

# Übung zu Betriebssystemtechnik

## Aufgabe 3: Paging in STUBSMI

---

21. Mai 2025

Dustin Nguyen, Maximilian Ott & Phillip Raffeck

Lehrstuhl für Informatik 4  
Friedrich-Alexander-Universität Erlangen-Nürnberg



Lehrstuhl für Informatik 4  
Systemsoftware



Friedrich-Alexander-Universität  
Technische Fakultät

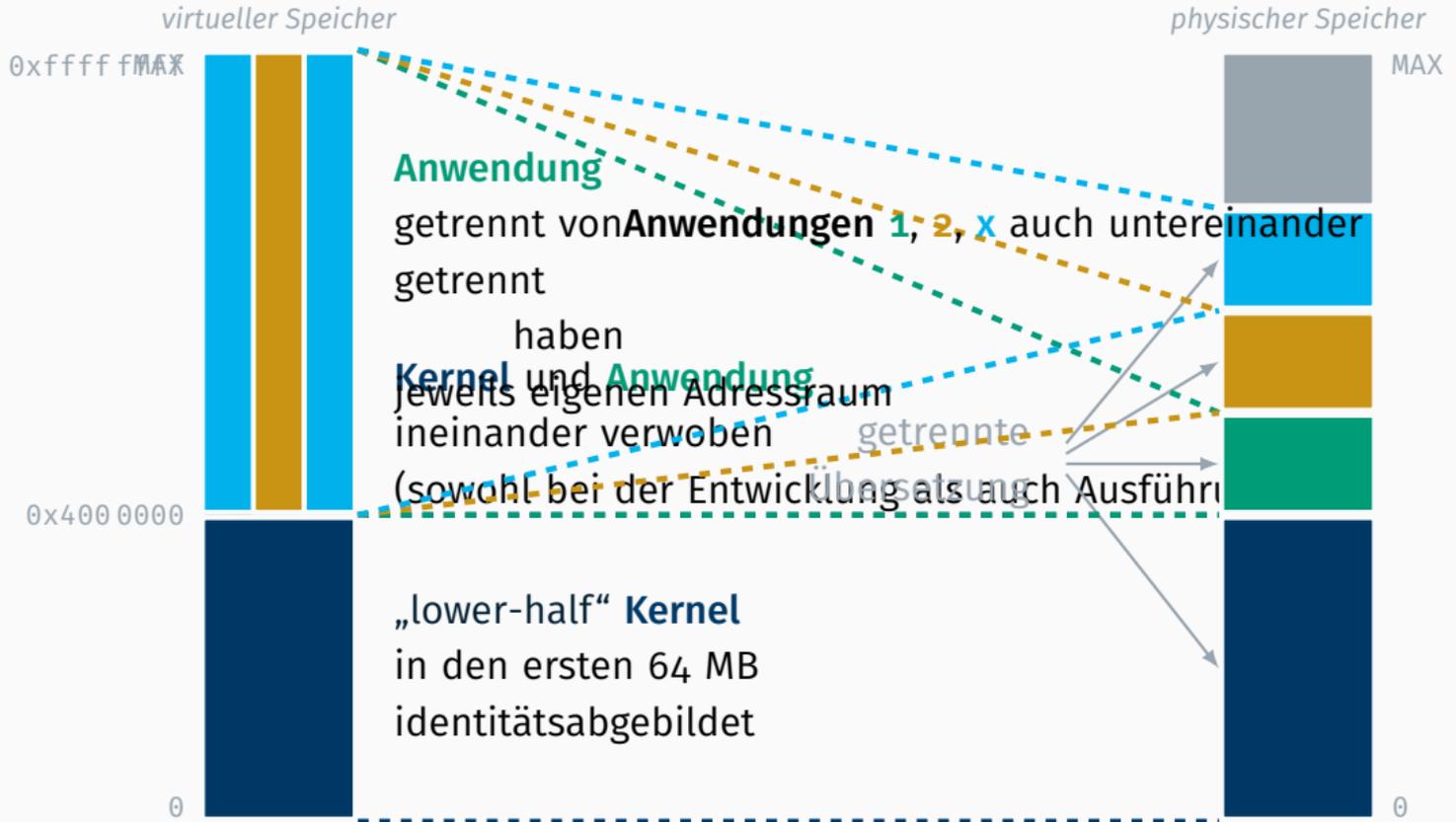
## Nachtrag: Systemaufrufbehandlung

```
extern "C" size_t syscall_handler(size_t p1, size_t p2, size_t p3,  
                                size_t p4, size_t p5, size_t sysnum) {  
    // Fast path for benchmarking  
    if (sysnum == SYSCALL_NOP)  
        return 0;  
  
    // Enter epilogue level  
    Guarded section;  
    // Enable Interrupts  
    Core::Interrupt::enable();  
  
    // Call syscall skeleton  
    switch(sysnum) {  
        case SYSCALL_WRITE:  
            return Skeleton::write(static_cast<int>(p1), reinterpret_cast<void*>(p2), p3);  
        // ...  
        default:  
            DBG << "Unknown SYSCALL " << sysnum << endl;  
            return static_cast<size_t>(-1);  
    }  
    return 0;  
}
```

