

Übung zu Betriebssystemtechnik

Aufgabe 4: Trennung von Kern & Anwendungen

11. Juni 2025

Dustin Nguyen, Maximilian Ott & Phillip Raffeck

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



Lehrstuhl für Informatik 4
Systemsoftware



Friedrich-Alexander-Universität
Technische Fakultät

Anwendungen sollen (sowohl bei Übersetzung als auch bei Ausführung) vom Kernel getrennt werden

Restrukturierung der Quellen

Organisation der Dateien

Bisher:

- Anwendungscode und Betriebssystemcode vermischt
→ alles in einer großen `system[64]`-Datei

Nun Aufteilung des Codes in unterschiedliche Verzeichnisse:

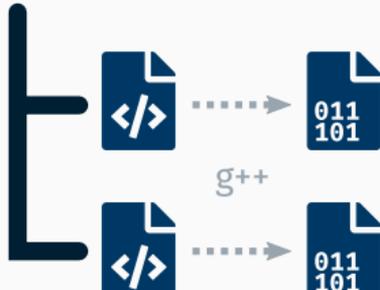
`kernel/` nur der Betriebssystemkern
→ `system[64]`

`libsys/` Unterstützung für Anwendungen
→ statische Bibliothek `libsys.a`

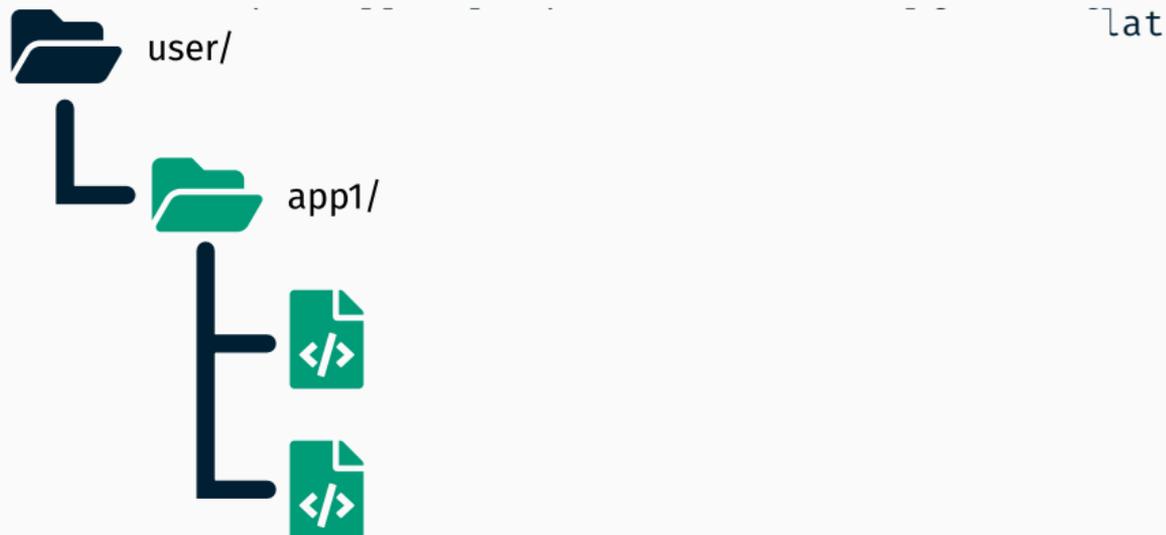
`user/` enthält **mehrere** Anwendungen
Jede Anwendung wird gegen `libsys` gelinkt
und zu einer eigenen Binärdatei kompiliert
→ Archiv `initrd.img` mit den einzelnen Binärdateien

libsys

```
$ g++ -fno-builtin -nodefaultlibs -nostdlib -nostdinc ... → $CXXFLAGS  
aus tools/build.mk berücksichtigen → oder gleich tools/common.mk in Makefile  
einbinden $ ar rcs libsys.a libsys/*.o
```



```
Linkerskript user/sections.ld mit 0x400 0000 (64 MiB) als Einsprungpunkt $  
g++ $CXXFLAGS -Wl,T sections.ld -o app1.elf $LDFLAGS init.o [...] $  
objcopy -O binary --set-section-flags \
```



