

X . s y s t e m s . p r e s s

X.systems.press ist eine praxisorientierte Reihe zur Entwicklung und Administration von Betriebssystemen, Netzwerken und Datenbanken.

Markus Zahn

Unix- Netzwerkprogrammierung mit Threads, Sockets und SSL

Mit 44 Abbildungen und 19 Tabellen

 Springer

Markus Zahn
unp@bit-oase.de
<http://unp.bit-oase.de>

Bibliografische Information der Deutschen Bibliothek
Die Deutsche Bibliothek verzeichnet diese Publikation in der Deutschen
Nationalbibliografie; detaillierte bibliografische Daten sind im Internet über
<http://dnb.ddb.de> abrufbar.

ISSN 1611-8618
ISBN-10 3-540-00299-5 Springer Berlin Heidelberg New York
ISBN-13 978-3-540-00299-4 Springer Berlin Heidelberg New York

Dieses Werk ist urheberrechtlich geschützt. Die dadurch begründeten Rechte, insbesondere die der Übersetzung, des Nachdrucks, des Vortrags, der Entnahme von Abbildungen und Tabellen, der Funksendung, der Mikroverfilmung oder der Vervielfältigung auf anderen Wegen und der Speicherung in Datenverarbeitungsanlagen, bleiben, auch bei nur auszugsweiser Verwertung, vorbehalten. Eine Vervielfältigung dieses Werkes oder von Teilen dieses Werkes ist auch im Einzelfall nur in den Grenzen der gesetzlichen Bestimmungen des Urheberrechtsgesetzes der Bundesrepublik Deutschland vom 9. September 1965 in der jeweils geltenden Fassung zulässig. Sie ist grundsätzlich vergütungspflichtig. Zuwiderhandlungen unterliegen den Strafbestimmungen des Urheberrechtsgesetzes.

Springer ist ein Unternehmen von Springer Science+Business Media
springer.de

© Springer-Verlag Berlin Heidelberg 2006

Die Wiedergabe von Gebrauchsnamen, Handelsnamen, Warenbezeichnungen usw. in diesem Werk berechtigt auch ohne besondere Kennzeichnung nicht zu der Annahme, dass solche Namen im Sinne der Warenzeichen- und Markenschutz-Gesetzgebung als frei zu betrachten wären und daher von jedermann benutzt werden dürften. Text und Abbildungen wurden mit größter Sorgfalt erarbeitet. Verlag und Autor können jedoch für eventuell verbliebene fehlerhafte Angaben und deren Folgen weder eine juristische Verantwortung noch irgendeine Haftung übernehmen.

Satz: Druckfertige Daten des Autors
Herstellung: LE-TeX, Jelonek, Schmidt & Vöckler GbR, Leipzig
Umschlaggestaltung: KünkelLopka Werbeagentur, Heidelberg
Gedruckt auf säurefreiem Papier 33/3100 YL - 5 4 3 2 1 0

Vorwort

Vernetzte Rechnersysteme und insbesondere das weltumspannende *Internet* haben unsere Welt verändert. Mit Hilfe der dabei entstandenen Technologien ist es heute nicht nur möglich, sondern sogar äußerst einfach, mit dem eigenen PC selbst Teil dieses riesigen Computernetzwerks zu werden. Natürlich ist allein ein Verbund vernetzter Rechnersysteme für den Normalverbraucher noch nicht sonderlich interessant. Erst die Fülle von Anwendungen, von *Online-Enzyklopädien* über *Online-Banking* und *Online-Shopping* bis hin zu *File-Sharing* und *Online-Spielen*, die seit den Anfängen des Internets entstanden sind, gestaltet dieses Netz anziehend für seine Nutzer.

Die Anziehungskraft vernetzter Rechnersysteme steht und fällt also mit der Attraktivität und Zuverlässigkeit der darüber verfügbaren Anwendungen. Das vorliegende Buch beschäftigt sich deshalb mit der Programmierung vernetzter Computersysteme, genauer gesagt mit der Entwicklung netzwerkfähiger Client-/Server-Programme für Unix-Systeme (oder Unix-ähnliche Computersysteme). Es hat zum Ziel, dem Leser einen fundierten Einstieg in die Welt der Unix-Netzwerkprogrammierung zu vermitteln, klammert aber auch fortgeschrittene Themen nicht aus. Die notwendigen Grundlagen der Unix-Systemprogrammierung werden demnach ebenso berücksichtigt wie die Absicherung des Datenverkehrs mittels *SSL (Secure Socket Layer)*. Zahlreiche Programmbeispiele mit typischen Implementierungsmustern stellen dem Leser darüber hinaus eine solide Codebasis für die Entwicklung zuverlässiger, leistungsfähiger und sicherer Netzwerkprogramme zur Verfügung.

Die Einschränkung auf Unix und Unix-ähnliche Systeme geht auf die gemeinsame Entwicklungsgeschichte des Unix-Betriebssystems und des Internets zurück. So fand z. B. in den 70'er Jahren die Implementierung von TCP/IP und der zugehörigen Socket-API zunächst auf Unix-Systemen statt. Unix-Systeme bilden aufgrund ihrer hohen Betriebsstabilität sowie ihrer seit langem etablierten Multiuser-, Multiprozeß- und Multithreading-Fähigkeiten auch heute die Plattform für die wichtigsten Netzwerkdienste im Internet. Selbst auf der Seite der Arbeitsplatzsysteme gewinnt momentan mit Linux wieder

ein Unix-ähnliches Betriebssystem mehr und mehr an Bedeutung. Nachdem andere Betriebssysteme ebenfalls die in diesem Buch beschriebene Socket-API für die Netzwerkprogrammierung adaptiert haben und auch das vorgestellte OpenSSL auf vielen Systemen verfügbar ist, sind die hier diskutierten Konzepte der Netzwerkprogrammierung mitsamt der Beispiele meist ohne größere Probleme auf nicht Unix-basierte Betriebssysteme übertragbar.

Sowohl durch meine eigenen praktischen Erfahrungen als auch durch meine Tätigkeit als Dozent mußte ich lernen, daß es für die Entwicklung effizienter und stabiler Netzwerkanwendungen keinesfalls ausreichend ist, allein die Grundlagen der klassischen Netzwerkprogrammierung zu beherrschen. Zum einen entstammen viele Problemstellungen weniger der Netzwerkprogrammierung als vielmehr der Unix-Systemprogrammierung. Auch typische Fehler beruhen oftmals auf grundlegenden Mißverständnissen (oder Informationsdefiziten) aus dem Bereich der Systemprogrammierung. Zum anderen steigt aber auch durch die wachsende Nutzung des Internets für private und geschäftliche Transaktionen wie Online-Banking und Online-Shopping der Bedarf an einer gesicherten Übertragung der transportierten Daten.

Insofern gliedert sich das Buch in drei Teile: der erste Teil widmet sich der Unix-Systemprogrammierung, der zweite Teil beschäftigt sich mit der klassischen Netzwerkprogrammierung und der dritte Teil beleuchtet die Absicherung des Datenverkehrs.

Nach einer allgemeinen Einführung in Kapitel 1 bereitet Kapitel 2, *Programmieren mit Unix-Prozessen*, die wichtigsten Unix-Grundlagen auf. Das Kapitel konzentriert sich ausschließlich auf die Ausschnitte der Systemprogrammierung, auf die wir später im Rahmen der Netzwerkprogrammierung zurückgreifen werden. Im Wesentlichen sind dies die Themen Ein-/Ausgabe, Signalbehandlung und Nebenläufigkeit auf Prozeßebene. Kapitel 3, *Programmieren mit POSIX-Threads*, zeigt, wie nebenläufige Handlungen auch innerhalb eines Prozesses implementiert werden können.

Die Kapitel 4, *Grundlagen der Socket-Programmierung*, und 5, *Netzwerkprogrammierung in der Praxis*, beschäftigen sich dann mit der Netzwerkprogrammierung auf Basis der Socket-API und besprechen fünf verschiedene Implementierungsmuster für typische Server-Programme.

