

Compilerbau

②

Literatur

- 1) Aho, A.V., Sethi, R., Ullman, J.D. (1988).
Compilerbau. Teil 1 und 2. ("Drachenbuch")
Addison-Wesley
- 2) Appel, A.W. (2002).
Modern Compiler Implementation in Java.
2nd Edition. Cambridge University Press
- 3) Grune, D., Bal, H.E., Jacobs, J.H., Langendoen, K.G.
(2001). Modern Compiler Design.
John Wiley & Sons
- 4) Jäger, M. (2011).
Compilerbau - eine Einführung.
Skript zur Vorlesung an der THM,
FB MNI. Im WWW

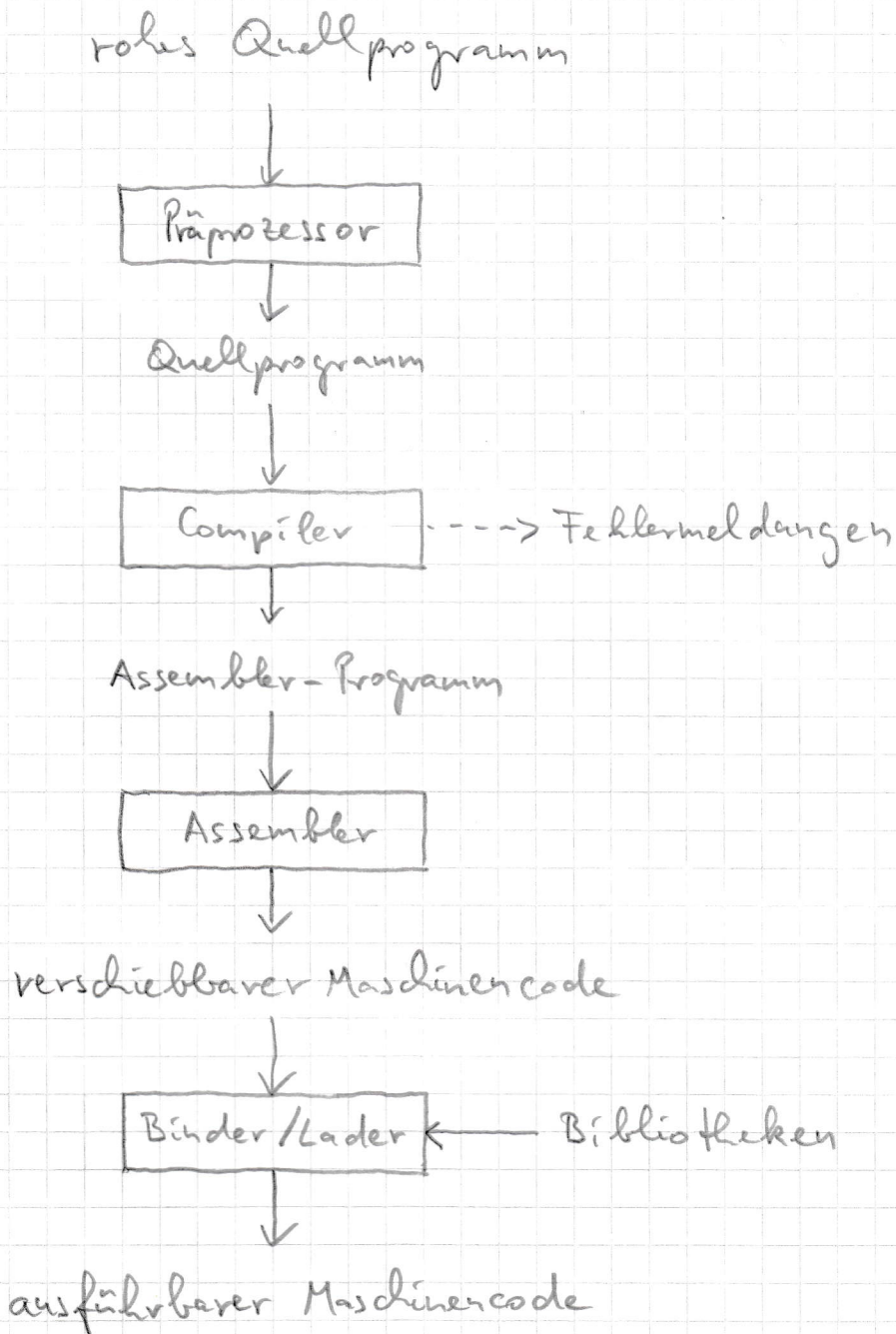
Überblick

③

- 1.5 1. Einführung = Aufgaben, Phasen, Läufe, Schnittstellen
- 1.5 2. Lexikalische Analyse = Reguläre Sprachen, endliche Automaten, Scanner-Generatoren
- 3 3. Syntaktische Analyse = Kontextfreie Sprachen, Abbildungen, LL- und LR-Parser, Parser-Generatoren
- 1 4. Abstrakte Syntax und Attributierung
- 2 5. Semantische Analyse = Typen und Typkonstrukturen, Typüberprüfung, Symboltabelle
- 1 6. Laufzeit-Organisation = Stack und Stack-Frames, Prozedur-Eintritt/Austritt, Parameter-Übergabe
- 2 7. Codegenerierung = Stackmaschinen, Registermaschinen, Übersetzung von Ausdrücken, Übersetzung von Kontrollanweisungen, Übersetzung von Prozeduranrufen
- 1 [8. Registerallokation]

1. Einführung

1.1. Der Compiler in der Werkzeugkette



1.2. Phasen, Läufe, Schnittstellen

Verteilung der Aufgaben auf Phasen.

Zusammenfassung von Phasen zu Läufen

(ein Lauf $\hat{=}$ ein Durchgang durch das Programm).

Interne Darstellung des Programms durch Schnittstellen.

Quellprogramm -
characters

Lexikalische
Analyse

Tokens

Syntaktische
Analyse

Abstrakte Syntax

Semantische
Analyse

Symbol-
tabellen

Attributierte Syntax

Zwischencode-
Erzeugung

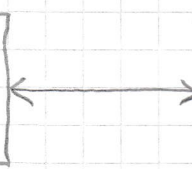
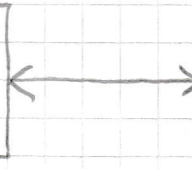
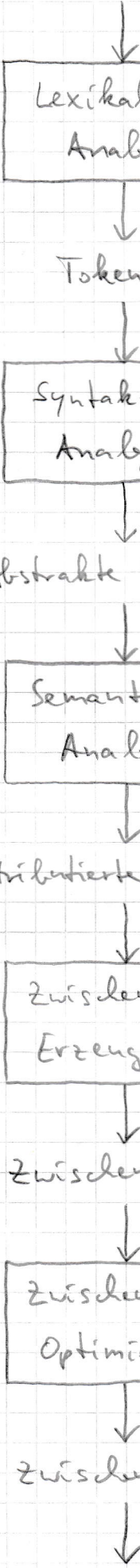
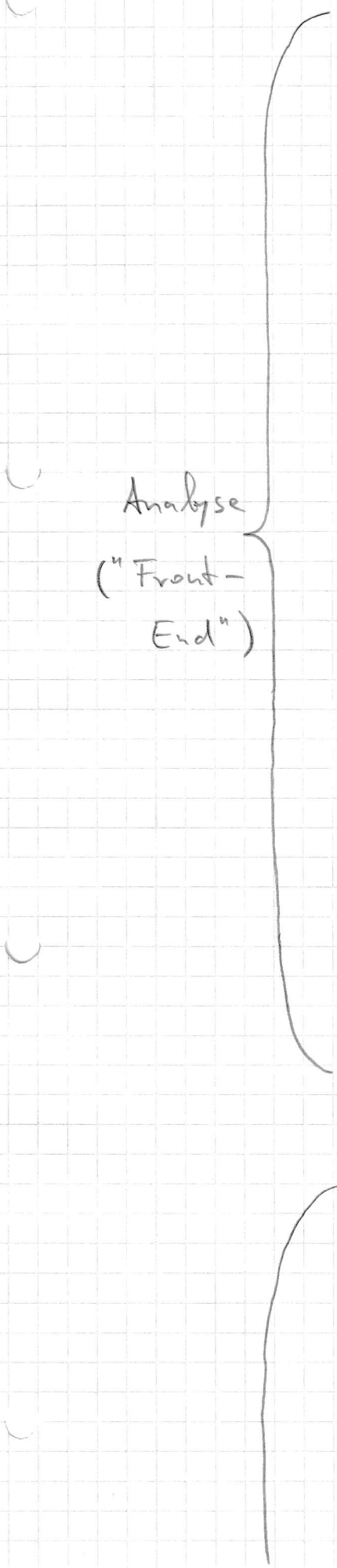
Frame
Layout