**Motivation für die nächsten  
beiden Aufgaben (3 & 4)**

---

# Nach Aufgabe 4



***Exkurs:* MULTIBOOT SPECIFICATION**  
*oder* **Wie stiefel ich meinen Kernel?**

---

# Überblick zur MULTIBOOT SPECIFICATION

- offener PC **Bootloader Standard**, ab 1995 entwickelt
- Betriebssystem muss als **32 bit ELF** oder **a.out** vorliegen
- übernimmt die [hässliche] Initialisierung eines x86 PCs in einen **wohl definierten Zustand**
  - 32 bit Protected Mode
  - nur BSP (*Bootstrap Processor*)
  - A20 Gate aktiviert
  - setzt optional auch Grafikmodus
- übergibt dem BS „**vitale**“ **Informationen** über das System
- lädt ggf. auch *Boot Module* (weitere Dateien wie die **initiale Ramdisk**) in den Speicher
- wird u.a. von **GRUB** (Referenzimplementierung) und **PXELINUX** (Netzwerkboot) unterstützt
- und wird in **STUBS** verwendet

# Laden einer MULTIBOOT-kompatiblen Binärdatei

## Ablauf im Bootloader

1. liest System ELF Wert `0x1bad b002` (und `0x1000000`)
2. kopiert Code & Daten in den ersten `8192` Bytes (der ELF) liest die `ELF` Header und erstellt `BSS` Bereich in `boot/multiboot/header.asm`
3. lädt ggf. **Boot Module** (via Flags) und behält Konfiguration (via Flags)
4. setzt `eax` auf `0x2bad b002` sowie `ebx` als Zeiger auf Struktur mit

Hauptspeicher



5. Springt an den **Einsprungpunkt** (und übergibt somit an das Betriebssystem)

```
[SECTION .multiboot_header]
; Constants included from boot/multiboot/config.inc
MULTIBOOT_HEADER_MAGIC_OS equ 0x1badb002 ; Magic Header

MULTIBOOT_PAGE_ALIGN      equ 1<<0      ; Align boot modules at 4K
MULTIBOOT_MEMORY_INFO     equ 1<<1      ; Request Memory Map info
MULTIBOOT_VIDEO_MODE      equ 1<<2      ; Configure video mode

MULTIBOOT_HEADER_FLAGS equ MULTIBOOT_PAGE_ALIGN | MULTIBOOT_MEMORY_INFO

MULTIBOOT_HEADER_CHKSUM equ -(MULTIBOOT_HEADER_MAGIC_OS + MULTIBOOT_HEADER_FLAGS)

align 4 ; Align section
multiboot_header:
    dd MULTIBOOT_HEADER_MAGIC_OS ; Magic Header Value
    dd MULTIBOOT_HEADER_FLAGS    ; Flags (affects following entries)
    dd MULTIBOOT_HEADER_CHKSUM   ; Header Checksum
    ; additional fields depending on flags
    ; (e.g. specifying the desired video mode)
```



# Laden einer MULTIBOOT-kompatiblen Binärdatei

## Ablauf im Bootloader

1. liest System ELF Wert `0x1bad b002` (und `0x1000000`)
2. kopiert **Code & Daten** in den ersten `8192` Bytes (der ELF) liest die **Symboltabelle** und erstellt **BSS** `boot/multiboot/header.asm`
3. lädt ggf. **Boot Module** (via Flags) `boot/multiboot/header.asm` behält Konfiguration (via Flags)
4. setzt `eax` auf `0x2bad b002` sowie `ebx` als Zeiger auf Struktur mit

Hauptspeicher

0x8c 838  
0x85 340

**MULTIBOOT Information**

Symboltabelle

Symbole (aus STUBSML Linkerskript):

