

Linux – Prinzipien und Programmierung

Dr. Klaus Höppner

Hochschule Darmstadt – Wintersemester 2011/2012

Linux-Kernel

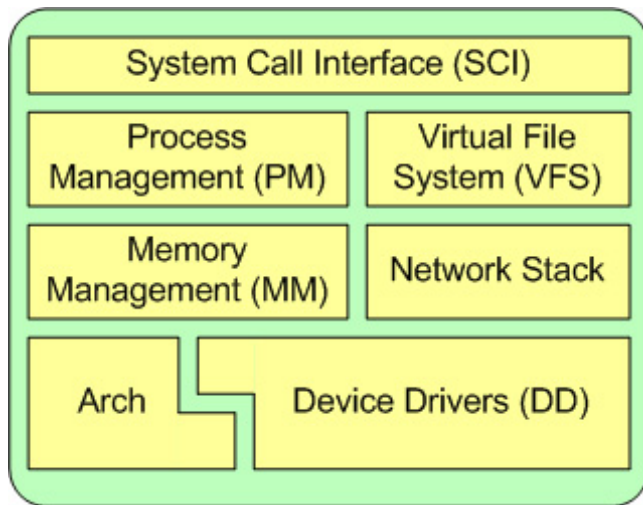
Kernel-Sources

Process-Scheduling

VFS

Memory Management

Schaubild: Komponenten des Kernels



Komponenten des Kernels

Neben dem Interface für Systemaufrufe (SCI – *System Call Interface*) besteht der Linux-Kernel aus folgenden grundlegenden Subsystemen:

Prozess-Management Ausführung von Prozessen (aus Kernelsicht werden die *Threads* genannt). Applikationen werden über das SCI Funktionen zur Verfügung gestellt, wie `fork`, `exec`, `kill`, ... Bedeutender Teil ist natürlich das Hin- und Herschalten zwischen den Prozessen, *Scheduling*.

Speicher-Management Verwaltung von *Memory Pages* (üblich 4 KB-Blöcke bei 32bit-Systemen, 8KB bei 64bit). Kritisch, da innerhalb des Kernels jedweder Speicher ungeschützt ist.

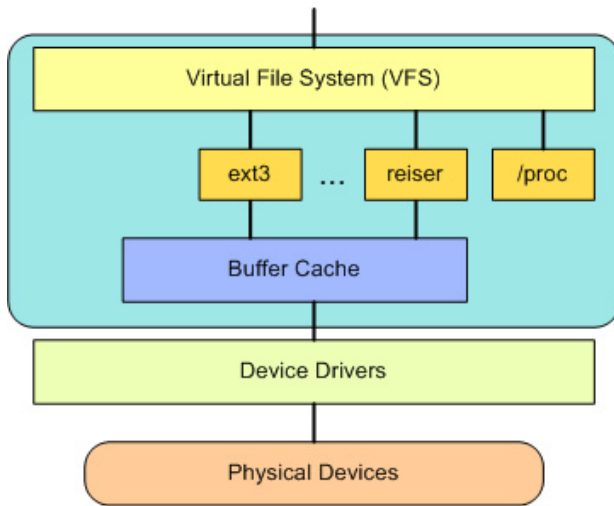
Komponenten des Kernels (Forts.)

Virtuelles Dateisystem Das VFS (*Virtual File System*) ist eine Abstraktionsschicht über den realen Dateisystemen (wie `ext2`, `ext3`, ...) und macht diese über eine allgemeine API zugänglich.

Netzwerk-Stack Der Kernel stellt die in den letzten beiden Vorlesungen beschriebenen Protokollschichten (IP – Netzwerkschicht; TCP, UDP – Transportschicht) bereit.

Gerätetreiber (*Device Drivers*) für z. B. Schnittstellen, Netzwerkkarten, Sound usw.
Architektur-spezifischer Code, z. B. für i386, PowerPC, M68k usw.

Schaubild: VFS



Herunterladen der Kernel-Sources

Zum Runterladen der Quellen stehen zwei Wege zur Verfügung:

- Archiv einer bestimmten Version runterladen, z. B.:
<http://www.kernel.org/pub/linux/kernel/v3.0/linux3.2.1.tar.bz2> (oder von einem der Mirrors, s. <http://www.kernel.org/mirrors/>)
- Sich eine lokale Kopie des Code-Repositories (git) erstellen:
`git clone`
`git://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git`
(in einer Zeile)
Letzteres ist natürlich bequemer, wenn Änderungen am Kernel jeweils nachvollzogen werden sollen.

Der Source-Tree

<code>arch</code>	Architektur-spezifische Quelldateien
<code>block</code>	Block I/O
<code>crypto</code>	Crypto API
<code>drivers</code>	Gerätetreiber
<code>firmware</code>	für bestimmte Geräte benötigt
<code>fs</code>	Dateisystem
<code>include</code>	Kernel-Headerdateien
<code>init</code>	Boot- und Initialisierungs-Code
<code>ipc</code>	Code für Interprozess-Kommunikation
<code>kernel</code>	Kern, z. B. Scheduler
<code>lib</code>	Hilfsbibliotheken
<code>mm</code>	Speichermanagement
<code>net</code>	Netzwerk
<code>samples</code>	Beispiele
<code>scripts</code>	zum Kompilieren benötigte Skripte
<code>security</code>	Sicherheit
<code>sound</code>	Audio
<code>tools</code>	Entwicklungstools
<code>usr</code>	für initramfs benutzt
<code>virt</code>	Infrastruktur für Virtualisierung

Kernel kompilieren

Bevor der Kernel kompiliert werden kann, muss er (passend) konfiguriert werden.