Die letzten beiden Kapitel widmen sich schließlich der Absicherung des Datenverkehrs über das SSL-Protokoll. Nachdem in Kapitel 6, *Netzwerkprogrammierung mit SSL*, sowohl die SSL-Grundlagen als auch die Basisfunktionalität der freien SSL-Implementierung OpenSSL erläutert wurden, widmet sich Kapitel 7, *Client-/Server-Programmierung mit OpenSSL*, v. a. dem sachgemäßen, sicheren Umgang mit SSL-Zertifikaten und damit der Entwicklung sicherer SSL-fähiger Anwendungen.

Bei der Erstellung des Manuskripts sind mir zahlreiche Helfer mit ihren kritischen Fragen, inhaltlichen Anregungen und durch ihr fleißiges Korrekturlesen zur Seite gestanden. Mein Dank gilt den Freunden und Kollegen Dr. Lars

Freund, Thomas Morper, Hella Seebach, Dr. Michael Westenburg und ganz besonders Dr. Harald Görl für zahlreiche fruchtbare Diskussionen und viele hilfreiche Anregungen. Dem Springer-Verlag danke ich in Person von Frau Jutta Maria Fleschutz und Herrn Frank Schmidt für die angenehme, problemlose und notwendigerweise auch geduldige Zusammenarbeit.

Sämtliche Beispielprogramme aus dem vorliegenden Buch stehen unter der Adresse <http://unp.bit-oase.de/> zum Download bereit. Konstruktive Kritik und Hinweise auf Fehler, die sich trotz sorgsamer Ausarbeitung des Manuskripts eingeschlichen haben könnten, nehme ich jederzeit gerne entgegen.

Augsburg, im Juni 2006

Dr. Markus Zahn
unp@bit-oase.de

Inhaltsverzeichnis

1	Einführung	1
1.1	TCP/IP-Grundlagen	2
1.1.1	Netzwerkschicht	3
1.1.2	Internet-Schicht	4
1.1.3	Transportschicht	5
1.1.4	Anwendungsschicht	6
1.2	Internet-Standards	6
1.3	Unix-Standards	7
2	Programmieren mit Unix-Prozessen	9
2.1	Unix-Prozesse	9
2.1.1	Prozeßgruppen und Sessions	10
2.1.2	Kontrollierendes Terminal	12
2.1.3	Verwaiste Prozesse und verwaiste Prozeßgruppen	13
2.1.4	Prozeßumgebung	14
2.1.5	Lebenszyklus	15
2.1.6	User- und Gruppen-ID	20
2.2	Ein- und Ausgabe	21
2.2.1	Dateideskriptoren	21
2.2.2	Elementare Ein- und Ausgabe	24
2.2.3	Standerdeingabe und -ausgabe	35
2.2.4	Ausgabe über den Syslog-Dienst	41
2.3	Buffer-Overflows und Format-String-Schwachstellen	45
2.3.1	Buffer-Overflows	47
2.3.2	Format-String-Schwachstellen	54
2.3.3	Geeignete Gegenmaßnahmen	57
2.4	Signale	58
2.4.1	Signale behandeln	59
2.4.2	Signale blockieren	68
2.4.3	Signale annehmen	72
2.4.4	Signale generieren	74

2.5	Prozeßkontrolle	76
2.5.1	Was bin ich? Prozeß-IDs und mehr	77
2.5.2	Neue Prozesse erzeugen	79
2.5.3	Prozesse synchronisieren	83
2.5.4	Zombie-Prozesse	89
2.5.5	Andere Programme ausführen	90
2.5.6	User- und Group-IDs wechseln	94
2.6	Dæmon-Prozesse	96
3	Programmieren mit POSIX-Threads	103
3.1	Grundlagen	104
3.2	Synchronisation	115
3.2.1	Race Conditions und kritische Bereiche	116
3.2.2	Gegenseitiger Ausschluß	120
3.2.3	Bedingungsvariablen	126
3.3	Pthreads und Unix-Prozesse	135
3.3.1	Threadsichere und eintrittsinvariante Funktionen	135
3.3.2	Fehlerbehandlung und errno	137
3.3.3	Signalverarbeitung	138
3.3.4	fork() und exec() in Pthreads-Programmen	144
4	Grundlagen der Socket-Programmierung	147
4.1	Erste Schritte mit telnet und inetd	147
4.1.1	Das telnet-Kommando als Netzwerk-Client	147
4.1.2	Einfache Netzwerkdienste mit dem inetd	152
4.2	IP-Namen und IP-Adressen	156
4.2.1	Das Domain Name System	157
4.2.2	IPv4-Adressen	159
4.2.3	IPv6-Adressen	164
4.2.4	Netzwerkdarstellung von IP-Adressen	168
4.3	Sockets	179
4.3.1	Socket anlegen	180
4.3.2	Socket-Strukturen	182
4.3.3	Client-seitiger TCP-Verbindungsaufbau	184
4.3.4	Socket-Adressen zuweisen	189
4.3.5	Annehmende Sockets	192
4.3.6	TCP-Verbindungen annehmen	194
4.3.7	Drei-Wege-Handshake und TCP-Zustandsübergänge ..	199
4.3.8	Kommunikation über UDP	205
4.3.9	Standardeingabe und -ausgabe über Sockets	212
4.3.10	Socket-Adressen ermitteln	213
4.3.11	Multiplexing von Netzwerkverbindungen	218
4.3.12	Socket-Optionen	223
4.4	Namensauflösung	227

5	Netzwerkprogrammierung in der Praxis	235
5.1	Aufbau der Testumgebung	236
5.1.1	Funktionsumfang der Testumgebung	237
5.1.2	Hilfsfunktionen für die Socket-Kommunikation	238
5.1.3	Der Test-Client	249
5.2	Iterative Server	255
5.2.1	Sequentielle Verarbeitung der Anfragen	256
5.2.2	Clientbehandlung	259
5.2.3	Hilfsfunktionen zur Laufzeitmessung	262
5.2.4	Eigenschaften und Einsatzgebiete	264
5.3	Nebenläufige Server mit mehreren Threads	266
5.3.1	Abgewandelte Signalbehandlung	267
5.3.2	Ein neuer Thread pro Client	268
5.3.3	Das Hauptprogramm als Signalverarbeiter	270
5.3.4	Eigenschaften und Einsatzgebiete	272
5.4	Nebenläufige Server mit Prethreading	274
5.4.1	Clientbehandlung mittels paralleler Accept-Handler	275
5.4.2	Das Hauptprogramm als Signalverarbeiter	277
5.4.3	Eigenschaften und Einsatzgebiete	279
5.5	Nebenläufige Server mit mehreren Prozessen	281
5.5.1	Anpassung der Signalbehandlung	282
5.5.2	Ein neuer Prozeß pro Client	284
5.5.3	Das Hauptprogramm	286
5.5.4	Eigenschaften und Einsatzgebiete	287
5.6	Nebenläufige Server mit Preforking	289
5.6.1	Buchführende Signalbehandlung	290
5.6.2	Parallele Accept-Handler in mehreren Prozessen	292
5.6.3	Preforking im Hauptprogramm	294
5.6.4	Eigenschaften und Einsatzgebiete	296
5.7	Zusammenfassung	298
6	Netzwerkprogrammierung mit SSL	301
6.1	Strategien zur Absicherung des Datenverkehrs	302
6.1.1	Datenverschlüsselung	304
6.1.2	Hashfunktionen und Message Authentication Codes	307
6.1.3	Digitale Signaturen	308
6.1.4	Zertifizierungsstellen und digitale Zertifikate	309
6.1.5	Praktische Absicherung des Datenverkehrs	310
6.2	SSL-Grundlagen	313
6.2.1	Datentransfer über SSL	314
6.2.2	Anwendungsprotokolle um SSL erweitern	317
6.2.3	SSL-Verbindungen interaktiv testen	321
6.3	OpenSSL-Basisfunktionalität	323
6.3.1	Das Konzept der BIO-API	324
6.3.2	Lebenszyklus von BIO-Objekten	325

