

ÜBERSETZERBAU VORLESUNGSSKRIPTUM

Manfred Brockhaus

Anton Ertl

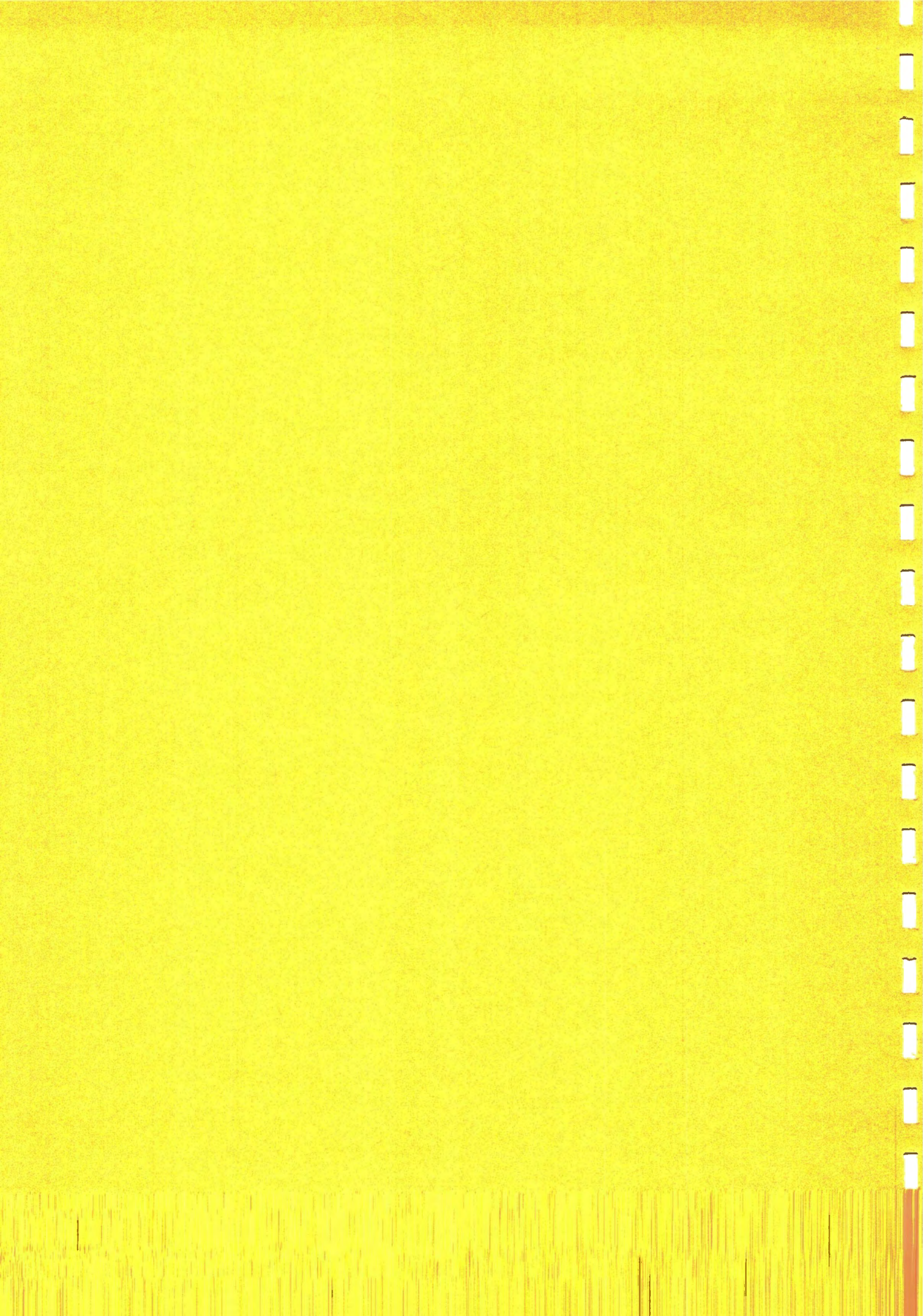
Andreas Krall

Institut für Computersprachen

Arbeitsbereich für Programmiersprachen und Übersetzerbau

Technische Universität Wien

©2014



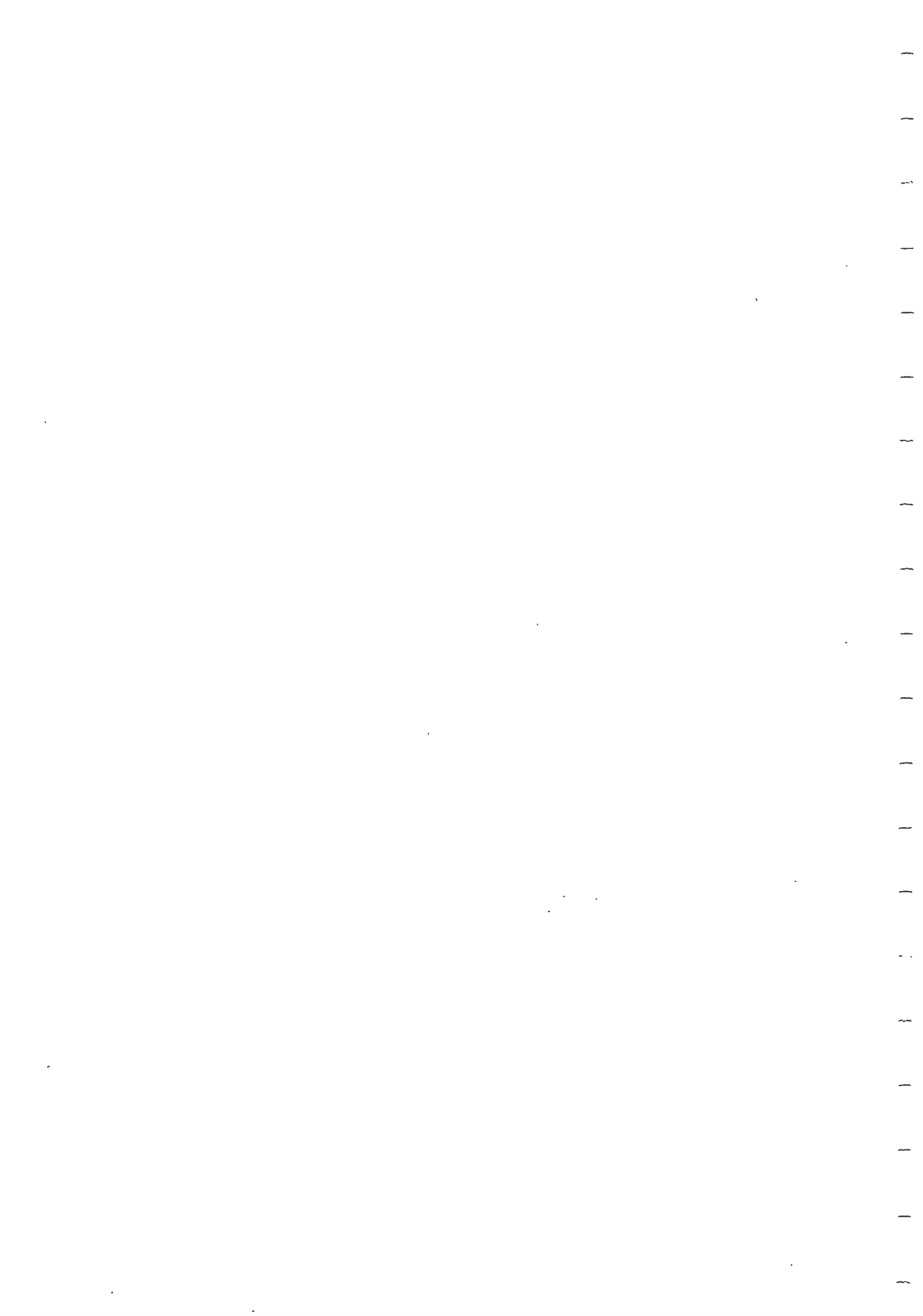
Vorwort

Dieses Skriptum behandelt grundlegende Begriffe und Methoden des Compilerbaus (Übersetzerbau). Der Inhalt entspricht dem Stoff der Vorlesung. Wir setzen Kenntnisse in höheren Programmiersprachen voraus. Grundlegende Kenntnisse über Automaten und Formale Sprachen sind für das Verständnis von Vorteil.

Der Compilerbau ist ein klassisches Gebiet der Informatik. Er gehört seit den Anfängen der Informatik bzw. Computer Science zu den Kernfächern und hat seine Aktualität bis heute nicht verloren. Dafür sind unter anderem folgende Gründe maßgebend:

- Compiler verarbeiten formale Sprachen und können daher weitgehend formal spezifiziert werden. Dabei lassen sich Ergebnisse aus der Theorie der Automaten und Formalen Sprachen praktisch anwenden. Dies hat insbesondere zur Entwicklung von Compiler-Generatoren geführt. Die Methoden und Werkzeuge des Compilerbaus sind aber nicht nur zur Konstruktion von Compilern, sondern auch zur Implementierung von Editoren, Textsystemen, Netzwerkprotokollen etc. anwendbar. Ganz allgemein eignen sie sich zur Entwicklung von Programmen, deren Operationen stark an die Struktur der Eingabedaten gebunden sind. Solche Programme können analog zu Compilern formal spezifiziert und generiert werden, indem man die Eingabe als Sprache und die Verarbeitung als Semantik der Sprache auffaßt.
- Ein Compiler ist ein Programm, das über andere Programme „nachdenkt“, d.h. diese analysiert und transformiert. Die Methoden des Compilerbaus kann man daher auch einsetzen, um universelle Programme in spezialisierte und effizientere Programme zu transformieren, sofern ein Teil der Eingabedaten bekannt ist. Diese Technik der „partiellen Auswertung“ ist im Bereich der funktionalen und logischen Sprachen von aktueller Bedeutung.
- Bei neuen Rechnerarchitekturen geht man zunehmend davon aus, daß der Compiler die Erzeugung des Maschinencodes genau auf die Anforderungen der Hardware abstimmt, und zwar nicht nur im Hinblick auf die Effizienz, sondern auch im Hinblick auf den korrekten zeitlichen (parallelen) Ablauf der Befehle. Es handelt sich also um integrierte Hardware/Software-Entwicklungen mit aktuellen Anforderungen an den Compilerbau.

Das Skriptum ist wie folgt aufgebaut: Im Kapitel 1 werden die Grundbegriffe Compiler und Interpreter definiert. Kapitel 2 bietet eine kurze Einführung in



Computerarchitektur und Assembler. Kapitel 3 gibt einen Überblick über die Teilaufgaben eines Compilers und zeigt an einem einfachen Beispiel, wie ein Programm schrittweise übersetzt wird. Die restlichen Kapitel behandeln die einzelnen Teilaufgaben des Compilers. Jedes Kapitel behandelt eine Teilaufgabe, wobei die Folge der Kapitel der logischen Folge der Teilaufgaben entspricht.

Algorithmen werden in einer Java-ähnlichen Notation angegeben. Übergeordnete Klassendefinition oder Befehle zum Importieren von *Packages* werden weggelassen.

Eine umfassende Behandlung der globalen Optimierung und der automatisierten Code-Erzeugung (Code-Generator-Generatoren) ist im gegebenen Rahmen nicht möglich. Wir verweisen auf entsprechende Speziallehrveranstaltungen und die wissenschaftliche Literatur.

Der Anhang des Skriptums enthält zwei verschiedene Implementationen von Übersetzer-Interpreter-Systemen für die Programmiersprache Mini, eine kleine Teilmenge von Java/C. Sie demonstrieren typische Formen von Zwischensprachen und Interpretern. Die erste Version ist in Java programmiert. Als Zwischensprache wird ein Postfix-Code verwendet, der einer Teilmenge des standardisierten Java VM Codes entspricht. Der Interpreter simuliert eine Stackmaschine und entspricht einer Java VM. Die zweite Version ist in Prolog mit DCGs programmiert. Als Zwischensprache wird die abstrakte Syntax verwendet und der Interpreter arbeitet rekursiv.

Für die Mitarbeit an diesem Skriptum danken wir Andreas Falkner, Thomas Berger, Alexander Forst-Rakoczy und Ulrich Neumerkel.

Inhaltsverzeichnis

| | | |
|----------|-----------------------------------------------------|-----------|
| 1 | Grundlagen | 7 |
| 1.1 | Compiler | 7 |
| 1.2 | Interpreter | 9 |
| 1.3 | Compiler-Interpreter | 10 |
| 2 | Computerarchitektur und Assembler | 12 |
| 2.1 | Speicher | 12 |
| 2.2 | Datenstrukturen | 13 |
| 2.3 | Register | 14 |
| 2.4 | Maschinencode und Assemblercode | 15 |
| 2.5 | Befehle | 15 |
| 2.6 | Funktionen | 19 |
| 2.7 | Was dieses Kapitel verschweigt | 21 |
| 3 | Struktur von Compilern | 22 |
| 3.1 | Phasen | 22 |
| 3.2 | Ablaufsteuerung | 28 |
| 4 | Lexikalische Analyse | 31 |
| 4.1 | Grundsymbole | 31 |
| 4.2 | Ausgabe-Aktionen | 37 |
| 4.3 | Analyse-Verfahren | 39 |
| 5 | Syntax-Analyse | 45 |
| 5.1 | Grundlagen | 45 |
| 5.2 | Top-Down-Analyse | 50 |
| 5.3 | Bottom-Up-Analyse | 58 |
| 5.4 | Hierarchie der Analyse-Verfahren | 68 |
| 5.5 | Wie schreibt man Grammatiken? | 70 |
| 5.6 | Qualitätskriterien für Grammatiken | 72 |
| 6 | Syntaxgesteuerte Übersetzung | 73 |
| 6.1 | Attributierte Grammatik | 73 |
| 6.2 | Wie attribuiert man Grammatiken? | 74 |
| 6.3 | Qualitätskriterien für die Attributierung | 78 |
| 6.4 | S-Attributierte Grammatik | 78 |

| | | |
|-----------|------------------------------------------------|------------|
| 6.5 | L-Attributierte Grammatik | 81 |
| 6.6 | Übersetzungsschema | 82 |
| 6.7 | Top-Down-Compiler-Generator | 82 |
| 6.8 | Bottom-Up-Compiler-Generator | 83 |
| 6.9 | Ein-Pass-Übersetzung | 88 |
| 6.10 | Mehrpas-Übersetzung | 89 |
| 7 | Semantische Analyse | 92 |
| 7.1 | Typen | 93 |
| 7.2 | Symboltabelle | 95 |
| 7.3 | Typ-Überprüfung | 97 |
| 7.4 | Typ-Konversion | 98 |
| 7.5 | Overloading | 99 |
| 8 | Zwischendarstellungen | 103 |
| 8.1 | Syntaxbaum | 103 |
| 8.2 | Postfix und Prefix | 104 |
| 8.3 | Quadrupel | 105 |
| 8.4 | Zuweisungen | 106 |
| 8.5 | Indizierte Variablen | 107 |
| 8.6 | Boolesche Ausdrücke | 108 |
| 8.7 | Kontrollanweisungen | 109 |
| 8.8 | Kontrollflußgraph und Datenflußgraph | 111 |
| 9 | Codeerzeugung | 116 |
| 9.1 | Befehlsauswahl | 117 |
| 9.2 | Befehlsanordnung | 122 |
| 9.3 | Registerbelegung | 126 |
| 10 | Laufzeitsystem | 131 |
| 10.1 | Speicherorganisation | 132 |
| 10.2 | Stack-Verwaltung | 133 |
| 10.3 | Nicht-lokale Variablen | 136 |
| 10.4 | Parameterübergabe | 139 |
| 10.5 | Dynamische Objekte | 139 |
| 10.6 | Heap-Verwaltung | 141 |
| 11 | Optimierungen | 147 |
| 11.1 | Analysen | 147 |
| 11.2 | Skalare Optimierungen | 149 |
| 11.3 | Codeoptimierungen | 149 |
| 11.4 | Array- und Schleifenoptimierungen | 150 |

| | |
|------------------------------------------------------|------------|
| 12 Übersetzung objektorientierter Konzepte | 156 |
| 12.1 Klassendarstellung und Methodenaufruf | 156 |
| 12.2 Dynamische Typüberprüfung | 162 |
| 12.3 Devirtualisierung | 166 |
| 12.4 Escapeanalyse | 167 |
| A Mini-Compiler in Java | 171 |
| B Mini-Compiler in Prolog | 182 |

Kapitel 1

Grundlagen

1.1 Compiler

Ein **Compiler** (*translator*) oder **Übersetzer** ist ein Programm, das Programme einer **Quellsprache** (*source language*) in semantisch äquivalente Programme einer **Zielsprache** (*target language*) umwandelt.



Abbildung 1.1: Blockdiagramm eines Compilers

Ein Compiler als Programm ist in einer Sprache formuliert. Diese Sprache nennen wir **Implementierungssprache** oder **Systemsprache**. Sie ist unabhängig von Quell- und Zielsprache. Mit einem **T-Diagramm** kann eine Übersetzung dargestellt werden. Abb. 1.2 beschreibt die Übersetzung eines Programms mit Bedeutung P von Q nach Z durch einen Compiler mit der Implementierungssprache S .

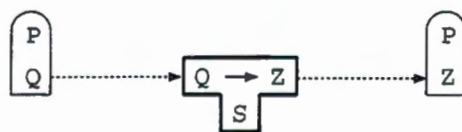


Abbildung 1.2: T-Diagramm eines Compilers

Als Implementierungssprache S eignet sich jede universelle (niedere oder höhere) Programmiersprache. Ein Compiler kann auf jedem Computer verwendet werden, der seine Implementierungssprache 'versteht'.

Beispiel 1.3: Gegeben sei ein Computer (eine Maschine) M , auf dem ein C-Compiler verfügbar ist. Gesucht sei ein Ada-Compiler für die Maschine M . Wir schreiben den

Ada-Compiler in C und übersetzen ihn wie ein gewöhnliches C-Programm mit dem verfügbaren C-Compiler. Danach ist der Ada-Compiler auf M verwendbar, um Ada-Programme zu übersetzen. Die T-Diagramme in Abb. 1.4 stellen den Zusammenhang bildlich dar.

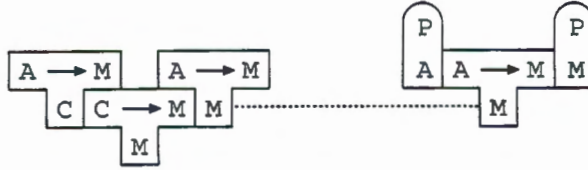


Abbildung 1.4: Übersetzung eines Compilers

Es wäre auch möglich, den Ada-Compiler von Bsp. 1.3 so zu schreiben, daß er nicht M, sondern C als Zielsprache hat. Dann müßte aber jedes Ada-Programm in zwei Stufen übersetzt werden ($A \rightarrow C \rightarrow M$).

Wenn ein Compiler in seiner Quellsprache Q geschrieben wird ($Q = S$), so kann er *sich selbst* übersetzen. Diese Möglichkeit läßt sich ausnutzen, um in ein oder mehreren Schritten modifizierte Versionen eines Compilers zu erstellen. Das Prinzip wird als **Bootstrapping** bezeichnet. Das folgende Beispiel zeigt eine typische Anwendung.

Beispiel 1.5: Gegeben sei ein Ada-Compiler auf einer Maschine M. Gesucht sei ein Ada-Compiler für eine andere Maschine N. Der Compiler $A \rightarrow N$ wird in Ada geschrieben und auf der Maschine M mit dem dort vorhandenen Compiler übersetzt. Er kann auf M bereits als Cross-Compiler verwendet werden, um Ada-Programme für die Maschine N zu übersetzen. Er kann aber auch sich selbst nach N übersetzen und dann auf der Maschine N verwendet werden (Abb. 1.6).



Abbildung 1.6: Selbst-Übersetzung eines Compilers

Der Compiler $A \rightarrow N$ mit der Implementierungssprache N wird erzeugt, ohne daß dazu die Maschine N selbst benötigt wird. N kann also eine beliebige, insbesondere auch eine abstrakte Maschine sein. Darauf wird im Kapitel 1.2 näher eingegangen. Die Erstellung von Compilern wird durch die Verwendung von höheren Programmiersprachen erleichtert. Eine noch weitergehende Unterstützung bieten **Compiler-Generatoren**. Compiler-Generatoren sind Werkzeuge zur Konstruktion von Compilern, die den Konstrukteur von Routinearbeit befreien.

Die Eingabe für einen Compilergenerator besteht aus einer formalen Definition der gewünschten Übersetzung $X \rightarrow Y$. Der Generator erzeugt aus der Definition einen Compiler, indem er die Definition in ein ablauffähiges Programm übersetzt. Das Prinzip wird in Abb. 1.7 durch T-Diagramme dargestellt.

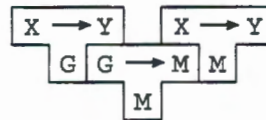


Abbildung 1.7: Compiler-Generator

Bei Verwendung eines Compiler-Generators wird der Compiler in einer **Generatorsprache G** geschrieben bzw. definiert. Das theoretische Modell für praktische Generatorsprachen ist die **Attributierte Grammatik**. Dabei wird die Syntax der zu übersetzenden Sprache X durch eine kontextfreie Grammatik definiert und die Übersetzung, d.h. die Zuordnung $X \rightarrow Y$ wird als Berechnung von Attributen der Grammatiksymbole spezifiziert (siehe Kapitel 6).

Besonders bekannt in der Praxis ist der Compiler-Generator **YACC** unter Unix. Die Generatorsprache ist hier eine Mischung bestehend aus einer speziellen Notation der Grammatik und aus der Programmiersprache C zur Berechnung der Attribute.

Die Logik-orientierte Programmiersprache **PROLOG** ist eine Generatorsprache, in der man sowohl die Syntax als auch die Attributberechnung, d.h. die Übersetzung einheitlich beschreiben kann.

1.2 Interpreter

Ein **Interpreter** ist ein Programm, das Programme ausführen kann, die in seiner **Interpretersprache I** formuliert sind. Ein Interpreter verhält sich wie eine Maschine bzw. ein Prozessor. Man bezeichnet ihn daher auch als **virtuelle** oder **abstrakte Maschine**.

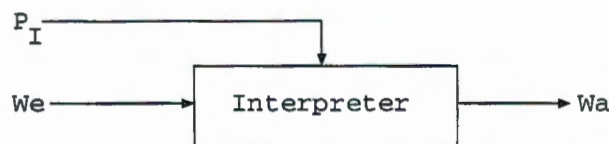


Abbildung 1.8: Blockdiagramm eines Interpreters

Der Interpreter (Abb. 1.8) führt das Programm P_I aus. Bei seiner Ausführung benötigt P_I im allgemeinen Eingangswerte W_e und liefert Ausgangswerte W_a .

Zum Vergleich von Compiler und Interpreter betrachten wir ein Quellprogramm, das eine lineare Folge von k Schritten beschreibt. Jeder Programmschritt i erfordert einen Decodierungsschritt D_i und einen Ausführungsschritt A_i .

- Bei der **Übersetzung** ergibt sich der Ablauf D_1, D_2, \dots, D_k , Optimieren, Assemblieren, Linken, A_1, A_2, \dots, A_k .
- Bei der **Interpretation** ergibt sich der Ablauf $D_1, A_1, D_2, A_2, \dots, D_k, A_k$.

Bei einem linearen Programm (jeder Schritt wird genau einmal ausgeführt) ist der gesamte Zeitbedarf bei beiden Verfahren etwa gleich. Falls das Programm eine Schleife enthält (Schritt i wird n -mal ausgeführt), so ist der Zeitbedarf beim Interpreter um $(n - 1) * D_i$ größer als beim Compiler. Falls das Programm eine Auswahl enthält (Schritt i besteht aus n Alternativen), so ist der Zeitbedarf beim Interpreter um $(n - 1) * D_i$ kleiner als beim Compiler.

Interpretation ist zweckmäßig, wenn ein Programm oft geändert wird, z.B. beim Debugging.

Übersetzung ist generell zweckmäßig, wenn ein Programm mehrmals ausgeführt werden soll.

Übersetzung ist notwendig, wenn die statische Korrektheit eines Programms vor der Ausführung überprüft werden soll.

Interpretation ist notwendig, wenn das Programm selbst während der Laufzeit manipuliert werden soll. Dynamische Übersetzung (*just-in-time compilation*) ist eine Alternative, wenn sowohl Programmmanipulationen als auch Effizienz gefordert sind.

1.3 Compiler-Interpreter

Oft werden Compiler und Interpreter kombiniert. Ein einfaches Beispiel ist ein Compiler, der Programme für einen Interpreter (abstrakte Maschine) erzeugt (siehe Abb. 1.9).

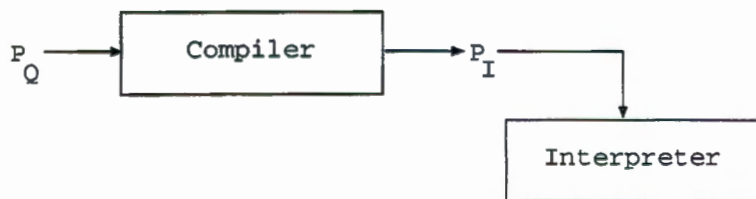


Abbildung 1.9: Blockdiagramm eines Compiler-Interpreter Systems

Diese Kombination kann man z.B. verwenden, um eine Quellsprache Q auf einer Maschine M zu implementieren, ohne daß man die Maschinenbefehle von M kennen und beachten muß. Als Zielsprache des Compilers wird eine Interpretersprache I

entwickelt und ein Interpreter dafür geschrieben. Zweckmäßig wählt man für I eine Sprache, die einfach und effizient decodiert werden kann (z.B. Postfix-Code). Der Compiler und der Interpreter können in einer beliebigen auf M verfügbaren Sprache X geschrieben werden. Abb. 1.10 zeigt ein solches System, wobei neben dem T-Diagramm auch ein I-Diagramm (für Interpreter) verwendet wird.

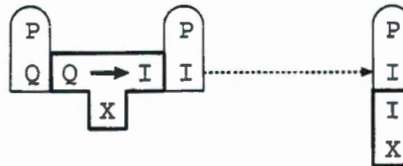


Abbildung 1.10: Implementierung von Q mittels Compiler-Interpreter

Der Nachteil dieser Implementierung ist, daß die erzeugten Programme durch den Interpreter langsamer als durch einen Hardware-Prozessor ausgeführt werden. Ein Vorteil liegt aber z.B. darin, daß das Zielprogramm in der Sprache I wesentlich kompakter dargestellt werden kann als in einem realen Maschinencode und daß die Übersetzung meist wesentlich schneller ist. Die Beispiele im Anhang sind Compiler-Interpreter Systeme nach dem Muster von Abb. 1.10. Eine Compiler-Interpreter Kombination kann man insbesondere verwenden, um die **Portabilität** eines Compilers zu erreichen. Dazu muß der Compiler $Q \rightarrow I$ in seiner Zielsprache I dargestellt werden. Er muß deswegen nicht in I geschrieben werden, sondern kann – falls er in der Quellsprache Q geschrieben wird – nach dem Prinzip von Abb. 1.6 durch Selbst-Übersetzung nach I übersetzt werden. Danach ist der Compiler auf allen Maschinen einsetzbar, auf denen die abstrakte Maschine I als Interpreter implementiert ist. In Abb. 1.11 wird der Einsatz des Systems auf einer Maschine X dargestellt.

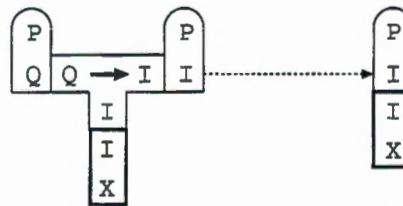


Abbildung 1.11: Portabler Compiler mittels Interpreter

Um das System auf einen neuen Computer X zu übertragen, muß also lediglich der Interpreter neu geschrieben werden. Der Nachteil dieses Verfahrens ist aber, daß der Compiler *und* die erzeugten Programme interpretiert werden (Bsp. UCSD-Pascal).

Das aktuelle Beispiel für Compiler-Interpreter Systeme ist Java. Der Zwischencode (Java VM Code) bringt nicht nur den Vorteil der Maschinenunabhängigkeit, sondern er verkürzt durch seine Kompaktheit auch die Übertragungszeiten im Internet.

Kapitel 2

Computerarchitektur und Assembler

Viele Compiler erzeugen (oft im Zusammenspiel mit Assembler und Linker) ausführbaren Maschinencode für eine bestimmte Computerarchitektur. Dieses Kapitel erklärt Computerarchitektur aus Sicht des Compilers bzw. eines User-level-Programms in Maschinensprache.

2.1 Speicher

Daten jedweder Art liegen im Speicher. Der Speicher besteht aus einer Menge von 8-bit-Bytes; jedes dieser Bytes hat eine eigene Adresse (eine 32-bit- oder 64-bit-Zahl), über die auf genau dieses Byte zugegriffen werden kann. Die Bytes des Speichers liegen nebeneinander zusammenhängend im Adressraum (siehe Abb. 2.1).

Der Prozessor kann auf einzelne Bytes zugreifen, aber auch auf Gruppen von zwei, vier, oder acht Bytes. Allerdings muss die Adresse dabei durch 2, 4, bzw. 8 teilbar sein (Ausrichtung, alignment), da der Zugriff sonst eine Exception auslöst oder zumindest viel langsamer durchgeführt wird als ein ausgerichteter Zugriff (je nach Architektur).

Was der Inhalt des Speichers bedeutet, wird nicht mitgespeichert, sondern ist allein eine Frage der Interpretation durch das Programm; häufige (und von vielen Architekturen direkt unterstützte) Interpretationen sind:

| Bytes | Interpretation |
|-------|--------------------------------------------------------------------------------------------------------------------------------------------------------|
| 1 | Zeichen, ganze Zahl im Bereich $-128...127$, natürliche Zahl < 256 , 1 Flag, 8 Flags, Teil eines IA32-Maschinenbefehls |
| 2 | UTF-16-Zeichen, ganze Zahl im Bereich $-2^{15}...2^{15} - 1$, natürliche Zahl $< 2^{16}$ |
| 4 | ganze Zahl im Bereich $-2^{31}...2^{31} - 1$, natürliche Zahl $< 2^{32}$, 32-bit-Adresse, RISC-Maschinenbefehl, IEEE single-precision Gleitkommazahl |
| 8 | ganze Zahl im Bereich $-2^{63}...2^{63} - 1$, natürliche Zahl $< 2^{64}$, 64-bit-Adresse, IEEE double-precision Gleitkommazahl |

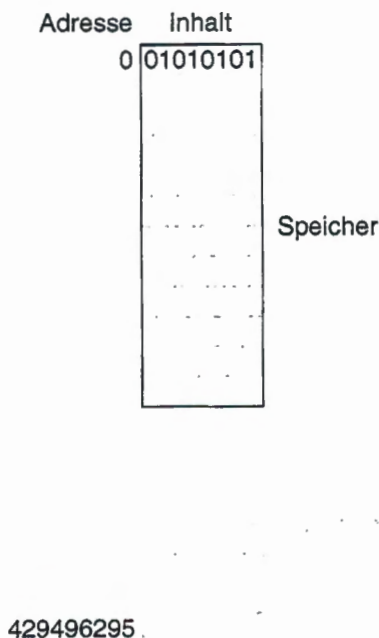


Abbildung 2.1: Adressraum und Speicher auf einer 32-bit-Maschine

So kann man zum Beispiel das Byte 11011111 als Zeichen “ß” (im Zeichensatz ISO-8859-1), als ganze Zahl -33 (in der Zweierkomplementdarstellung) oder als natürliche Zahl 223 interpretieren; andere Interpretationen durch die Software sind auch möglich, oder das Byte könnte Teil eines Datums mit mehr als einem Byte sein.

2.2 Datenstrukturen

Bei Strukturen aus mehreren Feldern liegen die einzelnen Felder direkt hintereinander im Speicher (siehe Abb. 2.2); allerdings ist oft ein Zwischenraum (padding) nötig (in unserem Beispiel nach Feld c), damit das nächste Feld (in unserem Beispiel m) ausgerichtet ist. Die Adressen der einzelnen Felder berechnen sich durch addieren eines Offsets zur Anfangsadresse der Struktur. Bei Arrays liegen die einzelnen Elemente hintereinander im Speicher (siehe Abb. 2.2). Die Adresse eines bestimmten Elements ist $a + s * i$, wobei a die Startadresse des Arrays ist, s die Größe eines Elements, und i der Index des Elements (angefangen mit 0).

Man kann diese Datenstrukturen miteinander verschachteln. Bei der Adressberechnung wird dabei zuerst die äussere Datenstruktur aufgelöst und dann sukzessive die inneren. In unserem Beispiel berechnen wir die Adresse von $a[2].m$, indem wir zuerst die Adresse von $a[2]$ berechnen (Ergebnis: $a+12$), und dann den Offset 2 für m dazuzählen; Gesamtergebnis: $a+14$ (die Anfangsadresse a hängt von der konkreten Plattform (Maschine, Compiler, Betriebssystem) ab, daher geben wir sie hier nur

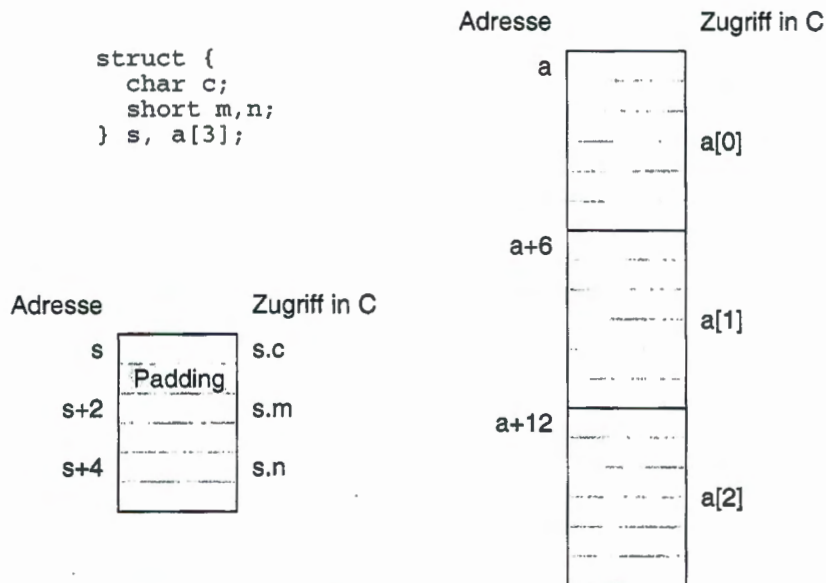


Abbildung 2.2: Repräsentation einer Struktur und eines Arrays (von Strukturen) im Speicher

symbolisch an.

Diese direkte Art der Verschachtelung funktioniert aber nur, wenn die Felder bzw. Elemente nicht größer sein können als der vorgesehene Platz. In den anderen Fällen (z.B. bei Objekten; Instanzen von Unterklassen können beliebig groß werden) legt man die eigentlichen Daten des Feldes/Elements woanders ab, und reserviert in der Datenstruktur nur Platz für die Adresse dieser Daten (alias Zeiger auf diese Daten). In C muß der Programmierer das explizit angeben, in vielen anderen Sprachen wird das aber automatisch gemacht.

2.3 Register

Neben dem Speicher kann der Prozessor auch noch auf Register zugreifen. Z.B. hat die AMD64-Architektur 16 Integerregister (*rax-rdx, rsp, rbp, rsi, rdi, r8-r15*), jeweils 64 bit breit und 16 SSE-Register (*xmm0-xmm15*), die 128 bit breit sind.

Register dienen dazu, Operanden von Befehlen zwischenspeichern, denn viele Befehle arbeiten nur mit Registern, oder nicht in jeder Form mit Operanden im Speicher, und Registerzugriffe sind jedenfalls schneller als Speicherzugriffe. Eine weitere Verwendung für Register ist daher auch, häufig verwendete Daten zu speichern, um schnelle Zugriffe zu erlauben.

Unterschied zum Speicher (indirekte Adressierung)

Ein wichtiger Unterschied zwischen Registern und Speicher ist, dass die Nummer des Registers immer fix im Befehl steht. Sie kann nicht zur Laufzeit geändert werden.

Beim Speicher dagegen kann die Adresse der Speicherstelle nicht nur direkt im Befehl angegeben werden (direkte Adressierung), sondern es kann auch z.B. auf eine Adresse zugegriffen werden, die in einem Register steht, und dort zur Laufzeit berechnet wurde (indirekte Adressierung).

2.4 Maschinencode und Assemblercode

Ausführbare Befehle der Architektur sind als bestimmte Bitmuster codiert; z.B. repräsentiert die Byte-Sequenz 01001000 10000011 11000111 00001000 einen bestimmten Befehl der AMD64-Architektur. Diese binäre Maschinensprache ist schwer zu lesen und zu schreiben; daher gibt es Disassembler und Assmblers, die zwischen Maschinensprache und einer lesbaren Repräsentation (Assemblersprache) hin- und herübersetzen. Die Assembler-Repräsentation des obigen Befehls ist `addq $8, %rdi` und erhöht das Register `rdi` um den Wert 8.

Assemblerbefehle bestehen aus einem Mnemonic (`addq`), das die Art der Operation angibt (addiere quadwords), und darauffolgenden Parametern (im Beispiel den Wert 8 (`$8`) und das Register `rdi` (`%rdi`)). Im allgemeinen entspricht ein Assemblerbefehl einem Maschinenbefehl.

Zusätzlich zu den Befehlen verstehen Assembler auch noch Assembleranweisungen, um z.B. Daten in Binärform abzulegen oder für diverse Verwaltungsaufgaben, die die Programmierung erleichtern oder die Wartbarkeit erhöhen.

2.5 Befehle

Verschiedene Architekturen haben verschiedene Befehle. Allerdings kann man die Befehle als Kombinationen von Grundoperationen ansehen, die zum Großteil für alle Architekturen gleich sind. Im folgenden betrachten wir als Beispiel einige Befehle der AMD64-Architektur. Die AMD64-Architektur ist eine 2-Adressarchitektur, d.h. jeder Befehl hat 2 Parameter, die die Operanden des Befehls beschreiben. Ein Befehl der Form

- `operation source, source_dest`

führt die Operation `source_dest = source_dest operation source` durch. Der zweite Parameter `source_dest` wird meist zweimal verwendet, einmal als Quell- und einmal als Zieloperand. Die Operanden können entweder Konstante, Register oder Werte im Speicher sein. Es gibt Einschränkungen, so kann z.B. eine Konstante nur ein Quelloperand sein und maximal ein Operand darf sich im Speicher befinden.

Ein einfacher, häufig gebrauchter Befehl ist `mov source, dest`, der einfach den Inhalt des Quelloperanden in den Zieloperanden kopiert. Ein einfaches Beispiel ist `mov %r8, %r9`, der den Inhalt des Registers `r8` in das Register `r9` schreibt.

Speicherzugriffe

Statt eines Registers kann man als Operand auch eine Speicherstelle angeben, z.B.:

```

mov 456, %r9
mov %r8, 456
mov (%r8), %r9
mov %r8, (%r9)

```

Der erste Befehl liest einen 8-Byte-Wert aus dem Speicher, der dort bei der Adresse 456 anfängt (und die folgenden Bytes bis inklusive 463 belegt), und kopiert ihn in das Register r9. Der zweite Befehl liest umgekehrt den Inhalt des Registers r8 und schreibt ihn in den Speicher bei Adresse 456(-463). Da hier die Adresse direkt angegeben ist, spricht man von direkter Adressierung. Sie spielt bei Datenzugriffen keine große Rolle, dafür aber beim Kontrollfluss. Normalerweise wird im Assemblercode die Adresse symbolisch (als Labelname) angegeben und nicht als Zahl.

Der dritte Befehl liest wie der erste einen 8-Byte-Wert aus dem Speicher und schreibt ihn in das Register r9. Allerdings steht hier die Adresse der Speicherstelle in Register r8; wenn in r8 z.B. 456 steht, dann ist dieser Befehl äquivalent zum ersten. Hier wird die Speicheradresse indirekt angegeben (in diesem Fall über ein Register), daher spricht man von indirekter Adressierung.

Der vierte Befehl liest wie der zweite den Inhalt des Registers r8 und schreibt ihn in den Speicher, und zwar bei der Adresse, die in r9 steht. Wenn in r9 also 456 steht, dann ist dieser Befehl äquivalent zum zweiten Befehl. Beachten Sie, dass dieser Befehl den Inhalt von r9 nicht ändert, sondern nur den Inhalt des Speichers.

Die AMD64-Architektur bietet noch weitere Möglichkeiten, die Speicheradresse anzugeben (Adressierungsarten); die allgemeinste Variante ist $o(r_1, r_2, f)$, wobei o eine Konstante ist, r_1 und r_2 Register, und f der Wert 1, 2, 4 oder 8; diese Adressierungsart berechnet folgende Speicheradresse: $o + r_1 + r_2 f$. Die meisten anderen Adressierungsarten lassen sich aus dieser ableiten, indem man verschiedene Teile weglässt. Beispiele:

```

mov 400(,%r8,8),%r9

```

Wenn r8 den Wert 7 enthält, liest dieser Befehl wieder aus der Speicherstelle 456.

Die AMD64-Architektur erlaubt bei fast jedem Befehl, dass einer der beiden Parameter einen Speicherzugriff durchführt. Liegt der Zieloperand im Speicher, so wird einmal lesend und einmal schreibend auf den Speicher zugegriffen. Beispiel:

```

addq %r8, (%r9)

```

Wenn z.B. in r9 der Wert 456 steht, dann wird der Wert von dieser Speicherstelle gelesen, r8 dazugaddiert, und das Ergebnis wieder an Speicherstelle 456 gespeichert; r9 bleibt unverändert.

Effektive Adressen

Es ist manchmal praktisch, die Adresse, die von einer Adressierungsart ausgerechnet wird (die *effektive Adresse*), nicht sofort zum Speicherzugriff zu benutzen, sondern zunächst in ein Register zu laden. Das macht der Befehl `leaq`. Beispiel:


```
leaq 400(,%r8,8),%r9
```

Wenn r8 den Wert 7 enthält, rechnet der Befehl die Adresse 456 aus, und schreibt sie dann (ohne Speicherzugriff) in das Register r9.

Wir verwenden den Adressberechnungsbefehl `leaq` (load effective address), der den Quelloperanden als Adresse interpretiert und im Zieloperand ablegt (Achtung, hier erfolgt kein Speicherzugriff), um einige der komplexen Addressierungsarten zu demonstrieren:

```
leaq -4(%ebp), %rax    /* rax = ebp-4      */
leaq (%rdi,%rsi), %rax /* rax = rdi+rsi    */
leaq 8(%rdi,%rsi), %rax /* rax = rdi+rsi+8  */
leaq 4(,%rdi,8), %rax  /* rax = rdi*8+4    */
leaq 2(%rax,%rax,2), %rax /* rax = rax+rax*2+2 */
```

Man kann damit auch manche Rechnungen leichter durchführen als mit den eigentlichen Rechenbefehlen.

Datenstrukturzugriff

Das folgende Programm entspricht der C-Anweisung `a[i].m++` (mit den Datenstrukturdefinitionen aus Abb. 2.2), wobei `i` in Register `rdi` liegt. Die Multiplikation mit 6 kann effizient durch das Ausnutzen der mächtigen Addressierungsarten berechnet werden ($((i+i)+(i+i)*2)$):

```
leaq (%rdi,%rdi), %rax /* rax = rdi+rdi    ; i + i    */
incw a+2(%rax,%rax,2) /* *(a+2+rax+rax*2)++; a[i].m++ */
```

Der Befehl `incw` erhöht dabei den Inhalt der vom Operanden adressierten 2-Byte-Speicherstelle um 1. Die Adresse von `a` wird hier symbolisch durch einen Label mit dem Namen `a` repräsentiert; der Assembler ersetzt diesen Namen dann durch die numerische Adresse. Weiters ist `a+2` ein konstanter Ausdruck, der vom Assembler zu einer Konstanten ausgerechnet wird.

Rechenbefehle

`addq` ist ein Beispiel für die Klasse der arithmetischen/logischen Befehle. Es gibt diese Befehle in drei Formen:

- `addq reg1, reg2/mem2`
- `addq const, reg2/mem2`
- `addq reg1/mem1, reg2`

wobei das Ergebnis in den zweiten Parameter (`reg2` oder `mem2`) geschrieben wird, und beide Parameter die zu addierenden Zahlen enthalten.

Folgendes Programm besteht aus arithmetischen/logischen Befehlen:

```

leaq  -1(%rdi), %rax  /* rax = rdi-1 */
imulq %rdi,%rax     /* rax *= rdi  */
sarq  $1,%rax       /* rax = rax>>1 */

```

Dieses Programm berechnet $x = y(y-1)/2$, wobei y in `rdi` liegt und das Ergebnis x in `rax` gespeichert wird. Die Division wird in diesem Beispiel durch Verschieben um ein Bit nach rechts implementiert.

Kontrollfluss

Um den Kontrollfluss steuern zu können, gibt es Verzweigungsbefehle (Sprungbefehle). Die unbedingte Verzweigung ist

- `jmp Ziel`

Diese Anweisung springt zum *Ziel* (direkte Adressierung).¹ Am Sprungziel muss man den Label natürlich definieren, mit *Ziel* :

Zusätzlich zur unbedingten Verzweigung gibt es auch noch die bedingte Verzweigung, z.B.:

- `jCC Ziel`

Diese Anweisung springt abhängig vom Konditioncode CC (`l` (less, signed), `ge` (greater or equal, signed), `le` (less or equal, signed), `g` (greater, signed), `e` (equal), `ne` (not equal), `b` (below, unsigned), `ae` (above or equal, unsigned), `be` (below or equal, unsigned), `a` (above, unsigned)) und dem Inhalt der Flagregister zum *Ziel*. Die Flagregister werden von den meisten arithmetischen Befehlen und dem Vergleichsbefehl (`cmp`) gesetzt:

```

cmp  %rcx,%rdi      /* flags = rdi - rcx      */
je   ziel          /* if (rdi==rcx) goto ziel */
subq $1, %rdi      /* rdi--                  */
ziel:

```

In diesem Beispiel wird `rdi` dekrementiert, wenn `rcx` ungleich `rdi` ist (entsprechend der C-Anweisung `if (a != b) b--;`).

Ein etwas größeres Beispiel:

```

/* C code: */
long vsum(long x[], long n)
{
    long s=0, *p;
    for (p=x; p<x+n; p++)
        s += *p;
    return s;
}

```

¹Der `jmp`-Befehl unterstützt auch die indirekte Adressierung, aber darauf gehen wir hier nicht ein.

```

/* Assembler: x/p in rdi, n in rsi, s in rax, x+n in rcx */
movq  $0, %rax          /* rax = 0 */
leaq  (%rdi,%rsi,8), %rcx /* rcx = rdi+rsi*8 */
cmpq  %rcx, %rdi        /* flags = rdi - rcx */
jae   exit              /* if (rdi>=rcx) goto exit */
loop:
addq  (%rdi), %rax      /* rax += *(rdi) */
addq  $8, %rdi          /* rdi += 8 */
cmpq  %rcx, %rdi        /* flags = rdi - rcx */
jb    loop              /* if (rdi<rcx) goto loop */
exit:

```

Dieses Beispiel berechnet die Summe der Elemente des Arrays `x`. Man sieht folgende Unterschiede zwischen C und Assembler:

- Verwendung von Verzweigungsbefehlen statt Kontrollstrukturen mit Blockstruktur.
- In C wird bei Adressarithmetik automatisch skaliert, in Assembler muss der Programmierer das explizit machen (z.B. Übersetzung von `p++` nach `addq $8, %rdi`).
- Verwendung von Registern statt Variablen.
- Zwischenergebnisse müssen in Assembler oft explizit in Registern abgelegt werden.

2.6 Funktionen

Damit der Assemblercode mit von Compilern erzeugtem Code zusammenarbeitet, ist es sinnvoll, den Assemblercode in Form von Funktionen zu schreiben, die sich an die Konvention zum Funktionsaufruf halten, an die sich auch die Compiler auf dieser Plattform halten.

Diese Konvention legt fest, welche Register von der aufgerufenen Funktion zerstört werden dürfen (der Inhalt dieser *caller-saved*-Register muss ggf. vom Aufrufer gesichert werden), und welche die aufgerufene Funktion unverändert lassen muss (bzw. abspeichern und danach wiederherstellen muss (*callee-saved* (`rbx`, `rbp`, `r12` – `r15`))). Sie legen auch fest, wie Argumente und Rückgabewerte übergeben werden, und ähnlich Details.

Die Aufrufkonvention ist so gestaltet, dass eine beliebige Funktion eine beliebige andere aufrufen kann, ohne viel über die andere Funktion zu wissen.

Auf der AMD64-Architektur werden z.B. die ersten sechs Integer-Argumente in den Registern `rdi`, `rsi`, `rdx`, `rcx`, `r8`, `r9` und weitere Argumente auf dem Stack übergeben, das Funktionsergebnis in `rax`, und die Rücksprungadresse auf dem Stack. Der Unterprogrammaufruf erfolgt mit dem Befehl `call`, der Rücksprung erfolgt mit dem Befehl `ret`.

Wir können also unser erstes Programmfragment durch Hinzufügen dieser Befehle und eines Einsprunglabels zu einer vollwertigen Funktion aufrüsten (die Register sind nicht ganz zufällig schon entsprechend gewählt):

```
.globl f /* Funktionsnamen für den Linker exportieren */
f:      /* Funktionsname */
leaq   -1(%rdi), %rax /* rax = rdi-1 */
imulq  %rdi,%rax     /* rax *= rdi */
sarq   $1,%rax       /* rax = rax>>1 */
ret    /* return to caller */
```

Diese Funktion kann zum Beispiel von C aus mit `n=f(10)`; aufgerufen werden. Die `.globl`-Anweisung exportiert den Namen `f` so, dass der Linker ihn verwenden kann, und der Assembler-Code mit C-Code in einem anderen File zusammengebunden werden kann. Andere Namen sind nur innerhalb des Sourcefiles sichtbar, entsprechend der Speicherklasse `static` in C. Komplizierter wird die Lösung, wenn eine Funktion eine weitere Funktion aufruft. Nehmen wir als Beispiel die Funktion `g` die `f` aufruft:

```
static long g(long a, long b) {
    return b + f(a);
}
```

Dann müssen die Rücksprungadresse (implizit mit dem Befehl `call`) und die *callee-saved*-Register (explizit) zu Beginn gesichert und am Ende wiederhergestellt werden. *Caller-saved*-Register werden direkt vor einem Unterprogrammaufruf (vor eventuellen Argumenten, die nicht in Argumentregistern Platz gefunden haben) gesichert und direkt nach dem Unterprogrammaufruf wiederhergestellt. Dazu werden auf dem Stack entsprechende Speicherzellen reserviert. Der Stack wächst von hohen zu niedrigen Adressen. Der Stackpointer (Zeiger auf die Spitze des Stacks) `rsp` wird bei den Befehlen `call` und `push` automatisch vermindert und bei den Befehlen `pop` und `ret` automatisch erhöht:

```
g:      /* Funktionsname g, a in rdi, b in rsi */
pushq  %rsi /* caller-saved Register rsi (b) sichern */
call   f   /* f aufrufen */
popq   %rsi /* rsi wiederherstellen */
addq   %rsi, %rax /* rsi (b) zu Ergebnis von f addieren */
ret    /* Rücksprung */
```

Die meisten bisher gezeigten Programmfragmente können so in vollwertige Funktionen umgewandelt werden. Weitere Feinheiten der AMD64-Aufrufkonvention werden im Übungsskriptum behandelt.

Neben `.globl` produzieren Compiler noch weitere Assembleranweisungen mit Informationen für den Debugger und andere Werkzeuge.

2.7 Was dieses Kapitel verschweigt

Dieses Kapitel kann natürlich nur eine stark vereinfachte Darstellung von Computerarchitektur geben, für eine fundierte Darstellung empfiehlt es sich, ein gutes Lehrbuch zu lesen [HP03]. Einige der ausgelassenen Themen:

Die Architektur ist die Spezifikation für das Interface zwischen Maschinenprogramm und Prozessor. Daneben gibt es noch die *Mikroarchitektur*, die die Ausführungszeit bestimmt; hierunter fallen Konzepte wie z.B. Caches, Pipelines, und superskalare Ausführung, die für das Programm normalerweise nicht sichtbar sind. Für eine Architektur (z.B. IA32) gibt es oft viele Mikroarchitekturen (z.B. 386, 486, P5, P6, K5, K6, K7, K8, K10, Netburst, Core 2), und von den Mikroarchitekturen oft wieder verschiedene Subvarianten (z.B. vom P6: Pentium Pro–Pentium III).

Die virtuelle Speicherverwaltung ist ebenfalls für User-Level-Programme kaum sichtbar, allerdings auf Betriebssystemebene schon; sie wird daher zur Architektur gezählt, aber trotzdem in diesem Kapitel ignoriert.

Dieses Kapitel geht auch nicht auf die Architektur im I/O-Bereich ein (z.B. Platten, Graphikkarte, etc.); I/O steht User-Level-Programmen nur über Betriebssystemaufrufe zur Verfügung.

Die Darstellung des Speichers ist eine Vereinfachung; neben den hier dargestellten weit verbreiteten byte-adressierten 32- und 64-bit-Maschinen mit flacher Adressierung gibt es noch andere.

Kapitel 3

Struktur von Compilern

Ein Compiler wandelt das Quellprogramm in mehreren **Phasen** (logischen Teilaufgaben) in das Zielprogramm um. Jede Phase bearbeitet eine abgeschlossene Teilaufgabe. Die Phasen sind durch den Datenfluß voneinander abhängig. In vielen Fällen können sie zeitlich verzahnt durchgeführt werden.

Phasen, die beim Ablauf nicht mit anderen Phasen verzahnt sind, bezeichnet man als einen „Durchlauf“ oder **Pass**. Es gibt logische und technische Gründe, den Ablauf eines Compilers in einen oder mehrere Pässe zu gliedern. Man beachte aber, daß die Phasen - die logischen Teilaufgaben - immer die gleichen sind.

3.1 Phasen

Jede Übersetzung hat zwei grundlegende Phasen: Analyse und Synthese.

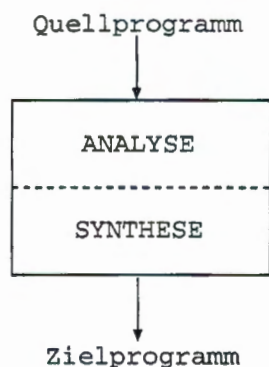


Abbildung 3.1: Grundstruktur eines Compilers

Die **Analyse** überprüft die **Syntax** des Quellprogramms und ermittelt seine **Struktur**. Außerdem wird die **statische Semantik** überprüft und ausgewertet. Die statische Semantik umfaßt semantische Eigenschaften des Programms, die zur Übersetzungszeit festliegen. Dies sind insbesondere die **Typen** von Objekten.

Die **Synthese** baut das Zielprogramm auf. Besondere Aspekte dabei sind einerseits **Optimierungen** zur Erhöhung der Effizienz des Zielprogramms und andererseits der Einbau **dynamischer Kontrollen** (z.B. Indexüberprüfung) zur Erhöhung der Sicherheit des Zielprogramms. Die Grundstruktur läßt sich weiter verfeinern:

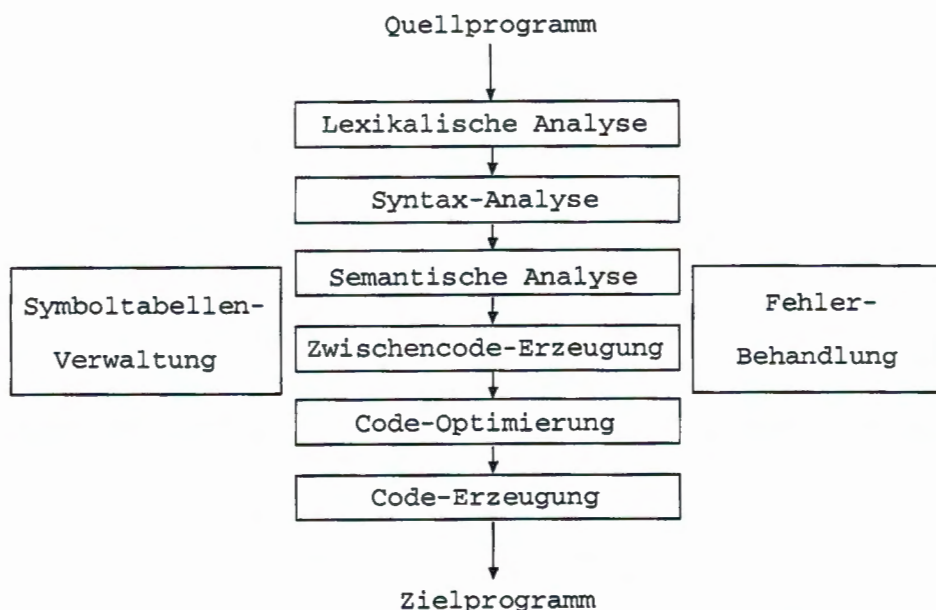


Abbildung 3.2: Verfeinerte Phasenstruktur

Gemäß Abb. 3.2 lassen sich Analyse und Synthese jeweils in drei Phasen gliedern. Dabei sind einzelne Phasen je nach Programmiersprache und Anspruch an die Qualität des erzeugten Codes mehr oder weniger stark ausgeprägt.

