

# Kernels Make The OS

<https://1.mafiasi.de/kernel-kbs>

---

Bene 17ostendo

Finn 17sell

15. Januar 2024

Well actually, Linux is just the Kernel. What you're talking about is the GNU/Linux Operating System...

**Warum eigentlich einen Kernel entwickeln?**

---

## Linux Syscalls

read write open close stat fstat lstat poll lseek  
mmap mprotect munmap brk rt\_sigaction  
rt\_sigprocmask rt\_sigreturn ioctl pread64 pwrite64  
readv writev access pipe select sched\_yield mmap  
msync mincore madvise shmget shmat shmctl dup  
dup2 pause nanosleep getitimer alarm setitimer  
getpid sendfile socket connect accept sendto  
recvfrom sendmsg recvmsg shutdown bind  
getsockname getpeername socketpair setsockopt  
getsockopt clone fork vfork execve exit wait4 kill

clone fork vfork execve exit wait4 kill uname semop semop semctl shmctl msgctl msgctl msgctl msgctl msgctl flock fsync fdatsync truncate ftruncate getdents getcwd chdir fchdir rename mkdir rmdir creat link unlink symlink readlink chmod fchmod chown fchown lchown umask gettimeofday getrlimit getrusage sysinfo times ptrace getuid syslog getgid setuid setgid geteuid getegid setpgid getppid getpgrp setsid setreuid setregid getgroups setgroups setresuid getresuid setresgid getresgid setpgid setsuid setfsgid getsid capget capset rt\_sigpending rt\_sigtimedwait rt\_sigqueueinfo rt\_sigsuspend sigaltstack utime mknod uselib ustat statfs fstatfs sysfs getpriority setpriority sched\_setparam sched\_getparam sched\_setscheduler sched\_getscheduler sched\_get\_priority\_max sched\_get\_priority\_min sched\_rr\_get\_interval mlock munlock munlockall vhangup modify\_ldt pivot\_root \_\_sysctl prctl arch\_prctl adjtimex setrlimit chroot sync acct settimeofday mount umount2 swapon swapoff reboot sethostname setdomainname iopl ioperm create\_module delete\_module get\_kernel\_syms nfsservctl readahead setattr lsetattr fsetattr getxattr lgetxattr fgetxattr listxattr llistxattr flistxattr removexattr lremovexattr fremovexattr tkill time futex sched\_setaffinity sched\_getaffinity set\_thread\_area io\_destroy io\_getevents io\_submit io\_cancel get\_thread\_area epoll\_ctl\_old getdents64 set\_tid\_address restart\_syscall semtimedop fadvise64 timer\_create timer\_settime timer\_gettime timer\_getoverrun timer\_delete clock\_settime clock\_gettime clock\_getres clock\_nanosleep exit\_group epoll\_wait epoll\_ctl tgkill utimes vserver set\_mempolicy get\_mempolicy mq\_open mq\_unlink mq\_timedsend mq\_timedreceive mq\_notify mq\_getsetattr kexec\_load waitid add\_key request\_key keyctl ioprio\_set ioprio\_get inotify\_init inotify\_add\_watch inotify\_rm\_watch migrate\_pages openat mkdirat mknodat fchownat futimesat newfstatat unlinkat renameat linkat readlinkat fchmodat faccessat pselect6 ppoll unshare set\_robust\_list get\_robust\_list splice tee sync\_file\_range vmsplice move\_pages utimensat epoll\_pwait signalfd timerfd\_create eventfd fallocate timerfd\_settime timerfd\_gettime accept4 signalfd4 eventfd2 epoll\_create1 dup3 pipe2 inotify\_init1 preadv pwritev rt\_tsigqueueinfo perf\_event\_open recvmmsg fanotify\_init fanotify\_mark prlimit64 name\_to\_handle\_at open\_by\_handle\_at clock\_adjtime syncfs sendmmsg setns getcpu process\_vm\_readv process\_vm\_writev kcmp finit\_module sched\_setattr sched\_getattr renameat2 seccomp getrandom memfd\_create kexec\_file\_load bpf execveat userfaultfd membarrier mlock2 copy\_file\_range preadv2 pwritev2 pkey\_mprotect pkey\_alloc pkey\_free statx io\_pgetevents rseq 't numbers 387 through 423, add new calls after the last " pidfd\_send\_signal io\_uring\_setup io\_uring\_enter io\_uring\_register open\_tree move\_mount fsopen fsconfig fsmount fspick pidfd\_open clone3 close\_range openat2 pidfd\_getfd faccessat2 process\_madvise epoll\_pwait2 mount\_setattr landlock\_create\_ruleset landlock\_add\_rule landlock\_restrict\_self memfd\_secret process\_mrelease futex\_waitv set\_mempolicy\_home\_node cachestat fchmodat2 map\_shadow\_stack futex\_wake futex\_wait futex\_requeue statmount listmount lsm\_get\_self\_attr lsm\_set\_self\_attr lsm\_list\_modules rt\_sigaction rt\_sigreturn ioctl readv writev recvfrom sendmsg recvmmsg execve ptrace rt\_sigpending rt\_sigtimedwait rt\_sigqueueinfo sigaltstack timer\_create mq\_notify kexec\_load waitid set\_robust\_list get\_robust\_list vmsplice move\_pages preadv pwritev rt\_tsigqueueinfo recvmmsg sendmmsg process\_vm\_readv process\_vm\_writev setsockopt getsockopt io\_setup io\_submit execveat preadv2 pwritev2