C/C++ Laufzeitumgebung

- `crti.o` und `crtn.o` für C Funktionsprologe
 - für Initialisierungs- (`_init`) und Beendigungsroutine (`_fini`)
(siehe `compiler/crti.asm` sowie `compiler/crtn.asm`)
- `init.o` setzt die Laufzeitumgebung auf
 1. Einsprungspunkt ist `void start()`
 2. Ausführung der C Startup Initialisierungsroutinen
(`__preinit_array`, `_init()` und `__init_array`, siehe `compiler/libc.cc`)
 3. Übergabe an `main` der App
 4. bei Rückkehr Deinitialisierung
(`__fini_array`, `_fini()` und Endlosschleife)
 - Bereitstellung der für C++ benötigten Funktionen
(Dummies für `__cxa_pure_virtual()` und `__cxa_atexit()`, siehe `compiler/libcxx.cc`)



Beim Linken muss `init.o` die erste Objektdatei sein, damit `start()` die erste Funktion (am Einsprungspunkt an Adresse `0x400 0000`) ist!

Flat Binaries und Image Builder *(für 5 ECTS)*

Erstellung einer flachen Binärdatei

Man nehme eine ELF-Datei

0. prüfe, ob **start** an der Einsprungsadresse liegt

```
$ nm app1.elf | grep -q "4000000 T start" || echo "Fehler"
```

1. integriere **BSS** in die **Datensektion**

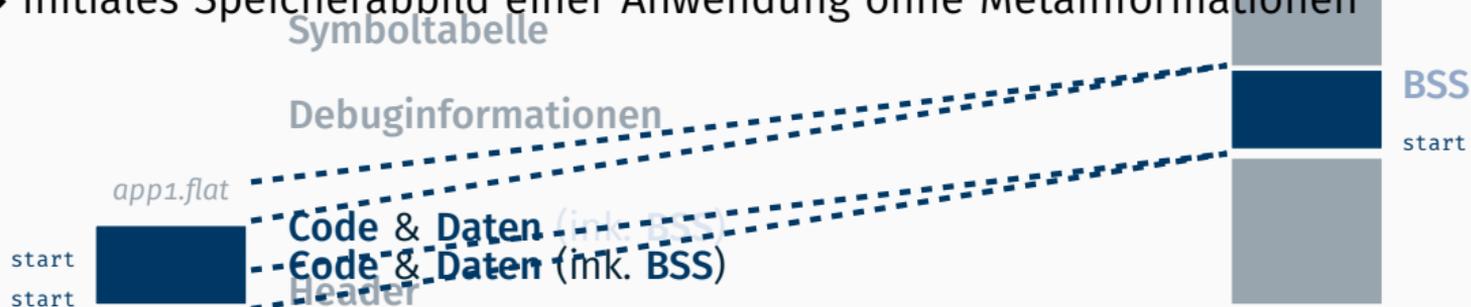
```
objcopy Argument: --set-section-flags .bss=alloc,load,contents
```

2. werfe alles außer **Code & Datensektion** weg

```
objcopy Argument: -O binary
```

Zack fertig: **Flat Binary**

→ initiales Speicherabbild einer Anwendung ohne Metainformationen



Anwendungsabbildarchiv als Boot Module

```
$ ./imgbuilder app1.flat app2.flat appx.flat > initrd.img
```