5. Springt an den **Einsprungpunkt** (und übergibt somit an das Betriebssystem) `boot/multiboot/header.asm` Auflösung Adresse zu Quelltext `boot/multiboot/header.asm`

0x10e04

0x1000

**Code, Header**

`initrd (1.1 MB)`

**MULTIBOOT Header**

MAX

(me)

ert

0x3b4 004c

0x3a4 0000 ←

0x101 9860

0x100 fe04

0x100 0000

0

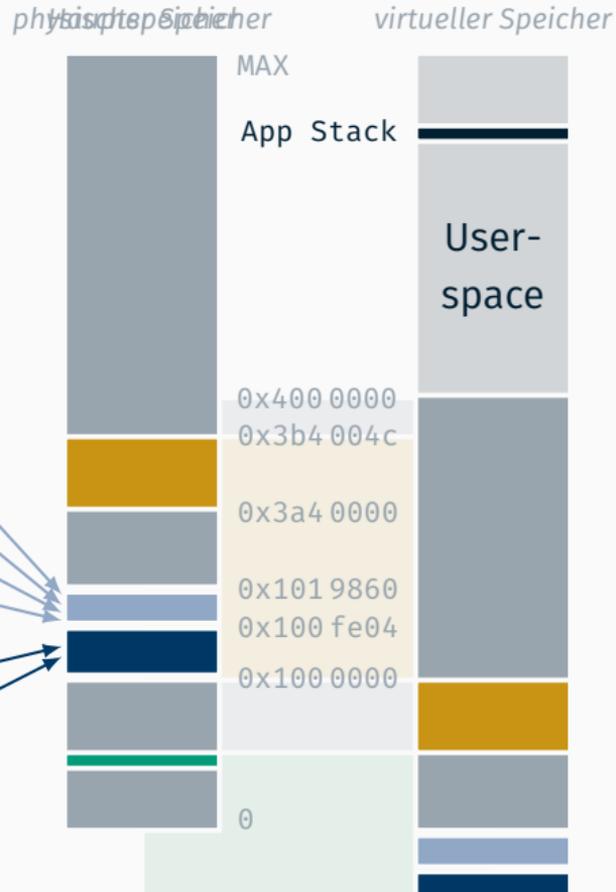
## Aufgabe 3

---

# Aufgabeninhalt

Ziel dieser Übung ist es, den zeitlichen (N/A/P/S) und statischen Kernel-Application-Kontext von stabiler and Kernel-Application-Kontext

- dynamisch Anwendungsstapel allokalieren
  - kein Hauptspeicherbelegung untersuchen
  - keine Isolation der Anwendungen
    - freien Speicher verwalten
- eigener Adressraum pro Anwendung
  - virtueller Adressraum mittels Paging
  - ersten 64 MB sind Kernelspace (lower-half) mit Identitätsabbildung
  - darüber liegt der Userspace mit Anwendungsstapel an einer fixen Startadresse (z.B. Top-of-Stack bei 0x8000 0000 0000)



