

# Übung zu Betriebssystemtechnik

## Privilegienisolation und Systemaufrufe

---

2. Mai 2019

Andreas Ziegler  
Bernhard Heinloth

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



Lehrstuhl für Verteilte Systeme  
und Betriebssysteme



FRIEDRICH-ALEXANDER  
UNIVERSITÄT  
ERLANGEN-NÜRNBERG

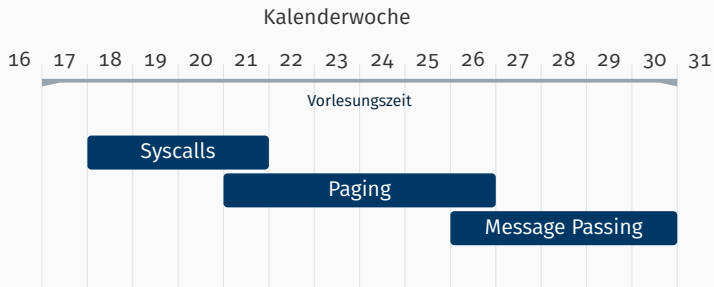
TECHNISCHE FAKULTÄT

# Organisation des Übungsbetriebs

---

**STUBSMI**  
single-core  
5 oder 7.5 ECTS Modul





## April 2019

1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30					

## Mai 2019

	1	<b>2</b>	3	4	5	
6	<b>7</b>	8	<b>9</b>	10	11	12
13	<b>14</b>	15	<b>16</b>	17	18	19
20	<b>21</b>	22	<b>23</b>	24	25	26
27	<b>28</b>	29	30	31		

**Tafelübung**  
im WinCIP (01.153)

**Rechnerübung**  
im WinCIP (01.153)

Abgabe der Aufgabe in  
der **Rechnerübung**

## Juni 2019

				1	2	
3	<b>4</b>	5	<b>6</b>	7	8	9
10	11	12	<b>13</b>	14	15	16
17	<b>18</b>	19	20	21	22	23
24	<b>25</b>	26	<b>27</b>	28	29	30

## Juli 2019

	1	<b>2</b>	3	<b>4</b>	5	6	7
8	<b>9</b>	10	11	12	13	14	
15	<b>16</b>	17	<b>18</b>	19	20	21	
22	<b>23</b>	24	25	26	27	28	
29	30	31					

- Abgabe nur in festen **2er Gruppen**
- eine (obligatorische) Tafelübung pro Aufgabe
- Anmeldung zur Tafelübung (bis 9. Mai) im Waffel auf **waffel.informatik.uni-erlangen.de/signup?course=375**
- Informationen und Aufgabenstellung auf **www4.cs.fau.de/Lehre/SS19/V\_BST/**
- Quelltextvorlage der Aufgaben im Gitlab  
`https://gitlab.cs.fau.de/i4/bs/stubsmi`

**Die Übung wird zwar durch Folien unterstützt –  
diese dienen jedoch ausschließlich als  
ergänzende Veranschaulichung des zu  
vermittelnden Stoffes und haben somit keinen  
Anspruch auf Vollständigkeit.  
Sie sind nicht zum Selbststudium geeignet und  
auch nicht auf Druck optimiert!**

## Selbsthilfe

- Aufpassen in der Tafelübung
- Internet (wiki.osdev.org, lowlevel.eu, Stack Overflow)
- Intel Handbuch

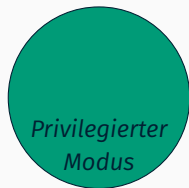
## Eskalationsstufen

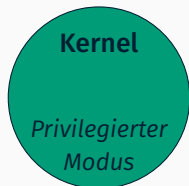
- Rechnerübung
- **#fau\_i4bs** im IRCnet
- XMPP-MUC **i4bs@conference.cs.fau.de**
- Mail an **bsstud@lists.informatik.uni-erlangen.de**
- **Raum 0.055** in der Martensstr. 1 (begründeter Notfall)

