_dispatch(registers_t * regs)
160 {
161     if (REG_VECTOR(regs) > countof(syscall_table))
162     {
163         REG_RETURN(regs) = -ENOSYS;
164         return;
165     }
166
167     syscall_t *func = syscall_table[regs->eax];
168
169     if (!func)
170     {
171         REG_RETURN(regs) = -ENOSYS;
172         return;
173     }
174
175     current_task->syscall_registers = regs;
176
177     REG_RETURN(regs) = func(regs);
178
179     current_task->syscall_registers = NULL;
180 }
```

```
1 use32
2
3 extern syscall_dispatch
4
5 global syscall_entry
6 global syscall_init
7
8 %define IA32_SYSENTER_CS 0x174
9 %define IA32_SYSENTER_ESP 0x175
10 %define IA32_SYSENTER_EIP 0x176
11
12 %define KERNEL_CODE 0x08
13 %define KERNEL_DATA 0x10
14
15 syscall_init:
16     xor edx, edx
17
18     mov ecx, IA32_SYSENTER_CS
19     mov eax, KERNEL_CODE
20     wrmsr
21
22     mov ecx, IA32_SYSENTER_ESP
23     mov eax, 0x0fffffc
24     wrmsr
25
26     mov ecx, IA32_SYSENTER_EIP
27     mov eax, syscall_entry
28     wrmsr
29
30     ret
31
32 syscall_entry:
33     pusha
34
35     push dword 0
36     push dword 0
37     mov ebp, esp
38
39     lea eax, [esp + 8]
40     push eax
41     call syscall_dispatch
42     add esp, 12 ; 4 bytes for the argument we passed to syscall_dispatch, and
43                 ; another 8 to compensate for the fake stack frame we pushed.
44
45     popa
46
47     ; STI apparently waits one instruction before enabling interrupts, so
48     ; despite how it appears, this return sequence should be race-free.
49     sti
50     sysexit
```

```
1 #ifndef SYSCALL_H
2 #define SYSCALL_H
3
4 #include "types.h"
5
6 typedef struct
7 {
8     uint32_t edi;
9     uint32_t esi;
10    uint32_t ebp;
11    uint32_t esp;
12    uint32_t ebx;
13    uint32_t edx;
14    uint32_t ecx;
15    uint32_t eax;
16 } registers_t;
17
18 void syscall_init();
19
20 #endif
```

```

1  #include "console.h"
2  #include "gdt.h"
3  #include "kernel_page.h"
4  #include "paging.h"
5  #include "panic.h"
6  #include "sched.h"
7  #include "string.h"
8  #include "task.h"
9  #include "util.h"
10
11 static_assert(tss_t.is_0x68_bytes_long, sizeof(tss_t) == 0x68);
12
13 static tss_t tss;
14
15 task_t *current_task;
16
17 static task_t *tasks[1024];
18
19 task_t *task_for_pid(uint32_t pid)
20 {
21     if (pid > countof(tasks))
22     {
23         return NULL;
24     }
25     return tasks[pid];
26 }
27
28 static task_t *alloc_empty_task()
29 {
30     for (uint32_t pid = 1; pid < countof(tasks); pid++)
31     {
32         if (tasks[pid])
33         {
34             continue;
35         }
36     }
37
38     tasks[pid] = kernel_page_alloc_zeroed();
39     tasks[pid]->state = TASK_READY;
40     tasks[pid]->pid = pid;
41     tasks[pid]->wait_queue.live.head = NULL;
42     tasks[pid]->wait_queue.live.tail = NULL;
43     return tasks[pid];
44 }
45
46 return NULL;
47 }
48
49 static void create_skeleton_page_directory(task_t * task)
50 {
51     // initialise task page directory
52     uint32_t *current_page_directory = (uint32_t *) CURRENT_PAGE_DIRECTORY;
53     uint32_t *task_page_directory = kernel_page_alloc_zeroed();
54
55     if (!task_page_directory)
56     {
57         panic("could not allocate page for new task's page directory");
58     }
59
60     for (size_t i = 0; i < KERNEL_STACK_BEGIN / (4 * 1024 * 1024); i++)
61     {
62         task_page_directory[i] = current_page_directory[i];
63     }
64
65     task->page_directory = task_page_directory;
66     task->page_directory_phys = virt_to_phys((virt_t) task_page_directory);
67     task->page_directory[1023] =
68         task->page_directory_phys | PE_PRESENT | PE_READ_WRITE;
69
70     uint32_t *kernel_stack_page_table = kernel_page_alloc_zeroed();
71     task->page_directory[KERNEL_STACK_BEGIN / (4 * 1024 * 1024)] =
72     virt_to_phys((virt_t) kernel_stack_page_table) | PE_PRESENT |

```

