

1 Verwaltungsdaten F.3 UNIX — Verwaltungsstrukturen des Systemkerns pro Prozeß (2)

- ◆ der Speicherbereich der *user area* des gerade aktiven Prozesses
 - wird an eine feste Adresse des virtuellen Adreßraums des Systemkerns abgebildet
 - ↳ Zugriff auf die *user area* erfolgt in den Kern-Funktionen über die Struktur *u* (z. B. *u.u_ofile*)
 - ↳ bei einem Prozeßwechsel wird die *user area* des neuen Prozesses auf den Speicherbereich der Struktur *u* abgebildet
 - ↳ bei jedem Prozeßwechsel muß die Abbildung virtuelle Adressen → physikalische Adressen geändert werden
 - oder wird über einen festen Zeiger im Systemkern adressiert
 - ↳ Zugriff auf die *user area* erfolgt in den Kern-Funktionen über den Zeiger *uunix* (z. B. *uunix->u_ofile*)
 - ↳ bei einem Prozeßwechsel wird der Zeiger *uunix* auf die neue *user area* gelegt
 - ↳ jeder Zugriff auf Elemente der *u-area* erfolgt indirekt

AKBP

Ausgewählte Kapitel der praktischen Betriebsprogrammierung
© Jürgen Kleinöder, Universität Erlangen-Nürnberg, IMMD IV, 1998

F-Prozesse.doc 1998-11-25 08.32

F.13

Reproduktion jeder Art oder Verwendung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

2 Globale F.3 UNIX — Verwaltungsstrukturen des Systemkerns Prozeßverwaltungsdaten (2)

- ... Prozeßtabelle
 - ◆ eine *proc*-Struktur enthält u. a. folgende Daten:
 - Prozeßzustand (laufend/bereit, blockiert, ausgelagert, ...)
 - Verweise auf die Segmente des Prozesses (meist über *page table*)
 - Verweis auf die *user area*
 - Verweis auf die *credential*-Struktur mit den Zugriffsrechten (*User-Ids*, Gruppen-Rechte)
 - Speicher für Daten zur Signalverarbeitung
 - Scheduling-Parameter (z. B. Priorität)
 - Information über Ereignis, auf das der Prozeß wartet
 - ◆ das Kommando *ps(1)*, gibt Auszüge der Prozeßtabelle aus

AKBP

Ausgewählte Kapitel der praktischen Betriebsprogrammierung
© Jürgen Kleinöder, Universität Erlangen-Nürnberg, IMMD IV, 1998

F-Prozesse.doc 1998-11-25 08.32

F.15

Reproduktion jeder Art oder Verwendung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

1 Verwaltungsdaten F.3 UNIX — Verwaltungsstrukturen des Systemkerns pro Prozeß (3)

- Kernel stack
 - ◆ Stack für Funktionen des Systemkerns, die im Kontext des Prozesses ablaufen (zur Abwicklung von Systemaufrufen)
 - ◆ bei Prozeßwechsel wird der Kernel Stack des neuen Prozesses entsprechend in den virt. Adreßraum des Kerns gemapped

2 Globale Prozeßverwaltungsdaten

- Prozeßtabelle
 - ◆ alle wichtigen Verwaltungsdaten über einen Prozeß, die immer verfügbar sein müssen, werden in *proc*-Strukturen gehalten
 - ◆ die *proc*-Strukturen bilden die **Prozeßtabelle** (bei Booten statisch als Feld angelegt oder dynamisch als Liste verwaltet)

AKBP

Ausgewählte Kapitel der praktischen Betriebsprogrammierung
© Jürgen Kleinöder, Universität Erlangen-Nürnberg, IMMD IV, 1998

F-Prozesse.doc 1998-11-25 08.32

F.14

Reproduktion jeder Art oder Verwendung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

