# x86 Memory management

*What was it you were talking about?*

Communications Security Establishment    Centre de la sécurité des télécommunications
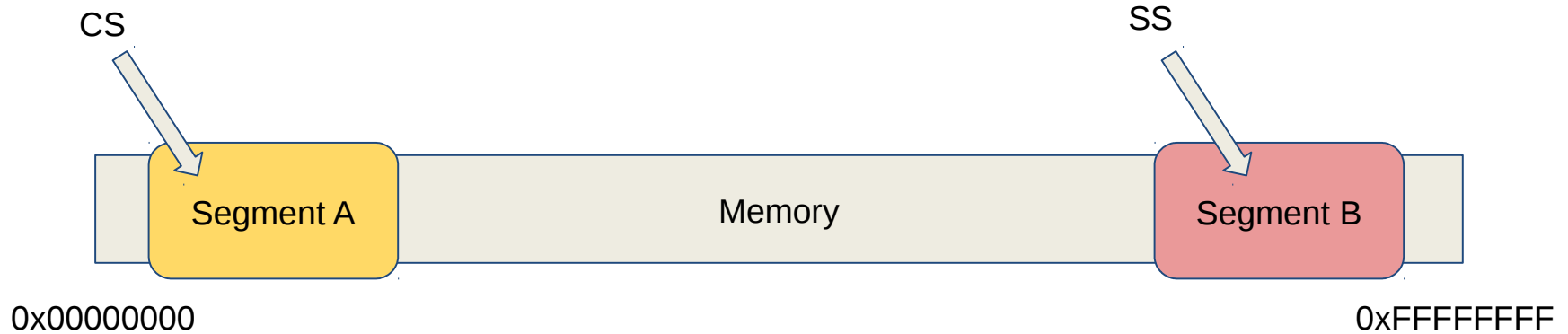
Canada

# So, exactly what is the difference between segmentation and virtual memory?

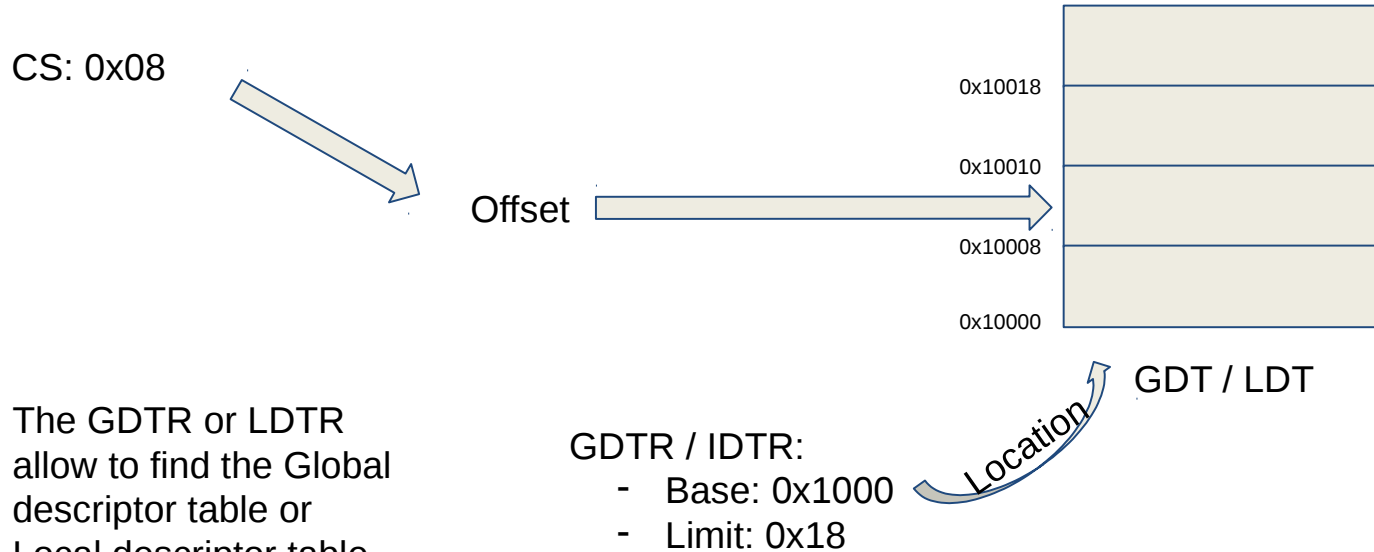# What is segmentation?