Die Pfeile geben die logische Abhängigkeit (Datenfluß) der Phasen an. Die seitlich dargestellten „Phasen“ stellen Teilaufgaben dar, die nicht in diesem Datenfluß liegen, sondern im wesentlichen mit allen Phasen Informationen austauschen.

Wir werden im folgenden anhand eines Programm-Beispiels einen Überblick über die Phasen einer Übersetzung geben. Als Quellsprache wird hier Modula-2 verwendet.

In Abb. 3.3 und Abb. 3.4 sind die Zwischenergebnisse der einzelnen Phasen dargestellt. Man beachte, daß es sich dabei um konzeptuelle Zwischendarstellungen handelt. Sie zeichnen die „Geschichte“ der erzeugten Informationen auf. Tatsächlich werden die Informationen oft nur in Teilen erzeugt. Es folgt nun eine kurze Beschreibung der Phasen.

Lexikalische Analyse

Die Lexikalische Analyse untersucht die Zeichenfolge des Quellprogramms und teilt sie in **Grundsymbole** ein. Diese werden durch Paare (Token, Attribut) dargestellt.

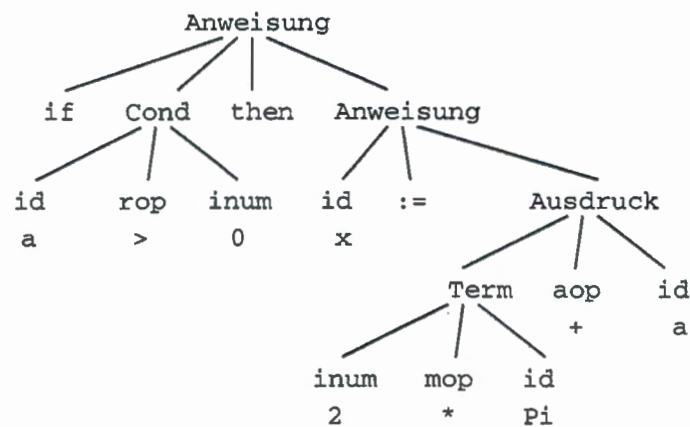
Ausschnitt eines Modula-2 Quellprogramms:

```
... CONST Pi=3.14;
    VAR a:INTEGER; x:REAL;
... IF a > 0 THEN
    x:=2*Pi+a ...
```

Folge von Symbolen:

```
... (const)(id,Pi)(=)(rnum,3.14)(;)
    (var)(id,a)(:)(int)(;)(id,x)(:)(real)(;)
... (if)(id,a)(rop,>)(inum,0)(then)
    (id,x)(:=)(inum,2)(mop,*)(id,Pi)(aop,+)(id,a) ...
```

Ableitungsbaum der IF-Anweisung:



Abstrakter Syntaxbaum:

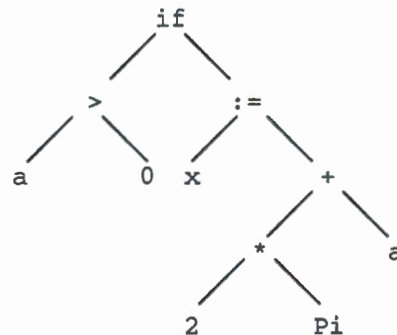
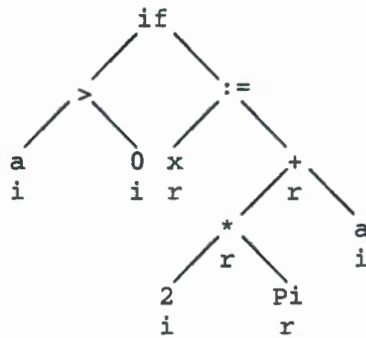


Abbildung 3.3: Vom Quellprogramm zum Syntaxbaum

Attributierter Syntaxbaum und Symboltabelle:



name art typ wert/(adr)

| name | art | typ | wert/(adr) |
|------|-----|-----|------------|
| Pi | c | r | 3.14 |
| a | v | i | (0) |
| x | v | r | (4) |

Linearer Zwischencode:

```
if a >i 0 goto L1 (*>i vergleicht integer-Werte *)
goto L2
```

L1: t1 =_r intreal(2) (*=_r bedeutet Zuweisung von real auf real *)

```
t2 =r t1 *r 3.14 (* *r multipliziert real-Werte *)
```

```
t3 =r intreal(a)
```

```
t4 =r t2 +r t3
```

```
x =r t4
```

L2:

Optimierter Zwischencode:

```
if a <=i 0 goto L
```

```
t1 =r intreal(a)
```

```
x =r 6.28 +r t1
```

L:

Symbolischer Maschinencode:

```
MOVi 0(A),R1
```

```
CMPi #0,R1
```

```
BLE L
```

```
FLT R1,F1 (* Konversionsbefehl von integer nach real *)
```

```
ADDR #6.28,F1
```

```
MOVr F1,4(A)
```

L:

Abbildung 3.4: Vom Syntaxbaum zum Maschinencode

Das Token kennzeichnet das Symbol und das Attribut die individuelle Ausprägung. Falls das Token eindeutig ist, bleibt das Attribut leer.

Die Lexikalische Analyse untersucht nur einfache (also keine geschachtelten) Strukturen. Diese können durch **Reguläre Ausdrücke** beschrieben werden. Die Analyse arbeitet nach dem Prinzip eines **Endlichen Automaten**. Durch die Lexikalische Analyse wird die nachfolgende Syntax-Analyse wesentlich vereinfacht.

Syntax-Analyse

Die Syntax-Analyse betrachtet die von der Lexikalischen Analyse gelieferten Token (nicht die Attribute) und erzeugt einen **Ableitungsbaum** entsprechend der Grammatik der Quellsprache, die in der Regel als **kontextfreie Grammatik** gegeben ist. Die Syntax-Analyse arbeitet nach dem Prinzip eines **Kellerautomaten**.

In den Endknoten des Ableitungsbaums stehen die Token als Terminalsymbole der Grammatik. Die Nonterminale der Grammatik dienen nur zum Aufbau der Struktur und können danach eliminiert werden. Man markiert dabei jeden Zwischenknoten durch den Operator, der den Teilbaum charakterisiert. Es entsteht ein Operatorbaum bzw. ein **(abstrakter) Syntaxbaum**.

Ein abstrakter Syntaxbaum kann auf verschiedene Weise implementiert werden: In imperativen Sprachen als dynamische Datenstruktur, in objektorientierten Sprachen als Instanz einer Klassenhierarchie, in funktionalen Sprachen als geschachtelte Liste und in logischen Sprachen als Term. Oft wird der Baum nicht explizit dargestellt, sondern nur implizit durch den Aufruf entsprechender Codestücke, die die semantische Analyse und die Zwischencodeerzeugung durchführen.

Semantische Analyse

Während die Syntax-Analyse nur strukturelle Beziehungen zwischen Symbolen untereinander analysiert, behandelt die Semantische Analyse die Bedeutung von Symbolen, d.h. die Beziehungen zwischen Symbolen und den von ihnen bezeichneten Objekten (z.B. Konstanten, Variablen).

Die Semantische Analyse trägt Art und Typ der von Symbolen bezeichneten Objekte in die **Symboltabelle** ein. Die zentrale Aufgabe der Semantischen Analyse ist die **Typ-Überprüfung** (*type checking*). Dabei ist zu prüfen, ob Verwendungen von Objekten mit ihren Deklarationen typkonsistent sind. Bei überladenen Operatoren (z.B. "+" für integer- und real-Addition) muß dabei aus dem Typ der Operanden der Typ des Operators abgeleitet werden. Man nennt dies **Operator-Identifizierung**. Im einfachsten Fall kann die Semantische Analyse den Syntaxbaum von den Blättern zur Wurzel hin untersuchen. Das Ergebnis der Analyse ist ein **attributierter Syntaxbaum**, der das syntaktisch und statisch-semantisch analysierte Programm repräsentiert.

Die Semantische Analyse behandelt **Kontextbeziehungen**, die nicht durch eine kontextfreie Grammatik definiert werden können - formal ist auch sie ein Teil der Syntax-Analyse, nämlich der nicht-kontextfreien. Als Basis für die Semantische Analyse eignet sich die **Attributierte Grammatik**.

Zwischencode-Erzeugung

Der attributierte Syntaxbaum kann bereits als eine Darstellung des Programms in einem Zwischencode angesehen werden, der von einer abstrakten Maschine ausgeführt werden könnte. Im Hinblick auf reale Maschinen muß der Baum aber linearisiert werden. Dabei sind insbesondere Sprungbefehle zu erzeugen.

In Abb. 3.4 ist für den linearen Zwischencode eine abstrakte Drei-Adreß-Maschine gewählt. Zur Linearisierung wird der Syntaxbaum in einem Tiefendurchlauf von links nach rechts abgearbeitet. Für den bedingten Ausdruck werden Sprungbefehle und Marken erzeugt. Bei der Zuweisungs-Anweisung wird für jeden Operator-knoten ein Befehl generiert. Zur Typ-Konversion werden zusätzliche Befehle eingefügt (z.B. `intreal`). Für jedes Zwischenergebnis wird eine temporäre Hilfsvariable t_i verwendet.

Alle bisherigen Phasen sind prinzipiell noch unabhängig von einer konkreten Zielmaschine. Diese ersten Phasen bezeichnet man auch als **front-end** des Compilers und die anschließenden Phasen als **back-end**.

Code-Optimierung

Die Code-Optimierung versucht, das Programm zu vereinfachen, um Speicherplatz und/oder Laufzeit einzusparen. Grundsätzlich kann man unterscheiden zwischen maschinenunabhängigen und maschinenabhängigen Optimierungen. Typische maschinenunabhängige Optimierungen sind z.B. die Vereinfachung von Indexberechnungen in Schleifen und die Auswertung konstanter Ausdrücke. Bei der maschinenabhängigen Optimierung dagegen werden vor allem spezielle Befehle der Zielmaschine ausgenutzt.

In Abb. 3.4 werden einfache maschinenunabhängige Optimierungen gezeigt: Der Sprungbefehl `goto L2` kann durch Umkehrung der Vergleichsoperation eliminiert werden. Die Konversionsoperation `intreal(2)` muß nicht erzeugt werden, sondern kann zur Übersetzungszeit ausgeführt werden. Dann läßt sich auch die Multiplikation $2.0 * 3.14$ berechnen. Insgesamt ergibt sich eine Optimierung von 7 auf 3 Befehle.

Code-Erzeugung

Die Code-Erzeugung generiert ein Programm, das auf einer realen Maschine ablaufen kann. Dieses Programm wird auch **Objektprogramm** genannt. Die in Abb. 3.4 gewählte Maschine entspricht einem 68000-Prozessor, der über Register R_i verfügt und um einen Koprozessor für real-Operationen mit Registern F_i erweitert ist. Wir nehmen an, daß ein integer-Wert 4 Bytes benötigt. Wenn a im Variablenspeicher die relative Adresse $d=0$ hat, bekommt x die relative Adresse $d=4$. Diese relativen Adressen werden in die Symboltabelle eingetragen. Das Adreßregister A enthält zur Laufzeit des Objektprogramms die Anfangsadresse des Variablenspeichers. Dann können die Variablen durch indizierte Adressierung $d(A)$ angesprochen werden.

In Abb. 3.4 ist der erzeugte Code als Assemblercode angegeben, um den Code lesbar zu machen. Dieser Code muß durch den **Assembler** in ausführbaren Maschinencode übersetzt werden.

Binder und Lader

Programme werden im allgemeinen als verschiebbarer (*relocatable*) Code erzeugt. D.h., daß sie nicht in einen bestimmten Speicherbereich geladen werden müssen. In vielen Programmiersprachen gibt es außerdem die Möglichkeit, ein Programm auf verschiedene Module (Dateien) aufzuteilen und getrennt zu übersetzen. Das Überprüfen der Schnittstellen ist Aufgabe des Compilers, während das Zusammenfügen von Programmteilen (Auflösen der Referenzen) und das Laden in den Speicher vom Betriebssystem, d.h. vom **Binder** (*linker*) und **Lader** (*loader*) durchgeführt werden. Dabei können auch Bibliotheksprozeduren (*libraries*), die schon in Objektcode vorliegen, eingebunden werden. Grundsätzlich kann Binden und Laden von Programmteilen vor oder während der Laufzeit des Objektprogramms erfolgen.

Laufzeitsystem

Die Zielmaschine eines Compilers ist im allgemeinen eine um Software erweiterte Hardware gemäß Abb. 3.5. Das Objektprogramm verwendet Funktionen der Hardware, Funktionen des Betriebssystems (z.B. E/A Routinen) und Funktionen des **Laufzeitsystems**. Das Laufzeitsystem umfaßt Unterprogramme, die zum Ablauf von Objektprogrammen generell benötigt werden (z.B. Speicherverwaltung). Es bildet eine Compiler-spezifische Softwareschicht über dem Betriebssystem und muß daher auf das verwendete Betriebssystem genau abgestimmt werden.

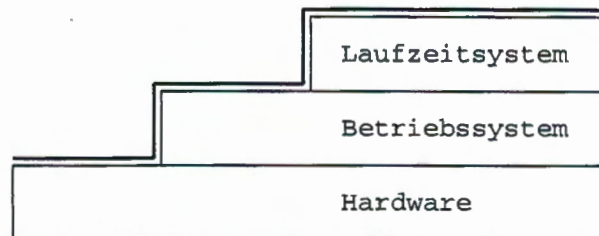


Abbildung 3.5: Schichtenmodell der Zielmaschine

Ein Interpreter kann als ein Laufzeitsystem angesehen werden, das Hardware und Betriebssystem voll überdeckt und einen eigenen Maschinencode interpretiert.

3.2 Ablaufsteuerung

Alle Phasen (Teilaufgaben) eines Compilers können unter bestimmten Umständen zeitlich verzahnt durchgeführt werden, d.h. das Quellprogramm kann in einem Pass (Durchlauf) übersetzt werden. Einen solchen Compiler nennt man daher einen Ein-Pass-Compiler. Wenn das Quellprogramm oder Zwischendarstellungen davon mehrmals von einem Compiler gelesen werden müssen, spricht man von einem Mehr-Pass-Compiler. Wir betrachten im folgenden einige typische Strukturen.

Ein-Pass-Compiler

Ein-Pass-Compiler sind besonders schnell, da das Programm nur einmal gelesen wird und keine Zwischendarstellungen des Programms erzeugt und weitergegeben werden müssen. Ein-Pass-Übersetzung ist aber nur möglich, wenn es die Programmiersprache zuläßt. Insbesondere müssen alle Deklarationen textlich *vor* ihrer Verwendung stehen, damit der Compiler zu jedem Zeitpunkt über vollständige Informationen verfügt. Optimierungen sind nur sehr beschränkt möglich, da keine globalen Informationen gesammelt und ausgewertet werden können. Abb. 3.6 zeigt die typische Ablaufstruktur.

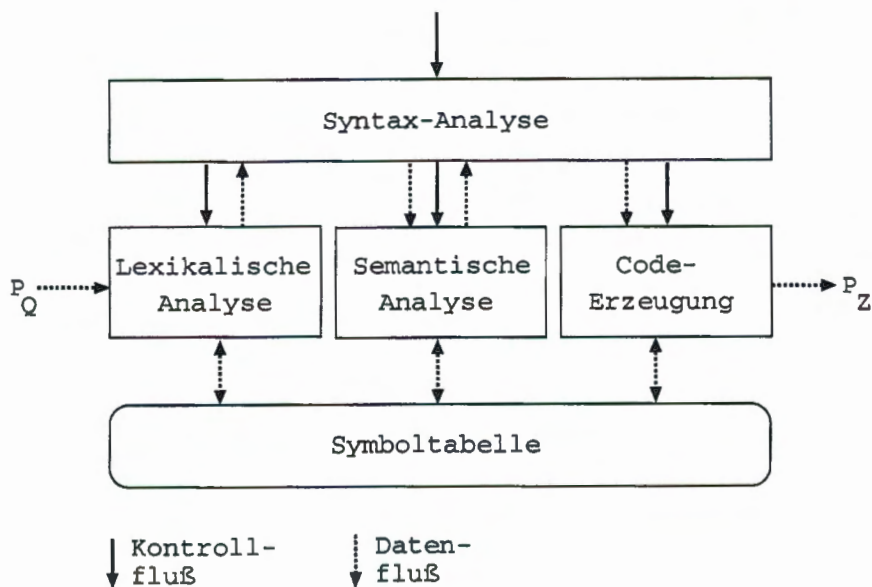


Abbildung 3.6: Ablaufstruktur eines Ein-Pass-Compilers

Der Ein-Pass-Compiler wird von der Syntax-Analyse gesteuert. Die anderen Teilaufgaben (Lexikalische Analyse, Semantische Analyse und Code-Erzeugung) werden als Prozeduren aufgerufen.

Mehr-Pass-Compiler

Falls die Quellsprache erlaubt, daß Deklarationen auch textlich *nach* ihrer Verwendung stehen können, muß der Compiler mit mindestens zwei Pässen arbeiten. Der erste Pass analysiert zunächst alle Deklarationen und der zweite Pass baut das Zielprogramm auf. Abb. 3.7 zeigt die typische Struktur.

Der erste Pass hat im Modell von Abb. 3.7 nur Analyseaufgaben. Er sammelt alle Deklarationen in der Symboltabelle und wandelt das Programm in eine geeignete Zwischendarstellung (Zwischencode) um. Nach Ablauf des ersten Passes wird der zweite Pass aufgerufen. Er verwendet die erzeugten Informationen und baut das Zielprogramm auf. Zwei-Pass-Übersetzung wird häufig für Assembler verwendet,

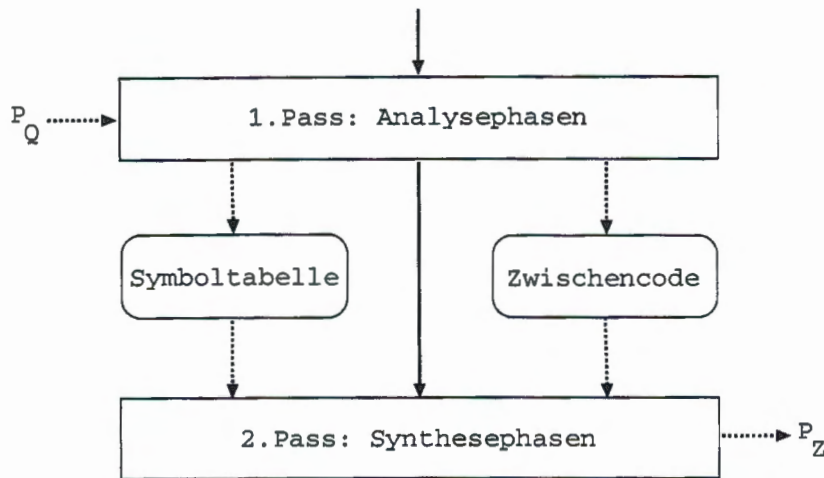


Abbildung 3.7: Ablaufstruktur eines Zwei-Pass-Compilers

wobei der erste Pass die Adressen der symbolischen Marken bestimmt und der zweite Pass den Maschinencode erzeugt.

Im Prinzip kann jede Phase eines Compilers als ein Pass implementiert werden. Der Steuerfluß der Pässe entspricht dann genau dem Datenfluß von Abb. 3.2. Der Compiler startet mit der Lexikalischen Analyse. Jeder Pass ruft seinen logischen Nachfolger auf. Zwischen den Pässen wird jeweils die Zwischendarstellung des Programms als Ganzes übergeben.

Ein Mehr-Pass-Compiler hat den Vorteil, daß die einzelnen Pässe klare, explizite Schnittstellen haben. Außerdem benötigt der Compiler wenig Hauptspeicher, da die Pässe nacheinander geladen werden können. Der Nachteil liegt darin, daß das Erzeugen und Verarbeiten von Zwischendarstellungen zeitaufwendig ist. Heutige Compiler können nicht so einfach eingeordnet werden. Oft zeigen sie im Back-End zwar im Prinzip eine Mehr-Pass-Struktur, nicht aber auf der Ebene von Eingabedateien. Die Zwischendarstellungen werden im Speicher übergeben.

Kapitel 4

Lexikalische Analyse

Die Lexikalische Analyse (*scanning*) liest das Quellprogramm als Zeichenfolge (z.B. ASCII) und teilt es in Symbole ein. Für den Compiler bedeutungslose Zeichen (z.B. Leerzeichen, Kommentare) werden dabei entfernt.

4.1 Grundsymbole

Ein oder mehrere aufeinanderfolgende Zeichen im Quellprogramm bilden ein **Grundsymbol** (auch kurz: Symbol). Nach ihrer Schreibweise im Quellprogramm kann man folgende Arten von Grundsymbolen unterscheiden:

- Wortsymbole (*keywords*), z.B. `if`, `while`
- Bezeichner (*identifier*) als Namen für definierte Objekte, z.B. `x1`
- Literale (*literals*) als Notationen für Werte, z.B. `3.14`
- Spezialsymbole (*delimiter*), z.B. `=`, `+`, `==`

Die Codierung der Symbole in der Ausgabe der Lexikalischen Analyse wird durch die weitere Verarbeitung im Compiler bestimmt. Die nachfolgende Syntax-Analyse basiert auf einer kontextfreien Grammatik, welche die Grundsymbole als *Terminale* verwendet. Entsprechend klassifiziert man die Symbole. Jedes Wortsymbol wird auf ein eigenes Terminal, alle Bezeichner werden auf ein einziges Terminal abgebildet. Literale werden nach dem Typ des Wertes (z.B. `integer`, `real`, `string`) auf Terminale abgebildet. Spezialsymbole entsprechen grundsätzlich jeweils einem Terminal. Man kann aber Operatoren auch zusammenfassen und entsprechend ihrer syntaktischen Rolle (z.B. `relational`, `arithmetisch`) auf Terminale abbilden.

Für die Syntax-Analyse ist allein die Abbildung der Symbole auf Terminale relevant, aber für weitere Phasen des Compilers ist auch die individuelle Ausprägung, z.B. die Identität der Bezeichner und der Wert von Literalen von Bedeutung. Zur Codierung eines Symbols benötigt man daher ein Tupel (Token, Attribut).

Token ist der Terminal-Code für die Syntax-Analyse, während das **Attribut** die Ausprägung des Symbols beschreibt. Als Attribut kann man im Prinzip das **Lexem**

(d.h. die Zeichenfolge des Symbols) aus dem Quellprogramm übernehmen. Abb. 3.3 in Kapitel 2 zeigt eine entsprechende Codierung von Symbolen. Meistens wird für die Fehlerbehandlung in späteren Phasen auch die Stelle, an der ein Symbol bzw. Lexem im Quellprogramm steht (Zeile, Spalte), als weitere Komponente codiert.

Spezifikation von Symbolen

Der Aufbau von Symbolen aus Zeichen kann durch **Reguläre Ausdrücke** oder **Reguläre Definitionen** exakt spezifiziert werden. Man bezeichnet in diesem Zusammenhang die Symbole auch als lexikalische Konstrukte.

$$[0-9]^+ \mid [0-9]^* ([0-9] \text{"."} \mid \text{"."}[0-9]) [0-9]^*$$

Beispiel 4.1: Regulärer Ausdruck zur Beschreibung von Zahlen. Wir betrachten integer- und real-Zahlen ohne Exponent, wobei vor oder nach dem Dezimalpunkt mindestens eine Ziffer stehen soll. (Die meisten Programmiersprachen verlangen, daß mindestens eine Ziffer vor dem Punkt steht.)

Ein Regulärer Ausdruck (Bsp. 4.1) beschreibt eine Menge von Zeichenfolgen. Er wird aus Zeichen und Operatoren (Metasymbolen) aufgebaut. Die Zeichen werden zur Unterscheidung von den Metasymbolen oder zur Verdeutlichung bei Bedarf zwischen Anführungsstriche ("t") gestellt.

Als Operatoren (Metasymbole) stehen zur Verfügung:

- Bereich von Eingabezeichen ($[t_1-t_n]$ für Zeichen zwischen t_1 und t_n),
- Iteration (r^* beliebige Vorkommnisse, r^+ mindestens ein Vorkommnis von r),
- Option ($r?$),
- Verkettung ($r_1 r_2$) und
- Alternative ($r_1 \mid r_2$)

Die Reihenfolge entspricht den Operator-Prioritäten (die Alternative bindet am schwächsten), auch Klammerung ist möglich.

$[0-9]$ spezifiziert z.B. eine Ziffer, $[0-9]^+$ eine oder mehrere Ziffern.

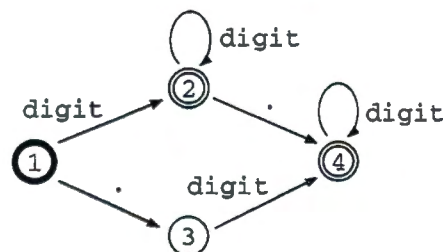
Reguläre Definitionen dienen zur Benennung und übersichtlichen Schreibweise von Regulären Ausdrücken (siehe Bsp. 4.2). Auf der rechten Seite der Definitionen stehen Reguläre Ausdrücke, die *Namen* von Regulären Ausdrücken enthalten können. Rekursionen sind dabei nicht erlaubt, d.h. rechts dürfen nur *vorher* definierte Namen stehen.

| | |
|------------|------------------------------|
| digit | = [0-9] |
| digits | = digit+ |
| pointgroup | = digit "." "."digit |
| integer | = digits |
| real | = digits? pointgroup digits? |
| number | = integer real |

Beispiel 4.2: Die reguläre Definition der Zahlen von Bsp. 4.1.

Erkennen von Symbolen

Zum Erkennen von Symbolen werden **Endliche Automaten** verwendet: Sie untersuchen Zeichen für Zeichen, ob der Anfang der vorliegenden Zeichenfolge einem Symbol entspricht. In Abhängigkeit von ihrem Zustand und dem zuletzt gelesenen Eingabezeichen gehen sie in einen neuen Zustand über, bis sie einen Endzustand erreichen und damit ein Symbol erkannt haben. Endliche Automaten können als **Zustandsdiagramm** (Bsp. 4.3) oder als **Zustandsmatrix** (Bsp. 4.4) dargestellt werden.



Beispiel 4.3: Zustandsdiagramm für die Zahlen-Analyse (Bsp. 4.1)

Die Zustände entsprechen den Knoten des Graphen (Anfangszustand: fester Rand, Endzustände: doppelter Rand). An den Kanten stehen die Eingabezeichen, bzw. Zeichenmengen, mit denen ein Übergang zwischen den Zuständen stattfindet.

Jede Zeile entspricht einem Zustand, jede Spalte einem Eingabezeichen. In der Matrix stehen die Folgezustände. Leereinträge (.) zeigen lexikalische Fehler an.

Der bisher betrachtete Automat ist ein deterministischer Endlicher Automat (DFA, *deterministic finite automaton*). Bei einem DFA gibt es für jede Kombination aus Zustand und Eingabezeichen nur maximal einen Folgezustand. Für einfache Reguläre Ausdrücke kann man den zugehörigen DFA intuitiv konstruieren, aber für komplexere Ausdrücke benötigt man ein algorithmisches Verfahren.

Im folgenden wird ein Verfahren vorgestellt, das Reguläre Definitionen in minimale DFA umformt. Die meisten Scanner-Generatoren (z.B. LEX unter Unix) arbeiten nach diesem Verfahren. Zunächst wird ein nichtdeterministischer Endlicher

| Zustand | Eingabezeichen | |
|---------|----------------|-----|
| | digit | "." |
| 1 | 2 | 3 |
| ◦ 2 | 2 | 4 |
| 3 | 4 | . |
| ◦ 4 | 4 | . |

Beispiel 4.4: Zustandsmatrix für die Zahlen-Analyse (vgl. Bsp. 4.3). Anfangszustand ist Zustand 1, Endzustände sind durch ◦ markiert. Die Zustandsmatrix ist hier optimiert: Durch Bilden der Zeichenmenge `digit` wurden 9 Spalten eingespart.

Automat (NFA) erzeugt, der dann in einen DFA umgewandelt wird. Wir beschreiben im folgenden das Verfahren in vier Schritten ausgehend von Bsp. 4.2.

(1) Umformung der Regulären Definition in Reguläre Ausdrücke

Durch wiederholte Textersetzung wird die Definition jedes Symbols in einen Regulären Ausdruck umgewandelt. Vereinfachungen sind dabei möglich, z.B. Ersetzen von `r+?` durch `r*` oder von `r* r` durch `r+`. Angewandt auf Bsp. 4.2 ergeben sich für `integer` und `real` die folgenden Regulären Ausdrücke, wobei zur besseren Lesbarkeit `digit` nicht ersetzt wurde.

```
digit+ | (* integer *)
digit* (digit "." | "."digit) digit* (* real *)
```

(2) Umformung der Regulären Ausdrücke in einen NFA

Bei einem NFA kann es für einen Zustand und ein Eingabezeichen mehrere Folgezustände geben. Er kann auch ϵ -Kanten enthalten. Der NFA erkennt ein Eingabewort (Symbol), wenn es einen Weg vom Start- zu einem Endzustand gibt, an dessen Kanten die Zeichen des Eingabeworts in entsprechender Folge stehen.

Für jeden Regulären Ausdruck, der einem Symbol entspricht, wird aus den Elementen von Abb. 4.5 ein NFA konstruiert (Bsp. 4.6). Die Anfangszustände der gewonnenen NFAs werden durch ϵ -Kanten mit einem gemeinsamen Startzustand verbunden (Bsp. 4.7). Es entsteht ein NFA mit einem Startzustand und mehreren Endzuständen, von denen jeder dem Erkennen eines bestimmten Symbols entspricht.

(3) Umformung des NFA in einen DFA

Der NFA läßt sich direkt zur Analyse einsetzen, indem bei mehreren möglichen Nachfolgezuständen alle Alternativen durchprobiert werden. Bei Anwendungen, die mit wechselnden Regulären Ausdrücken jeweils relativ wenig Text verarbeiten (z.B. Editoren), wird dies auch praktiziert. Für die Lexikalische Analyse im Compiler ist es effizienter, einen DFA zu verwenden. Jeder NFA kann in einen äquivalenten DFA (d.h. einen DFA, der die gleiche Sprache akzeptiert) umgewandelt werden. Diese

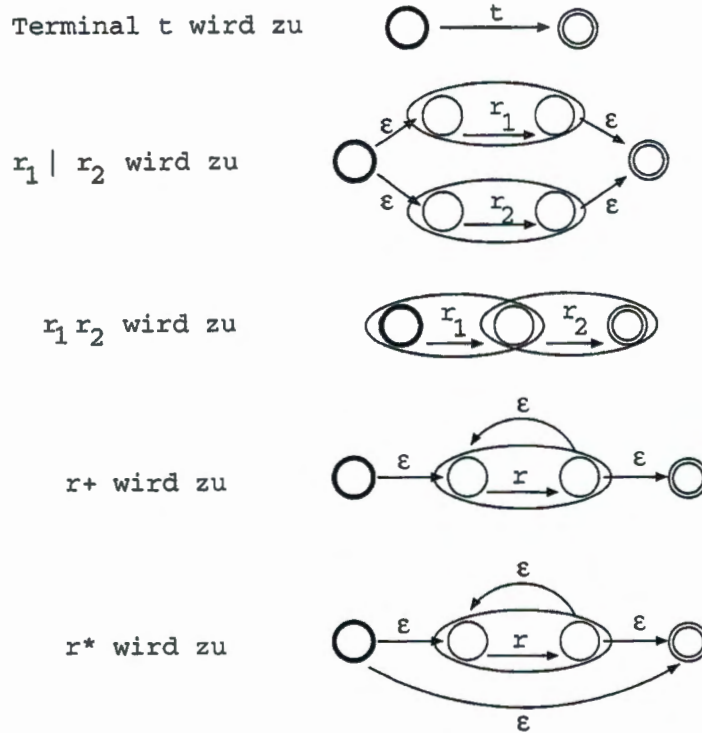
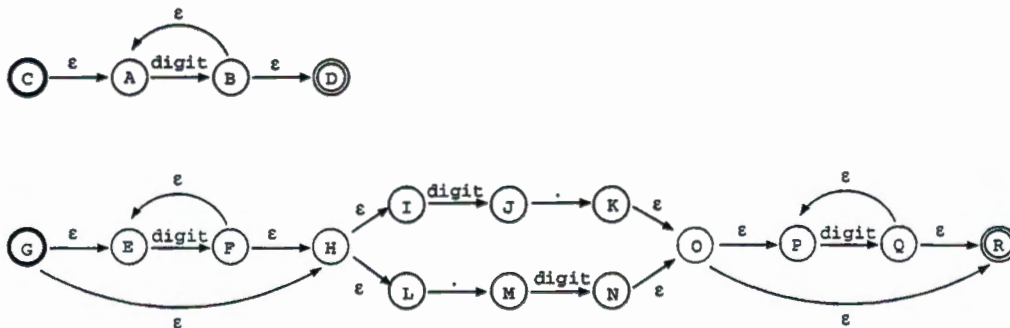


Abbildung 4.5: Konstruktion eines NFA aus einem regulären Ausdruck



Beispiel 4.6: Umformung der Regulären Ausdrücke für integer und real

| Zustand | digit | "." | ϵ |
|---------|-------|-----|------------|
| S | . | . | CG |
| A | B | . | . |
| B | . | . | AD |
| C | . | . | A |
| o D | . | . | . |
| E | F | . | . |
| F | . | . | EH |
| G | . | . | EH |
| H | . | . | IL |
| I | J | . | . |
| J | . | K | . |
| K | . | . | O |
| L | . | M | . |
| M | N | . | . |
| N | . | . | O |
| O | . | . | PR |
| P | Q | . | . |
| Q | . | . | PR |
| o R | . | . | . |

Man erkennt an den Doppeleinträgen und ϵ -Übergängen, daß der Automat nichtdeterministisch ist.

Beispiel 4.7: Die Zustandsmatrix des erzeugten NFA

Umwandlung ist ein häufig angewandtes Verfahren, das uns in diesem Skriptum noch mehrfach begegnen wird (z.B. in Kap. 5). Das Grundprinzip dabei ist, mehrere mögliche Nachfolgezustände zu einem einzigen neuen Zustand zusammenzufassen und Entscheidungen aufzuschieben, bis mehr von der Eingabe bekannt ist. Der Alg. 4.8 erzeugt nach diesem Prinzip aus einem NFA die Zustandsmatrix `dTransTab` des äquivalenten DFA.

Jedes Element in `dStates`, das einen NFA-Endzustand enthält, wird Endzustand des DFA (Bsp. 4.9). Der DFA kann durch Zusammenfassen von Zuständen verkleinert werden - eine derartige Minimierung ist aus Effizienzgründen immer sinnvoll.

(4) Minimierung eines DFA

Die Anzahl der Zustände wird minimiert, indem äquivalente Zustände zusammengefaßt werden (zwei Zustände sind äquivalent, wenn sie durch die gleiche Menge von Eingabeworten in einen Endzustand führen). Praktisch geht man so vor, daß man die Menge der Zustände iterativ in nicht äquivalente Teilmengen zerlegt (Alg. 4.10 und Bsp. 4.11).

```

public NStates move(NStates nStates, char input) liefert die Menge aller
    NFA-Zustände, zu denen Übergänge von einem Zustand aus der Menge
    nStates mittels des Eingabezeichens input führen;

public NStates epsclosure(NStates nStates) liefert die Menge aller NFA-
    Zustände, die über eine Folge von  $\epsilon$ -Übergängen von einem Zustand in nStates
    erreicht werden können (d.h. die  $\epsilon$ -Hülle). Darin sind alle Zustände aus
    nStates selbst enthalten

dStates = epsclosure( InitialState ) // Startzustand
do {
    Für jeden DFA-Zustand dState in dStates
    und jedes Eingabezeichen input
    newState = epsclosure(move(dState, input))
    if( !empty(newState) && !dStates.member(newState) )
        // neuer Zustand noch nicht enthalten
        füge newState zu dStates hinzu
    dTransTab[dState][input] = newState
} while(dStates verändert sich noch immer)

```

Algorithmus 4.8: Erzeugung der Zustandsmatrix eines äquivalenten DFA aus einem NFA.

4.2 Ausgabe-Aktionen

Das vorige Kapitel behandelte die Spezifikation bzw. Analyse von Symbolen (lexikalischen Konstrukten). Die eingelesene Zeichenfolge soll aber nicht nur geprüft, sondern auch umcodiert werden. Dazu erweitert man die Analyse um geeignete (Ausgabe-)Aktionen.

Aktionen werden prinzipiell erst nach dem Erkennen eines Symbols ausgeführt. Das kann zu Redundanz führen, wenn das Lexem nochmals durchsucht werden muß (z.B. bei der Behandlung von doppelten Hochkommas in Strings). Im allgemeinen aber kann der Scanner erst am Schluß entscheiden, welche Aktion(en) auszuführen sind.

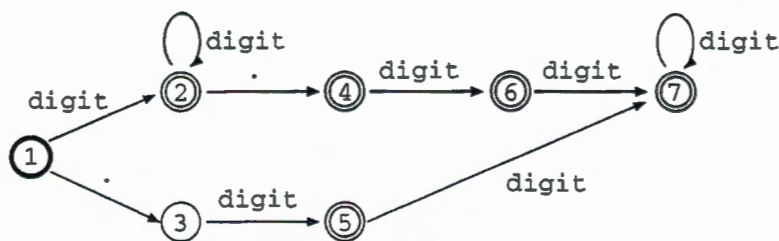
Bei Regulären Definitionen wird jedem Symbol eine Aktionsfolge oder eine Prozedur zugeordnet (Bsp. 4.12).

Bei Ausführung der Aktion muß die als Symbol erkannte Zeichenfolge (Lexem) des Quellprogramms verfügbar sein. Dazu kann man die Zeichen während der Analyse in einem Array sammeln, oder man verwendet zwei Zeiger in den Eingabepuffer: Der erste wird auf den Anfang des Lexems gesetzt und der zweite wird bei der Analyse fortgeschaltet, so daß er nach Erkennen des Symbols das Ende des Lexems anzeigt. Folgende Aktionen treten bei der Analyse üblicher Programmiersprachen auf:

Man beginnt mit der ϵ -Hülle (siehe Prozedur `eps closure`) von S: Sie enthält S selbst; von S führen ϵ -Kanten nach C und G, von da nach A bzw. nach E und H und von dort nach I und L. Von keinem dieser Zustände führen ϵ -Kanten weg. Der Startzustand des DFA ist also SACEGHIL (kurz für $\{S, A, C, E, G, H, I, L\}$). `move(SACEGHIL, digit) = BFJ` (wegen der Übergänge von A, E und I). Die Hülle davon (ABDEFHIJL) bildet den zweiten DFA-Zustand. Die erste Zeile der Zustandsmatrix wird mit der Hülle von `move(SACEGHIL, ".") = M` fertiggestellt. Mit den neu gewonnenen DFA-Zuständen wird die Matrix Zeile für Zeile erstellt:

| | DFA-Zustand | digit | "." |
|-------|-------------|-----------|-------|
| (1) | SACEGHIL | ABDEFHIJL | M |
| (2) ◦ | ABDEFHIJL | ABDEFHIJL | KMOPR |
| (3) | M | NOPR | . |
| (4) ◦ | KMOPR | NOPQR | . |
| (5) ◦ | NOPR | QPR | . |
| (6) ◦ | NOPQR | QPR | . |
| (7) ◦ | QPR | QPR | . |

Zustandsdiagramm des generierten DFA:



Beispiel 4.9: Determinisierung des NFA

Zurückliefern des Token. Token werden als integer-Werte in Form benannter Konstanten oder durch einen Aufzähltyp repräsentiert.

Eintragungen in Tabellen. Bezeichner und Konstanten werden aus Effizienzgründen oft in Tabellen gehalten. Bei ihnen muß also überprüft werden, ob das Lexem bereits in der entsprechenden Tabelle enthalten ist. Sollte das nicht der Fall sein, wird es eingetragen. In jedem Fall wird ein Verweis auf den Eintrag (Index) als Attribut zurückgeliefert. In Bsp. 4.12 erledigen das die Prozeduren `installInum` bzw. `installRnum`.

Behandlung von Wortsymbolen. Nach dem Erkennen eines Namens muß überprüft werden, ob dieser Name ein Wortsymbol ist. Das entsprechende Token soll ausgegeben werden (z.B. `if`, `then`, ... im Gegensatz zu `id`). Dazu wird entweder eine eigene Tabelle für Wortsymbole geführt oder die Bezeichner-Tabelle wird am Anfang

Anfangsaufteilung herstellen: teile die Zustandsmenge `dStates` des DFA in 2 Teilmengen: alle Endzustände, alle anderen Zustände

```
rest = dStates - finalStates
partition = SetOf( finalStates, rest )
endgültige Aufteilung iterativ generieren
```

```
do {
    teile jede Menge in partition so in Untermengen, daß zwei
    Zustände genau dann in dieselbe Untermenge kommen, wenn die
    Übergänge mittels der gleichen Eingabezeichen jeweils in dieselben
    ursprünglichen Mengen führen
    partition = die neuen Untermengen
} while(partition hat sich geändert)
```

Zustandsmatrix des minimalen Automaten erzeugen: Die neuen Zustände sind die Mengen der endgültigen Aufteilung. Das sind maximal soviele wie beim ursprünglichen DFA, wenn jede Menge einelementig ist.

```
states = partition
Übergänge (transTab) zwischen den neuen Zuständen (Mengen) genau
dann generieren, wenn es Übergänge zwischen ihren Elementen gibt.
```

Algorithmus 4.10: Minimierung eines DFA durch Zerlegung der Menge der Zustände in nicht äquivalente Teilmengen.

mit den Wortsymbolen initialisiert. Indizes kleiner als ein bestimmter Grenzwert sind dann für Wortsymbole reserviert, größere zeigen auf einen Bezeichner.

Ausblenden von Zeichen. In den meisten Programmen gibt es viele Zeichen, die nur die Lesbarkeit erhöhen (Zwischenraum, Tabulator, Zeilenwechsel, ...), für den Compiler also unwichtig sind. Manche benötigt der Scanner, um das Ende eines Lexems zu erkennen (Trennzeichen). Alle derartigen Leerzeichen können von der Lexikalischen Analyse ausgeblendet werden, d.h. ihre Ausgabeaktion ist leer. Das Mitzählen von Zeilenwechselln (und Spalten) ist jedoch für genaue Fehlermeldungen nötig. Ähnlich wie die Leerzeichen sind Kommentare für den Compiler bedeutungslos und können ausgeblendet werden. Geschachtelte Kommentare (wie z.B. in Modula erlaubt) lassen sich mit normalen regulären Methoden nicht behandeln.

4.3 Analyse-Verfahren

Die Analyse soll den jeweils längsten passenden Anfang der Eingabefolge als Symbol erkennen (*longest input match*). Beim Automaten in Bsp. 4.3 führen von den Endzuständen noch Kanten weg. Daher sind auch in einem Endzustand weitere Zeichen zu lesen, solange es Übergänge für sie gibt.

Anfangsaufteilung: 24567 (Endzustände), 13 (restliche Zustände)

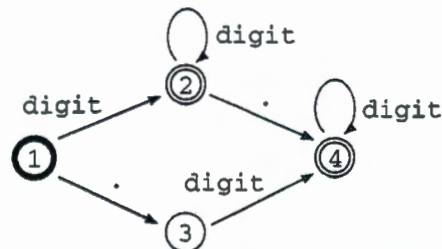
1. Iterationsschritt:

1. Menge: 2 (digit → 24567, "." → 24567)
 4 (digit → 24567, "." → {})
 5 (digit → 24567, "." → {})
 6 (digit → 24567, "." → {})
 7 (digit → 24567, "." → {})
2. Menge: 1 (digit → 24567, "." → 13)
 3 (digit → 24567, "." → {})
1. Aufteilung: 2, 4567, 1, 3

2. Iterationsschritt: 2, 1 und 3 sind einelementig – es ist keine weitere Aufteilung möglich. Bei der Menge 4567 führen bei jedem Element Übergänge mit digit nach 4567 selbst, Übergänge für "." existieren nicht. Es ist keine weitere Aufteilung nötig. In der Folge wird kurz 4 statt 4567 geschrieben.

2. Aufteilung = 1. Aufteilung = endgültige Aufteilung

Der minimale DFA hat das folgende Aussehen (vgl. Bsp. 4.3):



Beispiel 4.11: Minimierung des DFA: Zustandsmengen werden durch Ziffernfolgen dargestellt. Bei den Iterationsschritten stehen die möglichen Übergänge in Klammern, wobei der ursprüngliche Folgezustand innerhalb der Mengen unterstrichen ist.


```

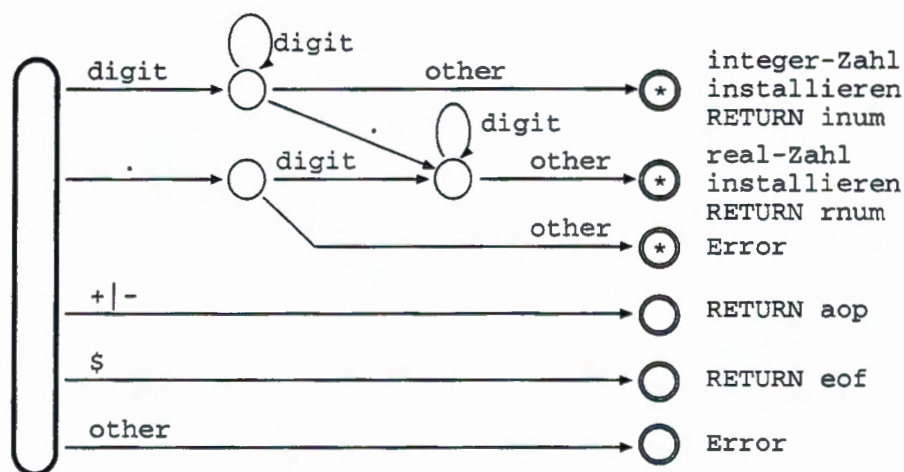
/* Definition von Symbolen */
integer = digits
real    = digits? pointgroup digits?

/* Aktionen */
integer : installInum(index); return( inum, index)
real    : installRnum(index); return( rnum, index)
    
```

| Zustand | Eingabezeichen | | Aktionen |
|---------|----------------|-----|-------------------------------|
| | digit | "." | |
| 1 | 2 | 3 | . |
| o 2 | 2 | 4 | install...; return(inum,...) |
| 3 | 4 | . | . |
| o 4 | 4 | . | install...; return(rnum,...) |

Beispiel 4.12: Bsp. 4.2 und Bsp. 4.4 werden um Aktionen erweitert

Die Analyse wird einfacher, wenn von Endzuständen keine Kanten wegführen. Bei Lexemen ohne erkennbare Endezeichen (z.B. Zahlen) muß man das erste Zeichen, das nicht mehr zum Lexem gehört, als Endezeichen (*other*) verwenden. Damit wird aber ein Zeichen „zuviel“ gelesen, das für das nächste Symbol wieder berücksichtigt werden muß. Im Bsp. 4.13 sind solche Endzustände mit * markiert.



Beispiel 4.13: Ein DFA zur Analyse von Zahlen (wie in Bsp. 4.3) und Additionsoperatoren. Der Automat kann z.B einen Ausdruck 17+3.14-.32\$ lexikalisch analysieren.

Bei einem erkannten Fehler (Error) wird im einfachsten Fall ein „Fehlertoken“ err zurückgeliefert und eventuell im Attribut zusätzliche Information übergeben.

Eine aufwendige Fehlerbehandlung in der Lexikalischen Analyse selbst ist nicht besonders sinnvoll, da der jeweilige Kontext nicht bekannt ist.

\$ wird als eof-Zeichen verwendet; es darf in der Sprachdefinition nicht vorkommen.

Tabellengesteuerte Analyse

Zugrunde liegt ein Kontroll-Algorithmus, der für alle Automaten gleichbleibt. Mittels der jeweiligen Zustandsmatrix (Tabelle) wird die Analyse gesteuert. Dieses Prinzip ist insbesondere für Scanner-Generatoren geeignet. Das Problem des „Zuweit-Lesens“ kann man dabei auf zwei Arten lösen:

- Das überschüssige Zeichen wird *zurückgeschrieben* (siehe Alg. 4.14).
- Es wird grundsätzlich ein Zeichen *vorausgelesen* (siehe Alg. 4.15).

```
public void scan1() {
    state.set(initialState); // Ausgangszustand
    while( !finalStates.member(state.current) ) {
        nextchar(); // nächstes Zeichen lesen
        state.set(transTab[state.current][input]); // neuer Zustand
    }
    if( state.marked ) writeback( input ); // zurückschreiben
}
```

Algorithmus 4.14: Ein überschüssiges Zeichen wird zurückgeschrieben.

```
public void scan2() {
    state.set(initialState);
    while( !finalStates.member(state.current) ) {
        state.set(transTab[state.current][input]); // neuer Zustand
        if( !state.marked ) nextchar(); // nächstes Zeichen
    }
}
```

Algorithmus 4.15: Ein Zeichen wird grundsätzlich vorausgelesen.

Die Hilfsprozedur `nextchar` liest das nächste Eingabezeichen in die Variable `input`. Die Prozedur in Alg. 4.15 erwartet im Gegensatz zu Alg. 4.14, daß vor ihrem Aufruf bereits ein Zeichen gelesen wurde.

Die angegebenen Prozeduren erkennen jeweils ein Symbol. Am Ende der Prozeduren kann daher die dem Endzustand `state` zugeordnete Aktionsfolge ausgeführt werden. Zur vollständigen Lexikalischen Analyse eines Quellprogramms sind die Prozeduren wiederholt aufzurufen. Bei Alg. 4.15 muß dann vor dem ersten Aufruf ein Zeichen gelesen werden.

Ausprogrammieren des Automaten

Aus einem Zustandsdiagramm (Tabelle) kann man ein Programm erzeugen, wobei Zustände durch Programmstellen repräsentiert werden, auf die in Abhängigkeit vom nächsten Eingabezeichen verzweigt wird. An denjenigen Stellen, die Endzuständen entsprechen, setzt man die zugehörigen Aktionen ein (Bsp. 4.16).

```
public Token scan3() {
    switch( input ) {
        case "0": ... case "9":
            while( isDigit( input )) nextchar();
            if( input == "." {
                nextchar();
                while( isDigit( input )) nextchar();
                return "rnum";
            } else return "inum";
            break;
        case ".":
            nextchar();
            if( isDigit( input ) {
                nextchar();
                while( isDigit( input )) nextchar();
                return "rnum";
            } else return "err";
            break;
        case "+": nextchar(); return "aop"; break;
        case "-": nextchar(); return "aop"; break;
        case "$": return "eof"; break;
        default: nextchar(); return "err";
    }
}
```

Beispiel 4.16: Aus dem Diagramm von Bsp. 4.13 programmierter Scanner. Zur besseren Übersicht ist hier nur die Ausgabe der Token (ohne Attribute) eingetragen.

Vor dem ersten Aufruf von `scan3` muß `nextchar()` aufgerufen werden, um das erste Zeichen bereitzustellen. Auch nach Erkennen der Operatoren (+,-) und beim letzten Fehlerzustand wird `nextchar()` aufgerufen, da diese Endzustände nicht markiert sind.

Eine vollständige ausprogrammierte Lexikalische Analyse nach dem obigen Muster findet man im Anhang A.

Kapitel 5

Syntax-Analyse

Die Syntax-Analyse (*parsing*) untersucht, ob die von der Lexikalischen Analyse gelieferte Symbolfolge der Grammatik der Programmiersprache entspricht und liefert Informationen über die Struktur des analysierten Programms. Die Attribute werden dabei nicht benötigt, sondern an die Semantische Analyse weitergereicht.

5.1 Grundlagen

Die Syntax von Programmiersprachen wird durch **kontextfreie Grammatiken** beschrieben. Die Grammatik spezifiziert nicht nur die Menge der syntaktisch richtigen Programme, sondern auch die *Struktur* dieser Programme.

-
- (1) $F \rightarrow id$
 - (2) $F \rightarrow id (L)$
 - (3) $L \rightarrow L , L$
 - (4) $L \rightarrow num$
 - (5) $L \rightarrow F$

Beispiel 5.1: Grammatik für Funktionsaufrufe in Backus-Naurform (BNF).

Eine kontextfreie Grammatik $G = (T, N, P, S)$ besteht aus:

Terminalen T (z.B. `id`, `(`, `)`, `...`). Sie können nicht weiter abgeleitet werden (daher der Name).

Nonterminalen N (z.B. `F` und `L`), mit großem Anfangsbuchstaben. Bei gut konstruierten Grammatiken entspricht jedes Nonterminal einer syntaktischen Klasse (in Bsp. 5.1 Funktionsaufruf bzw. Parameterliste).

Produktionen P (z.B. $L \rightarrow L, L$), auf ihrer linken Seite steht genau ein Nonterminal, die rechte Seite besteht aus einem Wort. Mehrere Produktionen für dasselbe Nonterminal können auch zusammengefaßt werden (z.B. $L \rightarrow L, L \mid num \mid id$).

Startsymbol S Wird das Startsymbol nicht explizit angegeben, so ist das Nonterminal der linken Seite der *ersten* Produktion das Startsymbol (im Bsp. 5.1: F).

Ein **Wort** ist eine beliebige Folge (Verkettung) von Nonterminalen und Terminalen, z.B. id L, num. Es wird mit griechischen Kleinbuchstaben (α, β, \dots) bezeichnet. Die leere Folge (ϵ) nennt man *Leerwort*. Eine **Ableitung** ist die Anwendung von Produktionen auf ein Wort derart, daß ein Nonterminal im Wort, das der linken Seite einer Produktion entspricht, durch deren rechte Seite ersetzt wird. Bei **Linksableitungen** wird jeweils das am weitesten links stehende Nonterminal des Wortes ersetzt.

Umgekehrt ist eine **Reduktion** die hintereinander ausgeführte Ersetzung von Teilworten, die rechten Seiten von Produktionen entsprechen, durch die zugehörigen Nonterminale. Die Schritte einer Linksreduktion entsprechen den Schritten einer Rechtsableitung in umgekehrter Folge.

| | |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| $\begin{aligned} F &\Rightarrow \\ \text{id} (\underline{L}) &\Rightarrow \\ \text{id} (\underline{L} , L) &\Rightarrow \\ \text{id} (\text{num} , \underline{L}) &\Rightarrow \\ \text{id} (\text{num} , \underline{F}) &\Rightarrow \\ \text{id} (\text{num} , \text{id}) & \end{aligned}$ | $\begin{aligned} \text{id} (\underline{\text{num}} , \text{id}) &\Leftarrow \\ \text{id} (L , \underline{\text{id}}) &\Leftarrow \\ \text{id} (L , \underline{F}) &\Leftarrow \\ \text{id} (\underline{L} , \underline{L}) &\Leftarrow \\ \underline{\text{id} (L)} &\Leftarrow \\ F & \end{aligned}$ |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|

Beispiel 5.2: Linksableitung $F \xRightarrow{*} \text{id} (\text{num} , \text{id})$ und Linksreduktion $\text{id} (\text{num} , \text{id}) \xleftarrow{*} F$ mit der Grammatik aus Bsp. 5.1. Die Teilworte, die im nächsten Schritt ersetzt werden, sind unterstrichen.

Das Ergebnis von Ableitungen bzw. Reduktionen kann durch einen **Ableitungsbaum** (*parse tree*) strukturiert dargestellt werden (siehe Abb. 5.3).

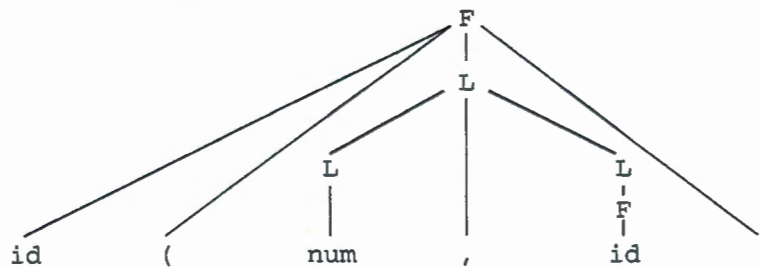


Abbildung 5.3: Ableitungsbaum zu Bsp. 5.2

Die Innenknoten entsprechen den Nonterminalen (Wurzel = Startsymbol). Ein Innenknoten mit seinen Nachfolgern entspricht der Anwendung einer Produktion.