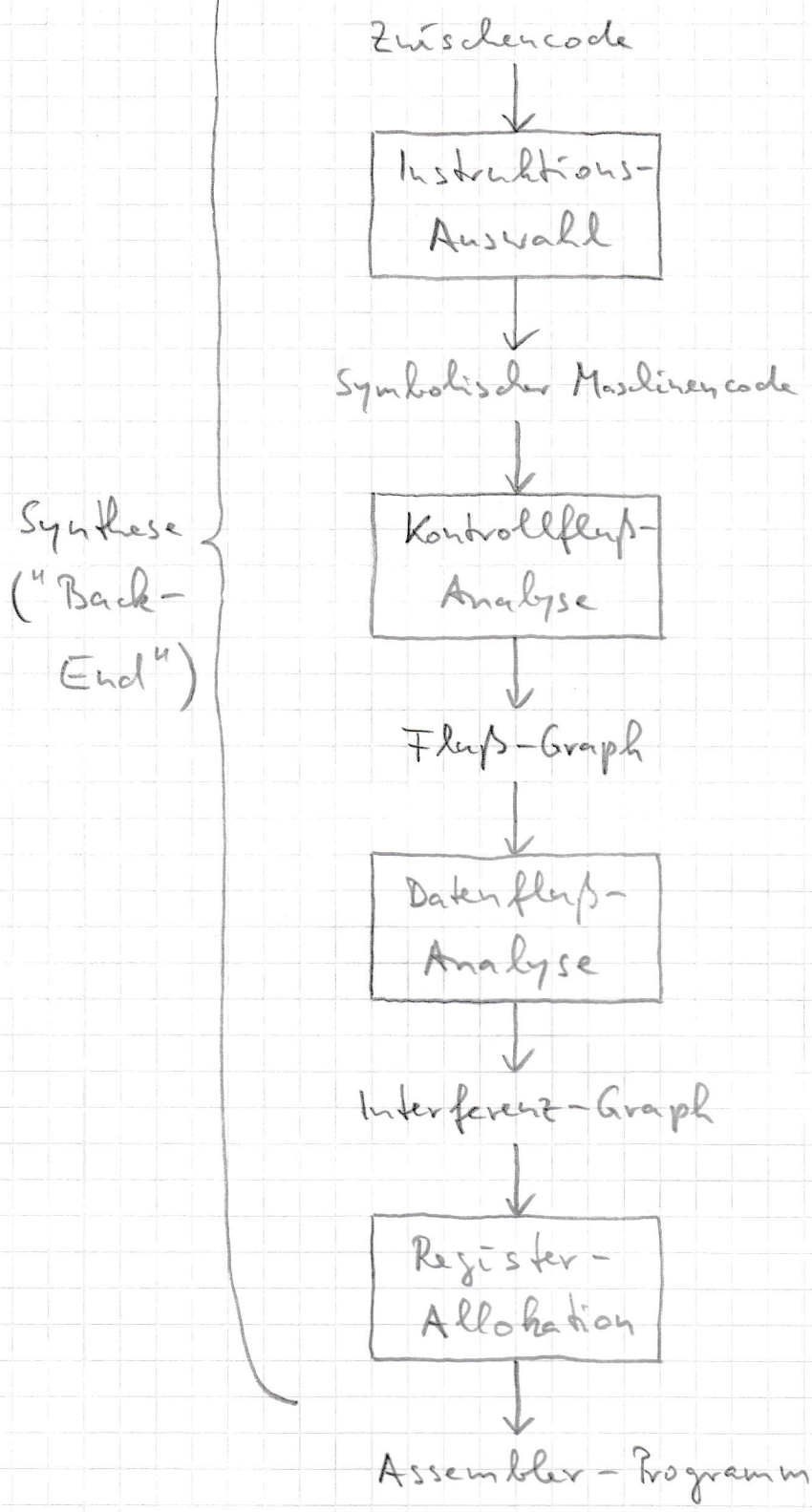
Zwischencode

Zwischen code-
Optimierung

Zwischencode

Analyse
("Front-
End")





Im Folgenden etwas vereinfacht:

1. Lexikalische Analyse
2. Syntaktische Analyse
3. Abstrakte Syntax
4. Semantische Analyse
5. Frame-Layout
6. Codegenerierung

Bsp.:

(7)

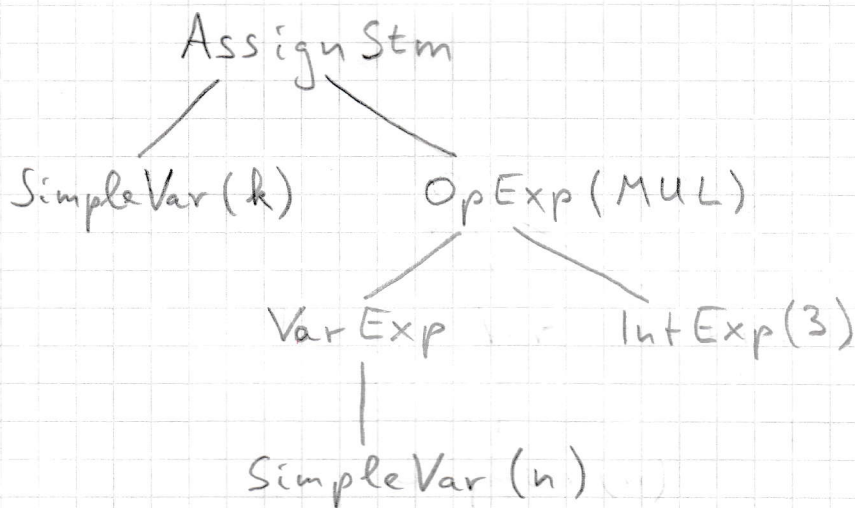
1) Quellcode

```
proc demo (n: int, ref k: int) {  
    k := n * 3;  
}
```

2) Tokens (für die Zuweisung)

IDENT("k") ASGN IDENT("n")
STAR INTLIT(3) SEMIC

3) Abstrakte Syntax (für die Zuweisung)



4) Symboltabelle (nach Analyse von "demo")

Level ϕ :
n \rightarrow var: int
k \rightarrow var: ref int

Level 1:
demo \rightarrow proc: (int, ref int)
int \rightarrow type: int
main \rightarrow proc: ()

⋮
(Bibliotheksfunktionen)

5) Variablenallokation (für "demo")

param 'n': $fp + \phi$ param 'k': $fp + 4$ local vars: ϕ

proc calls: none

6) Assemblercode (für die Zuweisung)

add	\$8, \$25, 4	; $fp + 4 =$ Adresse von k
ldw	\$8, \$8, ϕ	; Inhalt von k = Referenz
add	\$9, \$25, ϕ	; $fp + \phi =$ Adresse von n
ldw	\$9, \$9, ϕ	; Wert holen
add	\$10, ϕ , 3	; Konstante 3 nach \$10
mul	\$9, \$9, \$10	; $\$9 \leftarrow \$9 * \$10$
stw	\$9, \$8, ϕ	; $mem[\$8 + \phi] \leftarrow \9

2. Lexikalische Analyse

Aufgabe: Strom von Zeichen \rightarrow Strom von Token

Token: kleinste Bedeutung tragende Einheiten einer Prog. Spr.
(„Worte“ einer Prog. Spr.)

Bsp.: Token-Typ Beispiele aus einem konkreten Programm

IDENT	k ab123 das_Letzte
NUM	123 ϕ $\phi \times$ Affe
IF	if
COMMA	,
LE	\leq
RPAREN)

2.1. Reguläre Ausdrücke

Def.:

- c) Sprache — Menge von Strings
- b) String — endliche Folge von Zeichen
- a) Alphabet — alle zur Verfügung stehenden Zeichen (endlich viele)

Bsp.: Alphabet = $\{a, b\}$

Sprache = $\{b, ab, aab, aabb, \dots\}$

Beschreibung einer (regulären) Sprache durch „reguläre Ausdrücke“.

Jeder regulären Ausdruck steht für eine Menge von Strings.

Zeichen: a steht für $L(a) = \{a\}$

Alternative: β/γ steht für $L(\beta/\gamma) = L(\beta) \cup L(\gamma)$

Bsp: $L(a|b) = L(a) \cup L(b) = \{a\} \cup \{b\} = \{a, b\}$

Konkatenation: $\beta\gamma$ steht für $L(\beta\gamma) =$

$$\{st \mid s \in L(\beta) \wedge t \in L(\gamma)\}$$

Bsp: $L(a(b|c)) = \{st \mid s \in L(a) \wedge t \in L(b|c)\}$

$$= \{st \mid s \in \{a\} \wedge t \in \{b, c\}\} = \{ab, ac\}$$

Epsilon: ϵ steht für $L(\epsilon) = \{\epsilon\}$ (leerer String)

Kleene'scher Abschluss: β^* steht für $L(\beta^*) =$

$$L(\epsilon) \cup L(\beta) \cup L(\beta\beta) \cup L(\beta\beta\beta) \cup \dots$$

Bsp.: $L((ab)^*) = L(\epsilon) \cup L(ab) \cup L(abab) \cup \dots$

$$= \{\epsilon, ab, abab, ababab, \dots\}$$

Ein paar Beispiele:

a) Binärzahlen, die Vielfache von 2 sind: $(0|1)^*0$

b) Strings aus a und b, die keine 2 a's hintereinander enthalten:

$$b^*(a|b)^*(a|\epsilon)$$

c) Strings aus a und b, die mind. einmal 2 a's hintereinander enthalten:

$$(a|b)^*aa(a|b)^*$$

Weitere Abkürzungen:

$[abcd]$ bedeutet $(a|b|c|d)$

$[b-e]$ " $(b|c|d|e)$

$\beta^?$ " $(\beta|\epsilon)$

β^+ " $(\beta\beta^*)$

Bsp: Token einer Programmiersprache =

$if \dots \rightarrow IF$

\cdot bedeutet irgend ein Zeichen (außer Zeilenbruch)

$[a-z][a-z\phi-9]^*$ \rightarrow IDENT

$[\phi-9]^+$ \rightarrow NUM

$\phi^x[\phi-9a-fA-F]^+$ \rightarrow NUM

\backslash) \rightarrow RPAREN

(11)

Auflösen von Mehrdeutigkeiten:

a) Was ist "if8"? Ein IDENT oder IF NUM?

b) Was ist "if"? IDENT oder IF?

Zwei Regeln:

"Längste Übereinstimmung": "if8" ist IDENT

"Zuerst genannte Regel hat Vorrang": "if" ist IF

2.2. Deterministische Endliche Automaten (DEAs)

Bestandteile:

a) Endliche Menge von Zuständen (bezeichnet mit Zahlen)

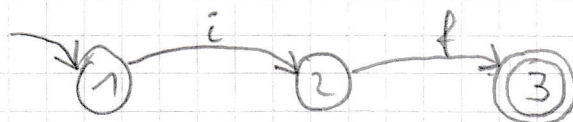
b) " " " Zustandsübergängen (markiert mit Eingabezeichen)

c) Ein ausgezeichnete Zustand, der "Startzustand"

d) Eine Teilmenge der Zustände sind "Endzustände"

"deterministisch": verschiedene Zustandsübergänge aus dem gleichen Zustand heraus sind mit verschiedenen Eingabezeichen markiert; kein Übergang ist mit dem leeren Eingabestring ϵ markiert.