CS

SS

Segment A

Memory

Segment B

0x00000000

0xFFFFFFFF

A segment defines a **zone of memory**. A segment register then points to that zone allowing access to it.

Communications Security Establishment

Centre de la sécurité des télécommunications

Canadä

# Segmentation setup x86



CS: 0x08

Offset

0x10018

0x10010

0x10008

0x10000

GDT / LDT

The GDTR or LDTR allow to find the Global descriptor table or Local descriptor table on the memory.

GDTR / IDTR:
- Base: 0x1000
- Limit: 0x18

Location

Communications Security Establishment

Centre de la sécurité des télécommunications

Canada

# On Linux

traps.c

```
1015  ▼ #ifdef CONFIG_X86_32
1016  »          set_system_intr_gate(IA32_SYSCALL_VECTOR, entry_INT80_32);
1017  »          set_bit(IA32_SYSCALL_VECTOR, used_vectors);
1018    #endif
1019
1020  ▼ »        /*
1021  »          * Set the IDT descriptor to a fixed read-only location, so that the
1022  »          * "sidt" instruction will not leak the location of the kernel, and
1023  »          * to defend the IDT against arbit
1024  »          * It will be reloaded in cpu_init    1576    void cpu_init(void)
1025  »          __set_fixmap(FIX_RO_IDT, __pa_symb  1577  ▼ {
1026  »          idt_descr.address = fix_to_virt(FI  1578  »            int cpu = smp_processor_id();
1027                                                1579  »            struct task_struct *curr = current;
1028  ▼ »        /*                                  1580  »            struct tss_struct *t = &per_cpu(cpu_tss, cpu);
1029  »          * Should be a barrier for any ext  1581  »            struct thread_struct *thread = &curr->thread;
1030  »          */                                 1582
1031  »          cpu_init();                        1583  »            wait_for_master_cpu(cpu);
                                                    1584
                                                    1585  ▼ »          /*
                                                    1586  »            * Initialize the CR4 shadow before doing anything that could
                                                    1587  »            * try to read it.
                                                    1588  »            */
                                                    1589  »            cr4_init_shadow();
                                                    1590
                                                    1591  »            show_ucode_info_early();
                                                    1592
                                                    1593  »            pr_info("Initializing CPU#%d\n", cpu);
                                                    1594
                                                    1595  »            if (cpu_feature_enabled(X86_FEATURE_VME) ||
                                                    1596  »               boot_cpu_has(X86_FEATURE_TSC) ||
                                                    1597  »               boot_cpu_has(X86_FEATURE_DE))
                                                    1598  »     »          cr4_clear_bits(X86_CR4_VME|X86_CR4_PVI|X86_CR4_TSD|X86_CR4_DE);
                                                    1599
                                                    1600  »            load_current_idt();
                                                    1601  »            switch_to_new_gdt(cpu);
                                                    1602
```

common.c

# On Linux

common.c

```
500     void switch_to_new_gdt(int cpu)
501   ▼ {
502   »       /* Load the original GDT */
503   »       load_direct_gdt(cpu);
504   »       /* Reload the per-cpu base */
505   »       load_percpu_segment(cpu);
506     }
507
```

```
475     void load_direct_gdt(int cpu)
476   ▼ {
477   »       struct desc_ptr gdt_descr;
478
479   »       gdt_descr.address = (long)get_cpu_gdt_rw(cpu);
480   »       gdt_descr.size = GDT_SIZE - 1;
481   »       load_gdt(&gdt_descr);
482     }
```