Von links nach rechts gelesen ergeben die Blätter das abgeleitete Wort. Eine **Satzform** ist ein Wort α , das vom Startsymbol S aus ableitbar ist ($S \xRightarrow{*} \alpha$), z.B. id (L, L). Ein **Satz** ist eine *Satzform*, die nur aus Terminalen besteht, also das 'Endergebnis' einer Ableitung, z.B. id (num, id). Bei der Syntax-Analyse sind die Sätze die syntaktisch richtigen Programme. Die von einer Grammatik erzeugte **Sprache** $L(G)$ ist die Menge all ihrer Sätze. Eine Programmiersprache ist damit die Menge aller ableitbaren Programme.

Äquivalenz von Grammatiken

Zwei verschiedene Grammatiken G_1 und G_2 , die dieselbe Sprache L beschreiben, nennt man **äquivalent**. Aufgrund derartiger Äquivalenzen können mittels *Umformungen* bestimmte, für das gewählte Analyseverfahren (siehe Kap. 5.2) oder für nachfolgende Phasen (z.B. Semantische Analyse) ungünstige Eigenschaften von Grammatiken entfernt werden.

Alle Grammatiken in Bsp. 5.4 außer G_1 sind **eindeutig** (d.h. für jede Satzform gibt es genau einen Ableitungsbaum). G_1 ist **mehrdeutig**: es gibt Satzformen (z.B. num, num, num) mit verschiedenen Ableitungsbäumen (Abb. 5.5).

Die Grammatik G_2 erzeugt „linksrekursive“ Bäume (Abb. 5.6 links), G_3 jedoch „rechtsrekursive“ Bäume (Abb. 5.6 rechts). Beide eignen sich zum Beschreiben von *listenartigen* Strukturen, bei denen nur die *Reihenfolge* der Blätter wichtig ist. Bei *Ausdrücken* mit links- bzw. rechtsassoziativen Operatoren soll die Links- bzw. Rechtsrekursion verwendet werden – der Baum bildet dann die Struktur richtig ab.

Die Grammatik G_4 (Abb. 5.7 verwendet Iterationen, dargestellt mit regulären Ausdrücken. Sie ist anschaulich, jedoch kann im 'Ableitungsbaum' die Anzahl der Nachfolger eines Knoten statisch nicht bestimmt und die Struktur von Ausdrücken nicht erkannt werden.

Bei der Grammatik G_5 wird die Iteration mit einer Rekursion in einem zusätzlichen Nonterminal („Restliste“) aufgelöst – das erste Terminal (num bei L) bekommt eine Sonderstellung. In G_6 kommt es zu einer indirekten Rekursion. Mit beiden können Iterationen durch Knoten mit fester Nachfolgerzahl dargestellt werden (Abb. 5.8).

Regular Right Part Grammars (RRPG)

Die Syntax vieler Programmiersprachen wird unter Verwendung von Regulären Ausdrücken auf den rechten Seiten sehr einfach und klar beschrieben. Die syntaktischen Klassen sind deutlich erkennbar (in Bsp. 5.9 Ausdruck, Term, Faktor). Im Unterschied zur Regulären Definition sind Rekursionen erlaubt. Es folgt eine Gegenüberstellung dieser Konstrukte und gleichwertiger BNF-Konstrukte. Die Umformung ist automatisch durchführbar; die Iterationen werden dabei in Rechtsrekursionen mit neuen Zusatz-Nonterminalen aufgelöst.

| | | | | |
|-----------|-----------------------|------------|------------------------|--------------------------------------------------------|
| Iteration | $\dots\alpha * \dots$ | \implies | $\dots N \dots$ | mit neuer Prod. $N \rightarrow \alpha N \mid \epsilon$ |
| Iteration | $\dots\alpha + \dots$ | \implies | $\dots \alpha N \dots$ | mit neuer Prod. $N \rightarrow \alpha N \mid \epsilon$ |
| Option | $\dots\alpha ? \dots$ | \implies | $\dots N \dots$ | mit neuer Prod. $N \rightarrow \alpha \mid \epsilon$ |

| | |
|-------------------------------|-----------------------------------------------------------------------------------------|
| G1 (doppelte Rekursion): | $L \rightarrow L, L \mid \text{num}$ |
| G2 (Linksrekursion): | $L \rightarrow L, \text{num} \mid \text{num}$ |
| G3 (Rechtsrekursion): | $L \rightarrow \text{num}, L \mid \text{num}$ |
| G4 (reguläre Schreibweise): | $L \rightarrow \text{num} (, \text{num})^*$ |
| G5 (Restlistenrekursion): | $L \rightarrow \text{num } L_r$ $L_r \rightarrow , \text{num } L_r \mid \varepsilon$ |
| G6 (alternierende Rekursion): | $L \rightarrow \text{num } L_r$ $L_r \rightarrow , L \mid \varepsilon$ |

Beispiel 5.4: Äquivalente Grammatiken (Zeilen (3) und (4) aus Bsp. 5.1)

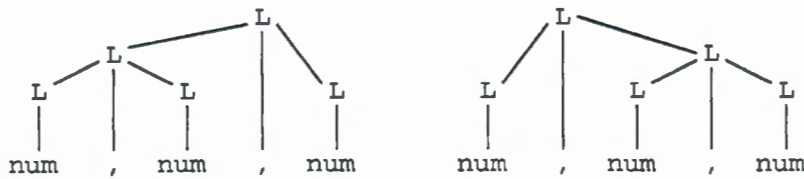


Abbildung 5.5: Ableitungsbäume für G1

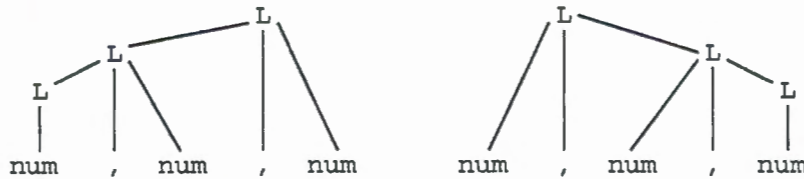


Abbildung 5.6: Ableitungsbäume für G2 und G3

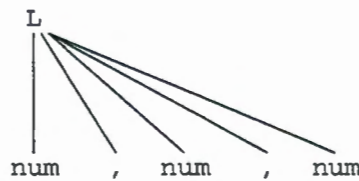


Abbildung 5.7: Ableitungsbaum für G4

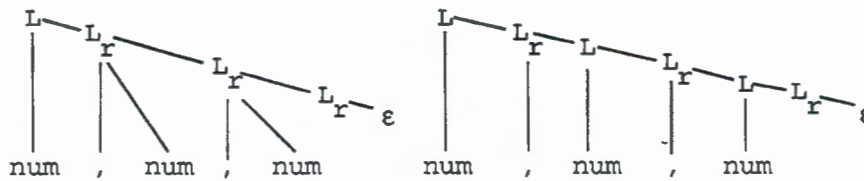


Abbildung 5.8: Ableitungsbäume für G5 und G6

Dabei besteht α aus einem Symbol oder einer geklammerten Folge von Symbolen, N ist jeweils ein neues Nonterminal.

$$\begin{array}{lcl}
 E \rightarrow T (+ T)^* & \Rightarrow & E \rightarrow T E_r \\
 & & E_r \rightarrow + T E_r \mid \varepsilon \\
 T \rightarrow F (* F)^* & \Rightarrow & T \rightarrow F T_r \\
 & & T_r \rightarrow * F T_r \mid \varepsilon \\
 F \rightarrow (E) \mid id & \Rightarrow & F \rightarrow (E) \mid id
 \end{array}$$

Beispiel 5.9: Automatische Umformung einer RRPg in eine BNF

In Bsp. 5.9 treten Klammern als Terminale und als Meta-Operatoren auf. Zur eindeutigen Unterscheidung werden Terminale daher oft zwischen Hochkommas (z.B. "+", "*", "(", etc.) geschrieben.

Eine andere Schreibweise von Regulären Ausdrücken auf der rechten Seite einer Produktion wird in der EBNF (erweiterte BNF) verwirklicht: Iterationen werden mit {}, Optionale mit [] beschrieben. Bekannt ist diese Schreibweise von der Modula-Spezifikation.

Operatoren und Trenner

In fast allen Programmiersprachen gibt es **Ausdrücke**, deren *Struktur*, beruhend auf der Assoziativität und Priorität ihrer **Operatoren**, wichtig ist (spätestens bei der Code-Erzeugung). Für solche Ausdrücke sollten die passenden Produktionen gewählt werden (z.B. Linksrekursion, s.o.).

Andererseits gibt es **Listen** und listenähnliche Konstrukte (z.B. Parameterlisten, Anweisungsfolgen), bei denen nur die *Reihenfolge* der einzelnen Elemente wichtig ist. Dabei werden die 'Operatoren' als **Trenner** bezeichnet, da sie nur zur Trennung der Elemente dienen, z.B. S;S;S.

Analyseverfahren

Die Grammatiken, die in erster Linie für die *Sprachbeschreibung* gedacht sind, können auch für die *Sprachanalyse* verwendet werden. Dabei wird untersucht, ob ein gege-

benes Wort in einer bestimmten Sprache enthalten ist, d.h. ob das Wort mit der Grammatik ableitbar ist; man sagt, das Wort wird akzeptiert. Ähnlich wie bei der Lexikalischen Analyse wird diese Aufgabe meist durch Automaten erledigt. Man kann sie als Interpreter sehen, die als (Steuer-)Programm eine Grammatik verwenden. Außerdem benötigen sie einen Keller zum Speichern der Geschichte der Ableitung (daher auch der Name **Keller-Automat**).

Um die Analyse zu vereinfachen, wird die Grammatik erweitert: Den Produktionen wird $S \rightarrow \text{Startsymbol } \$$ hinzugefügt (z.B. $(0) S \rightarrow F \$$ für Bsp. 5.1). S ist ein neues Startsymbol, das in der Grammatik sonst nicht vorkommt. $\$$ ist ein spezielles *eof-Symbol*, mit dem jedes Eingabewort abgeschlossen ist und das als zusätzliches neues Terminalsymbol gesehen werden kann. Diese einfache Erweiterung findet bei jedem Analysator statt, wird in der Folge aber nur dann explizit angeführt, wenn das zum besseren Verständnis nötig ist.

Für die im folgenden behandelten Analyseverfahren werden zwei grundlegende Funktionen benötigt:

$\text{First}(\alpha)$ gibt für eine Grammatik an, mit welchen Terminalsymbolen die aus α ableitbaren Satzformen beginnen können, bzw. ob α nach ϵ ableitbar ist. Die *First-Mengen* aller Grammatiksymbole (Terminale und Nonterminale) können gemeinsam iterativ berechnet werden (siehe Alg. 5.10).

$\text{First}(s_1 s_2 \dots s_k)$ wird äquivalent zur Produktion $N \rightarrow s_1 s_2 \dots s_k$ gebildet, wobei $\text{First}(N)$ durch $\text{First}(s_1 s_2 \dots s_k)$ zu ersetzen ist.

$\text{Follow}(N)$ gibt für eine Grammatik an, welche Terminalsymbole dem Nonterminal N in einer beliebigen Satzform unmittelbar folgen können, bzw. ob eof ($\$$) folgen kann. Die *Follow-Mengen* aller Nonterminale einer Grammatik können ebenfalls iterativ berechnet werden (siehe Alg. 5.11).

5.2 Top-Down-Analyse

Es wird versucht, mittels einer gegebenen Grammatik einen Ableitungsbaum für eine gegebene Symbolfolge zu finden. Dabei wird „**absteigend**“ (*top down*) vorgegangen (Abb. 5.12): Ausgehend vom Startsymbol, das die Wurzel des Ableitungsbaumes wird, und vom linken Rand der Symbolfolge werden solange **Linksableitungen** durchgeführt, bis die ganze Symbolfolge (d.h. das Eingabewort) abgeleitet wurde oder feststeht, daß sie nicht abgeleitet werden kann.

Tabellengesteuerte Top-Down-Analyse

Die Analyse wird durch einen allgemeinen Kontroll-Algorithmus (Alg. 5.14) mithilfe von Tabellen (wie in Bsp. 5.13) durchgeführt. Jeder Ableitungsschritt hängt vom *obersten* Kellersymbol und vom *vordersten* Eingabesymbol ab. Für jedes Nonterminal gibt es eine Zeile in der Tabelle. Sie gibt an, welche Produktion in Abhängigkeit von der Eingabe anzuwenden ist. Leereinträge (.) stellen Fehler dar.

```

∀ Terminalsymbole t: First(t) = {t}
∀ Nonterminale N: First(N) = {}
do {
  if( es gibt eine Produktion N → ε ) füge ε zu First(N) hinzu
  if( es gibt eine Produktion N → s1 s2 ... sk )
    füge die Terminalsymbole aus First(s1) zu First(N) hinzu
  i=1;
  while( i < k && ε ∈ First(si) ) {
    füge die Terminalsymbole aus First(si+1) zu First(N) hinzu
    i=i+1;
  }
  if( s1 s2 ... sk  $\xRightarrow{*}$  ε ) füge ε zu First(N) hinzu
} while( First-Menge hat sich geändert )

```

Algorithmus 5.10: Berechnung der First-Mengen aller Grammatiksymbole.

```

∀ alle Nonterminale N: DO Follow(N) = {}
Follow(Startsymbol) = {$}
do {
  if( es gibt eine Produktion A → α N oder A → N )
    füge die Elemente aus Follow(A) der Follow-Menge von N hinzu
  if( es gibt eine Produktion A → α N β oder A → N β {
    füge die Terminalsymbole aus First(β) zur Follow-Menge von N
    if( ε ∈ First(β) )
      füge die Elemente aus Follow(A) der Follow-Menge von N hinzu
  }
} while( Follow-Menge hat sich geändert )

```

Algorithmus 5.11: Berechnung der Follow-Mengen aller Nonterminale einer Grammatik.

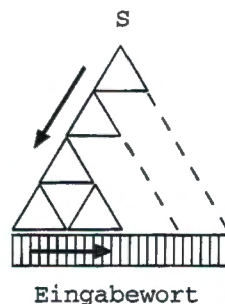


Abbildung 5.12: Top-Down-Analyse

Konstruktion der Tabelle. Die Top-Down-Tabelle wird durch die Matrix $M[N, t]$ mit den Indexbereichen Nonterminale = Kellersymbole (für die Zeilen) und Eingabesymbole (für die Spalten) dargestellt. Sie wird durch den Alg. 5.15 gefüllt. Bsp. 5.16 zeigt die Anwendung des Verfahrens.

Bei der Ableitung eines Nonterminalen N wird jene Produktion $N \rightarrow \alpha$ gewählt, die zum vordersten Eingabesymbol t hinführt. Es muß $t \in \text{First}(\alpha)$ sein; läßt sich allerdings α nach ϵ ableiten, so muß $t \in \text{Follow}(N)$ sein.

Rekursiver Abstieg (recursive descent)

Bei der Methode des Rekursiven Abstiegs (*recursive descent*) wird für jedes Nonterminal der Grammatik eine Prozedur gebildet (siehe Bsp. 5.17).

Als Optimierung läßt sich die Rekursion von R durch eine *Iteration* ersetzen und statt des Aufrufs direkt in E eintragen. Die Prozedur E hat dann das Aussehen:

```
... T(); while( sym == op ) { nextsym(); T(); }
```

Sie entspricht genau der RRP-Produktion $E = T (op T)^*$. So lassen sich RRP-Grammatiken ohne Umweg über BNF zu Parsern ausprogrammieren. Leider lassen sich nicht alle Grammatiken direkt mit den obigen Methoden behandeln. In der Tabelle kann es zu Mehrfacheinträgen kommen; im ausprogrammierten Parser kann der Fall auftreten, daß die Verzweigung für einen Wert von sym nicht eindeutig ist.

LL(1)-Grammatiken

Eine Grammatik, deren Top-Down-Tabelle keine Mehrfacheinträge hat, nennt man LL(1): sie analysiert Wörter von Links nach rechts mittels Linksableitung und einer Vorausschau in der Eingabe von 1 Symbol.

Eine Grammatik ist genau dann LL(1), wenn für die Alternativen α_i der Produktionen für jedes N gilt:

1. $\text{First}(\alpha_i) \cap \text{First}(\alpha_j) = \{\}$ f.a. i, j ($i \neq j$)

Grammatik:

$$\begin{aligned}
 E &\rightarrow T E_r \\
 E_r &\rightarrow + T E_r \mid \epsilon \\
 T &\rightarrow F T_r \\
 T_r &\rightarrow * F T_r \mid \epsilon \\
 F &\rightarrow (E) \mid id
 \end{aligned}$$

Tabelle:

| oberstes Kellersymbol | vorderstes Eingabesymbol | | | | | |
|--------------------------|--------------------------|---------|------------|------------|-----------|------------|
| | id | (|) | + | * | \$ |
| E | $T E_r$ | $T E_r$ | . | . | . | . |
| E_r | . | . | ϵ | $+ T E_r$ | . | ϵ |
| T | $F T_r$ | $F T_r$ | . | . | . | . |
| T_r | . | . | ϵ | ϵ | $* F T_r$ | ϵ |
| F | id | (E) | . | . | . | . |

Analyse des Wortes id + id * id \$:

| Keller | verbleibendes Eingabewort | anzuwendende Produktion |
|--------------------|---------------------------|----------------------------|
| \$ E | id + id * id \$ | $E \rightarrow T E_r$ |
| \$ E_r T | id + id * id \$ | $T \rightarrow F T_r$ |
| \$ E_r T_r F | id + id * id \$ | $F \rightarrow id$ |
| \$ E_r T_r id | id + id * id \$ | |
| \$ E_r T_r | + id * id \$ | $T_r \rightarrow \epsilon$ |
| \$ E_r | + id * id \$ | $E_r \rightarrow + T E_r$ |
| \$ E_r T + | + id * id \$ | |
| \$ E_r T | id * id \$ | $T \rightarrow F T_r$ |
| \$ E_r T_r F | id * id \$ | $F \rightarrow id$ |
| \$ E_r T_r id | id * id \$ | |
| \$ E_r T_r | * id \$ | $T_r \rightarrow * F T_r$ |
| \$ E_r T_r F * | * id \$ | |
| \$ E_r T_r F | id \$ | $F \rightarrow id$ |
| \$ E_r T_r id | id \$ | |
| \$ E_r T_r | \$ | $T_r \rightarrow \epsilon$ |
| \$ E_r | \$ | $E_r \rightarrow \epsilon$ |
| \$ | \$ | |

Beispiel 5.13: Tabellengesteuerte Top-Down-Analyse

```

for(;;) {
  sym = vorderstes Eingabesymbol
  switch( oberstes Kellersymbol ) {
    case $: if( sym == $ ) accept           // Analyse richtig beendet
           else Error ...                 // es gibt noch Eingabesymbole
           break;
    case Terminal t: if( t == sym ) skip    // verbrauche Terminal
                    else Error             // erwartetes Terminal nicht in der Eingabe
                    break;
    case Nonterminal N: if( table[N][sym] ==  $\alpha$  )
                        lösche N vom Keller
                        schreibe  $\alpha$  in umgekehrter Reihenfolge auf den Keller
                    else Error ...         // es gibt keine passende Produktion
  }
}

```

Algorithmus 5.14: Allgemeiner Kontroll-Algorithmus mit Tabellen.

Setze alle Einträge $M[N][t]$ auf error

Für jede Produktion $N \rightarrow \alpha$ und jedes Terminalsymbol t :

if(Terminalsymbol $t \in \text{First}(\alpha)$) füge α zu $M[N][t]$ hinzu

if($\epsilon \in \text{First}(\alpha)$ && $t \in \text{Follow}(N)$) füge α zu $M[N][t]$ hinzu

Algorithmus 5.15: Füllen der Top-Down-Tabelle.

2. höchstens ein α_i läßt sich nach ϵ ableiten
3. falls $\alpha_i \xRightarrow{*} \epsilon : \text{First}(\alpha_j) \cap \text{Follow}(N) = \{\}$ f.a. $j \neq i$

Die dritte Bedingung besagt: Läßt sich eine der Alternativen eines Nonterminals N zum Leerwort (ϵ) ableiten, so darf keine der anderen Alternativen für N mit einem Element aus $\text{Follow}(N)$ anfangen (siehe Bsp. 5.18).

Bei der Ableitung von $op\ id$ ist nicht eindeutig, welche Produktion von Opd_0 anzuwenden ist, wenn id das vorderste Eingabesymbol ist.

LL-Grammatiken haben den Vorteil, daß bei jedem Analyseschritt klar ist, welche Aktion als nächste zu treffen ist: welche Produktion anzuwenden ist, bzw. ob ein Eingabesymbol zu verbrauchen ist, bzw. ob die Analyse richtig oder falsch beendet ist. Die Analyse läuft *deterministisch* ab und ist daher zeiteffizient. LL(1)-Analysen sind viel speichereffizienter als solche mit mehr Vorausschau (LL(k) mit $k > 1$).

Grammatik:

$$\begin{array}{lll}
 (1) E \rightarrow T E_r & (4) T \rightarrow F T_r & (7) F \rightarrow (E) \\
 (2) E_r \rightarrow + T E_r & (5) T_r \rightarrow * F T_r & (8) F \rightarrow id \\
 (3) E_r \rightarrow \varepsilon & (6) T_r \rightarrow \varepsilon &
 \end{array}$$

First- und Follow-Mengen:

$$\begin{array}{ll}
 \text{First}(E) = \{id, (\} & \text{Follow}(E) = \{), \$ \} \\
 \text{First}(E_r) = \{+, \varepsilon \} & \text{Follow}(E_r) = \{), \$ \} \\
 \text{First}(T) = \{id, (\} & \text{Follow}(T) = \{+,), \$ \} \\
 \text{First}(T_r) = \{*, \varepsilon \} & \text{Follow}(T_r) = \{+,), \$ \} \\
 \text{First}(F) = \{id, (\} & \text{Follow}(F) = \{*, +,), \$ \}
 \end{array}$$

Einträge für das Nonterminal E: Zu betrachten ist nur die Produktion (1): $\text{First}(T E_r) = \{(, id\}$ – unter (und unter id wird die rechte Seite der Regel (1), d.h. $T E_r$, eingetragen. Die anderen Einträge der Zeile E (für), +, * und \$) bleiben leer – sie zeigen Syntaxfehler an.

Einträge für das Nonterminal E_r : Zuerst Produktion (2): $\text{First}(+ T E_r) = \{+\}$ – deshalb wird in $M[E_r, +]$ die rechte Seite ($+ T E_r$) eingetragen; dann (3): $\varepsilon \in \text{First}(\varepsilon)$ – also $\text{Follow}(E_r)$ bestimmen: das ist $\{), \$\}$. Also steht unter) und unter \$ die rechte Seite: ε . Die Einträge für die restlichen Nonterminale werden analog bestimmt.

Das Erstellen der Tabelle kann maschinell (mit einem Parser-Generator) erfolgen.

Beispiel 5.16: Erstellen der Tabelle von Bsp. 5.13.

Top-Down-Analyse mit nicht-LL(1)-Grammatiken

Steht oben im Keller stm und vorne in der Eingabe id (siehe Bsp. 5.19), so ist nicht klar, welche Produktion angewendet werden muß. Lösungsmöglichkeiten sind:

- Vorausschau um ein weiteres Eingabesymbol – also LL(2)-Analyse
- Verwenden semantischer Informationen (z.B. Symboltabelle) – das führt zu einer Vermischung mit der Semantischen Analyse
- Rücksetzen im Kontroll-Algorithmus

Rekursiver Abstieg mit Backtracking

Zum Speichern der Geschichte der bisherigen Analyse ist ein Keller nötig, der nicht nur die noch abzuleitenden Symbole (wie bei der tabellengesteuerten Analyse), sondern den *gesamten* Ableitungsbaum und in jedem Innenknoten die zuletzt versuchte

```

E → T R                                     // expression
R → op T R | ε                             // rest of expression
T → ( E ) | id                             // term

class RecursiveDescent {
    Token sym;
    public void nextsym() ... // setzt sym auf das nächste Eingabesymbol
    public void E() { // expression
        if( sym.member("id,(") ) { T(); R(); } else Error...
    }
    public void R() { // rest of expression
        if( sym.equals("op") ) { nextsym(); T(); R(); }
        else if( sym.member(",,$") /* ε (leeres statement) */
        else Error...
    }
    public void T() { // term
        if( sym.equals("(") ) { nextsym(); E(); nextsym(); } // ")"
        else if( sym.equals("id") ) nextsym();
        else Error...
    }
    void main() { // Recursive Descent
        nextsym();
        E();
    }
}

```

Beispiel 5.17: Recursive-descent-Parser für einen Teil der Grammatik von Bsp. 5.16

```

Opd → op Opd0 Opd | id
Opd0 → Opd | ε

```

```

First(Opd) ∩ Follow(Opd0) = {op,id} ≠ {}

```

Beispiel 5.18: Die folgende Grammatik für Präfix-Ausdrücke mit ein- und zweistelligen Operatoren ist nicht LL(1), da die 3. Bedingung nicht erfüllt wird.

| | |
|----------------------------------------------|-------------------------------------|
| (1) $\text{Stm} \rightarrow V = E$ | /* Zuweisung */ |
| (2) $\text{Stm} \rightarrow \text{id} (L)$ | /* Prozeduraufruf */ |
| (3) $\text{Stm} \rightarrow \text{id}$ | /* Prozeduraufruf ohne Parameter */ |
| (4) $V \rightarrow \text{id}$ | /* Variable */ |
| (5) $V \rightarrow \text{id} [L]$ | /* indizierte Variable */ |
| $E \rightarrow \dots$ | /* Ausdruck */ |
| $L \rightarrow \dots$ | /* Parameterliste */ |

Beispiel 5.19: Grammatik zur Beschreibung von Anweisungen

Alternative speichert. Derartige recursive-descent-Parser untersuchen prinzipiell alle Ableitungsmöglichkeiten (siehe Bsp. 5.20). Sie sind mit Sprachen, die Backtracking oder Parallelität erlauben (siehe Anhang B), besonders einfach zu schreiben, da der Keller bereits implizit in der Sprache vorhanden ist.

| | | |
|--------------------------------|-----------------------------------|---------------------------------------|
| $\text{Stm} \Rightarrow V = E$ | $\Rightarrow \text{id} = E$ | = paßt nicht \rightarrow rücksetzen |
| | $\Rightarrow \text{id} [L] = E$ | [paßt nicht \rightarrow rücksetzen |
| $\Rightarrow \text{id} (L)$ | | |

Beispiel 5.20: Analyse von $\text{id} (L) \$$ mittels Backtracking.

Bestimmte Grammatiken, die die LL_r -Bedingung nicht erfüllen, können durch folgende Algorithmen in $LL(1)$ -Grammatiken umgewandelt werden:

Linksfaktorisierung. Gleiche Anfänge der Alternativen eines Nonterminals werden „herausgehoben“, für den Rest wird ein neues Nonterminal eingeführt. Stehen dabei Nonterminale am Anfang, ist es manchmal nötig, sie durch ihre rechten Seiten zu ersetzen (siehe Bsp. 5.21).

Linksrekursion. Linksrekursive Grammatiken beschreiben linksassoziative Operatoren gut, sind aber für Top-Down-Verfahren nicht geeignet: bei Verfahren mit (beliebiger) Vorausschau ist der Schnitt der **First**-Mengen nicht leer, selbst Verfahren mit Rücksetzen fallen bei fehlerhaften Eingabeworten in Endlos-Linksrekursionen.

Linksrekursionen lassen sich auf folgende Weise in $LL(1)$ -Grammatiken umwandeln:

$$A \rightarrow A \alpha \mid \beta \implies \begin{aligned} A &\rightarrow \beta A_r \\ A_r &\rightarrow \alpha A_r \mid \epsilon \end{aligned}$$

Dabei können Probleme entstehen, wenn A auch in α oder β vorkommt.

$$\begin{array}{ll} V \rightarrow \text{id } V_r & V_r \rightarrow \varepsilon \\ & V_r \rightarrow [L] \end{array}$$

Auch die rechten Seiten der Produktionen für Stm haben gleiche Elemente in ihren **First-Mengen**. Zuerst wird V ersetzt, dann wird wieder faktorisiert:

$$\begin{array}{ll} \text{Stm} \rightarrow \text{id } \text{Stm}_r & \text{Stm}_r \rightarrow V_r = E \\ & \text{Stm}_r \rightarrow (L) \\ & \text{Stm}_r \rightarrow \varepsilon \end{array}$$

Nun ist die Grammatik LL(1). Die Produktionen für V werden nicht mehr gebraucht.

Beispiel 5.21: Linksfaktorisierung für V aus Bsp. 5.19.

Fehlerbehandlung

LL-Parser erkennen einen Fehler (bzw. sein Symptom) an der frühestmöglichen Stelle, d.h. wenn sich der Rest der Eingabe nicht mehr ableiten läßt (**valid-prefix-Eigenschaft**). Natürlich kann die Fehlerursache weitaus früher liegen: z.B. verursacht das Weglassen von $($ in $\text{id} * (\text{id} + \text{id})$ einen Fehler, der erst beim Lesen von $)$ erkannt wird. Im Gegensatz dazu können Verfahren mit Rücksetzen die Fehlerstelle nicht bestimmen.

Die einfachste Fehlerbehandlung ist der Abbruch der Analyse nach einer genauen Fehlermeldung. Tabellengesteuerte Verfahren unterstützen das mit Informationen aus ihrem Keller. Er enthält die 'erwarteten' Terminalsymbole bzw. Nonterminale und ermöglicht auch *phrase level recovery*, eine bessere Methode, die alle Fehler findet und 'verbessert', aber Teile überliest bzw. Folgefehler erzeugt.

Zur Vorbereitung setzt man alle **error**-Einträge jedes Nonterminalen mit einer ε -Alternative auch auf ε (siehe $(*)$ in Bsp. 5.23). Diese falschen ε -Ableitungen werden später abgefangen. Bei den anderen Nonterminalen N ändert man die **error**-Einträge für die Folgesymbole ($\in \text{Follow}(N)$) auf **stop**, die anderen auf **skip**. Die Analyse erfolgt durch Alg. 5.22 (vgl. Alg. 5.14).

Bei der ausprogrammierten LL-Analyse können die recursive-descent-Prozeduren ähnlich erweitert werden. Man kann dabei *dynamische* (d.h. vom Input abhängige) Informationen mittels Parametern weiterreichen.

Es gibt globale Strategien, die nicht nur die nächste Umgebung des Fehlers (vorderstes Eingabesymbol, oberstes Kellersymbol) betrachten. Dabei wird versucht, mit Kostenmatrizen die Fehlerbehebung zu optimieren.

5.3 Bottom-Up-Analyse

Mittels einer gegebenen Grammatik wird versucht, einen Ableitungsbaum für eine gegebene Symbolfolge zu finden. Dabei wird „aufsteigend“ (*bottom up*) vorgegan-

```

for(;;) {
  sym = vorderstes Eingabesymbol
  switch( oberstes Kellersymbol {
    case $:
      if( sym.equals("$") ) accept           // Analyse beendet
      break;
    case Terminal t:
      if( sym.equals(t) ) skip                // verbrauche Terminal t
      else lösche t vom Keller              // „Einfügen“ von t im Eingabewort
      break;
    case Nonterminal N:
      if( table[N][sym] ==  $\alpha$  ) ersetze N im Keller durch  $\alpha$ 
      else if( table[N][sym] == stop )
        lösche N vom Keller                // erzwingt die Ableitung von N
        else skip                          // Überlesen des Eingabesymbols
        break;
  }
}

```

Algorithmus 5.22: Analyse mit Fehlerbehandlung.

gen: Ausgehend vom Eingabewort (= die von der Lexikalischen Analyse gelieferte Symbolfolge) werden solange **Linksreduktionen** durchgeführt, bis zum Startsymbol (= Wurzel des Ableitungsbaumes) reduziert wurde oder feststeht, daß das nicht möglich ist. Dabei wird, umgekehrt wie bei der Top-Down-Analyse, jeder Knoten des Ableitungsbaumes *nach* seinen Nachfolgern erzeugt (Abb. 5.24).

Bottom-Up-Verfahren haben den Vorteil, daß die zugrundeliegende Grammatik auch Linksrekursionen und gleiche Anfänge enthalten darf. Umformungen (Faktorisierung, Elimination von Linksrekursionen) sind *nicht* nötig.

Grundbegriffe

Die Bottom-Up-Analyse heißt auch **shift-reduce-Analyse**. Die wesentlichen Aktionen des Parsers sind das „Schieben“ eines Eingabesymbols in den Keller (*shift*) und das Reduzieren der obersten Kellersymbole, die der rechten Seite einer Produktion entsprechen, zu einem Nonterminal (*reduce*). Im Keller stehen, im Gegensatz zur Top-Down-Analyse, die abgeleiteten Grammatiksymbole.

Bei jedem Analyseschritt stellen daher Kellerinhalt und verbleibendes Eingabewort aneinandergereiht (d.h. die Symbolfolge zwischen \$ und \$) eine *Satzform* dar, die mittels Rechtsableitungen gewonnen wird. Für eine entsprechende Rechtsableitung sind die Zeilen in Bsp. 5.25 von unten nach oben zu lesen.

Die unterstrichenen Teilworte in Bsp. 5.25 bezeichnet man als **Reduktionsansätze** (*handles*). Sie sind jene Symbolfolgen oben im Keller, die der rechten Seite

- (1) $E \rightarrow T E_r$ (4) $T \rightarrow F T_r$ (7) $F \rightarrow (E)$
 (2) $E_r \rightarrow + T E_r$ (5) $T_r \rightarrow * F T_r$ (8) $F \rightarrow id$
 (3) $E_r \rightarrow \varepsilon$ (6) $T_r \rightarrow \varepsilon$

Zugehörige Tabelle mit Fehler-Aktionen:

| oberstes Kellersymbol | vorderstes Eingabesymbol | | | | | |
|--------------------------|--------------------------|------------------|---------------|---------------|------------------|---------------|
| | id | (|) | + | * | \$ |
| E | $T E_r$ | $T E_r$ | stop | skip | skip | stop |
| E_r | $\varepsilon(*)$ | $\varepsilon(*)$ | ε | $+ T E_r$ | $\varepsilon(*)$ | ε |
| T | $F T_r$ | $F T_r$ | stop | stop | skip | stop |
| T_r | $\varepsilon(*)$ | $\varepsilon(*)$ | ε | ε | $* F T_r$ | ε |
| F | id | (E) | stop | stop | stop | stop |

Die Analyse von (+ id * + id \$ liefert (id * F + id) \$ als richtiges Wort:

| Keller | verbleibendes Eingabewort | Aktion |
|----------------------------------|---------------------------|------------------|
| \$ E | (+ id * + id \$ | |
| \$ E_r T | (+ id * + id \$ | |
| \$ E_r T_r F | (+ id * + id \$ | |
| \$ E_r T_r) E (| (+ id * + id \$ | |
| \$ E_r T_r) E | + id * + id \$ | skip über + |
| \$ E_r T_r) E | id * + id \$ | id ist richtiger |
| \$ E_r T_r) E_r T | id * + id \$ | Anfang von E |
| \$ E_r T_r) E_r T_r F | id * + id \$ | |
| \$ E_r T_r) E_r T_r id | id * + id \$ | |
| \$ E_r T_r) E_r T_r | * + id \$ | |
| \$ E_r T_r) E_r T_r F * | * + id \$ | |
| \$ E_r T_r) E_r T_r F | + id \$ | stop Ableitung |
| \$ E_r T_r) E_r T_r | + id \$ | von F erzwingen |
| ... | | |
| \$ E_r T_r) E_r | \$ | |
| \$ E_r T_r) | \$ |) einfügen |
| \$ E_r T_r | \$ | |
| \$ E_r | \$ | |
| \$ | \$ | |

Beispiel 5.23: Fehlerbehandlung mit LL(1)-Grammatik (vgl. Bsp. 5.16).

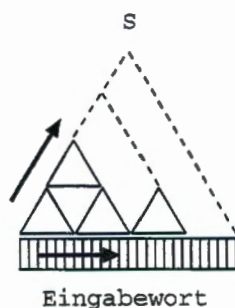


Abbildung 5.24: Bottom-Up-Analyse

$$\begin{array}{l}
 E \rightarrow E + T \mid T \\
 T \rightarrow T * F \mid F \\
 F \rightarrow (E) \mid id
 \end{array}$$

| Keller | verbleibendes Eingabewort | anzuwendende Aktion |
|----------------------|---------------------------|--------------------------------------|
| \$ | id + id * id \$ | shift |
| \$ <u>id</u> | + id * id \$ | reduce mittels $F \rightarrow id$ |
| \$ <u>F</u> | + id * id \$ | reduce mittels $T \rightarrow F$ |
| \$ <u>T</u> | + id * id \$ | reduce mittels $E \rightarrow T$ |
| \$ E | + id * id \$ | shift |
| \$ E + | id * id \$ | shift |
| \$ E + <u>id</u> | * id \$ | reduce mittels $F \rightarrow id$ |
| \$ E + <u>F</u> | * id \$ | reduce mittels $T \rightarrow F$ |
| \$ E + T | * id \$ | shift |
| \$ E + T * | id \$ | shift |
| \$ E + T * <u>id</u> | \$ | reduce mittels $F \rightarrow id$ |
| \$ E + <u>T * F</u> | \$ | reduce mittels $T \rightarrow T * F$ |
| \$ <u>E + T</u> | \$ | reduce mittels $E \rightarrow E + T$ |
| \$ E | \$ | accept (E ist Startsymbol) |

Beispiel 5.25: Eine shift-reduce-Analyse des Wortes $id + id * id \$$ mit einer linksrekursiven Grammatik (vgl. Bsp. 5.13)

einer Produktion entsprechen und deren Ersetzung einen *Schritt* der *entsprechenden Rechtsableitung* des Wortes darstellt.

In einer Satzform kann es Teilworte geben, die einer rechten Seite entsprechen, aber keine Reduktionsansätze sind (weil sie keinen Schritt in der Rechtsableitung darstellen). Im Bsp. 5.25 ist in der Satzform $E + T * id$ nicht $E + T$ der Ansatz, sondern id .

Eine Folge von Anfangssymbolen einer Satzform, die rechts nicht über den Ansatz hinausgeht, heißt **zulässiges Präfix** (*viabile prefix*). D.h. sie steht zur Gänze im Keller (E oder $E + id$ oder $E + T * id$ in Bsp. 5.25). Man kann jedes zulässige Präfix um Terminalsymbole ergänzen, sodaß eine (vollständige) Satzform entsteht. Das bedeutet, daß bei der Analyse kein Fehler entdeckt wird, solange man den Anfang der Eingabefolge zu einem zulässigen Präfix reduzieren kann.

Alle Sprachen, die *deterministisch* analysierbar sind, lassen sich durch **LR-Grammatiken** (Analyse von Links nach rechts mittels Rechtsableitungen (=Linksreduktionen)) beschreiben. Diese Grammatik-Klasse hat zudem die nützliche Eigenschaft, daß die LR-Grammatiken mit Vorausschau 1 genauso mächtig sind wie solche mit mehr Vorausschau. Das heißt, jede Sprache, die mit einer LR(k)-Grammatik (k beliebig) beschrieben werden kann, kann auch mit einer LR(1)-Grammatik beschrieben werden. LR(1)-Grammatiken sind bei der Analyse weitaus weniger aufwendig.

Bei einer LR(k)-Grammatik kann für jedes Eingabewort bei jedem Analyseschritt der Reduktionsansatz und damit die zu setzende Aktion aufgrund des Kellerinhalts und der nächsten k Eingabesymbole bestimmt werden. Weitere Informationen werden nicht benötigt. Hier sieht man (informell) den Grund für die größere Mächtigkeit der LR-Grammatiken: während bei LL(k)-Grammatiken die richtige Produktion allein mit ihren ersten k Terminalsymbolen bestimmt werden muß, stehen bei einer LR(k)-Grammatik die *ganze* rechte Seite der Produktion (im Keller) und noch weitere k Symbole zur Verfügung.

Weitere Vorteile der LR-Grammatiken sind, daß (mit ähnlichem Aufwand wie bei LL-Grammatiken) mächtigere Sprachen analysiert werden können, und daß die Erzeugung von tabellengesteuerten Parsern relativ einfach automatisiert werden kann.

LR-Analyse

Ein Kontroll-Algorithmus, der für alle LR-Parser gleich ist, verwendet Tabellen zur Durchführung der Analyse (Abb. 5.26). Die Tabellen bestimmen die zu analysierende Sprache und das genaue verwendete Analyseverfahren (LR, SLR, LALR, ...). Wir werden uns hier auf die Beschreibung des SLR-Verfahrens (*simple-LR*) beschränken.

Bei der Bottom-Up-Analyse wird generell im Keller das bisher analysierte Teilwort durch Zustände beschrieben.

Die Analysetabellen werden durch Zustände (Zeilen) und Grammatiksymbole inkl. \$ (Spalten) indiziert. Sie beinhalten Aktionen (action-Tabelle) bzw. Folgezustände (goto-Tabelle). Die Aktionen sind:

- shift (Eingabesymbol in den Keller schieben)
- reduce p (reduzieren; p ist die Nummer der Produktion)

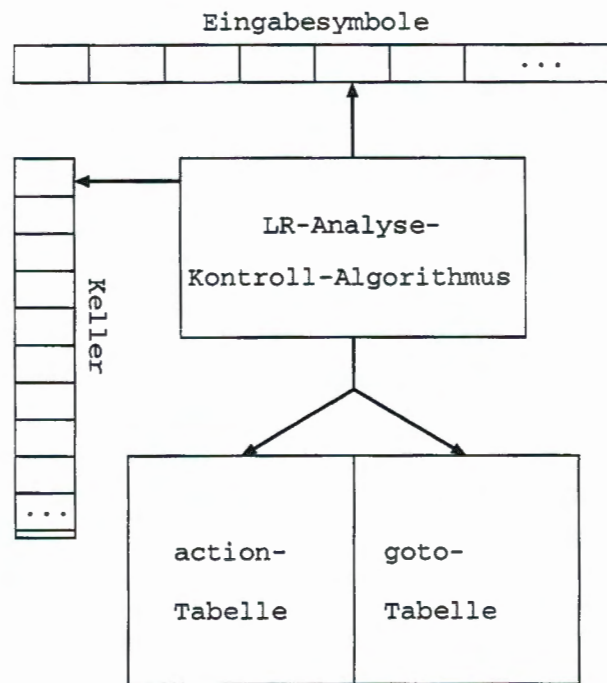


Abbildung 5.26: LR-Analyse-Modell

- accept (richtige Beendigung der Analyse)
- error (Fehlermeldung und -behandlung)

Generierung von Tabellen. Die Generierung der SLR-Tabellen wird anhand eines Beispiels vorgeführt. Zuerst wird die Grammatik durch Numerierung der Produktionen und Erweiterung um ein neues Startsymbol S vorbereitet (siehe Bsp. 5.27).

- (0) $S \rightarrow E$ (1) $E \rightarrow E \text{ op } T$ (3) $T \rightarrow (E)$
 (2) $E \rightarrow T$ (4) $T \rightarrow \text{id}$

Beispiel 5.27: Erweiterte Grammatik für einfache Ausdrücke

Generierung der goto-Tabelle. Jede Zeile der goto-Tabelle entspricht einem Zustand, repräsentiert durch eine Item-Menge. Ein **Item** ist eine Produktion, deren rechte Seite um eine Marke (".") erweitert ist: z.B.

$$T \rightarrow .\text{id} \text{ oder } E \rightarrow E \text{ op}.T$$

Die Marke trennt Keller und Eingabefolge, d.h. die bereits analysierten und noch zu analysierenden Symbole.

Enthält die Item-Menge I eines Zustands ein Item der Form $N \rightarrow \alpha . A \beta$ (Marke befindet sich unmittelbar links von einem Nonterminal), so sind auch alle Items der Form $A \rightarrow . \gamma$ in I enthalten (z.B. sind mit $E \rightarrow E \text{ op. } T$ auch $T \rightarrow . (T)$ und $T \rightarrow . id$ enthalten). Das entspricht der Tatsache, daß zuerst eine der rechten Seiten von A abgeleitet werden muß, bevor A reduziert werden kann. Die Funktion $\text{closure}(I)$ erweitert die Menge I zur vollständigen Item-Menge.

Ein Zustandsübergang findet statt, wenn die Marke über das Grammatiksymbol nach rechts geschoben wird: das entspricht der gerade erfolgten Reduktion eines Nonterminals bzw. dem Verbrauchen eines Terminalsymbols. Die Funktion $\text{goto}(I, G)$ liefert die closure-Menge jener Menge von Items, die von I aus über das Grammatiksymbol G erreichbar sind. Beschreibt I den zulässigen Präfix α , so beschreibt $\text{goto}(I, G)$ den zulässigen Präfix αG (siehe Alg. 5.28 und Bsp. 5.29).

```

I0 = closure({S → . Startsymbol }); M = {I0}
do {
  for( jede Item-Menge I in M und jedes Grammatiksymbol G ) {
    if( goto(I,G) != {} && goto(I,G) ∉ M )
      füge die Item-Menge goto(I,G) als In mit neuem Index n zu M hinzu
    if( goto(I,G) == Ik && I == Ij ) goto[j][G] = k;
  }
} while( neue Mengen wurden zu M hinzugefügt );

```

Algorithmus 5.28: Generierung der goto-Tabelle.

Generierung der action-Tabelle. Für die action-Tabelle werden die First- und Follow-Mengen der Grammatiksymbole und die goto-Tabelle benötigt. Das Verfahren ist in Alg. 5.30 zusammengefaßt. Bsp. 5.31 zeigt die Anwendung.

LR-Analyse-Algorithmus. Bei jedem Schritt wird in Abhängigkeit vom aktuellen Eingabesymbol a_i und vom obersten Kellerelement s_m eine Aktion $\text{action}[s_m, a_i]$ aus der action-Tabelle durchgeführt (siehe Alg. 5.32 und Bsp. 5.33).

Mehrdeutigkeiten

Bei bestimmten Sprachkonstrukten (Parameterlisten, etc.) ist die Struktur des Ableitungsbaumes unwichtig, sofern die Reihenfolge der Blätter gewahrt bleibt. Zu ihrer Beschreibung könnte man mehrdeutige Grammatiken verwenden und dem Analyse-Verfahren überlassen, welcher Baum generiert wird (siehe Abb. 5.5) .

Oft lassen sich auch andere Konstrukte (z.B. Ausdrücke) einfach mit mehrdeutigen Grammatiken beschreiben. Man braucht dann aber Zusatzregeln, mit denen die

| | |
|-------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| $I_0 = \{ \begin{array}{l} S \rightarrow .E \\ E \rightarrow .E \text{ op } T \\ E \rightarrow .T \\ T \rightarrow .(E) \\ T \rightarrow .id \end{array} \}$ | <p>Der Anfangszustand 0 enthält $S \rightarrow .E$, davon die <i>closure</i>-Menge I_0 bilden.</p> |
| $I_1 = \{ \begin{array}{l} S \rightarrow E. \\ E \rightarrow E. \text{ op } T \end{array} \}$ | <p>$I_1 = \text{goto}(I_0, E)$. Rechts von Marken keine Nonterminale \rightarrow keine <i>closure</i>-Operation.</p> |
| $I_2 = \{ E \rightarrow T. \}$ | <p>$I_2 = \text{goto}(I_0, T)$. Zustand 2 enthält nur fertige Items (die Marke steht ganz rechts) – hier muß reduziert werden.</p> |
| $I_3 = \{ \begin{array}{l} T \rightarrow (.E) \\ E \rightarrow .E \text{ op } T \\ E \rightarrow .T \\ T \rightarrow .(E) \\ T \rightarrow .id \end{array} \}$ | <p>$I_3 = \text{goto}(I_0, ($). Darin sind die Items von E enthalten, weil E im 1.Item unmittelbar nach einer Marke steht. Im Zustand 3 werden alle zulässigen Präfixe erkannt, die mit (enden, z.B. (, ((, E op (, ...</p> |
| $I_4 = \{ T \rightarrow id. \}$ | <p>$I_4 = \text{goto}(I_0, id)$. Jetzt sind alle Nachfolger des Zustands 0 bestimmt.</p> |
| $I_5 = \{ \begin{array}{l} E \rightarrow E \text{ op}.T \\ T \rightarrow .(E) \\ T \rightarrow .id \end{array} \}$ | <p>$I_5 = \text{goto}(I_1, \text{op})$.</p> |

Die Generierung wird mit den anderen Elementen von $M(I_2, I_3, \dots)$ fortgesetzt. Interessant sind $\text{goto}(I_3, T) = I_2$ und $\text{goto}(I_3, () = I_3$. Hier werden keine neuen Item-Mengen aufgebaut, sondern die bereits existierenden verwendet. Gleichzeitig mit der Generierung der Item-Mengen wird die *goto*-Tabelle Zeile für Zeile angelegt. Die Zustände werden dabei wie die Indizes der Item-Mengen nummeriert (I_0 – Zustand 0, I_1 – Zustand 1, ...).

Beispiel 5.29: Generierung der *goto*-Tabelle zu Bsp. 5.27

First- und Follow-Mengen der Grammatiksymbole bestimmen

\forall Zustände s und alle Terminalsymbole a (inkl. $\$$):

trage in $\text{action}[s][a]$ (a ist dabei das nächste Eingabesymbol) ein:

shift: wenn $\text{goto}(I_s, a) = I_j$

reduce p : wenn I_s das Item $N \rightarrow \alpha$ enthält und

$a \in \text{Follow}(N)$ und

die Produktion $N \rightarrow \alpha$ die Nummer p hat

accept: anstelle von reduce 0

error: sonst

Algorithmus 5.30: Generierung der action-Tabelle.

$\text{First}(S) = \{\text{id}, (\}$, $\text{First}(E) = \{\text{id}, (\}$, $\text{First}(T) = \{\text{id}, (\}$
 $\text{Follow}(S) = \{\$\}$, $\text{Follow}(E) = \{\$, \text{op}, \}$, $\text{Follow}(T) = \{\$, \text{op}, \}$

$\text{action}[0][\text{id}] = \text{shift}$, da $\text{goto}(I_0, \text{id}) = I_4$

$\text{action}[2][\text{op}] = \text{reduce } 2$, da Item $E \rightarrow T$. $\in I_2$ und
 $\text{op} \in \text{Follow}(E)$ und
 $E \rightarrow T$ die Produktion Nr.(2) ist.

$\text{action}[1][\$] = \text{accept}$, da $\$ \in \text{Follow}(S)$ und
Item $S \rightarrow E$. $\in I_1$

$\text{action}[0][\text{op}] = \text{error}$ (zur Kontrolle: $\text{op} \notin \text{First}(S)$)

vollständige Tabellen:

| state | action | | | | | goto | | | | |
|-------|--------|----|-----|----|-----|------|----|-----|---|-----|
| | id | op | () | \$ | | id | op | () | E | T |
| 0 | s | . | s | . | . | 4 | . | 3 | . | 1 2 |
| 1 | . | s | . | . | acc | . | 5 | . | . | . |
| 2 | . | r2 | . | r2 | r2 | . | . | . | . | . |
| 3 | s | . | s | . | . | 4 | . | 3 | . | 6 2 |
| 4 | . | r4 | . | r4 | r4 | . | . | . | . | . |
| 5 | s | . | s | . | . | 4 | . | 3 | . | 7 |
| 6 | . | s | . | s | . | . | 5 | . | 8 | . |
| 7 | . | r1 | . | r1 | r1 | . | . | . | . | . |
| 8 | . | r3 | . | r3 | r3 | . | . | . | . | . |

Beispiel 5.31: Generierung der action-Tabelle (zu Bsp. 5.29)

```

for(;;)
  if( action[sm][ai] == shift ) {
    verbrauche das Eingabesymbol ai
    schreibe den Inhalt von goto[sm][ai] in den Keller
  } else if( action[sm][ai] == reduce p ) {
    lösche soviel Elemente vom Keller, wie die rechte Seite von p Symbole hat
    sei s das neue oberste Kellersymbol und A die linke Seite von p,
    so schreibe den Inhalt von goto[s][A] auf den Keller
  } else if( action[sm][ai] == accept ) break;
else Error

```

Algorithmus 5.32: LR-Analyse-Algorithmus.

| Keller | verbleibende Eingabe | anzuwendende Aktion |
|---------|----------------------|------------------------------------------------|
| 0 | id op id op id \$ | s (* shift, push goto[0,id] = 4 *) |
| 0 4 | op id op id \$ | r4 (* T → id, 1x pop, push goto[0,T] *) |
| 0 2 | op id op id \$ | r2 (* E → T *) |
| 0 1 | op id op id \$ | s |
| 0 1 5 | id op id \$ | s |
| 0 1 5 4 | op id \$ | r4 (* T → id *) |
| 0 1 5 7 | op id \$ | r1 (* E → E op T, 3x pop, push goto[0,E] *) |
| 0 1 | op id \$ | s |
| 0 1 5 | id \$ | s |
| 0 1 5 4 | \$ | r4 (* T → id *) |
| 0 1 5 7 | \$ | r1 (* E → E op T *) |
| 0 1 | \$ | acc |

Beispiel 5.33: Analyse von id op id op id \$ (vgl. Bsp. 5.25).

Analyse eindeutig wird. Ergebnis der Syntax-Analyse soll in jedem Fall jener (einzige) Ableitungsbaum sein, der der gewünschten Struktur entspricht. Die mehrdeutige Grammatik

$$E \rightarrow E + E \mid E * E \mid (E) \mid id$$

kann z.B. mit den Zusatzregeln *left +*, *left ** Eingabewörter eindeutig analysieren. Die Zusatzregeln bedeuten, daß die Operatoren + und * linksassoziativ sind, und (aufgrund der Reihenfolge) daß * eine höhere Priorität als + besitzt. Auf diese Weise kann der Ausdruck $a+b*c$ richtig als $a+(b*c)$ analysiert werden.

Bei LR-Methoden entstehen durch mehrdeutige Grammatiken Mehrfach-Einträge in den action-Tabellen. Bei reduce/reduce-Konflikten ist nicht klar, mit welcher Produktion zu reduzieren ist; bei shift/reduce-Konflikten kann entweder sofort reduziert oder aber ein Eingabesymbol verbraucht werden.

Streichen der unerwünschten Einträge löst die Konflikte und macht die Analyse eindeutig (siehe Bsp. 5.34). Der Parser-Generator YACC erlaubt das automatische Streichen von Einträgen durch die Angabe von Zusatzregeln und durch Default-Annahmen. Bei shift/reduce-Konflikten streicht er den reduce-Eintrag (und wählt damit die 'längere' Variante); bei reduce/reduce-Konflikten bleibt der Eintrag für die früher definierte Produktion stehen.

So lassen sich auch *Spezialfälle* in Programmiersprachen einfach durch Streichen von mehrdeutigen Tabelleneinträgen behandeln, wenn der Sonderfall vor der allgemeinen Regel steht. Eine Grammatik zur Beschreibung mathematischer Formeln enthält z.B. die Produktionen

$$F \rightarrow F < F \mid F = F \mid \dots \mid F + F \mid \dots \mid id \mid num.$$

Um auch Variablen-Eingrenzungen darstellen zu können, wird als (längerer) Spezialfall als erste Alternative $F < id < F$ eingefügt. Die entstehenden Konflikte in den Tabellen werden automatisch richtig gelöst.

5.4 Hierarchie der Analyse-Verfahren

Das in Kapitel 5.3 vorgestellte Verfahren heißt **SLR(1)**: Simple-LR mit Vorausschau 1. Es kann zu Doppelseinträgen in der action-Tabelle führen, weil der SLR-Algorithmus zuwenig Information über die Geschichte der Ableitung mitführt. Die Verwendung der Follow-Mengen ist zu allgemein, da $Follow(N)$ alle Terminale enthält, die dem Nonterminal N folgen können. Diese Mengen sind unabhängig vom momentanen Zustand. Oft können aber einem Nonterminal N während einer bestimmten Ableitung (d.h. in einem bestimmten Zustand) nur einige der Follow-Symbole *tatsächlich* folgen.

Beim Generieren von **LR(1)**-Tabellen werden daher Items eingesetzt, die zusätzlich die ihnen tatsächlich folgenden Follow-Symbole (look-ahead-Mengen) beinhalten. Die Zustandsmenge des Automaten wird dadurch viel größer (viele Zustände sind in solche mit gleichen Items, aber verschiedenen look-ahead-Mengen aufgespalten); der Speicheraufwand, aber auch die Mächtigkeit steigt.