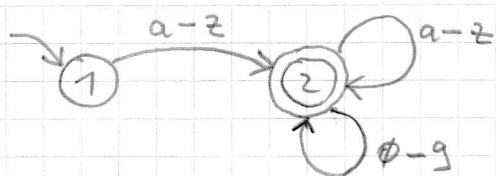
Bsp.:



IF

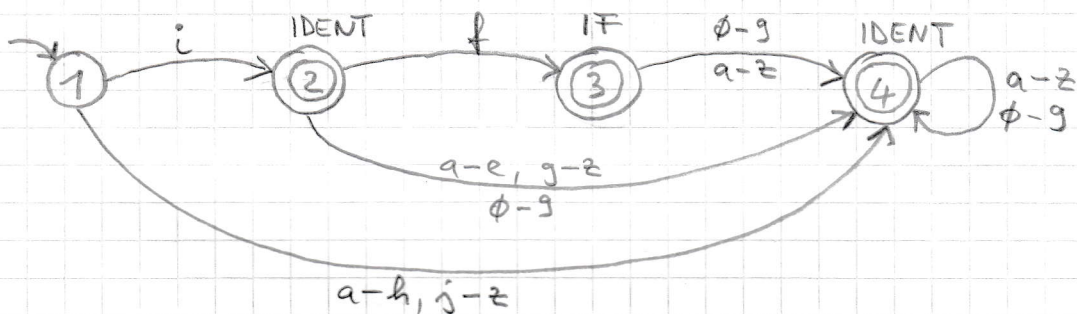
1: Startzustand, 3: Endzustand

Automat erkennt das Token IF.



Automat erkennt das Token IDENT.

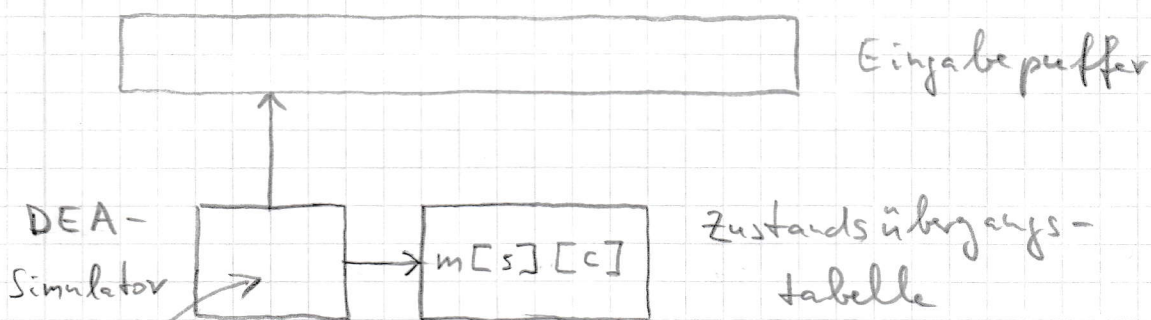
Kombinierter Automat zum Erkennen von IF und IDENT:



2 Fragen: Wie wird der Automat implementiert?

Wie wird der Automat konstruiert?

2.3. Tabellengesteuerter Scanner



Algorithmus:

$s := s_0$ // Startzustand

$c := \text{nextchar}()$; // nächstes Eingabezeichen

while ($c \neq \text{EOF}$) {

$s := m[s][c]$;

$c := \text{nextchar}()$;

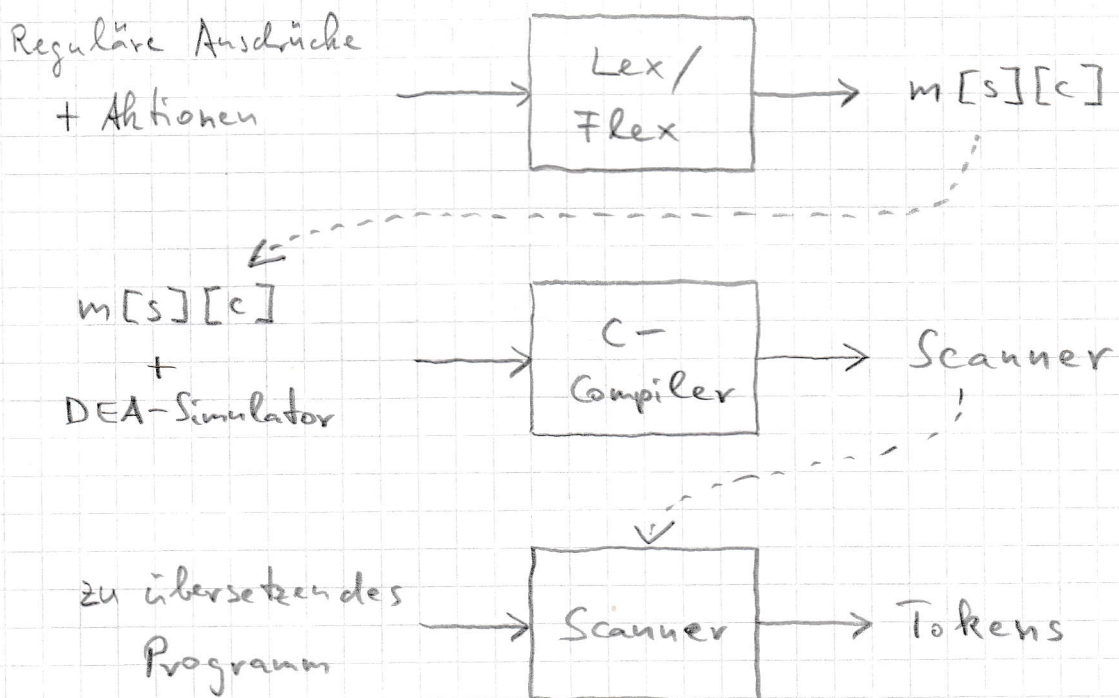
}

Die erkannte Sprache wird durch m festgelegt! (13)

3 Ergänzungen:

- Wenn ein Zustand s für ein Eingabezeichen c keinen Übergang besitzt, muß $m[s][c]$ eine spezielle Markierung bekommen ("Weiterlaufen nicht möglich")
- Es muß eine weitere Tabelle geben, die für jeden Zustand sagt, ob's ein Endzustand ist, und ggf. welches Token abgeliefert werden soll.
- Realisierung der Regel "längste Übereinstimmung":
Markieren des letzten erreichten Endzustandes und Weiterlaufen des Automaten bis "Weiterlaufen nicht mögl."

2.4. Scanner-Generatoren



Bsp.: $[\backslash \backslash t]^+$ { /* kein Token zurück */ }
if { return IF; }

$[a-z][a-z\phi-9]^* \{ \text{return IDENT}; \}$
 $\cdot \{ \text{error} ("..."); \}$

Bem.: Vermünftige Fehlerausgaben in späteren Phasen verlangen das Speichern einer "Koordinate" (mindestens der Zeilennummer im Quelltext) pro Token!

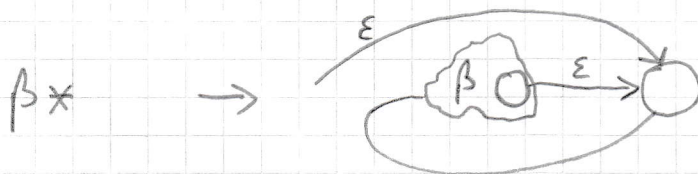
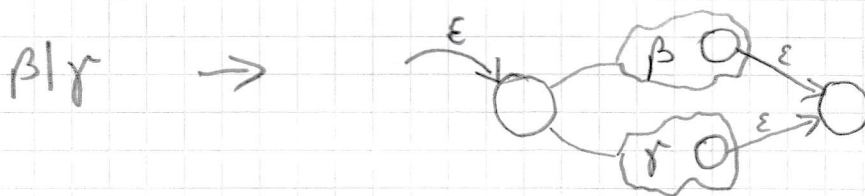
2.5. Konstruktion des DEA

Zwei Schritte:

a) reguläre Ausdrücke \rightarrow nichtdeterminist. endl. Automat ("NEA")

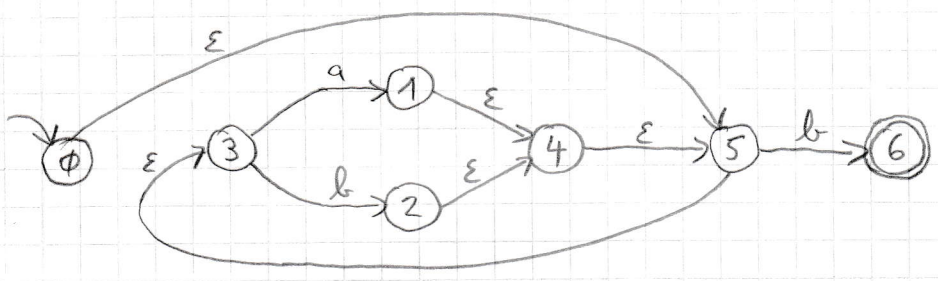
b) NEA \rightarrow DEA

zu a) "Thompson's Konstruktion":



Am Ende: Startzustand s_0 hinzufügen.

Bsp.: $(a|b)^* b$



Bem.: "nichtdeterministisch": der Automat muß bei gegeb. Zustand und Eingabezeichen raten, wo's langgeht, und er muß richtig raten!

zu b) "Teilmengenkonstruktion"

Idee: Menge aller NEA-Zustände, die für einen gegeb. Eingabestring erreicht werden können = ein DEA-Zustand

Def.:

ϵ -Abschluß (t): Menge der NEA-Zustände, die vom NEA-Zustand t allein über ϵ -Übergänge erreicht werden

ϵ -Abschluß (T): Menge der NEA-Zustände, die von irgendwelchen NEA-Zuständen $t \in T$ allein über ϵ -Übergänge erreicht werden

$move(T, c)$: Menge der NEA-Zustände, die von irgendwelchen NEA-Zuständen $t \in T$ durch das Eingabezeichen c erreicht werden

S : Menge der Zustände des DEA

$m[T][c]$: Übergangstabelle des DEA

Teilmengenkonstruktion:

Start: $S := \{\epsilon\text{-Abschluß}(s_0)\}$; // unmarkiert

Iteration: while (S enthält einen unmarkierten Zustand T) {

markieren T ;

16

for (jedes Eingabezeichen c) {

$X := \varepsilon$ -Abschluss (move (T, c));

if ($X \notin S$) $S := S \cup \{X\}$; // unmarkiert

$m[T][c] := X$;

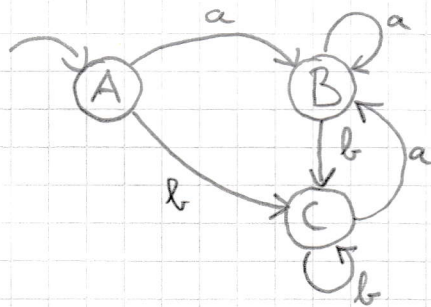
}

}

Bsp.: NEA wie oben

<u>NEA-Zustände</u>	<u>DEA-Zustand</u>	<u>Übergänge</u>
$\{\emptyset, 5, 3\}$	A	
$(A, a) = \{1, 4, 5, 3\}$	B	$m[A][a] := B$
$(A, b) = \{6, 2, 4, 5, 3\}$	C	$m[A][b] := C$
$(B, a) = \{1, 4, 5, 3\}$	B	$m[B][a] := B$
$(B, b) = \{6, 2, 4, 5, 3\}$	C	$m[B][b] := C$
$(C, a) = \{1, 4, 5, 3\}$	B	$m[C][a] := B$
$(C, b) = \{6, 2, 4, 5, 3\}$	C	$m[C][b] := C$

DEA:



Endzustand: C, weil $6 \in C$

3. Syntaktische Analyse

Reguläre Ausdrücke können nicht alle relevanten Konstrukte einer Progr. Sprache beschreiben, Bsp.: Korrekte Klammerung

Alphabet = $\{a, l, r\}$, Sprache = $\{a, lar, llarr, \dots\}$

Ursache: Keine Information über frühere Zustände des Automaten

Abhilfe: Automat wird mit Stack ausgestattet, Dann

können die "kontextfreien Sprachen" erkannt werden.

