



# Isolation und Schutz in Betriebssystemen

## 3. Systemaufrufe bei x86-64

Michael Schöttner

- Bisher können unsere User-Mode Threads (Ring 3) einfach alle Funktionen des Kernels direkt aufrufen.
  - Jedoch wirft der Prozessor eine General Protection Fault (GPF), falls in den aufgerufenen Kern-Funktionen eine privilegierte Instruktion verwendet wird
  - Normalerweise werden die Funktionen des Kerns über das Paging vor direkten Aufrufen aus dem Ring 3 geschützt.
- Wir wollen nun Systemaufrufe realisieren, wodurch nur noch bestimmte Kernel-Funktionen indirekt in kontrollierter Weise aufgerufen werden

## 3.2 Mechanismen für Systemaufrufe bei x86

### ■ Call Gate

- Deskriptor in der GDT
- Zugriff mit einer `call` oder `jmp` Instruktion
- Stack wird mithilfe des TSS umgeschaltet

### ■ Interrupt Trap Gate

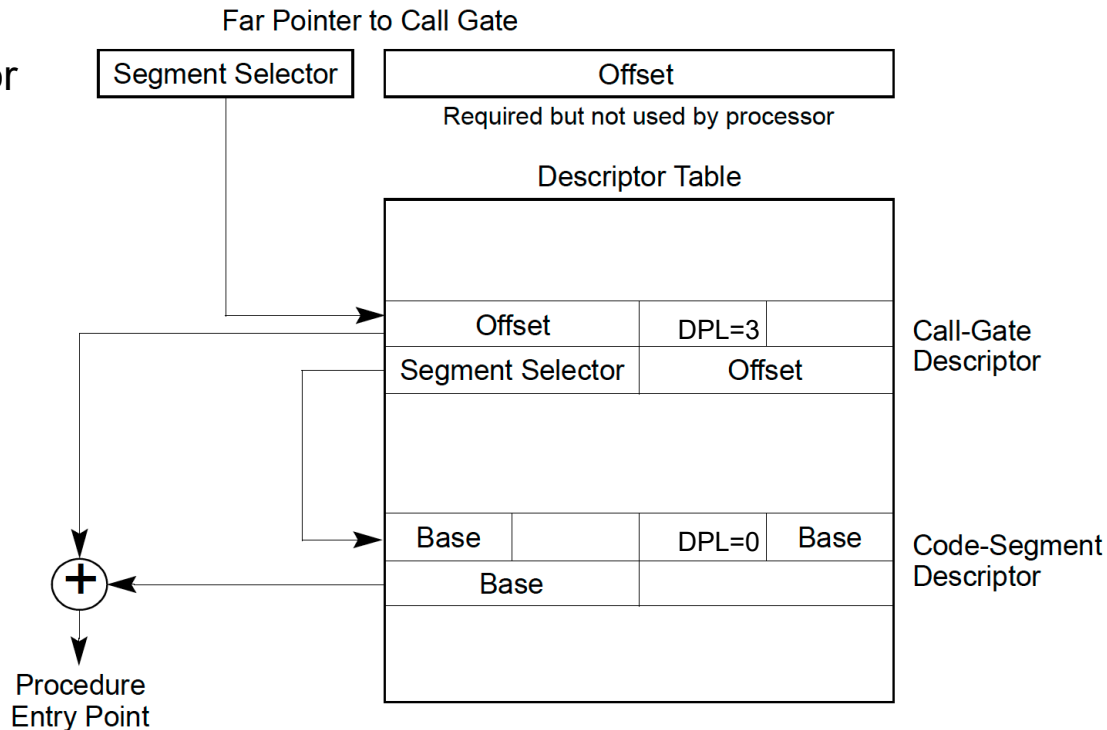
- Deskriptor in der IDT
- Zugriff mit einer `int` Instruktion
- Stack wird mithilfe des TSS umgeschaltet

### ■ Schnelle Systemaufrufe

- Mithilfe MSR (Model Specific Register)
- Zugriff mit einer `syscall/sysret`
- Stack muss in Software umgeschaltet werden