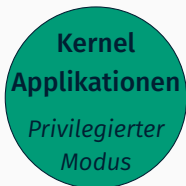


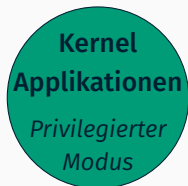
# Einleitung

---

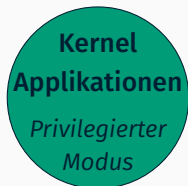






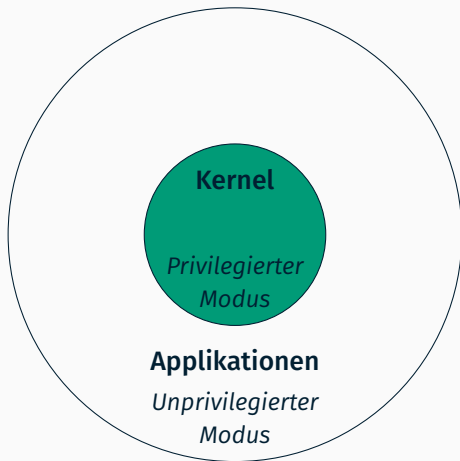


Probleme?



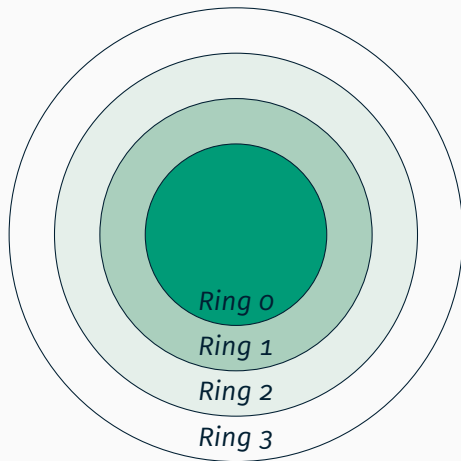
## Probleme?

- Zugriff auf privilegierte Befehle (`hlt`, `cli/sti`, `inb/outb`, ...)
- Zugriff auf CPU-Kontrollregister (z.B. MMU-Konfiguration)



## Probleme?

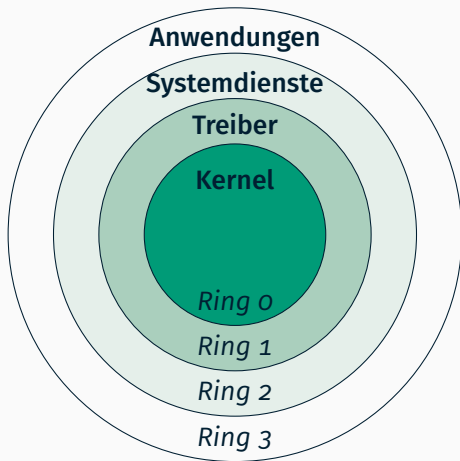
- Zugriff auf privilegierte Befehle (`hlt`, `cli/sti`, `inb/outb`, ...)
- Zugriff auf CPU-Kontrollregister (z.B. MMU-Konfiguration)



- Schutzringkonzept ursprünglich aus *Multics*
- Auf **x86** insgesamt 4 Ringe

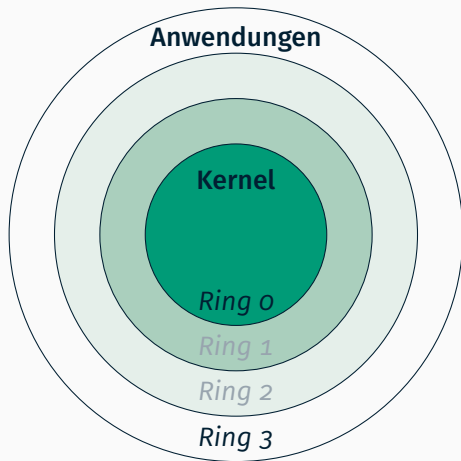


# Privilegientrennung auf x86: Schutzringe



- Schutzringkonzept ursprünglich aus *Multics*
- Auf **x86** insgesamt 4 Ringe

# Privilegientrennung auf x86: Schutzringe



- Schutzringkonzept ursprünglich aus *Multics*
- Auf **x86** insgesamt 4 Ringe
- Tatsächlich verwendet: Ring 0 (System) und 3 (Benutzer)