Aus der mehrdeutigen Grammatik für Ausdrücke

- (1) $E \rightarrow E + E$ (3) $E \rightarrow (E)$
 (2) $E \rightarrow E * E$ (4) $E \rightarrow id$

werden die folgenden SLR-Tabellen erstellt:

| handle | state | action | | | | | | goto |
|------------|-------|--------|-------|-------|-----|----|-----|------|
| | | id | + | * | () | \$ | E | |
| ϵ | 0 | s3 | . | . | s2 | . | . | 1 |
| E | 1 | . | s4 | s5 | . | . | acc | . |
| (| 2 | s3 | . | . | s2 | . | . | 6 |
| id | 3 | . | r4 | r4 | . | r4 | r4 | . |
| E + | 4 | s3 | . | . | s2 | . | . | 7 |
| E * | 5 | s3 | . | . | s2 | . | . | 8 |
| (E | 6 | . | s4 | s5 | . | s9 | . | . |
| E + E | 7 | . | s4/r1 | s5/r1 | . | r1 | r1 | . |
| E * E | 8 | . | s4/r2 | s5/r2 | . | r2 | r2 | . |
| (E) | 9 | . | r3 | r3 | . | r3 | r3 | . |

In den Zuständen 7 und 8 kommt es zu shift/reduce-Konflikten. Beide Operatoren sollen linksassoziativ sein – in $action[7,+]$ und $action[8,*]$ muß also reduziert werden (s4 bzw. s5 werden gestrichen). Wird im Zustand 7 das Symbol * gelesen, d.h. wurde $E + E *$ erkannt (siehe Spalte handle), so ist shift die richtige Aktion, da * die höhere Priorität hat und zuerst ausgeführt werden muß. Aus demselben Grund muß $action[8,+]$ = r2 gelten (das entspricht ... $E * E +$...).

Für eine eindeutige (und strukturell zutreffende) Analyse haben die Zeilen 7 und 8 also das Aussehen:

| handle | state | action | | | | | | goto |
|--------|-------|--------|----|----|-----|----|----|------|
| | | id | + | * | () | \$ | E | |
| E + E | 7 | . | r1 | s5 | . | r1 | r1 | . |
| E * E | 8 | . | r2 | r2 | . | r2 | r2 | . |

Die neue Tabelle ist auch mit LALR(1)- oder LR(1)-Verfahren nicht direkt (d.h. ohne Zusatzregeln) zu gewinnen.

Beispiel 5.34: Eindeutige Analyse durch das Streichen der unerwünschten Einträge aus der Tabelle.

Der hohe Speicherbedarf von LR(1)-Parsern wird vermieden, wenn man Zustände mit gleichen Items zusammenfaßt und dabei die Vereinigung ihrer look-ahead-Mengen bildet. Das führt zu Look-Ahead-LR-Analysatoren, deren Zustandszahlen der von SLR-Parsern entsprechen. Durch das Zusammenfügen der Zustände können reduce/reduce-, nicht jedoch shift/reduce-Konflikte entstehen. Bei richtigen Eingabeworten verhält sich der Parser wie der ursprüngliche LR-Parser, arbeitet aber wesentlich effizienter. Außerdem gibt es effiziente Verfahren, mit denen man LALR(1)-Tabellen direkt (also nicht über den Umweg der LR(1)-Zustände) aus einer Grammatik generieren kann.

Die Verfahren sind verschieden mächtig in Bezug auf die Grammatiken, die sie verarbeiten können.

$$\text{LR}(1) \supset \text{LALR}(1) \supset \text{SLR}(1)$$

Die weniger mächtigen Verfahren erzeugen u.U. Doppelseinträge beim Aufbau der Tabellen (wie z.B. SLR für Bsp. 5.19), die mächtigere Verfahren noch eindeutig behandeln können. Bei mehrdeutigen Grammatiken entstehen bei jedem LR-Verfahren Mehrfach-Einträge in den action-Tabellen.

$$\text{kontextfrei (auch mehrdeutig)} \supset \text{LR}$$

LR(1)-Parser können alle LL(1)-Grammatiken behandeln. Das gilt nicht für LALR und SLR, bei denen manche LL-Grammatiken mit ϵ -Produktionen zu Doppelseinträgen führen.

$$\text{LR}(1) \supset \text{LL}(1)$$

5.5 Wie schreibt man Grammatiken?

Das Erstellen von Grammatiken ähnelt sehr dem Schreiben von Programmen. Die zugrunde liegende Aufgabe läßt sich sehr oft auf mehrere Arten lösen. Man sollte versuchen, durch Aufspalten des Problems in kleine, überschaubare Module die Wiederverwendung von bewährten Teillösungen zu fördern.

Das Aufstellen einer Grammatik geht in drei Schritten vor sich:

1. Identifizieren von *syntaktischen Kategorien*, d.h. von wichtigen Bestandteilen der zu beschreibenden Sprache. Diesen Kategorien werden Nonterminale zugeordnet.
2. Suchen von Beziehungen zwischen syntaktischen Kategorien. Diese Beziehungen werden dann mit Grammatikregeln beschrieben. Eine Grammatikregel zeigt, wie sich die syntaktische Kategorie auf der linken Seite aus den syntaktischen Kategorien der rechten Seite zusammensetzt. Verschiedene Regeln mit gleicher linker Seite entsprechen verschiedenen Spezialisierungen der syntaktischen Kategorie der linken Seite.

3. Aus der *abstrakten Syntax* (das Ergebnis der ersten beiden Schritte) wird nun die *konkrete Syntax* festgelegt. Das betrifft vor allem die Reihenfolge der Symbole auf der rechten Seite einer Produktionsregel und eventuell zusätzliche Terminalsymbole, um die Struktur zu betonen und die Lesbarkeit zu verbessern.

Als Beispiel wollen wir eine Grammatik aufstellen, die Zahlen in verschiedenen Zahlensystemen beschreiben soll.

1. **Schritt:** Als Kandidaten für syntaktische Kategorien fallen die Basis und die Ziffernfolge ins Auge. Eine Ziffernfolge besteht naturgemäß aus Ziffern bzw. Einzelbuchstaben. Wir führen als Bezeichnungen die Nonterminale B, F und Z ein. Den ganzen Ausdruck bezeichnen wir mit A. Der Aufbau der Basis soll im Gegensatz zu den Ziffern der Ziffernfolge in der Grammatik nicht erklärt werden sondern als Einheit mit einem terminalen Symbol *num* betrachtet werden. Ihr Wertebereich liegt zwischen 2 und 36.

2. **Schritt:** Die Beziehungen zwischen den syntaktischen Kategorien sind leicht zu erstellen:

1. Ein Ausdruck A besteht aus einer Ziffernfolge F und der Basis B.
2. Eine Ziffernfolge F ist eine Aneinanderreihung von Ziffern Z.
3. Eine Ziffer ist eines der Zeichen '0' bis '9' oder einer der Buchstaben 'A' bis 'Z' (für Zahlensysteme bis zur Basis 36).
4. Eine Basis B ist eine Zahl *num*.

3. **Schritt:** Wir setzen die in den vorigen Schritten aufgestellte abstrakte Syntax aus Symbolen und Beziehungen nacheinander in konkrete Grammatikregeln um. Dabei wird die konkrete Syntax der Sprache festgelegt.

In unserem Beispiel muß die Ziffernfolge irgendwie von der Basis getrennt werden, damit wir diese beiden Gruppen von Zeichen auseinander halten können. Als Trennzeichen wählen wir das Zeichen '/'. Danach steht die Basis des Zahlensystems, z.B. A6E2/16, 10011/2, 755/8. Aus der ersten Beziehung ergibt sich die Regel

$$A \rightarrow F '/' B$$

Die Ziffernfolge aus der zweiten Beziehung wird z.B. durch

$$\begin{aligned} F &\rightarrow F Z \\ F &\rightarrow Z \end{aligned}$$

dargestellt.

Eine einfache Umsetzung der dritten Beziehung nimmt keine Rücksicht auf Klassen von Zeichen sondern listet ganz einfach alle Alternativen auf.

$$Z \rightarrow '0' \mid '1' \mid \dots \mid '9' \mid 'A' \mid 'B' \mid \dots \mid 'Z'$$

Die letzte Beziehung ergibt sich zu

$$B \rightarrow \text{num}$$

Die vollständige Grammatik hat damit die Form:

$$\begin{aligned} A &\rightarrow F \text{ '/' } B \\ F &\rightarrow F Z \\ F &\rightarrow Z \\ Z &\rightarrow '0' \mid '1' \mid \dots \mid '9' \mid 'A' \mid 'B' \mid \dots \mid 'Z' \\ B &\rightarrow \text{num} \end{aligned}$$

5.6 Qualitätskriterien für Grammatiken

Es gibt keine eindeutigen Kriterien für die Qualität von Grammatiken, wohl aber nützliche Stilregeln. Einige solcher Stilregeln sind:

- *Die Grammatik muß aufgabenmäßig abgefaßt sein.* Die syntaktischen Kategorien, die für die Verarbeitung von Bedeutung sind, müssen in der Grammatik festgelegt sein. Sind sie das nicht, dann läßt sich die Grammatik nur schwer attributieren. Die Beschreibung arithmetischer Ausdrücke hat z.B. als Forderung die unterschiedliche Priorität von Operatoren (Punktrechnung vor Strichrechnung). Sind diese Unterscheidungen schon in der Grammatik repräsentiert, dann vereinfacht sich die Attributierung ganz erheblich.
- *Die Grammatik muß analysegeeignet sein.* Dieser Punkt trifft besonders auf Grammatiken zu, die später mit Compilergeneratoren bearbeitet werden sollen. Für manche Syntaxanalyseverfahren sind z.B. Grammatiken mit linksrekursiven Symbolen nicht geeignet.
- *Grammatiken sollen eindeutig sein.* Die Eindeutigkeit erleichtert sowohl die Syntaxanalyse als auch die Zuordnung von Bedeutung durch die Attributierung.

Kapitel 6

Syntaxgesteuerte Übersetzung

In diesem Kapitel werden Methoden behandelt, die es erlauben, auf der Basis einer kontextfreien Grammatik nicht nur die Syntax-Analyse, sondern beliebige Übersetzungen zu definieren und zu implementieren. Die Methoden beruhen auf der **Attributierten Grammatik**. In einer Attributierten Grammatik werden den Grammatiksymbolen Attribute zugeordnet, und Übersetzungen werden als Berechnungen der Attribute beschrieben. Da die Berechnungen an die Grammatik gebunden sind, handelt es sich um eine **Syntaxgesteuerte Übersetzung**.

In ihrer allgemeinen Form kann man eine Attributierte Grammatik als „deklarative“ Spezifikation eines Compilers ansehen. Der Ablauf der Übersetzung wird nicht explizit, sondern nur implizit durch die Abhängigkeiten der Attribute definiert. Compiler-Generatoren können aufgrund der Attributabhängigkeiten den Ablauf der Übersetzung ermitteln und einen Compiler aus der Spezifikation generieren.

Unter bestimmten Voraussetzungen kann man eine Attributierte Grammatik auch als **Übersetzungsschema** formulieren. Dabei wird der Ablauf der Übersetzung explizit festgelegt. Ein Übersetzungsschema ist eine „prozedurale“ Spezifikation eines Compilers und als Eingabe für effiziente Compiler-Generatoren besonders geeignet. Wir werden zwei typische Vertreter solcher Generatoren behandeln.

6.1 Attributierte Grammatik

Eine **Attributierte Grammatik (AG)** ist eine kontextfreie Grammatik mit folgenden Erweiterungen:

- Jedem Grammatiksymbol können **Attribute** zugeordnet werden. Die Attribute sind in zwei Klassen eingeteilt, in **synthetisierte** (*synthesized*) und **ererbte** (*inherited*) Attribute.
- Jeder Produktion $A \rightarrow \dots X \dots$ können **Regeln** zugeordnet werden. Jede Regel berechnet ein Attribut und hat die allgemeine Form

$$b = f(c_1, \dots, c_k)$$

b ist ein **synthetisiertes** Attribut des Symbols A der linken Seite oder ein **ererbtes** Attribut eines Symbols X der rechten Seite.

c_1, \dots, c_k sind **ererbte** Attribute des Symbols der linken Seite oder **synthetisierte** Attribute der Symbole der rechten Seite.

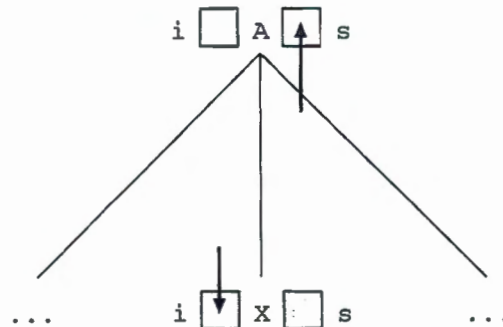


Abbildung 6.1: Attributberechnungen einer Produktion

Wir betrachten die Produktion $p: A \rightarrow \dots X \dots$ im Ableitungsbaum (Abb. 6.1), wobei angenommen ist, daß A und X jeweils ein ererbtes Attribut i und ein synthetisiertes Attribut s haben. Die Pfeile zeigen, welche Attribute durch Regeln der Produktion p berechnet werden. Die anderen Attribute (schraffiert) werden durch Regeln der angrenzenden Produktionen berechnet. Synthetisierte Attribute geben also entlang des Ableitungsbaums Informationen von unten nach oben weiter, ererbte Attribute von oben nach unten.

Wir verwenden im folgenden die Record-Schreibweise $Y.a$, um ein Attribut a des Symbols Y zu bezeichnen. Dabei wählen wir im allgemeinen den Attributnamen s für synthetisierte, i für ererbte und x für Attribute von Terminalen.

6.2 Wie attributiert man Grammatiken?

Für die Attributierung geht man von der kontextfreien Grammatik aus und geht wie folgt vor:

1. Identifizieren von *semantischen Kategorien*, d.h. von Eigenschaften bzw. Werten, die für die beabsichtigte Übersetzung bzw. Auswertung eine Rolle spielen. Zur Benennung dieser semantischen Kategorien werden Attribute eingeführt. Diese sind entweder synthetisierte oder ererbte Attribute, je nachdem, ob die betrachtete Information im Syntaxbaum von unten nach oben oder in umgekehrter Richtung fließen soll.

Dieser erste Schritt beeinflusst die entstehende Attributierung am nachhaltigsten. Man muß sich an dieser Stelle genau überlegen, welche Informationen man benötigt und in welcher Weise diese repräsentiert werden sollen.

2. *Zuordnung* von semantischen zu syntaktischen Kategorien, d.h. die Zuordnung von den Attributen zu den Nonterminalen. Es kann vorkommen, daß Nonterminale als „Zwischenlager“ für bestimmte Attribute eingesetzt werden müssen, die dort nicht benötigt werden. Nur so ist es möglich, Informationen weiter nach oben oder unten durchzureichen.
3. Bestimmen der *Attributauswertungsregeln*: Für jedes synthetisierte Attribut eines Nonterminals der linken Seite einer Regel und für jedes ererbte Attribut eines Nonterminals der rechten Seite einer Regel muß es eine Auswertungsregel geben.

Bei der Formulierung der Auswertungsregeln kann es notwendig sein, weitere Attribute einzuführen oder bestehende Attribute weiteren Nonterminalen zuzuordnen. Dies führt dann zu Iteration der Schritte 1 bis 3.

Synthetisierte und ererbte Attribute

Betrachten wir noch einmal unser Beispiel aus Kapitel 5.5. Wir wollen diese Grammatik attributieren, um den Wert einer Zahl zu berechnen.

1. **Schritt:** Als Ergebnis der Attributauswertung erwarten wir den Dezimalwert der Zahl in einem Attribut *val*. Bei der Betrachtung der Grammatik sehen wir, daß dieses Attribut Informationen in Richtung des Startsymbols zurückliefert, es wird also synthetisiert. Für die Basis brauchen wir auch ein Attribut *base*. Die Information wird jedoch bei der Berechnung des Wertes einer Ziffer gebraucht und muß daher als ererbtes Attribut zur entsprechenden Regel übergeben werden.
2. **Schritt:** Nun erfolgt die Zuordnung der semantischen zu den syntaktischen Kategorien. Alle Nonterminale erhalten das synthetisierte Attribut *val*. Die Ziffernfolge *F* erbt zusätzlich noch das Attribut *base*.
3. **Schritt:** Als letzten Schritt müssen wir alle konkreten Attributauswertungsregeln angeben. Dabei ist zu beachten, daß nur Attribute verwendet werden, die im jeweiligen Kontext auch definiert sind.

Wir beginnen mit den Definitionen der Regeln von unten, d.h. bei den Blättern. Die Reihenfolge ist rein willkürlich gewählt. Sie folgt nur dem Fluß der synthetisierten Attribute, die in unserem Beispiel häufiger auftreten.

Die Attribute werden nicht explizit deklariert, sondern implizit durch die Regeln eingeführt. Die Attributnamen sind frei wählbar und lokal für ein Grammatiksymbol gültig. Wenn in einer Produktion ein Grammatiksymbol mehrmals vorkommt (z.B. *F*), so wird zur Unterscheidung indiziert, damit sich die Regeln eindeutig auf die unterschiedlichen Vorkommen des Symbols in der Produktion beziehen können.

Den trivialen Abhängigkeiten der Regeln

$$\begin{aligned}
 F &\rightarrow Z \\
 Z &\rightarrow '0' \mid '1' \mid \dots \mid '9' \mid 'A' \mid 'B' \mid \dots \mid 'Z' \\
 B &\rightarrow \text{num}
 \end{aligned}$$

entsprechen häufig „Transferregeln“ wie

$$\begin{aligned}
 F.\text{val} &= Z.\text{val} \\
 Z.\text{val} &= 0 \\
 Z.\text{val} &= 1 \\
 &\vdots \\
 Z.\text{val} &= 35 \\
 B.\text{val} &= \text{value}(\text{num}.x)
 \end{aligned}$$

die Attributwerte unverändert weiterreichen bzw. bei Terminalsymbolen die Werte erst zur Verfügung stellen¹. Alle synthetisierten Attribute der linken Seiten haben damit einen Wert erhalten. Auf den rechten Seiten gibt es keine ererbten Attribute, denen ein Wert zugewiesen werden müßte.

Die eigentliche Berechnung des Zahlenwertes geschieht in der Regel

$$F_1 \rightarrow F_2 Z$$

Der Wert für jede neue Ziffernfolge ergibt sich aus der alten Ziffernfolge multipliziert mit der Basis des Zahlensystems. Zu diesem Wert wird die neue Ziffer addiert. Die Basis steht im ererbten Attribut $F_1.\text{base}$. Damit und mit den restlichen synthetisierten Attributen der rechten Seite können wir die Attributierungsregel für das einzige synthetisierte Attribut der linken Seite formulieren. Fehlt nur noch die Regel für das einzige ererbte Attribut auf der rechten Seite. $F_2.\text{base}$ bekommt den Wert von $F_1.\text{base}$.

Wir erhalten

$$\begin{aligned}
 F_1.\text{val} &= F_2.\text{val} * F_1.\text{base} + Z.\text{val} \\
 F_2.\text{base} &= F_1.\text{base}
 \end{aligned}$$

In der ersten Produktion müssen wir jetzt noch den Wert der Basis zur Verfügung stellen und die errechnete Zahl an das Startsymbol übergeben. Es ergeben sich die Regeln

$$\begin{aligned}
 A.\text{val} &= F.\text{val} \\
 F.\text{base} &= B.\text{val}
 \end{aligned}$$

Die vollständige attributierte Grammatik hat damit das folgende Aussehen:

¹Z.B. übersetzt die Funktion `value()` die Textrepräsentation (die im Attribut `x` abgelegt ist) der Zahl des Symbols in einen numerischen Wert.

| Produktion | Regeln |
|----------------------------------|----------------------------------------------------------------------|
| $A \rightarrow F \text{ '/' } B$ | $A.val = F.val$ $F.base = B.val$ |
| $F_1 \rightarrow F_2 Z$ | $F_1.val = F_2.val * F_1.base$ $+ Z.val$ $F_2.base = F_1.base$ |
| $F \rightarrow Z$ | $F.val = Z.val$ |
| $Z \rightarrow '0'$ | $Z.val = 0$ |
| $Z \rightarrow '1'$ | $Z.val = 1$ |
| \vdots | \vdots |
| $Z \rightarrow 'Z'$ | $Z.val = 35$ |
| $B \rightarrow \text{num}$ | $B.val = \text{value}(\text{num.x})$ |

Die durch eine attributierte Grammatik spezifizierte Übersetzung kann man exemplarisch durch einen **attributierten Ableitungsbaum** darstellen. Dabei werden die Attributwerte in die Ableitungsknoten eingetragen. Jede Übersetzung liefert genau eine Wertebelegung der Knoten. Abb. 6.2 zeigt den attributierten Ableitungsbaum für die Umwandlung von "3F7/16" in die Dezimalzahl 1015. Zur besseren Übersicht sind nur die synthetisierten Attribute dargestellt.

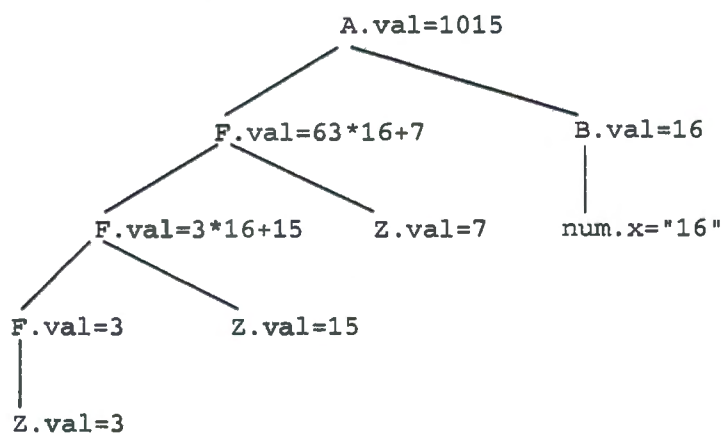


Abbildung 6.2: Attributierter Ableitungsbaum für die Umwandlung von "3F7/16" in die Dezimalzahl 1015. Nur synthetisierte Attribute sind dargestellt.

In Abb. 6.3 ist der **Abhängigkeitsgraph** des Ableitungsbaums dargestellt. Die Knoten des Abhängigkeitsgraphen sind als Punkte neben den Ableitungsknoten eingetragen. Sie repräsentieren beliebige Attributwerte. Ererbte Attribute stehen dabei links und synthetisierte Attribute rechts. Der Abhängigkeitsgraph des gesamten Ableitungsbaums wird aus den Abhängigkeitsgraphen der angewandten Produktionen zusammengesetzt.

Abb. 6.3 veranschaulicht den Informationsfluß bei der Ableitung der dreistelligen Hexadezimalzahl "3F7/16". In diesem Beispiel sieht man, wie ererbte Attribute dazu verwendet werden, Informationen „hinunter“ in die Zweige des Baumes zu liefern.

Die Basis der umzuwandelnden Zahl muß ja für die Berechnung des Dezimalwertes bereitgestellt werden.

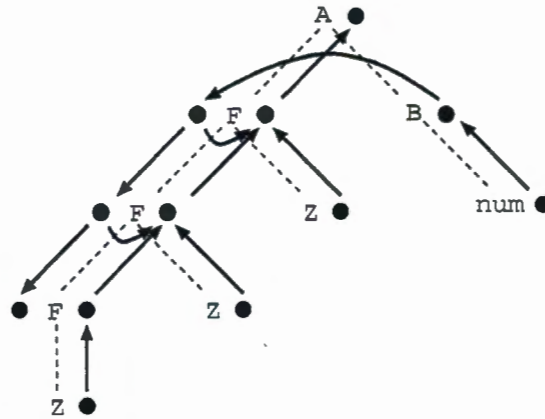


Abbildung 6.3: Abhängigkeitsgraph bei der Berechnung einer dreistelligen Zahl.

6.3 Qualitätskriterien für die Attributierung

Wir ergänzen die Stilregeln für Grammatiken um solche für Attributierungen:

- *Die Attributierung soll einfach und natürlich formuliert sein.* Dies empfiehlt sich wegen des Umfangs der Attributierung, der oft den der zugrundeliegenden Grammatik um ein Vielfaches übertrifft.
- *Die Attributierung muß vollständig sein,* d.h. in jedem Kontext muß für jedes zu berechnende Attribut eine Attributierungsregel zur Verfügung stehen.
- *Die Attributierung muß zyklensfrei sein.* Kein Attribut darf (direkt oder indirekt) von sich selbst abhängen.

6.4 S-Attributierte Grammatik

Eine AG, die nur synthetisierte Attribute enthält, heißt **S-attributiert**.

Die Nonterminale E und T in Bsp. 6.4 besitzen jeweils ein synthetisiertes Attribut s, in dem der Postfix-Code (P-Code) als String „berechnet“ wird.

Die Terminale op und id besitzen ein Attribut x, dessen Wert von der Lexikalischen Analyse geliefert wird (Tokenattribut). Wir nehmen hier an, daß die Operanden und Operatoren des zu übersetzenden Infix-Ausdrucks aus jeweils nur einem Zeichen bestehen und daß dieses Zeichen als Tokenattribut übergeben wird.

Die den Produktionen zugeordneten Regeln können als „Übersetzungsgleichungen“ gelesen werden:

| Produktion | Regeln |
|---------------------------------------|-----------------------------------|
| (1) $E \rightarrow E_1 \text{ op } T$ | $E.s = E_1.s + T.s + \text{op}.x$ |
| (2) $E \rightarrow T$ | $E.s = T.s$ |
| (3) $T \rightarrow (E)$ | $T.s = E.s$ |
| (4) $T \rightarrow \text{id}$ | $T.s = \text{id}.x$ |

Beispiel 6.4: Eine AG zur Übersetzung von Infix- in Postfix-Ausdrücke

1. Der P-Code für E ist die Verkettung der P-Codes für E_1 und T mit $\text{op}.x$.
2. Der P-Code für E ist der P-Code für T.
3. Der P-Code für T ist der P-Code für E (Klammern werden eliminiert).
4. Der P-Code für T ist $\text{id}.x$.

Abb. 6.5 zeigt auf der linken Seite den attributierten Ableitungsbaum für die Übersetzung von "a+b" in "ab+". Die eingetragenen Attributwerte sind als Strings aufzufassen.

Rechts in Abb. 6.5 ist der **Abhängigkeitsgraph** des Ableitungsbaums dargestellt. In diesem Beispiel erfolgt die Berechnung aller Attribute im Ableitungsbaum von unten nach oben.

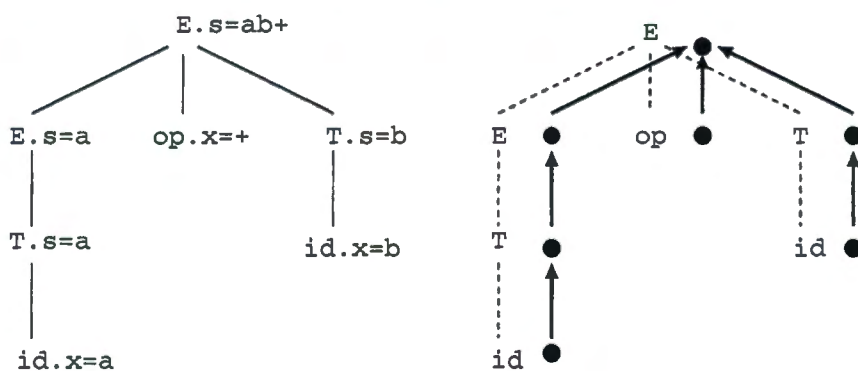


Abbildung 6.5: Attributierter Ableitungsbaum und Abhängigkeitsgraph

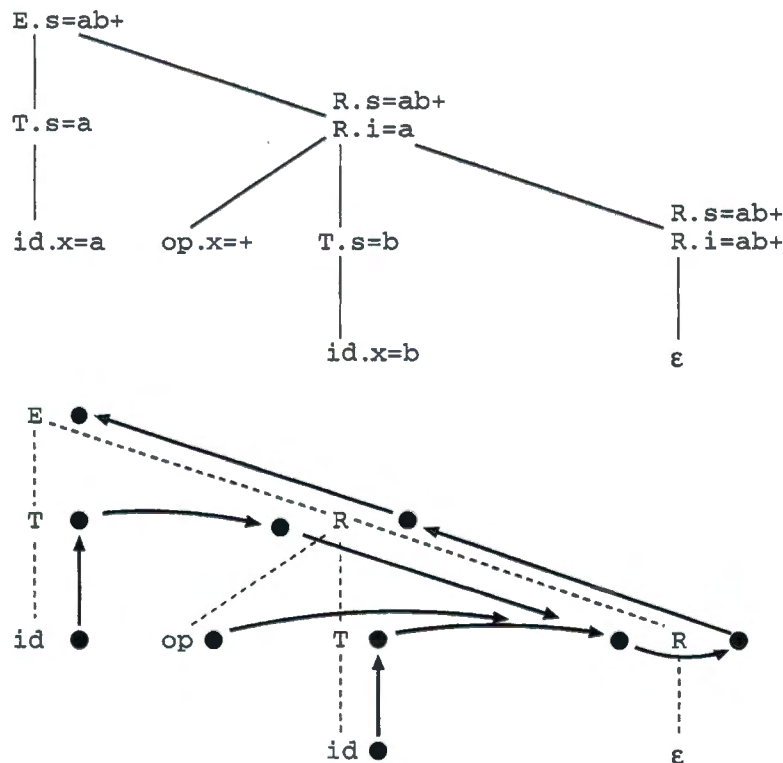
Durch einfache Modifikation kann man die AG von Bsp. 6.4 auch für andere Zwecke verwenden, z.B. würde durch $\text{op}.x + E_1.s + T.s$ ein **Prefix-Code** erzeugt. Wenn die Operanden nicht symbolisch, sondern als Werte gegeben sind, kann man anstelle der Verkettung zum Beispiel eine Funktion $\text{apply}(\text{op}.x, E_1.s, T.s)$ aufrufen, die den Operator auf die Operanden anwendet, d.h. man erhält die Spezifikation eines **Interpreters** für Infix-Ausdrücke.

Ereberte Attribute

Grundsätzlich dienen ererbte Attribute dazu, Werte im Ableitungsbaum von oben nach unten zu geben. Sie werden z.B. bei LL-Grammatiken wegen der Ersetzung von Linksrekursion durch Rechtsrekursion benötigt (Bsp. 6.6 und Kapitel 6.5).

| Produktion | Regel |
|-----------------------------|---------------------------------------------|
| $E \rightarrow T R$ | $R.i = T.s$ $E.s = R.s$ |
| $R \rightarrow op T R_1$ | $R_1.i = R.i + T.s + op.x$ $R.s = R_1.s$ |
| $R \rightarrow \varepsilon$ | $R.s = R.i$ |
| $T \rightarrow (E)$ | $T.s = E.s$ |
| $T \rightarrow id$ | $T.s = id.x$ |

Attributierter Ableitungsbaum und Abhängigkeitsgraph (ererbte Attribute stehen links, synthetisierte rechts vom Ableitungsknoten)



Beispiel 6.6: Attributierte LL-Grammatik mit ererbten Attributen zur Übersetzung von Infix in Postfix.

6.5 L-Attributierte Grammatik

Eine AG heißt **L-attribuiert** (LAG, genauer: LAG(1)), wenn alle Attribute in einem links-rechts Tiefendurchlauf (*left-right depth-first*) des Ableitungsbaums berechnet werden können, entsprechend der rekursiven Prozedur Alg. 6.7.

```
public void dfvisit(Node n) {
  for( alle direkten Nachfolger m von links nach rechts ) {
    berechne die ererbten Attribute von m
    dfvisit(m);
  }
  berechne die synthetisierten Attribute von n
}
```

Algorithmus 6.7: Berechnung der Attribute in einem links-rechts Tiefendurchlauf des Ableitungsbaums.

Aus der Prozedur `dfvisit` in Alg. 6.7 kann man ersehen, daß eine L-attributierte Grammatik folgende Bedingung erfüllen muß:

In jeder Produktion $A \rightarrow \dots X \dots$ dürfen die ererbten Attribute von X nur von Attributen jener Symbole abhängen, die **links** von X stehen (d.h. von ererbten von A und synthetisierten von anderen Symbolen links von X).

Da die einschränkende Bedingung der L-attribuierten Grammatik sich nur auf ererbte Attribute bezieht, ist jede S-attribuierte Grammatik auch L-attribuiert.

Die S-attribuierte Grammatik von Bsp. 6.4 ist also L-attribuiert und ebenfalls die Grammatik von Bsp. 6.6, da sie die obige Bedingung erfüllt. Man erkennt an den Pfeilen im Abhängigkeitsgraphen von Bsp. 6.6, daß die Attributberechnung in einem Tiefendurchlauf des Ableitungsbaums erfolgen kann.

Bsp. 6.8 zeigt eine nicht L-attribuierte Grammatik. Es handelt sich um die Übersetzung Boolescher Ausdrücke, wobei der linke Teilausdruck E_1 Information über die Art des Operators benötigt.

| Produktion | Regel |
|-----------------------------------|--------------------------|
| $E \rightarrow E_1 \text{ op } T$ | $E_1.i = f(\text{op}.x)$ |
| ... | |

Beispiel 6.8: Ausschnitt aus einer nicht-L-attribuierten Grammatik

6.6 Übersetzungsschema

Eine L-attributierte Grammatik kann man als **Übersetzungsschema** schreiben. Dabei werden die Regeln jeder Produktion geordnet und so in die rechte Seite der Produktion eingesetzt, daß sie beim Tiefendurchlauf des Ableitungsbaums von links nach rechts ausgeführt werden können. Die eingesetzten Regeln bezeichnet man als (semantische) **Aktionen**.

$$\begin{aligned}
 E &\rightarrow T \{R.i = T.s\} R \{E.s = R.s\} \\
 R &\rightarrow op T \{R_1.i = R.i + T.s + op.x\} R_1 \{R.s = R_1.s\} \\
 R &\rightarrow \varepsilon \{R.s = R.i\} \\
 T &\rightarrow (E) \{T.s = E.s\} \\
 T &\rightarrow id \{T.s = id.x\}
 \end{aligned}$$

Beispiel 6.9: LL-Übersetzungsschema der AG von Bsp. 6.6

In einem Übersetzungsschema muß jede Aktion, die ein ererbtes Attribut eines Symbols berechnet, links dieses Symbols stehen. Die Aktion darf nur Attributwerte von Symbolen verwenden, die links dieser Aktion stehen. Aktionen zur Berechnung des synthetisierten Attributs der linken Seite können generell am Ende der Produktion geschrieben werden.

Die in Bsp. 6.8 gezeigte Produktion kann man nicht als Übersetzungsschema schreiben, da die Aktion $E_1.i = f(op.x)$ zur Berechnung des ererbten Attributs vor E_1 plaziert werden muß. Die Aktion verwendet aber den Attributwert $op.x$ des Symbols op , das nicht links dieser Aktion steht.

Bei einer allgemeinen AG (z.B. in Bsp. 6.6) wird nicht spezifiziert, in welcher Reihenfolge die Regeln der Attributberechnung auszuführen sind. Die Reihenfolge ergibt sich nur implizit aus den Abhängigkeiten der Attribute. Es handelt sich also um eine „deklarative“ Spezifikation. Beim Übersetzungsschema wird dagegen die Ausführungsfolge der Regeln in Form von Aktionen explizit angegeben. Es handelt sich hier um eine „prozedurale“ Spezifikation, und es ist entsprechend einfach, aus einem Übersetzungsschema einen Compiler zu generieren.

6.7 Top-Down-Compiler-Generator

Aus jedem Übersetzungsschema, das auf einer LL-Grammatik basiert, kann man einen Ein-Pass-Top-Down-Compiler erzeugen, denn ein Übersetzungsschema setzt voraus, daß die Grammatik L-attribuiert ist. Die Attribute können daher in einem links-rechts Tiefendurchlauf berechnet werden, der genau dem Ablauf der Top-Down Syntax-Analyse entspricht.

Eingabe: Übersetzungsschema auf der Basis einer LL(1)-Grammatik

Ausgabe: Compiler; der nach dem Verfahren des rekursiven Abstiegs arbeitet.

Algorithmus: Der Generator erzeugt einen LL-Parser nach der Methode des rekursiven Abstiegs (siehe Kap. 5.2), erweitert um Parameter und Aktionen.

1. Für jedes Nonterminal A der Grammatik wird eine Funktion generiert, die für jedes ererbte Attribut von A einen Eingangs-Parameter hat und den Wert des synthetisierten Attributs von A als Funktionswert zurückgibt. Für jedes Attribut eines jeden Symbols, das in Produktionen von A auf der rechten Seite auftritt, wird eine lokale Variable deklariert.
2. Zur Behandlung der rechten Seite einer Produktion sind folgende drei Fälle zu unterscheiden:
 - (a) Für jedes Terminal t mit einem synthetisierten Attribut wird eine Wertzuweisung an die für das Attribut $t.s$ deklarierte lokale Variable generiert.
 - (b) Für jedes Nonterminal N wird ein Aufruf $s = N(i_1, \dots, i_n)$ generiert. Dabei sind i_1, \dots, i_n die lokalen Variablen für die ererbten Attribute von N , und s ist die lokale Variable für das synthetisierte Attribut von N .
 - (c) Für jede Aktion wird der Code der Aktion eingesetzt, wobei jede Referenz auf ein Attribut durch die lokale Variable für dieses Attribut ersetzt wird.

Falls ein Nonterminal mehrere synthetisierte Attribute besitzt, können diese zu einem Record zusammengefaßt werden. Dies ist bei einem Übersetzungsschema immer möglich, weil alle Attribute in einem Durchlauf ausgewertet werden (siehe Bsp. 6.10).

Optimierung. Die erzeugten Funktionen E und R lassen sich vereinfachen. Zunächst kann man alle unnötigen lokalen Variablen eliminieren (siehe Bsp. 6.11).

Die Prozedur R kann wegen der Endrekursion als Schleife in E eingesetzt werden (vgl. Kap. 5.2). Man beachte aber, daß der erste Aufruf von T vor der Schleife stehen muß, da die Verkettung nur bei nachfolgenden Aufrufen von T (innerhalb der Schleife) durchzuführen ist (Bsp. 6.12).

Die optimierte Version entspricht auf Grammatikebene einer RRPg (bzw. EB-NF). Man kann Übersetzungsschemata auch direkt mit RRPgs formulieren. Dann treten allerdings wegen der Schleifen Attributvariablen mit Mehrfach-Zuweisungen auf (entsprechend der Variablen r im obigen Programm).

6.8 Bottom-Up-Compiler-Generator

Aus einem Übersetzungsschema auf der Basis einer LR-Grammatik kann problemlos ein Bottom-Up-Compiler konstruiert werden, wenn die Grammatik **S-attribuiert** ist. In diesem Fall stehen Aktionen nur am Ende der Produktionen und können daher mit den syntaktischen Reduktionen gekoppelt werden.

```

class TopDown {
  String E() {
    String ri, ts, rs;
    ts = T(); ri = ts; rs = R(ri); return rs;
  }

  String R( String ri ) {
    String opx, ts, rli, rls;
    if( sym.equals("op") {
      opx = att; nextsym();
      ts = T(); rli = ri+ts+opx; rls = R(rli); return rls;
    } else return ri;
  }

  String T() {
    String es, idx; if( sym.equals("(") {
      nextsym(); es = E(); nextsym(); return es;
    } else {
      idx = att; nextsym(); return idx;
    }
  }
}

```

Der Aufruf `nextsym()` liefert das nächste Token und dessen Attribut in `sym` und `att`. Man startet den generierten Compiler durch den Aufruf `nextsym()`; `System.out.println(TopDown.E());`;

Beispiel 6.10: Aus dem Übersetzungsschema von Bsp. 6.9 erzeugt der beschriebene Generator folgenden Top-Down-Compiler (als Menge von Prozeduren)

Für die synthetisierten Attribute wird ein **Attributkeller** verwendet. Er nimmt das zu einem gekellerten Grammatiksymbol gehörige Attribut auf. Zustandskeller und Attributkeller können „parallel“ mit den gleichen Kellerzeigern verwaltet werden. Wir zeigen das Prinzip anhand des folgenden Übersetzungsschemas.

$$A \rightarrow X Y Z \{A.s = f(X.s, Y.s, Z.s)\}$$

Die Reduktion mit gekoppelter Aktion läuft in folgenden Schritten ab:

- Vor der Reduktion:

| | | | | |
|-----------------|-----|-----|-----|-----|
| | | | | top |
| | | | | ↓ |
| Zustandskeller: | ... | X | Y | Z |
| Attributkeller: | ... | X.s | Y.s | Z.s |

```
String E() {
    return R(T());
}

String R( String ri ) {
    String opx;
    if( sym.equals("op") ) {
        opx = att; nextsym();
        return( ri+T()+opx );
    } else return ri;
}
```

Beispiel 6.11: Vereinfachung von E und R aus Bsp. 6.10.

```
String E() {
    String opx, r;
    r = T();
    while( sym.equals("op") ) {
        opx = att; nextsym();
        r = r+T()+opx;
    }
    return r;
}
```

Beispiel 6.12: Optimierte Version zu Bsp. 6.11.

- Ermitteln von $ntop = top - r + 1$ (r = Länge der rechten Seite):

| | | | | |
|-----------------|-----|--------|-----|-------|
| | | $ntop$ | | top |
| | | ↓ | | ↓ |
| Zustandskeller: | ... | X | Y | Z |
| Attributkeller: | ... | X.s | Y.s | Z.s |

- Berechnen des Attributs **A.s** auf der Position $ntop$ (und Reduktion im Zustandskeller)
- Setzen von $top = ntop$

| | | |
|-----------------|-----|-------|
| | | top |
| | | ↓ |
| Zustandskeller: | ... | A |
| Attributkeller: | ... | A.s |

Die Anwendung dieses Prinzips wird in Bsp. 6.13 gezeigt.

| | |
|-----------------------------------|---------------------------------------|
| $E \rightarrow E_1 \text{ op } T$ | $\{E.s = E_1.s + T.s + \text{op}.x\}$ |
| $E \rightarrow T$ | $\{E.s = T.s\}$ |
| $T \rightarrow (E)$ | $\{T.s = E.s\}$ |
| $T \rightarrow \text{id}$ | $\{T.s = \text{id}.x\}$ |

Daraus ergibt sich der folgende Code für die Aktionen:

| Produktion | Code |
|-----------------------------------|----------------------------------------------------------------------|
| $E \rightarrow E_1 \text{ op } T$ | $a[\text{ntop}] = a[\text{top}-2] + a[\text{top}] + a[\text{top}-1]$ |
| $E \rightarrow T$ | |
| $T \rightarrow (E)$ | $a[\text{ntop}] = a[\text{top}-1]$ |
| $T \rightarrow \text{id}$ | |

Beispiel 6.13: Gegeben sei ein S-attribuiertes Übersetzungsschema zur Übersetzung von Infix nach Postfix gemäß der AG von Bsp. 6.4.

Der Attributkeller ist mit a bezeichnet. top zeigt auf den letzten Kellereintrag und wird bei shift-Operationen entsprechend erhöht. Immer wenn ein Symbol von der Eingabe in den Keller gelangt, wird der Wert des Attributs – den die Lexikalische Analyse liefert – in den Attributkeller gespeichert. Dies gilt z.B. für die Symbole aop und id . Falls ein Symbol kein Attribut besitzt, wie z.B. die Symbole $($ und $)$, bleibt die zugehörige Position des Attributkellers leer.

Die Zeiger top und ntop beziehen sich auf die Situation beim Schritt (2) im vorher skizzierten Ablauf. Bei der Produktion $E \rightarrow T$ ist kein Code für die Zuweisung notwendig, da die Attribute $T.s$ und $E.s$ die gleiche Position im Attributkeller einnehmen ($\text{top} = \text{ntop}$).

Man sieht, daß ein Compilergenerator den Code für die Aktionen generieren kann. Allerdings haben wir bisher nur synthetisierte Attribute berücksichtigt.

Erebt Attribute

Da die Aktionen zur Berechnung ererbter Attribute im Übersetzungsschema nicht am Ende der Produktion stehen, können sie nicht mit der Reduktion gekoppelt werden. Man wendet daher eine Transformation der Grammatik an (siehe Bsp. 6.14).

Bottom-Up-Compiler-Generator. Der Generator zur Erzeugung eines Ein-Pass-Bottom-Up-Compilers aus einem Übersetzungsschema arbeitet insgesamt wie folgt:

Eingabe: Übersetzungsschema auf der Basis einer LR(1)-Grammatik. Jede Aktion zur Berechnung eines ererbten Attributs $X.i$ muß unmittelbar links von X stehen, die Aktion zur Berechnung des synthetisierten Attributs $A.s$ muß am

$$\begin{aligned} S &\rightarrow A B \{C.i = f(A.s, B.s)\} C \{S.s = g(A.s, C.s)\} \\ C &\rightarrow D E \{C.s = h(C.i, E.s)\} \end{aligned}$$

Die Aktion, die das ererbte Attribut $C.i$ berechnet, wird durch eine eindeutige **Markierung**, d.h. ein neues Grammatiksymbol M mit der Produktion $M \rightarrow \varepsilon$ ersetzt. Das Symbol M belegt im Attributkeller vor C einen Platz für den Wert von $C.i$. Der Wert von $C.i$ kann somit als Wert eines synthetisierten Attributs von M berechnet werden, d.h. die Aktion zur Berechnung von $C.i$ wird der Produktion $M \rightarrow \varepsilon$ als Code zur Berechnung des synthetisierten Attributs $M.s$ zugeordnet:

| Produktion | Code |
|-----------------------------|---------------------------------|
| $S \rightarrow A B M C$ | $a[ntop] = g(a[top-3], a[top])$ |
| $M \rightarrow \varepsilon$ | $a[ntop] = f(a[top-1], a[top])$ |

Vor der Reduktion von M zeigt top auf B und es gilt $ntop = top+1$, weil $r = 0$ ist. Nach der Reduktion ist der Wert von $C.i$ auf der Position M (direkt vor C) gespeichert. Obwohl C selbst noch nicht im Keller steht, kann jede Produktion von C den Wert $C.i$ bereits im Keller ansprechen:

| Produktion | Code |
|---------------------|---------------------------------|
| $C \rightarrow D E$ | $a[ntop] = h(a[top-2], a[top])$ |

Beispiel 6.14: Ein Übersetzungsschema mit einem ererbten Attribut.

Ende stehen. Beide Forderungen stellen keine prinzipiellen Einschränkungen dar.

Ausgabe: Compiler, der nach dem Verfahren der LR-Analyse arbeitet.

Algorithmus: Der Generator erzeugt einen LR-Parser (siehe Kap. 5.3) mit Erweiterung des Zustandskellers um einen **Attributkeller**.

Für jede Produktion $A \rightarrow \dots X \dots$ des Übersetzungsschemas sind die auftretenden Aktionen folgendermaßen zu behandeln:

- Aus der Aktion zur Berechnung des **synthetisierten** Attributs $A.s$ wird der entsprechende Code generiert und dem Reduktionsschritt der Produktion zugeordnet.
- Jede Aktion zur Berechnung eines **ererbten** Attributs $X.i$ wird durch eine eindeutige Markierung M ersetzt und die Grammatik um die Produktion $M \rightarrow \varepsilon$ ergänzt. Aus der Aktion selbst wird der entsprechende Code generiert und dem Reduktionsschritt der Produktion $M \rightarrow \varepsilon$ zugeordnet.

Die folgenden zwei Arten von Kopieraktionen können ignoriert werden, da der Wert, der an $X.i$ zugewiesen werden soll, bereits auf der richtigen Position für $X.i$ (d.h. direkt vor X) im Attributkeller steht:

$$\begin{aligned} A &\rightarrow \{X.i = A.i\}X \dots \\ A &\rightarrow \dots Y\{X.i = Y.s\}X \dots \end{aligned}$$

- Für die geänderte Grammatik wird ein LR-Parser erzeugt.

Das Problem der Transformation durch Markierungen ist aber, daß die Grammatik unter Umständen die LR-Eigenschaft verliert, sodaß Konflikte in den Parser-Tabellen auftreten. Bsp. 6.15 zeigt einen solchen Fall.

6.9 Ein-Pass-Übersetzung

Notwendige Voraussetzung zur vollständigen Attributberechnung während der Syntax-Analyse ist eine L-attributierte Grammatik (LAG). Die LAG muß entweder auf einer LL-Grammatik basieren (LL-attributierte Grammatik, LLAG) oder auf einer LR-Grammatik, die durch die oben angegebene Transformation ihre LR-Eigenschaft nicht verliert (LR-attributierte Grammatik, LRAG).

Jede LLAG ist auch LRAG, da durch eingefügte ε -Produktionen ohne Alternativen die LL-Eigenschaft nicht zerstört wird.

Bsp. 6.15 zeigt drei mögliche Lösungen eines Problems entsprechend den AG-Klassen von Abb. 6.16.

LLAG:

$$\begin{aligned} S &\rightarrow id \{L.i = 1\} L \{S.s = L.s\} \\ L &\rightarrow id \{L_1.i = L.i+1\} L_1 \{L.s = L_1.s\} \\ L &\rightarrow \varepsilon \{L.s = L.i\} \end{aligned}$$

LRAG:

$$\begin{aligned} S &\rightarrow L \{S.s = L.s\} \\ L &\rightarrow L_1 id \{L.s = L_1.s+1\} \\ L &\rightarrow id \{L.s = 1\} \end{aligned}$$

LAG, aber nicht LRAG:

$$\begin{aligned} S &\rightarrow \{L.i = 0\} L \{S.s = L.s\} \\ L &\rightarrow \{L_1.i = L.i+1\} L_1 id \{L.s = L_1.s\} \\ L &\rightarrow id \{L.s = L.i+1\} \end{aligned}$$

Beispiel 6.15: Die Anzahl der Bezeichner in einer Liste soll berechnet werden.

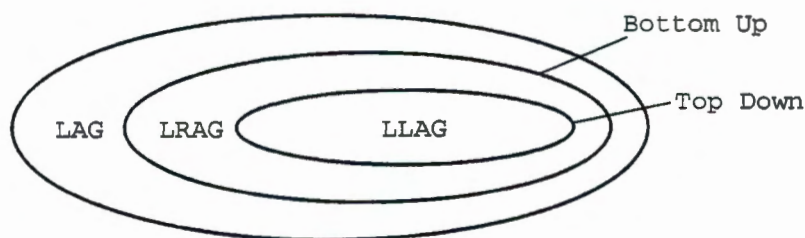


Abbildung 6.16: AG-Klassen für Ein-Pass-Übersetzung

6.10 Mehrpass-Übersetzung

Falls die Einpass-Übersetzung nicht möglich ist, erzeugt die Syntax-Analyse einen expliziten Ableitungsbaum (Syntaxbaum), der von einem **Attribut-Auswerter** beliebig durchlaufen werden kann. Im allgemeinsten Fall wird für jeden Syntaxbaum der Abhängigkeitsgraph konstruiert und die Attributauswertung aufgrund der Abhängigkeiten vorgenommen. Diese Methode ist sehr aufwendig, da sie die statisch bekannten Abhängigkeiten der Grammatik nicht ausnützt. Wir behandeln im folgenden zwei effizientere Methoden.

Pass-orientierte Auswerter durchlaufen den Baum ein oder mehrmals nach einer *globalen Strategie*, z.B. mit links-rechts-Tiefensuche (LAG-Auswerter). Bei jedem Durchgang werden Attribute ausgewertet, die aufgrund der bis dahin verfügbaren Attributwerte berechnet werden können. Durch Analyse der Attributabhängigkeiten wird vor der Auswertung bestimmt, welche Attribute in den einzelnen Durchgängen zu berechnen sind.

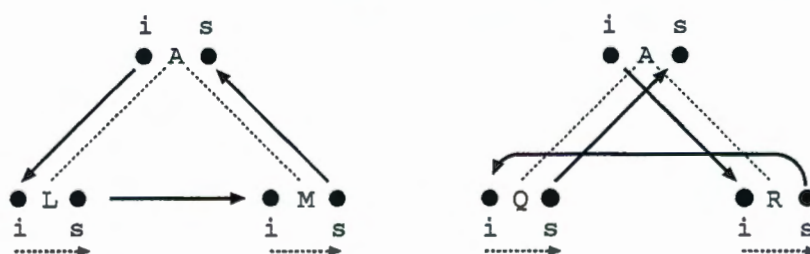


Abbildung 6.17: Eine L-attributierte und eine nicht-L-attributierte Produktion für A

Abb. 6.17 zeigt zwei alternative Produktionen für ein Nonterminal A und die Attributabhängigkeiten. Die zweite Produktion erfordert zwei Durchläufe eines LAG-Auswerter: Im ersten Durchgang wird R.s und im zweiten Q.s berechnet.

In Abb. 6.18 ist eine Situation dargestellt, die bei der Semantischen Analyse des Overloading in Ada auftritt: Das ererbte Attribut X.i ist vom synthetisierten Attribut X.s abhängig. Der LAG-Auswerter berechnet im ersten Durchgang X.s

und im zweiten Durchgang X.c.

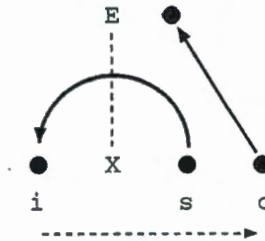


Abbildung 6.18: Nicht-L-attribuierte Produktion $E \rightarrow X$

Visit-orientierte Auswerter durchlaufen den Baum nicht nach einer globalen Strategie, sondern entsprechend den Attributabhängigkeiten der Teilbäume.

Ein Visit-orientierter Auswerter kann z.B. konstruiert werden, indem man jedem Nonterminal eine eigene Auswertungsfunktion zuordnet.

```
public SType A( Node n, IType ai ) {
  switch( Produktion am Knoten n ) {
    case 'A → L M':
      // L vor M aufrufen
      li = ...ai...
      ls = L(child(n,1),li);
      mi = ...ls...
      ms = M(child(n,2),mi)
      return ...ms...
    case 'A → Q R':
      // R vor q aufrufen
      ri = ...ai...
      rs = R(child(n,2),ri)
      qi = ...rs...
      qs = Q(child(n,1),qi)
      return ...qs...
  }
}
```

Algorithmus 6.19: Auswertungsfunktion für das Nonterminal A.

Die Abhängigkeiten in Abb. 6.17 zeigen, daß bei der Produktion $A \rightarrow LM$ die Auswertungsfunktion L vor M und bei der Produktion $A \rightarrow QR$ die Funktion R vor Q aufgerufen werden muß. Die unterschiedlichen Aufruffolgen können in der Funktion A berücksichtigt werden. Die Funktion A ist in Alg. 6.19 dargestellt. Dabei ist angenommen, daß die Knoten des Baums durch eine eindeutige Nummer n identifiziert sind und die Funktion $child(n,k)$ die Nummer des k-ten direkten Nachfolgers von

n liefert. Außerdem sei zu jedem Knoten gespeichert, welche Produktion angewandt wurde.

Die Auswertungsfunktion A bringt gegenüber dem LAG-Auswerter den Vorteil, daß jeder Teilbaum A in *einem* Durchlauf ausgewertet werden kann.

Im Beispiel von Abb. 6.18 benötigt man für X zwei Funktionen X1 und X2, die in E nacheinander aufzurufen sind: X1 berechnet X.s und X2 berechnet X.c. Der Teilbaum X wird also zweimal besucht, d.h. der Visit-orientierte Auswerter verhält sich in diesem Beispiel wie der Pass-orientierte LAG-Auswerter.

Kapitel 7

Semantische Analyse

Dieses Kapitel behandelt die Semantische Analyse für Programmiersprachen mit einem strengen statischen Typsystem wie z.B. Pascal, Modula und Ada. Die wesentliche Aufgabe der Semantischen Analyse ist die **Typ-Überprüfung** (*type checking*). Dabei muß überprüft werden, ob die Verwendung von Objekten mit ihrer Deklaration verträglich ist.

Grundlage für die Typ-Überprüfung ist das **Typsystem** der jeweiligen Programmiersprache. Als Typsystem bezeichnet man die Menge der Regeln, die den Elementen einer Programmiersprache Typen zuordnen. Typen lassen sich formal durch Typ-Ausdrücke beschreiben und durch Graphen darstellen. Die zweckmäßige interne Darstellung von Typen ist für eine effiziente Typ-Überprüfung von großer Bedeutung. In Pascal und Modula beeinflusst die interne Darstellung sogar das Sprachkonzept bezüglich der Typ-Äquivalenz.

Der Zusammenhang zwischen Deklaration und Verwendung läßt sich nicht durch eine kontextfreie Grammatik beschreiben. Die Semantische Analyse ist daher ein klassisches Anwendungsfeld der Attributierten Grammatik. Zentrale Datenstruktur der Semantischen Analyse ist die **Symboltabelle** als 'Datenbasis' für die Typ-Überprüfung.