## 3.3 Call Gate

- Offset im Register sowie im Code Segment Deskriptor werden ignoriert



- Interrupt/Trap Gate**

The diagram illustrates the structure of the Interrupt/Trap Gate, which is 32 bits wide. It is divided into several fields:

  - Reserved:** 12 bits (bits 31-20), shaded gray.
  - Offset 63..32:** 8 bits (bits 19-12).
  - Control and Type Fields:** 4 bits (bits 11-8), divided into:
    - P:** Privilege level (1 bit, bit 11).
    - DPL:** Descriptor Privilege Level (2 bits, bits 10-9).
    - 0:** Reserved zero (1 bit, bit 8).
    - TYPE:** Interrupt type (2 bits, bits 7-6).
  - IST:** Interrupt Stack Table index (2 bits, bits 5-4).
  - Segment Selector:** 16 bits (bits 15-0), the left half of the bottom section.
  - Offset 15..0:** 16 bits (bits 15-0), the right half of the bottom section.

### Interrupt Stack Table

## 3.4 Interrupt/Trap Gate

### ■ Interrupt Gate

- Interrupt Enable Flag wird gelöscht
- Verwendet für Interrupt Handler → sequentielle Abarbeitung (auf Single Core)

### ■ Trap Gate

- Interrupt Enable Flag wird nicht gelöscht
- Verwendet für Systemaufrufe

- Externe- oder Hardware-Interrupts
  - von einem Gerät, z.B. Timer-Interrupt
  - Kommen von außerhalb, aus Sicht der CPU
- Interne- oder Software- Interrupts:
  - Kommen von der CPU selbst
    - Exceptions (siehe nächste Seite)
    - Oder durch die Assemblerinstruktion `INT <nr>`
      - Verwendet für Systemaufrufe (Linux, Windows NT, MacOS, MSDOS)

- **Fault** (dt. Störung):
  - kann behoben werden, z.B. Page Fault
  - CPU-Zustand wird gesichert & Adresse der Instruktion, die Fault ausgelöst hat
  
- **Trap** (~ dt. Falle):
  - ausgelöst durch speziellen Befehl, z.B. INT 3 (Breakpoint)
  - Programm kann fortgeführt werden
  
- **Abort** (dt. Abbruch):
  - bei schwerem Fehler
  - Auslöser oft nicht genau lokalisierbar
  - führt zum Restart (z.B. Double Fault)

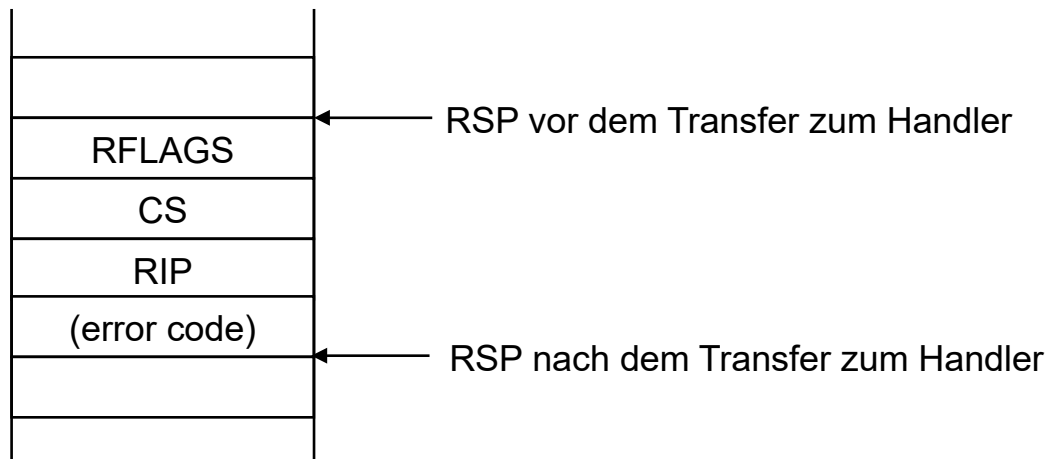


- Interrupts und Exceptions werden durch eine Vektornummer identifiziert
- 0 – 31 ist reserviert für x86 Exceptions
- 32 – 255 steht zur freien Verfügung
- Vektoren und ihre Bedeutung (Auszug)