# Umsetzung

---

## Schutzringe

- Konfiguration des Prozessors
- Wechsel zwischen den Modi

## Schutzringe

- Konfiguration des Prozessors
- Wechsel zwischen den Modi

## Systemaufrufe

- Erlauben von Traps aus Ring 3
- Übergabe von Argumenten

# Interrupt Deskriptor

63		<b>Offset (high):</b> oberer Teil der Einsprungsadresse für die Interruptbehandlung (z.B. <code>irq_entry_128</code> )
48		
47		<b>Present:</b> Eintrag aktiv (1) oder inaktiv (0)
46		
45		<b>Descriptor Privilege Level:</b> Erlaubter Ring
44		<b>Storage Segment:</b> 0 für Interrupt und Traps
43		<b>Mode:</b> 16-bit (0) oder 32-bit (1)
42		
40		<b>Type:</b> Task (5), Interrupt (6) oder Trap (7)?
39		
32		<b>Unused</b> – muss 0 sein
31		
16		<b>Selector:</b> Codesegment, in das beim Interrupt gewechselt wird (i.d.R. Kernel-Codesegment)
15		
0		<b>Offset (low):</b> unterer Teil der Einsprungsadresse für die Interruptbehandlung

# Interrupt Deskriptor

63		<b>Offset (high):</b> oberer Teil der Einsprungsadresse für die Interruptbehandlung (z.B. <code>irq_entry_128</code> )
48		
47		<b>Present:</b> Eintrag aktiv (1) oder inaktiv (0)
46		
45		<b>Descriptor Privilege Level:</b> Erlaubter Ring
44		<b>Storage Segment:</b> 0 für Interrupt und Traps
43		<b>Mode:</b> 16-bit (0) oder 32-bit (1)
42		
40		<b>Type:</b> Task (5), Interrupt (6) oder Trap (7)?
39		
32		<b>Unused</b> – muss 0 sein
31		
16		<b>Selector:</b> Codesegment, in das beim Interrupt gewechselt wird (i.d.R. Kernel-Codesegment)
15		
0		<b>Offset (low):</b> unterer Teil der Einsprungsadresse für die Interruptbehandlung

# Interrupt Deskriptor (Beispiel)

63		<b>Offset (high):</b> oberer Teil der Einsprungsadresse für die Interruptbehandlung (z.B. <code>irq_entry_128</code> )
48		
47		<b>Present:</b> Eintrag aktiv (1) oder inaktiv (0)
46		
45		<b>Descriptor Privilege Level:</b> Erlaubter Ring
44		<b>Storage Segment:</b> 0 für Interrupt und Traps
43		<b>Mode:</b> 16-bit (0) oder 32-bit (1)
42		<b>Type:</b> Task (5), Interrupt (6) oder Trap (7)?
40		
39		
32		<b>Unused</b> – muss 0 sein
31		
		<b>Selector:</b> Codesegment, in das beim Interrupt gewechselt wird (i.d.R. Kernel-Codesegment)
16		
15		
		<b>Offset (low):</b> unterer Teil der Einsprungsadresse für die Interruptbehandlung
0		

für `100 0c68 <irq_entry_128>`



# Interrupt Deskriptor (Beispiel)

63	0x0100	<b>Offset (high):</b> oberer Teil der Einsprungsadresse für die Interruptbehandlung (z.B. <code>irq_entry_128</code> )
48		
47	1	<b>Present:</b> Eintrag aktiv (1) oder inaktiv (0)
46	3	<b>Descriptor Privilege Level:</b> Erlaubter Ring
45		
44	0	<b>Storage Segment:</b> 0 für Interrupt und Traps
43	1	<b>Mode:</b> 16-bit (0) oder 32-bit (1)
42		
40	6	<b>Type:</b> Task (5), Interrupt (6) oder Trap (7)?
39		
32	0	<b>Unused</b> – muss 0 sein
31		
16	8	<b>Selector:</b> Codesegment, in das beim Interrupt gewechselt wird (i.d.R. Kernel-Codesegment)
15		
0	0x0c68	<b>Offset (low):</b> unterer Teil der Einsprungsadresse für die Interruptbehandlung

für `100 0c68 <irq_entry_128>` → `0x0100 ee00 0008 0c68`

# Segmente

- Segmentierung tief in X86 verwurzelt
- Segmentregister (ss, cs, ds, ...) beinhalten Selektoren
- Selektoren (16-bit Werte) werden direkt in Register geladen:

```
mov ax, 0x10  
mov ds, ax
```

- Segmentierter Zugriff:

```
mov [fs:eax], 42
```

- Selektoren bestehen aus Segmentdeskriptorindex in GDT und Requested Privilege Level (RPL)



