

# Operating Systems

Betriebssystem: Abstraktion von Hardware, Ressourcenverwaltung, Sicherheit/Schutz

<u>User Mode</u>	unprivilegierte Befehle teilweiser Zugriff auf Speicher	schützt Betriebssystem und Nutzerprozesse vor anderen Nutzerprozessen
<u>Kernel Mode</u>	alle Befehle kompletter Zugriff auf Speicher	bspw. privilegierte Instruktionen: ↳ Kontrollregister lesen/schreiben ↳ Interrupts de-/maskieren ↳ Gerätezugriff

Position-independent-Code / dynamische Linking:  
↳ direkt mit GOT/PLT

CPU-bound: Prozess arbeitet viel ohne I/O

I/O-bound: kurze Berechnung + viel I/O

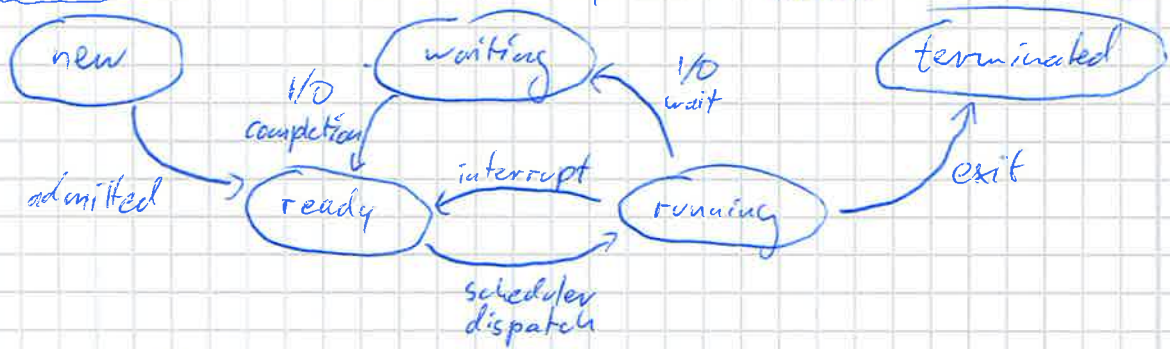
Betriebssystem schaltet sich ein:	Interrupts	Exceptions	System Call
freiwillig	nein	nein	ja
synchron	nein	ja	ja

Trap-Instruktion: Prozessor wechselt in Kernel Mode  
Syscall-Nr wird übergeben

Syscall-Parameter: Register / Stack / Pointer im Hauptspeicher

Scheduling : Prozesse  $\leftrightarrow$  Ressourcen  
 CPU scheduling  
 I/O scheduling

long-term : Prozess zur ready-Queue hinzufügen (admit)  
short-term : Prozess aus ready-Queue ~~aus~~ auf CPU ausführen (dispatch)



Metriken für scheduling policy:

- Prozessor-Ausnutzung (processor utilization)
- Durchsatz (throughput)
- Zeit von Einreichen bis Fertigstellung (turn around time)
- Wartezeit (waiting time)
- Zeit von Einreichen bis erstem Laufen (response time)

nicht-kooperatives Scheduling braucht Hardware-Unterstützung (Interrupts), weil Prozesse nicht auf CPU laufen (direct execution)

kurze Zeitscheiben

- viele Kontext-Switches
- Interaktivität
- kurze Ladezeiten

lange Zeitscheiben

- wenige Kontext-Switches
- wenig Overhead
- höherer Durchsatz (throughput)

Scheduling Policies Time-Slices werden wiederholt an Prozesse  
Round Robin : Prozesse werden in der Reihenfolge gegeben, wie sie gekommen sind.  
Shortest Job First (präemptiv: periodisch; bei neuen Job)

Priority Scheduling Prozess  $\rightarrow$  Prio, Prozess mit max Prio läuft  
 Starving, wenn mehrere Prozesse selbe Prio.

Multi-Level Feedback Queue

mehrere Warteschlangen mit verschiedenen Prios  
 Prios (hohe Warteschlange  $\rightarrow$  hohe Prio)  
 Zeitscheiben (hohe Warteschlange  $\rightarrow$  kurze Zeitscheibe)

Wille Prozess aus höchster Warteschlange (round robin)  
 Prozess steigt in niedrige Warteschlange (prio), wenn  
 Zeitscheibe aufgebraucht, selbst er wurde blockiert.

kooperatives Scheduling: Prozesse geben durch yield() CPU freiwillig ab.