3.1. Kontextfreie Grammatiken

Bestandteile:

a) Menge von Terminalsymbolen T (das sind die Tokens aus Kap 2)

b) Menge von Nichtterminalsymbolen N ("Variablen")

c) Menge von Produktionen der Form $l \rightarrow r$, wobei $l \in N$ und $r \in (T \cup N)^*$ ("String von Grammatiksymbolen")

d) Ein ausgezeichnetes $S \in N$, das "Startsymbol".

Bsp: Korrekte Klammerung

$T = \{a, l, r\}$, $N = \{S\}$

Produktionen: $S \rightarrow a$

$S \rightarrow lS r$

Abkürzung dafür: $S \rightarrow a \mid lS r$

Bem: Jede reguläre Sprache ist eine kontextfreie Sprache.

a $S \rightarrow a$

ϵ $S \rightarrow \epsilon$

$\beta | \gamma$ $S \rightarrow \beta | \gamma$

(18)

 $\beta \gamma$ $S \rightarrow \beta \gamma$ β^* $S \rightarrow \varepsilon | \beta S$

("rechtsrekursiv")

alternativ: $S \rightarrow \varepsilon | S \beta$

("linksrekursiv")

Also: Reguläre Sprachen \subsetneq Kontextfreie Sprachen

Warum dann überhaupt Scanner? Effizienz, Einfachheit!

3.2. Ableitungen, Parse-Bäume, Mehrdeutigkeiten

Grammatik für arithmet. Ausdrücke:

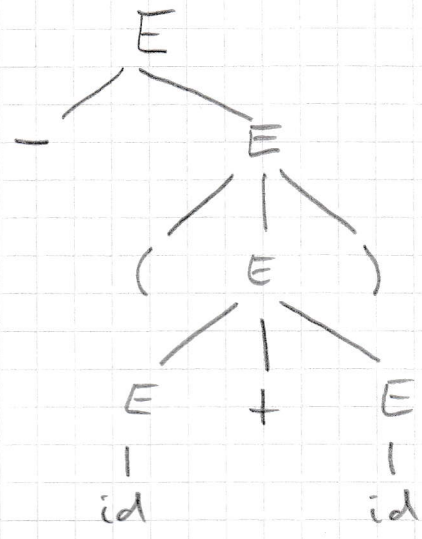
 $E \rightarrow E + E \mid E - E \mid E * E \mid E / E \mid (E) \mid -E \mid id$ Ist $-(id + id)$ ein Satz der Grammatik?Ja, denn es gibt eine "Ableitung": $E \Rightarrow -E \Rightarrow -(E)$ $\Rightarrow -(E + E) \Rightarrow -(id + E) \Rightarrow -(id + id)$

In jedem Herleitungsschritt zwei Entscheidungen:

a) Welches Nichtterminal wird ersetzt?

b) Welche Alternative für dieses Nichtterminal wird gewählt?

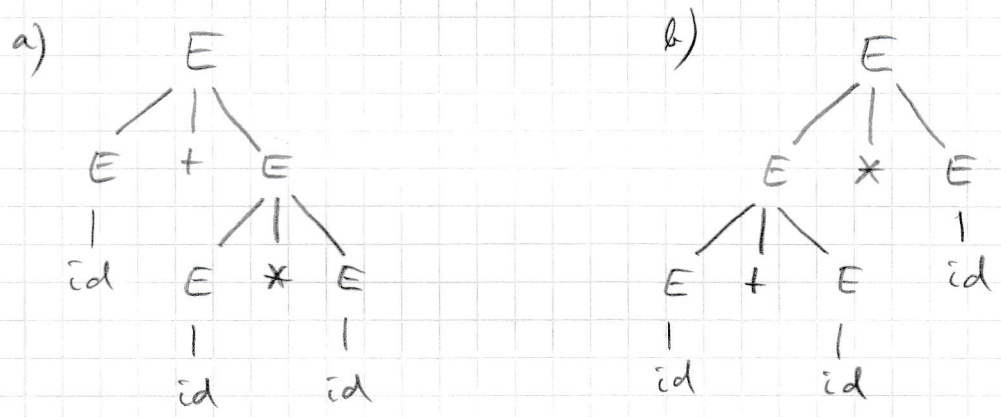
Linksableitungen ersetzen das jeweils am weitesten links stehende Nichtterminalsymbol, Rechtsableitungen das am weitesten rechts stehende Nichtterminalsymbol.Parse-Baum: grafische Darstellung einer Ableitung ohne Berücksichtigung der Reihenfolge.Bsp.: $-(id + id)$



Syntaxanalyse = Finden einer Ableitung = Konstr. d. Parse-Baums bei gegeb. Satz

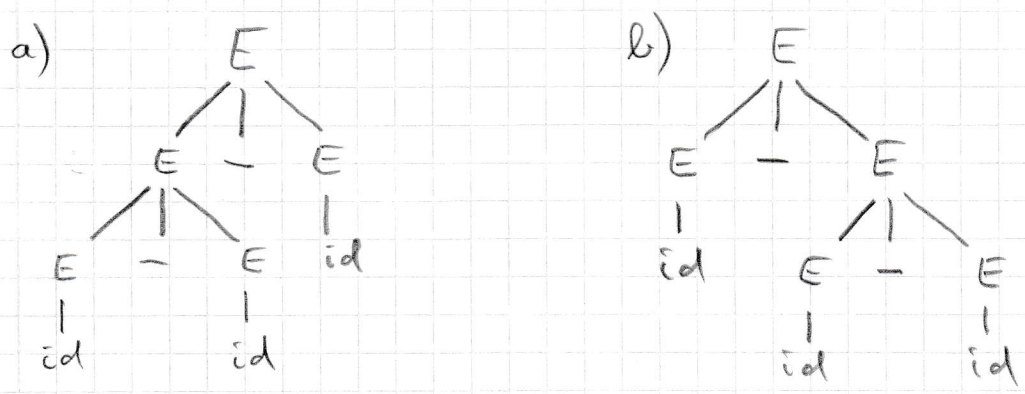
Eine Grammatik, bei der es mind. einen Satz gibt, der durch verschiedene Parse-Bäume repräsentiert wird, heißt mehrdeutig.

Bsp.: Grammatik von oben, Satz $id + id * id$



Wir wollen haben: a), wg. "Priorität"

Bsp.: Grammatik von oben, Satz $id - id - id$



Wir wollen haben: a), wg. "Assoziativität"

Grammatik wird eindeutig durch neue Nichtterminale: (20)

$E \rightarrow E + T \mid E - T \mid T$

E : "Expression"

$T \rightarrow T * F \mid T / F \mid F$

T : "Term"

$F \rightarrow -F \mid \text{id} \mid (E)$

F : "Faktor"

Bem.: ^(Links-) Priorität durch $E/T/F$, ^(Links-) Assoziativität durch Rekursion

Die beiden Beispiele behandeln!

3.3. Prädiktive Parser (Top-Down-Parser, LL(1)-Parser)

Verfahren: "Rekursiver Abstieg" - jedes Nichtterminalsymbol wird durch eine rekursive Prozedur dargestellt, jede Produktion ist ein Zweig in der betr. Prozedur

Eigenschaften: leicht zu programmieren, wenig mächtig

Bsp.: $S \rightarrow \text{if } E \text{ then } S$

$S \rightarrow \text{while } E \text{ do } S$

Annahme: in t steht das jeweils nächste Token

```
void S(void) {
```

```
    switch (t) {
```

```
        case IF: t = getToken();
```

```
                E();
```

```
                if (t != THEN) error(...);
```

```
                t = getToken();
```

```
                S();
```

```
                break;
```

```
        case WHILE: ;
```

```
    }
```

```
}
```

Aber: $E \rightarrow E+T \mid E-T \mid T$ läßt sich so nicht (21)

parse! Warum? Erstes Terminalsymbol muß genug Information liefern, um die richtige Produktion auszuwählen!

3.3.1. First- und Follow-Mengen

Sei $\beta \in (T \cup N)^*$ eine bel. Folge von Grammatiksymbolen.

Def.: $FIRST(\beta) = \{t \in T \mid \beta \Rightarrow \dots \Rightarrow t\gamma\}$

d.h. diejenigen Terminalsymbole, mit denen ein aus β abgeleiteter

Satz beginnen kann. Gilt $\beta \Rightarrow \dots \Rightarrow \varepsilon$, dann ist $\varepsilon \in FIRST(\beta)$.

1. Forderung für Top-Down-Parser:

Fasse alle Produktionen für ein Nichtterminalsymbol A zusammen zu $A \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$. Dann muß

$FIRST(\beta_i) \cap FIRST(\beta_j) = \{\}$ f. alle $i \neq j$ sein.

Das reicht nicht: Falls $\beta \Rightarrow \dots \Rightarrow \varepsilon$, dann kommt's drauf an, was dem β nachfolgt!