# Global Deskriptor Table Entry

63		<b>Base Address (high):</b> Oberer Teil der Startadresse
56		
55		<b>Granularity:</b> Byte (0) oder 4KiB-Block (1)
54		<b>Operation Size:</b> 16-bit (0) oder 32-bit (1)
53		<b>64-bit Segment:</b> Bei uns 0 für Protected Mode
52		<b>Available</b> – zur freien Verwendung
51		
48		<b>Segment Limit (high):</b> Oberer Teil der Segmentgröße
47		<b>Present:</b> Eintrag aktiv (1) oder inaktiv (0)
46		<b>Descriptor Privilege Level:</b> Erlaubter Ring
45		<b>Deskriptortyp:</b> System (0) oder Code/Daten (1)
44		
43		<b>Type:</b> Legt Berechtigungen fest (R/W/X)
40		
39		
		<b>Base Address (low):</b> Unterer Teil der Startadresse
16		
15		
		<b>Segment Limit (low):</b> Unterer Teil der Segmentgröße
0		



## Auslösen eines Systemaufrufs

Was passiert bei `int 0x80` aus Ring 3?

Hardware:

# Auslösen eines Systemaufrufs

Was passiert bei `int 0x80` aus Ring 3?

**Hardware:**

- Instruction Pointer setzen
- Ring wechseln

# Auslösen eines Systemaufrufs

Was passiert bei `int 0x80` aus Ring 3?

Hardware:

- Instruction Pointer setzen (aus **IDT**)
- Ring wechseln (aus **IDT**)

Was passiert bei `int 0x80` aus Ring 3?

Hardware:

- Instruction Pointer setzen (aus **IDT**)
- Ring wechseln (aus **IDT**)
  
- Zustand auf Stack sichern



Was passiert bei `int 0x80` aus Ring 3?

Hardware:

- Instruction Pointer setzen (aus IDT)
- Ring wechseln (aus IDT)
  
- Zustand auf Stack sichern

Woher kommt der Stack-Pointer (und -Segment)?

# Auslösen eines Systemaufrufs

Was passiert bei `int 0x80` aus Ring 3?

Hardware:

- Instruction Pointer setzen (aus IDT)
- Ring wechseln (aus IDT)
  
- Zustand auf Stack sichern

Woher kommt der Stack-Pointer (und -Segment)?

```
mov esp, 0x0  
int 0x80
```

# Auslösen eines Systemaufrufs

Was passiert bei `int 0x80` aus Ring 3?

Hardware:

- Instruction Pointer setzen (aus IDT)
- Ring wechseln (aus IDT)
  
- Zustand auf Stack sichern

Woher kommt der Stack-Pointer (und -Segment)?

```
mov esp, 0x0
```

```
int 0x80 
```

# Auslösen eines Systemaufrufs

Was passiert bei `int 0x80` aus Ring 3?

Hardware:

- Instruction Pointer setzen (aus IDT)
- Ring wechseln (aus IDT)
- **Auf Kernel-Stack wechseln**
- Zustand auf Stack sichern

Woher kommt der Stack-Pointer (und -Segment)?

```
mov esp, 0x0  
int 0x80
```



→ Task State Segment (TSS)

# Task State Segment

	CR4	40
	EFLAGS	36
	EIP	32
	CR3 (PDBR)	28
	SS2	24
	ESP2	20
	SS1	16
	ESP1	12
	SS0	8
	ESP0	4
	Previous Task Link	0

Intel unterstützt Hardware-Tasks  $\Rightarrow$  Zustandssicherung

Ein Deskriptor in der **GDT**, Base = Adresse des **TSS**

Wird mittels `mov ax, <offset>; ltr ax` geladen

Bei uns: **ein** globales **TSS**, jeder Task hat eigenen Kernel Stack

## **Syscall-Nummer und Parameter in Registern übergeben**

Reihenfolge (z. B.): `eax, ecx, edx, ebx, esi, edi`

```
int sys_foo(int p1, int p2, int p3);
```

## Syscall-Nummer und Parameter in Registern übergeben

Reihenfolge (z. B.): eax, ecx, edx, ebx, esi, edi

```
int sys_foo(int p1, int p2, int p3);
```

**Entweder:** alles in purem Assembler schreiben

```
[GLOBAL sys_foo]
sys_foo:
    push ebx
    mov eax, FOO_SYSNUM
    mov ecx, [esp + 8]
    mov edx, [esp + 12]
    mov ebx, [esp + 16]
    int 0x42
    pop ebx
    ret
```

### Syscall-Nummer und Parameter in Registern übergeben

Reihenfolge (z. B.): `eax`, `ecx`, `edx`, `ebx`, `esi`, `edi`

```
int sys_foo(int p1, int p2, int p3);
```

**Oder:** Systemaufruf (fast) in C mit *Inline Assembly*

```
int sys_foo(int p1, int p2, int p3) {
    int retval;
    asm volatile (
        "int $0x42\n\t"      // Befehle
        : "=a"(retval)      // Outputs
        : "a"(FOO_SYSNUM), "c"(p1), "d"(p2),
          "b"(p3)           // Inputs
        : "memory");       // Clobber
    return retval;
}
```



Was fehlt noch?