2 Globale F.3 UNIX — Verwaltungsstrukturen des Systemkerns Prozeßverwaltungsdaten (3)

- Prozeßidentifikatoren und Prozeßgruppen
 - pid* **process id** - eindeutiger Prozeßidentifikator ($0 \leq pid \leq 30000$) **getpid(2)**
 - ppid* **parent process id** = *pid* des "Vaterprozesses" **getppid(2)**
 - pgrp* **process group id** = *pid* des Prozeßgruppenführers **setpgrp(2), getpgrp(2), setpgid(2)** (POSIX) (siehe auch Abschnitt "Sessions und Prozeßgruppen")
 - real user id (ruid)** und **real group id (rgid)**
uid und *gid* des Benutzers, der den Prozeß erzeugt hat (festgelegt durch *login* des Benutzers)
setuid(), getuid(), setgid(), getgid()

AKBP

Ausgewählte Kapitel der praktischen Betriebsprogrammierung
© Jürgen Kleinöder, Universität Erlangen-Nürnberg, IMMD IV, 1998

F-Prozesse.doc 1998-11-25 08.32

F.16

Reproduktion jeder Art oder Verwendung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

3 Kommando ps F.3 UNIX — Verwaltungsstrukturen des Systemkerns

F	S	UID	PID	PPID	C	PRI	NI	ADDR	SZ	WCHAN	TTY	TIME	CMD
19	T	0	0	0	0	0	SY	10416288	0	?	?	0:00	sched
8	S	0	1	0	0	41	20	604df608	85	604df600	?	0:03	init
19	S	0	2	0	0	0	SY	604def48	0	10432c34	?	0:00	pageout
19	S	0	3	0	0	0	SY	604de888	0	10435cdc	?	34:37	fsflush
8	S	0	239	1	0	51	20	604dd448	253	604dd440	?	0:11	sendmail
8	S	10402	3705	1	1	51	20	608c3620	4808	6098490e	console	28:16	makecs5k.
8	S	0	122	1	0	40	20	60693610	257	600a8b8e	?	0:15	rpcbnd
8	S	0	124	1	0	94	20	60692f50	249	600a8aee	?	0:00	keyserv
8	S	0	135	1	0	41	20	60691b10	231	600a8bb6	?	0:00	ypbind
8	S	0	179	1	0	41	20	60691450	479	609848e6	?	8:06	autocmount
8	S	0	238	1	0	40	8	606906d0	230	606908c8	?	1:01	xmtpd
8	S	0	201	1	0	51	20	60759b18	206	6002deb8	?	0:02	cron
8	S	10402	1776	1	0	50	20	6098a028	396	60985ade	console	0:01	xterm
8	S	0	269	1	0	41	20	60758018	210	604fe060	?	0:12	ashd
8	S	10402	3579	1	0	58	20	608c1b20	159	608c1d18	console	0:00	csd
8	S	0	313	1	0	40	20	608c1460	726	600a8206	?	0:00	dtlogin
8	S	10402	3638	3629	0	40	20	60aac0b8	180	60a690be	pts/3	0:01	elm
8	S	10402	3597	3579	0	77	20	608c21e0	121	608c2250	console	0:00	openwin
8	O	0	8150	3633	0	41	20	609ec1e8	114	pts/2	?	0:00	ps.exe
8	S	10402	3601	3597	0	51	20	60aaadb8	321	60aaaae28	console	0:00	init
8	S	10402	3693	3668	0	51	20	60a921f0	1963	609853fe	console	0:00	netscape
8	S	10402	1777	1776	0	40	20	6098a6e8	188	60ade4c6	pts/4	0:00	rlogin
8	S	10402	3613	3603	0	40	20	60a90db0	381	600a813e	console	0:45	vtvwm
8	S	10402	27875	1	0	51	20	60758d98	399	60985ba6	console	0:18	xterm
8	S	10402	3603	3601	0	49	20	60aaab78	121	60aaab4e8	console	0:00	xinitrc
8	S	10402	3669	1	0	50	20	608c2f60	708	6098585e	console	0:04	cm
8	Z	10402	3610	3603	0	0	0	0	0	0	0	0:00	<defunct>
8	K	10402	3602	3601	2	43	20	60a92f70	23544	console	63:44	Xsun	
8	S	10402	3633	3622	0	41	20	6098ecf8	276	609ed160	pts/2	0:02	tcsh-6.0
8	S	10402	3623	3618	0	40	20	608c28a0	258	60a689be	pts/0	0:01	tcsh-6.0
8	S	10402	3668	1	2	51	20	60b5edc0	4165	600a809e	console	125:56	netscape
8	S	10402	3622	1	0	40	20	60a928b0	401	600a804e	console	0:03	xterm
8	S	10402	3680	1	0	40	20	60b55b40	373	600a818e	console	0:00	xclipboa
8	S	10402	27876	27875	0	40	20	60b60200	286	60b905de	pts/1	3:15	elm
8	T	10402	3799	3633	0	52	20	60b5e700	188	pts/2	?	0:00	rlogin
8	Z	10402	3614	3613	0	0	0	0	0	0	0	0:00	<defunct>
8	S	10402	3019	1	0	50	20	60a91b30	397	6098463e	console	0:01	xterm