Sei $A \in N$ ein Nichtterminalsymbol und S das Startsymbol.

Def.: $FOLLOW(A) = \{t \in T \mid S \Rightarrow \dots \Rightarrow \beta A t \gamma\}$

d.h. diejenigen Terminalsymbole, die in einer Satzform direkt rechts neben A stehen können.

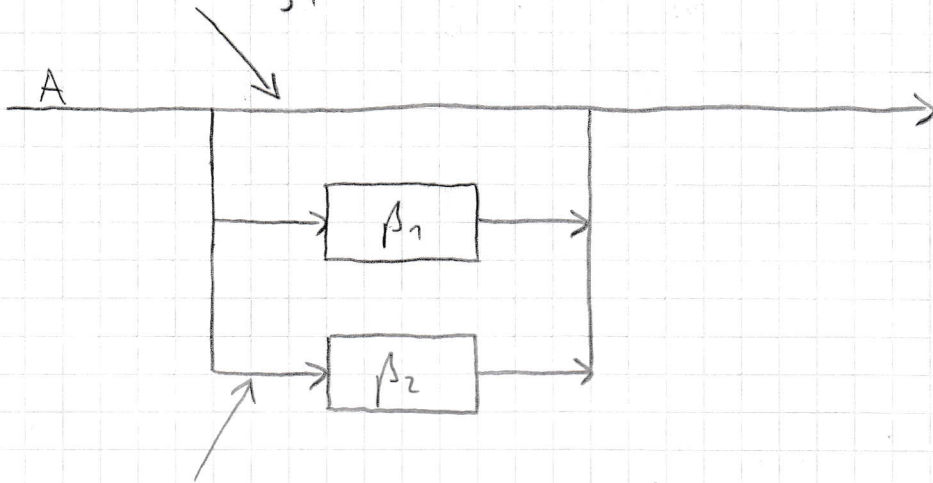
2. Forderung für Top-Down-Parser:

Für jedes Nichtterminalsymbol A mit $A \Rightarrow \dots \Rightarrow \varepsilon$ muß

$FIRST(A) \cap FOLLOW(A) = \{\}$ sein.

Ansichtlich in Syntax-Diagrammen:

Hier entlang, wenn nächstes Token $\in \text{FOLLOW}(A)$ ist! (22)



Hier entlang, wenn nächstes Token $\in \text{FIRST}(\beta_2)$ ist!

Grammatiken, die beide Forderungen erfüllen, heißen LL(1).

1. "L" $\hat{=}$ Eingabe wird von links nach rechts gelesen

2. "L" $\hat{=}$ Beim Parsen wird eine Links-Ableitung erzeugt

"(1)" $\hat{=}$ 1 Token Vorausschau ("Lookahead")

Eine mehrdeutige oder linksrekursive Grammatik ist nicht LL(1).

3.3.2. Elimination von Links-Rekursion, Links-Faktorisierung

a) Transformation von $A \rightarrow A\beta \mid \gamma$ zu:

$$A \rightarrow \gamma A'$$

$$A' \rightarrow \beta A' \mid \epsilon$$

Bsp.: $E \rightarrow E+T \mid T$ umformen zu:

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \epsilon$$

b) Transformation von $A \rightarrow \beta\gamma_1 \mid \beta\gamma_2$ zu:

$$A \rightarrow \beta A'$$

$$A' \rightarrow \gamma_1 \mid \gamma_2$$

3.3.3. Wo ist beim Top-Down-Parser mit rek. Abstieg der Stack? \rightarrow

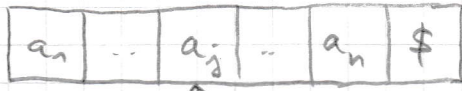
→ Das ist der Laufzeit-Stack der Programmiersprache,
in der der Compiler geschrieben ist!

3.4. Shift-Reduce-Parser (Bottom-Up-Parser, LR(1)-Parser)

Schwäche beim LR(1)-Parser: Entscheidung, welche Produktion benutzt wird, fällt aufgrund eines Tokens!

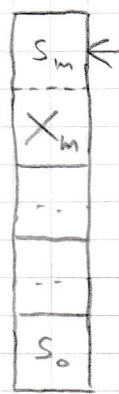
Idee: Zurückschaltung dieser Entscheidung, bis die gesamte rechte Seite der Produktion gelesen wurde (und noch ein Token mehr).

3.4.1. Tabellengesteuerter Parser

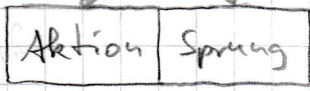
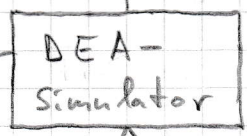


Eingabe

Stack



gehören jeweils zusammen!



($\$$: Ende der Eingabe,
 a_j : Eingabe-Token,
 s_m : Zustand,
 x_m : Grammatiksymbol)

DEA-Simulator:

```

push( $s_0$ );
accept == false;
while (!accept) {
  sei  $s$  der Zustand oben auf dem Stack;
  sei  $a$  das nächste Token;
  switch (aktion[ $s$ ][ $a$ ]) {
    case shift( $n$ ):
      push( $a, n$ );
      getToken();
      break;
  }
}

```

case reduce (m):

sei $A \rightarrow \beta$ die Produktion mit der Nummer m;

sei k die Anzahl der Grammatiksymbole von β ;

pop (k Elemente);

sei t der Zustand, der jetzt oben auf dem Stack ist;

push (A, sprung [t][A]);

break;

case accept:

accept := true;

break;

case error:

error ();

break;

}

}

Bsp.:

r0: $S \rightarrow E \$$

r2: $E \rightarrow T$

r1: $E \rightarrow T + E$

r3: $T \rightarrow a$

Zustand	Aktion			Sprung	
	a	+	\$	E	T
1	s5			2	3
2			acc		
3		s4	r2		
4	s5			6	3
5		r3	r3		
6			r1		

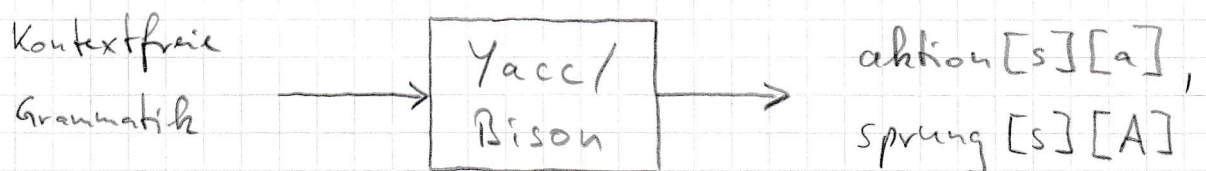
leer $\hat{=}$ error

<u>Stack</u>	<u>Eingabe</u>	<u>Aktion / Sprung</u> (25)
1	a + a \$	s 5
1 a 5	+ a \$	r 3 → 3
1 T 3	+ a \$	s 4
1 T 3 + 4	a \$	s 5
1 T 3 + 4 a 5	\$	r 3 → 3
1 T 3 + 4 T 3	\$	r 2 → 6
1 T 3 + 4 E 6	\$	r 1 → 2
1 E 2	\$	acc

Bem.:

- Die Grammatiksymbole (a, T, +, E) müssen nicht wirklich auf dem Stack stehen; der Zustand genügt!
- LR(1) bedeutet: Eingabe von links nach rechts lesen, Rechtsableitung konstruieren, mit einem Token Vorausschau.
Stack + Eingabe = Rechtsableitung in umgekehrter Reihenfolge:
 $E \$ \Rightarrow T + E \$ \Rightarrow T + T \$ \Rightarrow T + a \$ \Rightarrow a + a \$$
- Die Tabelleneinträge bestimmen sowohl die erkannte Sprache als auch die Parsing-Methode (LR(0), SLR(1), LR(1), LALR(1))

3.4.2. Parser-Generatoren



Weiteres Vorgehen wie bei Scanner-Generatoren.

Bsp.:

% token PLUS IDENT

% start expr

%%

expr : term PLUS expr

| term

;

term : IDENT

;

%%

3.4.3. Konstruktion von LR(ϕ)-Tabellen

LR(ϕ): Schiebe / reduziere - Entscheidung ohne Vorausschau.

Def.: Eine Produktion $A \rightarrow \beta$ mit einem Punkt irgendwo auf der rechten Seite heißt "LR(ϕ)-Element" (kurz: Element).

Bsp.: Die Produktion $A \rightarrow aBc$ bringt 4 Elemente hervor:

$A \rightarrow \cdot aBc$ $A \rightarrow a \cdot Bc$ $A \rightarrow aB \cdot c$ $A \rightarrow aBc \cdot$

Def.: Eine Menge von Elementen heißt "Zustand" des DEA.

Algorithmus zur Tabellenkonstruktion benutzt zwei Hilfsfunktionen.

Sei I eine Menge von Elementen und X ein Grammatiksymbol.

a) Hülle (I) {

do {

Der Punkt kennzeichnet den Stand der Analyse.

for (jedes Element $A \rightarrow \alpha.X\beta$ in I) { (27)

for (jede Produktion $X \rightarrow \gamma$) {

$I := I \cup \{X \rightarrow \gamma\}$

}
} while (I hat sich in dieser Iteration geändert);
return I ;

}

b) Sprung (I, X) {

$J := \{\}$;

for (jedes Element $A \rightarrow \alpha.X\beta$ in I) {

$J := J \cup \{A \rightarrow \alpha.X.\beta\}$;

}

return Hülle (J);

}

Tabellenkonstruktion:

Ergänze die Grammatik um die neue Startproduktion

$S' \rightarrow S \$$. Sei T die Menge der Zustände des DEA und

E die Menge seiner Kanten.

$T := \{ \text{Hülle}(\{S' \rightarrow .S \$\}) \}$;

$E := \{\}$;

do {

for (jeden Zustand I in T) {

for (jedes Element $A \rightarrow \alpha.X\beta$ in I) {

$J := \text{Sprung}(I, X)$;

$$T := T \cup \{ \} ;$$

$$E := E \cup \{ I \xrightarrow{x} \} ;$$

$$\}$$

$$\}$$

while (T oder E haben sich in dieser Iteration geändert);

Aber: Für das Symbol \$ wird kein Sprung berechnet;

jeder Zustand mit Element $S' \rightarrow S. \$$ akzeptiert bei Token = \$!