6.3.3	Ein-/Ausgabe über BIO-Objekte	326
6.3.4	BIO-Quellen/Senken und BIO-Filter	329
6.3.5	Fehlerbehandlung	342
6.3.6	Thread-Support	345
6.3.7	Pseudozufallszahlengenerator	352
7	Client-/Server-Programmierung mit OpenSSL	357
7.1	Initialisierung der ssl-Bibliothek	357
7.2	Der SSL-Kontext	359
7.2.1	Ein unvollständiger SSMTP-Client	360
7.2.2	SSL-Optionen, SSL-Modi und Chiffrenfolgen	364
7.3	Sicherer Umgang mit X.509-Zertifikaten	368
7.3.1	Zertifikatsüberprüfung aktivieren	372
7.3.2	Zertifikatsüberprüfung per Callback nachbereiten	374
7.3.3	Identitätsabgleich mit digitalen Zertifikaten	380
7.3.4	SSL-Kommunikation mit eigener Identität	387
7.4	Client-/Server-Beispiel: SMTP mit SARTTLS	389
7.4.1	Ein SMTP-Client mit STARTTLS	389
7.4.2	Ein SMTP-Server mit STARTTLS	397
7.5	Zusammenfassung	406
A	Anhang	409
A.1	Zertifikate erstellen mit OpenSSL	409
A.1.1	Aufbau einer Zertifizierungsstelle	409
A.1.2	Neue Zertifikate ausstellen	412
A.1.3	Vertrauenswürdige Zertifizierungsstellen	414
A.2	Barrieren mit POSIX-Threads	415
	Literaturverzeichnis	423
	Sachverzeichnis	427

Beispielprogramme

2.1	exit-test.c	18
2.2	write-data.c	27
2.3	read-data.c	29
2.4	read-data-stdin.c	30
2.5	readn.c	32
2.6	writen.c	33
2.7	readn-data-stdin.c	34
2.8	logfile.c	42
2.9	syslog.c	45
2.10	overflow.c, Teil 1	47
2.11	overflow.c, Teil 2	48
2.12	cat1.c	54
2.13	cat2.c	55
2.14	quiz1.c	62
2.15	quiz2.c	65
2.16	signal-block.c	70
2.17	signal-wait.c	73
2.18	show-ids.c	77
2.19	fork-ids.c	79
2.20	wait-child.c	87
2.21	exec-test.c	93
2.22	daemon.c, Teil 1	97
2.23	daemon.c, Teil 2	99
3.1	pthread-lifecycle.c	106
3.2	pthread-exit1.c	110
3.3	pthread-exit2.c	113
3.4	average.c	117
3.5	average-mutex.c	123
3.6	average-mutex-cv.c, Teil 1	131
3.7	average-mutex-cv.c, Teil 2	133
3.8	pthread-signal.c, Teil 1	141

3.9	pthread-signal.c, Teil 2	143
4.1	quiz3.c	154
4.2	bitwise.c	171
4.3	ntoatest.c	175
4.4	byteorder.c	177
4.5	timeclient.c	186
4.6	timeserver.c	196
4.7	timeclientudp.c	207
4.8	timeclientudpconn.c	211
4.9	quiz4.c	215
4.10	select.c	221
4.11	hostaddr.c	231
5.1	server.h	238
5.2	tcp-listen.c	241
5.3	tcp-connect.c	243
5.4	readwrite.c, Teil 1	246
5.5	readwrite.c, Teil 2	248
5.6	test-cli.c, Teil 1	250
5.7	test-cli.c, Teil 2	251
5.8	test-cli.c, Teil 3	253
5.9	iter-srv.c	257
5.10	handle-client.c	260
5.11	srv-stats.c	263
5.12	thread-srv.c, Teil 1	267
5.13	thread-srv.c, Teil 2	269
5.14	thread-srv.c, Teil 3	270
5.15	pthread-srv.c, Teil 1	275
5.16	pthread-srv.c, Teil 2	278
5.17	fork-srv.c, Teil 1	283
5.18	fork-srv.c, Teil 2	284
5.19	fork-srv.c, Teil 3	286
5.20	prefork-srv.c, Teil 1	291
5.21	prefork-srv.c, Teil 2	292
5.22	prefork-srv.c, Teil 3	294
6.1	bio-timeclient.c	333
6.2	openssl-thread-init.c, Teil 1	348
6.3	openssl-thread-init.c, Teil 2	349
6.4	openssl-thread-init.h	352
6.5	openssl-rand-seed.c	354
7.1	openssl-lib-init.c	358
7.2	openssl-lib-init.h	359
7.3	bio-ssl-smtpcli1.c	361
7.4	openssl-util.c, Teil 1	372
7.5	openssl-util.c, Teil 2	378
7.6	openssl-util.c, Teil 3	383

7.7	openssl-util.c, Teil 4	385
7.8	openssl-util.h	386
7.9	bio-ssl-smtpcli2.c, Teil 1	390
7.10	bio-ssl-smtpcli2.c, Teil 2	393
7.11	bio-ssl-smtpsrv.c, Teil 1	398
7.12	bio-ssl-smtpsrv.c, Teil 2	400
7.13	bio-ssl-smtpsrv.c, Teil 3	403
A.1	barrier.h	416
A.2	barrier.c	418

Einführung

Das *Internet* ist in seiner heutigen Form ein weltumspannendes Rechnernetz, das sich selbst aus einer Menge einzelner, voneinander unabhängiger Netzwerke zusammensetzt. Der Begriff Internet ist dabei als Kurzschreibweise für *Interconnected Networks*, also miteinander verbundene Netzwerke, entstanden. Angefangen bei einzelnen Rechnersystemen, die zu einem *Local Area Network (LAN)* zusammengeschlossen sind, werden mehrere LANs über größere Distanzen hinweg zu einem *Wide Area Network (WAN)* verknüpft. Der verschachtelte Verbund vieler derartiger WANs ergibt schließlich das Internet.

Natürlich hatte zu den Anfangszeiten des Internets niemand die kühne Vision, ein Rechnernetz bestehend aus knapp 400 Millionen¹ miteinander vernetzter Computersysteme zu initiieren. Vielmehr sollte gegen Ende der 60'er Jahre im Auftrag der *Advanced Research Projects Agency (ARPA)*, die damals für das US-amerikanische Verteidigungsministerium Forschungsprojekte förderte, ein Netzwerk entstehen, über das die damals knappen Rechenkapazitäten der Hochschulen durch den Austausch von Daten besser ausgenutzt werden sollten. Das resultierende Netzwerk hieß zu dieser Zeit auch noch nicht Internet, sondern *ARPANET*, und wurde Ende 1969 von der *University of California, Los Angeles*, der *University of California, Santa Barbara*, der *University of Utah* und dem *Stanford Research Institute* in Betrieb genommen.² Das Netz verband zunächst genau vier Rechner der vier beteiligten Universitäten.

Um eine möglichst hohe Ausfallsicherheit zu erreichen, ersannen die beteiligten Forscher ein *paketvermitteltes Netzwerk*. In paketvermittelten Netzen werden die zu übertragenden Daten vom Absender in einzelne Pakete zerlegt und vom Empfänger nach dem Eintreffen wieder zusammengesetzt. In einem komplexen Netzwerk können die einzelnen Pakete bei der Übertragung

¹ Stand: Januar 2006, siehe <http://www.isc.org/ds/>

² Der erste Datenaustausch soll laut [ZEI01] am 29. Oktober 1969 stattgefunden haben, um die Buchstabenkombination *LOG* zu übermitteln. Von den gleichzeitig telefonierenden Technikern werden die Worte „Hast du das L?“ – „Ja!“ – „Hast du das O?“ – „Ja!“ – „Hast du das G?“ überliefert, dann sei der Rechner abgestürzt.

durchaus unterschiedliche Wege gehen. Absender und Empfänger sind also, im Gegensatz zu *leitungsvermittelten Netzwerken* mit fester Verbindung, lediglich lose miteinander verbunden.

Weitere Forschungsarbeiten führten bis 1974 zur Entwicklung von TCP/IP, einer Familie von Netzwerkprotokollen, die bis heute die Basis für den Datenaustausch über das Internet bilden. Um die Akzeptanz von TCP/IP zu forcieren, wurde die *University of California at Berkeley* damit beauftragt, TCP/IP in *Berkeley Unix* zu integrieren. Mit der zunehmenden Verbreitung von TCP/IP wuchs gleichzeitig die Anzahl der Rechner im ARPANET rapide an. Der große Verbund von Netzen, den das ARPANET nun langsam aber sicher darstellte, wurde schließlich als *das Internet* bekannt. TCP/IP wird seitdem auch als Internet-Protokoll-Familie bezeichnet.