In den meisten Programmiersprachen treten Operatoren auf, deren Typ nicht eindeutig deklariert ist. Die Semantische Analyse muß den Typ solcher Operatoren bzw. den Typ entsprechender Teilausdrücke aus dem jeweiligen Kontext ableiten. Dabei sind auch die Regeln bezüglich impliziter Typ-Konversionen zu beachten.

Die Semantische Analyse soll die **statische Korrektheit** des Quellprogramms sicherstellen. Neben der Typanalyse fallen daher unter Umständen noch einige Nebenaufgaben an, z.B. in Modula die Überprüfung, ob der Bezeichner von Modulen und Prozeduren an ihrem Ende korrekt wiederholt wird. Solche Prüfungen benötigen aber keine speziellen Verfahren und werden daher in diesem Kapitel nicht explizit behandelt.

7.1 Typen

Typsysteme

Ein Typsystem ist die Menge der Regeln einer Programmiersprache, die den Programmelementen Typen zuordnen und die Äquivalenz und Kompatibilität von Typen definieren. Man kann **statische** und **dynamische Typsysteme** unterscheiden, je nachdem, ob die Typregeln zur Übersetzungszeit (=statisch) oder zur Laufzeit (=dynamisch) zur Anwendung kommen.

Statische Typsysteme haben den Vorteil, daß sie die Korrektheit (Sicherheit) und die Effizienz des Programms erhöhen. Der Compiler kann die Typen überprüfen und im erzeugten Programm die genau passenden Operatoren und Speicherbereiche festlegen. Bei dynamischen Typsystemen muß die Typ-Überprüfung generiert bzw. in das Laufzeitsystem eingebaut werden. Das Laufzeitverhalten ist daher entsprechend weniger effizient. Beide Systeme können auch kombiniert werden.

Ein Typsystem kann man als 'Filter' ansehen, das die Menge von Programmen bzw. Prozessen auf typkorrekte Teilmengen einschränkt. Dies wird in Abb. 7.1 anschaulich dargestellt.

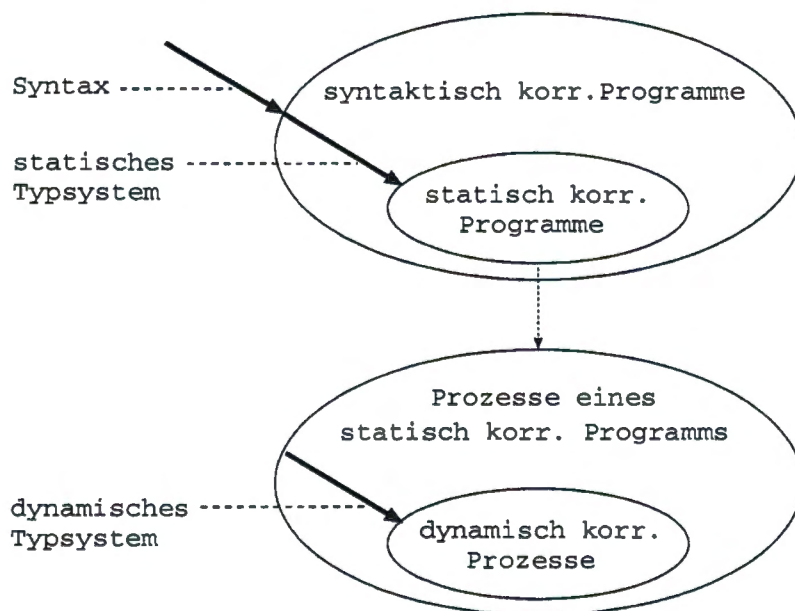


Abbildung 7.1: Statisches und dynamisches Typsystem

Typ-Ausdrücke

Ein Typ kann formal durch einen **Typ-Ausdruck** beschrieben werden. Typ-Ausdrücke bilden die Grundlage zur internen Darstellung von Typen. Daher werden wir Typ-Ausdrücke für eine 'typische' Programmiersprache definieren:

Ein **Basistyp** ist ein Typ-Ausdruck. Basistypen sind z.B. `bool`, `char`, `int` und `real`.

Ein **Typ-Konstruktor**, angewendet auf einen oder mehrere Typ-Ausdrücke, ist ein Typ-Ausdruck. Konstruktoren sind z.B.:

Array: `array(n,T)` mit `T` als Elementtyp und `n` = Anzahl der Elemente.

Produkt: $T_1 \times T_2$ (z.B. für Parameterlisten)

Record: `record((n1 : T1) × ... × (nk : Tk))` (n_i: Feldnamen)

Zeiger: `pointer(T)`

Funktion: $T_1 \rightarrow T_2$

Ein Typ-Ausdruck kann als Baum dargestellt und implementiert werden. Endknoten sind mit Basistypen, alle anderen Knoten mit Typkonstruktoren markiert. Abb. 7.2 zeigt als Beispiel den Baum für einen Funktionstyp.

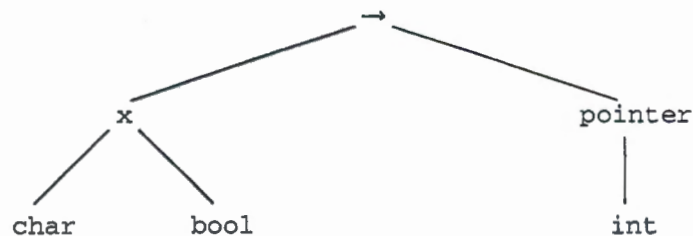


Abbildung 7.2: Typ-Ausdruck `char × bool → pointer(int)`

Typnamen und Typ-Äquivalenz

In vielen Programmiersprachen können Typnamen definiert werden. Dabei ist die Frage der Typ-Äquivalenz zu klären (siehe Bsp. 7.3).

```

TYPE link = POINTER TO cell;
VAR next : link;
    last : link;
    p :    POINTER TO cell;
    q,r :  POINTER TO cell;
  
```

Beispiel 7.3: Eine Typ- und Variablendeklaration (Modula-2).

Man unterscheidet im wesentlichen zwei Möglichkeiten, um die Äquivalenz von Typen zu definieren:

Struktur-Äquivalenz bedeutet, daß Typen mit gleicher Struktur gleich sind. Im Bsp. 7.3 haben dann alle Variablen den gleichen Typ.

Namen-Äquivalenz bedeutet, daß Typen mit gleichem Namen gleich sind. Jeder anonymen Typ-Konstruktion wird ein eindeutiger impliziter Name zugeordnet. Im Beispiel haben dann `next` und `last`, sowie `q` und `r` jeweils den gleichen Typ. Diese Art der Typ-Äquivalenz findet man in Pascal und Modula.

Die Namen-Äquivalenz kann anhand des Typbaums dargestellt und implementiert werden, indem man den Knoten (explizite oder implizite) Typnamen zuordnet. Zwei Typen sind dann gleich, wenn sie durch denselben Knoten repräsentiert werden.

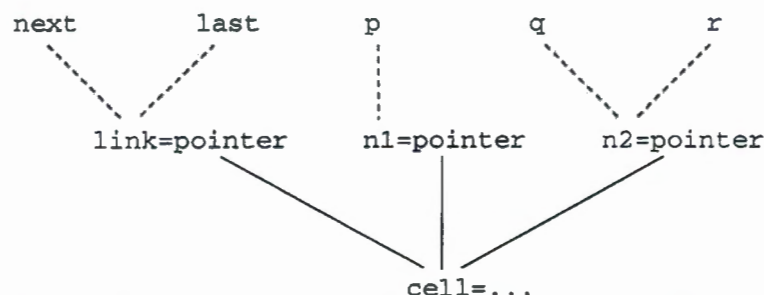


Abbildung 7.4: Namen-Äquivalenz

7.2 Symboltabelle

Die Namen aller im Quellprogramm deklarierten Objekte und deren Typen werden in einer **Symboltabelle** gesammelt. Jeder Eintrag enthält den Namen und den Typ des deklarierten Objekts. Namen werden durch das Attribut `id.x` repräsentiert, Typen werden als Zeiger auf den Typbaum dargestellt.

Die Symboltabelle wird bei der Verarbeitung der Deklarationen aufgebaut und bei der Verarbeitung der Anweisungen verwendet. In einer Attributierten Grammatik (ohne Seiteneffekte) tritt die Symboltabelle daher als Attribut auf, dessen Wert im Deklarationsteil berechnet und an den Anweisungsteil übergeben wird. Die AG in Bsp. 7.5 zeigt, wie die Symboltabelle aus den Deklarationen aufgebaut werden kann. Das Nonterminal `D` besitzt ein ererbtes und ein synthetisiertes Attribut. Der Wert des ererbten Attributs ist jeweils die Symboltabelle *vor* der Verarbeitung von `D`, der Wert des synthetisierten Attributs ist die Symboltabelle *nach* der Verarbeitung von `D`. Auf der obersten Ebene (erste Produktion) wird mit der leeren Tabelle begonnen und die vollständige Tabelle zurückgeliefert. Jede Deklaration einer Variablen (dritte Produktion) erzeugt einen neuen Eintrag in der Tabelle. Im Beispiel fassen wir die Symboltabelle als Liste auf, wobei die Operation `||` ein Element am Ende der Liste anhängt.

| Produktion | Regeln |
|----------------------------------------------------------|---------------------------------------------------|
| $P \rightarrow D \text{ begin } S \text{ end}$ | $D.i = ()$ $S.i = D.s$ |
| $D \rightarrow D_1 ; D_2$ | $D_1.i = D.i$ $D_2.i = D_1.s$ $D.s = D_2.s$ |
| $D \rightarrow \text{id} : T$ | $D.s = D.i \parallel (\text{id}.x, T.type)$ |
| $T \rightarrow \text{int}$ | $T.type = \text{int}$ |
| $T \rightarrow \text{real}$ | $T.type = \text{real}$ |
| $T \rightarrow \text{arr} [\text{inum}] \text{ of } T_1$ | $T.type = \text{array}(\text{inum}.x, T_1.type)$ |

Beispiel 7.5: AG zur Erzeugung einer Symboltabelle aus Variablendeklarationen (Typen: integer, real und array mit Untergrenze 0).

Basistypen sind vorgegebene Endknoten von Typbäumen. Für strukturierte Typen muß aus der Deklaration der entsprechende Typbaum konstruiert werden. Die letzte Zeile der AG in Bsp. 7.5 definiert die Konstruktion des Typbaums für mehrdimensionale Arrays. Der Typbaum wird im synthetisierten Attribut *T.type* aufgebaut. Abb. 7.6 veranschaulicht die Konstruktion am attributierten Ableitungsbaum für ein zweidimensionales Array.

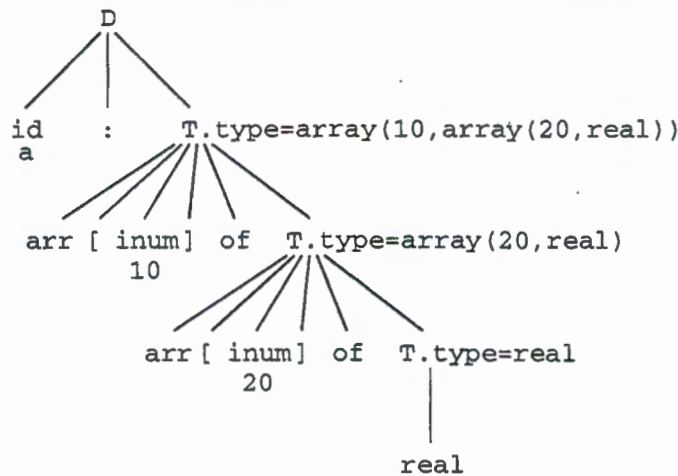


Abbildung 7.6: Attributierter Ableitungsbaum für die Deklaration *a:ARRAY[10] OF ARRAY[20] OF REAL*

Beim Aufbau der Symboltabelle müssen die Gültigkeitsbereiche von geschachtelten Prozedurdeklarationen berücksichtigt werden, d.h. innerhalb jeder Prozedurde-

klaration ist die Symboltabelle um die lokalen Deklarationen zu erweitern. Dies läßt sich auf einfache Weise durch eine AG beschreiben (siehe Bsp. 7.7). Die Prozedurdeklaration D erweitert die Tabelle zunächst um den Eintrag für den Prozedurnamen. Die so erweiterte Tabelle wird einerseits nach außen zurückgegeben (da der Name außen sichtbar ist) und andererseits an den inneren Prozedurblock P weitergegeben. Dort wird die Tabelle durch die lokalen Deklarationen D ergänzt und an den Anweisungsteil S der Prozedur übergeben.

| Produktion | Regeln |
|------------------------------------------------|-------------------------------------------------------------------|
| $D \rightarrow \text{proc id } P$ | $h = D.i \parallel (\text{id.x, proc})$ $D.s = h$ $P.i = h$ |
| $P \rightarrow D \text{ begin } S \text{ end}$ | $D.i = P.i$ $S.i = D.s$ |

Beispiel 7.7: AG zur Erzeugung der Symboltabelle bei geschachtelten Prozedurdeklarationen

7.3 Typ-Überprüfung

Bei imperativen Sprachen ist zwischen Anweisungen und Ausdrücken zu unterscheiden. Die Typ-Überprüfung bezieht sich nur auf Ausdrücke. Da Ausdrücke aber innerhalb von Anweisungen auftreten, muß die Symboltabelle als Attribut über die Anweisungen zu den Ausdrücken bis zu den elementaren Operanden weitergegeben werden (siehe Bsp. 7.8). Wir nehmen an, daß die Zugriffsfunktion $\text{looktype}(n, t)$ zu einem gegebenen Namen n den in der Tabelle t eingetragenen Typ liefert. Wenn die Symboltabelle in der vorher beschriebenen Form als Liste aufgebaut ist, kann man die Funktion looktype durch sequentielle Rückwärtssuche implementieren, sodaß die lokalen Deklarationen zuerst gefunden werden.

Die eigentliche Typ-Überprüfung erfolgt in Bsp. 7.8 mithilfe des synthetisierten Attributs `type` aufsteigend im Ableitungsbaum. Dabei wird der Typ aller Teilausdrücke bestimmt, und mehrdeutige Operatoren werden eindeutig identifiziert. In Modula hat "+" zum Beispiel die zwei „überladenen“ Bedeutungen $\text{int} \times \text{int} \rightarrow \text{int}$ und $\text{real} \times \text{real} \rightarrow \text{real}$.

Die Typ-Überprüfung verläuft ähnlich wie eine Berechnung des Ausdrucks. Man kann daher die Semantische Analyse als Interpreter ansehen, der Ausdrücke auf der 'Typebene' interpretiert. Anstelle der konkreten Werte, die zur Übersetzungszeit noch nicht vorliegen, werden nur die Typen als abstrakte Werte 'berechnet'. Diese Berechnung liefert als Ergebnis den Typ des Ausdrucks und in den Zwischenschritten die Typen der Teilausdrücke und Operatoren – oder eine Fehlermeldung.

Zur besseren Übersicht wird in den nachfolgenden Beispielen die Weitergabe der Symboltabelle als Attribut nicht mehr explizit angegeben.

| Produktion | Regeln |
|------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| $P \rightarrow D \text{ begin } S \text{ end}$ | $D.i = ()$ $S.i = D.s$ |
| $S \rightarrow S_1 ; S_2$ | $S_1.i = S.i$ $S_2.i = S.i$ |
| $S \rightarrow id = E$ | $E.i = S.i$ $\text{if(looktype(id.x,S.i) \neq E.type) error}$ |
| $E \rightarrow E_1 + E_2$ | $E_1.i = E.i$ $E_2.i = E.i$ $E.type =$ $\text{if}(E_1.type == \text{int} \ \&\& \ E_2.type == \text{int}) \ \text{int}$ $\text{else if}(E_1.type == \text{real} \ \&\& \ E_2.type == \text{real}) \ \text{real}$ else error |
| $E \rightarrow id$ | $E.type = \text{looktype}(id.x,E.i)$ |
| $E \rightarrow \text{inum}$ | $E.type = \text{int}$ |
| $E \rightarrow \text{rnum}$ | $E.type = \text{real}$ |

Beispiel 7.8: AG zur Typ-Überprüfung (Modula).

7.4 Typ-Konversion

Man unterscheidet explizite und implizite Typ-Konversion. Bei der **expliziten** Typ-Konversion muß eine Konversionsfunktion im Programm geschrieben werden (z.B. `FLOAT(i)` in Modula). Solche Funktionen können von der Semantischen Analyse wie jede andere Funktion behandelt werden.

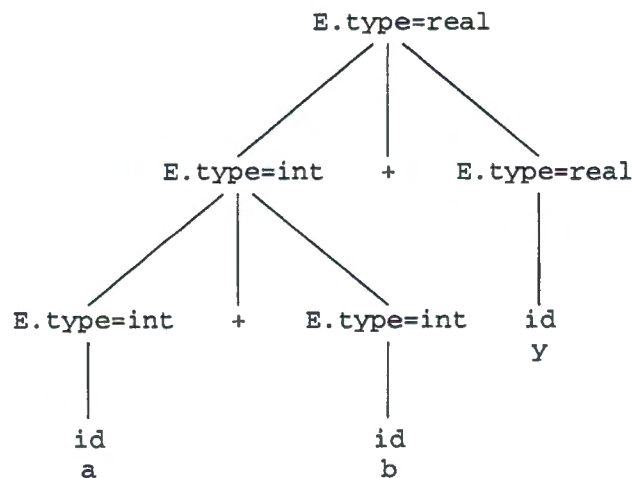
Implizite Konversionen dagegen müssen von der Semantischen Analyse besonders berücksichtigt werden, z.B. in Pascal bei der Verknüpfung von integer- und real-Operanden, wobei integer nach real implizit konvertiert wird.

Die AG in Bsp. 7.9 überprüft die Typen unter Berücksichtigung der Regeln für implizite Konversionen in Pascal. Die Konversions-Operationen müssen bei der Code-Erzeugung an den entsprechenden Stellen generiert werden. Der Ablauf einer Typ-Überprüfung kann am attributierten Ableitungsbaum dargestellt werden. Die Abb. 7.10 basiert auf den Regeln von Bsp. 7.9. Die verwendeten Variablen seien wie folgt deklariert: $a, b: \text{INT}; y: \text{REAL}$.

Typ-Überprüfung und Operator-Identifizierung können bei den meisten Programmiersprachen in einem Durchlauf des Ausdrucksbaums mit einem synthetisierten Attribut erfolgen. Voraussetzung dafür ist, daß der Ergebnistyp jeder Operation eindeutig durch die Operandentypen bestimmt ist. Das gilt z.B. auch für C++, da die Operandentypen überladener Operatoren eindeutig sein müssen. Im Gegensatz dazu werden in Ada überladene Operatoren durch die Operandentypen *und den*

| Produktion | Regeln |
|---------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| $E \rightarrow E_1 + E_2$ | $E.type =$ if($E_1.type == int \ \&\& \ E_2.type == int$) int else if($E_1.type == real \ \&\& \ E_2.type == real$) real else if($E_1.type == int \ \&\& \ E_2.type == real$) real else if($E_1.type == real \ \&\& \ E_2.type == int$) real else error |

Beispiel 7.9: AG zur Typ-Überprüfung (Pascal).

Abbildung 7.10: Typ-Überprüfung für $a+b+y$

Ergebnistyp identifiziert. Dadurch wird die Analyse wesentlich komplizierter.

7.5 Overloading

Overloading bedeutet, daß ein Symbol (Bezeichner, Operator) mit mehreren Bedeutungen 'überladen' wird. Die aktuelle Bedeutung muß jeweils aus dem Kontext ermittelt werden. In den meisten Programmiersprachen ist z.B. der Operator "+" überladen:

1. $int \times int \rightarrow int$
2. $real \times real \rightarrow real$

Der Typ von "+" kann hier eindeutig aus dem Typ der Operanden bestimmt werden. In Ada können jedoch zusätzlich im Programm weitere Typen für "+" definiert werden, und zwar auch solche mit den gleichen Operandentypen, z.B.


```
function "+"(i,j:integer) return real;
```

Damit ist der Operator "+" dreifach überladen:

1. $\text{int} \times \text{int} \rightarrow \text{int}$
2. $\text{real} \times \text{real} \rightarrow \text{real}$
3. $\text{int} \times \text{int} \rightarrow \text{real}$

Die beiden Operanden allein identifizieren den Operator nicht mehr eindeutig. Der Typ eines Ausdrucks $a+b$ mit zwei int -Operanden ist z.B. int oder real , d.h. ein Element aus der Menge $\{\text{int}, \text{real}\}$. Der Typ des Operators "+" ist ein Element aus der Menge $\{\text{int} \times \text{int} \rightarrow \text{int}, \text{int} \times \text{int} \rightarrow \text{real}\}$.

Welcher Operator anzuwenden ist, muß durch den jeweiligen Kontext, d.h. durch den in den Kontext passenden Ergebnistyp des Ausdrucks, entschieden werden. Falls auch dadurch keine eindeutige Entscheidung möglich ist, liegt gemäß Sprachdefinition von Ada ein Typfehler vor.

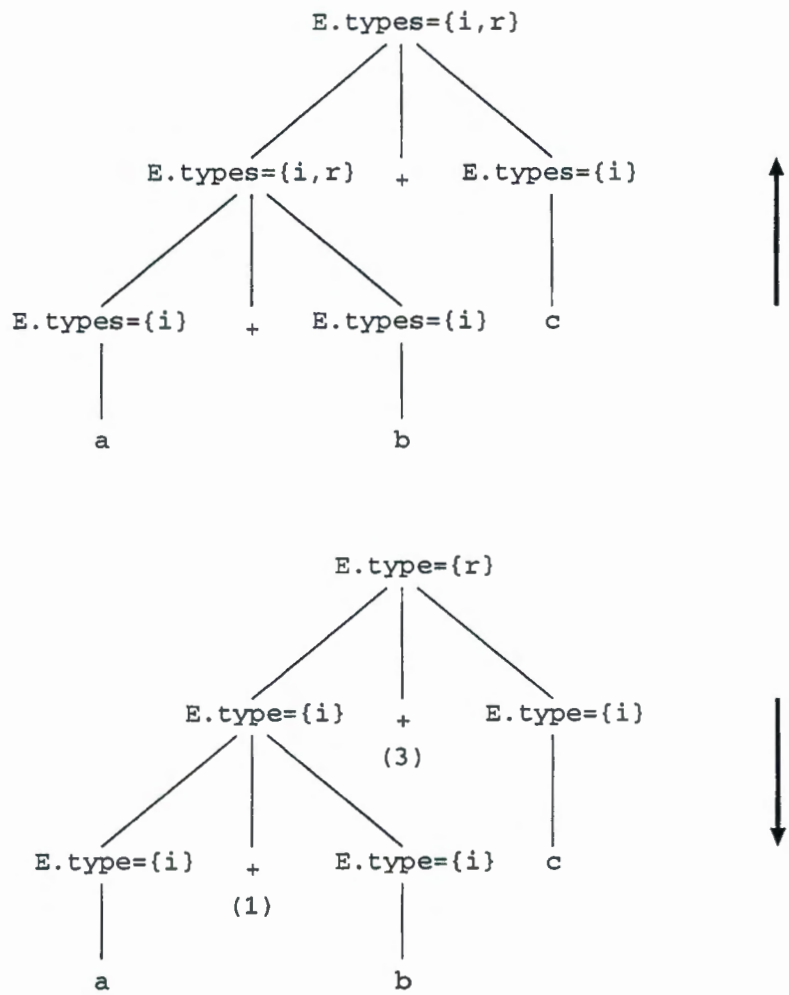
Zur Analyse eines Ausdrucks muß der zugehörige Baum in zwei Richtungen durchlaufen werden:

- Aufsteigende Berechnung der Typmengen an den Knoten.
- Absteigende Berechnung des eindeutigen Typs für jeden Knoten.

Bei der absteigenden Berechnung muß sich die Typmenge jedes Knotens auf genau ein Element reduzieren, damit die anzuwendenden Operatoren eindeutig bestimmt sind. Das Verfahren wird in Bsp. 7.11 veranschaulicht. Zunächst werden im Attribut `types` die Typmengen der Knoten *aufsteigend* anhand der definierten Typen der Bezeichner und des Operators "+" berechnet. Im Beispiel ergibt sich an der Wurzel ein eindeutiger Typ. Dieser wird in das Attribut `type` übernommen, und *absteigend* werden jene Typen gestrichen, für die es keine passende Operatordefinition gibt. Im Beispiel ist der Typ i aus der Menge $\{i, r\}$ zu streichen, da $i \times r \rightarrow r$ für den Operator "+" nicht definiert ist.

Die AG in Bsp. 7.12 beschreibt das Verfahren formal, wobei angenommen ist, daß die Funktion `looktypes` die Typmengen der deklarierten Bezeichner bzw. Operatoren aus der Symboltabelle liefert.

Aus der Sicht der Attributierten Grammatik erkennt man die Notwendigkeit für zwei Baumdurchläufe an der Produktion $E' \rightarrow E$ (Wurzel des Ableitungsbaumes). Das ererbte Attribut `type` des Knotens E hängt vom synthetisierten Attribut `types` desselben Knotens ab. Die AG von Bsp. 7.12 ist also nicht L-attributiert.



Beispiel 7.11: Gegeben seien die Variablen a , b und c vom Typ i (integer) und y vom Typ r (real). Der Operator "+" sei mit $\{(1) i \times i \rightarrow i, (2) r \times r \rightarrow r, (3) i \times i \rightarrow r\}$ überladen (Ada). Die Typ-Überprüfung des Ausdrucks $y = a+b+c$ wird durch die gezeigten attributierten Ableitungsbäume veranschaulicht

| Produktion | Regeln |
|-------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| $S \rightarrow id = E$ | $E.type = looktypes(id.x) \cap E.types$ |
| $E \rightarrow E_1 \text{ op } E_2$ | $E.types = \{t \mid (t_1 \times t_2 \rightarrow t) \in looktypes(op.x)$ $\wedge t_1 \in E_1.types \wedge t_2 \in E_2.types\}$ $E_1.type = \{t_1 \mid (t_1 \times t_2 \rightarrow t) \in looktypes(op.x)$ $\wedge t_1 \in E_1.types \wedge t_2 \in E_2.types$ $\wedge t = E.type\}$ $E_2.type = \{t_2 \mid (t_1 \times t_2 \rightarrow t) \in looktypes(op.x)$ $\wedge t_1 \in E_1.types \wedge t_2 \in E_2.types$ $\wedge t = E.type\}$ |
| $E \rightarrow id$ | $E.types = looktypes(id.x)$ |

Falls für das Attribut $E.type$ mehrelementige Mengen auftreten wird ein Fehler ausgelöst.

Beispiel 7.12: AG zur Typ-Überprüfung (Ada)

Kapitel 8

Zwischendarstellungen

Ein Quellprogramm, das Syntax-Analyse und Semantische Analyse durchlaufen hat, kann prinzipiell (durch einen Interpreter) ausgeführt werden. Die Darstellung eines Programms in diesem Stadium bezeichnen wir als **Zwischendarstellung**.

Die allgemeinste Form der Zwischendarstellung ist der (attributierte) **Syntaxbaum** oder **Operatorbaum**, der das Programm als hierarchische Struktur von Operatoren und Operanden darstellt.

Bei einem **linearen Zwischencode** wird das Programm als lineare Folge von Operatoren und Operanden dargestellt. Einfache Codes dieser Art sind **Postfix-** und **Prefix-Code**. Der Postfix-Code hat als Zwischencode für Interpreter besondere Bedeutung, da er nur einen Stack für die Operanden erfordert. Als Zwischencode im Compiler wird häufig der **Quadrupel-Code** verwendet, ein linearer Zwischencode, dessen Elemente Quadrupel (4-Tupel) sind.

Wir werden in diesem Kapitel zunächst die elementaren Zusammenhänge zwischen Syntaxbaum und linearen Codes zeigen und dann den Quadrupel-Code eingehend behandeln. Insbesondere wird unter Verwendung Attributierter Grammatiken gezeigt, wie aus typischen Konstrukten der Quellsprache entsprechende Folgen von Quadrupeln erzeugt werden können.

Der Quadrupelcode eignet sich als Grundlage für einfache Verfahren der Codeerzeugung. Er kann aber auch verwendet werden, um Kontroll- und Datenflußgraphen zu konstruieren, die man zur optimalen Codeerzeugung benötigt.

8.1 Syntaxbaum

Ein Syntaxbaum (Operatorbaum) stellt das Quellprogramm als Struktur von Operatoren und Operanden dar. Jeder Knoten des Baums bestimmt eine Operation, deren Operanden die anhängenden Teilbäume sind. Abb. 8.1 skizziert zwei Beispiele.

Bei LR-Grammatiken erhält man den Syntaxbaum auf einfache Weise durch Vergrößerung (Abstraktion) des Ableitungsbaums. Bei LL-Grammatiken wird wegen der fehlenden Linksrekursion die Programmstruktur weniger deutlich durch den Ableitungsbaum dargestellt, so daß die Erzeugung des Syntaxbaums aufwendiger ist. In jedem Fall aber ist es möglich, den Syntaxbaum während der Syntax-Analyse explizit zu erzeugen.

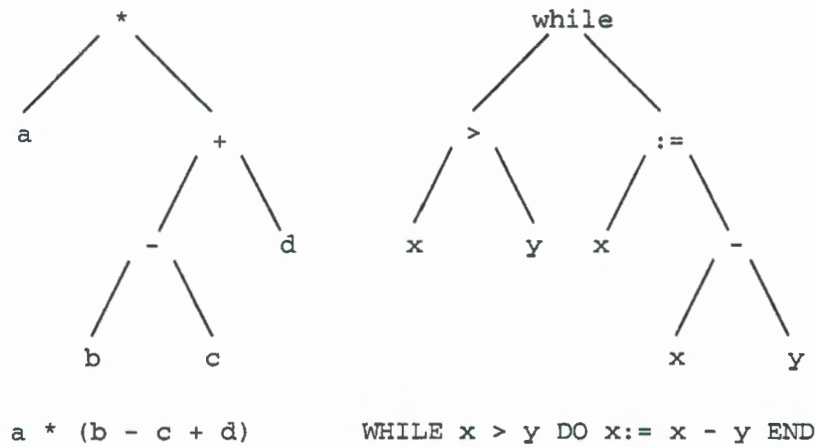


Abbildung 8.1: Syntaxbäume

Der Syntaxbaum ist als universelle Zwischendarstellung des Quellprogramms für weitere Transformationen in Richtung eines Zielprogramms gut geeignet. Er kann prinzipiell aber auch als Programm durch einen rekursiven Interpreter ausgeführt werden (mit entsprechenden Eingabedaten). Bei einem WHILE-Operator müßte der rekursive Interpreter z.B. den Schleifen-Unterbaum wiederholt aufrufen, bis die Wiederholungsbedingung nicht mehr erfüllt ist. Im Anhang B wird die Implementierung eines rekursiven Interpreters in Prolog gezeigt.

Zur Erzeugung von Maschinencode ist jedenfalls eine Linearisierung des Baums erforderlich. Im folgenden zeigen wir den Zusammenhang zwischen Syntaxbaum und den linearen Formen Postfix-, Prefix- und Quadrupel-Code.

8.2 Postfix und Prefix

Postfix- und Prefix-Code sind **klammerfreie Folgen** von Operatoren und Operanden. Dabei ist zu beachten, daß die Stelligkeit jedes Operators eindeutig ist (z.B. müssen unäres und binäres "-" unterschiedlich dargestellt werden).

Die **Postfix-Folge** (*postorder*) kann wie folgt aus dem Baum erzeugt werden: Bei einem Tiefendurchlauf wird jeder Knoten *nach* dem Besuch seiner Nachfolger ausgegeben, also jeder Operator nach seinen Operanden (Abb. 8.2).

Die **Prefix-Folge** (*preorder*) kann wie folgt aus dem Baum erzeugt werden: Bei einem Tiefendurchlauf wird jeder Knoten *vor* dem Besuch seiner Nachfolger ausgegeben, also jeder Operator vor seinen Operanden (Abb. 8.3).

Die Prefix-Folge ist eine klammerfreie Form der Funktions-Schreibweise

$$/(*(a, b), +(-(c, d), e))$$

Die geklammerte Version kann man durch folgende Ergänzung leicht erzeugen: Bei jedem Innenknoten wird nach Ausgabe des Operators vor dem Besuch des ersten Nachfolgers eine öffnende Klammer, vor jedem weiteren Nachfolgerbesuch ein

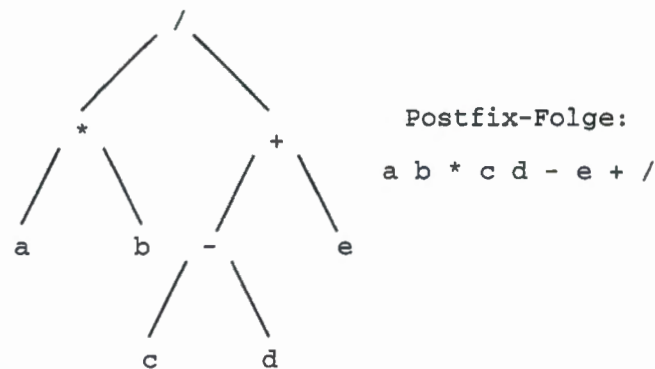


Abbildung 8.2: Postfix-Folge zum Syntaxbaum

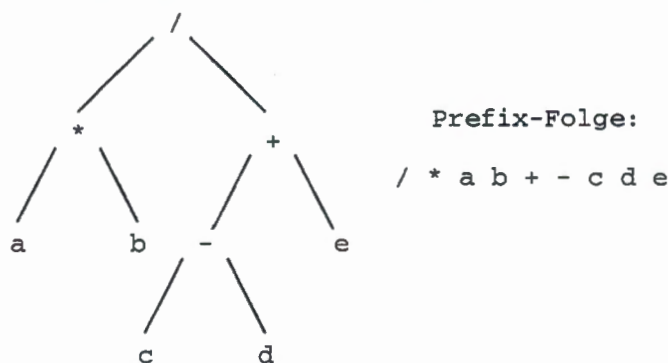


Abbildung 8.3: Prefix-Folge zum Syntaxbaum

Beistrich und nach dem letzten Nachfolgerbesuch eine schließende Klammer ausgegeben.

8.3 Quadrupel

Ein Quadrupel besteht aus folgenden vier Elementen:

(Operator, 1.Operand, 2.Operand, Ergebnis).

Da Quadrupel drei Adressen enthalten, nennt man sie auch **Drei-Adreß-Code**.

Die Erzeugung aus dem Syntaxbaum erfolgt ähnlich wie bei der Postfix-Folge mit dem Unterschied, daß eindeutige symbolische Namen für die Operatorknoten generiert und ausgegeben werden (Abb. 8.4).

Zur besseren Lesbarkeit werden Quadrupel oft im Stil einer Programmiersprache als elementare Anweisungen geschrieben, also z.B.

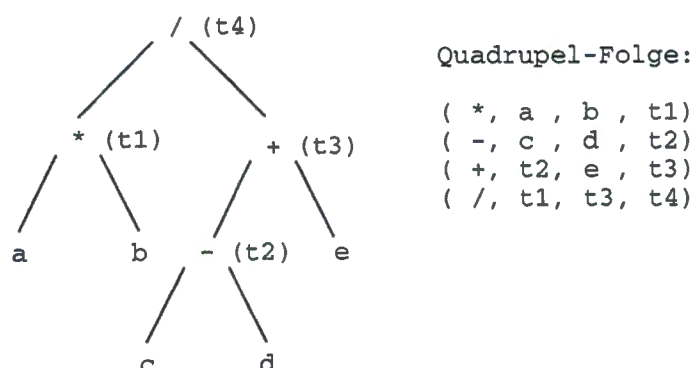


Abbildung 8.4: Quadrupel-Folge zum Syntaxbaum

$t1 = a + b$ statt $(+, a, b, t1)$

Man beachte, daß Quadrupel im Unterschied zum Postfix- oder Prefix-Code nur maximal zweistellige Operatoren enthalten können. Der Vorteil von Quadrupeln liegt vor allem darin, daß sie einen symbolischen Namen für das Ergebnis enthalten und daher innerhalb einer Folge verschiebbar sind. Quadrupel-Code eignet sich daher für Optimierungen.

Wir verwenden folgenden Code (für die Quadrupel (op, x, y, z)):

| | |
|-------------------------------|----------------------------------------------------------------------------------|
| $z = x$ | z erhält den Wert von x |
| $z = op\ x$ | op ist ein unärer Operator |
| $z = x\ op\ y$ | op ist ein binärer Operator |
| <code>goto m</code> | unbedingter Sprung zur Marke m |
| <code>if x op y goto m</code> | Sprung zur Marke m , wenn $x\ op\ y$ TRUE ist, op ist ein Vergleichsoperator |
| $z = @\ x$ | z erhält den Wert des Speicherplatzes, dessen Adresse in x steht |
| $@\ z = x$ | der Speicherplatz, dessen Adresse in z steht, erhält den Wert von x |

Im folgenden werden wir mit Attributierten Grammatiken die Erzeugung von Quadrupeln für typische Konstrukte einer Programmiersprache beschreiben.

8.4 Zuweisungen

Die AG in Abb. 8.5 definiert die Übersetzung von Zuweisungen in Quadrupel. Der Operator "+" steht stellvertretend für binäre arithmetische Operatoren.

Die Funktion `lookadr(id.x)` liefert die Speicheradresse von `id` (aus der Symboltabelle). Die Funktion `newtemp` erzeugt eine neue Hilfsvariable `ti` und liefert deren

| Produktion | Regeln |
|----------------------------|-------------------------------------------------------------------------------------------------------------------------|
| $S \rightarrow V = E$ | $S.c = E.c \parallel \text{gen}(V.a \text{ '=' } E.a)$ |
| $E \rightarrow E_1 + E_2$ | $E.a = \text{newtemp}$ $E.c = E_1.c \parallel E_2.c \parallel \text{gen}(E.a \text{ '=' } E_1.a \text{ '+' } E_2.a)$ |
| $E \rightarrow V$ | $E.a = V.a$ $E.c = V.c$ |
| $E \rightarrow \text{num}$ | $E.a = \text{newconst}(\text{num}.x)$ $E.c = \text{''}$ |
| $V \rightarrow \text{id}$ | $V.a = \text{lookadr}(\text{id}.x)$ $V.c = \text{''}$ |

Abbildung 8.5: Erzeugung von Quadrupeln für die Zuweisungs-Anweisung

Adresse zurück. Die Funktion `newconst(num.x)` speichert den Wert der Konstanten `num` ab und liefert dessen Adresse zurück. Aus Gründen der leichten Lesbarkeit ist im Quadrupelcode jedoch der Wert der Konstanten angegeben und nicht dessen Adresse. Die Adressen werden im Attribut `E.a` weitergegeben.

Die Folge der Quadrupel wird im Attribut `E.c` generiert. Der Operator `"||"` verkettet Quadrupel bzw. Quadrupelfolgen miteinander. Die Funktion `gen` generiert ein einzelnes Quadrupel. Argument von `gen` ist eine Folge von Stringkonstanten und Stringvariablen. In dieser Folge ersetzt `gen` alle Variablen durch ihren Wert und liefert den gesamten String als Funktionswert. Bsp. 8.6 zeigt die für eine Zuweisung generierten Quadrupel.

```

t1 = a + b
t2 = t1 + c
a  = t2

```

Beispiel 8.6: Quadrupel für die Zuweisung `a=a+b+c`. Die Namen stehen hier für die (in der Symboltabelle eingetragenen) Speicheradressen.

8.5 Indizierte Variablen

Wir setzen voraus, daß die Untergrenze aller Indexbereiche 0 ist. Die Adresse einer indizierten Variablen `a[i1, i2, ..., ik]` wird durch folgende Formel berechnet.

$$(((...((i_1 n_2 + i_2) n_3 + i_3) \dots) n_k + i_k) w + \text{adr}(a))$$

mit

n_j Länge (=Anzahl der Elemente) der j -ten Dimension
 w Anzahl der Speichereinheiten für ein Element
 $\text{adr}(a)$ Anfangsadresse des Arrays a im Speicher

Die Längen n_j der Indexbereiche seien statisch deklariert – sie können in die Symboltabelle zu a eingetragen werden. Die Größe w ergibt sich aus dem Typ der Elemente und ist bei einem statischen Typsystem auch zur Übersetzungszeit bekannt. Bsp. 8.7 zeigt die für eine indizierte Variable generierten Quadrupel.

Formel: $(u*5+v)*4 + \text{adr}(a)$

Quadrupel: $t1 = 5 * u$
 $t2 = t1 + v$
 $t3 = 4 * t2$
 $t4 = t3 + \text{adr}(a)$
 $t5 = @ t4$
 $x = t5$

Beispiel 8.7: Gegeben ist eine Modula-2-Deklaration $a:\text{ARRAY}[10,5] \text{ OF INTEGER}$ und eine Zuweisung an eine indizierte Variable $x = a[u,v]$. Ein integer-Element benötigt 4 Bytes.

8.6 Boolesche Ausdrücke

Bei der Übersetzung Boolescher Ausdrücke ist zunächst die Frage zu klären, wie der Wert eines solchen Ausdrucks dargestellt werden soll. Es gibt zwei prinzipielle Möglichkeiten:

Numerische Methode: Die Werte werden durch Zahlen dargestellt, z.B. $\text{FALSE} = 0$ und $\text{TRUE} \neq 0$. Diese Werte können durch Maschinenbefehle verarbeitet und bei der Zuweisung an eine Variable direkt verwendet werden.

Kontrollfluß-Methode: Die Werte werden durch Positionen im Code dargestellt. Jedem Booleschen Ausdruck werden zwei Marken zugeordnet, die als Sprungziele die beiden möglichen Werte des Ausdrucks repräsentieren.

Wir zeigen das Prinzip beider Methoden an einem einfachen Ausdruck in Abb. 8.8. Dabei bedeutet $**n$ die relative Adresse n im Quadrupel-Code. Der Operator ' $<$ ' wird hier stellvertretend für alle Vergleichsoperatoren verwendet.

Bei der numerischen Methode wird ein Muster von Sprungbefehlen mit festen Relativadressen und einer temporären Variable verwendet, deren Wert (0 oder 1) als Wert des Ausdrucks weiterverarbeitet werden kann. Bei Verwendung einer AG wird

| Numerische Methode | Kontrollfluß-Methode |
|--------------------|----------------------|
| if a<b goto **3 | if a<b goto Ltrue |
| t1 = 0 | goto Lfalse |
| goto **2 | |
| t1 = 1 | |

Abbildung 8.8: Quadrupel-Darstellungen des Ausdrucks a<b

die Adresse der temporären Variable als synthetisiertes Attribut des Ausdrucks der Umgebung bekanntgegeben.

Die Kontrollfluß-Methode verwendet ein parametrisiertes Muster: Ltrue und Lfalse sind formale Wertparameter, denen symbolische Marken als aktuelle Sprungziele übergeben werden. Dem Ausdruck wird also „mitgeteilt“, wohin im Falle true bzw. false gesprungen werden soll. Aus der Sicht einer AG sind Ltrue und Lfalse zwei ererbte Attribute des Ausdrucks.

Die Übersetzung Boolescher Ausdrücke nach der Kontrollfluß-Methode wird durch die AG von Abb. 8.9 formal beschrieben. Die Funktion newlabel erzeugt – ähnlich wie die Funktion newtemp – eindeutige symbolische Namen. Das Grundmuster von Abb. 8.8 findet man als Regel zur Vergleichsoperation rop in Abb. 8.9 wieder. Man beachte, daß die Operatoren and und or im erzeugten Code nicht vorkommen. Für den Operator not wird überhaupt kein Code generiert, es werden nur die Sprungziele vertauscht.

Für die Operatoren and und or wird „jumping code“ generiert, d.h. die Ausdrücke werden nicht unbedingt zur Gänze ausgewertet. Bsp. 8.10 zeigt den Code für einen Booleschen Ausdruck mit dem Operator and. Falls der linke Operand a<b den Wert 'false' ergibt, wird das externe Sprungziel L2 verwendet, d.h. der zweite Operand c<d wird nicht mehr ausgewertet.

Abb. 8.11 veranschaulicht, wie die Sprungziele als Attribute gemäß der AG von Abb. 8.9 übergeben werden, um den „jumping code“ zu erzeugen. Jedem Booleschen Ausdruck B sind die beiden ererbten Attribute B.true und B.false zugeordnet. Ihre Werte sind die beiden Sprungziele des Ausdrucks im Falle 'true' bzw. 'false'. Bei einem and-Operator wird das 'false'-Sprungziel direkt an das Attribut B.false des linken Teilausdrucks weitergegeben, während das Attribut B.true des linken Teilausdrucks eine neu generierte Marke als Sprungziel auf den zweiten Teilausdruck erhält. Ausdrücke mit or werden analog behandelt (siehe AG in Abb. 8.9).

8.7 Kontrollanweisungen

Die Übersetzung von Kontrollanweisungen hängt davon ab, ob der steuernde Ausdruck mit der numerischen oder der Kontrollfluß-Methode arbeitet. Wir werden daher beide Methoden berücksichtigen und jeweils als typische Vertreter von Kontrollanweisungen die IF- und die WHILE-Anweisung behandeln.

| Produktion | Regeln |
|--------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| $B \rightarrow B_1 \text{ or } B_2$ | $B_1.\text{true} = B.\text{true}$ $B_1.\text{false} = \text{newlabel}$ $B_2.\text{true} = B.\text{true}$ $B_2.\text{false} = B.\text{false}$ $B.c = B_1.c \parallel \text{gen}(B_1.\text{false} \text{ ':'}) \parallel B_2.c$ |
| $B \rightarrow B_1 \text{ and } B_2$ | $B_1.\text{true} = \text{newlabel}$ $B_1.\text{false} = B.\text{false}$ $B_2.\text{true} = B.\text{true}$ $B_2.\text{false} = B.\text{false}$ $B.c = B_1.c \parallel \text{gen}(B_1.\text{true} \text{ ':'}) \parallel B_2.c$ |
| $B \rightarrow \text{not } B_1$ | $B_1.\text{true} = B.\text{false}$ $B_1.\text{false} = B.\text{true}$ $B.c = B_1.c$ |
| $B \rightarrow E_1 \text{ rop } E_2$ | $B.c = E_1.c \parallel E_2.c \parallel$ $\text{gen}(\text{'if' } E_1.a \text{ rop.x } E_2.a \text{ 'goto' } B.\text{true}) \parallel$ $\text{gen}(\text{'goto' } B.\text{false})$ |
| $B \rightarrow \text{true}$ | $B.c = \text{gen}(\text{'goto' } B.\text{true})$ |
| $B \rightarrow \text{false}$ | $B.c = \text{gen}(\text{'goto' } B.\text{false})$ |

Abbildung 8.9: Übersetzung Boolescher Ausdrücke nach der Kontrollfluß-Methode

```

    if a<b goto L3
    goto L2
L3:  if c<d goto L1
    goto L2

```

Beispiel 8.10: Quadrupel für den Booleschen Ausdruck $(a < b)$ and $(c < d)$. Die externen Sprungziele für die Fälle 'true' und 'false' seien L1 und L2.

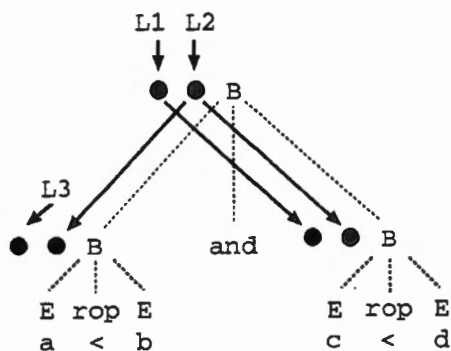


Abbildung 8.11: Syntaxbaum und Attributübergabe zu Bsp. 8.10

Für die **numerische Methode** sind die zu erzeugenden Sprungstrukturen und eine AG zur Erzeugung dieser Strukturen in Abb. 8.12 zusammengefaßt.

Für die **Kontrollfluß-Methode** sind die zu erzeugenden Sprungstrukturen und die AG in Abb. 8.13 zusammengefaßt.

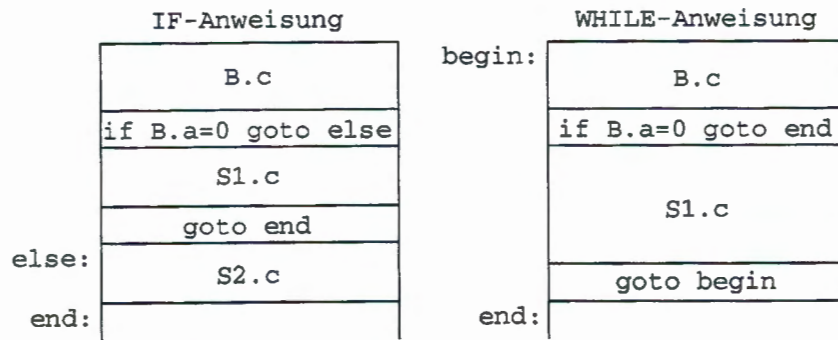
Bsp. 8.14 vergleicht die numerische und die Kontrollfluß-Methode bei der Übersetzung des Programmstücks zur Berechnung des ggT. Die Kontrollfluß-Methode ist für dieses Beispiel günstiger.

Der Code kann noch weiter vereinfacht werden (Nachoptimierung), indem man die beiden Vergleichsoperatoren umkehrt, d.h. \neq durch $=$, und $>$ durch \leq ersetzt. Bsp. 8.15 zeigt das optimierte Programmstück zur Kontrollflußmethode aus Bsp. 8.14. Sprünge zu einer Marke, an der wieder ein unbedingter Sprung steht, können durch diesen ersetzt werden, wie im Beispiel `goto L6` durch `goto L1`. Ein Sprung auf die dem Sprungbefehl unmittelbar folgende Stelle kann ebenfalls eliminiert werden.

8.8 Kontrollflußgraph und Datenflußgraph

Der **Kontrollflußgraph** beschreibt alle strukturell möglichen Abläufe eines Programms. Jeder Knoten stellt einen Grundblock (*basic block*) und jede Kante einen möglichen Kontrollübergang dar. Ein Grundblock ist eine Anweisungsfolge, die genau dann ausgeführt wird, wenn die erste Anweisung der Folge ausgeführt wird. Der Kontrollflußgraph kann aus dem Quadrupelcode eines Programms konstruiert werden. Abb. 8.16 zeigt den Graphen für das Programm ggT von Bsp. 8.14.

Jeder Grundblock kann vollständig durch einen **Datenflußgraphen** beschrieben werden. Die Knoten sind mit Operatoren markiert. Die Kanten repräsentieren Daten, die am Ausgangsknoten erzeugt und am Zielknoten verwendet werden. Knoten ohne Vorgänger sind als nullstellige Operatoren aufzufassen, d.h. jeder Knoten des Graphen erzeugt bzw. repräsentiert einen Wert. Datenflußgraphen lassen sich aus Quadrupelfolgen konstruieren, indem man die Teilgraphen der einzelnen Quadrupel aneinanderhängt. Abb. 8.17 zeigt den Datenflußgraphen für die Quadrupelfolge von



| Produktion | Regeln |
|------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| $S \rightarrow \text{if } B \text{ then } S_1 \text{ else } S_2$ | <pre> else = newlabel end = newlabel S.c = B.c gen('if' B.a '=0 goto' else) S1.c gen('goto' end) gen(else ':') S2.c gen(end ':') </pre> |
| $S \rightarrow \text{while } B \text{ do } S_1$ | <pre> begin = newlabel end = newlabel S.c = gen(begin ':') B.c gen('if' B.a '=0 goto' end) S1.c gen('goto' begin) gen(end ':') </pre> |

Abbildung 8.12: Übersetzung von Kontrollanweisungen bei der numerischen Methode

Bsp. 8.7.

Insgesamt bietet der Quadrupelcode also die Möglichkeit, Programmdarstellungen zu erzeugen, bei denen einerseits der Kontrollfluß und andererseits der Datenfluß betrachtet wird. Diese Aufteilung ist zweckmäßig für Methoden der Optimierung und Codeerzeugung.

| Produktion | Regeln |
|------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| $S \rightarrow \text{if } B \text{ then } S_1 \text{ else } S_2$ | <pre> B.true = newlabel B.false = newlabel end = newlabel S.c = B.c gen(B.true ':') S1.c gen('goto' end) gen(B.false ':') S2.c gen(end ':') </pre> |
| $S \rightarrow \text{while } B \text{ do } S_1$ | <pre> begin = newlabel B.true = newlabel B.false = newlabel S.c = gen(begin ':') B.c gen(B.true ':') S1.c gen('goto' begin) gen(B.false ':') </pre> |

Abbildung 8.13: Übersetzung von Kontrollanweisungen bei der Kontrollfluß-Methode

```

WHILE x≠y DO
  IF x>y THEN
    x = x-y
  ELSE
    y = y-x
  END
END
END

```

| Numerische Methode | Kontrollfluß-Methode |
|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------|
| L1: if x≠y goto *+3 t1 = 0 goto *+2 t1 = 1 if t1=0 goto L2 if x>y goto *+3 t2 = 0 goto *+2 t2 = 1 if t2=0 goto L3 t3 = x-y x = t3 goto L4 | L1: if x≠y goto L2 goto L3 L2: if x>y goto L4 goto L5 L4: t1 = x-y x = t1 goto L6 L5: t2 = y-x y = t2 L6: goto L1 L3: |
| L3: t4 = y-x y = t4 | |
| L4: goto L1 | |
| L2: | |

Beispiel 8.14: Zum Vergleich werden die numerische und die Kontrollfluß-Methode zur Übersetzung eines Quellprogramms (ggT) angewendet.

```

L1: if x=y goto L3
    if x≤y goto L5
    t1 = x-y
    x = t1
    goto L1
L5: t2 = y-x
    y = t2
    goto L1
L3:

```

Beispiel 8.15: Durch Anwendung der Nachoptimierung auf den mit der Kontrollfluß-Methode entstandenen Code aus Bsp. 8.14 entsteht dieses Programm.

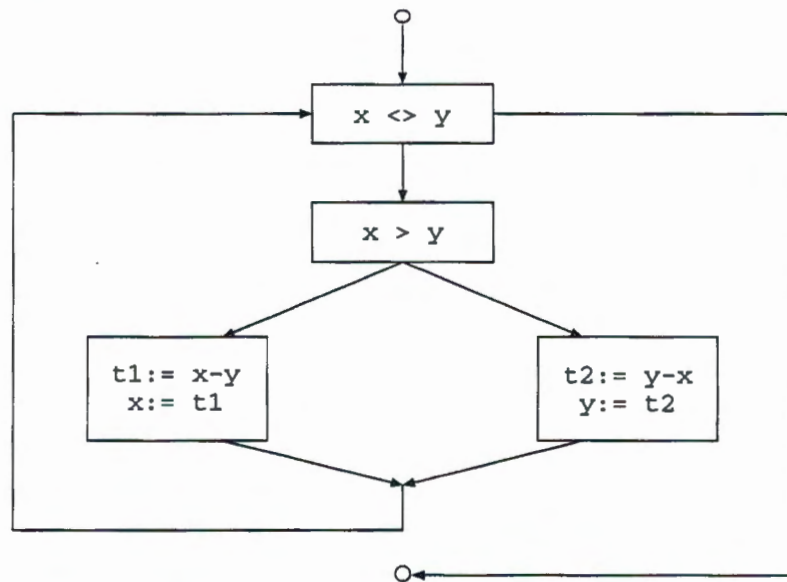


Abbildung 8.16: Kontrollflußgraph für das Programm ggT von Bsp. 8.14.

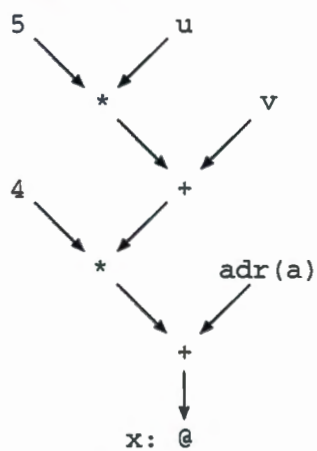


Abbildung 8.17: Datenflußgraph für die Quadrupelfolge von Bsp. 8.7.

Kapitel 9

Codeerzeugung

Die Codeerzeugung wandelt die maschinenunabhängige Zwischendarstellung in Assembler- oder Maschinencode um. Sie besteht aus drei Teilaufgaben:

Befehlsauswahl (*instruction selection, code selection*). Die Operationen der Zwischendarstellung werden zu Befehlen der Zielmaschine zusammengefaßt.

Befehlsanordnung (*instruction scheduling*). Die Befehle werden im Hinblick auf eine schnelle Ausführung auf Pipeline- oder superskalaren Prozessoren angeordnet.