appx.flat



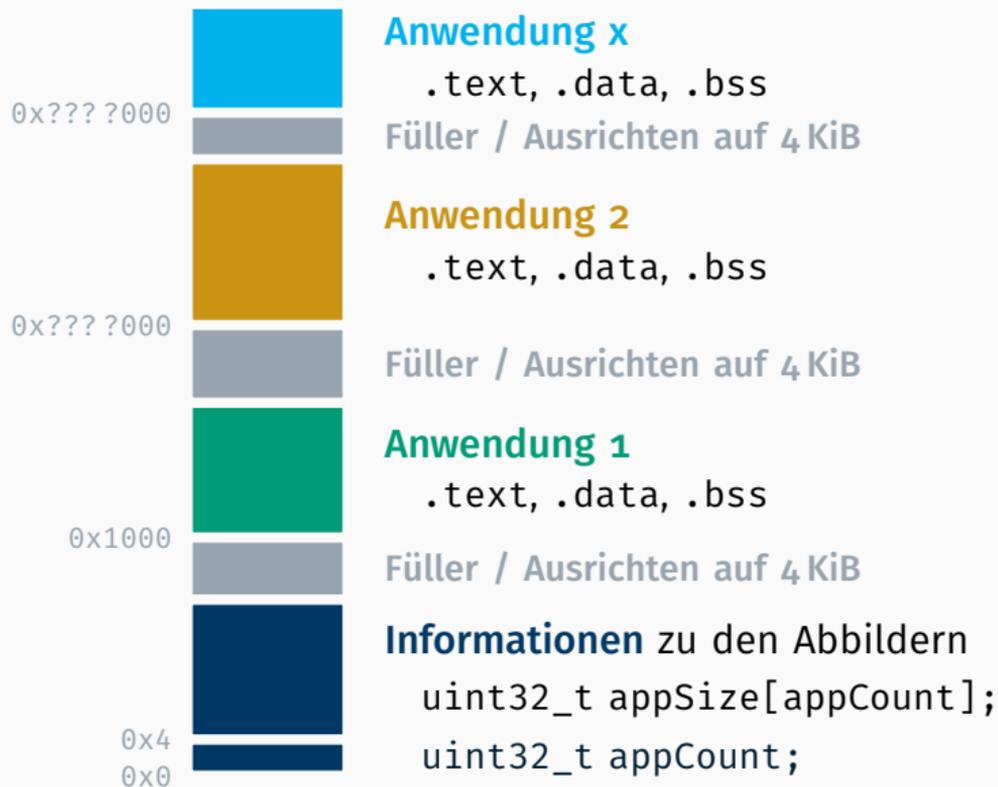
app2.flat



app1.flat



Image Builder Format



ELF und TAR (*für 7.5 ECTS*)

Executable and Linking Format

- 1993 vom *Tool Interface Standard (TIS)* Committee) spezifiziert
- Nachfolger von *a.out* (*Assembler output*, ab 1971) und *coff* (*Common Object File Format*, ab 1983)
- Standardformat auf unixoiden Systemen
- inzwischen plattformübergreifender Einsatz
- vielfältige Einsatzgebiete
 - relocatable** Objektdateien (.o)
 - executable** (statische) Programme
 - dynamic** dynamische Programme und [geteilte] Bibliotheken (.so)
- durch Sektionen flexibel erweiterbar
- bei Bedarf Einbettung von *DWARF* Debuginformationen

Die Vorgabe enthält unter `utils/` bereits das Grundgerüst (Parser) für einen einfachen Loader. Dieser soll

Program Header Table

- beschreibt die einzelnen Programmsegmente
- **Recht** **LOAD** von statischen Binärdateien (ohne Relok.) unterstützen
- jeweils mit Position & Größe in ELF-Datei und im Zielspeicher
- Berechtigungen (Writable / Executable) berücksichtigen
- verschiedene Typen: **LOAD**, **DYNAMIC**, **INTERP**
- (ELF) Quelldaten in **helle (Userspace) Zielseiten kopieren**
- und Attribute wie **les-, schreib- und ausführbar**
- optional unter bestimmten Bedingungen auch **Quellseite einblenden**