Präemptiv = erzwungenes Scheduling  
 regelmäßiger Interrupts  
 Nutzer muss nicht um Scheduling kümmern (yield)  
 verändert nicht automatisch Scheduling.



# Process Switching

Process Control Block: instruction pointer, stack pointer  
 general purpose registers, Adressraum,  
 Scheduler-Info (Pri), offene Dateien

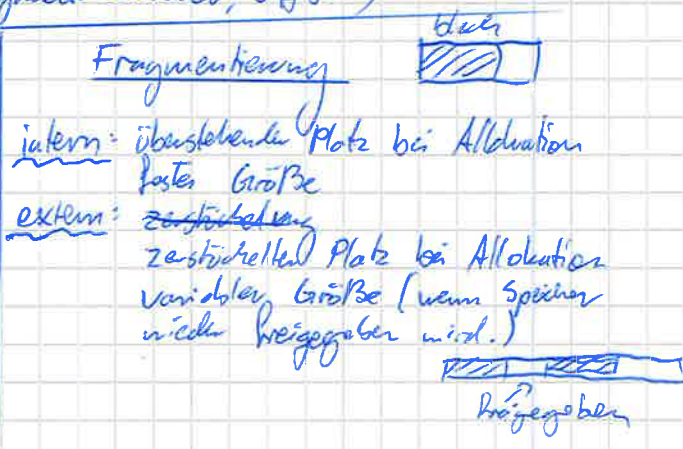
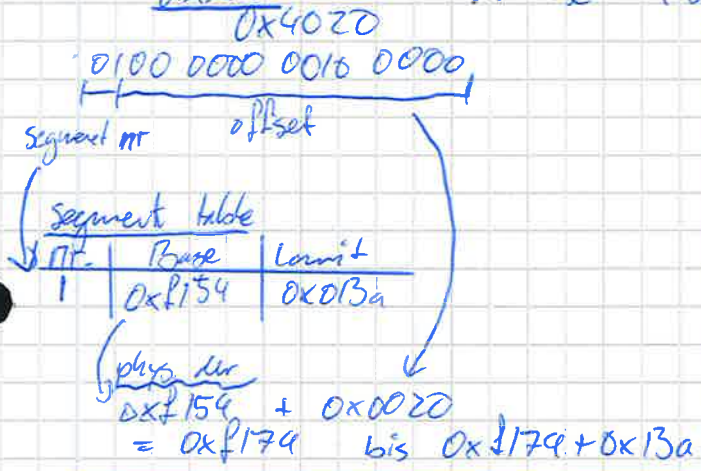
Threads	Kernel-level	User-level	Kernel-Mode
	= One-to-One	= Many-to-One	
	Kernel hat Threads Process + Thread Tabelle	Userspace hat Threads + Thread-tabelle Kernel hat Prozess-tabelle kein Parallelismus Blaßieren als Threads bei Syscall	

Many-to-Many (= Hybrid Threads)  
 $N$  User-Threads  $\leftrightarrow$   $M$  Kernel-Threads ( $N \geq M$ )  
 Parallelismus, kein komplettes Blockieren

## Segmentation

Virtuelle Adressen: genutzt von Anwender / Betriebssystem  
Physische Adressen: Hauptspeicher, MMU

Segmenttabelle: Segment (Base, Limit)  
 virt. adr. Adresse (Segmentnummer, Offset)