## Linux Kernel Features

- Dateisystem (-e)
- Nutzer (-Rechte)
- PAM
- Netzwerk
- Firewall x2

und vieles mehr...



Linux

13 IN 1

Memory Management

Time Management

CPU Isolation

File Systems

Sound

Rendering

M-Access-Control

Networking

Logging

Kexec

FRESH ENERGY

Users (as in people)

Crypto API

+5 vitamins

eBPF

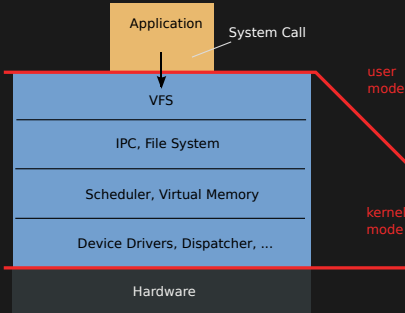
## Micro vs. Monolithic

- Principle of Least Privilege
- Principle of Least Scope

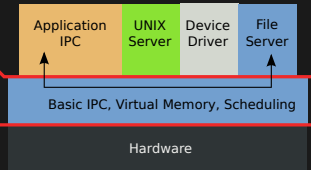
Ein Kernel **sollte** irgendwann *fertig* sein!



## Monolithic Kernel based Operating System



## Microkernel based Operating System



# Delegation!

### Permissions?

→ Reicht nicht!

→ Mandatory-Access-Control auch nicht! (AppArmor, seccomp, pledge)

→ DenyListing: Alles doof...

**Capability** Token, das Resource identifiziert  
und Zugriff erlaubt<sup>1</sup>

---

<sup>1</sup>Think: File Descriptor



Everything  
is a File



Everything  
is a  
Capability

# Intro to Kernel Development

---

## Was ist nochmal ein Computer?

- + CPU
- + RAM
- + Diverse Geräte
- + ROM

## Bootstrapping

1. Gerät lädt Firmware vom ROM
2. Firmware lädt Bootloader
3. Bootloader lädt Kernel
4. Kernel lädt OS

- **Ein** Computer vs **mehrere** Programme
- **Ein** Computer vs **mehrere** Benutzer

Der Kernel macht aus einem physikalischen Computer mehrere logische Computer.



### Multiplexing!

- Rechenzeit
- Prozessorzustand
- Speicherplatz
- Gerätezugriff

## Hardwareunterstützung

- Privilegierter CPU-Modus
- Speichersegmentierung
- Virtuelle Addressierung
- Preemption
- Timer

## Kernel

- Gott-Zugriff
- konfiguriert Speicherisolierung
- wechselt regelmäßig Prozesse
- Syscalls

not\_kernel.rs

```
import std;  
malloc()  
println!("Hi")  
let x = new MyObject()  
let y = List::new()
```

kernel.rs

```
# boot  
todo!()
```

```
while true:  
    todo!()
```

# Time-Management

---

```
# boot
```

```
todo!()
```

```
while True:
```

```
    proc = schedule()
```

```
    exec_user_proc(proc)
```

# RISC-V CPU State

**Regs** 32x 64bit Register

**Float Regs** 32x 64bit IEEE Floating Point

**PC** Aktuelle Stelle im Code



## Switch to Userspace

```
func exec_user_proc(proc):  
    cpu.regs = proc.load_state()  
    cpu.pc = proc.pc  
    cpu.mode = "user"  
    unreachable!()
```

Problem: Gierige Prozesse!