Berechnung der Menge der Reduzier-Aktionen R:

$$R := \{ \} ;$$

for (jeder Zustand I in T) {

 for (jedes Element $A \rightarrow \alpha.$ in I) {

$$R := R \cup \{ (I, A \rightarrow \alpha) \} ;$$

 }

}

↑ ↑

Zustand Reduktion mit dieser Produktion

Bsp.:

rφ: $S \rightarrow E \$$ r2: $E \rightarrow T$

r1: $E \rightarrow T + E$ r3: $T \rightarrow a$

① (Startzustand) = Hülle ({ $S \rightarrow . E \$$ })

Elemente: $S \rightarrow . E \$$

$E \rightarrow . T + E$

$E \rightarrow . T$

$T \rightarrow . a$

} ①

Füllen der Tabellen: →

a) Für jede Kante $I \xrightarrow{x} J$:

Falls X ein Terminalsymbol ist,

setze Aktion $[I][X] := \text{shift } J$

Falls X ein Nichtterminalsymbol ist,

setze Sprung $[I][X] := J$

b) Für jeden Zustand I , der das Element

$S' \rightarrow S. \$$ enthält, setze Aktion $[I][\$] := \text{accept}$.

c) Für jeden Zustand I , der ein Element

$A \rightarrow \alpha.$ enthält, setze Aktion $[I][y] := \text{reduce } u$

für alle Token y ($L R(\phi)$, keine Voraussetzungen!),

wobei u die Produktion ist, die $A \rightarrow \alpha.$ hervorbringt hat.

d) Setze alle anderen Einträge auf error.

① \xrightarrow{E} Sprung (①, E): $S \rightarrow E \cdot \$$ } ②

① \xrightarrow{T} Sprung (①, T): $E \rightarrow T \cdot + E$
 $E \rightarrow T \cdot$ } ③

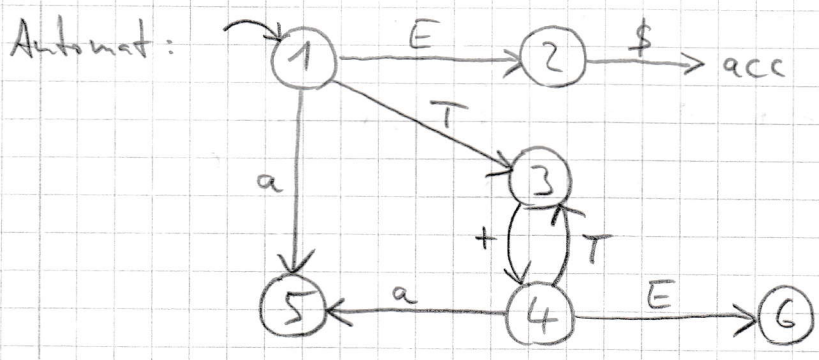
③ $\xrightarrow{+}$ Sprung (③, +): $E \rightarrow T + \cdot E$
 $E \rightarrow \cdot T + E$
 $E \rightarrow \cdot T$
 $T \rightarrow \cdot a$ } ④

① \xrightarrow{a} Sprung (①, a): $T \rightarrow a \cdot$ } ⑤

④ \xrightarrow{E} Sprung (④, E): $E \rightarrow T + E \cdot$ } ⑥

④ \xrightarrow{T} Sprung (④, T): $E \rightarrow T \cdot + E$
 $E \rightarrow T \cdot$ } ③

④ \xrightarrow{a} Sprung (④, a): $T \rightarrow a \cdot$ } ⑤



Reduzier - Aktionen:

(③, $E \rightarrow T$)

(⑤, $T \rightarrow a$)

(⑥, $E \rightarrow T + E$)

LR(Φ)-Tabelle:

Zustand	Aktion			Sprung	
	a	+	\$	E	T
1	s5			2	3
2			acc		
3	r2	s4, r2	r2		
4	s5			6	3
5	r3	r3	r3		
6	r1	r1	r1		

③, + : "Schiebe-Reduziere-Konflikt"

Grammatik ist nicht LR(Φ)!

Bem.: Parsergeneratoren lösen solche Konflikte nach zwei Regeln auf:

- a) Bei Schiebe-Reduziere-Konflikten wird immer geschoben. Das ist in selteneren Fällen akzeptabel ~~aber nicht ideal~~ \rightarrow
- b) Bei Reduziere-Reduziere-Konflikten hat die zuerst längeschriebene Produktion Vorrang. Das ist praktisch immer unbrauchbar!

Konsequenz: Bei Konflikten ist der Automat zu analysieren und die Grammatik umzuformen!

3.4.4. Konstruktion von SLR(1)-Tabellen

Der Konflikt in der LR(Φ)-Tabelle kann leicht aufgelöst werden: Reduktion mit $A \rightarrow \gamma$ nur dann, wenn nächstes Token $\in FOLLOW(A)$ ist!

Bsp.: Das "hängende else" ist eine Mehrdeutigkeit, die sich beim if-Statement ergeben kann:

```
if_stmt      : IF LPAREN exp RPAREN stmt
              | IF LPAREN exp RPAREN stmt ELSE stmt
              ;
```

Der Parsergenerator meldet einen Schiebe/Reduziere-Konflikt in einem der Zustände des Parsers, wenn das Lookahead-Token ELSE ist:

if_stmt → IF LPAREN exp RPAREN stmt.

if_stmt → IF LPAREN exp RPAREN stmt. ELSE stmt

Warum schiebt er bei ELSE nicht einfach?

```
if (x < 5)
    if (y > 10)
        ...
    else
        ...
```

Wozu gehört das else...?

Schieben → else gehört zum inneren if

reduzieren → " " " äußeren "

$R := \{\};$

for (jeden Zustand I in T) {

for (jedes Element $A \rightarrow \alpha$ in I) {

for (jedes Token X in $FOLLOW(A)$) {

$R := R \cup \{(I, X, A \rightarrow \alpha)\};$

}
}
}

↑ Zustand ↑ lookahead ↑ Reduktion mit dieser Produktion

Unser Bsp.: Automat bleibt derselbe, aber Tabelle enthält

weniger Reduktionen:

$FOLLOW(E) = \{\$ \}$

$FOLLOW(T) = \{+, \$ \}$

Zustand	Aktion			Sprung	
	a	+	\$	E	T
1	SS			2	3
2			acc		
3		S4	r2		
4	SS			6	3
5		r3	r3		
6			r1		

3.4.5. Konstruktion von LR(1)- bzw. LALR(1)-Tabellen

LR(1)-Elemente sind Paare: (LR(0)-Element, Token)

Mehr Information in den Zuständen des DEA \Rightarrow viel mehr Zustände; weniger Konflikte in der Tabelle.

LALR(1): Wie LR(1), aber zusammenfassen der Zustände, die sich nur im Token unterscheiden: ~ Faktor 10 weniger Zustände als LR(1)!

4. Abstrakte Syntax

(32)

Compiler soll Code (oder Datenstruktur) erzeugen: Verknüpfung von Aktionen mit Grammatik-Produktionen notwendig.

4.1. Semantische Aktionen in LL-Parsern

Einbau von Aktionen in die Parser-Routinen. Bsp.:

$E \rightarrow TE'$ $E' \rightarrow +TE' \mid \varepsilon$ $T \rightarrow n$

```
void E(void) {  
    T();  
    E'();  
}
```

```
void E'(void) {  
    if (token == PLUS) {  
        getToken();  
        T();  
        printf("add");  
        E'();  
    }  
}
```

```
void T(void) {  
    if (token == NUM) {  
        printf("push %d", token.value);  
        getToken();  
    } else {  
        error("Zahl erwartet");  
    }  
}
```

Eingabe: $1 + 2 + 3$

(33)

Ausgabe: push 1

push 2

add

push 3

add

4.2. Semantische Aktionen in LR-Parsern

Einbau von Aktionen in die Reduktionen. Bsp.:

$E \rightarrow E + T \mid T \quad T \rightarrow n$

$E : E + T$

{ printf("add"); }

| T

;

$T : n$

{ printf("push %d", n.value); }

;

Eingabe / Ausgabe wie oben! Warum?

Rechtsableitung: $E \Rightarrow E + T \Rightarrow E + 3 \Rightarrow E + T + 3 \Rightarrow E + 2 + 3$

$\Rightarrow T + 2 + 3 \Rightarrow 1 + 2 + 3$. In umgekehrter Reihenfolge:

$T \rightarrow 1$ push 1

$E \rightarrow T$

$T \rightarrow 2$ push 2

$E \rightarrow E + T$ add

$T \rightarrow 3$ push 3

$E \rightarrow E + T$ add

Möglichkeit: ganzer Compiler eingebettet in semantischen (34)
Aktionen. Nachteile: schwer zu lesen und zu warten (syntakt.
Analyse, semant. Analyse und Codeerzeugung miteinander verweben),
Informationsfluß ist auf die Parse-Reihenfolge beschränkt.
Besser: Erzeugen einer Datenstruktur ("Abstrakte Syntax").

4.3. Semantische Werte in LL-Parsern

Jede Parser-Routine gibt den "semantischen Wert" (z. B. Baum-
knoten) des von ihr erkannten Nichtterminals zurück. Bsp.:

ifstm \rightarrow IF cond THEN stm

```
Node * ifstm (wid) {  
    Node * c, * s;  
    if (token == IF) getToken(); else error();  
    c = cond();  
    if (token == THEN) getToken(); else error();  
    s = stm();  
    return newNode (IFSTM, c, s);  
}
```

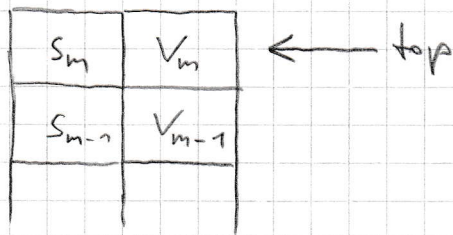
4.4. Semantische Werte in LR-Parsern

Jede Reduktion produziert den semantischen Wert ihrer linken
Seite. Bsp.:

```
ifstm : IF cond THEN stm  
      { $$ = newNode (IFSTM, $2, $4); }  
      ;
```

Technische Realisierung: "Semantik-Stack"

parallel zum Zustands-Stack.