# Speicherverwaltung

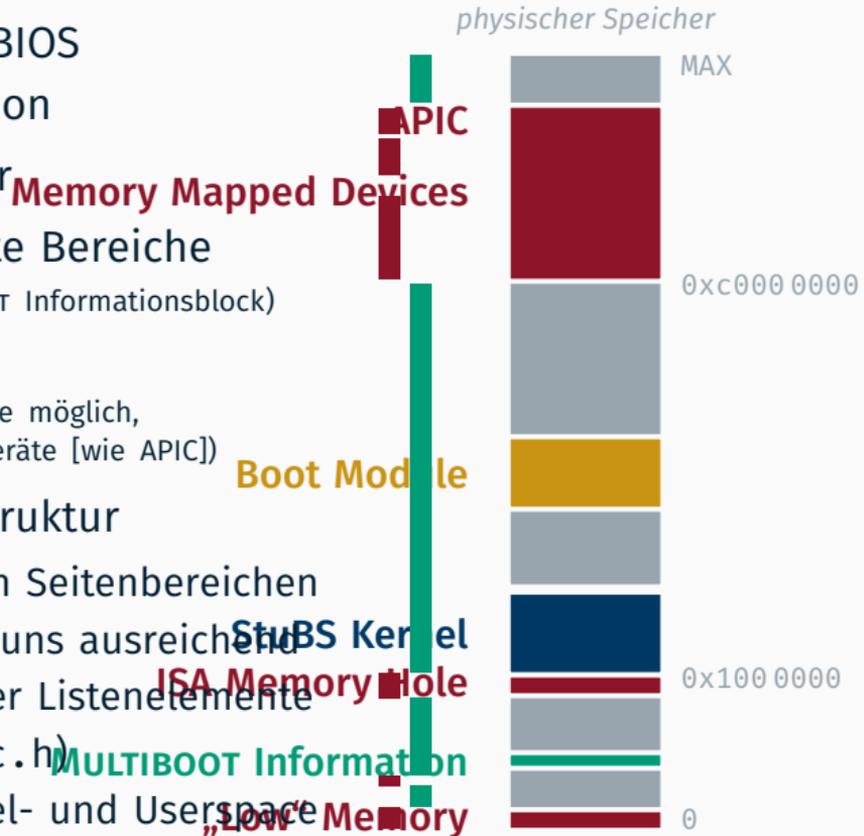
---

# Freien Speicher finden

Abfrage der **Memory Map** über BIOS

Ergebnis in MULTIBOOT Information

- freier und belegter Speicher
- ignoriert aber später belegte Bereiche  
(wie Kernel, initrd und den MULTIBOOT Informationsblock)
- besser defensiv auswerten  
(überlappende/widersprüchliche Bereiche möglich, ggf. fehlen im Speicher eingeblendete Geräte [wie APIC])
- Verwaltung in geeigneter Struktur
  - verkettete Liste mit freien Seitenbereichen (Startadresse, Länge) für uns ausreichend
  - dynamische Allokation der Listenelemente zulässig (→ `utils/alloc.h`)
  - Unterscheidung zw. Kernel- und Userspace





Der *Page Frame Allocator* sollte unbedingt vor dem nachfolgenden Schritt ausgiebig getestet werden!

```
void *addr;
void *prev = nullptr;
const size_t page_size = 4096;
while ((addr = reinterpret_cast<void*>(alloc_page())) != nullptr) {
    // longmode.asm maps only the first 4 GiB
    if (reinterpret_cast<uintptr_t>(addr) >= 0x100000000)
        continue;

    DBG << "Checking " << addr << endl;

    // Fill full page with 0b01011010 pattern
    memset(addr, 0x5a, page_size);
    // Check contents of page with previous filled one
    assert(prev == nullptr || memcmp(prev, addr, page_size) == 0);

    prev = addr;
}
```

# Paging

---

Mit 64 bit sind (theoretisch) bis zu 16 EiB adressierbar, aber der von uns verwendete x64 unterstützt „nur“

- maximal 52 bit (= 4 PiB) **physischen** Speicher
  - 40 bits durch Adressumsetzung + 12 bits Seite
  - MAXPHYADDR im Intel Manual (vgl. ISDMv3 4.1.4)
- standardmäßig 48 bit (= 256 TiB) **virtuellen** Speicher
  - über 4-stufige Adressumsetzung
  - die oberen 17 Bits einer Adresse müssen identisch sein (= *canonical*)  
→ valide Adressen sind 0x0 – 0x7fff ffff ffff sowie  
0xffff 8000 0000 0000 – 0xffff ffff ffff ffff
- neuere Architekturen 57 bit (= 128 PiB) **virt.** Speicher
  - über 5-stufige Adressumsetzung
  - muss extra aktiviert werden
  - die oberen 8 Bits müssen identisch sein (= *canonical*)

Für **STUBSMI** verwenden wir in dieser Aufgabe

- eine **Seitengröße von 4 KiB** wie im Beispiel  
(Hardware unterstützt 2 MiB sowie ggf. 1 GiB Seiten)
- eine **4-stufige Adressumsetzung** → 48 bit Adressen
- bei 7.5 ECTS auch die Möglichkeit Seiten als **nicht-ausführbar** zu markieren  
(dafür muss im *Extended Feature Enable Register* [MSR\_EFER] das 11. Bit [MSR\_EFER\_NXE] gesetzt sein)
- **keine weiteren Features** wie *Protection Keys* (ISDMv3 4.6.2)

# 4-stufige Adressumsetzung (48 bit) am Beispiel

ISDMv3 4.5

Virtuelle Adresse: 0x791bf4f2daffe

011110010001101111110100111100101101101011111110

Bits 39 ... 47

Bits 0 ... 11

Bits 30 ... 38

Bits 20 ... 29

Page Table

0xffff

target data

Page Directory

PDP Table

PML4

511	000000004b1d3000
...	...
243	000000003410001
242	000000003229007
241	0000000032c0800
...	...
0	0000000000000000