→ Regelmäßige Unterbrechungen!

→ Interrupts!

## RISC-V Privileged Registers

**ie** Interrupt Enabled

Timer, External, Software

**ip** Interrupt Pending

**tvec** Interrupt Handler Address

```
# boot  
todo!()
```

```
while True:
```

```
    proc = schedule()
```

```
    setup_timer(5ms) ← new
```

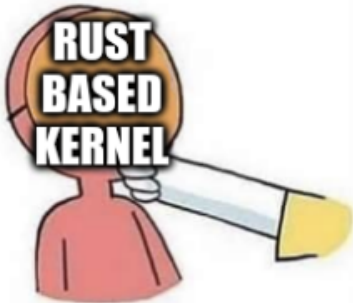
```
    interrupt = exec_user_proc(proc)
```

```
    handle_interrupt(interrupt) ← new
```

HEY **RUST BASED KERNEL** WHY  
DO YOU ALWAYS  
WEAR THAT MASK?



**RUST  
BASED  
KERNEL**



**ASSEMBLY  
TRAP  
HANDLER**



LET'S KEEP  
THIS ON.



## RISC-V Privileged Registers. Cont

**scratch** Scratch

**epc** Exception Program Counter

**cause** Trap Cause

```
# boot
setup_trap_handler()           ← new

while True:
    proc = schedule()
    setup_timer(5ms)
    exec_user_proc(proc)
    label!("loop_resume")
```

```
#[asm]
```

```
func trap_handler():
```

```
    proc = cpu.scratch
```

```
    proc.save_state(cpu.regs)
```

```
    proc.pc = cpu.epc
```

```
    goto!("loop_resume")
```



# Multiplexing!

- Rechenzeit
- Prozessorzustand
- Speicherplatz
- Gerätezugriff

# Memory-Management

---

## Segmentation (-fault)

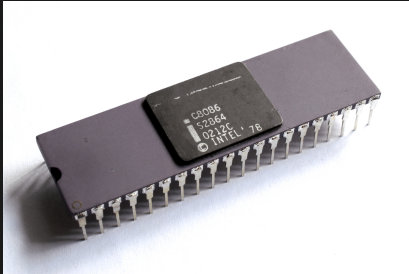


Abbildung 1: Intel iAPX 86

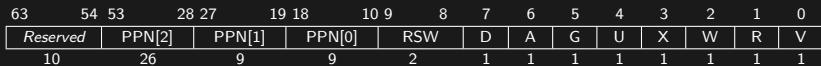
- + Introduces Memory Access Rights
- Fixed number of memory segments
- Fragmentation

## RISC-V Privileged Registers. Cont

**atp** Address Translation and Protection

**status** Verschiedene Status Bits

# Address Translation

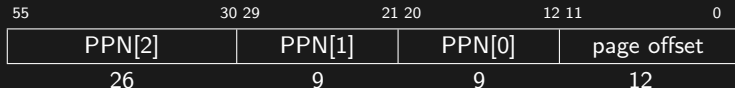


Sv39 page table entry.

# RISC-V Virtual Addressing



Sv39 virtual address.



Sv39 physical address.

```
func exec_user_proc(proc):
    cpu.regs = proc.load_state()
    cpu.atp = proc.pagetable    ← new
    cpu.pc = proc.pc
    cpu.mode = "user"
    unreachable!()
```

# Multiplexing!

- Rechenzeit
- Prozessorzustand
- Speicherplatz
- Gerätezugriff



# Devices

---

Geräte sind memory-mapped  
→ Memory können wir schon!

# Multiplexing!

- Rechenzeit
- Prozessorzustand
- Speicherplatz
- Gerätezugriff

Kernel fertig?