Registerbelegung (*register allocation*). Die während früherer Phasen verwendeten Pseudoregister werden durch Maschinenregister ersetzt.

Die Codeerzeugung wird oft als maschinenabhängiger Teil des Compilers bezeichnet. Die im folgenden vorgestellten Methoden lassen sich jedoch auf eine große Klasse von Prozessorarchitekturen anwenden, zu der fast alle heute in PCs und größeren Computern verwendeten Prozessoren gehören. Die Methoden müssen nur mit den entsprechenden Prozessordaten parametrisiert werden. In diesem Kapitel zeigen wir die Probleme der Codeerzeugung anhand der Prozessorarchitekturen 68000 und MIPS.

Ein Befehl in 68000-Assembler hat folgende Syntax:

Name.Größe Operand1,Operand2

Der Name des Befehls ist meist selbsterklärend, die Größe gibt an, auf welchen Daten der Befehl arbeitet. *b* steht für 1 byte, *w* für 2, und *l* für 4. Von den Operanden ist der rechte der Zieloperand, viele Befehle (z.B. *add*) verwenden ihn auch als Eingabeoperanden. Für die Operanden kann eine Vielfalt von Adressierungsarten verwendet werden, auf die wir hier nicht eingehen. Der 68000 hat zwei Klassen von Registern, Adreßregister (*An*) und Datenregister (*Dn*).

Befehle im MIPS-Assembler haben folgende Form:

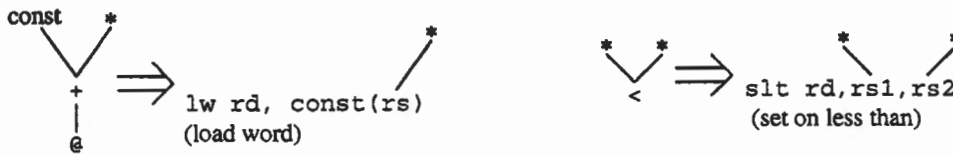
Name Operand1,Operand2,Operand3

Dabei ist im allgemeinen der linke Operand der Zielooperand, die anderen beiden die Eingabeoperanden. Die Befehle arbeiten auf Registern und mit Konstanten, Zugriffe auf den Speicher erfolgen über spezielle Lade- und Speicherbefehle der Form *Name Register, const(Register)*. Der zweite Operand gibt dabei die Adresse an.

9.1 Befehlsauswahl

Die Befehlsauswahl erfolgt auf Grundblockebene, wobei Grundblöcke als Datenflußgraphen dargestellt werden. Zur Vereinfachung betrachten wir zunächst aber nur Bäume¹, also Datenflußgraphen, bei denen jedes Datum nur einmal verwendet wird.

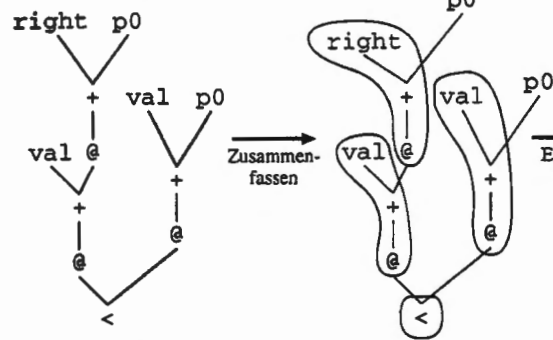
Befehlsbäume:



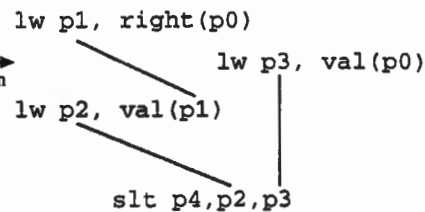
rd ist das Ergebnisregister, rs* sind die Eingaberegister

C-Code: `p->right->val < p->val`

Datenflußgraph (Zwischencode):



Datenflußgraph (Maschinencode):



p0 ist das Pseudoregister der Variablen p, p1-p4 sind Pseudoregister für Zwischenergebnisse
right und val sind Konstanten

Abbildung 9.1: Befehlsauswahl für die MIPS-Architektur

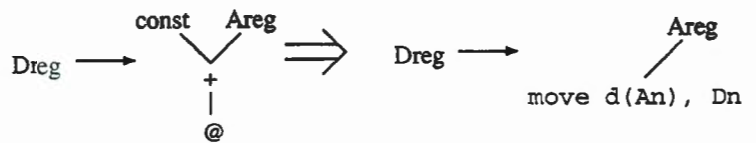
Die Knoten des Datenflußgraphen sind die Operatoren der Zwischensprache. Maschinenbefehle entsprechen Bäumen der Zwischensprache (Ausnahmen werden später besprochen). Die Aufgabe der Befehlsauswahl ist es, den Datenflußgraphen

¹Datenflußgraphen werden üblicherweise so gezeichnet, daß die Daten von oben nach unten fließen. Bei Reduktion auf Bäume ergibt sich dabei die etwas ungewöhnliche Darstellung mit der Wurzel nach unten. Bei einer Expedition in den Resselpark können Sie sich überzeugen, daß diese Darstellung den Grundsätzen der Natur nicht widerspricht.

in Bäume aufzuteilen, die Maschinenbefehlen entsprechen. Die Bäume werden dann durch die entsprechenden Befehle ersetzt. Das Ergebnis ist ein Datenflußgraph, in dem die Knoten Maschinenbefehle sind. Abb. 9.1 zeigt beispielhaft die Befehlsauswahl für die MIPS-Architektur.

Im Baum eines Befehls steht dort ein *, wo der Baum eines weiteren Befehls eingesetzt werden kann. Allerdings können nicht immer alle Befehle so einfach kombiniert werden. So kann z.B. auf dem 68000 der Befehl „move (An), Dn“ das Ergebnis des Befehls „and #const, Dn“ nicht unmittelbar verwenden, da es in einem Datenregister liegt, aber in einem Adreßregister gebraucht wird. Beim 68000 müssen wir also zwei Arten von Sternen verwenden. Die eine (nennen wir sie Dreg) steht für Bäume, die ihr Ergebnis in einem Datenregister abliefern, die andere (Areg) steht entsprechend für Bäume mit dem Ergebnis in einem Adreßregister.

Die Beziehung zwischen dem Zwischencode und dem Maschinencode kann mit Regeln wie der folgenden beschrieben werden:



Die beiden Teile der Regeln schauen ähnlich wie Grammatikproduktionen aus. Man spricht daher von Baumgrammatiken, Areg und Dreg werden als Nonterminale bezeichnet. Die Beziehung zwischen Zwischencodebaum und Maschinencodebaum ist: Wenn ein Baum auf der einen Seite mit einer bestimmten Folge von Produktionen dieser Seite abgeleitet werden kann, dann erhält man durch die entsprechende Folge von Produktionen auf der anderen Seite einen äquivalenten Baum in der anderen Darstellung.

In der Praxis werden die Regeln etwas anders geschrieben:

$$Dreg \rightarrow @(+(const, Areg)) \text{ move } const(An), Dn$$

Bäume werden linearisiert als prefix-Ausdrücke dargestellt. Und die Produktion im Maschinenteil der Regel wird einfach durch den Befehl ersetzt. Die Produktion kann folgendermaßen rekonstruiert werden: Die linke Seite ist das gleiche Nonterminal wie im Zwischendarstellungsteil und die rechte Seite enthält ebenfalls dieselben Nonterminale wie die rechte Seite des Zwischendarstellungsteils; diese Nonterminale werden als Kinder direkt an den Maschinenbefehl gehängt.

Die Einführung verschiedener Nonterminale anstatt des * hat einen weiteren Vorteil: Prozessoren wie der 68000 können die Befehle mit vielen Adressierungsarten kombinieren. Anstatt jede Kombination von Befehl und Adressierungsart hinzuschreiben, kann man eine Regel für eine Klasse von Adressierungsarten bilden, und sie in allen entsprechenden Befehlen verwenden.

Baumgrammatiken für die Befehlsauswahl müssen jeden Zwischencode-Baum beschreiben können. Dabei stört es nicht, wenn die Grammatik mehrdeutig ist; Die verschiedenen Ableitungen für den Baum entsprechen den verschiedenen Möglichkeiten, den Baum in Maschinensprache zu übersetzen.

Von diesen Möglichkeiten sollte eine optimale bezüglich Ausführungszeit und/oder Speicherverbrauch ausgewählt werden. Daher fügen wir zu jeder Regel noch eine Zahl, die die Kosten des entsprechenden Maschinenbefehls beschreibt, also die Ausführungszeit oder den Platzbedarf des Befehls. Der im folgenden beschriebene Algorithmus liefert dann ein optimales Ergebnis, wenn die Gesamtkosten sich aus der Summe der Kosten für die einzelnen Regeln ergeben. Das ist bei heutigen Prozessoren nicht immer der Fall, da Befehle überlappt ausgeführt werden können. Trotzdem ist die Annahme eines additiven Verhaltens eine gute Annäherung.

Für jeden Knoten n im Baum, von den Blättern ausgehend:

t ist der Teilbaum, dessen Wurzel n ist

Für jedes Nonterminal A :

$cost[n][A] = \infty$

Für jede Regel R der Form „ $A \rightarrow$ Baummuster“ (ausgenommen Kettenregeln), sodaß Baummuster auf t paßt:

$cost1 = R.cost$

$tree1.op = n.op$

$childno = 1$

Für jedes Nonterminal B in Baummuster:

u ist der Knoten an der Stelle in t ,

die B in Baummuster einnimmt

$cost1 = cost1 + cost[u][B]$

$tree1.child[childno] = tree[u][B]$

$childno = childno + 1$

Wenn $cost1 < cost[n][A]$ dann

$cost[n][A] = cost1$

$tree[n][A] = tree1$

Wende Kettenregeln an wie oben die anderen bis sich nichts mehr ändert

Algorithmus 9.2: Befehlsauswahl-Algorithmus

Der Befehlsauswahl-Algorithmus (siehe Alg. 9.2) geht von den Blättern aus zur Wurzel vor und ermittelt für jeden Knoten n und jedes Nonterminal A eine Regel R mit folgenden Eigenschaften: Aus A kann mit der Regel der Teilbaum abgeleitet werden, dessen Wurzel n ist; und diese Ableitung hat minimale Kosten. Diese Berechnung greift auf die optimale Befehlsauswahl und die zugehörigen Kosten für untergeordnete Baum-Nonterminal-Kombinationen zurück.

Konkret geschieht folgendes: Die minimalen Kosten für eine Baum-Nonterminal-Kombination bei Verwendung einer bestimmten Regel an der Wurzel berechnen sich als Summe der Kosten der Regel und der Kosten der von der Regel verwendeten Teilbaum-Nonterminal-Kombinationen. Die Kosten aller Teilbäume sind schon bekannt, da sich der Algorithmus von den Blättern zur Wurzel vorarbeitet. Zu jedem

Knoten werden für jedes Nonterminal die Gesamtkosten und der Maschinencodebaum mit den niedrigsten Gesamtkosten gespeichert.

Ein Problem stellen dabei die Kettenregeln dar, also Regeln der Form *Nonterminal1* \rightarrow *Nonterminal2*. Bei der Berechnung von Kosten mit Hilfe einer Kettenregel würde auf die Kosten eines Nonterminals desselben Teilbaumes zurückgegriffen, die eventuell noch nicht oder nicht vollständig berechnet sind. Eine einfache Lösung ist es, Kettenregeln nach allen anderen zu berücksichtigen, und, wenn sich die Kosten des Baums ändern, noch einmal alle Kettenregeln zu berücksichtigen, solange bis sich nichts mehr ändert. Diese Wiederholung der Kettenregelberechnung dient für den Fall, daß eine Kettenregel von einer anderen abhängt.

| Nr. | Regel | Befehl | Kosten |
|-----|---------------------------------------------------|----------------------|--------|
| 1 | Dreg \rightarrow Areg | move.l An, Dn | 4 |
| 2 | Areg \rightarrow Dreg | move.l Dn, An | 4 |
| 3 | Areg \rightarrow @(Areg) | move.l (An), An | 12 |
| 4 | Dreg \rightarrow @(Areg) | move.l (An), Dn | 12 |
| 5 | Dreg \rightarrow +(Areg, Dreg) | add.l An, Dn | 8 |
| 6 | Dreg \rightarrow +(Dreg, Dreg) | add.l Dn, Dn | 8 |
| 7 | Areg \rightarrow +(Dreg, Areg) | add.l Dn, An | 8 |
| 8 | Areg \rightarrow +(Areg, Areg) | add.l An, An | 8 |
| 9 | Dreg \rightarrow +(@(Areg), Dreg) | add.l (An), Dn | 14 |
| 10 | Areg \rightarrow +(@(Areg), Areg) | add.l (An), An | 14 |
| 11 | Dreg \rightarrow @ _b (Areg) | move.b (An), Dn | 8 |
| 12 | Dreg \rightarrow @ _b (+(Areg, Dreg)) | move.b 0(An, Dn), Dn | 14 |
| 13 | Dreg \rightarrow @ _b (+(Areg, Areg)) | move.b 0(An, An), Dn | 14 |
| 14 | Areg \rightarrow an | ϵ | 0 |
| 15 | Dreg \rightarrow dn | ϵ | 0 |

Abbildung 9.3: Teil einer Befehlsauswahl-Baumgrammatik für den 68000

Abb. 9.3 zeigt einige Regeln, die sich im Beispiel in Abb. 9.4 anwenden lassen. an ist ein Terminalsymbol, das für ein Adreßregister steht, entsprechend steht dn für ein Datenregister. Die Kosten einer Regel sind die Laufzeit des Befehls in Maschinenzyklen. Der Zwischencode-Baum in Abb. 9.4 entspricht dem C-Ausdruck (*p) [x], wobei p auf einen Zeiger auf ein Zeichen-Array zeigt und in einem Adreßregister liegt; x liegt in einem Datenregister. Das Ergebnis soll in einem Datenregister landen. Bei jedem Knoten sind für jedes Nonterminal die Kosten der optimalen Befehlsauswahl für den entsprechenden Unterbaum eingetragen. „tree“ enthält einen optimalen Maschinencodebaum für den entsprechenden Teilbaum und das Nonterminal, der allerdings aus Platzgründen linearisiert und mit Regelnummern statt der Maschinenbefehle dargestellt ist.

Betrachten wir nun, wie die Berechnung für den Wurzelknoten vor sich geht: Das Nonterminal Areg kann den Baum nur über die Kettenregel 2 ableiten; wir

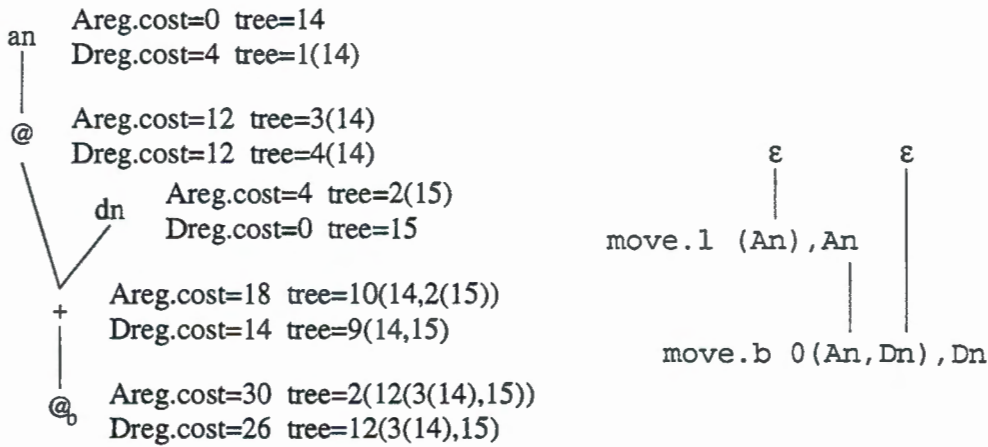


Abbildung 9.4: Optimale Befehlsauswahl

betrachten diese Variante daher später. Für Dreg kommen die Regeln 11, 12, und 13 in Frage. Die Berechnung der Kosten für die Regel 12 verläuft folgendermaßen: Die Regel verwendet den @-Knoten des Baums als Areg (Kosten: 12) und den dn-Knoten des Baums als Dreg (Kosten: 0), die Regel selbst kostet 14 Zyklen, die Gesamtkosten sind daher 26 Zyklen. Für Regel 11 ergeben sich ebenfalls Kosten von 26 Zyklen, für Regel 13 dagegen 30 Zyklen. Die Regeln 11 und 12 sind daher optimal. Wir wählen willkürlich Regel 12. Der zugehörige Maschinencodebaum ergibt sich wie folgt: Der Operator ist der zur Regel gehörige Befehl `move.b 0(An,Dn),Dn`. Der Baum hat für jedes Nonterminal der rechten Seite der Zwischencodeproduktion ein Kind, und zwar für Areg den optimalen Maschinencodebaum für Areg am @-Knoten: `move.l (An),An(ε)` und für Dreg den optimalen Maschinencodebaum für Dreg am Knoten dn: `ε`. Insgesamt ergibt sich also der in Abb. 9.4 rechts gezeigte Baum.

Nun können die Kettenregeln angewandt werden, z.B. die Regel 2, um die optimalen Kosten für das Nonterminal Areg des Wurzelknotens zu berechnen: Sie ergeben sich aus den 26 Zyklen von Dreg im Wurzelknoten und 4 Zyklen für die Regel selbst, also insgesamt 30 Zyklen. Der Maschinencodebaum hat als Operator den Befehl `move.l Dn, An` und nur ein Kind: den optimalen Maschinencodebaum für Dreg am Knoten @_b, den wir gerade vorher ermittelt haben.

Jetzt kann auch die Regel 1 für die Berechnung von Dreg angewandt werden, dadurch ergibt sich aber natürlich keine Verbesserung (Kosten: 34); die Regeln 1 und 2 fügen ja nur Kopierbefehle ein, deren wiederholte Anwendung das Programm wohl nicht schneller macht.

Die Eingabe der Befehlsauswahl sind meist nicht Bäume, sondern azyklische Graphen (DAGs). Leider gibt es für DAGs im allgemeinen kein ähnlich einfaches, optimales und effizientes Verfahren wie das oben beschriebene für Bäume. Ein Lösungsweg ist, den DAG in Bäume aufzuteilen, indem alle Kanten entfernt werden, die zwischen einem Knoten mit mehreren Eltern und diesen Eltern verlaufen. Der Code für die Teilbäume ist optimal, der Code für den gesamten DAG nicht unbedingt,

aber oft.

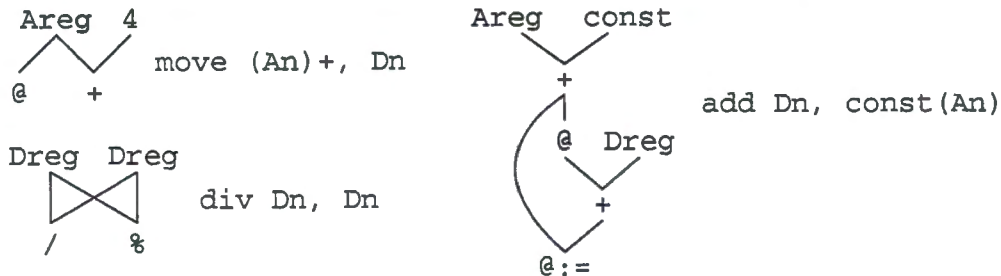


Abbildung 9.5: Befehle, deren Operatordarstellung kein Baum ist

Ein weiteres Problem stellen Befehle dar, deren Operatordarstellung kein Baum ist (siehe Abb. 9.5). Dazu zählen Befehle mit mehreren Resultaten wie z.B. Divisionsbefehle, die auch noch den Rest liefern, und Befehle mit der Autoincrement-Adressierungsart, aber auch viele Befehle, die auf Speicheroperanden arbeiten (read-modify-write-Befehle). Solche Befehle können nicht mit Verfahren erzeugt werden, die auf Baumustererkennung beruhen. Die einfachste Lösung ist, sich damit zufriedenzugeben, daß der Compiler diese Befehle nicht verwendet. Eine andere Möglichkeit ist, nachträglich geeignete Befehlskombinationen zusammenzufassen, wobei diese Befehle entstehen.

9.2 Befehlsanordnung

Die Umformung des Datenflußgraphen eines Grundblocks in eine lineare Form, wie sie der Prozessor erwartet, ist Aufgabe der Befehlsanordnung. Bei Prozessoren, die durch Pipelining und/oder superskalare Architektur Befehle parallel abarbeiten, ist das Ziel der Befehlsanordnung eine gute Ausnutzung des Parallelismus durch Minimierung von ungenutzten Wartezyklen (*delay slots, pipeline bubbles*).

Auf dem MIPS R4000-Prozessor kann zum Beispiel einen Zyklus nach dem Start eines Ladebefehls ein anderer Befehl ausgeführt werden. Das Ergebnis des Ladebefehls steht aber erst nach drei Zyklen zur Verfügung. Wird ein Befehl gestartet, der das Ergebnis verwendet, dann muß er so lange warten und es entstehen ungenutzte Wartezyklen. Wenn die Befehle so angeordnet werden, daß nach einem Ladebefehl zwei Befehle folgen, die das Ergebnis nicht verwenden, wird der Prozessor optimal ausgenutzt. Auch bei superskalaren Prozessoren wie dem 586 tritt das Problem auf: Nur unabhängige Befehle können im gleichen Zyklus ausgeführt werden. Ein verwandtes Problem gibt es auch bei Verzweigungsbefehlen in einigen RISC-Architekturen. Bei der MIPS-Architektur wird ein Verzweigungsbefehl erst mit Verspätung ausgeführt, vorher wird noch der unmittelbar hinter dem Verzweigungsbefehl liegende Befehl ausgeführt (*delayed branch*). Dieser Branch-delay-slot sollte natürlich mit einem sinnvollen und unabhängigen Befehl gefüllt werden.

Zunächst stellt sich die Frage, welche Befehlsanordnungen überhaupt korrekt sind. Ein Befehl a muß in folgenden Fällen vor einem anderen Befehl b ausgeführt werden:

Datenflußabhängigkeit (*flow dependence, RAW dependence*): Wenn b von einem Register oder Speicherplatz liest, der von a beschrieben wurde.

Lese-Schreib-Abhängigkeit (*anti dependence, WAR dependence*): Wenn b auf eine Stelle schreibt, die zuvor von a gelesen wurde.

Schreib-Schreib-Abhängigkeit (*output dependence, WAW dependence*): Wenn b auf eine Stelle schreibt, die zuvor von a beschrieben wurde.

potentielle Abhängigkeit Wenn nicht sicher ist, daß zwei Zugriffe auf verschiedene Stellen erfolgen, z.B. bei Zugriffen über Zeiger, muß die Annahme getroffen werden, daß eine Abhängigkeit zwischen den zwei Zugriffen besteht.

Die Abhängigkeiten zwischen den Befehlen bilden einen Datenabhängigkeitsgraphen (*data dependence graph*), das ist ein um Lese-Schreib-, Schreib-Schreib-, und potentielle Abhängigkeitskanten erweiterter Datenflußgraph. Jede Befehlsanordnung, die die Abhängigkeiten berücksichtigt, ist korrekt.² Die Abb. 9.6 zeigt den Abhängigkeitsgraphen für eine Schleife in MIPS-Assembler; Integerregister werden in der Form $\$n$ geschrieben, Fließkommaregister in der Form $\$fn$. $\$32$ ist ein Label.

Wenn es um die beste Lösung geht, ist das Problem der Befehlsanordnung im allgemeinen NP-vollständig, d.h. alle bekannten Algorithmen weisen exponentielles Laufzeitverhalten auf. Daher werden in der Praxis heuristische Algorithmen verwendet, die in vertretbarer Zeit gute, wenn auch nicht immer optimale Befehlsanordnungen finden.

Der Standard-Algorithmus für Befehlsanordnung heißt *List Scheduling* und ist in Alg. 9.7 dargestellt: Einer der Befehle ohne Vorgänger im Graphen (ein sogenannter Anführer) wird ausgewählt und aus dem Graphen entfernt. Dieser Schritt wird wiederholt bis der Graph leer ist. Die Reihenfolge, in der die Befehle entfernt wurden, ist die Befehlsanordnung.

Dieser Algorithmus läßt offen, welcher Befehl ausgewählt wird. Dies ist die Aufgabe der heuristischen Auswahlfunktion „wähle“. Eine typische Auswahlfunktion zur Ausnutzung des Parallelismus lautet:

- Wähle einen Befehl, der als nächstes ausgeführt werden kann; also einen, der möglichst kurz oder am besten gar nicht warten muß, bis andere Befehle seine Eingabedaten berechnet haben. Wenn es mehrere Kandidaten gibt,
- Wähle einen Befehl mit maximaler Pfadlänge zum Ende des Datenabhängigkeitsgraphen. Zur Berechnung der Pfadlänge geben wir jeder Kante eine Länge: Die Länge einer Abhängigkeitskante ist die Zeit, die nach dem Start des Befehls am Anfang der Kante vergehen muß, bis der Befehl am Ende der Kante ausgeführt werden kann. Diese Heuristik sorgt dafür, daß Delay-slots früh freigelegt werden, solange sie noch gefüllt werden können.


```

for (i=0; i<n; i++) {
    y[iy] = y[iy]+a*x[ix];
    ix += incx;
    iy += incy;
}

```

```

$32:
A mul    $14, $4, 8
B addu   $6, $5, $14
C mul    $15, $7, 8
D addu   $24, $15, $8
E l.d    $f4, 0($24)
F mul.d  $f6, $f4, $f12
G l.d    $f8, 0($6)
H add.d  $f10, $f6, $f8
I s.d    $f10, 0($6)
J addu   $7, $7, $9
K addu   $4, $4, $10
L addu   $2, $2, 1
M bne    $11, $2, $32

```

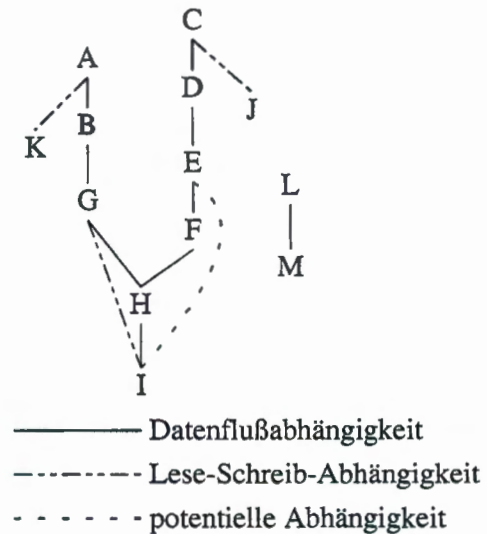


Abbildung 9.6: Eine Schleife in C und MIPS-Assembler und ihr Datenabhängigkeitsgraph

```

list_scheduling(Graph)
  Anordnung=leer
  while Graph ≠ leer
    Anführer={Knoten im Graph ohne Vorgänger}
    Nächster-Befehl=wähle(Anführer)
    Hänge Nächster-Befehl an Anordnung an
    Entferne Nächster-Befehl aus Graph

```

Algorithmus 9.7: List Scheduling

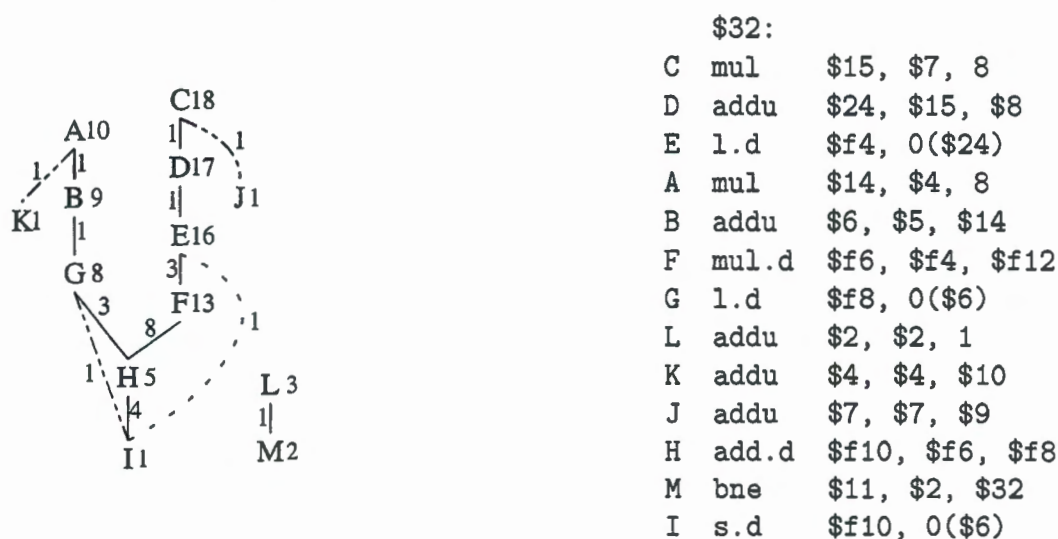


Abbildung 9.8: Der Datenabhängigkeitsgraph mit Kanten- und Pfadlängen und der Code nach der Befehlsanordnung

Abb. 9.8 zeigt den Datenabhängigkeitsgraphen dekoriert mit Kantenlängen für den R4000. Weiters ist bei jedem Knoten die Pfadlänge bis zum Ende dargestellt. Für den Verzweigungsbefehl M ist dabei eine Pfadlänge von zwei Zyklen bis zum Ende angegeben, da es von der Ausführung des Verzweigungsbefehls bis zum Ende des Grundblocks eben zwei Zyklen dauert (*delayed branch*).

Am Anfang der Befehlsanordnung sind A, C, und L die Anführer und sie könnten alle sofort ausgeführt werden. Aufgrund der Pfadlänge wird C gewählt. Dadurch werden D und J zu Anführern, die ebenfalls sofort ausgeführt werden könnten. Als nächstes werden die Befehle D und E ausgewählt. Dann gibt es mit F erstmals einen Anführer, der nicht sofort ausgeführt werden kann. Aber noch gibt es genügend andere Kandidaten. Im weiteren werden A, B, F, G und L ausgewählt. Jetzt wird M zu einem Anführer, aber als Verzweigungsbefehl darf er erst als vorletzter oder letzter Befehl gewählt werden. Daher lautet die weitere Folge K, J (oder J, K) und H. H muß zwar noch drei Zyklen auf das Ergebnis von F warten, aber H ist der einzige wählbare Anführer und wird daher ausgewählt. Nun kann endlich M gewählt werden, und in den Delay-slot von M kommt der Befehl I, der weitere zwei Zyklen auf das Ergebnis von H warten muß. Der Grundblock benötigt bei dieser Befehlsanordnung 18 Zyklen, bei der ursprünglichen 25. Die Befehlsanordnung ist sogar optimal: Da der kritische Pfad des Datenabhängigkeitsgraphen 18 Zyklen lang ist, gibt es sicher keine kürzere.

²Graphentheoretiker nennen solche Anordnungen eines Graphen „topologisch sortiert“.

9.3 Registerbelegung

Moderne Compiler erzeugen zunächst Code, der beliebig viele Pseudoregister verwendet. Die Aufgabe der Registerbelegung ist es, die Pseudoregister möglichst alle in Maschinenregistern unterzubringen, bzw. wenn das nicht gelingt, weniger wichtige Pseudoregister in den Speicher auszulagern. Im Gegensatz zur Befehlsauswahl und der Befehlsanordnung, die auf der Grundblockebene („lokal“) arbeiten, erfolgt die Registerbelegung heute für eine ganze Prozedur („global“ im irreführenden Compilerbau-Jargon).

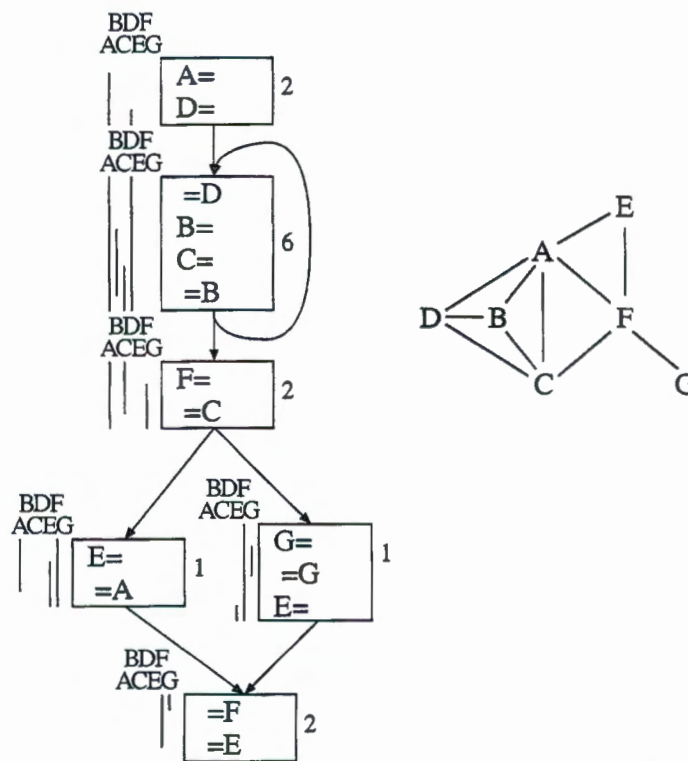


Abbildung 9.9: Ein Programm, die Aktivitätsbereiche seiner Pseudoregister, und sein Konfliktgraph

Ausgangspunkt der Registerbelegungsalgorithmen ist die Erkenntnis, daß zwei Pseudoregister im selben Maschinenregister untergebracht werden können, wenn sie nicht zur gleichen Zeit aktiv (*live*) sind. Daraus ergibt sich die grundlegende Datenstruktur für die Registerbelegung: Der Konfliktgraph (*interference graph, conflict graph*) enthält für jedes Pseudoregister einen Knoten. Zwei Knoten sind durch eine Kante verbunden, wenn die beiden Pseudoregister zur gleichen Zeit aktiv sind. Abb. 9.9 zeigt ein Programm, und zwar die für die Registerbelegung wichtigen Teile, nämlich die Zuweisungen ($x =$) an und die Verwendungen von ($= x$) Pseudoregistern und den Kontrollfluß. Weiters werden die Aktivitätsbereiche der Pseudoregister (links von den Grundblöcken), die erwartete Ausführungshäufigkeit der Grundblöcke (rechts) und der Konfliktgraph dargestellt.

Das Problem der Registerbelegung läßt sich umformulieren in das Problem, den Knoten Maschinenregister zuzuordnen, sodaß zwei Knoten, die durch eine Kante verbunden sind, verschiedene Register bekommen. Dieses Problem³ ist im Operations Research als Graphenfärben (*graph colouring*) bekannt und wie das Problem der Befehlsanordnung NP-vollständig, daher können wir die optimale Lösung nicht in vertretbarer Zeit finden. Wir greifen wieder auf heuristische Algorithmen zurück, um eine gute Lösung zu finden.

Das Grundprinzip der im Compilerbau verwendeten Algorithmen ist einfach: Nimm einen Knoten und weise ihm ein Register zu, das von den Nachbarn des Knotens noch nicht verwendet wird; wiederhole diesen Schritt für alle Knoten. Folgende Fragen bleiben offen:

Was passiert, wenn für einen Knoten kein Register mehr übrig ist?

Eine weit verbreitete Methode ist es, das entsprechende Pseudoregister in den Speicher auszulagern (*spill*). Viele Maschinenbefehle arbeiten jedoch nur mit Registern und nicht mit Operanden im Speicher. Die Pseudoregister müssen daher kurzfristig in einem Maschinenregister gehalten werden, in der Nähe jener Maschinenbefehle, die das Pseudoregister verwenden. Das Auslagern ersetzt also ein langlebiges Pseudoregister durch viele kurzlebige.

In der Praxis werden Knoten, für die kein Register übrig ist, zunächst einmal ignoriert. Nach dem Ende der Registerzuweisung wird der Auslagerungscode (*spill code*) eingefügt, der Konfliktgraph neu aufgebaut und erneut versucht, eine Registerzuweisung zu finden.

In welcher Reihenfolge werden die Knoten abgearbeitet?

Die Reihenfolge der Abarbeitung bestimmt, welche Pseudoregister ein Maschinenregister bekommen und welche nicht. Sie ist daher sehr wichtig für die Qualität der Registerbelegung.

| Pseudoreg | Kosten |
|-----------|--------|
| A | 4 |
| B | 18 |
| C | 10 |
| D | 14 |
| E | 6 |
| F | 6 |
| G | 3 |

Abbildung 9.10: Die Auslagerungskosten der Pseudoregister aus Abb. 9.9, wenn ein Speicherbefehl einen Zyklus und ein Ladebefehl zwei Zyklen kostet

³Genauer gesagt, das Problem, die minimale Anzahl von Registern zu finden, für die es eine solche Zuordnung gibt.

Ein wichtiges Kriterium ist, wieviel an Laufzeit⁴ es kosten würde, das Pseudoregister auszulagern (siehe Abb. 9.10). Ein weiterer wichtiger Faktor ist der Grad des Knotens, also die Anzahl der Kanten, die von ihm ausgehen: Ein hochgradiger Knoten hat viele Nachbarn, denen er ein Register wegnehmen kann und sollte daher eher ausgelagert werden als ein niedergradiger Knoten. Diese beiden Kriterien können z.B. folgendermaßen zu einer Prioritätsfunktion $p(n)$ kombiniert werden:

$$p(n) = \frac{\text{Kosten}(n)}{\text{Grad}(n)}$$

Die Knoten könnten nun einfach in der Reihenfolge absteigender Priorität abgearbeitet werden. Es gibt aber noch eine bessere Variante: Wenn ein Knoten weniger Nachbarn (also einen kleineren Grad) hat als k , die Anzahl der zur Verfügung stehenden Register, dann bleibt auf jeden Fall ein Register für den Knoten übrig, auch wenn zuerst alle seine Nachbarn ein Register bekommen. Wir entfernen daher zunächst einmal alle Knoten n mit $\text{Grad}(n) < k$ aus unserer Betrachtung und aus dem Konfliktgraphen. Dadurch reduziert sich eventuell der Grad weiterer Knoten unter die k -Grenze. Diese erhalten ebenfalls garantiert ein Register, allerdings unter der Bedingung, daß die früher entfernten Knoten erst nach diesen Knoten ihr Register erhalten. Auf diese Weise kann der Graph solange reduziert werden, bis der Graph leer ist oder nur mehr Knoten mit $\text{Grad}(n) \geq k$ übrigbleiben. Ist der Konfliktgraph leer, dann kann eine Registerbelegung gefunden werden, indem die Knoten in der umgekehrten Reihenfolge des Entfernens abgearbeitet werden.

Bleiben noch Knoten übrig, so wird zunächst die Prioritätsfunktion berechnet und der Knoten mit der niedrigsten Priorität entfernt, damit er dann bei der Abarbeitung so spät wie möglich drankommt. Ergeben sich dadurch wieder Knoten mit $\text{Grad}(n) < k$, werden die zuerst entfernt. Anhand der Prioritätsfunktion wird nun der nächste Knoten bestimmt usw., bis der Graph leer ist. Auch hier wird die Registerbelegung in der umgekehrten Reihenfolge der Entfernung aus dem Graphen vorgenommen. In diesem Fall ist der Erfolg aber nicht garantiert und es müssen unter Umständen Knoten ausgelagert werden.

Welches der verfügbaren Register wird einem Knoten zugewiesen?

Bei Knoten mit $\text{Grad}(n) \geq k$ hängt es von den Registern, die seine Nachbarn zugewiesen bekommen haben, ab, ob ein Register für den Knoten übrigbleibt oder nicht. Man könnte annehmen, daß dadurch die Qualität der Registerbelegung stark beeinflusst wird. Versuche mit ausgeklügelten Registerzuweisungsverfahren haben aber enttäuschende Resultate erbracht. In der Praxis wird daher das erstbeste Register verwendet.

Beispiel

Wir wollen eine Registerbelegung für das Programm in Abb. 9.9 mit drei Registern finden. Der erste Schritt ist, die Knoten zu ordnen. E und G können sofort aus dem

⁴Normalerweise muß man sich hier mit Abschätzungen behelfen, da zur Übersetzungszeit nicht bekannt ist, wie oft die entsprechenden Programmteile zur Laufzeit ausgeführt werden.

Graphen entfernt werden, da sie weniger als drei Nachbarn haben. Dadurch hat F auch weniger als drei Nachbarn und kann ebenfalls entfernt werden. Damit bleiben A, B, C und D mit jeweils drei Nachbarn übrig. Wir entfernen A, da seine Priorität mit $4/3$ am niedrigsten ist. Dann haben B, C, und D nur mehr zwei Nachbarn und können in beliebiger Reihenfolge entfernt werden. Die Reihenfolge der Registerzuweisung ist nun umgekehrt: B, C, D, A, F, E, G. Wir weisen nun B das Register 1 zu, was für A, C und D nur mehr die Register 2 und 3 läßt. Nach der Zuweisung von 2 an C und 3 an D ist für A kein Register mehr übrig. Daher setzen wir mit F fort (verfügbar: 1, 3) und geben ihm Register 1. E und G erhalten beide das Register 2.

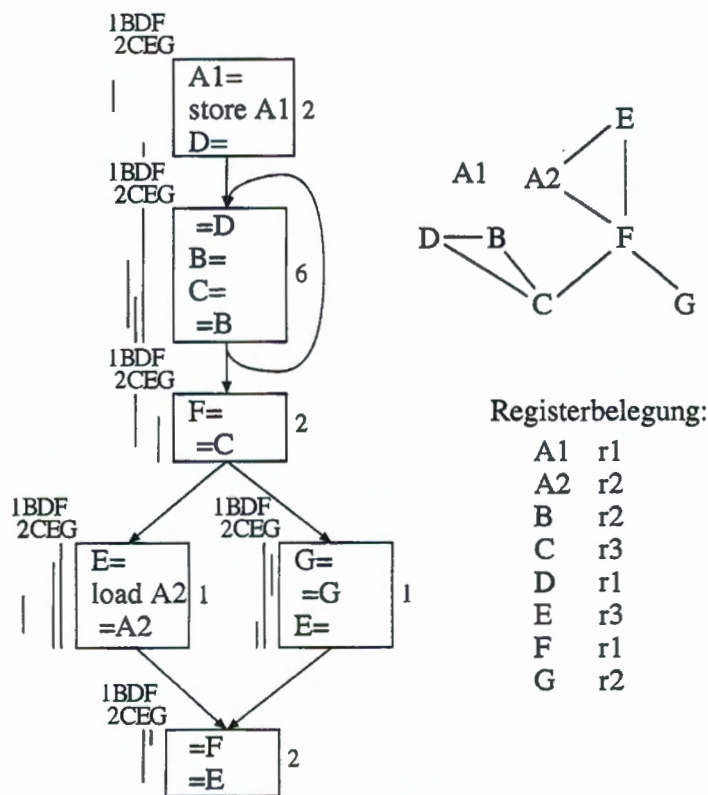


Abbildung 9.11: Auslagerung des Pseudoregisters A, der neue Konfliktgraph, und eine Registerbelegung

Da für das Pseudoregister A kein Maschinenregister gefunden wurde, wird A ausgelagert. Dies erfolgt durch Einfügen eines Speicherbefehls nach der Zuweisung an A und eines Ladebefehls vor der Verwendung von A. Die beiden Aktivitätsbereiche von A werden auf zwei Pseudoregister A1 bzw. A2 aufgeteilt, da sie nicht dasselbe Register erhalten müssen. Dadurch ergibt sich ein neuer Konfliktgraph (siehe Abb. 9.11). Diesmal können bei der Bestimmung der Reihenfolge alle Knoten entfernt werden, ohne einen Knoten mit drei Nachbarn entfernen zu müssen. Daher gibt es natürlich auch eine komplette Registerbelegung.

Entfernung von Kopierbefehlen

Der Konfliktgraph enthält auch alle Informationen, um Kopierbefehle von einem Register in ein anderes wegzuoptimieren. Ein Kopierbefehl von einem Pseudoregister in ein anderes kann wegoptimiert werden, wenn es keinen Konflikt zwischen den beiden Pseudoregistern gibt. Dann werden die beiden Knoten im Konfliktgraphen einfach zu einem verschmolzen (*coalesced*). Dabei werden natürlich auch die Kanten der beiden Knoten vereinigt. Diese Optimierungsphase wird unmittelbar nach dem Aufbau des Konfliktgraphen ausgeführt, noch bevor der erste Versuch der Registerzuteilung unternommen wird.

Es stellt sich natürlich die Frage, ob es genug Kopierbefehle gibt, daß sich das überhaupt auszahlt. Nun, viele Optimierungen lassen sich einfacher formulieren und programmieren, wenn bedenkenlos Kopierbefehle verwendet werden können. Diejenigen davon, die unnötig sind, werden später wieder wegoptimiert.

Damit sich diese Optimierung in mehr Fällen anwenden läßt, kann noch das Konfliktkonzept verfeinert werden: Zwei Pseudoregister können nicht dasselbe Register verwenden, wenn sie zur gleichen Zeit aktiv sind und (während dieser Zeit) verschiedene Werte enthalten. Das heißt, hinter einem Kopierbefehl sind zwei Pseudoregister zunächst einmal nicht in Konflikt, auch wenn sie beide aktiv sind.

Zwei-Adreß-Maschinen

Während der Befehlsauswahl und der Registerbelegung werden alle Befehle so behandelt, als wären die Operanden der Befehle unabhängig voneinander. Doch in sogenannten Zwei-Adreß-Maschinen (z.B. x86, 68000) muß bei den meisten Befehlen das Ergebnis im gleichen Register stehen wie eine Eingabe. Die bisher verwendeten Drei-Adreß-Befehle können folgendermaßen in Zwei-Adreßbefehle umgewandelt werden: Aus

$$r1 = r2+r3$$

wird

$$\begin{aligned} r1 &= r2 \\ r1 &= r1+r3 \end{aligned}$$

Diese Transformation wird zweckmäßigerweise vor dem Entfernen der Kopierbefehle durchgeführt.

Kapitel 10

Laufzeitsystem

Das Laufzeitsystem eines Compilers umfaßt alle Komponenten, die zum Ablauf eines erzeugten Programms notwendig sind (und vom Betriebssystem nicht bereitgestellt werden). Das Laufzeitsystem muß einerseits auf den Compiler und andererseits auf das Betriebssystem abgestimmt werden. Die Konstruktion von Laufzeitsystemen ist ein wesentlicher Teil des Compilerbaus. Dieser Bereich wird derzeit von Compilergeneratoren nicht unterstützt.

Eine der wichtigsten Komponenten des Laufzeitsystems ist die Verwaltung des Speichers für Variablen, Prozeduren und Parameter. Wir behandeln in diesem Kapitel die zwei grundlegenden Methoden der **Speicherverwaltung**: Die **Stack-** und die **Heap-Verwaltung**. Die Speicherverwaltung hängt von den Eigenschaften der Quellsprache ab. Compiler und Speicherverwaltung müssen daher in einem „integrierten Entwurf“ aufeinander abgestimmt werden.

Ein weiterer wichtiger Teil des Laufzeitsystems sind die Ein/Ausgabe-Prozeduren. Diese Prozeduren basieren auf dem verwendeten Betriebssystem. In den meisten Programmiersprachen werden Ein/Ausgabe-Prozeduren wie gewöhnliche Prozeduren behandelt. Daher ist eine besondere Abstimmung auf den Compiler nicht notwendig.

Das Laufzeitsystem ist auch dafür zuständig, daß das Objektprogramm in seinem Ablauf überwacht wird und daß im Fehlerfall Meldungen gegeben werden, die für den Benutzer verständlich sind. In der Praxis trifft man häufig auf Systeme, die in dieser Hinsicht unzulänglich sind.

Ein optimales Fehlerverhalten zur Laufzeit erreicht man durch verstärkten Informationsaustausch zwischen Compiler und Laufzeitsystem. Der Compiler muß dem Laufzeitsystem Informationen verfügbar machen, die dieses für den normalen Ablauf nicht benötigen würde. Dazu zählen z.B. die Bezeichner von Variablen und die Indexgrenzen von Feldern. Es ist daher zweckmäßig, die Symboltabelle an das Laufzeitsystem zu übergeben. Dann kann im Fehlerfall ein „post-mortem dump“ ausgegeben werden, der die aktuellen Werte der Variablen (mit ihren Bezeichnern) enthält.

10.1 Speicherorganisation

Das von einem Compiler erzeugte Maschinenprogramm (Objektprogramm) besteht zur Laufzeit aus folgenden Komponenten:

- Code (Befehle)
- Statische Daten
- Dynamische Daten

Der Code wird vom Compiler generiert. Speicherbedarf und Adressierung der statischen Daten können ebenfalls vom Compiler bestimmt werden. Die Adressen der dynamischen Daten werden erst zur Laufzeit des Objektprogramms bekannt. Daher muß der Speicher für diese Daten vom Laufzeitsystem verwaltet werden. Abb. 10.1 zeigt eine typische Speicheraufteilung zur Laufzeit.

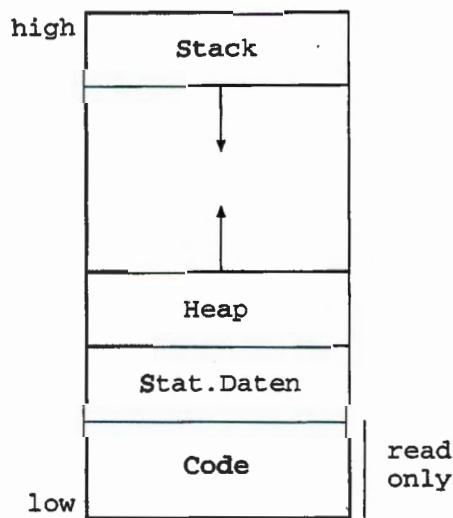


Abbildung 10.1: Speicheraufteilung zur Laufzeit

Man unterscheidet zwei Methoden zur Verwaltung der dynamischen Daten:

Die **Stack-Verwaltung** wird verwendet, um die Datenbereiche von Prozeduren bei geschachtelten und rekursiven Aufrufen zu verwalten. Die Speicherzuteilung ist streng an jeden Prozeduraufruf gekoppelt: Der Datenbereich wird beim Aufruf zugeteilt und bei der Rückkehr wieder freigegeben.

Die **Heap-Verwaltung** ist notwendig, wenn Speicher zu beliebigen Zeitpunkten angefordert und freigegeben werden soll.

Meistens benötigt man sowohl Stack- als auch Heap-Verwaltung. Es ist üblich, für beide einen gemeinsamen Speicherbereich zu reservieren und diesen wie in Abb. 10.1 nach Bedarf dynamisch in Stack und Heap aufzuteilen.

10.2 Stack-Verwaltung

Rekursive Prozeduren benötigen bei jedem Aufruf einen eigenen Speicherbereich, der als **Activation Record (AR)** oder **Frame** bezeichnet wird (Abb. 10.2) .

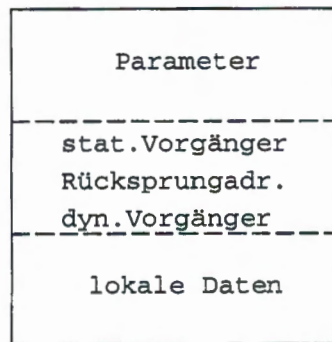


Abbildung 10.2: Activation Record einer Prozedur

Die **Parameter** werden von der aufrufenden Prozedur im AR der aufgerufenen Prozedur abgelegt. Bei einer Funktion wird im Parameterbereich ein Platz für den Funktionswert reserviert – falls die Rückgabe nicht über ein Register erfolgt. Der **statische Vorgänger** (*static link, access link*) ist ein Zeiger auf den AR der statisch umschließenden Prozedur. Er wird nur bei Programmiersprachen mit geschachtelten Prozedurdeklarationen benötigt (siehe Kap. 10.3). Die **Rücksprungadresse** gibt an, wo die aufrufende Prozedur nach Beendigung der aufgerufenen Prozedur fortgesetzt werden soll. Der **dynamische Vorgänger** (*dynamic link, control link*) ist ein Zeiger auf den AR der aufrufenden Prozedur. Er wird verwendet, um nach Beendigung der aufgerufenen Prozedur den Stack zurückzusetzen. Die **lokalen Daten** sind die Werte der lokal deklarierten Variablen der Prozedur. Dazu kommen unter Umständen noch Zwischenergebnisse für Ausdrücke.

Die Stack-Verwaltung beruht darauf, daß beim Aufruf einer Prozedur der zugehörige AR angelegt und beim Verlassen der Prozedur wieder freigegeben wird. Zur Verwaltung des Stack verwenden wir zwei Register (Abb. 10.3):

- S (Stack-Zeiger) zeigt auf die Spitze des Stack.
- A (AR-Zeiger) zeigt auf den aktuellen, zuletzt erzeugten AR.

Alle Komponenten im AR der aktuellen Prozedur können mit $r(A)$ indiziert adressiert werden, wobei die Abstände r zur Übersetzungszeit aus dem Typ und dem Speicherbedarf der Komponenten berechnet werden. Dadurch ist auch der gesamte Speicherbedarf d der lokalen Daten jeder Prozedur bekannt.

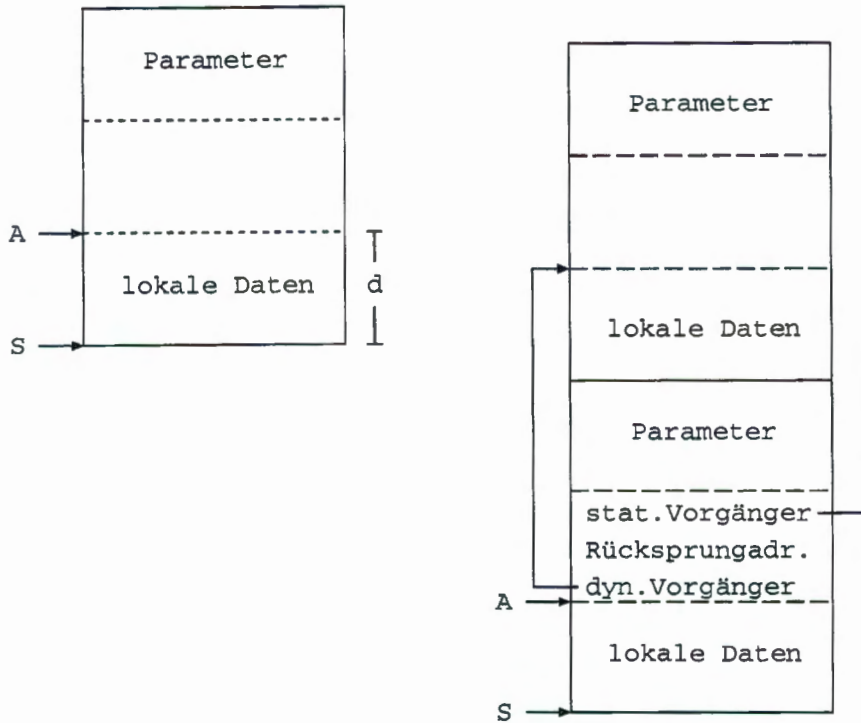


Abbildung 10.3: Prozeduraufruf

Aufrufende Prozedur:

aktuelle Parameter berechnen und kellern
 stat. Vorgänger berechnen und kellern
 Unterprogrammaufruf

stat. Vorgänger entkellern
 aktuelle Parameter entkellern

Aufgerufene Prozedur:

Rücksprungadresse kellern
 dyn. Vorgänger (A) kellern
 $A=S$
 $S=S+d$
 ... Prozedurrumpf ...
 $S=S-d$
 $A=\text{dyn. Vorgänger}$
 dyn. Vorgänger entkellern
 Rücksprungadresse entkellern
 Rücksprung

Abbildung 10.4: Sequenz beim Prozeduraufruf.

Aufruf-Sequenz (calling sequence)

Die Operationen „kellern“ und „entkellern“ in Abb. 10.4 beinhalten entsprechende Veränderungen des Stack-Zeigers S. Im allgemeinen muß auch der Inhalt von Registern beim Prozedur-Aufruf gesichert werden. Sie werden ebenfalls im AR abgespeichert, und zwar entweder von der aufrufenden oder von der aufgerufenen Prozedur.

Falls eine Prozedur mit variabler Parameteranzahl aufgerufen werden soll, muß auch diese Anzahl bekannt sein, damit die aufgerufene Prozedur die Parameterliste abarbeiten kann. Da A auf den Anfang der lokalen Daten und nicht auf den Anfang des AR zeigt, bleiben die Distanzadressen aller anderen Komponenten unabhängig von der Anzahl der Parameter konstant.

Dynamische Arrays werden im Bereich der lokalen Daten jeweils durch einen Zeiger repräsentiert und im Anschluß an den aktuellen AR auf den Stack gelegt.