Vektor	Bedeutung	Vektor	Bedeutung
0	Division by 0	11	Segment fault
1	Debug	12	Stack overflow
2	NMI	13	General protection fault
3	Break	14	Page fault
4	Overflow	16	-
5	Bounds range exceeded	18	Machine check
6	Illegal instruction	...	...
8	Double fault		

- Falls keine Privilegstufe gewechselt wird, so wird auch der Stack nicht umgeschaltet
- D.h. der Interrupt-Handler verwendet den Stack des unterbrochenen Threads
- Dies entspricht dem Ablauf bei hhuTOS (unser gesamter Code läuft im Ring 0)
- Einen error code gibt es nur bei manchen Exceptions

**Stack des unterbrochenen Threads**



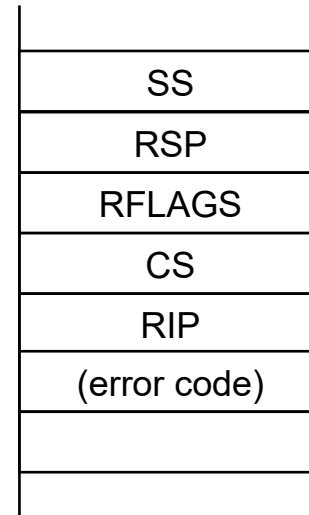
- Neuer Stack wird aus dem Task State Segment ermittelt (siehe letztes Kapitel)