Reproduktion jeder Art oder Vervielfältigung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

1 Überblick (2) F.4 UNIX - Erzeugen eines neuen Prozesses

- wesentliche Unterschiede zwischen Vater- und Sohnprozess:
- fork-Ergebnis:** Vaterprozess: *pid* des Sohnes, Sohnprozess: 0
- PID:** Sohnprozess erhält nächste freie *pid*
- Eintrag in Prozeßtafel:** Sohnprozess erhält eigene proc-Struktur, Inhalt wird weitgehend vom Vaterprozess kopiert

Reproduktion jeder Art oder Vervielfältigung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

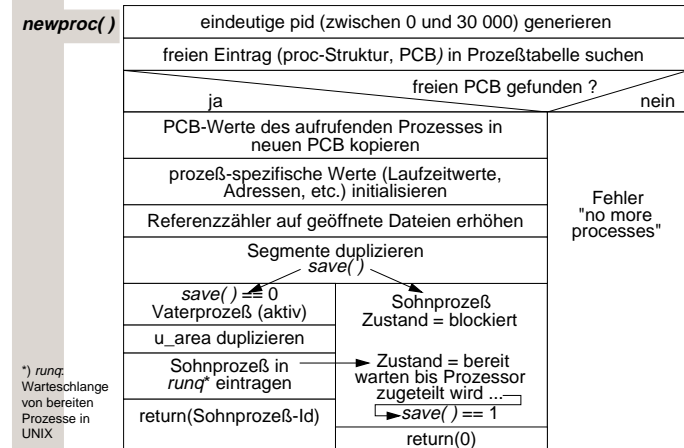
F.4 UNIX - Erzeugen eines neuen Prozesses

1 Überblick

- jeder UNIX-Prozess kann mit dem Systemaufruf **fork(2)** einen neuen Prozess erzeugen
- fork()** erzeugt eine nahezu identische Kopie des aufrufenden Prozesses
- der **fork()** aufrufende Prozess wird **Vaterprozess (parent process)** genannt
- der durch **fork()** neu erzeugte Prozess wird als **Sohnprozess (child process)** bezeichnet
- der Sohnprozess erbt alle Rechte und Einschränkungen vom Vaterprozess (einschließlich aller Daten der *user area*, z. B. offene Dateien)

Reproduktion jeder Art oder Vervielfältigung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

2 Ablauf F.4 UNIX - Erzeugen eines neuen Prozesses



*) runq: Warteschlange von bereiten Prozessen in UNIX

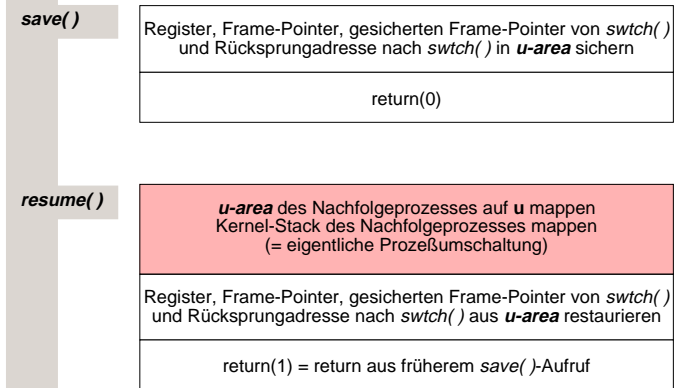
Reproduktion jeder Art oder Vervielfältigung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