1.1 TCP/IP-Grundlagen

Um die komplexen Zusammenhänge, die in einem Rechnernetz bestehen, besser strukturieren und beschreiben zu können, wurden sogenannte Referenzmodelle eingeführt. Das bekannteste dieser Modelle ist das von der *International Organization for Standardization (ISO)* standardisierte *OSI-Referenzmodell (Open Systems Interconnection Reference Model)*. Nachdem die Internet-Protokoll-Familie bereits vor dem OSI-Referenzmodell entwickelt wurde, liegt der TCP/IP-Architektur ein anderes, vereinfachtes Referenzmodell zugrunde. Die Erfahrungen mit dem TCP/IP-Referenzmodell sind dann bei der Ausarbeitung des OSI-Referenzmodells eingeflossen.

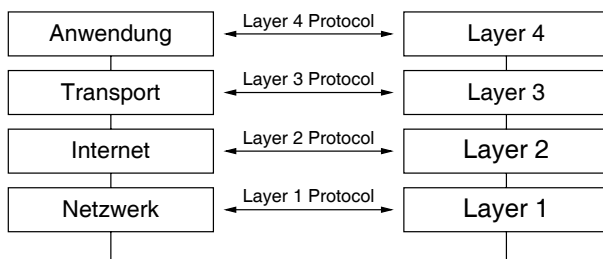


Abb. 1.1. Kommunikation im TCP/IP-Referenzmodell

Das TCP/IP-Referenzmodell unterteilt die einzelnen Aufgabenstellungen der Netzwerkkommunikation, angefangen bei der physikalischen Übertragung der Datensignale über die Vermittlung von Datenpaketen bis hin zu speziellen, anwendungsspezifischen Aufgaben, in einen Stapel von vier Schichten. Die einzelnen Schichten des Modells sind in Abb. 1.1 zu sehen. Jede der Schichten im Referenzmodell definiert bestimmte Funktionen, die in Form von Diensten und Protokollen implementiert werden:

- Jede Schicht kann die Dienste der darunterliegenden Schicht nutzen, ohne dabei konkrete Kenntnisse darüber besitzen zu müssen, wie die in Anspruch genommenen Dienstleistungen genau erbracht werden bzw. implementiert sind. Sämtliche Daten, die von einer Anwendung über das Netzwerk verschickt werden, durchlaufen von der Anwendungs- bis zur Netzwerkschicht alle Schichten des Modells (und auf dem empfangenden System wieder in der umgekehrten Reihenfolge).
- Die gleichnamigen Schichten zweier kommunizierender Systeme kooperieren miteinander jeweils über spezielle, schichtspezifische Protokolle. Dazu werden den von einer Anwendung verschickten Daten in jeder Schicht protokollspezifische Informationen hinzugefügt (und von der gleichnamigen Schicht auf dem empfangenden System wieder entnommen).

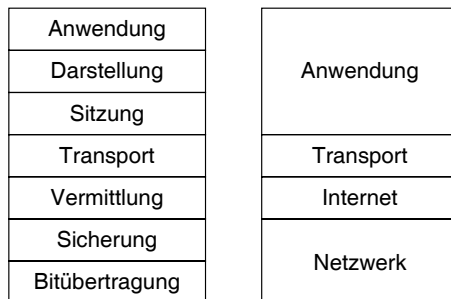


Abb. 1.2. OSI- und TCP/IP-Referenzmodell im Vergleich

Im direkten Vergleich zwischen OSI- und TCP/IP-Referenzmodell (vgl. dazu Abb. 1.2) zeigt sich die etwas feinere Untergliederung des OSI-Modells in einen Stapel von insgesamt sieben Schichten zusammen mit der korrespondierenden Aufteilung des TCP/IP-Modells.

1.1.1 Netzwerkschicht

Die beiden untersten Schichten des OSI-Referenzmodells repräsentieren grob gesprochen die Netzwerk-Hardware und die zugehörigen Gerätetreiber. Sie kapseln die eingesetzten Netzwerktechnologien wie z. B. *Ethernet*, *Token Ring* oder *FDDI* und liefern den übergeordneten Schichten eine davon unabhängige Schnittstelle zur Kommunikation. Diese beiden Schichten sind in der *Netzwerkschicht* des TCP/IP-Referenzmodells zusammengefaßt und spielen bei der Entwicklung netzwerkfähiger Client-/Server-Programme im Allgemeinen keine Rolle.

1.1.2 Internet-Schicht

Der Vermittlungsschicht des OSI-Modells entspricht im TCP/IP-Modell die *Internet-Schicht*. Diese Schicht kümmert sich in paketvermittelten Netzwerken wie dem Internet um die Weitervermittlung der einzelnen Datenpakete. Die Internet-Schicht etabliert damit eine Rechner-zu-Rechner-Verbindung, unabhängig von der zugrundeliegenden Netzwerkstruktur (bestehend aus zahlreichen LANs und WANs). Dadurch befreit die Schicht die übergeordneten Schichten von den Details der Datenübertragung.

Internet Protocol (IP)

Die Internet-Schicht wird bei TCP/IP vom *Internet Protocol (IP)* besetzt. Die wesentliche Aufgabe dieses in RFC 791 [IP81] spezifizierten Protokolls ist die Adressierung von Netzteilnehmern über *IP-Adressen* und die Datenübertragung von und zu anderen Teilnehmern. Die Daten werden dabei in einem vom Internet Protocol definierten Paketformat als sogenannte *Datagramme* an die tieferliegenden Schichten übergeben und über das Netzwerk übertragen. Bei Bedarf werden die Datagramme vom Internet Protocol fragmentiert, d. h. vor dem Versand in kleinere Fragmente zerlegt und nach dem Empfang wieder zusammengesetzt.

Beim Internet Protocol handelt es sich um ein *verbindungsloses Protokoll*, d. h. es existiert auf der IP-Ebene keine Ende-zu-Ende-Verbindung zwischen den kommunizierenden Anwendungen. Die einzelnen IP-Datagramme werden unabhängig von den jeweils anderen Datagrammen zugestellt, insbesondere kann sich dabei die Empfangsreihenfolge der Datagramme von der Reihenfolge beim Versand unterscheiden. Darüber hinaus wird der Empfang der einzelnen IP-Datagramme nicht bestätigt. Das Internet Protocol bietet also bezüglich der Datenübertragung keine Zuverlässigkeitsgarantie, weshalb das Protokoll auch oft als *nicht zuverlässig* oder sogar als *unzuverlässig* bezeichnet wird.³ Das Internet Protocol kann demnach die tatsächliche Zustellung der verschickten Datagramme nicht garantieren, die einzelnen IP-Datagramme werden lediglich bestmöglich verschickt.

Internet Control Message Protocol (ICMP)

Das in RFC 792 [Pos81] spezifizierte *Internet Control Message Protocol (ICMP)* ist ebenfalls auf der Internet-Schicht angesiedelt. Das Protokoll ist

³ Das Attribut *unzuverlässig* bedeutet in diesem Fall nicht, daß beim Internet Protocol ständig etwas schief läuft und man sich deshalb generell nicht darauf verlassen kann. Die Bezeichnung weist vielmehr darauf hin, daß das Protokoll selbst keine Mechanismen zur Erkennung und Behebung von Paketverlusten besitzt, so daß diese Aufgaben bei Bedarf von höheren Schichten übernommen werden müssen.

integraler Bestandteil jeder IP-Implementierung und transportiert Status-, Fehler- und Diagnoseinformationen für das Internet Protocol. Die ICMP-Pakete werden dazu als IP-Datagramme über das Netzwerk übertragen.

1.1.3 Transportschicht

Die *Transportschicht* ist die erste (niedrigste) Schicht, die den übergeordneten anwendungsorientierten Schichten, also den Schichten 5–7 des OSI-Modells bzw. Schicht 4 des TCP/IP-Modells, eine vollständige Ende-zu-Ende-Kommunikation zwischen zwei Anwendungen zur Verfügung stellt. Die Transportschicht leitet den Datenfluß von der Anwendung zur Internet-Schicht und umgekehrt. Die Adressierung der Anwendung erfolgt dabei über einen sogenannten *Port*, an den die Anwendung zuvor gebunden wurde.