```

73     PE_READ_WRITE;
74
75     task->kernel_stack = kernel_page_alloc_zeroed();
76     kernel_stack_page_table[1023] =
77     virt_to_phys((virt_t) task->kernel_stack) | PE_PRESENT | PE_READ_WRITE;
78 }
79
80 void task_init()
81 {
82     gdt_set_tss(GDT_TSS, (uint32_t) & tss, sizeof(tss));
83     gdt_reload();
84
85     tss.ss0 = GDT_KERNEL_DATA;
86     tss.esp0 = KERNEL_STACK_END;
87
88     // pointer to the IO permission bitmap is beyond the end of the segment
89     tss.iopb = sizeof(tss);
90
91     __asm__ volatile ("ltr ax::\"a\" ((uint16_t) GDT_TSS | 3));
92
93     task_t *init_task = alloc_empty_task();
94     create_skeleton_page_directory(init_task);
95
96     // user stack
97     uint32_t *user_stack_page_table = kernel_page_alloc_zeroed();
98     init_task->page_directory[1022] =
99     virt_to_phys((virt_t) user_stack_page_table) | PE_PRESENT |
100     PE_READ_WRITE | PE_USER;
101     user_stack_page_table[1023] =
102     page_alloc() | PE_PRESENT | PE_READ_WRITE | PE_USER;
103
104     init_task->ppid = 0;
105
106     current_task = init_task;
107 }
108
109 void task_boot_init(const char *text, size_t size)
110 {
111     task_t *init_task = task_for_pid(1);
112
113     set_page_directory(init_task->page_directory_phys);
114
115     for (size_t i = 0; i < size; i += PAGE_SIZE)
116     {
117         phys_t page = page_alloc();
118         page_map(USER_BEGIN + i, page, PE_PRESENT | PE_USER);
119         size_t copy_size = size - i;
120         if (copy_size > PAGE_SIZE)
121         {
122             copy_size = PAGE_SIZE;
123         }
124         memcpy((void *) (USER_BEGIN + i), (void *) (text + i), copy_size);
125     }
126
127     sched_begin();
128 }
129
130 static void copy_userland_pages(task_t * new_task)
131 {
132     uint32_t *current_page_directory = (uint32_t *) CURRENT_PAGE_DIRECTORY;
133
134     for (size_t dir_i = USER_BEGIN / (PAGE_SIZE * 1024); dir_i < 1023; dir_i++)
135     {
136         if (!current_page_directory[dir_i])
137         {
138             continue;
139         }
140
141         uint32_t *current_page_table =
142             (uint32_t *) (CURRENT_PAGE_TABLE_BASE + dir_i * PAGE_SIZE);
143         uint32_t *new_page_table = kernel_page_alloc();
144

```