DISSECTED FILE

```
-S uname -m
x86_64
-S ./simple64.elf
Hello world!
```

SIMPLE64.ELF

HEADER^{2/2}

TECHNICAL DETAILS FOR IDENTIFICATION AND EXECUTION

SECTIONS

CONTENTS OF THE EXECUTABLE

HEADER^{2/2}

TECHNICAL DETAILS FOR LINKING AND/OR FOR EXECUTION

ELF HEADER

IDENTIFICATION AND EXECUTION

PROGRAM HEADER TABLE

EXECUTION ORION

CODE

EXECUTION PROMPT

DATA

INFORMATION USED BY THE CODE

SECTIONS' NAMES

SYMBOL TABLE

SECTION HEADER TABLE

HEXADECIMAL DUMP

ASCII DUMP

FIELDS

VALUES

EXPLANATION

HEXADECIMAL DUMP	ASCII DUMP	FIELDS	VALUES	EXPLANATION
7F 45 4C 46 82 01 01 00 00 00 00 00 00 00 00 00	..ELF.....	e_ident	0x7F, "ELF"	CONSTANT SIGNATURE #BITS, LITTLE-ENDIAN
02 00 3E 00 01 00 00 00 00 00 10 00 00 00 00	e_class, e_data	2, 1	ALWAYS 1 EXECUTABLE
40 00 00 00 00 00 00 00 F0 00 00 00 00 00 00	e_machine	0x3E	AMD64 (X86-64)
00 00 00 00 40 00 10 00 01 00 40 00 03 00	e_version	0x10000000	ALWAYS 1
		e_entry	0x10000000	ADDRESS WHERE EXECUTION STARTS
		e_shoff	0x40	PROGRAM HEADERS OFFSET
		e_shstrndx	1	SECTION HEADERS OFFSET
		e_shsize	0x40	ELF HEADERS SIZE
		e_shentsize	0x30	SIZE OF A SINGLE PROGRAM HEADER
		e_shnum	4	COUNT OF PROGRAM HEADERS
		e_shndx	0	COUNT OF SECTION HEADERS
		e_strndx	0	INDEX OF THE FIRST SECTION IN THE TABLE

HEXADECIMAL DUMP	ASCII DUMP	FIELDS	VALUES	EXPLANATION
01 00 00 00 05 00 00 00 00 00 00 00 00 00 00	2_e_type	1	THE SEGMENT SHOULD BE LOADED IN MEMORY
00 00 00 00 00 00 00 00 10 00 00 00 00 00	2_e_flags	0x0	REMOVABLE AND EXECUTABLE
00 00 00 00 00 00 00 00 00 00 00 00 00 00	2_e_offset	0	OFFSET WHERE IT SHOULD BE READ
00 00 00 00 00 00 00 00 00 00 00 00 00 00	2_e_addr	0x10000000	VIRTUAL ADDRESS WHERE IT SHOULD BE LOADED
00 00 00 00 00 00 00 00 00 00 00 00 00 00	2_e_align	0x10000000	PHYSICAL ADDRESS WHERE IT SHOULD BE LOADED
00 00 00 00 00 00 00 00 00 00 00 00 00 00	2_e_entsize	0x30	SIZE OF ENTRY
00 00 00 00 00 00 00 00 00 00 00 00 00 00	2_e_estrndx	0x0	SIZE OF MEMORY

HEXADECIMAL DUMP	ASCII DUMP	FIELDS	VALUES	EXPLANATION
40 00 00 00 00 00 00 00 00 00 C7 C6 C0 00 00	3_0x64 ASSEMBLY		EQUIVALENT C CODE
10 40 C7 C7 01 00 00 00 48 C7 C6 31 00 00 00	3_1		
05 48 C7 C7 01 00 00 00 48 C7 C6 3C 00 00 00	3_2		
05 48 C7 C7 01 00 00 00 48 C7 C6 3C 00 00 00	3_3		

HEXADECIMAL DUMP	ASCII DUMP	FIELDS	VALUES	EXPLANATION
00 2E 73 68 73 74 72 74 61 62 00 2E 74 65 78 74_shstrtab..text	SECTION NAMES		
00 2E 72 6F 64 61 74 61 00_rodatab..rodatab			

INDEX	NAME	TYPE	FLAGS	CLASS	ADDRESS	OFFSET	SIZE
0	.text	PROGBITS	0	0	0x10000000	0x0	0x31
1	.rodatab	PROGBITS	2	0	0x10000000	0x40	0x0
2	.shstrtab	PROGBITS	0	0	0x10000000	0x40	0x31

THIS IS THE WHOLE FILE, HOWEVER, MOST ELF FILES CONTAIN MANY MORE ELEMENTS. ELEMENTS NOT DISPLAYED HERE CORRESPOND TO THE CONVENTIONS.