F.5 UNIX — Prozeßumschaltungen

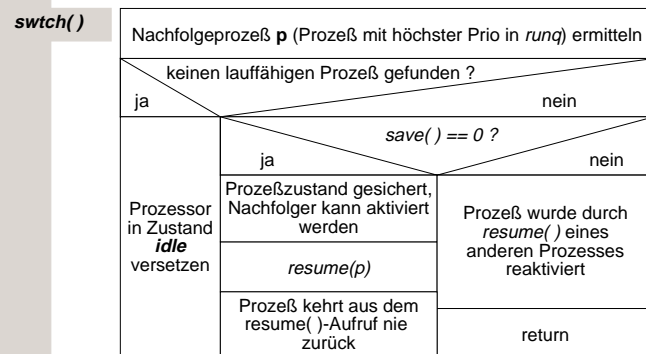
1 Überblick

- Prozeßumschaltungen erfolgen nur an wenigen Stellen im Systemkern:
 - ↳ wenn sich ein Prozeß blockiert (*sleep()*)
 - Prozeß aktiviert höchstprioren, lauffähigen Prozeß
 - wenn das Ereignis, auf das gewartet wird, eintritt, wird der Prozeß durch *wakeup()* in die Warteschlange der bereiten Prozesse (*runq*) eingetragen und kann reaktiviert werden
 - ↳ vor Rückkehr aus einem Systemaufruf, wenn ein höherpriorer Prozeß bereitsteht
 - Prozeß trägt sich in die *runq* ein
 - Prozeß aktiviert höchstprioren, lauffähigen Prozeß
- die Prozeßumschaltung erfolgt in den Funktionen *swtch()*, *save()* und *resume()* des Systemkerns (stark abhängig von Hardware und UNIX-Version!)

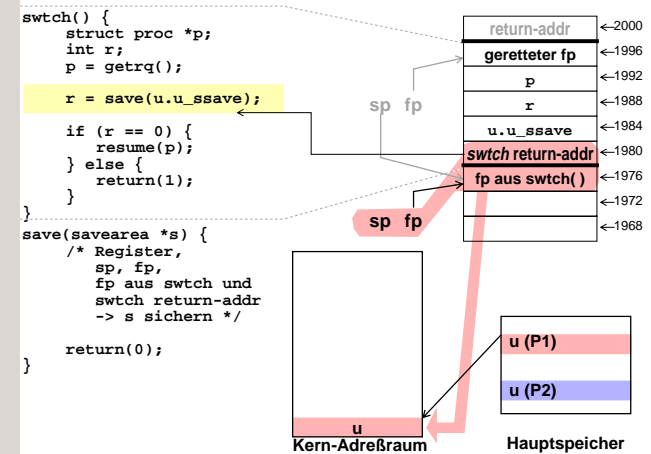
3 Ablauf (2)



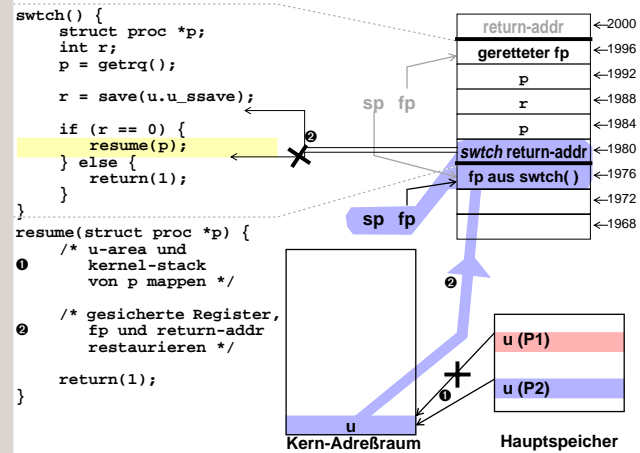
2 Ablauf (Beispiel aus XELOS)