User Datagram Protocol (UDP)

Genau wie das zugrundeliegende Internet Protocol ist auch das *User Datagram Protocol (UDP)* ein verbindungsloses Transportprotokoll ohne Zuverlässigkeitsgarantie.⁴ Die Spezifikation des Protokolls findet sich in RFC 768 [Pos80]. Das User Datagram Protocol erweitert die bereits vom Internet Protocol erbrachten Leistungen lediglich um die Portnummern der sendenden und empfangenden Anwendung. Die von UDP transportierten Pakete werden deshalb in Anlehnung an IP auch *UDP-Datagramme* genannt. UDP hat demzufolge auch nur einen minimalen *Protokoll-Overhead*, was das Protokoll v. a. für Anwendungen, die nur wenige Daten übertragen müssen, interessanter als das nachfolgend beschriebene TCP macht.

Soll auf Basis von UDP eine zuverlässige, reihenfolgetreue Datenübertragung erfolgen, so muß sich eine der übergeordneten Schichten (sprich: die Anwendung selbst) um diese Aufgabe kümmern.

Transmission Control Protocol (TCP)

Das ursprünglich in RFC 793 [TCP81] spezifizierte *Transmission Control Protocol (TCP)* ist, im Gegensatz zu UDP, ein verbindungsorientiertes und zuverlässiges Transportprotokoll. Das Protokoll stellt je zwei Kommunikationspartnern eine virtuelle Ende-zu-Ende-Verbindung (im Sinne einer festen Punkt-zu-Punkt-Verbindung in einem leitungsvermittelten Netzwerk) zur Verfügung. Die Datenübertragung erfolgt bei TCP deshalb in drei Phasen: dem Verbindungsaufbau, der eigentlichen Datenübertragung und dem

⁴ Bezüglich der (fehlenden) Zuverlässigkeitsgarantie gelten die selben Anmerkungen wie beim Internet Protocol (vgl. dazu Abschnitt 1.1.2).

abschließenden Verbindungsabbau. Die von TCP übermittelten Datenpakete werden *TCP-Segmente* genannt. Die Reihenfolge der über eine solche virtuelle Ende-zu-Ende-Verbindung verschickten TCP-Segmente bleibt im Rahmen der Übertragung erhalten. Nachdem das Protokoll zudem den Empfang eines TCP-Segments gegenüber dem Sender quittiert und der Sender beim Ausbleiben einer solchen Empfangsbestätigung die Übertragung des betreffenden Segments wiederholt, wird das Transmission Control Protocol als zuverlässiges Transportprotokoll bezeichnet.

TCP nimmt Daten von den übergeordneten Schichten als *Datenstrom* an und teilt diesen Datenstrom auf einzelne TCP-Segmente auf. Jedes dieser TCP-Segmente wird dann über die Internet-Schicht als IP-Datagramm verschickt. Umgekehrt werden die einer TCP-Verbindung zugeordneten IP-Datagramme bzw. TCP-Segmente auf der Empfängerseite wieder zum ursprünglichen Datenstrom zusammengesetzt. Nachdem das zugrundeliegende Internet Protocol bezüglich der Datenübertragung keine Zuverlässigkeitsgarantie übernimmt, muß sich das Transmission Control Protocol selbst um diese Aufgabe (wiederholtes Senden verlorengegangener Pakete, Einhalten der Reihenfolge) kümmern. Die TCP-Segmente erhalten dazu sogenannte *Sequenz- und Bestätigungsnummern*, mit deren Hilfe die Segmente quittiert, ggf. neu übertragen und in der korrekten Reihenfolge wieder zu einem Datenstrom zusammengesetzt werden können.

1.1.4 Anwendungsschicht

Die oberste Schicht des TCP/IP-Referenzmodells ist die *Anwendungsschicht*. Sie faßt die Schichten 5–7 des OSI-Modells in einer Schicht zusammen. Für diese Schicht ist inzwischen eine Fülle von Anwendungsprotokollen spezifiziert, über welche heute die gängigen Internetdienste miteinander kommunizieren. Stellvertretend seien an dieser Stelle die weithin bekannten Anwendungsprotokolle *Hypertext Transfer Protocol (HTTP)* für das *World Wide Web (WWW)* oder das *Simple Mail Transfer Protocol (SMTP)* zum Austausch *elektronischer Post (E-Mail)* genannt.

Die Anwendungsprotokolle bzw. die Anwendungen, die diese Anwendungsprotokolle implementieren, greifen ihrerseits über die *Socket-API* auf die Dienste der Transportschicht zu. Im weiteren Verlauf dieses Buchs kümmern wir uns nun im Wesentlichen um die Funktionalität der Socket-API und ihre Anwendung im Rahmen eigener Client-/Server-Programme.

1.2 Internet-Standards

Die sogenannten *RFCs (Request for Comments)*, von denen wir in diesem Kapitel bereits mehrere zitiert haben, spielten und spielen bei der Entstehung

und Weiterentwicklung des Internets eine gewichtige Rolle. Dabei handelt es sich um eine fortlaufend ergänzte Reihe von technischen und organisatorischen Dokumenten, über welche die Standards für das Internet festgelegt werden. Sie beschreiben die Dienste und Protokolle des Internets und legen Regeln und Grundsätze für dieses Netzwerk fest. Die Sammlung aller RFCs wird an zahlreichen Stellen im Internet publiziert, u. a. auf der Homepage des *RFC-Editors*.⁵

Neue RFC-Dokumente werden von einer Arbeitsgruppe bzw. dem RFC-Editor geprüft und durchlaufen dabei bestimmte Reifestufen. Ein einmal veröffentlichter RFC wird nie wieder verändert oder aktualisiert. Stattdessen wird er bei Bedarf durch einen neuen RFC ersetzt und erhält den Zusatz, daß er durch den neuen RFC abgelöst wurde. Der neue RFC enthält seinerseits einen Hinweis auf den RFC, den er abgelöst hat.

Neben ihrer verantwortungsvollen Aufgabe als eine Art Standardisierungsgremium beweisen RFC-Autoren mitunter auch sehr feinen Humor, bevorzugt am 1. April jeden Jahres. Stellvertretend für viele amüsante Exkursionen sei an dieser Stelle auf den lesenswerten und durchaus zu diesem Buch passenden RFC 1925 [Cal96], *The Twelve Networking Truths*, hingewiesen.

1.3 Unix-Standards

Die Entstehungsgeschichte von Unix reicht wie die Geschichte des Internets ins Jahr 1969 zurück. Ken Thompson begann damals bei den Bell Laboratories mit der Entwicklung des neuen Betriebssystems *UNICS*, einer in Assembler geschriebenen, abgespeckten Version des Betriebssystems *MULTICS*.⁶ Der Name UNICS wandelte sich im weiteren Projektverlauf zu Unix.⁷ Das Unix-System wurde dann Anfang der 70'er Jahre in der, im Rahmen der Unix-Entwicklung entstandenen, Programmiersprache C neu implementiert und wenig später zusammen mit einem C-Compiler kostenfrei an verschiedene Universitäten verteilt. In der Folgezeit entstanden viele variierende Systemlinien von Unix mit jeweils unterschiedlichen Kommandos, Kommandooptionen und Systemschnittstellen. So wurde z. B. von der *University of California at Berkeley* neben anderen Anpassungen auch TCP/IP in das dort entstandene *Berkeley Unix* integriert.

⁵ <http://www.rfc-editor.org/>

⁶ Übrigens: Der Name UNICS enthält gleich ein zweifaches Wortspiel: Zum einen dokumentiert *UNI* im Gegensatz zu *MULTI*, daß es sich bei UNICS um ein abgespecktes MULTICS handelt. Zum anderen wird UNICS in etwa so ausgesprochen wie das Wort *Eunuchs*, was UNICS als kastriertes MULTICS darstellt.

⁷ Streng genommen steht die Schreibweise *UNIX* für Unix-Systeme, deren Implementierung auf den ursprünglichen Unix-Quellen der Bell Laboratories basiert, während durch die Schreibweise *Unix* zusätzlich auch Unix-ähnliche Systeme, wie z. B. Linux, mit eingeschlossen werden.

Die vielen verschiedenen Unix-Derivate mit ihren gegenseitigen Inkompatibilitäten führten schließlich zu den ersten Standardisierungsbemühungen und es entstanden die sogenannten *POSIX-Standards*.⁸ Die aktuelle Version [SUS02], auf die sich die Ausführungen des vorliegenden Buchs beziehen, ist unter den drei Namen *IEEE Std 1003.1-2001*, *ISO/IEC 9945:2002* und *Single UNIX Specification, Version 3* bekannt. Wir werden im weiteren Verlauf den Namen IEEE Std 1003.1-2001 oder einfach nur kurz POSIX-Standard verwenden.