Bandarchivierer (TAPE ARCHIVER)

- 1979 eingeführte Archivdatei
 - *klebt* Dateien (unkomprimiert/unmodifiziert) aneinander
 - ausgerichtet an 512 Byte Grenzen
 - jeweils mit Kopfeintrag (beinhaltet Dateiname, Größe, Berechtigungen)
 - aber kein Index über alle Inhalte!
- für Bandlaufwerke (= sequentielle Ein-/Ausgabe) gedacht
- 1988 POSIX-Standardisierung und Einführung der *UStar* (= *Unix Standard*) Erweiterung (mehr Attribute, längere Dateinamen)
 - zusätzliche Informationen im Kopfeintrag
 - Format bleibt aber weiterhin kompatibel
- 1997 *pax* (= *portable archive formats*) Erweiterung
 - in großen Teilen kompatibel zu *UStar*
 - zur Befriedung der *Tar-Wars* (mit dem konkurrierenden *cpio*)

Bandarchivierer (TAPE ARCHIVER) Verwendung

```
$ tar -cf initrd.tar app1.elf app2.elf appx.elf
```



appx.elf



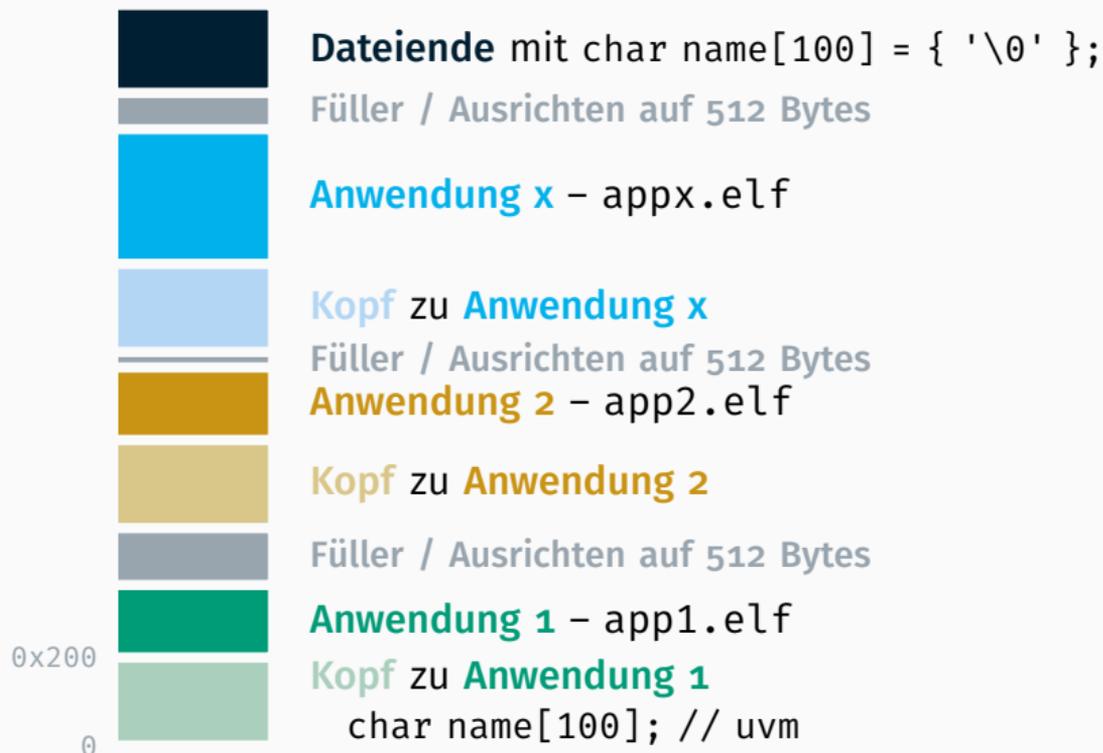
app2.elf



app1.elf



TAR Format



→ (rudimentärer) TAR-Parser bereits in Vorgabe (unter `utils/`)
enthalten!