Memory Allocation Policies: First-fit, Best-fit, Worst-fit

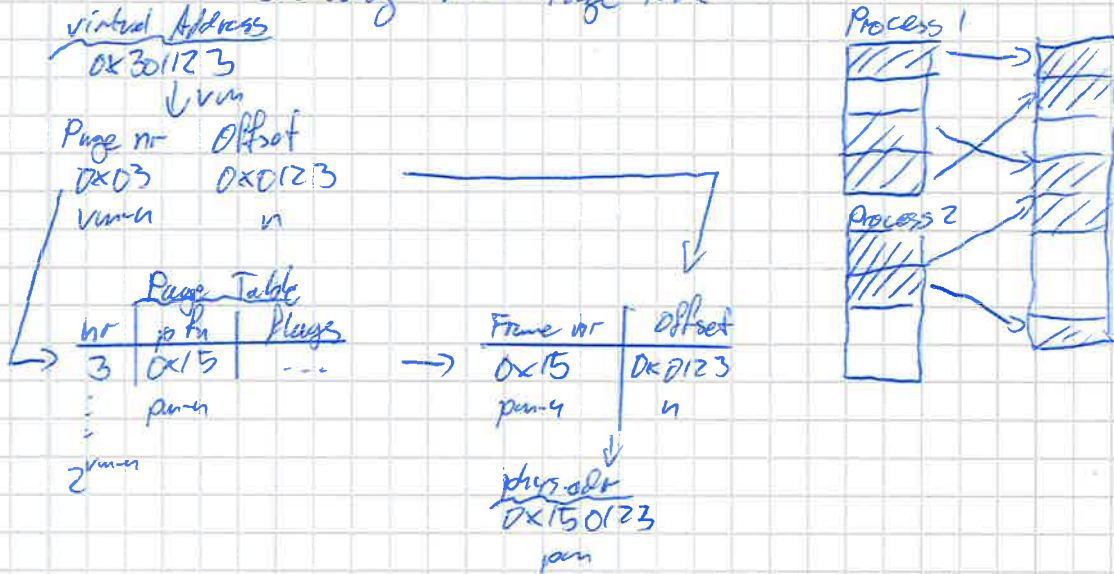
Buddy-Allocator: Speicherbereiche der Größe  $2^n$ ; Große Blöcke werden rekursiv geteilt  
 Kombinieren von freien Buddies; Schneller Zugriff durch Free-Lists



Slab-Cache: verwaltet Objekte selber Größe  
 keine interne Fragmentierung  
 keine externe Fragmentierung



Paging    virtuellen  $\rightarrow$  physischen Speicher  
 Pages  $\rightarrow$  Frames  
 jede Prozess hat zusammenhängende virt. Adressraum  
 Zuordnung mehr Page Table



Multilevel Page Table: virt. adr. hat mehrere Page nr  
 Tabelle muss nicht vollst. im Speicher sein  
 Adressübersetzung komplizierter

Inverted Page Table: eine Page Table im gesamten System  
 Zuordnung: physisch  $\rightarrow$  virtuell  
 PTE  
 Eintrag: VPN + PID  
 platz sparend aber lineare Suche

Hashed Inverted Page Table

Translation-Lookaside-Buffer (TLB): Cache für Page-Table  
 sucht erst in TLB nach  
 erst dann in PT und updated TLB

Copy-on-write: bei hoch geladenen beide Prozesse  
 auf Speicher zu, bis ein Prozess  
 auf Speicher schreibt, erst dann wird  
 kopiert

Page-Replacement-Policy: First-In; First-Out; Optimal (Theory)  
 Least-Recently-Used;  
 Clock  
 Random



## Coales

Zeitliche Kohärenz  
Räumliche Kohärenz

## Interprocess Communication (IPC)

- geteilter Speicher
- über OS: MessageQueue, Pipes, Sockets

Synchron: Sender wartet bis Nachricht bei Empfänger

Asynchron: Sender wartet bis Nachricht im Nachrichtensystem

## Race Conditions

Critical Section: Code-Abschnitt, in dem auf geteilte Ressourcen zugegriffen wird. Nur je ein Thread darf zugreifen.

Entry Section: Eintritts-Anfrage

Critical Section: geteilte Daten modifizieren

Exit Section: Abschnitt beenden.

## Anforderung an Synchronisations-Lösung

- Mutual Exclusion (gegenseitiger Ausschluss)
- Progress (Fortschritt)
- Bounded Waiting (begrenzte Warten)
- Performance (Leistung)

## Spinlock

```
lock() { while (locked) { wait(); } locked = true; }
```

~~lock() { while (locked) { wait(); } locked = true; }~~

benötigt atomare xchg Operation!

User-space, kurze Wartezeiten

## Mutex

lege thread schlafen und wecke ihn wenn unlocked wurde.

~~Kernel~~ Kernel-Erweiterung.

für lange Wartezeiten.