511	0000000000000000
...	...
112	00000000035b7001
111	00000000035bb007
110	0000000000000000
...	...
0	0000000003e11001

511	00000004b1d32000
...	...
424	000000000323c001
423	0000000003233007
422	000000000323a800
...	...
0	0000000000000000

511	0000000000000000
...	...
302	00000000ea5e5007
301	00001badbeefc003
300	00001badbeefb003
...	...
0	0000d0d0d0d0d003

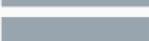
0x0

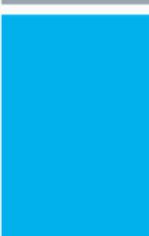
Page

Index ↑ Base Flags

%cr3

→ Physische Adresse: 0x1badbeefcafe

63		<b>Execute Disable:</b> 1 verhindert Ausführung von Code
62		<i>ignoriert</i>
52		<b>(Physische) Adresse der PDP Table,</b> welche an einer 4 KiB-Grenze ausgerichtet sein muss
51		
12		<i>ignoriert</i>
11		<i>ignoriert</i>
9		<b>Global:</b> Bei 4 KiB Seiten ignoriert
8		<i>reserviert, muss 0 sein</i>
7		<i>0</i>
6		<i>ignoriert</i>
5		<b>Accessed:</b> 1 falls die Zielseite verwendet wurde
4		<b>Page-Level Cache Disable:</b> 1 deaktiviert Caching
3		<b>Page-Level Write Through:</b> 1 aktiviert WT Caching
2		<b>User Mode:</b> 1 um Zugriff aus Ring 3 zu erlauben
1		<b>Writeable:</b> nur lesender (0) oder auch schreibender (1) Zugriff
0		<b>Present:</b> Eintrag aktiv (1) oder inaktiv (0)

63		<b>Execute Disable:</b> 1 verhindert Ausführung von Code
62		<i>ignoriert</i>
52		
51		<b>(Physische) Adresse des Seitenverzeichnisses (Page-Directory),</b> welche an einer 4 KiB-Grenze ausgerichtet sein muss
12		<i>ignoriert</i>
11		<i>ignoriert</i>
9		<b>Global:</b> Bei 4 KiB Seiten ignoriert
8		<b>Page Size:</b> Adresse zeigt auf Page-Directory (0) oder 1 GiB Seite (1)
7		
6		<i>ignoriert</i>
5		<b>Accessed:</b> 1 falls die Zielseite verwendet wurde
4		<b>Page-Level Cache Disable:</b> 1 deaktiviert Caching
3		<b>Page-Level Write Through:</b> 1 aktiviert WT Caching
2		<b>User Mode:</b> 1 um Zugriff aus Ring 3 zu erlauben
1		<b>Writeable:</b> nur lesender (0) oder auch schreibender (1) Zugriff
0		<b>Present:</b> Eintrag aktiv (1) oder inaktiv (0)

63		<b>Execute Disable:</b> 1 verhindert Ausführung von Code
62		<i>ignoriert</i>
52		
51		<b>(Physische) Adresse der Seitentabelle (Page Table),</b> welche an einer 4 KiB-Grenze ausgerichtet sein muss
12		<i>ignoriert</i>
11		<i>ignoriert</i>
9		<b>Global:</b> Bei 4 KiB Seiten ignoriert
8		<b>Page Size:</b> Adresse zeigt auf Page-Table (0) oder 2 MiB Seite (1)
7		
6		<i>ignoriert</i>
5		<b>Accessed:</b> 1 falls die Zielseite verwendet wurde
4		<b>Page-Level Cache Disable:</b> 1 deaktiviert Caching
3		<b>Page-Level Write Through:</b> 1 aktiviert WT Caching
2		<b>User Mode:</b> 1 um Zugriff aus Ring 3 zu erlauben
1		<b>Writeable:</b> nur lesender (0) oder auch schreibender (1) Zugriff
0		<b>Present:</b> Eintrag aktiv (1) oder inaktiv (0)