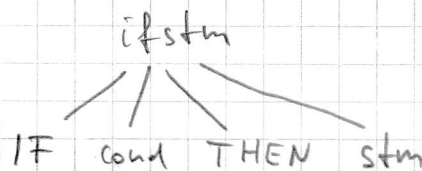


4.5. Abstrakte Syntax

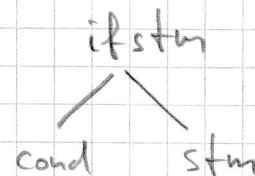
Zweck: sauberes Interface zwischen Parser und späteren Phasen.

Entsteht aus Parse-Baum durch Weglassen unwichtiger Details.

Bsp.: a) Parse-Baum



Syntax-Baum



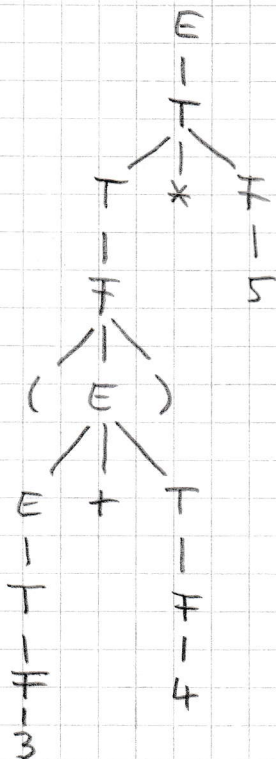
b) $E \rightarrow E + T \mid T$

Eingabe $(3+4) * 5$

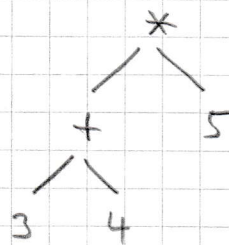
$T \rightarrow T * F \mid F$

$F \rightarrow n \mid (E)$

Parse-Baum



Syntax-Baum



Errichtet wird das durch z.B.:

$T : F$

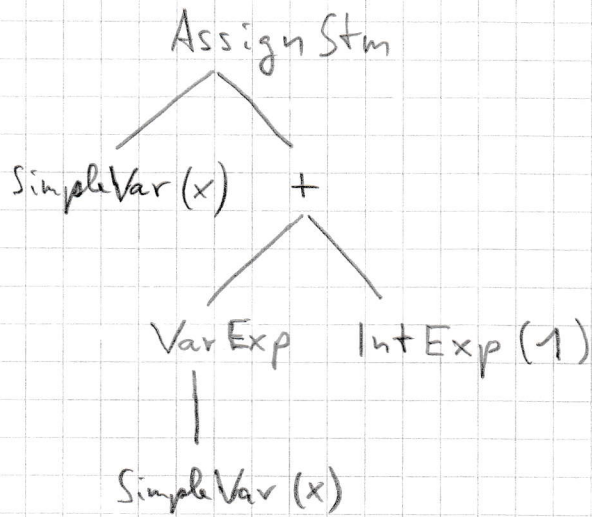
$\{ \$\$ = \$1 ; \}$

;

(ist häufig die "Default-Aktion" im Parser-Generator)

c) $x := x + 1$

35a

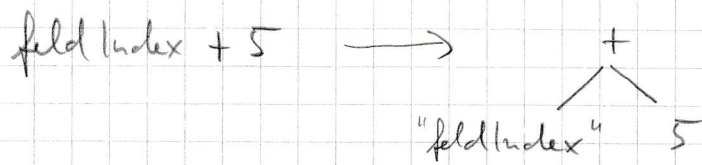


SimpleVar $\hat{=}$ Berechnung d. Adresse einer einfachen Var.

VarExp $\hat{=}$ Berechnung d. Inhalts aus einer Adresse

4.5.1. Effiziente Handhabung von Bezeichnern

(36)



Das ist ganz ineffizient, da bei jedem Durchgang durch diesen Knoten ein Stringvergleich stattfinden muß! Besser: Darstellung des Bezeichners durch "Symbol" (= eindeutige Zahl bzw. eindeutiger Zeiger). Funktionen d. Symbolverwaltung:

Erzeugen: String \rightarrow Symbol

ext. Darstellung: Symbol \rightarrow String

Anordnung: Symbol1 < Symbol2 ?

Technische Realisierung durch Hashtabelle.

4.6. Attributierung

Bedeutet: Anheften von Eigenschaften (Typ, Größe, Speicherort, etc.) an die Knoten des Syntax-Baumes.

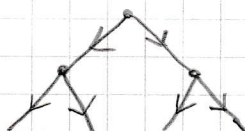
Realisierung: durch rekursive Attribut-Auswertung

"Synthetisierte Attribute" (Blätter \rightarrow Wurzel):



Bsp.: Auswertung konstanter Ausdrücke

"Vererbte Attribute" (Wurzel \rightarrow Blätter):



Bsp.: char *p1, *p2;

Im Allgemeinen: Attribute beider Sorten, berechnet in möglichen mehreren Durchläufen durch den Baum

5. Semantische Analyse

Aufgabe: Sammeln von Informationen über Bezeichner, überprüfen von Einschränkungen bei der Benutzung von Bezeichnern ("Typprüfungen").

5.1. Typen {, Literalen und Ausdrücken

Typ: Information (zur Übersetzungszeit), wie Bitmuster zu interpretieren sind (zur Laufzeit)

5.1.1. Darstellung von Typen

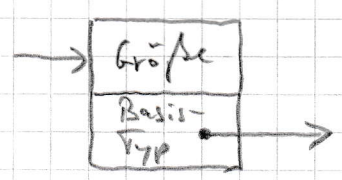
Typen können Bezug auf andere Typen nehmen (Bsp.: Array)

⇒ Darstellung durch "Typgraph"

a) Primitive ("eingebaute") Typen:

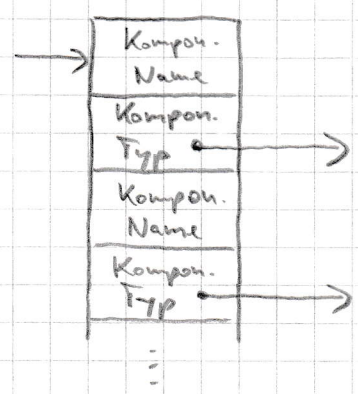


b) Arrays:

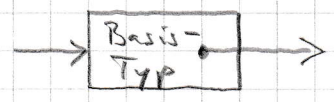


Bem.: Manche Sprachen betrachten die Größe als nicht zum Typ gehörend!

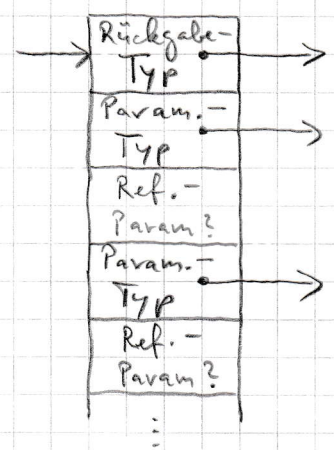
c) Records:



d) Pointer

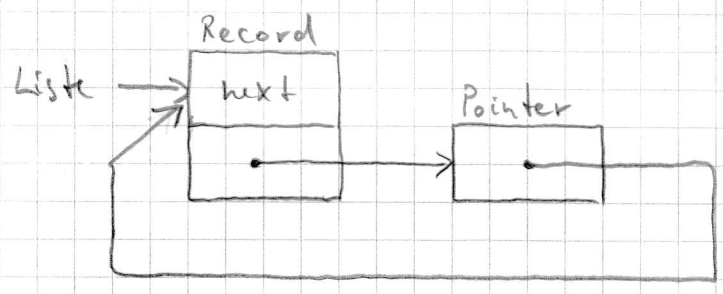


e) Funktionen



Achtung: Typgraphen können Zyklen enthalten:

```
typedef struct liste {
    struct liste *next;
} Liste;
```



5.1.2. Gleichheit von Typen

a) Namensgleichheit: Zwei Typen sind genau dann gleich, wenn sie denselben Namen haben.

```
Bsp.: type a = array [3] of int;
      type b = a;
      type c = array [3] of int;
```

Alle 3 Typen sind verschieden.

b) Ausdrucksgleichheit: Zwei Typen sind genau 39 dann gleich, wenn sie durch denselben Typausdruck konstruiert wurden (d.h. durch denselben Typgraphen repräsentiert werden).

Bsp.: Typdefinitionen wie oben. Die Typen a und b sind gleich, aber verschieden von c. (So macht es SPL.)

c) Strukturgleichheit: Zwei Typen sind genau dann gleich, wenn ihre Typgraphen (rekursiv) die gleiche Struktur aufweisen.

Bsp.: Typdefinitionen wie oben. Alle 3 Typen sind gleich.

5.2. Symboltabellen

Realisieren Abbildungen: Bezeichner \mapsto Attribute
(bzw. wegen Effizienz: Symbol \mapsto Attribute)

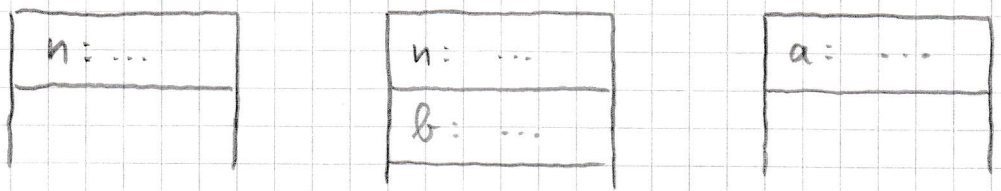
5.2.1. Mehrstufige Symboltabellen

Geschichtete Gültigkeitsbereiche von Bezeichnern \Rightarrow
verschiedene Attributsätze für den gleichen Bezeichner \Rightarrow
mehrstufige Symboltabelle.

Bsp.:

```
proc a () {  
  var n: int;  
  proc b () {  
    var n: array [10] of int;  
    ⋮  
  }  
}
```

Sym. Tab. von b() → Sym. Tab. von a() → glob. Sym. Tab.



(SPL benutzt zweistufige Symboltabellen.)