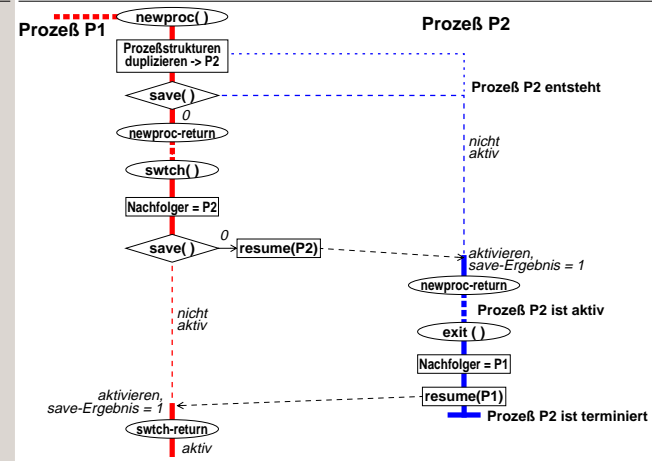
3 Ablauf (3)



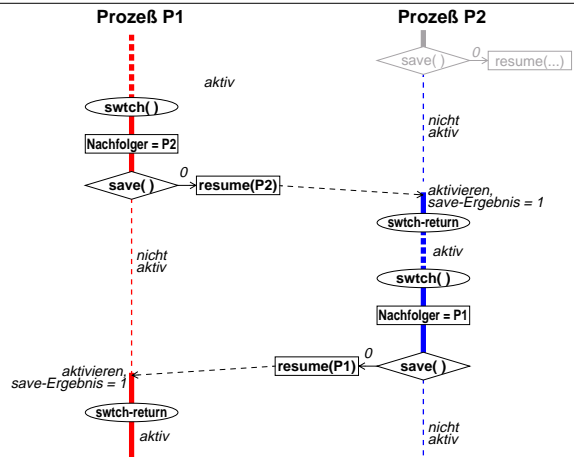
3 Ablauf (4)



5 Erzeugen eines neuen Prozesses



4 Wechsel zwischen zwei Prozessen



F.6 UNIX — Ausführen eines Programms

- ein UNIX-Prozess kann die Ausführung eines Programms durch ein anderes Programm ersetzen — **exec(2), execve(2)**
 - durch Laden eines neuen Programms werden die zuvor von dem Prozess bearbeiteten Programm-Datenstrukturen zerstört
 - im Erfolgsfall gibt es kein return aus exec()
 - durch Laden eines neuen Programms entsteht kein neuer Prozess
- execve() werden die Argumente (argv) und ein Environment (envp) für das neue Programm mitgegeben
- Details:
 - programmspezifische Daten des Prozesses (Segmente, signal-handler-Funktionen) werden initialisiert
 - effective uid / effective gid des Prozesses werden auf Eigentümer/Gruppe der Programm-Datei geändert, wenn bei dieser ein s-bit gesetzt ist

F.6 UNIX — Ausführen eines Programms (2)

■ ... Details:

- ◆ alle anderen Prozeßparameter (in *proc*-Struktur und *user area*) bleiben unverändert (*pid*, *ppid*, *pgrp*, *session*, *controlling tty*, *current working directory*, *current root directory*, ...)
- ◆ **offene Dateideskriptoren bleiben erhalten**, es sei denn, sie wurden als *close-on-exec* (siehe ***fcntl(2)***) markiert

■ Programmier-Beispiel:

```
if ( (pid = fork()) < 0 ) {
    perror("fork");
    exit(1)
} else if (pid == 0) {           /* child process */
    execl("/bin/cp", "cp",    "/tmp/a", "/tmp/b", (char
*)0);
    perror("exec");
    exit(1);
} else {                       /* parent process */
    ...
}
```

AKBP |

Ausgewählte Kapitel der praktischen Betriebsprogrammierung
© Jürgen Kleinöder, Universität Erlangen-Nürnberg, IMMD IV, 1998