Stack des unterbrochenen Threads



RSP vor  
dem Transfer  
zum Handler

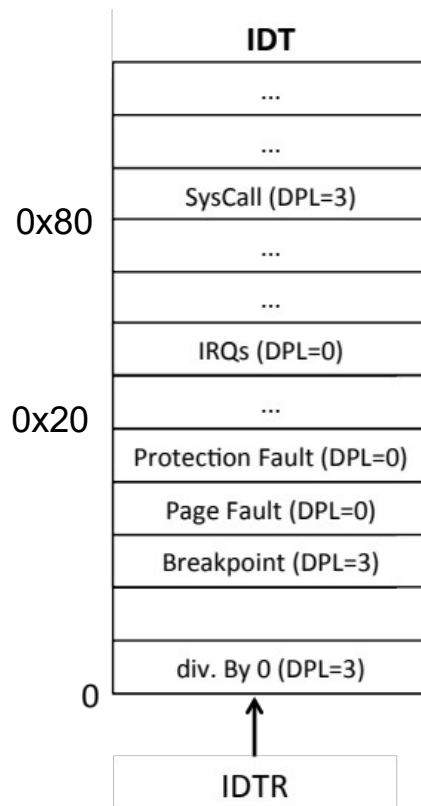
Handler Stack



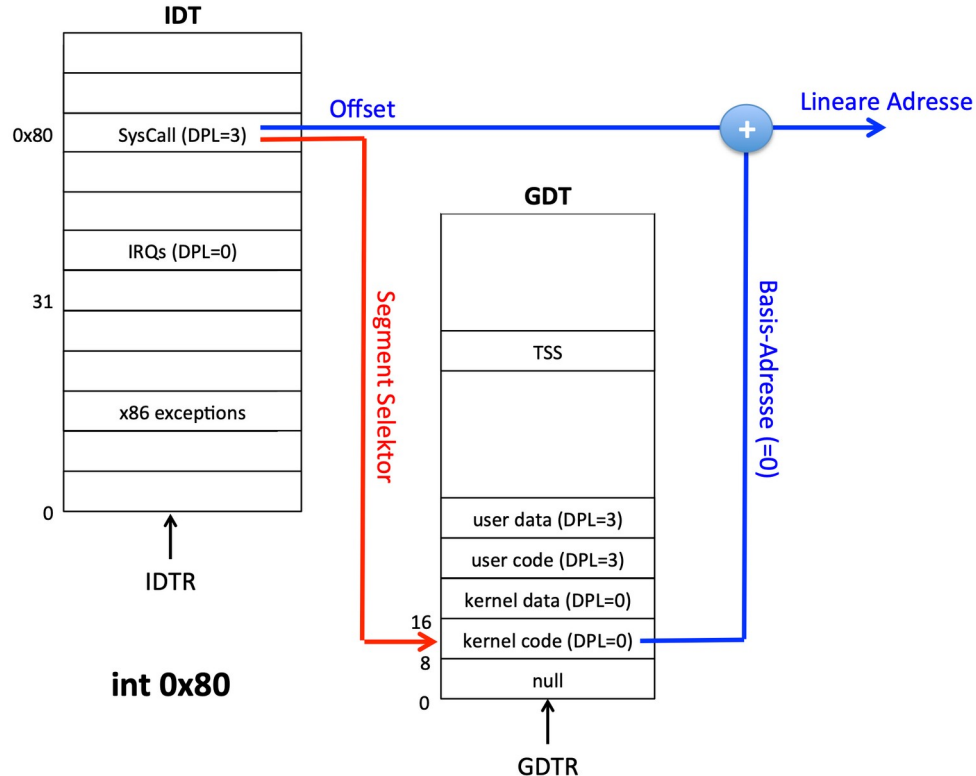
RSP nach  
dem Transfer  
zum Handler

# Beispiel unserer IDT

- Einträge 0-31 sind durch x86 reserviert für Exceptions
  - Interrupt Gates, alle DPL = 0
- Einträge 32 – 47 für externe Interrupts
  - Interrupt Gates, alle DPL = 0
  - IRQ0 != Vektor0
- Eintrag 0x80 für System-Aufrufe
  - Trap-Gate
  - DPL=3

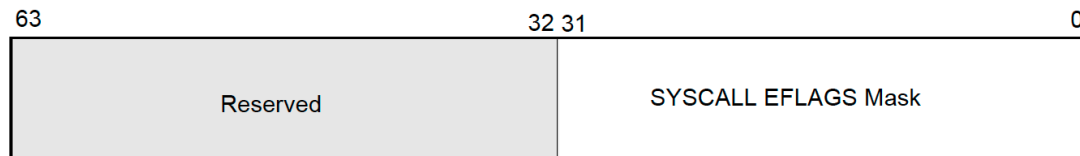


# Systemaufruf über ein Trap-Gate

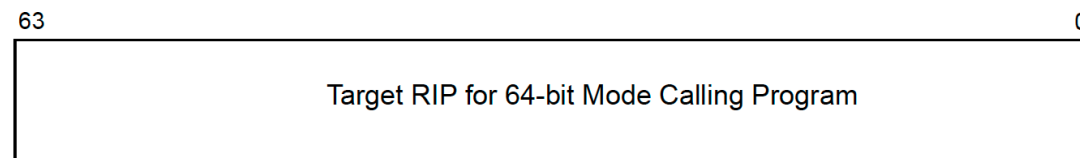


## 3.6 Syscall/sysret

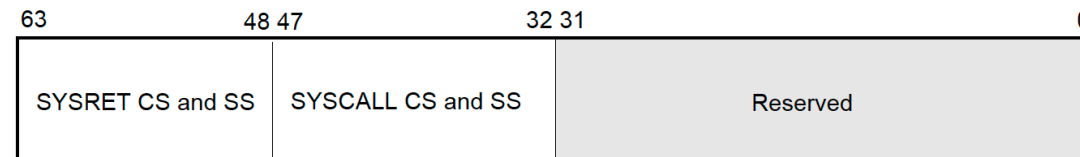
- Verwendet Model Specific Registers (MSRs) →
- Benötigt bestimmtes Layout in der GDT
- Stack muss händisch im Syscall-Handler umgeschaltet werden
- Verwenden FS und GS
  - Hier gibt es noch eine Basisadresse
  - Typischerweise speichert die Basisadresse in GS den Zeiger auf den Kernel-Stack



IA32\_FMASK



IA32\_LSTAR



IA32\_STAR