Hierfür stehen mehrere Wege zur Verfügung:

- `make defconfig` Voreinstellung für vorhandene CPU und
- ggfs. danach `make menuconfig` zum Anpassen der Konfiguration oder
- `make oldconfig` zum neuen Einlesen der (ggfs. manuell geänderten) Konfigurationsdatei.

Danach wird mit `make` das komprimierte Kernel-Image `arch/cpu/boot/bzImage` erzeugt.

Module werden mit `make modules_install` installiert.

Der Scheduler

Im Laufe der Entwicklung des Kernels gab es verschiedene Änderungen am *Scheduling*, die jeweils zu einer grundlegenden Neuimplementierung des *Schedulers* führten.

bis Linux 2.4 Scheduler mit $O(n)$ (bei jedem Task-Switch wurde die gesamte Prozessliste durchsucht, welcher als nächster drankommen soll).

Linux 2.5, frühes Linux 2.6 $O(1)$ -Scheduler, sortierung der Prozesse in Runqueues, Zeit zum Auswahl des nächsten Prozesses unabhängig von Gesamtzahl der Prozesse, teilweise mit heuristischen Elementen bei der Bewertung von Prozessen.

Der Scheduler (Forts.)

seit Linux 2.6.23 *Completely Fair Scheduler*, CFS Sortierung der Prozesse in einem balancierten binären Baum (red-black tree) anhand der Abweichung zwischen der real erhaltenen CPU-Zeit und der fairen CPU-Zeit, die dem Prozess aufgrund seiner Priorität zugestanden hätte.
Autor: Ingo Molnár

Prinzipielle Betrachtung des Scheduling

Wenn ein Scheduler an n Prozesse eine CPU-Zeit t_{ges} verteilen soll, so ist die faire Zeit pro Prozess

$$t_i = t_{\text{ges}} \cdot \frac{w_i}{\sum_i w_i}$$

In einem idealen Prozessor würde diese Formel für jede (beliebig kleine) Zeit t_{ges} gelten. Für einen realen Prozessor ist dies natürlich nicht realisierbar, jeder Prozesswechsel dauert eine gewisse Zeit \rightarrow bei zu kleiner Zeitscheibe (*timeslice*) CPU-Verbrauch für Scheduling statt für die eigentlichen Prozesse.

Konflikt: Timeslice groß, effizient aber interaktive Prozesse scheinen zu „hängen“ vs. Timeslice klein, also ineffizient.

Granularität: Mindestzeit zwischen zwei Taskswitchen (üblich 1 ms)

Unterschiedliches Vorgehen beim CFS-Scheduler

Der CFS-Scheduler benutzt weder Timeslices noch Runqueues, sondern pro Prozess wird (auf ns genau) die Abweichung zwischen realer CPU-Zeit und fairer Zeit protokolliert.

Nachdem der Scheduler einen Prozess zum Ausführen ausgewählt hat, wird dieser so lange ausgeführt, bis ein anderer Prozess nach der Abweichung zwischen realer und fairer Zeit „dringender“ wird.

Bei Linux ergibt sich das Gewicht eines Prozesses für die CPU-Zeitverteilung aus dem so genannten Nice-Level.

Das Virtuelle File System

Linux unterstützt eine Vielzahl von Dateisystemen:

- Physische Dateisysteme, wie `ext2`, `ext3`, `reiserfs`, `iso9660`, `(v)fat`,
- Netzwerk-Dateisysteme, insbes. `NFS`,
- Virtuelle Dateisysteme, wie das `proc`-Filesystem.

Merksatz: *Everything is a file*

Als Abstraktionsschicht stellt der Kernel das Virtuelle Dateisystem `VFS` zur Verfügung.

Wie kommen File Systems in den Kernel?

Für jedes Dateisystem gibt es unterhalb von `fs` ein Unterverzeichnis mit dem Namen des Dateisystems, also z. B. `ext2`.

Die Datei `fs/ext2/super.c` definiert nun eine Funktion, die das Dateisystem registriert:

```
int init_ext2_fs(void)
{
    return register_filesystem(&ext2_fs_type);
}
```

mit

```
static struct file_system_type ext2_fs_type = {
    ext2_read_super, "ext2", 1, NULL
};
```

Die Struktur `file_system_type`

Die Struktur `file_system_type` enthält:

- einen Zeiger auf eine Funktion zum Lesen des *Superblocks* des Dateisystems.
Der Superblock enthält Informationen wie Name des Dateisystems, einen Zeiger auf das Device (bei physikalischen FS), Blockgröße, Status, ...
- Name des Dateisystems (so wie er dann beim Befehl `mount` bei der Option `-t` angegeben wird),
- Die Angabe, ob das Dateisystem auf einem physikalischen Device aufsetzt (dann `1`), so wie sich z. B. `ext2` ein Device (z. B. `/dev/hda1`) braucht.
- Einen Nullpointer (warum?)