Abb. 10.5 zeigt den Stack für eine Prozedur mit zwei dynamischen Arrays X und Y. Man sieht, daß die Anfangsadresse des Array Y von der Größe des Array X abhängt und daher vom Compiler nicht als Relativadresse bestimmt werden kann.

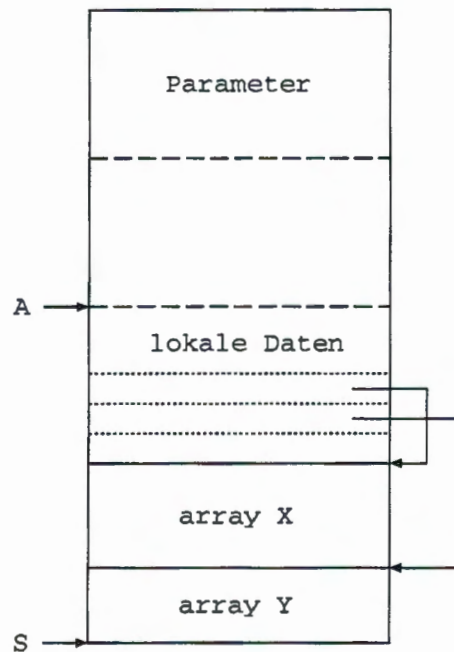


Abbildung 10.5: Dynamische Arrays im Stack

Bei einer Stack-Verwaltung ohne dynamische Arrays kann man den AR-Zeiger einsparen und nur mit dem Stack-Zeiger S adressieren (z.B. MIPS Compiler). Bei „Blatt-Prozeduren“ (*leaf procedures*), die selbst keine Prozeduren aufrufen, kann man das Anlegen eines AR im Stack überhaupt einsparen.

10.3 Nicht-lokale Variablen

Die lokalen Variablen der aktuell aufgerufenen Prozedur werden mit A als Basis adressiert. Bei Programmiersprachen mit geschachtelten Prozedurdeklarationen müssen auch die jeweils gültigen nicht-lokalen Variablen adressiert werden.

Adressierung mit statischen Ketten

Zu diesem Zweck wird in jedem AR ein Zeiger auf den **statischen Vorgänger**, d.h. auf den AR der umschließend deklarierten Prozedur, gespeichert. Ausgehend von A gibt dann die Kette der statischen Vorgänger (**statische Kette**) die auf jeder Schachtelungstiefe adressierbaren Variablen an. Der Anfangszeiger der statischen Kette (hier A) wird auch Umgebungszeiger (*environment pointer*) genannt.

Abb. 10.6 zeigt ein Programm mit geschachtelten Prozedurdeklarationen: Die Schachtelung ist durch Einrückungen dargestellt: Wenn ein Prozedurbezeichner X die Tiefe n hat, so haben die in X deklarierten Objekte und Anweisungen die Tiefe $n+1$. Prozeduraufrufe sind zur Verdeutlichung mit einem Pfeil gekennzeichnet.

Jede Deklaration und jede Referenz besitzt eine eindeutige Schachtelungstiefe. Wenn man P die Tiefe $n=0$ zuordnet, so hat die Deklaration von Q die Tiefe $n=1$ und der Aufruf (Referenz) $\rightarrow Q$ die Tiefe $n=2$. Die Differenz der Schachtelungstiefen zwischen Referenz und Deklaration bezeichnen wir mit dn . Diesen Wert kann der Compiler für jede Referenz bestimmen.

Angenommen, in jedem AR des Stack sei der statische Vorgänger eingetragen und eine Prozedur X soll aufgerufen werden: Dann findet man den statischen Vorgänger von X , indem man die statische Kette von A ausgehend dn -mal dereferenziert.

Als Beispiel betrachten wir Abb. 10.6. Der Stack wird mit P initialisiert (genauer: mit dem AR für P). A zeigt zunächst auf P . Dann ist T aufzurufen und der statische Vorgänger zu bestimmen. Aufruf und Deklaration von T haben die gleiche Tiefe, also ist $dn=0$. Daher zeigt A selbst (als Anfang der statischen Kette) auf den statischen Vorgänger, hier auf P . A wird in T als statischer Vorgänger eingetragen. Danach wird A auf den AR von T gesetzt.

Als nächstes muß Q aufgerufen und der statische Vorgänger bestimmt werden, wobei $dn=1$ ist. Die statische Kette wird ausgehend von A (jetzt auf T zeigend) einmal dereferenziert und liefert wiederum einen Zeiger auf P . Der Zeiger wird in Q eingetragen und A auf Q gesetzt.

Schließlich wird R mit $dn=0$ aufgerufen. A (jetzt auf Q zeigend) zeigt wieder direkt auf den statischen Vorgänger.

Abb. 10.6 zeigt den Inhalt des Stacks, wenn im Programm die Stelle (*) erreicht ist. Ausgehend von A (auf R zeigend) werden über die statische Kette alle momentan adressierbaren Variablen erreicht. Den statischen Gültigkeitsregeln entsprechend sind dies b in Q und a in P . Der AR von T wird von der statischen Kette „umgangen“.

Um nicht-lokale Variablen zu adressieren, wird wiederum die Differenz dn der Schachtelungstiefen zwischen Referenz und Deklaration verwendet. Der Compiler generiert für jeden Zugriff eine dn -malige Dereferenzierung der statischen Kette (siehe Bsp. 10.7).


```

PROGRAM P
  a:REAL
  PROCEDURE Q
    b:REAL
    PROCEDURE R
      c:REAL
      (*)
    END R
    →R
  END Q
  PROCEDURE T
    b:REAL
    →Q
  END T
  →T
END P

```

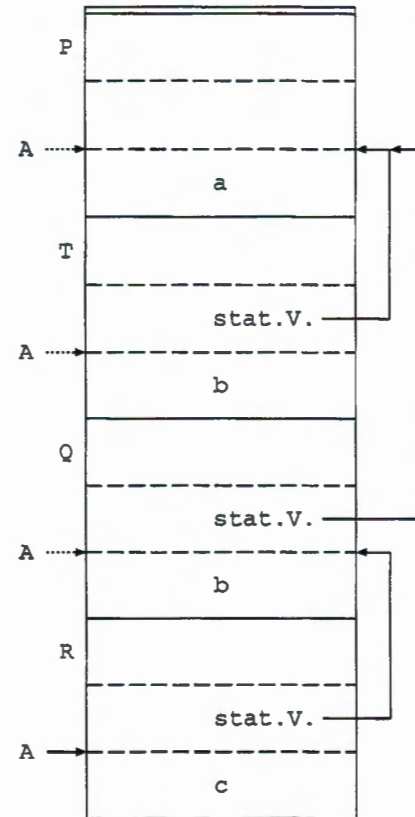


Abbildung 10.6: Modula-2-Programm und Stack mit statischen Ketten

```

Nicht-lokale Variable a:  move sv(A),R
                          move sv(R),R
                          ra(R)
Nicht-lokale Variable b:  move sv(A),R
                          rb(R)
Lokale Variable c:       rc(A)

```

Beispiel 10.7: Adressierung lokaler und nicht-lokaler Variablen. Angenommen, an der Position (*) im Programm von Abb. 10.6 werden die Variablen a, b und c verwendet. Die Relativadressen seien r_a , r_b und r_c . Die Relativadresse der Speicherstelle des statischen Vorgängers sei sv.

Adressierung mittels Display

Der Zugriff auf eine nicht-lokale Variable mittels statischer Kette wird umso aufwendiger, je größer dn ist. Dies kann durch ein **Display** vermieden werden. Ein Display ist ein Block von Indexregistern $D1 \dots Dm$, in dem eine Kopie der statischen Kette verwaltet wird (Abb. 10.8). Die statische Kette enthält zu jedem Zeitpunkt für jede Schachtelungstiefe genau einen Zeiger. Daher kann jeder Schachtelungstiefe ein Indexregister zugeordnet werden, das diesen Zeiger aufnimmt. Die Adressierung von Variablen erfolgt dann ausschließlich aufgrund der Schachtelungstiefe n ihrer *Deklaration* über das Indexregister D_i mit $i=n$ (vgl. Bsp. 10.7):

Nicht-lokale Variable a: $r_a(D1)$
 Nicht-lokale Variable b: $r_b(D2)$
 Lokale Variable c: $r_c(D3)$

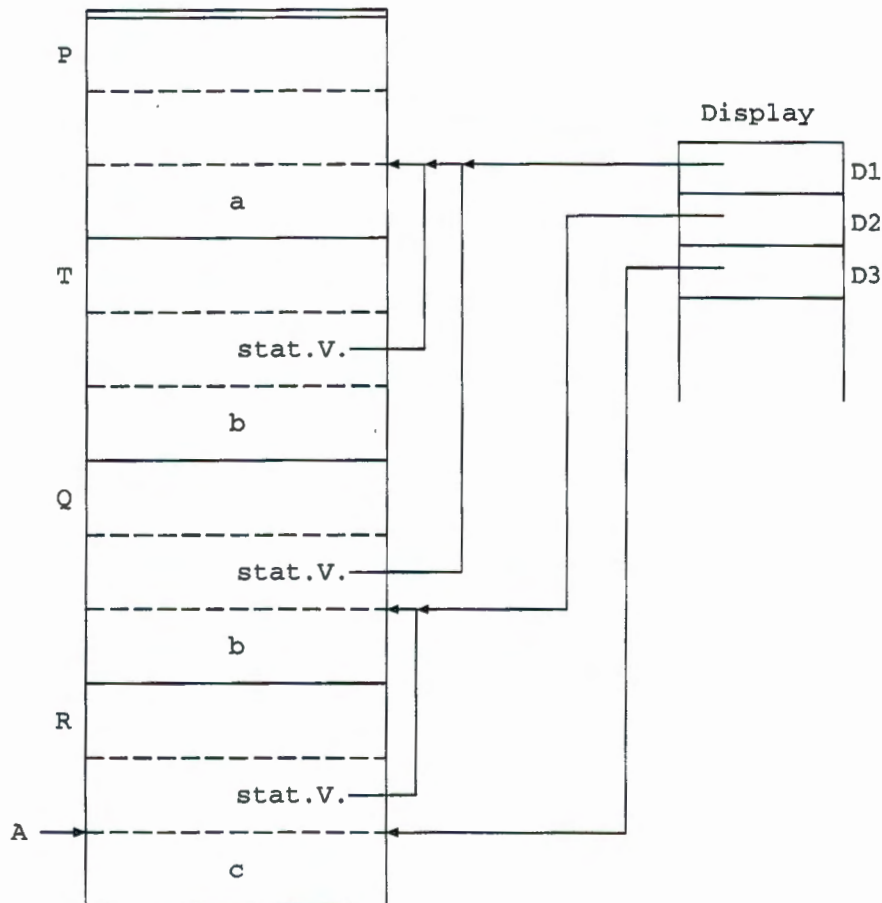


Abbildung 10.8: Stack mit Display für das Programm von Abb. 10.6

Die einfachste Display-Organisation besteht darin, bei jedem Prozeduraufruf alle Zeiger der statischen Kette in das Display zu kopieren (Abb. 10.8). Durch Unterscheiden der Aufrufsituationen läßt sich der Aufwand aber reduzieren. Man sieht

z.B., daß beim Aufruf von Q nur D2, aber nicht D1 geändert werden muß. Beim Aufruf von R ist nur D3 zu laden.

Während bei der statischen Kette die Rückkehr aus einer Prozedur keinen Aufwand verursacht, muß ein Display unter Umständen auch bei der Rückkehr reorganisiert werden, und zwar immer dann, wenn eine Prozedur mit $dn > 0$ aufgerufen wurde. Dies gilt z.B. für jede rekursive Prozedur.

Zusammenfassend ist festzustellen, daß ein Display die Adressierung von Variablen vereinfacht, aber einen höheren Aufwand bei Aufruf und Rückkehr der Prozeduren bedingt.

10.4 Parameterübergabe

Der Bereich am Beginn eines AR ist für die Parameter reserviert. Dieser Bereich liegt im Stack unmittelbar anschließend an den AR der aufrufenden Prozedur und kann daher von der aufrufenden Prozedur gefüllt werden. Die Parameterübergabe besteht generell darin, die Werte der aktuellen Parameter zu berechnen und in den Parameterbereich zu übergeben. In diesem Sinne ist jede Parameterübergabe eine Wert-Übergabe. Allerdings hängt die Art des berechneten und übergebenen Wertes folgendermaßen von der Art des Parameters ab:

Wert-Parameter Der Wert ist das Ergebnis der Auswertung des aktuellen Parameters. Ein formaler Wert-Parameter kann wie eine lokale Variable adressiert werden.

Variablen-Parameter Der Wert ist die absolute Adresse des aktuellen Parameters. Ein formaler Variablen-Parameter muß indirekt adressiert werden.

Prozedur-Parameter Der Wert ist das Tupel $\langle P, \text{statischer Vorgänger} \rangle$. P repräsentiert die Startadresse. Sie wird i.a. bei der Übersetzung festgelegt. Der statische Vorgänger wird bei der Übergabe der Prozedur bestimmt, so als ob die Prozedur an der Übergabestelle aufgerufen würde. Beim Aufruf der formalen Prozedur muß der statische Vorgänger nicht mehr ermittelt werden, da er übergeben wurde.

In Abb. 10.9 ist ein Programm mit einem Prozedur-Parameter dargestellt sowie der Zustand des Stacks, wenn der Ablauf die Stelle (*) erreicht hat.

Für alle Parameterarten gilt folgendes: Wenn ein formaler Parameter selbst wieder als aktueller Parameter verwendet wird, so kann der Wert des formalen Parameters direkt verwendet, d.h. kopiert werden.

10.5 Dynamische Objekte

In vielen prozeduralen Sprachen gibt es dynamische Objekte (Variablen), die zu einem beliebigen Zeitpunkt angelegt (z.B. mit `new`) und wieder freigegeben werden können (explizit oder implizit). Der Speicherbereich für diese Objekte kann nicht in

```

PROGRAM P
  a:REAL
  PROCEDURE T(X)
    b:REAL
    →X
  END T
  PROCEDURE Q
    b:REAL
    PROCEDURE R
      c:REAL
      (*)
    END R
    →T(R)
  END Q
  →Q
END P

```

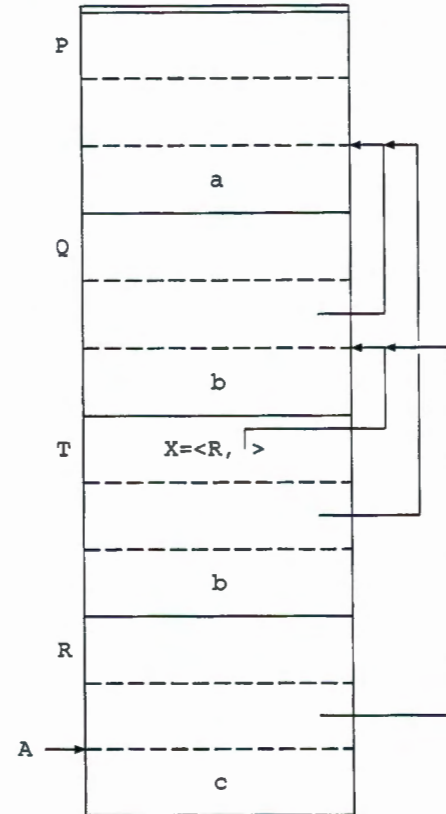


Abbildung 10.9: Modula-2-Programm und Stack mit Prozedur-Parameter

einem Stack verwaltet werden, sondern man benötigt dazu die allgemeinere Organisation eines Heap (siehe Kapitel 10.6).

In objektorientierten Sprachen treten dynamische Objekte als Instanzen von Klassen auf. Im Unterschied zu prozeduralen Sprachen enthalten diese Objekte nicht nur Werte und Referenzen von Variablen, sondern auch Referenzen auf die zugeordneten Methoden.

Abb. 10.10 zeigt die grundlegende Struktur von Objekten einer objektorientierten Sprache mit einfacher Vererbung. Die Unterklasse B erbt von der Oberklasse A die Variable *a* und die Methode *f*. Die Methode *g* wird in B überschrieben. An die deklarierte Variable *v* vom Typ A können Objekte der Klasse A oder B zugewiesen werden. Beim Aufruf *v.g()* bestimmt die Klassenzugehörigkeit des zugewiesenen Objekts, ob die Methode *g* von A oder B aufgerufen wird. Der Compiler kann die Methode statisch nicht bestimmen, da im allgemeinen Zuweisungen vom Ablauf des Programms abhängen. Um die Methode zur Laufzeit auswählen zu können, muß das Objekt selbst Informationen über seine Klassenzugehörigkeit bzw. die zugehörigen Methoden beinhalten. Daher enthält jedes Objekt einen Zeiger auf den Klassen-Deskriptor, in dem die Adressen der Methoden der Klasse stehen. Methoden gleichen Namens haben in den Deskriptoren die gleiche relative Position (z.B. *A.g* und *B.g*),

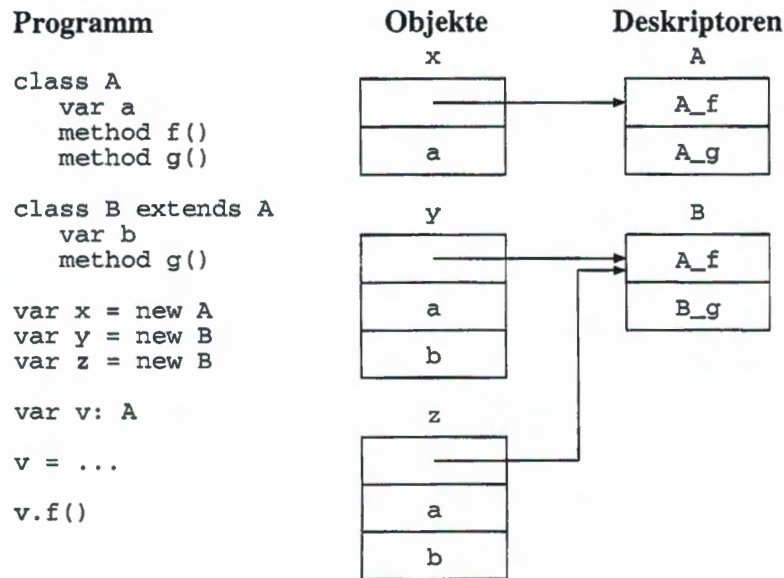


Abbildung 10.10: Grundlegende Struktur von Objekten einer objektorientierten Sprache mit einfacher Vererbung.

die statisch vom Compiler bestimmt werden kann.

In diesem grundlegenden Implementierungsmodell erfordert der Methodenauf-ruf gegenüber einem gewöhnlichen Prozeduraufruf zwei zusätzliche Speicherzugriffe (Adresse des Deskriptors, Adresse der Methode).

10.6 Heap-Verwaltung

In einem Heap werden Speicherobjekte (Zellen) in beliebiger Folge belegt und freigegeben. Heaps sind für Programmiersprachen ebenso relevant wie für Betriebs- bzw. Datenbanksysteme (z.B. Fileverwaltung). Das Laufzeitsystem muß für jede Zelle folgende Operationen bereitstellen (Abb. 10.11):

Allokierung Vor Verwendung einer Zelle muß Speicher für sie allokiert und ggf. initialisiert werden.

Manipulation Während ihrer Lebensdauer müssen Referenzen auf die Zelle manipu-liert werden können.

Deallokierung Nach der letzten Verwendung einer Zelle kann sie deallokiert wer-den – der Speicher ist für eine neue Zelle wiederverwendbar.

Die Wahl einer geeigneten Heap-Verwaltung wird durch Anforderungen der An-wendung sowie Vorgaben der Programmiersprache und des Systems (z.B. virtueller Speicher) bestimmt:

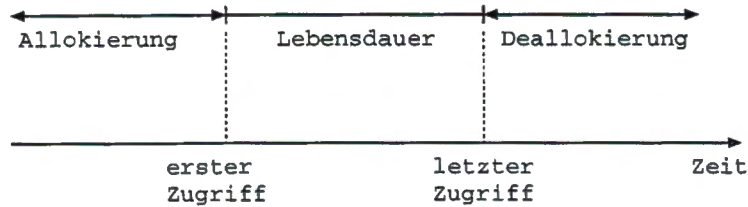


Abbildung 10.11: Lebensdauer einer Zelle

Anwendung

Meistens wird für Programme eine möglichst geringe Gesamtlaufzeit gefordert. Der durch die Gesamtkosten für Allokierung, Manipulation und Deallokierung bestimmte **Durchsatz** der Heap-Verwaltung soll also möglichst groß sein. Bei Programmen, die jedoch eine geringe maximale **Antwortzeit** erfordern (interaktive Systeme, Echtzeitsysteme), werden oft Verfahren mit schlechterem Durchsatz in Kauf genommen.

Programmiersprache

Allokierung/Deallokierung

In den meisten prozeduralen Programmiersprachen können Zellen nur **explizit** durch Anweisungen im Programm angefordert werden (z.B. `NEW` in Modula, `malloc` in C). Der Programmierer ist für die explizite Freigabe verantwortlich (z.B. `DISPOSE`, `free`). Dadurch können *ungültige Zeiger* (*dangling references*) und *Speicherverluste* (*garbage*) auftreten. Zugriffe über ungültige Zeiger auf freigegebene Zellen können das Laufzeitsystem zerstören.

Auch für prozedurale Sprachen ist eine **implizite** Freigabe denkbar, die oft sogar effizienter als eine explizite ist. Lokal in einer Prozedur angelegte Zellen (z.B. `alloca` in C) können bei Rückkehr der Prozedur ohne eigene Heap-Verwaltung implizit freigegeben werden. Falls die Lebensdauer einer Zelle unabhängig von Gültigkeitsbereichen sein soll, muß allerdings eine Heap-Verwaltung verwendet werden, die eine Speicherbereinigung (*garbage collection*) durchführt.

Operationen auf Referenzen

Modula erlaubt nur sehr eingeschränkte Operationen auf Zeiger, während in C Zeigerarithmetik erlaubt ist. Folgende Möglichkeiten sind zu unterscheiden:

- Test auf Gleichheit von Zeigern erlaubt (z.B. in Modula)
- Zeigerrichtung feststellbar (z.B. in C: `if (cp > dp) ...`)
- Zeigerarithmetik erlaubt (z.B. in C: `i = *(ip+5)`)

Falls der Zeigertyp der Programmiersprache zu mächtig ist (z.B. C, C++) oder Typinformationen zur Laufzeit nicht erkannt werden können, kann man einen sicheren „Zeigertyp“ als abstrakten Datentyp realisieren. In C++ können z.B. durch eine Klasse alle mit der Heap-Verwaltung verbundenen Operationen (Initialisierung, Zuweisung, Deallokierung) gekapselt werden.

Referenzstrukturen

Im einfachsten Fall enthalten die Zellen im Heap nur Werte und keine Referenzen auf weitere Zellen. Im allgemeinen aber können die Zellen Referenzen auf weitere Zellen mit Referenzen enthalten. Dabei kann man zwischen azyklischen Strukturen (z.B. Liste, Baum) und zyklischen Strukturen (z.B. Warteschlange) unterscheiden.

Kosten der Heap-Verwaltung

Zur Abschätzung von Antwortzeit und Durchsatz werden die Zellen des Heaps nach ihrer Lebensdauer bezüglich eines Zeitintervalls (t_1, t_2) klassifiziert (Abb. 10.12).

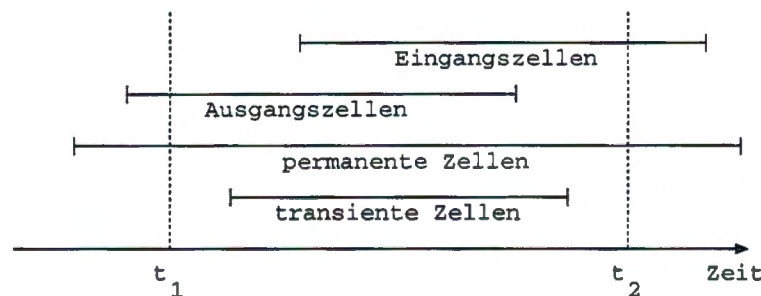


Abbildung 10.12: Klassifizierung von Zellen

Im betrachteten Intervall (t_1, t_2) gibt es N_{ein} Eingangszellen, N_{aus} Ausgangszellen, N_{perm} permanente Zellen und N_{trans} transiente Zellen:

- vor t_1 wurden $N_{perm} + N_{aus}$ Zellen vom Programm benötigt
- innerhalb von (t_1, t_2) wurden $N_{ein} + N_{trans}$ Zellen neu angelegt
- nach t_2 verbleiben $N_{perm} + N_{ein}$ vom Programm benötigte Zellen

Der Aufwand der Verfahren zur Speicherbereinigung ist im allgemeinen linear. Bei einem Verfahren, das z.B. nur von den permanenten und den Eingangszellen abhängt, beträgt der Aufwand $(N_{perm} + N_{ein})$, wobei jeder Summand noch mit einem systemabhängigen Faktor zu versehen ist.

Die Werte N_{ein} , N_{aus} , N_{perm} und N_{trans} sind vom konkreten Programm abhängig: z.B. ist bei File-Verwaltungen N_{perm} sehr groß, bei Verwendung des Heaps für lokale Variablen einer Prozedur ist N_{trans} groß. Man hat meistens mit einer Häufung

von sehr kurzlebigen (transienten) und sehr langlebigen (permanenten) Zellen zu rechnen. Vereinfachend betrachtet man oft nur diese beiden Klassen.

Algorithmen

Heap-Verwaltungen erfüllen im wesentlichen zwei Aufgaben: Zuerst sind die vom Programm noch benötigten (d.h. referenzierten) Zellen festzustellen (bei expliziter Freigabe fällt diese Aufgabe dem Programmierer zu), dann zu reorganisieren. Die entsprechenden Algorithmen sind weitgehend kombinierbar.

Feststellung der Referenzierbarkeit

Markierungsverfahren. Von den adressierbaren Programmvariablen aus werden alle erreichbaren Zellen und ihre Nachfolger markiert. Man benötigt pro Zelle zusätzlich ein Markierungsbit. Bei tief geschachtelten Referenzstrukturen wird der Stack der Markierungsprozedur sehr groß. Deswegen verwendet man oft nichtrekursive Varianten, die ohne Stack auskommen. Der Aufwand beträgt $(N_{perm} + N_{ein})$.

Referenzzähler. Um nicht mehr referenzierte Zellen möglichst früh zu erkennen, müßte man nach jeder Referenzmanipulation ein Markierungsverfahren anwenden. Es ist einfacher, die Veränderungen direkt zu beobachten. Dazu wird pro Zelle ein Zähler verwendet, dessen Wert die Anzahl der Referenzen auf diese Zelle ist. Bei jeder Manipulation einer Referenz auf die Zelle wird ihr Zähler entsprechend erhöht oder erniedrigt. Ist der Zähler = 0, wird die Zelle nicht mehr benötigt. Allerdings werden durch dieses Verfahren zyklische Strukturen nicht erkannt, da die Zähler der Zellen eines nicht referenzierten Zyklus $\neq 0$ bleiben.

Der Aufwand beträgt mindestens $(N_{trans} + N_{aus})$, da pro Zuweisung eines Zeigers zusätzlich der Zähler verändert wird. Dieser Algorithmus kommt daher für Systeme, für die N_{perm} sehr groß (z.B. Files) oder eine geringe maximale Antwortzeit gefordert ist, in Frage.

Werden viele Zellen (in funktionalen Sprachen rund 70%) nur einmal referenziert, kann ein aus nur einem Bit bestehender Referenzzähler verwendet werden, der im *Zeiger* auf die Zelle abgespeichert wird. Beim Überschreiben des Zeigers wird die referenzierte Zelle freigegeben und ihr Speicherplatz kann zugleich für eine neue Zelle verwendet werden. Wird ein Zeiger aber einmal kopiert, entsteht ein Überlauf (Zählerbit=0) und die Zelle kann nur mehr durch andere Verfahren freigegeben werden. Diese Variante wird von einigen Lisp-Maschinen verwendet.

Reorganisation

Reorganisation der freien Zellen. Der Aufwand für diese Verfahren beträgt mindestens $(N_{trans} + N_{aus})$. Die frei gewordenen Zellen werden z.B. über **Freispeicherlisten** verwaltet. Nachteile sind der bei Zellen unterschiedlicher Größe entstehende „Speicherverschnitt“ und der Verlust an Lokalität, weil neue Zellen an einem zufälligen Speicherplatz angelegt werden. Das ist bei virtuellen Speichern besonders ungünstig.

Mit dem Markierungsverfahren kombiniert ergibt diese Freispeicher-Verwaltung im einfachsten Fall den Algorithmus **mark and sweep**:

- Initialisiere alle Markierungsbits
- Markiere alle benötigten Zellen
- Hänge alle unmarkierten Zellen in die Freispeicherliste

Dieses Verfahren ist das älteste (1960 für Lisp) und auf heutigen Systemen (viel realer und virtueller Speicher) das langsamste. Es ist in dieser Form nur noch bei nicht verschiebbaren Zellen (z.B. C) sinnvoll. Geringfügig besser sind inkrementelle bzw. parallelisierte Varianten.

Reorganisation der benötigten Zellen. Die folgenden beiden Klassen von Algorithmen kompaktieren die benötigten Zellen durch Verlegung in einen geschlossenen Bereich. Die Zellen müssen verschiebbar sein. Bei virtuellem Speicher sind Kompaktieralgorithmen aus Gründen der **Lokalität** vorteilhaft - die Arbeitsmenge (*working set*) ist nach jeder Kompaktierung minimal. Der Aufwand hängt nur von den benötigten Zellen ab ($N_{perm} + N_{ein}$).

Ordnungsbewahrende Kompaktierer. Diese verschieben die Zellen unter Bewahrung ihrer Ordnung. Durch sie werden auch Stacks, die redundante Einträge besitzen, kompaktiert. Derartige Stacks treten z.B. in Prolog und Lisp auf.

Der folgende Algorithmus benötigt zu jeder Zelle zusätzlich Platz für einen Hilfszeiger. Weitere Algorithmen benötigen weniger Speicherplatz bzw. nur einen Pass.

1. Verschiebe alle markierten Zellen an ihre neue Position, speichere die neue Adresse in einem Hilfszeiger der alten Zelle ab.
2. Setze alle Zeiger in den kompaktierten Zellen auf die neuen Positionen, die in den Hilfszeigern der alten Zellen stehen.

Ordnungszerstörende Kompaktierer – Kopieralgorithmen. Die folgenden Algorithmen sind einsetzbar, wenn Zellen verschiebbar sind und die Reihenfolge zwischen ihnen keine Rolle spielt. Sie zählen zu den schnellsten Verfahren, da die Feststellung der Referenzierbarkeit und die Reorganisation in einem Pass erfolgen.

Die grundlegende Idee ist, den Heap in zwei Hälften zu teilen: den *Arbeitsspeicher* und den *Freispeicher*. Benötigte Zellen werden im Arbeitsspeicher sequentiell allokiert. Ist dieser voll, werden alle erreichbaren Zellen in den Freispeicher kopiert. Dabei durchwandert man – wie beim Markierungsverfahren – alle erreichbaren Zellen. Jede Zelle wird in den Freispeicher kopiert, an ihrer Stelle im Arbeitsspeicher wird die neue Position der Zelle eingetragen, damit andere auf sie zeigende Zellen ihre Referenzen entsprechend umsetzen können. Danach wird der Freispeicher zum Arbeitsspeicher und umgekehrt.

Die Unterbrechung durch ($N_{perm} + N_{ein}$) Kopiervorgänge bestimmt die maximale Antwortzeit. Zur Reduktion bieten sich zwei kombinierbare Vorgehensweisen an:

- inkrementelles Kopieren
- Unterteilung des Arbeitsspeichers in Segmente (durch Anwender oder System definiert) bzw. Generationen (nach bisherigem Alter einer Zelle), die unterschiedlich oft reorganisiert werden

Kopieralgorithmen werden vor allem für Smalltalk und Lisp verwendet. Einige Prozessoren (z.B. Explorer) unterstützen inkrementelle Kopieralgorithmen.

Kapitel 11

Optimierungen

Optimierungen transformieren ein Programm in ein äquivalentes, effizienteres Programm. Äquivalent bedeutet, daß sich das optimierte Programm bei allen Eingabedaten genauso verhält, wie das nicht optimierte Programm. Um die Korrektheit der Transformationen zu garantieren und die dafür notwendigen Informationen zur Verfügung zu stellen, müssen semantische Programmanalysen durchgeführt werden.

Bei Optimierungen unterscheidet man weiters zwischen maschinenabhängigen und maschinenunabhängigen Optimierungen. Die meisten Optimierungen bringen auf allen Architekturen einen Effizienzgewinn. Diese werden auf der Zwischendarstellung (abstrakter Syntaxbaum oder Quadrupelcode) durchgeführt. Maschinenabhängige Optimierungen werden erst nach oder während der Codeerzeugung durchgeführt und arbeiten meist auf einer abstrakten Darstellung des Maschinencodes (Quadrupelcode). Viele Optimierungen kann man parameterisieren, um Maschinenabhängigkeiten portabel zu implementieren (z.B. Cachegröße oder Registeranzahl). Zum besseren Verständnis sind die Beispiel in C angegeben, manche Optimierungen lassen sich aber nur auf Maschinenbefehlen durchführen.

Optimierungen können entweder lokal (innerhalb eines Grundblocks), global (innerhalb einer Funktion oder Prozedur) oder interprozedural (über das ganz Programm) wirken. Interprozedurale Optimierungen sind wesentlich aufwendiger als lokale Optimierungen, liefern aber auch bessere Ergebnisse. Skalare Optimierungen arbeiten auf einzelnen Werten, Schleifenoptimierungen arbeiten auf Arrays und Codeoptimierungen verbessern die Befehlsabarbeitung.

11.1 Analysen

Die Kontrollflußanalyse bestimmt die Ausführungsreihenfolge der Befehle. Die Datenflußanalyse analysiert den Wertefluß in einem Programm. Die Abhängigkeitsanalyse stellt die Abhängigkeiten zwischen den Befehlen fest. Die Aliasanalyse überprüft, ob zwei Zeiger auf die selbe Speicherzelle zeigen können.

Während der Kontrollflußanalyse wird das Programm in Grundblöcke zerlegt und der Kontrollflußgraph aufgebaut (siehe Abschnitt 8.8). Eine Schleifenanalyse (*loop analysis*) erkennt Schleifen und die Schachtelungstiefe von Schleifen. Weiters wird noch Dominanzinformation (*dominator tree*) berechnet, die angibt, welche beliebi-

gen Grundblöcke immer vor bzw. nach einem bestimmten Grundblock ausgeführt werden. Die Schleifeninformationen werden z.B. dazu genutzt, Schätzwerte für die Ausführungshäufigkeit von Grundblöcken zu berechnen. Alternativ dazu können reale Werte gespeichert werden, die durch Testläufe des Programms berechnet werden (*feedback directed compilation*).

Viele Optimierungen benötigen Informationen, wo Werte definiert und wo sie verwendet werden (*def-use information, reaching definitions*). Diese und ähnliche Informationen lassen sich mittels Datenflußanalyse berechnen. Die iterative Datenflußanalyse durchwandert entlang des Kontrollflußgraphen mehrmals die einzelnen Grundblöcke und Befehle. Dabei werden die Informationen mittels Ein- und Ausgangsgleichungen in Mengen gesammelt, bis sich die Inhalte in diesen Mengen nicht mehr ändern. Eine weitere Analysetechnik ist die abstrakte Interpretation. Anstelle der echten Programmbeefehle werden abstrakte Operationen ausgeführt, die die benötigten Informationen berechnen. Die gesammelten Informationen müssen endlich sein, um eine Terminierung der Analyse sicherzustellen.

Manche Analysen lassen sich einfacher durchführen, wenn für jede Variable maximal eine einzige Zuweisung existiert. Dazu wird das Programm in *static single assignment* (SSA) Form gebracht. Bei jeder Zuweisung an eine Variable wird eine neue Variable durch Indizierung des Variablennames angelegt. Bei der Vereinigung von Programmpfaden ist es dann möglich, daß von der selben Variablen mehrere Kopien mit unterschiedlichem Index existieren. Diese Variablen werden nun durch einen sogenannten Phi-Knoten zusammengefaßt. Ein Programm in SSA-Form kann wieder in ein übliches Programm zurücktransformiert werden, indem diese Phi-Knoten durch Kopierbefehle ersetzt werden. Unnötige Kopierbefehle können in einer weiteren Optimierungsphase eliminiert werden.

| | |
|-----------|-----------------------------------|
| a = | a ₀ = |
| if (cond) | if (cond) |
| a = | a ₁ = |
| else | else |
| a = | a ₂ = |
| | a ₃ = $\phi(a_1, a_2)$ |
| = a | = a ₃ |

Beispiel 11.1: *static single assignment*

Die Datenabhängigkeitsanalyse bestimmt die Datenabhängigkeiten der Befehle. Wird diese Analyse nur innerhalb eines Grundblocks durchgeführt, so ist der resultierende Datenabhängigkeitsgraph ein zyklensfreier gerichteter Graph (siehe Abschnitt 8.8). Für globale Optimierungen müssen auch Schleifen berücksichtigt werden und der resultierende Datenabhängigkeitsgraph ist ein zyklischer Graph. Für manche Optimierungen wird zusätzlich die Länge der Abhängigkeit als Anzahl der Iterationen berechnet. Alle Kontroll- und Datenabhängigkeiten können im *program dependency graph* zusammengefaßt werden.

11.2 Skalare Optimierungen

Die einfachste Optimierung ist die Auswertung von konstanten Ausdrücken (*constant folding*). Diese wird meistens noch mit einer Vereinfachung der Berechnung basierend auf den algebraischen Gesetzen (*algebraic simplification*) kombiniert. Der Ausdruck $(4 - 3) * a$ kann auf die Variable a reduziert werden, da der konstante Ausdruck Eins ergibt und die Multiplikation mit Eins eliminiert werden kann. *Constant propagation* (Konstantenverbreitung) verteilt Konstante über ein Unterprogramm oder über das ganze Programm. Das kann wieder zu neuen Möglichkeiten für *constant folding* führen. Daher werden diese beiden Optimierungen oft mehrmals angewendet. Bei der *copy propagation* werden Kopierbefehle eliminiert, wenn der Originalwert in einer anderen Variablen zur Verfügung steht.

Strength reduction reduziert die Komplexität von arithmetischen Operationen. So kann eine Multiplikation mit Zwei durch eine Addition, eine Multiplikation mit 16 durch eine vierfache Schiebeoperation und eine Division mit einer Konstanten durch eine Multiplikation mit dem Reziprokwert ersetzt werden. Manchmal wird auch die Vereinfachung von Addressausdrücken als *strength reduction* bezeichnet (siehe Abschnitt 11.4).

Eine wichtige Optimierung ist die Vermeidung von Doppelberechnungen. Die Elimination gemeinsamer Teilausdrücke (*common subexpression elimination*) entfernt solche Doppelberechnungen. Oft treten solche Doppelberechnungen nur auf manchen Pfaden durch das Programm auf. Eine Verallgemeinerung dieser Optimierung ist die *partial redundancy elimination*. Das Beispiel 11.2 zeigt, wie durch die Entfernung des Befehls im *else*-Zweig die Doppelberechnung vermieden werden kann.

| | |
|------------------------------------------------------------------|----------------------------------------------|
| <pre> if (cond) a = b; else x = a + c; x = a + c; </pre> | <pre> if (cond) a = b; x = a + c; </pre> |
|------------------------------------------------------------------|----------------------------------------------|

Beispiel 11.2: *partial redundancy elimination*

11.3 Codeoptimierungen

Konstante Ausdrücke in Bedingungen können dazu führen, daß Teile des Programms nicht mehr ausgeführt werden. Diese Teile können dann entfernt werden (*dead code elimination*). Grundblöcke können in beliebiger Reihenfolge angeordnet werden. Daher werden sie üblicherweise so angeordnet, daß die Anzahl der Sprungbefehle minimiert wird. Dabei werden teilweise die Sprungbedingungen invertiert.

Ein Unterprogrammaufruf ist recht aufwendig. Kurze Unterprogramme können direkt an die Aufrufstelle kopiert werden, um die Effizienz zu steigern (*function*

oder *method inlining*). Bei objektorientierten Programmiersprachen muß bei virtuellen Methodenaufrufen sichergestellt sein, daß nur eine einzige Methode aufgerufen werden kann. In eingebetteten Systemen ist Programmspeicher sehr begrenzt. *Procedural abstraction* sucht idente Programmteile und ersetzt sie durch einen Unterprogrammaufruf. *Procedural abstraction* ist damit die inverse Optimierung zu *inlining*. Rekursive Unterprogrammaufrufe können durch Schleifen ersetzt werden, wenn der Unterprogrammaufruf am Ende des Unterprogramms steht (*tail recursion elimination*). Beispiel 11.3 zeigt wie die rekursive Funktion zur Fakultätsberechnung in eine Schleife transformiert werden kann.

| | |
|------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------|
| <pre>int fak(int n) { if (n == 1) return n; return n * fak(n - 1); }</pre> | <pre>int fak(int n) { int f; for (f = 1; n > 1; n--) f = f * n; return f; }</pre> |
|------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------|

Beispiel 11.3: Rekursionseliminierung

Eine der einfachsten Optimierungen ist die *peephole* (Guckloch) Optimierung. Dabei werden mehrere aufeinanderfolgende Befehle auf Optimierungsmöglichkeiten untersucht. Z.B. erzeugt ein einfacher Codegenerator mit einer einfachen Registerzuteilung manchmal Befehlssequenzen, wo ein Register in den Speicher geschrieben wird und direkt danach der Wert wieder in das selbe Register geladen wird. In diesem Fall kann der Ladebefehl eliminiert werden.

11.4 Array- und Schleifenoptimierungen

Der überwiegende Teil der Rechenzeit wird normalerweise in Schleifen verbraucht. Daher zählen Schleifenoptimierungen zu den wichtigsten Optimierungen. Diese werden sowohl für klassische Architekturen wie Mikroprozessoren, als auch für Vektorrechner oder parallele Systeme benötigt. Das Ziel dabei ist, sowohl die Anzahl der ausgeführten Operationen zu reduzieren, als auch die Speicherzugriffe so zu gestalten, daß Caches oder Vektorregister effizient genützt werden können.

Ein typisches Beispiel für eine Schleifenoptimierung ist die *loop invariant code motion*, wo Berechnungen, die sich innerhalb einer Schleife nicht verändern, aus der Schleife hinaus verschoben werden. Besonders häufig treten solche Berechnungen implizit bei Adressberechnungen auf.

Weitere Optimierungen sind eine Reihe von Schleifentransformationen, die geschachtelte Schleifen aufspalten, einzelne Schleifen kombinieren, Schleifen vertauschen oder die Zugriffsreihenfolge auf einzelne Zellen von Arrays verändern. Diese Optimierungen haben teilweise inverse Transformationen und müssen daher abhängig

vom Optimierungsziel eingesetzt werden. Alle diese Optimierungen sind nur anwendbar, wenn sie die meistens zyklischen Datenabhängigkeiten berücksichtigen.

Elimination der Induktionsvariablen

Die Multiplikationen, die bei der Adressberechnung bei Zugriffen auf mehrdimensionale Arrays auftreten, sind aufwendige Operationen. Daher ist eine der wichtigsten Schleifenoptimierungen die Elimination der Induktionsvariablen (*induction variable elimination*) mit gleichzeitiger Vereinfachung der Adressberechnung (*strength reduction*). Im Bsp. 11.4 wird ein zweidimensionales Array mit 0 initialisiert.

```

for (i = 0; i < N; i++)
  for (j = 0; j < M; j++)
    a[i][j] = 0;

```

Beispiel 11.4: Beispielschleife

Für den doppelt indizierten Zugriff muß der erste Index mit der Größe der Dimension multipliziert werden, der zweite Index muß aufaddiert werden und (implizit im C Code) dieser Wert noch mit der Größe eines Arrayelements multipliziert werden, bevor er zur Anfangsadresse des Arrays aufaddiert werden kann (siehe Bsp. 11.5).

| | |
|------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------|
| <pre> for (i = 0; i < N; i++) for (j = 0; j < M; j++) *(a + i * M + j) = 0; </pre> | <pre> aptr = a; eptr = a + N * M; while (aptr < eptr) *aptr++ = 0; </pre> |
|------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------|

Beispiel 11.5: Elimination der Induktionsvariablen

Variablen, die in den einzelnen Schleifeniterationen mit gleichen arithmetischen Berechnungen fortschreitende Werte erhalten, heißen Induktionsvariablen. Im Beispiel 11.5 sind die Indexvariablen i und j die Induktionsvariablen. Diese werden nun entfernt und die komplexe Adressberechnung wird durch eine einfache Addition ersetzt. Die Abbruchbedingung für die Schleife wird so umgeformt, daß diese mit den Adressvariablen durchgeführt werden kann. In diesem Beispiel wird auch noch die doppelt geschachtelte Schleife auf eine einfache umgesetzt und die Matrix als linearer Speicher aufgefaßt.

Loop unrolling

Loop unrolling (Schleifenentrollen) vervielfacht den Schleifenrumpf. Dadurch wird der Schleifenverwaltungsaufwand reduziert. Im allgemeinen Fall muß vor oder nach der Schleife geprüft werden, ob die Iterationszahl durch den Vervielfachungsfaktor

teilbar ist und die überzähligen Iterationen ausgeführt werden. Im Bsp. 11.6 wird angenommen, daß die Anzahl der Iterationen immer ein Vielfaches von 2 ist.

```

assert(N%2 == 0);
for (i = 0; i < N; i++) {
    a[i] = b[i];
}

```

| |
|------------------------------------------------------------------------------------|
| <pre> for (i = 0; i < N; i+=2) { a[i] = b[i]; a[i+1] = b[i+1]; } </pre> |
|------------------------------------------------------------------------------------|

Beispiel 11.6: *loop unrolling*

Vektorisierung

Vektorisierung (*vectorization*) formt Schleifen mit Zugriffen auf die einzelnen Felder eines Arrays auf Vektoroperationen um. Im Bsp. 11.7 wird dabei ein Vektor durch das Paar [Untergrenze:Obergrenze] als Index eines Arrays dargestellt.

```

for (i = 0; i < N; i++) | a[0:N-1] = b[0:N-1] + c[0:N-1];
                        | a[i] = b[i] + c[i];

```

Beispiel 11.7: *vectorization*

Parallelisierung

Parallelisierung (*parallelization, concurrentization*) teilt eine Schleife auf mehrere Prozessoren auf und führt Iterationen parallel durch. Das funktioniert allerdings nur dann, wenn keine Abhängigkeiten zwischen den Iterationen existieren. Bei verteiltem Speicher müssen auch die Daten auf die zu den Prozessoren gehörenden Speicher verteilt werden. Bsp. 11.8 zeigt die Aufteilung einer geschachtelten Schleife auf N Prozessoren, wobei jeder Prozessor eine Iteration der äußeren Schleife ausführt.

```

for (i = 0; i < N; i++)
    for (j = 0; j < M; j++)
        a[i][j] = b[i][j];

```

| |
|----------------------------------------------------------------------------------------------------------------------------------------------|
| <pre> for (j = 0; j < M; j++) // CPU 0 a[0][j] = b[0][j]; ... for (j = 0; j < M; j++) // CPU N-1 a[N-1][j] = b[N-1][j]; </pre> |
|----------------------------------------------------------------------------------------------------------------------------------------------|

Beispiel 11.8: *loop parallelization*

(Loop Distribution, Loop Fission)

Loop distribution, loop fission (Schleifenaufteilung oder Schleifenaufspaltung) teilt eine einzige Schleife auf mehrere einzelne Schleifen auf. Dadurch können entweder Caches besser genutzt werden oder die Vektorisierung erleichtert werden. Bsp. 11.9 zeigt eine geschachtelte Schleife, die auf zwei Schleifen aufgeteilt wurde und leichter vektorisierbar ist.

| | |
|--------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------|
| <pre>for (i = 0; i < N; i++) { a[i] = 0; for (j = 0; j < N; j++) b[i][j] = a[i]; }</pre> | <pre>for (i = 0; i < N; i++) a[i] = 0; for (i = 0; i < N; i++) for (j = 0; j < N; j++) b[i][j] = a[i];</pre> |
|--------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------|

Beispiel 11.9: *loop distribution***Schleifenaustausch**

Schleifenaustausch (*loop interchange*) vertauscht in geschachtelten Schleifen die äußere mit der inneren Schleife. Dadurch werden die Arrays in verschiedenen Reihenfolgen abgearbeitet. Schleifenaustausch wird angewendet, um lineare Speicherzugriffe und damit eine bessere Nutzung von Caches zu ermöglichen. Bsp. 11.10 vertauscht die beiden Schleifen und ermöglicht einen linearen Speicherzugriff.

| | |
|----------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------|
| <pre>for (j = 0; j < N; j++) for (i = 1; i < N; i++) a[i][j] = a[i-1][j] + b[i];</pre> | <pre>for (i = 1; i < N; i++) for (j = 0; j < N; j++) a[i][j] = a[i-1][j] + b[i];</pre> |
|----------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------|

Beispiel 11.10: *loop interchange***Schleifenvereinigung**

Schleifenvereinigung (*loop fusion*) vereinigt mehrere Schleifen zu einer einzigen. Dabei werden die Schleifenmehrkosten (Zähler und Sprungbefehle) reduziert. Im Bsp. 11.11 kann weiters das zweifache Laden des Wertes `a[i]` eliminiert werden.

Strip Mining (Cache Blocking)

Strip mining (Abarbeitung in Streifen) formt eine einfach geschachtelte Schleife in eine doppelt geschachtelte Schleife um. Dabei werden immer kleine Teile des Arrays

| | |
|-------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------|
| <pre>for (i = 0; i < N; i++) a[i] = b[i]; for (i = 0; i < N; i++) c[i] = a[i] + d[i];</pre> | <pre>for (i = 0; i < N; i++) { a[i] = b[i]; c[i] = a[i] + d[i]; }</pre> |
|-------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------|

Beispiel 11.11: *loop fusion*

in Blöcken abgearbeitet. Als Blockgröße (im Bsp. 11.12 ist sie 32) wird dabei die Vektorlänge oder die Cacheblockgröße gewählt.

| | |
|------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------|
| <pre>for (i = 0; i < N*32; i++) { a[i] = b[i] + 1; c[i] = b[i] - 1; }</pre> | <pre>for (j = 0; j < N*32; j+=32) for (i = j; i < j+32; i++) { a[i] = b[i] + 1; c[i] = b[i] - 1; }</pre> |
|------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------|

Beispiel 11.12: *strip mining*

Loop collapsing

Loop collapsing (Schleifen zusammenklappen) transformiert eine doppelt geschachtelte Schleife in eine einfach geschachtelte Schleife. Diese Transformation wird verwendet, um die nutzbare Vektorlänge zu erhöhen. Im Bsp. 11.13 wird dabei ein zweidimensionales Array durch ein eindimensionales Array ersetzt.

| | |
|-------------------------------------------------------------------------------------------|--------------------------------------------------------|
| <pre>for (i = 0; i < N; i++) for (j = 0; j < N; j++) a[i][j] = b[i][j];</pre> | <pre>for (i = 0; i < N*N; i++) a[i] = b[i];</pre> |
|-------------------------------------------------------------------------------------------|--------------------------------------------------------|

Beispiel 11.13: *loop collapsing*

Loop Peeling

Loop peeling (Schleifen schälen) entfernt Iterationen am Anfang oder am Ende einer Schleife. Das ist dann sinnvoll, wenn diese Iterationen Optimierungen wie zum Beispiel eine Vektorisierung behindern. Bsp. 11.14 zeigt eine Schleife, bei der auf eine ringförmige Datenstruktur zugegriffen wird.

| | |
|--------------------------------------------------------------------------------|-------------------------------------------------------------------------|
| <pre> j = N - 1; for (i = 0; i < N; i++) { a[i] = b[j]; j = i; } </pre> | <pre> a[0] = b[N-1]; for (i = 1; i < N; i++) a[i] = b[i-1]; </pre> |
|--------------------------------------------------------------------------------|-------------------------------------------------------------------------|

Beispiel 11.14: *loop peeling*

Software Pipelining

Aktuelle Prozessoren können mehrere Befehle parallel ausführen. Manche Befehle (z.B. Ladebefehle) benötigen mehrere Zyklen. Kurze Schleifen können dann diese Prozessoren nicht ausnützen, da zu wenig Möglichkeiten für eine Befehlsumordnung vorhanden sind. Beim *software pipelining* werden mehrere Iterationen einer Schleife überlappend ausgeführt. Der Prolog startet die erste Schleifeniteration, der Epilog schließt die letzte Iteration ab. Im Beispiel 11.15 gehen wir davon aus, daß die Schleife mindestens einmal durchlaufen wird. Im Prolog wird die Variable *t* gelesen, im Epilog wird *t* gespeichert. In der Schleife kann das Lesen und Speichern gleichzeitig durchgeführt werden.

| | |
|---------------------------------------------------------------------|-----------------------------------------------------------------------------------------|
| <pre> for (i = 0; i < N; i++) { t = a[i]; b[i] = t; } </pre> | <pre> t = a[0]; for (i = 1; i < N; i++) b[i-1] = t; t = a[i]; } b[N-1] = t; </pre> |
|---------------------------------------------------------------------|-----------------------------------------------------------------------------------------|

Beispiel 11.15: *software pipelining*

Kapitel 12

Übersetzung objektorientierter Konzepte

Dieses Kapitel gibt einen Überblick über Übersetzungstechniken, die notwendig sind, um objektorientierte Programmiersprachen wie Java oder C++ zu implementieren. Die wichtigsten Bereiche sind Objekt- und Klassendarstellung, Methodenaufruf und dynamische Typüberprüfung. Weiters wird kurz auf Optimierungen wie Klassenhierarchieanalyse und *method inlining* eingegangen.

12.1 Klassendarstellung und Methodenaufruf

Für einfache Vererbung gibt es für den Methodenaufruf im wesentlichen nur zwei verwendete Techniken: direkter Methodenaufruf, der durch einen Typtest gesichert ist, und indirekter Methodenaufruf mittels einer Methodenzeigertabelle (*virtual function table, vtbl*). Für Mehrfachvererbung gibt es viele verschiedene Techniken, die unterschiedliche Kompromisse eingehen: Einbettung der Superklassen, *trampolines*, komprimierte Methodenzeigertabellen.

Einfachvererbung

Die Felder eines Objekts werden nach einer Speicherzelle, die den Typ eines Objekts beschreibt, hintereinander im Speicher angeordnet (siehe Abb. 12.1). Der Typ des Objekts kann entweder durch eine ganze Zahl oder durch einen Zeiger auf eine Klassenbeschreibungstabelle dargestellt werden. Bei Einfachvererbung werden im Speicher die zusätzlichen Felder der Unterklassen entsprechend der Hierarchie hintereinander nach der Superklasse angeordnet (siehe Abb. 12.1). Der Zugriff auf ein Feld des Objekts benötigt genauso wie der Zugriff auf ein Feld einer Struktur in einer nicht objektorientierten Programmiersprache genau einen Lade- oder Schreibbefehl.