desc.h

```
125     #define load_gdt(dtr)»   »         »         »         native_load_gdt(dtr)
```

```
43  ▼ struct gdt_page {
44  »       struct desc_struct gdt[GDT_ENTRIES];
45  } __attribute__((aligned(PAGE_SIZE)));
46
47  DECLARE_PER_CPU_PAGE_ALIGNED(struct gdt_page, gdt_page);
48
49  /* Provide the original GDT */
50  static inline struct desc_struct *get_cpu_gdt_rw(unsigned int cpu)
51  ▼ {
52  »       return per_cpu(gdt_page, cpu).gdt;|
53  }
54
```

```
236     static inline void native_load_gdt(const
237   ▼ {
238   »       asm volatile("lgdt %0"::"m" (*dt
239     }
```

# On Linux

Bolded bits shows the DPL (descriptor privilege level).
0 = Kernel mode, 3 = User mode.
Definite proof that Linux only uses 2 of the 4 rings available on x86!

common.c

```
101 ▾ DEFINE_PER_CPU_PAGE_ALIGNED(struct gdt_page, gdt_page) = { .gdt = {
102 ▾ #ifdef CONFIG_X86_64
103 ▾ »      /*
104   »       * We need valid kernel segments for data and code in long mode too
105   »       * IRET will check the segment types  kkeil 2000/10/28
106   »       * Also sysret mandates a special GDT layout
107   »       *
108   »       * TLS descriptors are currently at a different place compared to i386.
109   »       * Hopefully nobody expects them at a fixed place (Wine?)
110   »       */
111   »      [GDT_ENTRY_KERNEL32_CS]»»»      = GDT_ENTRY_INIT(0xc09b, 0, 0xfffff),
112   »      [GDT_ENTRY_KERNEL_CS]» »        = GDT_ENTRY_INIT(0xa09b, 0, 0xfffff),
113   »      [GDT_ENTRY_KERNEL_DS]» »        = GDT_ENTRY_INIT(0xc093, 0, 0xfffff),
114   »      [GDT_ENTRY_DEFAULT_USER32_CS]   = GDT_ENTRY_INIT(0xc0fb, 0, 0xfffff),
115   »      [GDT_ENTRY_DEFAULT_USER_DS]     = GDT_ENTRY_INIT(0xc0f3, 0, 0xfffff),
116   »      [GDT_ENTRY_DEFAULT_USER_CS]     = GDT_ENTRY_INIT(0xa0fb, 0, 0xfffff),
117 #else
118   »      [GDT_ENTRY_KERNEL_CS]» »        = GDT_ENTRY_INIT(0xc09a, 0, 0xfffff),
119   »      [GDT_ENTRY_KERNEL_DS]» »        = GDT_ENTRY_INIT(0xc092, 0, 0xff
120   »      [GDT_ENTRY_DEFAULT_USER_CS]»    = GDT_ENTRY_INIT(0xc0fa, 0, 0xff
121   »      [GDT_ENTRY_DEFAULT_USER_DS]»    = GDT_ENTRY_INIT(0xc0f2, 0, 0xff
```

Looking at the flags for the 32 bits version, you should see the following binary pattern for Kernel mode code:

0b 1100 0000 **100**1 1010

Looking at the flags for the 32 bits version, you should see the following binary pattern for User mode code:

0b 1100 0000 1**111** 1010

```
14 ▾ /*
15   * FIXME: Accessing the desc_struct through its fields is more elegant,
16   * and should be the one valid thing to do. However, a lot of open code
17   * still touches the a and b accessors, and doing this allow us to do it
18   * incrementally. We keep the signature as a struct, rather than a union,
19   * so we can get rid of it transparently in the future -- glommer
20   */
21   /* 8 byte segment descriptor */
22 ▾ struct desc_struct {
23 ▾ »     union {
24 ▾ »     »     struct {
25   » »     »           unsigned int a;
26   » »     »           unsigned int b;
27   » »     };
28 ▾ » »     struct {
29   » »     »           u16 limit0;
30   » »     »           u16 base0;
31   » »     »           unsigned base1: 8, type: 4, s: 1, dpl: 2, p: 1;
32   » »     »           unsigned limit: 4, avl: 1, l: 1, d: 1, g: 1, base2: 8;
33   » »     };
34   » };
35 } __attribute__((packed));
36
37 #define GDT_ENTRY_INIT(flags, base, limit) { { { \
38   »      .a = ((limit) & 0xffff) | (((base) & 0xffff) << 16), \
39   »      .b = (((base) & 0xff0000) >> 16) | (((flags) & 0xf0ff) << 8) | \
40   »           ((limit) & 0xf0000) | ((base) & 0xff000000), \
41   »      } } }
42
```