⁸ Das Akronym *POSIX* steht für *Portable Operating System Interface*, ergänzt um die Unix-typische Endung *IX*.

Programmieren mit Unix-Prozessen

Prozesse bilden die Grundlage der Datenverarbeitung auf Unix-Systemen. Es ist also nicht weiter überraschend, daß vertiefte Programmierkenntnisse mit Prozessen die beste Basis für die Implementierung solider Netzwerkanwendungen sind. Aus diesem Grund vermittelt dieses Kapitel zunächst die notwendigen Grundkenntnisse der Programmierung von Unix-Prozessen, auf die im weiteren Verlauf dieses Buchs aufgebaut wird. In [Ste92, Her96] finden Sie bei Bedarf ausführlichere Informationen zur Unix-Systemprogrammierung. Die Grundlage für dieses Kapitel sowie für den weiteren Verlauf dieses Buchs bilden die in IEEE Std 1003.1-2001 [SUS02] vereinbarten Unix-Standards.

2.1 Unix-Prozesse

Unter einem *Prozeß* wird in einer Unix-Umgebung ein *in Bearbeitung befindliches Programm* bezeichnet. Wird also ein ausführbares Programm, z. B. ein zuvor mit einem C-Compiler übersetztes C-Programm, gestartet, so nennt man die gerade laufende Instanz dieses Programms einen Prozeß. Formal definiert sich ein Prozeß als eine Einheit, die aus

- einem (zeitlich invarianten) Programm,
- einem Satz von Daten, mit denen der Prozeß initialisiert wird, und
- einem (zeitlich varianten) Zustand besteht.

Im klassischen Kontext liegt jedem Prozeß ein sequentielles Programm zugrunde. Dies bedeutet, daß die einzelnen Anweisungen des Programms und damit die einzelnen Schritte des implementierten Algorithmus' immer in einer bestimmten Reihenfolge nacheinander durchlaufen werden.

In allen modernen Unix-Systemen ist diese Annahme in der Zwischenzeit überholt. Die bislang beschriebenen Prozesse werden dort als *schwergewichtige Prozesse* bezeichnet, welche sich dann ihrerseits aus einem oder mehreren *leichtgewichtigen Prozessen*, sogenannten Threads, zusammensetzen. Der Standard IEEE Std 1003.1-2001 definiert deshalb wie folgt:

Eine Prozeß ist eine Einheit bestehend aus

- einem Adreßraum,
- mit einem oder mehreren in diesem Adreßraum ablaufenden Threads und
- den für die Threads benötigten Systemressourcen.

In Kapitel 3 werden wir ausführlich auf die Programmierung mit Threads eingehen. Für den weiteren Verlauf dieses Kapitels können wir allerdings unter einem Prozeß meist problemlos die klassische, sequentielle Version verstehen.

2.1.1 Prozeßgruppen und Sessions

Jedem Unix-Prozeß wird bei seinem Start vom Betriebssystem eine eindeutige Prozeßnummer (Prozeß-ID, kurz: PID) zugewiesen. Neue Prozesse können ausschließlich von bereits vorhandenen Prozessen erzeugt werden. Der erzeugende Prozeß wird als Elternprozeß, der neue Prozeß als Kindprozeß bezeichnet. Beim Systemstart wird unter Unix ein ausgezeichneter Prozeß, der sogenannte *init*-Prozeß, gestartet. Diesem Prozeß wird *immer* die PID 1 zugewiesen. Der *init*-Prozeß bildet den Ursprung aller Prozesse, alle weiteren Prozesse stammen mehr oder weniger direkt von diesem Prozeß ab.

Die verschiedenen Prozesse werden in Unix zu *Prozeßgruppen* zusammengefaßt und jede Prozeßgruppe gehört wiederum zu einer *Session*. Auch Prozeßgruppen und Sessions bekommen vom Betriebssystem eindeutige Identifikationsnummern (Prozeßgruppen-ID und Session-ID bzw. PGID und SID) zugewiesen.

Eine neue Session wird z.B. dann gestartet, wenn sich ein Benutzer über ein Terminal erfolgreich am System anmeldet. Für den Benutzer wird bei der Anmeldung eine Login-Shell, etwa die Korn-Shell (ksh), gestartet. Für diese Shell wird gleichzeitig eine neue Prozeßgruppe erzeugt. Werden nun von der Shell aus weitere Programme gestartet, gehören alle diese Prozesse inklusive der Login-Shell standardmäßig zur gleichen Session. Der erste Prozeß in einer neuen Session, in unserem Beispiel die Korn-Shell, wird als *Anführer der Session* (*Session Leader*) bezeichnet.

Um das beschriebene Szenario zu veranschaulichen, melden wir uns an einem Unix-System an und geben auf der Kommandozeile die beiden folgenden, an sich nicht besonders sinnvollen, Befehlsfolgen ein:

```
$ sleep 60 | tail &
$ ps | cat | head
```

Die beiden Kommandos der ersten Zeile werden von der Shell im Hintergrund gestartet. Während diese beiden Prozesse noch laufen, führt die Shell die nächsten drei Prozesse im Vordergrund aus. Abbildung 2.1 zeigt den Zusammenhang zwischen Session, Prozeßgruppen und Prozessen.

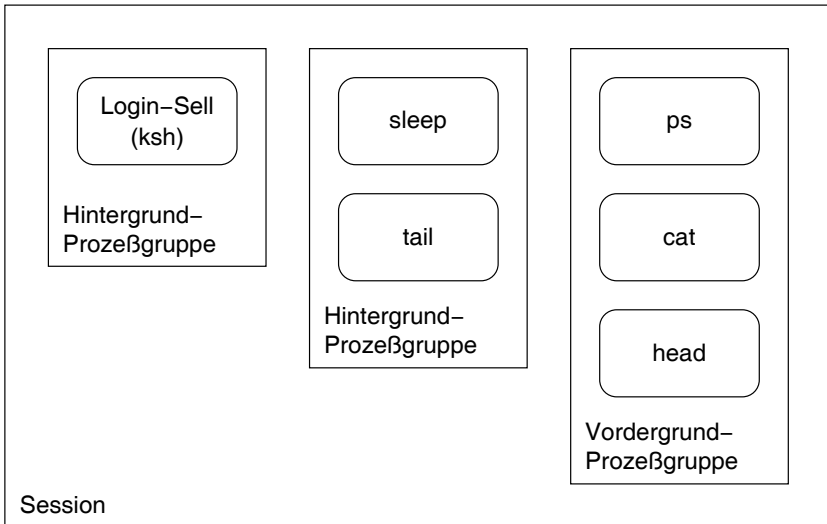


Abb. 2.1. Zusammenhang zwischen Session, Prozeßgruppen und Prozessen

Wie beschrieben gehören alle Prozesse zur gleichen Session. Neben der Login-Shell bilden die zuerst gestarteten Prozesse *sleep* und *tail* sowie die danach gestarteten Prozesse *ps*, *cat* und *head* zusammen je eine eigenständige Prozeßgruppe. Die Prozeßgruppe, die im Vordergrund ausgeführt wird – davon kann es pro Session maximal eine geben – wird als Vordergrund-Prozeßgruppe bezeichnet. Analog heißen alle anderen Prozeßgruppen der Session Hintergrund-Prozeßgruppen. Der nachfolgende Auszug aus der Prozeßtabelle zeigt die Details:

PID	PPID	PGID	SID	TTY	TPGID	COMMAND
2701	2698	2701	2701	tty0	3271	-ksh
3269	2701	3269	2701	tty0	3271	sleep
3270	2701	3269	2701	tty0	3271	tail
3271	2701	3271	2701	tty0	3271	ps
3272	2701	3271	2701	tty0	3271	cat
3273	2701	3271	2701	tty0	3271	head

Die sechs Prozesse (PID 2701 bis 3273) bilden drei verschiedene Prozeßgruppen (PGID 2701, 3269 und 3271). Die Elternprozeß-ID (PPID) zeigt an, daß die Login-Shell (PID 2701) die anderen Prozesse (PID 3269 bis 3273) gestartet hat. Die Login-Shell selbst wurde wieder von einem anderen Prozeß initiiert.

