

# Übung zu Betriebssystemtechnik

## Privilegienisolation und Systemaufrufe

---

23. April 2018

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

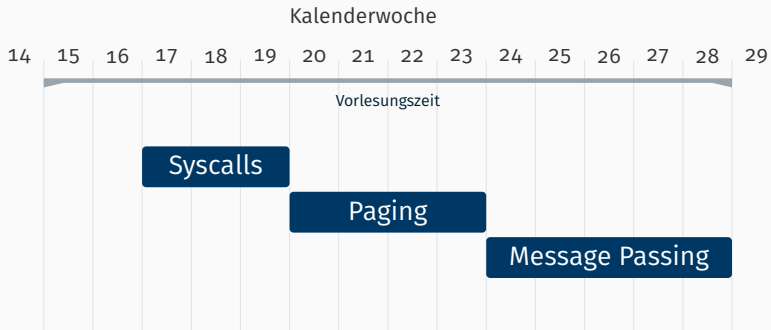
# Organisation des Übungsbetriebs

---

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



# Zeitplan





- Abgabe nur in festen **2er Gruppen**
- eine (obligatorische) Tafelübung pro Aufgabe
- Anmeldung zur Tafelübung (bis 24. April) im Waffel auf **waffel.informatik.uni-erlangen.de/signup?course=338**
- Informationen und Aufgabenstellung auf **www4.cs.fau.de/Lehre/SS18/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
- **#faii4bs** 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

---

# Ziel der Übung

**Auftrennung Usermodus  $\Leftrightarrow$  Systemmodus**

Unterschied?

## Auftrennung Usermodus $\Leftrightarrow$ Systemmodus

Unterschied?

- Erlaubte Befehle (hlt, cli/sti, iret, inb/outb, ...)
- Zugriff auf Kontrollregister: cr{0,2,3,4}  
⇒ CPU-/MMU-Konfiguration

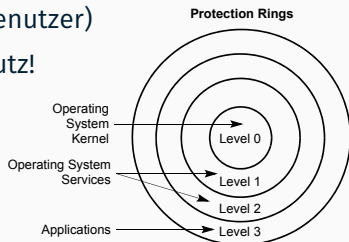
## Auftrennung Usermodus $\Leftrightarrow$ Systemmodus

Unterschied?

- Erlaubte Befehle (hlt, cli/sti, iret, inb/outb, ...)
- Zugriff auf Kontrollregister: cr{0,2,3,4}  
⇒ CPU-/MMU-Konfiguration

## x86: Ringmodell

- Insgesamt 4 Ringe: 0 (System) – 3 (Benutzer)
- Beinhaltet noch keinen Speicherschutz!
- Auf Ring 3: Privilegierte Befehle und Register geschützt



# Umsetzung

---

## Schutzringe

- Konfiguration des Prozessors
- Wechsel zwischen den Modi

## Schutzringe

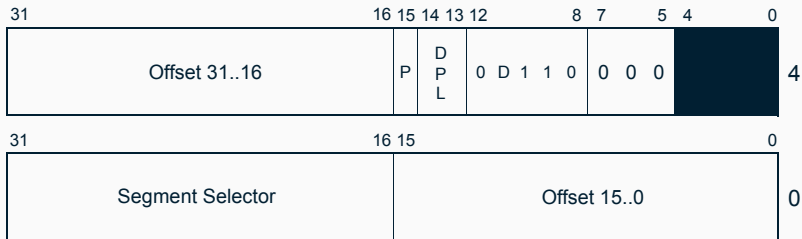
- Konfiguration des Prozessors
- Wechsel zwischen den Modi

## Systemaufrufe

- Erlauben von Traps aus Ring 3
- Übergabe von Argumenten

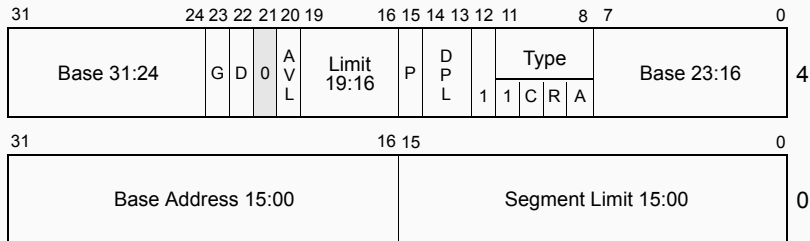
# Interrupt Descriptor Table (IDT)

## Deskriptoren für Interrupts



- Descriptor Privilege Level (DPL): Erlaubter Ring
- Offset: Angesprungene Funktion