63	█	<b>Execute Disable:</b> 1 verhindert Ausführung von Code
62	█	<i>ignoriert</i>
52	█	<b>Physische Adresse</b> der 4 KiB Zielseite
51	█	
12	█	<i>ignoriert</i>
11	█	<i>ignoriert</i>
9	█	<i>ignoriert</i>
8	█ 0	<b>Global:</b> Verhindert TLB Aktualisierung
7	█ 0	<b>Page Attribute Table:</b> 1 aktiviert feingranulare Cacheeinstellung
6	█	<b>Dirty:</b> 1 falls auf die Zielseite geschrieben wurde
5	█	<b>Accessed:</b> 1 falls die Zielseite verwendet wurde
4	█ 0	<b>Page-Level Cache Disable:</b> 1 deaktiviert Caching
3	█ 0	<b>Page-Level Write Through:</b> 1 aktiviert WT Caching
2	█	<b>User Mode:</b> 1 um Zugriff aus Ring 3 zu erlauben
1	█	<b>Writeable:</b> nur lesender (0) oder auch schreibender (1) Zugriff
0	█ 1	<b>Present:</b> Eintrag aktiv (1) oder inaktiv (0)

63	0	<b>Reserviert</b>
32		
31	1	<b>Paging</b> aktiv (1) oder inaktiv (0)
30	0	Cache Disable: 1 deaktiviert Caching
29	0	Not Write Through: 1 deaktiviert WT Caching
28		
19		<b>Reserviert</b>
18	0	Alignment Mask: 1 aktiviert Prüfung der Ausrichtung
17		<b>Reserviert</b>
16	0 / 1	<b>Write Protect:</b> 0 erlaubt schreiben in ro-Seiten im Ring 0
15		<b>Reserviert</b>
6		
5	0	Numeric Error: 1 aktiviert FPU Ausnahmebehandlung
4	0	Extension Type: für Koprozessor (Modelabhängig)
3	0	Task Switched
2	0	Emulation
1	0	Monitor Coprocessor
0	1	<b>Protection Enable:</b> Real (0) oder Protected (1) Mode

} für FPU Kontextsicherung



**Reserviert.** Für was auch immer.

63

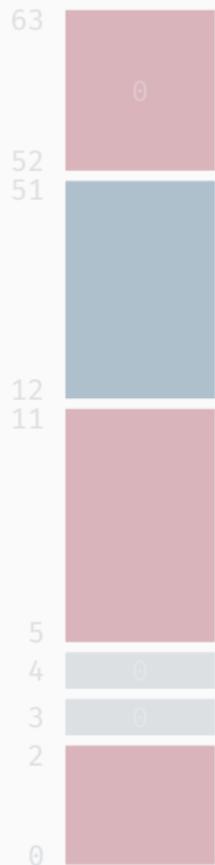


0

## Page-Fault Linear Address

Beinhaltet bei einem Seitenfehler die virtuelle Adresse, die den Fehler verursacht hat.

*(Noch) nicht notwendig in dieser Übung,  
aber kann das Entkäfern deutlich vereinfachen!*



Reserviert

(Physische) Adresse der PML4 (*Page-Map Level 4*) Tabelle, welche an einer 4 KiB-Grenze ausgerichtet sein muss



Reserviert

**Beim Schreiben von %cr3 wird TLB gespült**

Page-Level Cache Disable  
Page-Level Write Through

Reserviert

**%cr4** Steuerung von architekturabhängigen Erweiterungen wie **Page Size Extension** (4 MiB große Seiten) oder **Physical Address Extension** (erlaubt mehr als 4 GiB Speicher unter 32 Bit).

**%cr5** *reserviert*

**%cr6** *reserviert*

**%cr7** *reserviert*

**%cr8** steuert Zugriff auf *Task Priority Register*

**Aber:** Nicht wichtig für uns, wir ignorieren diese in der Übung.

# Implementierungshinweise

- Einträge in den Tabellen als Struktur/Klasse abbilden
  - Methoden zum Nachschlagen nützlich (nachgebildete MMU)
  - 4 KiB Ausrichtung der Tabellen nicht vergessen
- Codeduplikation ist eine hervorragende Quelle für Leichtsinnsfehler
  - ggf. sind hier C++ Templates hilfreich
  - virtuelle Methoden nur mit Bedacht einsetzen  
(vtable vergrößert Struktur → `static_assert` ist hilfreich)
- die ersten 64 MB (Kernel space) sollen identitätsabgebildet sein
  - **Ausnahme:** erste Seite im Speicher (Adresse `0x0`) nicht mappen
  - in dieser Aufgabe auch noch aus Userspace les- & schreibbar
  - für 7.5 ECTS: Die Seiten mit `Kernel .text` müssen ausführbar sein.
- an im Speicher eingeblendete Geräte denken
  - entweder anderweitige Verwendung & Zugriff im Userspace verhindern
  - oder im Kernel space einblenden

# Fragen?

---

Am 21.05. ist keine BST-Rechnerübung („bergfrei“)