TAPE ARCHIVE



ANGE ALBERTINI

<http://www.corkami.com>



```
$ tar -xOf hello.tar hello.txt
Hello World!
```

```
  0  1  2  3  4  5  6  7  8  9  A  B  C  D  E  F
0000: .h .e .l .l .o . . .t .x .t
0060:      .0 .0 .0 .0 .6 .4 .4 00 .0 .0 .0 .0
0070: .7 .6 .4 00 .0 .0 .0 .1 .0 .4 .0 00 .0 .0 .0 .0
0080: .0 .0 .0 .0 .0 .1 .5 00 .1 .2 .4 .2 .0 .0 .1 .0
0090: .5 .3 .2 00 .0 .1 .4 .6 .3 .6 00 20 .0
0100:  .u .s .t .a .r 00 .0 .0 .A .n .g .e
0120:                                     .A .d .m .i .n .i .s
0030: .t .r .a .t .o .r .s
-----
0200: .H .e .l .l .o 20 .W .o .r .l .d .! 0A
2800: ]
```

FILE HEADER

CONTENTS

FIELDS

VALUES

file name	hello.txt
file mode	0000644
owner user ID	0000764
group user ID	0001040
file size	0000013
timestamp	2014-10-16 20:41
checksum	014636 \0\x20
type flag	00 REGTYPE
magic	ustar\x00
version	"00"
owner user name	Ange
owner group name	Administrators
contents	Hello World!\n

TAR WAS INITIALLY DESIGNED FOR TAPE DRIVES, IN 1979:

- NO COMPRESSION, BLOCK ALIGNED

- NUMERIC VALUES ARE STORED IN OCTAL, ENCODED IN ASCII

TAR IS OFTEN COMBINED WITH GZIP, BZIP2 OR LZMA.

THE TAR FORMAT EVOLVED:

THIS EXAMPLE IS A "USTAR" FILE, AS DEFINED IN 1988

Unsere Makefile unterstützt bereits ein Bootmodul, dessen Pfad in der Variable INITRD definiert werden muss.

Beim Target netboot wird diese Datei nach /proj/i4stubs/student/LOGIN/initrd.img kopiert und automatisch vom Bootloader geladen (d.h. die Angabe in Multiboot::Module::getCommandLine() ist nicht sonderlich aussagekräftig, sondern beim Benutzer immer gleich)

Dynamisch wachsender Stapel

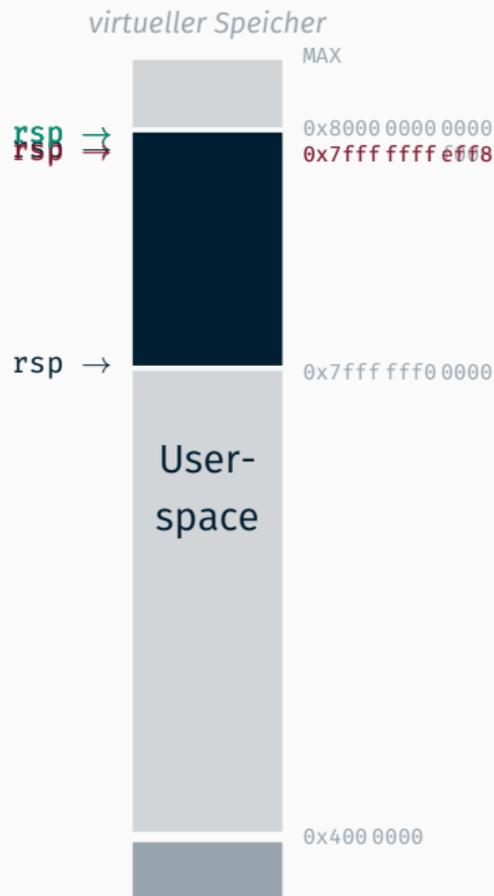
Ablauf:

- initial z.B. 4 KiB Userspace Stack
- sobald Stack voll → **Seitenfehler**
- Behandlungsroutine allokiert mehr Speicher (und setzt Ausführung fort)
- bis zu einem vorgegebenen Maximum

Umsetzung mittels neuem pagefault_handler

- Verwendung von `cr2` und `error_code`
- bei Änderung des aktiven Mappings an das Spülen des TLBs (z.B. via `invlpg`) denken!