# Segment-Selektoren? ⇒ GDT



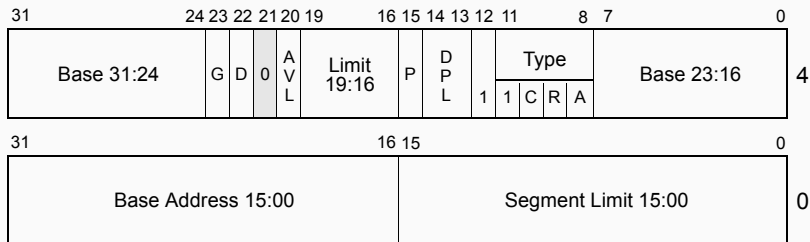
## Segmentierung tief in x86 verwurzelt

**Segmentregister** (SS, CS, DS, ...) beinhalten Offset in GDT

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

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

# Segment-Selektoren? ⇒ GDT



## Segmentierung tief in x86 verwurzelt

**Segmentregister** (SS, CS, DS, ...) beinhalten Offset in GDT

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

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

**Was bedeutet** `FF FF 00 00 | 00 9A CF 00`

## Auslösen eines Systemaufrufs

Was passiert bei `int 0x42` aus Ring 3?

Hardware

Was passiert bei `int 0x42` aus Ring 3?

## Hardware

- Instruction Pointer setzen
- Ring wechseln

Was passiert bei `int 0x42` aus Ring 3?

## Hardware

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

Was passiert bei `int 0x42` aus Ring 3?

## Hardware

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

Was passiert bei `int 0x42` aus Ring 3?

## Hardware

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

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

Was passiert bei `int 0x42` aus Ring 3?

## Hardware

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

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

```
mov esp, 0x0; int 0x42;
```

# Auslösen eines Systemaufrufs

Was passiert bei `int 0x42` aus Ring 3?

## Hardware

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

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

```
mov esp, 0x0; int 0x42; 
```

# Auslösen eines Systemaufrufs

Was passiert bei `int 0x42` 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 (+ -Segment)?

```
mov esp, 0x0; int 0x42; 
```

⇒ **Task State Segment**

# 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

## Parameterübergabe

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);
```

## Parameterübergabe

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 Assembly 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
```

Alternative: 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;
}
```

**Erinnerung:** Hardware wechselt Stack und Ring, setzt `eip`

**Erinnerung:** Hardware wechselt Stack und Ring, setzt `eip`

- Alle Argumentregister auf den (Kernel!) Stack legen
- System-Call-Handler (in C++) aufrufen
- `iret`

⇒ System-Call-Handler (in C++) ist `if-else`-Kaskade bzw. großes `switch/case`

Was fehlt noch?

**Was fehlt noch?**

⇒ initialer Wechsel in Ring 3

## Systemaufruf auf Kernel-Seite (cont.)

### 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 (cont.)

### Was fehlt noch?

⇒ initialer Wechsel in Ring 3

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

SS	11
ESP	
EFLAGS	
CS	11
EIP	

(Kernel-Stack)

**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

## Alternative zu Traps: `sysenter`

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

Eigentlich für uns nur Ringwechsel nötig

⇒ `sysenter` & `sysexit`

## Alternative zu Traps: `sysenter`

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

Eigentlich für uns nur Ringwechsel nötig

⇒ `sysenter` & `sysexit`

`sysenter` sichert keinen Zustand ⇒ auch kein `eip/esp`!

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

## Alternative zu Traps: `sysenter`

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

Eigentlich für uns nur Ringwechsel nötig

⇒ `sysenter` & `sysexit`

`sysenter` sichert keinen Zustand ⇒ auch kein `eip/esp`!

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

⇒ 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

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

%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

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

%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

Evtl. Hardware-Kontext nachbauen (→ später fork( ))

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)

⇒ esp bei Task-Wechsel 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)
```

⇒ esp bei Task-Wechsel 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)`

⇒ esp bei Task-Wechsel 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!

Intel bietet Zählerregister an, in jedem Takt inkrementiert

Auslesen mittels `rdtsc`-Instruktion

clang/gcc-Builtin: `uint64_t __builtin_ia32_rdtsc()`

Zu beachten:

- eigentlich™ unsauber

⇒ Out-of-order Ausführung

⇒ Serialisierende Instruktion (`cpuid`) und/oder `rdtsc`

- Whitepaper “How to Benchmark Code Execution Times on Intel® IA-32 and IA-64 Instruction Set Architectures”

⇒ für uns Builtin ausreichend (oder `cpuid` nutzen)

# **Erinnerung: Entwicklungsumgebung**

---

## Makefile Targets

**make qemu** QEMU **ohne** Hardware-Virtualisierung

**make kvm** QEMU **mit** Hardware-Virtualisierung

**make qemu-gdb** starte in QEMU und verbinde zu integrierten GDB-Stub → Fehlersuche mit gdb.

**make netboot** für Boot am Test-Rechner ins NFS kopieren

Suffix `-noopt` um Optimierungen auszuschalten (sonst `-O3`)

**Fragen?**