```
145 new_task->page_directory[dir_i] =
146     virt_to_phys((virt_t) new_page_table |
147     (current_page_directory[dir_i] & PE_FLAG_MASK);
148
149 for (size_t tab_i = 0; tab_i < 1024; tab_i++)
150 {
151     uint32_t current_entry = current_page_table[tab_i];
152     void *current_virt =
153         (void *)((dir_i * PAGE_SIZE * 1024) + (tab_i * PAGE_SIZE));
154
155     if (! (current_entry & PE_PRESENT))
156     {
157         continue;
158     }
159
160     phys_t new_page = page_alloc();
161
162     void *new_page_mapping = page_temp_map(new_page);
163     memcpy(new_page_mapping, current_virt, PAGE_SIZE);
164     page_temp_unmap();
165
166     new_page_table[tab_i] = new_page | (current_entry & PE_FLAG_MASK);
167 }
168
169 }
170
171 task_t *task_fork_inner()
172 {
173     task_t *new_task = alloc_empty_task();
174     create_skeleton_page_directory(new_task);
175
176     memcpy(&new_task->fpu_state, &current_task->fpu_state,
177     sizeof(new_task->fpu_state));
178     memcpy(new_task->kernel_stack, current_task->kernel_stack, PAGE_SIZE);
179
180     copy_userland_pages(new_task);
181
182     new_task->ppid = current_task->pid;
183
184     new_task->syscall_registers = current_task->syscall_registers;
185
186     return new_task;
187 }
188
189 void task_kill(task_t * task, uint8_t status)
190 {
191     task->exit_status = status;
192
193     task_t *parent = task_for_pid(task->ppid);
194
195     // if this task has any children in the wait queue, just clean them up as
196     // nothing is interested in them anymore.
197     task_t *waitable_child = task->wait_queue.live.head;
198     while (waitable_child)
199     {
200         task_t *next_waitable_child = waitable_child->wait_queue.dead.next;
201         task_destroy(waitable_child);
202         waitable_child = next_waitable_child;
203     }
204
205     // reparent children - T000 don't scan all processes...
206     for (uint32_t i = 2; i < countof(tasks); i++)
207     {
208         task_t *child = task_for_pid(i);
209         if (child && child->ppid == task->pid)
210         {
211             child->ppid = 1;
212         }
213     }
214
215     // insert this task into the parent's wait queue
216     if (parent->wait_queue.live.tail)
```

```
217 {
218     parent->wait_queue.live.tail->wait_queue.dead.next = task;
219     parent->wait_queue.live.tail = task;
220 }
221
222 else
223 {
224     parent->wait_queue.live.head = task;
225     parent->wait_queue.live.tail = task;
226 }
227
228 task->wait_queue.dead.next = NULL;
229
230 // wake parent up if it's blocked in wait()
231 if (parent->state == TASK_BLOCK_WAIT)
232 {
233     parent->state = TASK_READY;
234 }
235
236 // set state to EXITING so the scheduler never reschedules this task
237 current_task->state = TASK_EXITING;
238
239 void task_destroy(task_t * task)
240 {
241     tasks[task->pid] = NULL;
242
243     kernel_page_free(task->kernel_stack);
244
245     // T000 - free pages allocated for the task that are referenced from its
246     // page directory
247
248     kernel_page_free(task->page_directory);
249     kernel_page_free(task);
250 }
251
252 task_t *sched_next()
253 {
254     uint32_t current_pid = current_task->pid;
255
256     for (size_t i = current_pid + 1; i < countof(tasks); i++)
257     {
258         task_t *task = task_for_pid(i);
259
260         if (!task)
261         {
262             continue;
263         }
264
265         if (task->state == TASK_READY)
266         {
267             return task;
268         }
269     }
270
271     for (size_t i = 1; i <= current_pid; i++)
272     {
273         task_t *task = task_for_pid(i);
274
275         if (!task)
276         {
277             continue;
278         }
279
280         if (task->state == TASK_READY)
281         {
282             return task;
283         }
284     }
285
286     panic("no tasks ready to schedule!");
287 }
```