Protipp: Initial gar keinen Stack allokierten!





Der Stack muss bei Funktionsaufrufen (d.h. vor `call`) an 16 Byte ausgerichtet sein!

- gemäß System V ABI 3.2.2
- für FPU & SSE Instruktionen wichtig
- einfache Überprüfung mittels `cmpxchg16b` Instruktion:

```
__int128 x = 0;  
__sync_val_compare_and_swap(&x, 0, 1);
```

und mit zusätzlichem GCC Parameter `-mcx16` übersetzen
→ *Protection Fault* bei unausgerichtetem Stack!

***Optional:* Fließkommazahlen**

Geschichte der Fließkommazahlen auf dem x86

Gleitkomma-Operationen beim **x86** entweder in Software oder

- Anfangs Koprozessor (*Intel 8087* → **x87**)
- ab 486er Gleitkommaeinheit integriert in CPU
 - acht 80 bit Datenregister (16 bit Exponent + 64 bit Mantisse)
→ organisiert als Stack (st0 – st7)
 - zudem weitere Kontrollregister
 - Befehle für Stack- und Rechenoperationen
→ Parallele Ausführung, Fehler als Exceptions
- später Erweiterungen *Multi Media Extension (MMX)*, *Streaming SIMD Extensions (SSE)* [1...4] und *Advanced Vector Extensions (AVX)*
 - x64 unterstützen immer (mindestens) SSE2
 - min. sechzehn 128 bit Register (xmm0 – xmm15)
 - Nutzung (nach System V ABI 3.2.3) auch für Parameterübergabe

→ Zustandssicherung beim Kontextwechsel notwendig!

Früher:

Sicherung der zusätzlichen Register beim Kontextwechsel ist teuer, aber viele Programme brauchen eigentlich gar keine FPU...

- OS schaltet FPU standardmäßig aus und sichert beim Kontextwechsel nur die üblichen *general purpose* Register
- falls ein Prozess die FPU benutzt, gibt es eine Exception
- OS aktiviert FPU, führt die Instruktion nochmal aus und kümmert sich bei diesem Programm von nun an auch um die FPU Register

Heute:

- Übersetzer nutzen standardmäßig auch xmm-Register
- schnelle Spezialbefehle (wie `fxsave` & `fxrstor`)
 - hochoptimiert (84 & 44 Zyklen)
 - für alle FPU / MMX / SSE Register (512 bytes)

→ standardmäßig im Userspace aktiviert

FPU-Umsetzung in STUBSML

- FPU / MMX / SSE bei Übersetzung von Userspace Apps aktivieren (→ `CXXFLAGS_NOFPU`) – aber Kernel soll weiterhin ohne FPU laufen!
- beim Systemstart pro Kern FPU initialisieren (→ `FPU::init()`)
 1. Bits für *Software Emulation* (`CR0_EM`) und *Task Switched* (`CR0_TS`) in Kontrollregister 0 (`cr0`) löschen, *Monitor coprocessor* (`CR0_MP`) setzen
 2. FPU selbst auf Standardwerte initialisieren (`fninit`), Statuswort (muss 0 sein) sowie Kontrollwort prüfen (muss `0x37f` sein)
 3. Bits für *Enable SSE Instructions* (`CR4_OSFCSR`) und *Enable SSE Exceptions* (`CR4_OSXMMEXCPT`) in Kontrollregister 4 (`cr4`) setzen
- pro Thread Speicher für FPU-Zustand reservieren (→ `FPU::State`)
- beim Kontextwechsel aktuellen Zustand sichern (→ `FPU::State::save()`) und neuen laden (→ `FPU::State::restore()`)
- optional auch Unterstützung für `float` und `double` im Ausgabestrom der `libsys`

Fragen?