Die Vordergrund-Prozeßgruppe erkennt man in der Prozeßtabelle an der Übereinstimmung zwischen der Prozeßgruppen-ID (PGID) und der Terminal-Prozeßgruppen-ID (TPGID). Die Prozesse *ps*, *cat* und *head* bilden im vorangehenden Beispiel die Vordergrund-Prozeßgruppe. Alle Prozesse gehören zur selben Session (SID 2701). Die Prozesse, bei denen die Prozeß-ID mit der Prozeßgruppen-ID übereinstimmt, werden als Anführer der Prozeßgruppe (process group leader) bezeichnet. Im obigen Beispiel sind das die Login-Shell *ksh* sowie die beiden Prozesse *sleep* und *ps*.

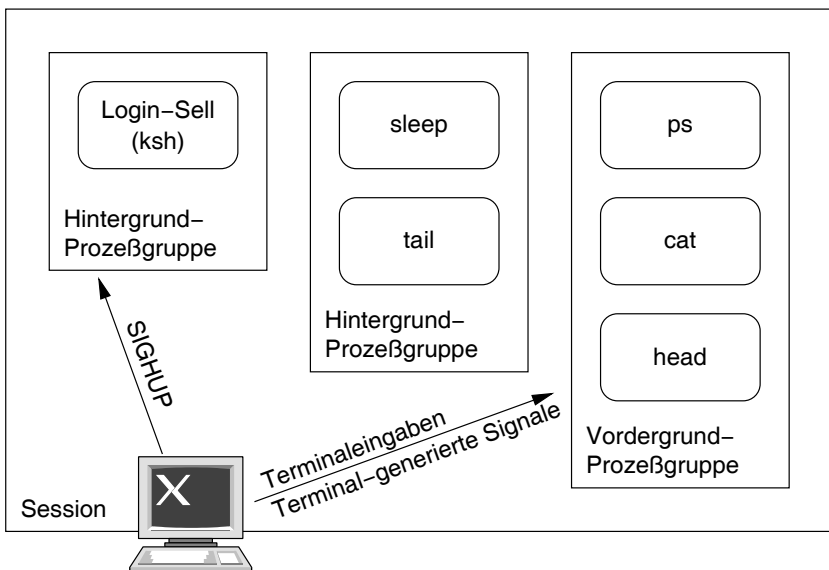


Abb. 2.2. Session und kontrollierendes Terminal

2.1.2 Kontrollierendes Terminal

Einer Session kann genau ein *kontrollierendes Terminal* zugeordnet sein. Dies ist z. B. immer dann der Fall, wenn die Session bei einem erfolgreichen Login erzeugt wurde. Über das kontrollierende Terminal können Daten an die Prozesse der Vordergrund-Prozeßgruppe geschickt werden. In unserem Beispiel bedeutet dies, daß alle Tastatureingaben am Login-Terminal an die Prozesse mit der PGID 3271 gesendet werden. Außerdem werden die über dieses Terminal generierten Signale (etwa SIGSTOP zum Stoppen oder SIGTERM zum Ab-

bruch eines Programms) an alle Prozesse der Vordergrund-Prozeßgruppe geschickt. In unserem Fall würden durch das Abbruchsignal die drei Kommandos *ps*, *cat* und *head* beendet. Danach existiert die Prozeßgruppe 3271 nicht mehr und die Prozeßgruppe der Login-Shell wird zur Vordergrund-Prozeßgruppe, an die dann bis auf weiteres alle Tastatureingaben gehen.

Der Anführer einer Session, der die Verbindung zum kontrollierenden Terminal herstellt, wird als *kontrollierender Prozeß* (*Controlling Process*) bezeichnet. Der kontrollierende Prozeß erhält das *SIGHUP*-Signal (*Hang Up*), sobald das kontrollierende Terminal die Verbindung trennt. Abbildung 2.2 veranschaulicht diesen Zusammenhang zwischen einer Session und dem assoziierten kontrollierenden Terminal.

2.1.3 Verwaiste Prozesse und verwaiste Prozeßgruppen

Beendet sich ein Prozeß, von dem noch Kindprozesse aktiv sind, so verlieren diese Kindprozesse natürlich ihren Elternprozeß. In Analogie zum richtigen Leben werden diese Kindprozesse dann als *verwaiste Prozesse* bezeichnet. Die Elternprozeß-ID solcher Waisen wäre damit nicht mehr gültig, sie referenziert keinen aktiven Prozeß mehr. Der Standard sieht deshalb vor, daß die Elternprozeß-ID in diesem Fall automatisch auf einen (von der Unix-Implementierung vorgegebenen) Systemprozeß gesetzt wird. In der Praxis ist dies der *init*-Prozeß.

Auch Prozeßgruppen können unter Unix verwaisen. Dieser Fall tritt genau dann ein, wenn für alle Elternprozesse aller Mitglieder einer Prozeßgruppe folgendes zutrifft:

- Der Elternprozeß ist entweder selbst Mitglied dieser Prozeßgruppe oder
- der Elternprozeß gehört nicht zur Session dieser Prozeßgruppe.

Mit anderen Worten ist eine Prozeßgruppe also solange *nicht* verwaist, solange ein Prozeß aus der Gruppe einen Elternprozeß besitzt, der zu einer anderen Prozeßgruppe innerhalb der gleichen Session gehört.

Würde sich im Beispiel aus Abb. 2.1 die Login-Shell beenden, bevor sich die Prozesse der beiden anderen Prozeßgruppen beendet haben, so wären alle Prozesse dieser Gruppen verwaist. Ihre Elternprozeß-ID würde in diesem Fall jeweils automatisch auf die Prozeß-ID von *init* gesetzt. Der *init*-Prozeß wäre damit für alle der neue Elternprozeß. Da für die Login-Shell eine eigene Session erzeugt wurde, hätten nun alle Prozesse aus den beiden Prozeßgruppen mit *init* einen Elternprozeß, der zu einer anderen Session (und damit selbstverständlich auch zu einer anderen Prozeßgruppe) gehört. Die beiden Prozeßgruppen wären damit verwaist.

2.1.4 Prozeßumgebung

Intern bestehen Prozesse im wesentlichen aus ihren Code- und Datenbereichen. Abbildung 2.3 zeigt die verschiedenen Segmente eines Prozesses:

- In das *Text-Segment* wird der Maschinencode des Programms geladen,
- im *Data-Segment* werden die explizit initialisierten globalen Variablen gespeichert,
- das *BSS-Segment* enthält alle uninitialisierten globalen Variablen (welche dann beim Programmstart mit 0 initialisiert werden),
- im nach oben wachsenden *Heap-Segment* (der Halde) werden die dynamisch allozierten Speicherbereiche angelegt und
- das nach unten wachsende *Stack-Segment* (der Stapel), enthält die automatischen und temporären Variablen, die Rücksprungadressen bei Funktionsaufrufen, zwischengespeicherte Registerwerte und einiges mehr.

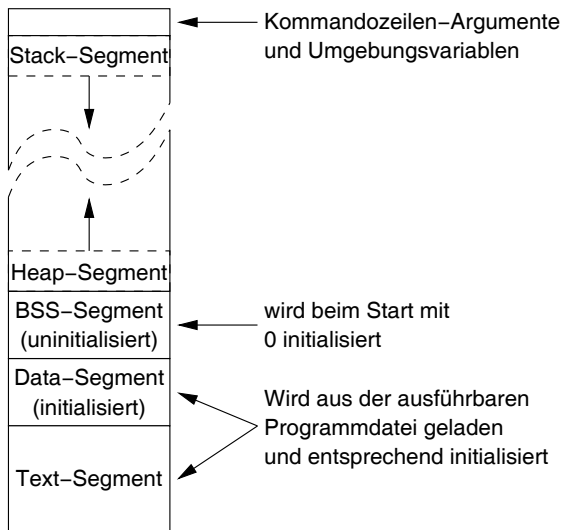


Abb. 2.3. Speicher-Layout von Prozessen

Der Systemkern hält in einer *Benutzerstruktur* weitere charakteristische Informationen für jeden Prozeß, u. a.:

- Registerinhalte die bei einem Systemaufruf dort abgelegt werden,
- Parameter und Ergebnisse des aktuellen Systemaufrufs,
- die Tabelle der Dateideskriptoren (siehe Abschnitt 2.2) sowie

- Abrechnungsdaten (verbrauchte CPU-Zeit im Benutzer und Systemmodus, Systembeschränkungen, wie z. B. maximale CPU-Zeit, die maximale Stackgröße, etc.)