## Futex

Spinlock als Anfang, dann Mutex

## Semaphore

wait(): counter -- ; thread schlüft, wenn counter < 0  
signal(): counter ++ ; sonst thread wecken

## Condition Variable

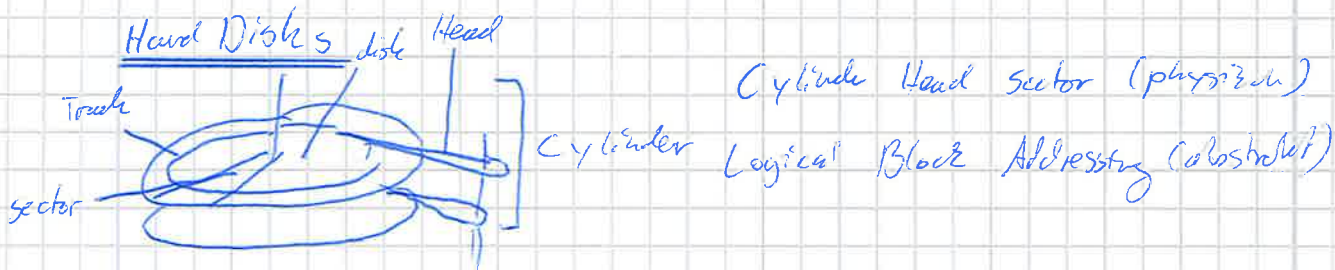
Mutex mit Bedingung

~~wait~~ sign

wait(cond, mutex), signal(cond), broadcast(cond)

Deadlock: kein Fortschritt, weil mehrere Threads auf gegenseitig zugewiesene Ressourcen warten

- Voraussetzungen:
- Mutual Exclusion: Ressourcen können nicht geteilt werden
  - Hold and Wait: Gleichzeitiges Halten und Warten auf Ressourcen
  - No preemption: Ressourcen dürfen Prozess nicht genommen werden
  - Circular Wait: keine zirkulären Abhängigkeiten



Sector Spanning: ersche kaputte sektoren durch neue (Lokalität zerstört)

Sektor Slipping: kaputte sektoren überspringen

Stingel: mehr Daten, was veraltet/überlagert geschrieben, langsames schreiben

Solid State Drive

- beim schreiben muss ein kompletter Block <sup>gelöscht und</sup> neu geschrieben werden
- durch remapping kann auf einen leeren Block geschrieben werden und später gelöscht werden

Trim = SSD mitteilen, welche Puffer nicht mehr genutzt werden.

Raid

SLED (Single Large Expensive Disk) = hohe Kapazität + Geschwindigkeit  
Robust + Fehlertoleranz

Raid 0: Block-Striping = Blöcke werden einfach auf Disks verteilt

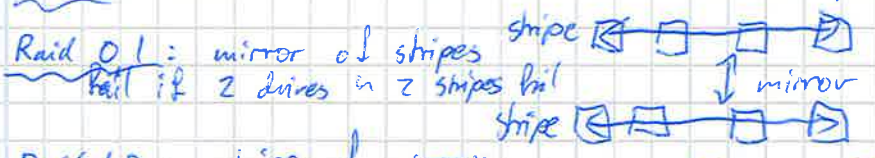
Raid 1: Mirroring

Raid 2: Für Disks werden ~~einige~~ Hamming-Codes auf anderen Disk gespeichert

Raid 3: word-interleaved-parity: parity-disk, high throughput

Raid 4: ~~block~~ block-interleaved-parity: parity-disk für mehrere Blöcke  
z.B. "or" = update braucht 2 reads + writes  
parity in Bottle neck

Raid 5: block-level distributed parity: distribute parity on all disks



Raid 10: stripe of mirror  
fail if 2 drives in same group fail  
lower probability of failure





## Disk Space Allocation

Contiguous Allocation: Datei hat Start-Adresse und Länge  
hat interne/externe Fragmentierung

Chained / Linked List Allocation:

wachsende Dateien variabler Größe  
keine externe Fragmentierung  
schlechte random R/W Performance

File Allocation Table: Tabelle mit Zeigern zu mehreren  
Einträgen

variable Dateigröße  
kann sehr groß werden  $\rightarrow$  hohe RAM Kosten

Indexed Allocation: Ein Block hat Liste von Pointern auf  
Datenblöcke

variable Dateigröße  
schnelles zufälliges Lesen/Schreiben  
geringe RAM-Kosten  
Platzverschwendung für kleine Dateien

Inode

direct	blocks
single	indirect
double	indirect
triple	indirect

Inhalt: Größe  
Typ = regulär, dir, link  
Zugriffsrechte  
Timesstamps  
Anzahl Hardlinks  
Nummern der Datenblöcke