desc_defs.h

Communications Security Establishment  Centre de la sécurité des télécommunications

# Alright, let's write an OS…
## Just kidding, let's talk about virtual memory...

Communications Security Establishment    Centre de la sécurité des télécommunications
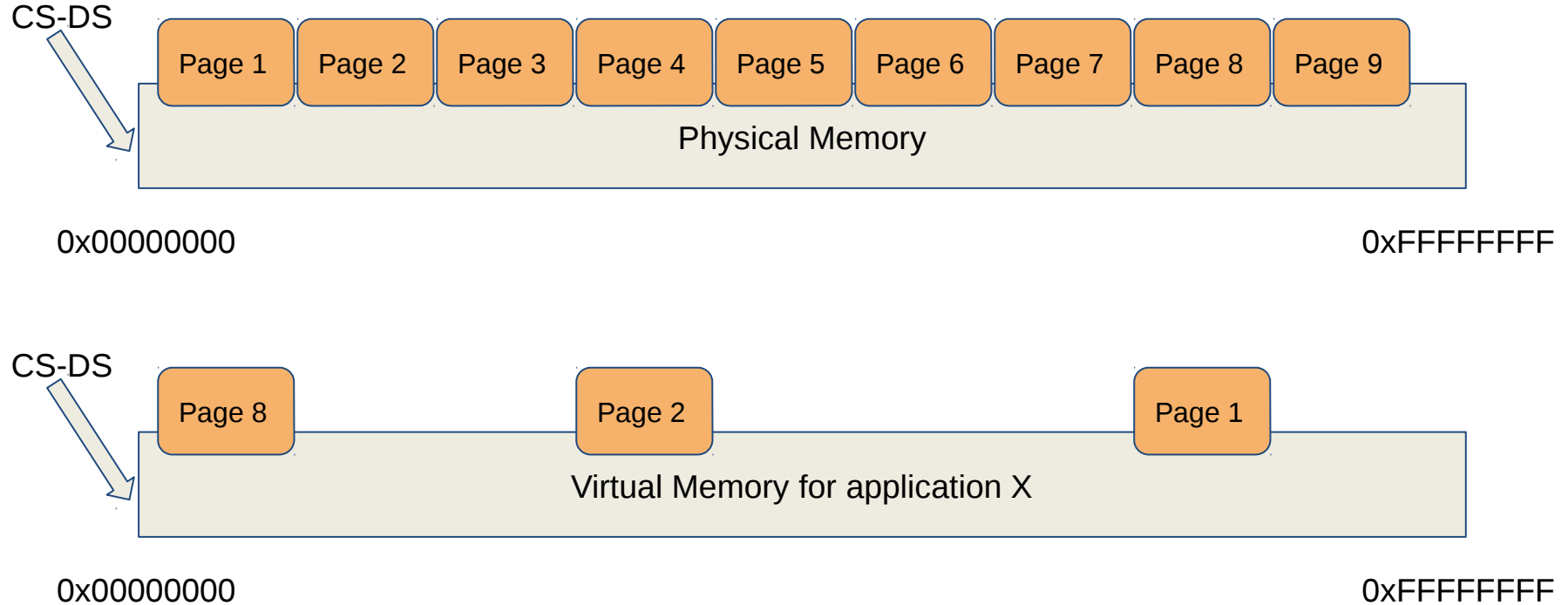
On x86, granularity in memory management is made possible by the use of virtual memory.