Bei der *Umgebung* eines Prozesses, also bei den genannten Attributen der Benutzerstruktur sowie den Code- und Datenbereichen, handelt es sich um geschützte Ressourcen. D. h. das Betriebssystem schottet diese Bereiche vor dem lesenden und schreibenden Zugriff durch andere Prozesse ab.

Unter dem *Kontext* eines Prozesses versteht man die aktuellen Registerwerte des ausführenden Prozessors, dazu gehören insbesondere der Befehlszähler und der Zeiger auf den Stack.

Ein Prozeßwechsel auf einem Unix-System bedingt demzufolge immer einen Umgebungswechsel *und* einen Kontextwechsel. In Kapitel 3 kommen wir nochmals auf die Bedeutung von Umgebungs- und Kontextwechsel zurück.

2.1.5 Lebenszyklus

Sobald ein neuer Prozeß gestartet wird, übergibt der Systemkern die weitere Ausführung des Programms über eine Startup-Routine der Laufzeitumgebung an die Hauptfunktion des Programms. Der sogenannten `main()`-Funktion werden dabei die Kommandozeilen-Argumente übergeben:

```
int main( int argc, char *argv[] );
```

Der neue Prozeß hat nun – nachdem er seine Aufgaben hoffentlich erfolgreich bewältigen konnte – drei verschiedene Möglichkeiten, sich selbst mit Anstand zu beenden:

1. Er verläßt mit `return` die `main()`-Funktion,
2. er ruft die Funktion `exit()` auf oder
3. er beendet sich durch Aufruf der Funktion `_exit()`.

Verläßt ein Prozeß die `main()`-Funktion durch `return`, so liegt die Kontrolle wieder bei der Startup-Routine, die nun ihrerseits die `exit()`-Funktion anspringt. Abbildung 2.4 veranschaulicht diesen Lebenszyklus. Darüber hinaus kann ein Prozeß von außerhalb – und daher meist unfreiwillig – durch ein Signal beendet werden. Eine letzte Möglichkeit: Der Prozeß beendet sich selbst durch ein Signal, das er durch Aufruf von `abort()` selbst auslöst. Die letzten beiden Varianten sind in Abb. 2.4 nicht dargestellt.

Die Funktionen `exit()` und `_exit()` sind durch ANSI/ISO C bzw. IEEE Std 1003.1-2001 spezifiziert. Der Standard für ANSI/ISO C legt fest, daß

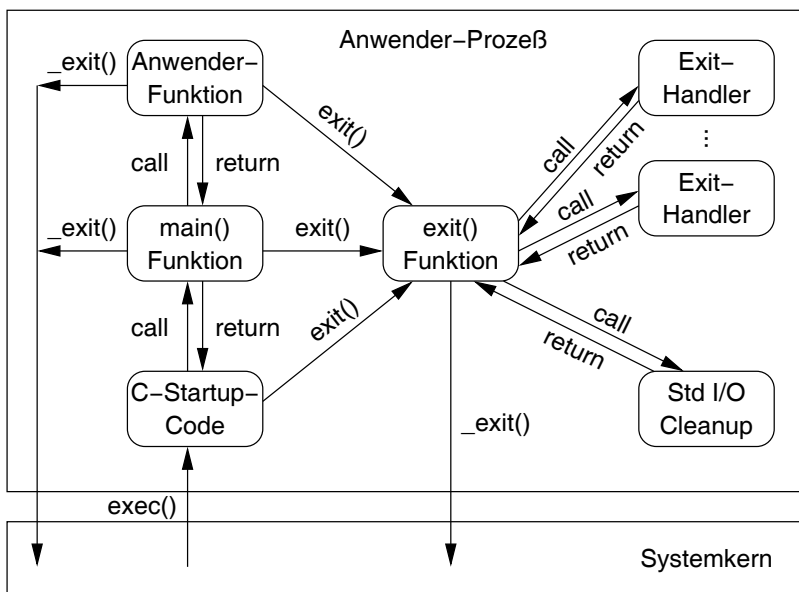


Abb. 2.4. Lebenszyklus von Prozessen

`exit()` die Puffer aller noch offenen Dateien schreibt, offene Datenströme mit `fclose()` abschließt, alle durch `tmpfile()` erzeugten temporären Dateien löscht, und die Kontrolle an das System zurück gibt. Ganz zu Beginn der `exit()`-Funktion müssen zudem noch alle mit `atexit()` hinterlegten *Exit-Handler* aufgerufen werden. Dies erfolgt in der umgekehrten Reihenfolge ihrer Registrierung mit `atexit()`.

```
#include <stdlib.h>

void exit( int status );
```

Nach ANSI/ISO C bedeutet der Statuswert Null bzw. `EXIT_SUCCESS` ein erfolgreiches Programmende. Hat `status` den Wert `EXIT_FAILURE`, wird dadurch ein nicht erfolgreiches Programmende angezeigt. Wie genau der aufrufenden Umgebung Erfolg oder Mißerfolg mitzuteilen ist und wie mit anderen Werten zu verfahren ist, überläßt ANSI/ISO C der jeweiligen Implementierung. Für Unix wird deshalb durch IEEE Std 1003.1-2001 festgelegt, daß der Status in Form der niederwertigen acht Bits des Statuswerts (d. h. der Wert `status & 0377`) zurück geliefert wird.

Genau genommen gibt die `exit()`-Funktion unter Unix den Rückgabewert und die Kontrolle nicht direkt an das System zurück, sondern ruft zu diesem

Zweck am Ende die Funktion `_exit()` auf, welche dann die Terminierung des Prozesses ordnungsgemäß abschließt.

```
#include <unistd.h>

void _exit( int status );
```

Die durch IEEE Std 1003.1-2001 festgelegte Funktion `_exit()` schließt u. a. alle noch offenen Dateideskriptoren (vgl. dazu Abschnitt 2.2) des aktuellen Prozesses und informiert gegebenenfalls den Elternprozeß über das „Ableben“ und den Rückgabe- bzw. Statuswert seines Kindes. Wir kommen in den Abschnitten 2.5 und 2.6 noch ausführlicher auf die Effekte und die z. T. damit verbundenen Probleme zu sprechen, die ein Prozeßende auf andere Prozesse wie Eltern- und Kindprozesse haben kann.

Handelt es sich bei dem terminierenden Prozeß um einen kontrollierenden Prozeß, also um den Anführer seiner Session, der die Verbindung zum kontrollierenden Terminal hergestellt hat, so wird an *alle* Prozesse in der Vordergrund-Prozeßgruppe (vgl. Abschnitt 2.1.2) des kontrollierenden Terminals das `SIGHUP`-Signal ausgeliefert. Außerdem wird die Assoziation zwischen der Session und dem kontrollierenden Terminal aufgelöst.

Verwaist durch den terminierenden Prozeß eine Prozeßgruppe und ist einer der Prozesse dieser Gruppe gestoppt, dann wird an alle Prozesse dieser Prozeßgruppe das `SIGHUP`-Signal, gefolgt vom `SIGCONT`-Signal ausgeliefert. Auch dazu später noch mehr.

```
#include <stdlib.h>

int atexit( void (*func)( void ) );
```

Die Funktion `atexit()` hinterlegt die benutzerdefinierte Funktion `func()`, die dann ausgeführt wird, wenn das Programm normal endet. Ein Rückgabewert ungleich Null signalisiert, daß der angegebene Exit-Handler nicht registriert werden konnte.

Exit-Handler dienen im Allgemeinen dazu, unmittelbar vor Programmende noch spezielle Aufräumarbeiten durchzuführen, die durch die `exit()`-Funktion nicht abgedeckt werden. Exit-Handler können sogar neue Exit-Handler hinterlegen. Es ist jedoch streng darauf zu achten, daß alle registrierten Funktionen auch zurückkehren.

Das nachfolgende Programm zeigt die Funktionalität und die Unterschiede der beiden Exit-Funktionen auf.

- 1-4 Zunächst binden wir die notwendigen Headerfiles in unser Programm ein. Das Headerfile `<stdlib.h>` enthält unter anderem die Prototypen für die in