F-Prozesse.doc 1998-11-25 08.32

F.33

Reproduktion jeder Art oder Verwendung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

2 Ablauf von exit im Kern

F.7 UNIX — Terminieren von Prozessen

- ◆ System V:
wenn der Prozeß ein Prozeßgruppenführer mit zugeordnetem Terminal war (typischerweise *login-shells*), wird ein *SIGHUP*-Signal an alle Prozesse dieser Prozeßgruppe gesendet
- ◆ POSIX:
wenn der Prozeß ein *controlling process* (*session leader* mit zugeordnetem *controlling tty*) war, wird allen Prozessen in der Vordergrund-Prozeßgruppe ein *SIGHUP*-Signal gesendet, das *controlling tty* wird von der *session* abgekoppelt
- ◆ alle offenen Dateideskriptoren werden geschlossen
- ◆ alle Segmente (Text, Daten, Stack) werden freigegeben; *shared memory* wird abgekoppelt, aber nicht gelöscht
- ◆ *user area* (enthält Kernel-Stack auf dem *exit* operiert!) und *proc*-Struktur bleiben noch erhalten
- ◆ alle noch existierenden Sohnprozesse erhalten den *init*-Prozeß (*pid* = 1) als neuen Vaterprozeß
- ◆ Statistikdaten und *exit*-Status werden in der *proc*-Struktur notiert

AKBP |

Ausgewählte Kapitel der praktischen Betriebsprogrammierung
© Jürgen Kleinöder, Universität Erlangen-Nürnberg, IMMD IV, 1998

F-Prozesse.doc 1998-11-25 08.32

F.35

Reproduktion jeder Art oder Verwendung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

F.7 UNIX — Terminieren von Prozessen

1 Überblick

- Prozesse terminieren, wenn
 - ◆ sie den Systemdienst ***exit(2)*** aufrufen
 - ◆ ein Signal an den Prozeß zugestellt wurde, für das keine Signalbearbeitungsfunktion vorgesehen ist
- die Startumgebung für C-Programme ruft nach einem *return* aus der Funktion *main* automatisch *exit(0)* auf
- dem *exit*-Aufruf kann ein Status-Wert (1 Byte) mitgegeben werden, der durch den Vaterprozeß abgefragt werden kann (→ *wait(2)*)
- ein Vaterprozeß kann auf das terminieren von Sohnprozessen warten und deren *exit*-Status abfragen — ***wait(2)***

AKBP |

Ausgewählte Kapitel der praktischen Betriebsprogrammierung
© Jürgen Kleinöder, Universität Erlangen-Nürnberg, IMMD IV, 1998

F-Prozesse.doc 1998-11-25 08.32

F.34

Reproduktion jeder Art oder Verwendung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

2 Ablauf von exit im Kern (2) F.7 UNIX — Terminieren von Prozessen

- ◆ neuer Prozeßzustand: **ZOMBIE**
- ◆ Vaterprozeß wird informiert (*wakeup()*), falls in *wait()* blockiert, und Signal *SIGCLD*/*SIGCHLD*)
- ◆ Nachfolgeprozeß bestimmen und Prozeßwechsel
- ZOMBIE-Prozesse bleiben so lange bestehen, bis der Vaterprozeß mit dem Systemaufruf ***wait(2)*** den *exit*-Status abfragt
 - ◆ mit *wait()* werden Statistikdaten an den Vaterprozeß übertragen und können danach mit ***times(2)*** abgefragt werden
 - ***times(2)*** liefert die Summe der Laufzeiten aller Sohnprozesse, auf die mit ***wait(2)*** gewartet wurde
 - ◆ alternativ können die Statistikdaten des Sohnprozesses mit ***wait3(2)*** oder ***wait4(2)*** direkt abgefragt werden
 - ◆ der Eintrag in Prozeßtablelle wird freigegeben
 - ◆ die *user area* des ZOMBIE-Prozesses ist jetzt passiv (Kernel-Stack wird nicht mehr genutzt) und kann freigegeben werden