Durch das `register_filesystem` wird das Dateisystem einer *Linked List* von bekannten Dateisystem hinzugefügt (beim Anhängen eines neuen Eintrags in die Liste wird der NULL-Pointer des bisherigen letzten FS dann überschrieben).

Mounten eines Dateisystems

Beim Mounten eines Dateisystems laufen folgende Schritte ab:

1. In der Linked-List der bekannten Dateisysteme wird das entsprechende Dateisystem gesucht.
2. Aus der Struktur `file_system_type` wird die Funktion zum Lesen des Superblocks aufgerufen.

Soll nun in dem Dateisystem eine Datei gefunden werden, so werden die Pfadkomponenten beginnend vom ersten Inode aus rekursiv gesucht. Dabei unterstützt jeder Inode eine Menge von Inode-Operationen, die in einem Zeiger auf eine Struktur des Typs `inode_operations` beschrieben sind (und die sich z. B. zwischen Inodes, die Dateien oder Verzeichnisse bezeichnen, unterscheiden können),

Die Struktur `inode_operations`

```
struct inode_operations {
    struct file_operations * default_file_ops;
    int (*create) (struct inode *,const char *,int,int,struct inode **);
    int (*lookup) (struct inode *,const char *,int,struct inode **);
    int (*link) (struct inode *,struct inode *,const char *,int);
    int (*unlink) (struct inode *,const char *,int);
    int (*symlink) (struct inode *,const char *,int,const char *);
    int (*mkdir) (struct inode *,const char *,int,int);
    int (*rmdir) (struct inode *,const char *,int);
    int (*mknod) (struct inode *,const char *,int,int,int);
    int (*rename) (struct inode *,const char *,int,struct inode *,const char *,int);
    int (*readlink) (struct inode *,char *,int);
    int (*follow_link) (struct inode *,struct inode *,int,int,struct inode **);
    int (*readpage) (struct inode *, struct page *);
    int (*writepage) (struct inode *, struct page *);
    int (*bmap) (struct inode *,int);
    void (*truncate) (struct inode *);
    int (*permission) (struct inode *, int);
    int (*smap) (struct inode *,int);
};
```

Memory Management

Der Linux-Kernel unterteilt den Arbeitsspeicher in Blöcke (*Memory Pages*) von 4 bzw. 8 KB.

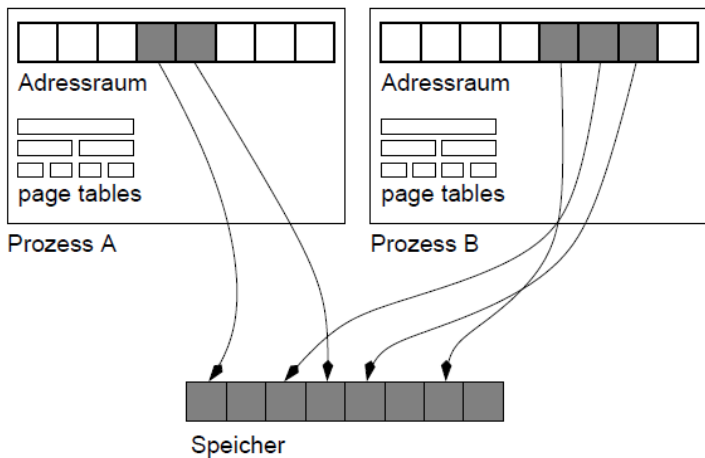
Hierbei erhält jeder Prozess (im Userspace) einen eigenen, geschützten Speicherbereich für

- Code
- Daten (Stack und Heap)

Naiver Ansatz: Jeder Prozess bekommt ein zusammenhängendes Segment im Speicher. Simpel, aber schlecht skalierbar (dynamische Speicherallokation?)

Prozesse arbeiten mit *Virtuellem Speicher*, der auf die physikalischen Pages gemappt wird.

Schaubild



Analyse

Vorteile:

- Beliebige Zuordnung von Memory Pages zu Prozessen
- Adressen prozessweit, nicht mehr global (perfekter Schutz)
- Prozesse sehen keine Speicherfragmentierung

Wichtige Begriffe

Page fault/Segmentation fault Nicht gestatteter Zugriff auf Speicher, z. B. Schreibzugriff auf RO-Section. Führt dazu, dass der Kernel dem Prozess das Signal SIGSEV schickt (Folge i. A. Core Dump)

Bus Error Zugriff auf Speicher, der physikalisch nicht zugänglich ist (Beispiel: Shared Memory über VME-Bus, wenn im adressierten Slot gar keine Karte steckt).

Copy on Write Wichtige Maßnahme des Linux Kernels beim Fork: Der virtuelle Speicher wird dupliziert, verweist aber noch auf den originalen physikalischen Speicher. Erst wenn Eltern- oder Kindprozesse in Speicherblöcken schreiben, werden diese vorher kopiert.