Beim üblichen indirekten Methodenaufruf wird eine Methodenzeigertabelle (*virtual function table, vtbl*) verwendet (siehe Abb. 12.2). Jedes Objekt enthält einen Zeiger auf die Methodenzeigertabelle (*vtblptr*, wird auch als Typinformation verwendet). Die Methodenzeigertabelle einer Unterklasse erweitert die Tabelle der Oberklasse. Falls in der Unterklasse die Implementierung einer Methode der Oberklasse

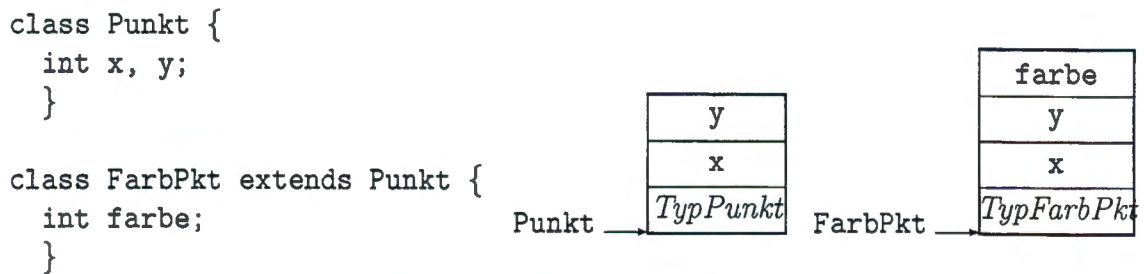


Abbildung 12.1: Objektdarstellung bei Einfachvererbung

geerbt wird, so sind die Methodenzeiger in den Methodenzeigertabelle der Unter- und Oberklasse identisch. Überschriebene Methoden ändern den Eintrag, neu hinzugekommene Methoden erweitern die Methodenzeigertabelle (siehe Abb. 12.2).

```

class Punkt {
  int x, y;
  void schiebe(int x, int y) {...}
  void zeichne() {...}
}

class FarbPkt extends Punkt {
  int farbe;
  void zeichne() {...}
  void faerbe(int f) {...}
}

```

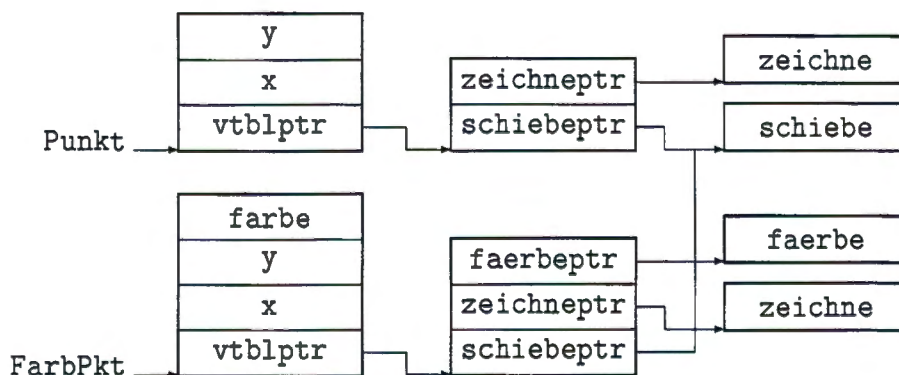


Abbildung 12.2: Methodenzeigertabelle bei Einfachvererbung

Jede Methode hat einen fixen Offset in der Methodenzeigertabelle. Der Methodenaufruf besteht aus drei Maschinenbefehlen (zwei Ladebefehle und ein indirekter Sprung):

```
LDQ vtblptr, (obj)      ; lade vtbl-Zeiger
LDQ mptr, method(vtblptr) ; lade Methoden Zeiger
JSR (mptr)              ; Methodenaufruf
```

Auf modernen Architekturen sind Ladebefehle und indirekte Sprünge teuer. Man kann einen Ladebefehl sparen, wenn der Methodencode direkt in der Methodentabelle steht (fat dispatch tables). Das Problem dabei ist, dass verschiedene Implementierungen einer Methode immer den selben Offset in der Tabelle haben müssen, was zu Kopien von identischem Code, ungenutzten Speicherbereichen und Sprungbefehlen zu Überlaufcode führt.

Direkter Methodenaufruf

Methodenzeigertabellen sind die gängigste Technik für den Methodenaufruf, jedoch sind indirekte Sprünge auf modernen Prozessoren relativ teuer. Der SmallEiffel Compiler verwendet eine Technik, die mit *polymorphic inline caches* verwandt ist. Der Methodenaufruf wird durch eine binäre Suche nach dem Typ der aufgerufenen Methode und einem direkten Methodenaufruf oder dem kopierten Methodencode implementiert (*method inlining*).

SmallEiffel verwendet als Typbezeichner keinen Zeiger (*vtblptr*), sondern eine ganze Zahl, die den Typ des Objekts repräsentiert. Dieser Typbezeichner wird zur Typbestimmung genutzt. SmallEiffel verwendet eine binäre Suche, eine lineare Suche gewichtet nach der Häufigkeit der Verwendung wäre aber auch möglich. Der Maschinencode zur Typbestimmung wird von verschiedenen Aufrufstellen genutzt, wenn die statisch bestimmte Menge der Typen für diese Aufrufstellen gleich sind. Unter der Annahme, dass die Typbezeichner T_A , T_B , T_C und T_D aufsteigend sortiert sind, sieht der Maschinencode für den Aufruf $x.f$ folgendermaßen aus:

```
if  $id_x \leq T_B$  then
  if  $id_x \leq T_A$  then  $f_A(x)$ 
  else  $f_B(x)$ 
else if  $id_x \leq T_C$  then  $f_C(x)$ 
  else  $f_D(x)$ 
```

In Sprachen, die dynamischen Nachladen von Klassen unterstützen, müßte der Maschinencode zur Laufzeit angepaßt werden.

Mehrfachvererbung

Als Stroustrup Mehrfachvererbung für C++ entwarf, schlug er auch verschiedene Implementierungen vor. Eine Erweiterung der Oberklasse wie in Einfachvererbung ist nicht möglich, da die Offsets der Felder bei mehreren Oberklassen nicht konstant gehalten werden können. Die Felder der Oberklasse werden als ein zusammenhängender Block in die Unterklasse eingefügt (siehe Abb. 12.3). Das Einfügen der Oberklasse ermöglicht den gleich schnellen Zugriff auf die Instanzvariablen wie

```

class Punkt {
    int x, y;
}

class Faerbig {
    int farbe;
}

class FarbPkt extends Punkt, Faerbig {}

```

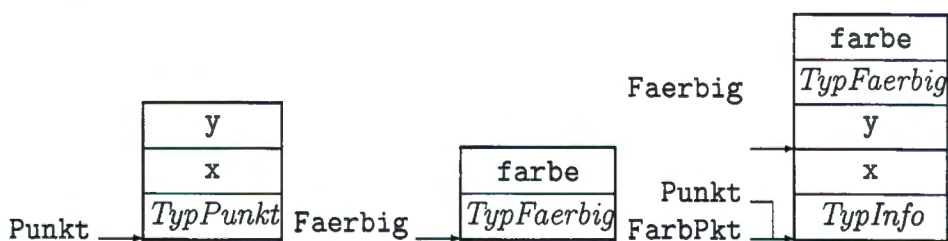


Abbildung 12.3: Objektdarstellung bei Mehrfachvererbung

bei Einfachvererbung. Jedesmal, wenn der Zeigertyp konvertiert wird (explizit oder implizit, Zuweisung, Casts, Parameterübergabe, Ergebnisse), muss der Zeiger korrigiert werden. Die Korrektur des Zeigers muss für Konvertierungen des null-Zeigers unterdrückt werden.

```

Faerbig f;
FarbPkt fp;
f = fp;          // f=fp; if (fp!=null)f=(Faerbig)((int*)fp+2)

```

In C++ versagt die Zeigerkorrektur, wenn Zeigerkonvertierungen ausserhalb der Klassenhierarchie durchgeführt werden (z.B. Konvertierungen nach `void*`). Die Speicherverwaltung wird komplexer, da Zeiger auch in die Mitte von Objekten zeigen können.

Der Methodenaufruf kann mittels Einfügen der Oberklassen erfolgen. Für jede Oberklasse müssen Methodenzeigertabellen angelegt werden und mehrere *vtbl*-Zeiger werden im Objekt gespeichert (siehe Abb. 12.4). Problematisch sind implizite Typkonvertierungen beim Methodenaufruf von der konkreten zur formalen Methode. Die aufrufende Methode kennt nicht den konkreten Typ der aufgerufenen Methode und die aufgerufene Methode kennt nicht den formalen Typ des Aufrufers. Daher kann der Offset für die Korrektur des Zeigers nicht im Programmcode angegeben werden, sondern wird als Korrekturoffset in der Methodenzeigertabelle gespeichert (siehe Abb. 12.4). Auch wenn der Korrekturoffset oft den Wert 0 hat, muss die Addition durchgeführt werden. Der Methodenaufruf benötigt abhängig vom Prozessor 4 bis 5 Maschinenbefehle:

```

class Punkt {
    int x, y;
    void schiebe(int x, int y) {...}
    void zeichne() {...}
}

class Faerbig {
    int farbe;
    void faerbe(int c) {...}
}

class FarbPkt extends Punkt, Faerbig {
    void zeichne() {...}
}

```

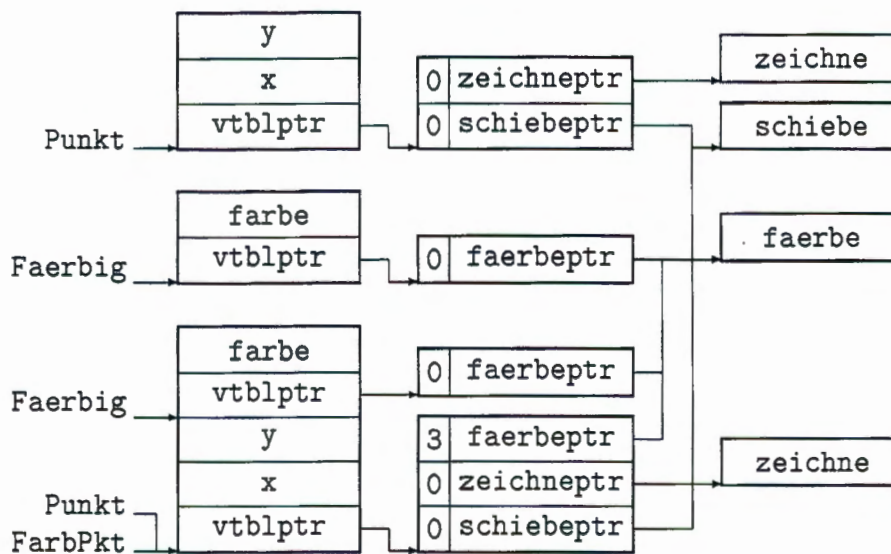


Abbildung 12.4: Methodenzeigertabelle bei Mehrfachvererbung

```

LD vtblptr,(obj)          ; lade vtbl-Zeiger
LD mptr,method_ptr(vtblptr) ; lade Methodenzeiger
LD off,method_off(vtblptr) ; lade Korrekturwert
ADD obj,off,obj          ; korrigiere Objektzeiger
JSR (mptr)               ; rufe Methode auf

```

Dieser Mehrverbrauch an Speicher für Tabellen und Maschinencode ist auch dann notwendig, wenn Mehrfachvererbung nicht genutzt wird. Weitere Anpassungen für die restlichen Parameter oder das Ergebnis sind nicht möglich.

Eine Optimierung ist eine kurze Befehlssequenz (*trampoline*), die über den Zeiger in der Methodentabelle angesprungen wird und die die Korrektur des Offsets

durchführt und dann auf die eigentliche Methode springt. Die Vorteile sind eine kleinere Tabelle (kein Speichern des Offsets nötig) und schneller Methodenaufruf (ident zur Einfachvererbung), wenn Mehrfachvererbung nicht benützt wird. Im Beispiel von Abbildung 12.4 würde `faerbeptr` in der Methodenzeigertabelle von `FarbPkt` auf ein `trampoline` zeigen, das 3 zum Objektzeiger addiert, bevor auf die Methode `faerbe` gesprungen wird:

```
ADD obj,3,obj          ; korrigiere Objektzeiger
BR faerbe             ; rufe Methode auf
```

Wenn Instanzvariablen von einer mehrfach vererbten Oberklasse als eine Variable genutzt werden sollen, dann müssen die Offsets zu diesen Variablen in der Methodenzeigertabelle gespeichert werden. Bei jedem Zugriff auf diese gemeinsam genutzte Variable muss der Offset geladen und addiert werden.

Klassendarstellung und Methodenaufruf in Java

Java und einige andere objektorientierte Programmiersprachen verwenden nicht Mehrfachvererbung, sondern Einfachvererbung mit mehrfachen Untertypen (*multiple subtyping*). Dieser kleine, aber wichtige Unterschied vereinfacht die Klassendarstellung und den Methodenaufruf. Die virtuelle Maschine CACAO legt eine Interfacetabelle mit negativen Offsets in die Methodenzeigertabelle (siehe Abb. 12.5). Für jedes implementierte Interface enthält die Interfacetabelle einen Zeiger auf eine eigene Methodenzeigertabelle für dieses Interface. Nicht implementierte Interfaces werden in der Interfacetabelle durch null-Zeiger dargestellt.

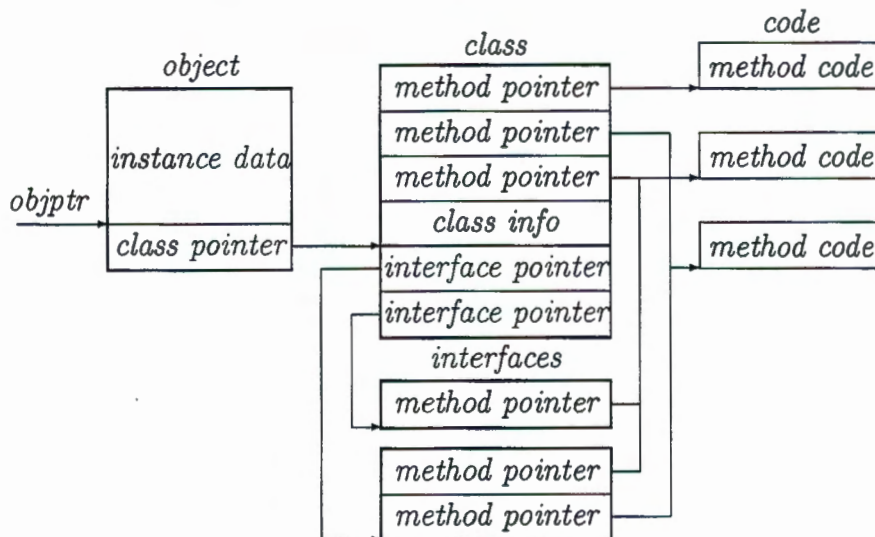


Abbildung 12.5: CACAO Objekt- und Klassendarstellung

Wie bei Einfachvererbung sind zwei Speicherzugriffe (lade Klassenzeiger, lade Methodenzeiger) gefolgt von einem indirekten Sprung nötig, um eine virtuelle Me-

thode aufzurufen. Der Aufruf einer Interfacemethode benötigt eine zusätzliche Indirektion (lade Interface).

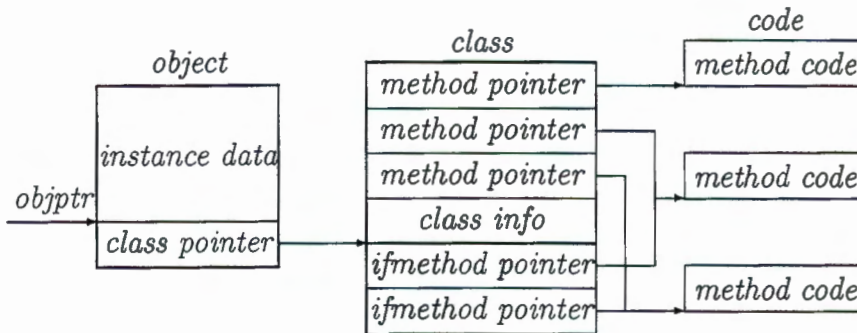


Abbildung 12.6: Komprimierte Objekt- und Klassendarstellung

Die zusätzliche Indirektion kann eingespart werden, wenn alle Methoden, die in irgendeiner Klasse in Interfaces definiert sind, mit negativen Offsets in der Methodenzeigertabelle gespeichert werden (siehe Abbildung 12.6). Der Speicherverbrauch für Methodenzeigertabellen mit Interfacemethoden wäre $Anzahl (Klassen + Interfaces) \times Anzahl\ unterschiedlicher\ Methoden$. Klassen, die keine Interfaces implementieren, benötigen keine Interfacemethodenzeigertabelle. Die Interfacemethoden-tabelle hat sehr viel leere Einträge. Graphenfärben kann verwendet werden, um die Interfacemethodenzeigertabelle zu komprimieren. Dynamisches Nachladen von Klassen kann aber eine Reorganisation der Komprimierung und Änderungen am bereits erzeugten Maschinencode verlangen.

Anstelle einer Komprimierung der Interfacemethodenzeigertabelle verwendet die JikesRVM Hashing. Die Interfacemethodenzeigertabelle wird mit einer kleinen fixen Größe angelegt. Wenn zwei Interfacemethoden den selben Hashwert ergeben, wird anstelle der Interfacemethoden ein kleines Programmstück aufgerufen, das die richtige Methode auswählt. Die Methodenauswahl erfolgt über einen Index, der als Parameter mitgegeben wird. Folgende vier Befehle sind für einen Interfacemethodenaufwurf notwendig:

```
LD  vtblptr, (obj)           ; lade Klassenzeiger
LD  mptr, hash(method_ptr)(vtblptr) ; lade Methodenzeiger
MV  mindex, idreg           ; lade Methodenindex
JSR (mptr)                  ; rufe Methode oder Konfliktstub
```

Die Anzahl der Maschinenbefehle ist gleich wie bei der CACAO Klassendarstellung, allerdings ist die Indirektion eliminiert, die teurer ist, als eine Konstante zu laden.

12.2 Dynamische Typüberprüfung

In einer statisch getypten Sprache kann eine Zuweisung eine Typüberprüfung zur Laufzeit erfordern, um die Korrektheit zu garantieren. Wenn z.B. B eine Unterklasse

von A ist, dann benötigt folgende Zuweisung an die Variable b

```
A a = new B();
...
B b = a;
```

eine Überprüfung, um zu garantieren, dass a einen Wert vom Typ B (oder einer Unterklasse von B) enthält. Weiters gibt es explizite Typüberprüfungen wie z.B. `instanceof` in Java. Diese Überprüfungen können durch statische Analysen nur selten eliminiert werden. Durch die homogene Übersetzung der Generizität in Java werden viele zusätzliche Typüberprüfungen notwendig. Daher sind effiziente Typüberprüfungen sehr wichtig.

Eine naheliegende Implementierung der Typüberprüfung ist das Traversal der Klassenhierarchie um den Obertyp in den Oberklassen zu suchen. Das Traversal ist für Einfachvererbung unkompliziert und für Mehrfachvererbung aufwendiger. Der Nachteil ist, dass der Laufzeitaufwand mit der Größe der Klassenhierarchie wächst. Der Aufwand kann durch Zwischenspeicherung von ein oder zwei Typen beschleunigt werden, besser sind allerdings Verfahren mit konstanter Laufzeit.

So ein Verfahren mit konstanten Kosten und in gleicher Weise für Einfach- und Mehrfachvererbung geeignet ist die binäre Matrix. Die Indizes der Matrix sind die Klassen, ein gesetzter Eintrag in der Matrix gibt eine gültige Untertypbeziehung an. Der Nachteil dieses Verfahrens ist, dass der Speicherverbrauch quadratisch mit der Anzahl der Klassen steigt.

Relative numbering

Relative numbering ist eine bekannte Technik, um in einem Baum zu überprüfen, ob ein Knoten ein Vorgänger eines anderen ist. Dazu wird jeder Knoten mit einem Paar ganzer Zahlen annotiert. Ein Beispiel mit einem annotierten Baum findet sich in Abbildung 12.7.

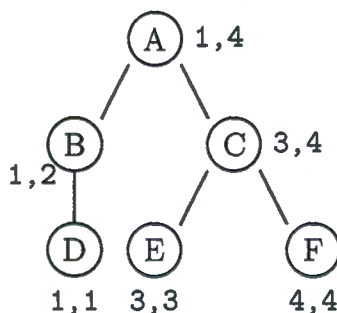


Abbildung 12.7: *Relative numbering* für eine Hierarchie

Um zu überprüfen, ob ein Knoten $n1$ ein Nachfolger von $n2$ (oder gleich $n2$) ist, wird der logische Ausdruck $n1.low \geq n2.low$ and $n1.high \leq n2.high$ aus-

gewertet, wobei *low* die linke bzw. *high* die rechte Zahl bei jedem Knoten entspricht. Dieses Verfahren läßt sich leider nicht einfach für Mehrfachvererbung erweitern.

Display

Bei der Typüberprüfung ist ein *display* ein Vektor von Typbezeichnern aller Oberklassen. Dazu erhalten alle Typen einen eindeutigen Typbezeichner, z.B. den Zeiger auf die Methodentabelle (*vftbl*) oder eine eindeutige Durchnummerierung. In der Datenstruktur, die auch die Methodentabelle enthält, wird ein Typbezeichnervektor (*display* in Abb. 12.8) mit fixer Größe und die Anzahl der Typbezeichner (*depth*, entspricht der Tiefe in der Typhierarchie) gespeichert. Die überzähligen Displayelemente werden mit einem ungültigen Typbezeichner initialisiert (siehe Abb. 12.8).

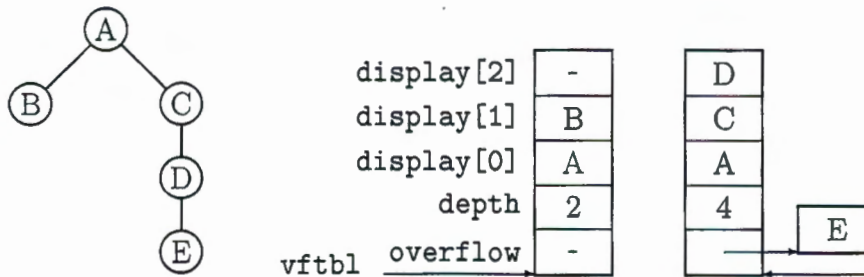


Abbildung 12.8: *Display* für Typüberprüfungen

Die Typüberprüfung ist dann in konstanter Zeit möglich. Das *display* des Untertyps wird mit der Tiefe des Obertyps indiziert und der gespeicherte Typbezeichner wird mit dem Typbezeichner des Obertyps verglichen. Bei den meisten Typüberprüfungen sind die Daten des Obertyps zur Übersetzungszeit bekannt und daher konstante Werte. Die Typüberprüfung besteht dann aus einem Speicherzugriff und einem Vergleich wie folgender C-Code zeigt:

```
return sub->vftbl->display[super->vftbl->depth] == super->vftbl
```

In dynamischen Programmiersprachen wie Java ist die Tiefe der Typhierarchie nicht statisch bekannt. Daher kann das *display* beliebig groß werden. Im Normalfall ist die Tiefe der Typhierarchie klein. Daher werden die ersten Elemente des *display*s in einem fixen Vektor gespeichert und die überzähligen Elemente in einem Überlaufbereich (*overflow*). Auch hier sind meistens zur Übersetzungszeit die Daten des Obertyps bekannt, daher auch ob auf den fixen Vektor oder den Überlaufbereich zugegriffen werden muss. Die Typüberprüfung besteht dann aus einer Bereichsüberprüfung, einer Indirektion, dem Speicherzugriff und einem Vergleich (siehe folgenden C-Code):

```
return sub->vftbl->depth >= super->vftbl->depth &&
       sub->vftbl->overflow[super->vftbl->depth - DISPLAY_SIZE]
       == super->vftbl;
```


Partitionierung der Klassenhierarchie

Da Typüberprüfungen für Bäume platzsparender als Typüberprüfungen für gerichtete Graphen sind, ist eine mögliche Lösung die Partitionierung des Graphen in einen Baumteil und den übrigen Graphen. Für Programmiersprachen mit Einfachvererbung und *multiple subtyping* (Java) wird diese Partitionierung bereits von der Sprache selbst vorgenommen.

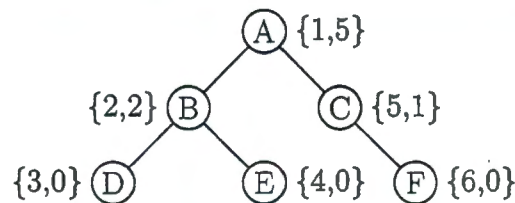


Abbildung 12.9: *Relative numbering* mit $\{\text{baseval}, \text{diffval}\}$ Paaren

Die CACAO virtuelle Maschine verwendet eine Typüberprüfung basierend auf *relative numbering* für Klassen und eine Art binäre Matrix für Interfaces. Zwei Zahlen *low* und *high* sind in der Methodenzeigertabelle für jede Klasse gespeichert. Eine Tiefensuche der Klassenhierarchie erhöht einen Zähler für jede Klasse und weist den Zähler *low* zu, wenn der Knoten das erste Mal beim Abstieg erreicht wird, und weist den Zähler *high* zu, wenn der Knoten beim Aufstieg verlassen wird. Wenn Klassen nachgeladen werden, müssen ab der Stelle, wo die Klasse in der Hierarchie eingefügt wurde, die Werte für *low* und *high* neu gesetzt werden. Eine Klasse *sub* ist ein Untertyp einer anderen Klasse *super*, falls $\text{super.low} \leq \text{sub.low} \leq \text{super.high}$. Da eine Bereichsüberprüfung effizienter durch einen vorzeichenlosen Vergleich implementiert wird, speichert CACAO die Differenz zwischen den *low* und *high* Werten und vergleicht die Differenz der *low* Werte der beiden Klassen (siehe Abb. 12.9). C-Code für `instanceof` sieht folgendermaßen aus:

```
return (unsigned) (sub->vftbl->baseval - super->vftbl->baseval) <=
        (unsigned) (super->vftbl->diffval);
```

Für *final* Klassen (Blattknoten) der Klassenhierarchie ist *diffval* Null, was zu einem simplen Vergleich der *baseval* Werte führt.

CACAO speichert die Interfacetabelle mit negativen Offsets in der Methodenzeigertabelle. Diese Tabelle wird zusätzlich für die Typüberprüfung von Interfaces verwendet. Falls ein Eintrag in der Interfacetabelle leer ist, ist keine Untertypbeziehung vorhanden. C-Code für `instanceof` sieht folgendermaßen aus:

```
return (sub->vftbl->interfacetable[-super->index] != NULL);
```

Beide Typüberprüfungen können in wenigen Maschinenbefehlen ohne teure Sprungbefehle implementiert werden. Aufgrund der Synchronisationsprobleme bei der Neu Nummerierung nach dem Laden einer Klasse wurde in der neuesten Version von CACAO die Typüberprüfung auf *display* umgestellt.

12.3 Devirtualisierung

Devirtualisierung ist eine Technik, um den Aufwand des Aufrufs virtueller Methoden zu reduzieren. Ziel der Devirtualisierung ist die statische Bestimmung aller Methoden, die von einer Aufrufstelle aus aufgerufen werden können. Falls nur eine Methode an einer Aufrufstelle aufgerufen wird, kann der Maschinencode dieser Methode an Stelle des Aufrufs hinkopiert werden (*method inlining*) oder der virtuelle Methodenaufruf durch einen statischen ersetzt werden. Bleiben nur sehr wenige Methoden über, kann eine Typüberprüfung verwendet werden, um *method inlining* zu unterstützen. Die für die Devirtualisierung nötigen Analysen verbessern auch die Genauigkeit des Methodenaufufrgraphens (*call graph*) und damit auch weiterer interprozeduraler Programmanalysen. Abbildung 12.10 zeigt eine Klassenhierarchie und den dazugehörigen Aufrufgraphen zu einem Ausschnitt aus folgendem Java Programm:

```
class A extends Object {
    void m1() {...}
    void m2() {...}
}

class B extends A {
    void m1() {...}
}

class C extends A {
    void m1() {...}
    public static void main(...) {
        A a = new A();
        B b = new B();

        a.m1(); b.m1(); b.m2();
    }
}
```

Schnelle Typanalyse (*rapid type analysis*)

Eine einfache, effiziente und trotzdem relativ genaue Analyse ist die schnelle Typanalyse (*rapid type analysis*, RTA). Die RTA nützt die Tatsache, dass eine Methode m einer Klasse c nur dann aufgerufen werden kann, wenn auch ein Objekt dieser Klasse c angelegt wurde. Der optimistische Algorithmus beginnt die Analyse in der Methode `main` und einer leeren Klassenhierarchie (kein Objekt wurde angelegt) (siehe Algorithmus 12.11). RTA arbeitet den Aufrufgraphen ab und ignoriert anfänglich virtuelle Methodenaufrufe (sie werden nur als potenzieller Aufruf markiert). Wenn dann das erste Mal ein Objekt einer Klasse angelegt wird, werden alle vorher

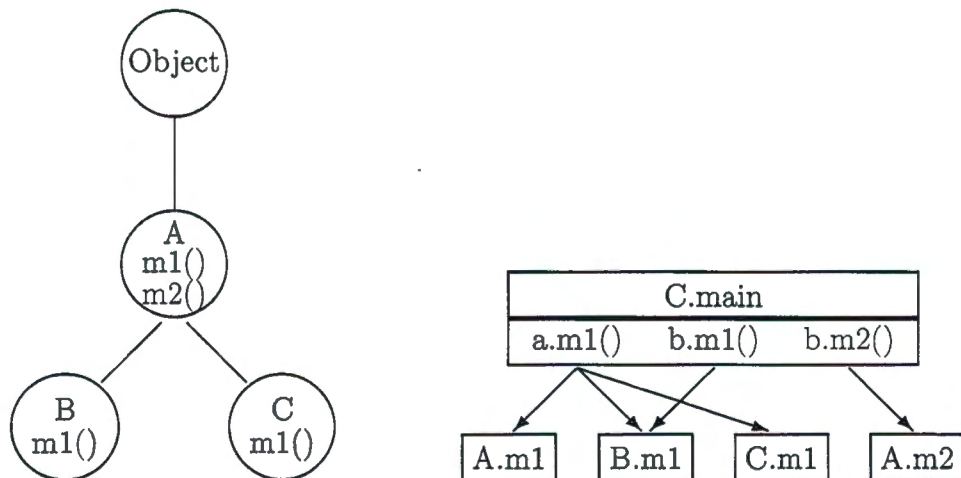


Abbildung 12.10: Klassenhierarchie und Aufrufgraph

markierten Methoden dieser Klasse dem Aufrufgraphen hinzugefügt und anschließend analysiert. Der Aufrufgraph und die Klassenhierarchie wachsen so gemeinsam während der Analyse.

12.4 Escapeanalyse

Üblicherweise werden Objekte dynamisch auf dem Heap angelegt. In Sprachen wie C oder C++ ist der Programmierer für die Speichieranforderung und Speicherfreigabe verantwortlich, in Sprachen wie Java erledigt dies die Heap-Verwaltung.

Das Ziel der Escapeanalyse ist zu erkennen, welche Objekte einen Aktivitätsbereich haben, der einem umschließenden Gültigkeitsbereich (z.B. einer Methode) entkommt (*escape*). Bleibt der Aktivitätsbereich innerhalb des Aufrufs einer Methode, so kann das Objekt auf dem Stack angelegt werden (für C Programmierer entspricht diese Transformation dem Ersetzen von `malloc` Funktionen durch `alloca` Funktionen). Entkommt ein Objekt nicht einem Thread, so können alle Synchronisationen für dieses Objekt entfallen, da kein anderer Thread darauf zugreifen kann.

Algorithmen für Escapeanalyse basieren auf abstrakter Interpretation. Sie unterscheiden sich dabei hauptsächlich im Analyseaufwand und der Genauigkeit der Analyse. Bei Javaprogrammen konnten Effizienzsteigerungen bis zu 44% erreicht werden.

```

main           // main-Methode des Programms
new c          // Erzeugung eines Objects der Klasse c
marked(c)     // Menge der markierten Methoden der Klasse c
m()          // Aufruf der statischen Methode m
type(e)       // der deklarierte Typ des Ausdrucks e
e.m()        // Aufruf der virtuellen Methode m im Ausdruck e
subtype(c)    // c und alle Unterklassen der Klasse c
method(m, c) // Methode m, die in der Klasse c definiert ist
mark(m, c)    // markiert Methode m der Klasse c

callgraph := main
hierarchy := {}
for each m ∈ callgraph do
  for each new c occurring in m do
    if c ∉ hierarchy then
      add c to hierarchy
      for each mmark ∈ marked(c) do
        add mmark to callgraph
  for each mstat() occurring in m do
    if mstat ∉ callgraph then
      add mstat to callgraph
  for each e.mvir() occurring in m do
    for each c ∈ subtype(type(e)) do
      mdef := method(c, mvir)
      if mdef ∉ callgraph then
        if c ∉ hierarchy then
          mark(mdef, c)
        else
          add mdef to callgraph

```

Algorithmus 12.11: *Rapid type analysis*

Literaturverzeichnis

- [AK01] Randy Allen and Ken Kennedy. *Optimizing Compilers for Modern Architectures: A Dependence-based Approach*. Morgan Kaufmann, 2001.
- Dieses Buch geht besonders auf Optimierungen für parallele Architekturen ein.
- [ALSU06] Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman. *Compilers. Principles, Techniques, and Tools*. Addison-Wesley, 2006.
- Dieses neu aktualisierte internationale Standardwerk bietet eine umfassende Einführung in den Übersetzerbau und eine Vertiefung vieler Themenbereiche.
- [App98] A.W. Appel. *Modern Compiler Implementation in Java*. Cambridge University Press, 1998.
- Das Buch beschreibt umfassend und klar die aktuellen Konzepte und Verfahren. Es werden auch objektorientierte Sprachen behandelt.
- [CT04] Keith D. Cooper and Linda Torczon. *Engineering a Compiler*. Elsevier, 2004.
- Die aktuellste und detaillierteste Einführung in den Übersetzerbau die es zur Zeit gibt.
- [FH95] Ch. Fraser and D. Hanson. *A Retargetable C compiler: Design and Implementation*. Benjamin/Cummings Publishing, 1995.
- Das Buch enthält einen kompletten ANSI C Compiler mit Codegeneratoren für Intel, SPARC und MIPS Prozessoren
- [HP03] John L. Hennessy and David A. Patterson. *Computer Architecture. A Quantitative Approach*. Morgan Kaufman Publishers 2003.
- [Kan89] Gerry Kane. *MIPS RISC Architecture*. Prentice Hall, 1989.
- [Kas90] U. Kastens. *Übersetzerbau*. R. Oldenbourg Verlag, 1990.
- Das Buch ist vergleichsweise kurz und vollständig

- [Mot84] Motorola. *M68000, Programmer's Reference Manual*. Prentice Hall, 1984.
- [Par92] Th.W. Parsons. *Introduction to Compiler Construction, Computer Science Press*. W.H.Freeman, 1992.
- Der Verfasser legt besonderen Wert auf Klarheit und weniger auf Tiefe. Das Buch behandelt auch theoretische Grundlagen sowie die Generatoren Lex und Yacc.
- [SF85] A.T. Schreiner and G. Friedman. *Compiler bauen mit UNIX - Eine Einführung*. Carl Hanser Verlag, 1985.
- Mithilfe der Übersetzer-Generatoren Lex und Yacc wird ein Compiler für ein C-Subset erstellt.
- [SS94] L. Sterling and E. Shapiro. *The Art of Prolog*. MIT Press, 1994.
- Das Buch enthält einen kleinen in Prolog geschriebenen Compiler (Kapitel 24).
- [StGDC05] Richard M. Stallman and the GCC Developer Community. *GNU Compiler Collection Internals (GCC)*. Free Software Foundation, 2005.
- Das Manual beschreibt die internen Details eines der besten portablen C-Compiler (Version 4.3.0). Er als freie Software verfügbar.
- [WG92] N. Wirth and J. Gutknecht. *Project Oberon, The Design of an Operating System and Compiler*. Addison-Wesley, 1992.
- Der Oberon-Compiler und die gesamte interaktive Systemumgebung werden beschrieben. Das Buch enthält die vollständigen Programmlisten.
- [Wir86] N. Wirth. *Compilerbau - Eine Einführung*. Teubner, 1986.
- Es wird ein kleiner Compiler-Interpreter für ein Pascal-Subset entwickelt.
- [WM92] R. Wilhelm and D. Maurer. *Übersetzerbau - Theorie, Konstruktion, Generierung*. Springer Verlag, 1992.
- Das Buch enthält theoretische Grundlagen und praktische Verfahren nach dem aktuellen Stand der Technik. Auch funktionale und logische Programmiersprachen werden berücksichtigt.

Anhang A

Mini-Compiler in Java

```
/* ===== MiniVM.java =====
```

The Mini language is a modified small subset of Java/C.
A Mini program consists of a single possibly recursive function.
The language has no declarations (implicit type is integer).

MiniVM.java is a compiler-interpreter for Mini written in Java.
The compiler generates code for a virtual machine, which is a modified
small subset of the Java VM (integer code instead of byte code).

The one-pass compiler is implemented by a top-down recursive descent
parser calling the methods of lexical analysis and code generation.
The parser routines correspond to the grammar rules in EBNF notation.
The regular right parts of EBNF are suitable to postfix generation.
Lexical analysis takes advantage of the Java class StreamTokenizer.

```
===== source language syntax (EBNF) =====
```

```
Program      = Function
Function     = identifier "(" identifier ")" Block
Block        = "{" [Statements] "}"
Statements   = Statement Statements
Statement    = identifier "=" Expression ";" |
              "if" Condition Statement "else" Statement |
              "while" Condition Statement |
              "return" Expression ";" |
              Block |
              ";"
Condition    = "(" Expression ("=="|"!="|">"|"<") Expression ")"
Expression   = Term {"+"|" -"} Term
Term         = Factor {"*"|" /"} Factor
Factor       = number |
              identifier |
              "(" Expression ")" |
              identifier "(" Expression ")"
```

```
===== VM code =====
```

```
0    do nothing
```

```

1 c push constant c onto stack
2 v load variable v onto stack
3 v store stack value into variable v
4 add two top elements of stack, replace by result
5 subtract ...
6 multiply ...
7 divide ...
8 a jump to a if the two top elements of stack are equal
9 a jump if ... not equal
10 a jump if ... less or equal
11 a jump if ... greater or equal
12 a unconditional jump to a
13 a jump to subroutine start address a
14 return from function
15 stop ececution

```

```
===== example =====
```

```
source file "fac.mini":
```

```

-----
fac(n) {
    if (n == 0)
        return 1;
    else
        return n * fac(n-1);
}

```

```
run:
```

```
----
```

```
<Java Compiler> MiniVM.java
```

```
<Java VM> MiniVM fac.mini 8
```

```
VMCode: 13 3 15 2 1 1 0 9 14 1 1 14 12 25 2 1 2 1 1 1 5 13 3 6 14 0
```

```
Result: 40320
```

```
===== */
```

```
////////////////////////////////////
```

```
import java.io.*;
```

```
public class MiniVM
```

```
{
```

```
    static int code_max = 1000;
    static int stack_max = 10000;
```

```
    public static void main(String args[]) {
```

```
        try {
```

```

            Lexer.init(args[0]);           // init lexer
            Gen.init(code_max);           // init generator
            Parser.program();             // call parser
            Gen.show();                   // show VM code
            VM.init(Gen.code,stack_max);  // init VM
            int x = Integer.parseInt(args[1]); // input data

```



```

        int y = VM.exec(x);           // call VM
        System.out.println("Result: "+y); } // print result
    catch (Error e) {
        System.out.println("error "+e.getMessage()); } }
}

```

```

class Error extends Exception
{
    public Error(String msg) {
        super(msg); }
}

```

```

////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////

```

```

interface Token
{
    // token
    final static int T_num = 1; // number
    final static int T_id = 2; // identifier
    final static int T_eql = 3; // ==
    final static int T_neq = 4; // !=
    final static int T_grt = '>';
    final static int T_les = '<';
    final static int T_add = '+';
    final static int T_sub = '-';
    final static int T_mul = '*';
    final static int T_div = '/';
    final static int T_lbr = '(';
    final static int T_rbr = ')';
    final static int T_clb = '{';
    final static int T_crb = '}';
    final static int T_com = ',';
    final static int T_sem = ';';
    final static int T_ass = '=';
    final static int T_eof = '$';

    final static String kw[] = {"if","else","while","return"};
    final static int kw0 = 10;
    final static int T_if = 10;
    final static int T_else = 11;
    final static int T_while = 12;
    final static int T_return = 13;
}

```

```

interface Code
{
    // target code (VM)
    final static int M_nop = 0;
    final static int M_push = 1;
    final static int M_load = 2;
    final static int M_store = 3;
    final static int M_add = 4;
}

```

```

final static int M_sub      = 5;
final static int M_mul      = 6;
final static int M_div      = 7;
final static int M_if_cmpeq = 8;
final static int M_if_cmpne = 9;
final static int M_if_cmple = 10;
final static int M_if_cmpge = 11;
final static int M_goto     = 12;
final static int M_jsr      = 13;
final static int M_return   = 14;
final static int M_stop     = 15;
}

////////////////////////////////////

// symbol table

class Symtab
{
    static String t[] = new String[100];    // table array
    static int n = 0;                       // number of variables

    static int enter(String s) {
        int i;
        for (i=0; i < n && !s.equals(t[i]); i++) {}
        if (i == n) {
            t[i] = s;
            n++; }
        return i; }
}

////////////////////////////////////

// lexical analysis

class Lexer implements Token
{
    static int num_val;    // attribute of number
    static int id_val;    // attribute of id (index of symbol table)

    private static Reader source;
    private static StreamTokenizer st;
    private static int tok;    // tokenizer token

    static void init(String file_name) throws Error {
        try {
            source = new BufferedReader(new FileReader(file_name)); }
        catch (FileNotFoundException e) {
            throw new Error("file not found"+" "+e.getMessage()); }
        st = new StreamTokenizer(source);
        st.ordinaryChar('/');
        st.ordinaryChar('-'); }

    static int scan() throws Error {

```

```

try {
    tok = st.nextToken();
    switch(tok) {
    case StreamTokenizer.TT_EOF:
        return T_eof;
    case StreamTokenizer.TT_NUMBER:
        num_val = (int)st.nval;
        return T_num;
    case StreamTokenizer.TT_WORD:
        int i = look_kw(st.sval);
        if (i >= 0)
            return kw0 + i;
        else {
            id_val = Syntab.enter(st.sval);
            return T_id; }
    default:
        char c = (char)tok;
        switch(c) {
        case '=':
            if ((char)st.nextToken() == '=')
                return T_eq;
            else
                st.pushBack();
            break;
        case '!':
            if ((char)st.nextToken() == '=')
                return T_neq;
            else
                st.pushBack();
            break; } } }
    catch (IOException e) {
        throw new Error ("IO"+" "+e.getMessage()); }
    return tok; }

private static int look_kw(String s) {
    int i;
    for (i=0; i < kw.length && !s.equals(kw[i]); i++) {}
    if (i < kw.length)
        return i;
    else
        return -1; }
}

////////////////////////////////////

// parser (generates VM code)

class Parser implements Token, Code
{
    private static int token;

    static void program() throws Error {
        next();
        Gen.start();
    }
}

```

```
function(); }

private static void function() throws Error {
    match(T_id);
    match(T_lbr);
    match(T_id);
    match(T_rbr);
    block();
    Gen.instr(M_nop); }

private static void block() throws Error {
    match(T_clb);
    statements();
    match(T_crb); }

private static void statements() throws Error {
    if (token != T_crb) {
        statement();
        statements(); } }

private static void statement() throws Error {
    int p1,p2;
    switch(token) {
    case T_id:
        int adr = Lexer.id_val;
        next();
        match(T_ass);
        expression();
        Gen.instr(M_store,adr);
        match(T_sem);
        break;
    case T_if:
        next();
        condition();
        p1 = Gen.pc-1;
        statement();
        Gen.instr(M_goto,0);
        p2 = Gen.pc-1;
        match(T_else);
        Gen.setjump(p1);
        statement();
        Gen.setjump(p2);
        break;
    case T_while:
        next();
        p1 = Gen.pc;
        condition();
        p2 = Gen.pc-1;
        statement();
        Gen.instr(M_goto,p1);
        Gen.setjump(p2);
        break;
    case T_return:
        next();
```



```

        expression();
        Gen.instr(M_return);
        match(T_sem);
        break;
    case T_clb:
        block();
        break;
    case T_sem:
        next();
        break;
    default:
        throw new Error("statement "+token); } }

private static void condition() throws Error {
    match(T_lbr);
    expression();
    int rop = token;
    next();
    expression();
    match(T_rbr);
    switch(rop) {
    case T_eql:
        Gen.instr(M_if_cmpne,0);
        break;
    case T_neq:
        Gen.instr(M_if_cmpeq,0);
        break;
    case T_grt:
        Gen.instr(M_if_cmple,0);
        break;
    case T_les:
        Gen.instr(M_if_cmpge,0);
        break;
    default:
        throw new Error("condition "+token); } }

private static void expression() throws Error {
    term();
    while (token == T_add || token == T_sub) {
        switch (token) {
        case T_add:
            next();
            term();
            Gen.instr(M_add);
            break;
        case T_sub:
            next();
            term();
            Gen.instr(M_sub);
            break; } } }

private static void term() throws Error {
    factor();
    while (token == T_mul || token == T_div) {

```

```

        switch (token) {
        case T_mul:
            next();
            term();
            Gen.instr(M_mul);
            break;
        case T_div:
            next();
            term();
            Gen.instr(M_div);
            break; } } }

private static void factor() throws Error {
    switch(token) {
    case T_num:
        Gen.instr(M_push, Lexer.num_val);
        next();
        break;
    case T_id:
        int id = Lexer.id_val;
        next();
        if (token != T_lbr)
            Gen.instr(M_load, id);
        else {
            next();
            expression();
            match(T_rbr);
            Gen.instr(M_jsr, Gen.start_adr); }
        break;
    case T_lbr:
        next();
        expression();
        match(T_rbr);
        break;
    default:
        throw new Error("expression "+token); } }

private static void next() throws Error {
    token = Lexer.scan(); }

private static void match(int x) throws Error {
    if (token == x)
        next();
    else
        throw new Error("syntax "+token); }
}

////////////////////////////////////

// code generator

class Gen implements Code
{
    static int code[];    // target code

```

```

static int pc;          // program counter
static int start_adr; // start address

static void init(int code_max) {
    code = new int[code_max];
    pc = 0; }

static void start() {
    instr(M_jsr,3);
    instr(M_stop);
    start_adr = pc; }

static void instr(int operator) {
    code[pc] = operator;
    pc = pc+1; }

static void instr(int operator,int operand) {
    code[pc] = operator;
    code[pc+1] = operand;
    pc = pc+2; }

static void setjump(int adr) {
    code[adr] = pc; }

static void show() {
    System.out.print("VMCode: ");
    for (int i=0;i < pc;i++)
        System.out.print(code[i]+" ");
    System.out.println(); }
}

////////////////////////////////////

// virtual machine

class VM implements Code
{
    private static int p[]; // program code
    private static int ip; // instruction pointer
    private static int s[]; // stack
    private static int sp; // stack pointer
    private static int fp; // frame pointer
    private static int fs; // frame size

    static void init(int code[],int stack_max) {
        p = code;
        ip = 0;
        s = new int[stack_max];
        sp = 0;
        fp = 0;
        fs = Symtab.n; }

    static int exec(int arg) throws Error {
        s[0] = arg;

```

```
sp++;
while (p[ip] != M_stop) {
    switch(p[ip]) {
        case M_push:
            s[sp] = p[ip+1];
            sp++;
            ip = ip+2;
            break;
        case M_load:
            s[sp] = s[fp+p[ip+1]];
            sp++;
            ip = ip+2;
            break;
        case M_store:
            s[fp+p[ip+1]] = s[sp-1];
            sp--;
            ip = ip+2;
            break;
        case M_add:
            sp--;
            s[sp-1] = s[sp-1] + s[sp];
            ip++;
            break;
        case M_sub:
            sp--;
            s[sp-1] = s[sp-1] - s[sp];
            ip++;
            break;
        case M_mul:
            sp--;
            s[sp-1] = s[sp-1] * s[sp];
            ip++;
            break;
        case M_div:
            sp--;
            s[sp-1] = s[sp-1] / s[sp];
            ip++;
            break;
        case M_if_cmpeq:
            sp = sp-2;
            ip = s[sp] == s[sp+1] ? p[ip+1]:ip+2;
            break;
        case M_if_cmpne:
            sp = sp-2;
            ip = s[sp] != s[sp+1] ? p[ip+1]:ip+2;
            break;
        case M_if_cmple:
            sp = sp-2;
            ip = s[sp] <= s[sp+1] ? p[ip+1]:ip+2;
            break;
        case M_if_cmpge:
            sp = sp-2;
            ip = s[sp] >= s[sp+1] ? p[ip+1]:ip+2;
            break;
    }
}
```



```
case M_goto:
    ip = p[ip+1];
    break;
case M_jsr:
    s[sp] = ip+2;        // save return address
    s[sp+1] = fp;       // save fp
    fp = sp+2;         // set fp
    sp = fp+fs;        // set sp
    s[fp+1] = s[fp-3]; // copy argument
    ip = p[ip+1];      // goto start address
    break;
case M_return:
    s[fp-3] = s[sp-1]; // copy return value
    sp = fp-2;         // reset sp
    fp = s[sp+1];      // reset fp
    ip = s[sp];        // goto return address
    break;
default:
    throw new Error("illegal vm code "+p[ip]); } }
return s[0]; }
```

```
////////////////////////////////////
```

Anhang B

Mini-Compiler in Prolog

```
/* ===== mini.pl =====
```

The Mini language is a modified small subset of Java/C.
A Mini program consists of a single possibly recursive function.
The language has no declarations (implicit type integer).

mini.pl is a compiler-interpreter for Mini written in Prolog.
The compiler generates a syntax tree which is executed by a recursive interpreter.

The compiler has two passes (lexical analysis and syntax analysis).
Both are implemented with definite clause grammars (DCG).
The target language (syntax tree) is specified by Prolog predicates.

```
===== source language syntax (DCG) =====
```

```
program    --> function
function   --> identifier,"(",identifier,")",block.
block      --> "{",statements,"}".
statements --> statement,statements.
statements --> "".

statement  --> identifier,"=",expression,";".
statement  --> "if",condition,statement,"else",statement.
statement  --> "while",condition,statement.
statement  --> "return",expression,";".
statement  --> block.
statement  --> ";".

condition  --> "(",expression,("==" ; "!=" ; ">" ; "<"),expression,")".

expression --> term,exp_rest.
exp_rest   --> ("+" ; "-"),term,exp_rest.
exp_rest   --> "".

term       --> factor,term_rest.
term_rest  --> ("*" ; "/" ),factor,term_rest.
term_rest  --> "".
```

```

factor    --> number.
factor    --> identifier.
factor    --> identifier,"(",expression,")".
factor    --> "(",expression,")".

===== target language ===== */

% syntax tree

is_int_prog(p(F)) :-
    is_fun(F).

is_fun(fun(Id,Pm,S)):-
    is_id(Id),
    is_id(Pm),
    is_stm(S).

is_stm(seq(S1,S2)) :-
    is_stm(S1),
    is_stm(S2).
is_stm(asm(Id,E)) :-
    is_id(Id),
    is_exp(E).
is_stm(ifs(E,S1,S2)) :-
    is_exp(E),
    is_stm(S1),
    is_stm(S2).
is_stm(whs(E,S)) :-
    is_exp(E),
    is_stm(S).
is_stm(ret(E)) :-
    is_exp(E).
is_stm(emp).

is_exp(con(N)) :-
    is_con(N).
is_exp(var(Id)) :-
    is_id(Id).
is_exp(fct(Id,E)) :-
    is_id(Id),
    is_exp(E).
is_exp(bin(Op,E1,E2)) :-
    is_opr(Op),
    is_exp(E1),
    is_exp(E2).

is_id(Id):- atom(Id).
is_con(N):- number(N).
is_opr(Op):- member(Op,[=,#,>,<,+,-,*,/]).

/* ===== example =====

source file "fac.mini":
-----

```



```

lex_analysis(SourceFile,Tokens) :-
    reader(SourceFile,Chars),
    phrase(tokens(Tokens),Chars).

tokens(Ts)    --> whitespace,
               tokens(Ts).
tokens([T|Ts]) --> tok(T),
                  !, % single solution:longest input match
                  tokens(Ts).
tokens([])    --> "".

tok('if')    --> "if".
tok('else')  --> "else".
tok('while') --> "while".
tok('return') --> "return".
tok('{')     --> "{".
tok('}')     --> "}".
tok(';')     --> ";".
tok(',')     --> ",".
tok('(')     --> "(".
tok(')')     --> ")".
tok(rop(=))  --> "==".
tok(rop(#))  --> "!=".
tok(rop(>))  --> ">".
tok(rop(<))  --> "<".
tok(aop(+))  --> "+".
tok(aop(-))  --> "-".
tok(mop(*))  --> "*".
tok(mop(/))  --> "/".
tok('=')     --> "=".
tok(id(I))   --> ident(Cs),{name(I,Cs)}.
tok(num(N))  --> number(Cs),{name(N,Cs)}.

ident([C|Cs]) --> letter(C),identr(Cs).
identr([C|Cs]) --> letter(C),identr(Cs).
identr([C|Cs]) --> digit(C),identr(Cs).
identr([])    --> "".

number([C|Cs]) --> digit(C),number(Cs).
number([C])    --> digit(C).

letter(C)     --> [C],{0'A=<C,C=<0'Z;0'a=<C,C=<0'z}.
digit(C)      --> [C],{0'0=<C,C=<0'9}.
whitespace    --> [C],{C=<0' }.

% -----
reader(F,P) :-
    open(F,read,S),
    read_chars(S,P),
    close(S),
    !.

```

```

read_chars(S, []) :-
    at_end_of_stream(S).
read_chars(S, [C|Cs]) :-
    \+at_end_of_stream(S),
    get0(S,C),
    read_chars(S,Cs).

%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%

% syntax analysis

syn_analysis(Tokens,SynTree) :-
    phrase(prog(SynTree),Tokens).

prog(p(F))          --> func(F).

func(fun(Id,Pm,S))  --> [id(Id)],['('],[id(Pm)],[')'],block(S).

block(S)            --> ['{'],stms(S),['}'].

stms(S)             --> stm(S1),stmr(S1,S).
stms(emp)           --> [].
stmr(S1,seq(S1,S)) --> stm(S2),stmr(S2,S).
stmr(S,S)           --> [].

stm(asm(Id,E))      --> [id(Id)],['='],exp(E),[';'].
stm(ifs(E,S1,S2))   --> ['if'],cond(E),stm(S1),['else'],stm(S2).
stm(whs(E,S))       --> ['while'],cond(E),stm(S).
stm(ret(E))         --> ['return'],exp(E),[';'].
stm(S)              --> block(S).
stm(emp)            --> [';'].

cond(bin(Op,E1,E2)) --> ['('],exp(E1),[rop(Op)],exp(E2),[')'].

exp(E)              --> term(E1),expr(E1,E).
expr(E1,E)          --> [aop(Op)],term(E2),expr(bin(Op,E1,E2),E).
expr(E,E)           --> [].

term(E)             --> factor(E1),termr(E1,E).
termr(E1,E)         --> [mop(Op)],factor(E2),termr(bin(Op,E1,E2),E).
termr(E,E)          --> [].

factor(con(N))      --> [num(N)].
factor(var(Id))     --> [id(Id)].
factor(fct(Id,E))   --> [id(Id)],['('],exp(E),[')'].
factor(E)           --> ['('],exp(E),[')'].

%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%

% execute target program (syntax tree)
% Env..Environment, V..Value
% return statements produce a return value replacing the environment

execution(p(F),Vin,Vout) :-

```

```

exec(F,Vin,Vout).

exec(fun(Id,Pm,S),Vin,Vout) :-
    new_env(fun(Id,Pm,S),Env0),
    put_env(Pm,Vin,Env0,Env1),
    exec(S,Env1,Vout),
    number(Vout). % fail if no return value

exec(seq(S1,S2),Env0,Env) :-
    exec(S1,Env0,Env1),
    (number(Env1) -> % return value ?
     Env = Env1 % don't execute S2
    );
    exec(S2,Env1,Env).

exec(asm(Id,E),Env0,Env) :-
    exec(E,Env0,V),
    put_env(Id,V,Env0,Env).

exec(ifs(E,S1,_),Env0,Env) :-
    exec(E,Env0,1),
    exec(S1,Env0,Env).

exec(ifs(E,_S2),Env0,Env) :-
    exec(E,Env0,0),
    exec(S2,Env0,Env).

exec(whs(E,S),Env0,Env) :-
    exec(E,Env0,1),
    exec(S,Env0,Env1),
    exec(whs(E,S),Env1,Env).

exec(whs(E,_),Env,Env) :-
    exec(E,Env,0).

exec(ret(E),Env0,V) :-
    exec(E,Env0,V).

exec(emp,Env,Env).

exec(con(V),_,V).

exec(var(Id),Env,V) :-
    get_env(Id,Env,V).

exec(fct(_,E),Env,Vout) :-
    exec(E,Env,Vin),
    Env = env(F,_),
    exec(F,Vin,Vout).

exec(bin(Op,E1,E2),Env,V) :-
    exec(E1,Env,V1),
    exec(E2,Env,V2),
    exec(Op,V1,V2,V).

exec(+,V1,V2,V) :- V is V1+V2.
exec(-,V1,V2,V) :- V is V1-V2.
exec(*,V1,V2,V) :- V is V1*V2.
exec(/,V1,V2,V) :- V is V1/V2.
exec(=,V1,V2,V) :- V1=V2 -> V=1.;V=0.
exec(≠,V1,V2,V) :- V1=V2 -> V=0.;V=1.
exec(>,V1,V2,V) :- V1>V2 -> V=1.;V=0.
exec(<,V1,V2,V) :- V1<V2 -> V=1.;V=0.

```