AKBP |

Ausgewählte Kapitel der praktischen Betriebsprogrammierung
© Jürgen Kleinöder, Universität Erlangen-Nürnberg, IMMD IV, 1998

F-Prozesse.doc 1998-11-25 08.32

F.36

Reproduktion jeder Art oder Verwendung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

3 Systemaufruf `wait()` F.7 UNIX — Terminieren von Prozessen

- ein Prozeß kann warten, bis ein Sohnprozeß terminiert oder gestoppt wird und dabei den Status des Sohnprozesses abfragen
 - übliche Arbeitsweise einer Shell bei Vordergrundprozessen
- **`wait(2)`** blockiert den aufrufenden Prozeß so lange, bis ein Sohnprozeß im Zustand ZOMBIE existiert oder ein Sohnprozeß gestoppt wird
 - ◆ `pid` dieses Sohnprozesses wird als Ergebnis geliefert
 - ◆ als Parameter kann ein Zeiger auf einen `int`-Wert mitgegeben werden, in dem der Status (16 Bit) des Sohnprozesses abgelegt wird

Sohnprozeß ist	Status & 0xff00	Status & 0x0fff
gestoppt	Signalnummer	0xff
terminiert durch <code>exit()</code>	exit-code	0x00
terminiert durch Signal	0x00	Signalnummer
terminiert durch Signal + <code>coredump</code>	0x00	Signalnummer + 0x80

4 typische `wait()`-Anwendung F.7 UNIX — Terminieren von Prozessen

- ein Prozeß führt eine Liste aller existierenden Sohnprozesse und interessiert sich für deren `exit`-Status, darf selbst aber nicht durch `wait()` blockiert werden
 - Lösung: Signal SIGCLD (SystemV) bzw. SIGCHLD (Berkeley) auswerten
 - ◆ großer Unterschied zwischen UNIX SystemV und Berkeley-UNIX!
- ▲ Verfahren bei Berkeley-UNIX
 - ◆ Prozeß setzt `signalhandler`-Funktion für Signal SIGCHLD
 - ◆ in der `signalhandler`-Funktion wird so lange `wait3()` mit Option `WNOHANG` aufgerufen, bis `wait3()` meldet, daß keine weiteren Söhne in gestopptem oder terminiertem Zustand existieren

3 Systemaufruf `wait()` (2) F.7 UNIX — Terminieren von Prozessen

- existiert kein Sohnprozeß, liefert `wait()` eine Fehlermeldung
- in BSD und POSIX gibt es zusätzlich die Möglichkeit, das Blockieren beim `wait`-Aufruf zu verhindern, wenn keine terminierten oder gestoppten Sohnprozesse existieren
 - ◆ `wait3(2)`, `waitpid(2)` oder `wait4(2)` mit Option `WNOHANG`
- in POSIX kann man mit `wait4()` oder `waitpid()` die Menge der Prozesse, auf die gewartet wird einschränken

4 typische `wait()`-Anwendung F.7 UNIX — Terminieren von Prozessen

- ▲ Verfahren bei SystemV
 - ◆ Prozeß setzt `signalhandler`-Funktion für Signal SIGCLD
 - ◆ in der `signalhandler`-Funktion wird **genau ein** `wait()`-Aufruf ausgeführt, durch den der Prozeß nicht blockiert werden kann
 - ◆ da die `signalhandler`-Funktion bei Eintreffen eines Signals zurückgesetzt wird, muß sie am Ende der Signalbearbeitung neu aufgesetzt werden
 - ◆ beim Setzen einer `Handler`-Funktion für SIGCLD stellt der Systemkern ein neues SIGCLD zu, wenn noch Sohnprozesse im ZOMBIE-Zustand existieren
- weitere Besonderheiten von SIGCLD in SystemV
 - ◆ SIGCLD führt (im Gegensatz zu anderen Signalen) per default nicht zum Abbruch eines Prozesses - es wird einfach ignoriert
 - ◆ wird für SIGCLD explizit vereinbart, daß es ignoriert werden soll, werden Sohnprozesse automatisch - ohne `wait()`-Aufruf - vom Kern aus dem ZOMBIE-Zustand befreit