**Was fehlt noch?**

→ initialer Wechsel in Ring 3

# Systemaufruf auf Kernel-Seite

## Was fehlt noch?

→ initialer Wechsel in Ring 3

Ähnlich zu `toc_settle()`: Kernel-Stack passend "faken"

SS
ESP
EFLAGS
CS
EIP

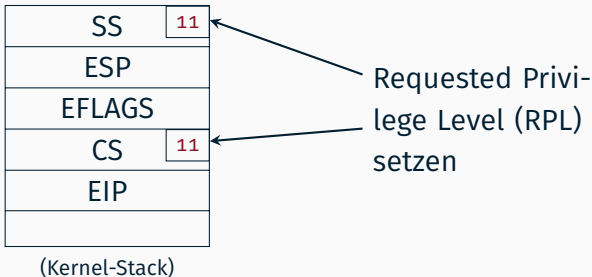
(Kernel-Stack)

# Systemaufruf auf Kernel-Seite

## Was fehlt noch?

→ initialer Wechsel in Ring 3

Ähnlich zu `toc_settle()`: Kernel-Stack passend “faken”



**Wichtig:** Segmentregister (ds, es, fs, gs) setzen

Danach: `iret`

## Was ist (noch) nötig?

- Neue Segment-Deskriptoren anlegen (Usermode & TSS)
- Adresse des TSS → Deskriptor & `ltr <offset>`
- User-Stacks für Anwendungen einführen  
(→ `toc_settle` weiterhin auf Kernel-Stack)

## Was ist (noch) nötig?

- Neue Segment-Deskriptoren anlegen (Usermode & TSS)
- Adresse des TSS → Deskriptor & ltr <offset>
- User-Stacks für Anwendungen einführen  
(→ `toc_settle` weiterhin auf Kernel-Stack)
  
- Dispatcher: `tss.esp0` (um)setzen (`go & dispatch`)
- `kickoff` macht jetzt Wechsel in Ring 3  
→ `kickoff_user(Thread *obj)` mit `obj->action()`  
**Woher kommt Parameter?**

## Was ist (noch) nötig?

- Neue Segment-Deskriptoren anlegen (Usermode & TSS)
- Adresse des TSS → Deskriptor & ltr <offset>
- User-Stacks für Anwendungen einführen  
(→ `toc_settle` weiterhin auf Kernel-Stack)
  
- Dispatcher: `tss.esp0` (um)setzen (`go & dispatch`)
- `kickoff` macht jetzt Wechsel in Ring 3  
→ `kickoff_user(Thread *obj)` mit `obj->action()`  
**Woher kommt Parameter?** → User-Stack präparieren

## Was ist (noch) nötig?

- Neue Segment-Deskriptoren anlegen (Usermode & TSS)
- Adresse des TSS → Deskriptor & ltr <offset>
- User-Stacks für Anwendungen einführen  
(→ `toc_settle` weiterhin auf Kernel-Stack)
  
- Dispatcher: `tss.esp0` (um)setzen (`go & dispatch`)
- `kickoff` macht jetzt Wechsel in Ring 3  
→ `kickoff_user(Thread *obj)` mit `obj->action()`  
**Woher kommt Parameter?** → User-Stack präparieren
  
- **IDT** & **IRQ-Makro** in `startup.asm` anpassen
- Systemaufruf-Stümpfe und -Zuteiler schreiben



## **Schnellere Systemaufrufe (7.5 ECTS)**

---

Traps sind langsam (Zustandssicherung, Segmentierung, ...)  
Eigentlich für uns nur Ringwechsel nötig

Traps sind langsam (Zustandssicherung, Segmentierung, ...)

Eigentlich für uns nur Ringwechsel nötig

→ `sysenter` & `sysexit`

Traps sind langsam (Zustandssicherung, Segmentierung, ...)