Be careful, Intel calls that "paging". Not to be confused with paging/swapping.
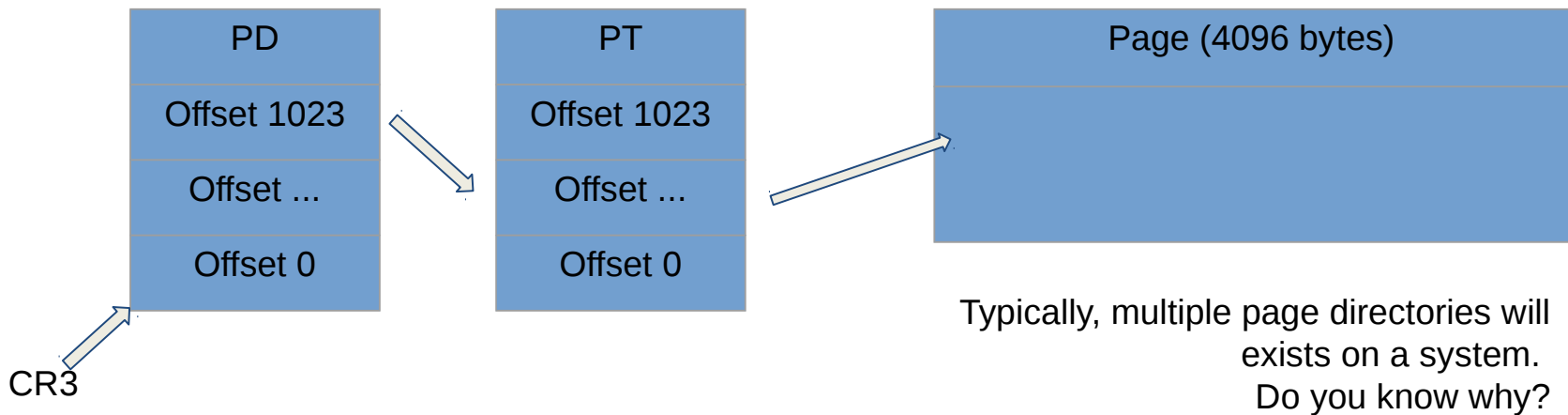
# Virtual memory

CS-DS

| Page 1 | Page 2 | Page 3 | Page 4 | Page 5 | Page 6 | Page 7 | Page 8 | Page 9 |

Physical Memory

0x00000000

0xFFFFFFFF

CS-DS

| Page 8 | | Page 2 | | Page 1 |

Virtual Memory for application X

0x00000000

0xFFFFFFFF

Communications Security Establishment

Centre de la sécurité des télécommunications

Canada

# Virtual Address

| 31 - 22 | 21 - 12 | 11 - 0 |
|---|---|---|
| Page Directory Offset | Page Table Offset | Page Offset |

| PD |
|---|
| Offset 1023 |
| Offset ... |
| Offset 0 |

| PT |
|---|
| Offset 1023 |
| Offset ... |
| Offset 0 |

| Page (4096 bytes) |
|---|
| |

CR3

Typically, multiple page directories will exists on a system.
Do you know why?

# On Linux

head_32.s

```
128    #endif
129
130    »        /* Create early pagetables. */
131    »        call   mk_early_pgtbl_32
132
```

```
290    /*
291     * Enable paging
292     */
293    »        movl $pa(initial_page_table), %eax
294    »        movl %eax,%cr3»  »        /* set the page table pointer.. */
295    »        movl $CR0_STATE,%eax
296    »        movl %eax,%cr0»  »        /* ..and set paging (PG) bit */
297    »        ljmp $__BOOT_CS,$1f»      /* Clear prefetch and normalize %eip */
298    1:
299    »        /* Shift the stack pointer to a virtual address */
300    »        addl $__PAGE_OFFSET, %esp
```

Before line 296, the system is running using physical addresses. Once x86 paging is activated, the system is running on virtual address.
The jump and esp modifications are there to make sure pointers shows virtual addresses.