```
1 #ifndef TASK_H
2 #define TASK_H
3
4 #include "types.h"
5 #include "syscall.h"
6
7 typedef struct
8 {
9     // 0x00
10    uint16_t link;
11    uint16_t _res_1;
12    // 0x04
13    uint32_t esp0;
14    // 0x08
15    uint16_t ss0;
16    uint16_t _res_2;
17    // 0x0c
18    uint32_t esp1;
19    // 0x10
20    uint16_t ss1;
21    uint16_t _res_3;
22    // 0x14
23    uint32_t esp2;
24    // 0x18
25    uint16_t ss2;
26    uint16_t _res_4;
27    // 0x1c
28    uint32_t cr3;
29    // 0x20
30    uint32_t eip;
31    // 0x24
32    uint32_t eflags;
33    // 0x28
34    uint32_t eax;
35    // 0x2c
36    uint32_t ecx;
37    // 0x30
38    uint32_t edx;
39    // 0x34
40    uint32_t ebx;
41    // 0x38
42    uint32_t esp;
43    // 0x3c
44    uint32_t ebp;
45    // 0x40
46    uint32_t esi;
47    // 0x44
48    uint32_t edi;
49    // 0x48
50    uint16_t es;
51    uint16_t _res_5;
52    // 0x4c
53    uint16_t cs;
54    uint16_t _res_6;
55    // 0x50
56    uint16_t ss;
57    uint16_t _res_7;
58    // 0x54
59    uint16_t ds;
60    uint16_t _res_8;
61    // 0x58
62    uint16_t fs;
63    uint16_t _res_9;
64    // 0x5c
65    uint16_t gs;
66    uint16_t _res_10;
67    // 0x60
68    uint16_t ldtr;
69    uint16_t _res_11;
70    // 0x64
71    uint16_t _res_12;
72    uint16_t iobp;
```

```
73 // 0x68
74 } __attribute__((packed)) tss_t;
75
76 typedef enum
77 {
78     TASK_READY = 1,
79     TASK_BLOCK_WAIT = 2,
80     TASK_EXITING = 3,
81 }
82 task_state_t;
83
84 // sched.asm refers to hardcoded offsets within this struct.
85 // make sure to change it when changing anything here.
86 typedef struct task
87 {
88     /* 0 */ uint8_t fpu_state[512];
89     /* 512 */ uint32_t esp;
90     /* 516 */ uint32_t eip;
91     /* 520 */ void *kernel_stack;
92     // allocated within the kernel's identity-mapped region:
93     /* 524 */ uint32_t *page_directory;
94     /* 528 */ phys_t page_directory_phys;
95
96     // struct members past this point may be freely rearranged without needing
97     // to update any assembly source.
98
99     registers_t *syscall_registers;
100
101     uint32_t pid;
102     uint32_t ppid;
103
104     task_state_t state;
105     uint8_t exit_status;
106
107     union
108     {
109         struct
110         {
111             struct task *head;
112             struct task *tail;
113         } live;
114         struct
115         {
116             struct task *next;
117         } dead;
118     } wait_queue;
119
120     task_t;
121
122     extern task_t *current_task;
123
124     void task_init();
125
126     void task_boot_init(const char *init_bin, size_t size) __attribute__((noreturn));
127
128     task_t *task_fork();
129
130     void task_kill(task_t * task, uint8_t status);
131
132     void task_destroy(task_t * task);
133
134     task_t *task_for_pid(uint32_t pid);
135
136 #endif
```



```
1 #ifndef TYPES_H
2 #define TYPES_H
3
4 typedef unsigned char uint8_t;
5 typedef unsigned short uint16_t;
6 typedef unsigned int uint32_t;
7 typedef unsigned long long uint64_t;
8
9 typedef signed char int8_t;
10 typedef signed short int16_t;
11 typedef signed int int32_t;
12 typedef signed long long int64_t;
13
14 typedef uint32_t size_t;
15
16 typedef _Bool bool;
17 #define true 1
18 #define false 0
19
20 #define NULL ((void*)0)
21
22 typedef uint32_t virt_t;
23 typedef uint32_t phys_t;
24
25 #endif
```

```
1 #ifndef UTIL_H
2 #define UTIL_H
3
4 #include "types.h"
5
6 #define static_assert(name, expr) static char __static_assert_##name[(expr) ? 1 : -1]
7 __attribute__((unused))
8
9 #define countof(x) (sizeof(x) / sizeof(*(x)))
10
11 static inline size_t round_down(size_t val, size_t divisor)
12 {
13     return val - val % divisor;
14 }
15
16 static inline size_t round_up(size_t val, size_t divisor)
17 {
18     return round_down(val + divisor - 1, divisor);
19 }
20 #endif
```