Eigentlich für uns nur Ringwechsel nötig

→ `sysenter` & `sysexit`

`sysenter` sichert keinen Zustand (somit auch kein `eip/esp!`)

Traps sind langsam (Zustandssicherung, Segmentierung, ...)

Eigentlich für uns nur Ringwechsel nötig

→ `sysenter` & `sysexit`

`sysenter` sichert keinen Zustand (somit auch kein `eip/esp`!)

`sysexit` erwartet aber `esp` in `ecx` und `eip` in `edx`

Traps sind langsam (Zustandssicherung, Segmentierung, ...)

Eigentlich für uns nur Ringwechsel nötig

→ `sysenter` & `sysexit`

`sysenter` sichert keinen Zustand (somit auch kein `eip/esp`!)

`sysexit` erwartet aber `esp` in `ecx` und `eip` in `edx`

⇒ müssen vom Aufruf-Stumpf mitgegeben werden!

## sysenter: Aufruf-Stumpf und Kern-Wrapper

```
sysenter_foo:  
    mov edx, post_label  
    mov ecx, esp  
    mov eax, FOO_SYSNUM  
    sysenter  
post_label:  
    ret
```

## sysenter: Aufruf-Stumpf und Kern-Wrapper

```
%macro syscall 2
[GLOBAL %1]
%1:
    mov edx, %%post_label
    mov ecx, esp
    mov eax, %2
    sysenter
%%post_label:
    ret
%endmacro
syscall sysenter_foo, FOO_SYSNUM
```



## sysenter: Aufruf-Stumpf und Kern-Wrapper

```
%macro syscall 2
[GLOBAL %1]
%1:
    mov edx, %%post_label
    mov ecx, esp
    mov eax, %2
    sysenter
%%post_label:
    ret
%endmacro
syscall sysenter_foo, FOO_SYSNUM
```

**Kernel-Seite:** eigener Wrapper, Parameter holen, C++, `sysexit`

Woher kommen Kern-cs/eip und -ss/esp?

Woher kommen Kern-cs/eip und -ss/esp?

→ Model-specific registers (MSRs)

`cs 0x174`

`esp 0x175`

`eip 0x176`

Woher kommen Kern-cs/eip und -ss/esp?

→ Model-specific registers (MSRs)

cs 0x174

esp 0x175

eip 0x176

(ss [Wert in 0x174] + 8)

Woher kommen Kern-cs/eip und -ss/esp?

→ Model-specific registers (MSRs)

`cs 0x174`

`esp 0x175`

`eip 0x176`

`(ss [Wert in 0x174] + 8)`

Somit esp bei Threadwechsel schreiben, Rest nur initial

**Woher kommen Kern-cs/eip und -ss/esp?**

→ Model-specific registers (MSRs)

`cs 0x174`

`esp 0x175`

`eip 0x176`

`(ss [Wert in 0x174] + 8)`

Somit esp bei Threadwechsel schreiben, Rest nur initial

**Woher kommen User-cs/ss bei sysexit?**

Woher kommen Kern-cs/eip und -ss/esp?

→ Model-specific registers (MSRs)

`cs 0x174`

`esp 0x175`

`eip 0x176`

`(ss [Wert in 0x174] + 8)`

Somit esp bei Threadwechsel schreiben, Rest nur initial

Woher kommen User-cs/ss bei `sysexit`?

**Code-Segment** Wert in MSR `0x174 + 16`

**Stack-/Daten-Segment** Wert in MSR `0x174 + 24`

⇒ richtiges **GDT**-Layout wichtig!





Auslesen mittels rdtsc-Instruktion

```
static inline uint64_t rdtsc() {
    uint32_t low, high;
    asm volatile (
        "rdtsc\n\t"
        : "=a" (low), "=d" (high)
        :: );
    return ((uint64_t)(high) << 32) | low;
}
```

oder einfach CLANG/GCC-Builtin:

```
uint64_t __builtin_ia32_rdtsc()
```

## Zu beachten:

- eigentlich™ unsauber
  - wegen *out-of-order* Ausführung
  - Lösung mittels serialisierender Instruktion (`cpuid`) oder `rdtscp`
- Whitepaper “How to Benchmark Code Execution Times on Intel® IA-32 and IA-64 Instruction Set Architectures”

⇒ für uns Builtin ausreichend

**Fragen?**