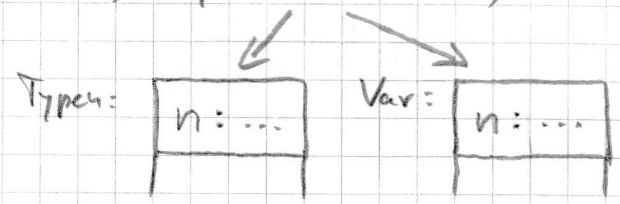
5.2.2. Parallele Symboltabellen

Mehrere Namensräume ⇒ verschiedene Attributsätze für den gleichen Bezeichner (auf derselben Stufe!) ⇒ parallele Symboltabellen (heißt zusammen: "Environment")

```

Bsp.: proc a() {
        type n = array [10] of int;
        var n: n;
    }
  
```

... → Sym. Tab. von a() → ...



einer Stufe

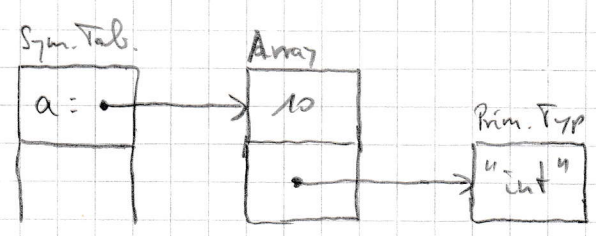
(SPL hat einen gemeinsamen Namensraum f. alle Bezeichner.)

5.2.3. Attribute für Bezeichner in SPL

a) Typ-Bezeichner ↦ Typ

```

Bsp.: type a = array [10] of int;
  
```



b) Variablen-Bezeichner \mapsto (Typ, ist Referenz)

(41)

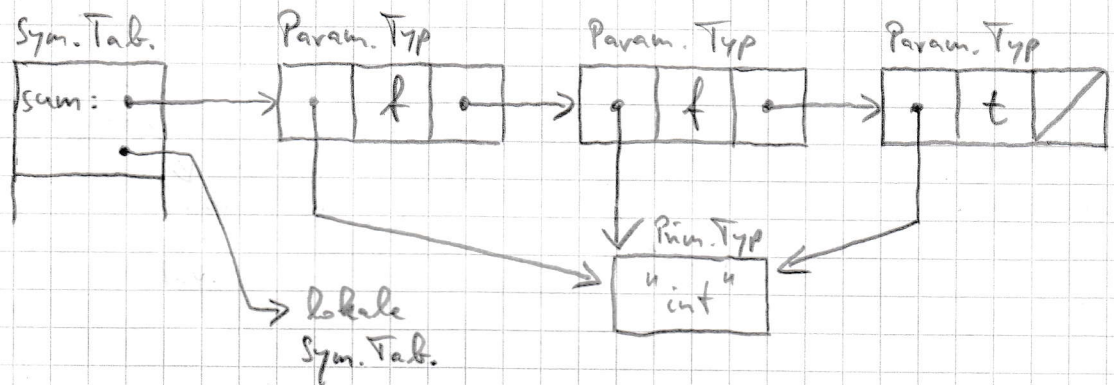
lokale Variablen: ist Ref. = false

Wertparameter: ist Ref. = false

Referenzparameter: ist Ref. = true

c) Prozedurbezeichner \mapsto (Parameter-Typen, lokale Sym.Tab.)

Bsp.: `proc sum (i: int, j: int, ref k: int) { ... }`



5.3. Type-Checking

Der Type-Checker ist eine rekursive Funktion, die die Abstrakte Syntax inspiziert und einen Typ zurückgibt:

`Type * checkNode (Absyn * node, Table * symTable) {`

`switch (node \rightarrow type) {`

`case ABSYN_NAME_TYP:`

`return checkNameType (node, symTable);`

`case ABSYN_PROCDECL:`

`return checkProcDecl (node, symTable);`

`case ABSYN_OP:`

`return checkOp (node, symTable);`

`}`

`:`

5.3.1. Type-Checking von Ausdrücken

(42)

a) Binäre Operatoren

```
Type * checkOp ( Absyn * node, Table * symTable ) {
    Type * leftType, * rightType, * type;
    leftType = checkNode ( node -> u.opNode.left, symTable );
    rightType = " " " " " .right, " " );
    if ( leftType != rightType ) error (...);
    switch ( node -> u.opNode.op ) {
        :
        case ABSYN_OP_ADD :
            if ( leftType != intType ) error (...);
            type = intType;
            break;
        :
    }
    return type;
}
```

b) Einfache Variable

```
Type * checkSimpleVar ( Absyn * node, Table * symTable ) {
    Entry * entry;
    entry = lookup ( symTable, node -> u.simplevarNode.name );
    if ( entry == NULL ) error (...);
    if ( entry -> kind != ENTRY_KIND_VAR ) error (...);
    return entry -> u.varEntry.type;
}
```

c) Array-Variable

(43)

- überprüfen, daß wirklich ein Array indiziert wird
- überprüfen, daß Index von Typ Integer ist
- Basistyp zurückgeben! Bsp.: type A = array [5] of int;
var a: A; \Rightarrow Typ von a ist A, Typ von a[z] ist int

Bem.: Es ist eine gute Idee, bei Ausdrücken den Typ in der Abstrakten Syntax als Attribut festzuhalten.

5.3.2. Type-Checking von Anweisungen

Liefern grundsätzlich den Typ "void" zurück (kann bei uns durch NULL repräsentiert werden).

a) Zuweisungen

- überprüfen, daß Typ(lhs) und Typ(rhs) "zuweisungskompatibel" (bei uns: gleich) sind

b) If- und While-Statements

- überprüfen, daß die Tests Boole'sche Werte liefern
- Rekursiven Abstieg in die Teilanweisungen (then-Teil, else-Teil, Schleifenrumpf) nicht vergessen!

c) Prozeduraufrufe

- Gibt's den Bezeichner? Ist es eine Prozedur?
- Stimmen die aktuellen Argumente in Zahl und Typ mit den formalen Parametern überein?
- Ist das aktuelle Argument eines Referenzparameters eine Variable? (Andere Ausdrücke sind hier verboten!)

5.3.3. Type-Checking von Typen

(44)

Das ist einfach: entspr. Typgraphen zurückgeben!

a) Benannte Typen:

- In Symboltabelle nachschauen (vorhanden?, Typ?)
- Typ zurückgeben

b) Array-Typen:

- Basistyp ermitteln (Rekursion!)
- Array-Typ konstruieren und zurückgeben

5.3.4. Type-Checking von Deklarationen

Deklarationen bewirken Einträge in der Symboltabelle.

a) Typ-Deklarationen

```
Type * checkTypedecl (Absyn * node, Table * symTable) {
    Type * type;
    Entry * entry;
    type = checkNode (node->u.typedeclNode.type, symTable);
    entry = new TypeEntry (type);
    if (enter (symTable, node->u.typedeclNode.name,
              entry) == NULL) error (...);
    return NULL;
}
```

b) Variablen-Deklarationen

Wie a), aber mit new VarEntry (type, FALSE)

($\hat{=}$ keine Referenz-Variablen)

c) Prozedur-Deklarationen

45

- Parameter-Typen ermitteln
- Prozedur in globale Symboltabelle eintragen
- Parameter in lokale Symboltabelle eintragen
- lokale Deklarationen bearbeiten
- Prozedur-Rumpf bearbeiten
- ggf. Symboltabellen ausgeben (Debugging!)

5.3.5. Vorwärtsreferenzen und Rekursion

Vorwärtsreferenz: Benutzung eines Bezeichners vor seiner Deklaration. Das ist unvermeidlich der Fall bei wechselseitig rekursiven Typen oder Prozeduren.

Schwierigkeit: Information über den Bezeichner wird benötigt, ist aber (noch) nicht vorhanden.

Lösung = Type-Checking in zwei Durchläufen.

1. Durchlauf: Sammeln von Informationen über Bezeichner aus den "Köpfen"; keine Behandlung der Benutzung von Bezeichnern in den "Rümpfen".

2. Durchlauf: Jetzt sind alle Bezeichner in der Symboltabelle eingetragen → Type-Checking der "Rümpfe".

Bsp: Prozedur-Deklarationen

1. Durchlauf: Type-Checking der Parameter und Eintrag der Prozedur in die globale Symboltabelle (Rumpf der Prozedur wird ignoriert)

2. Durchlauf: alle anderen o.g. Aktionen

(46)

Bem.: Bei der Typprüfung von rekursiven Typdefinitionen geht man genauso vor. Meist wird aber zusätzlich das Aufdecken von illegalen Zyklen verlangt. Bsp.:

type a = b; type b = c; type c = a;

5.3.6. Weitere Aufgaben des Type-Checkers

- Erzeugen der Typgraphen für die primitiven und die "eingebauten" Typen (Unterschied?)
- Erzeugen der globalen Symboltabelle
- Eintragen der vordefinierten Bezeichner (Typen, Prozeduren, etc.) in die globale Symboltabelle
- Überprüfen globaler semantischer Bedingungen (z.B. "ist main() vorhanden?")
- ggf. Ausgabe der ~~globalen~~ Symboltabelle(n)

6. Laufzeit - Organisation

(47)

Jeder Funktions / Prozeduraufruf hat seinen eigenen Satz von lokalen Variablen (und Parametern). Viele solcher Aufrufe können zur gleichen Zeit existieren (Rekursion!).

In vielen Sprachen (einschließlich SPL) sollen die Speicherplätze für lokale Variablen (und Parameter) automatisch freigegeben werden, wenn die Ausführung die Funktion verlässt \Rightarrow Stack ist geeignete Datenstruktur hierfür.

Bem: In Sprachen, die geschachtelte Funktionsdefinitionen und Funktionen als Rückgabewerte zulassen (sog. "Funktionen höherer Ordnung") geht das nicht!

Dort leben lokale Variablen u.U. länger, als der Funktionsaufruf, der sie erzeugt. Bsp.: Funktionsgenerator

```
int (*f(int x))(int) {  
    int g(int y) { return x+y; }  
    return g;  
}
```

```
int (*h)(int) = f(3); // h == "Erhöhen um 3"
```

```
int (*j)(int) = f(4); // j == "Erhöhen um 4"
```

```
int z = h(5); // z == 8
```

```
int w = j(7); // w == 11
```

Offensichtlich müssen die Speicherplätze für 3 und 4 über den Aufruf von f hinaus erhalten bleiben! \Rightarrow Stack ungeeignet!

6.1. Stack-Frames

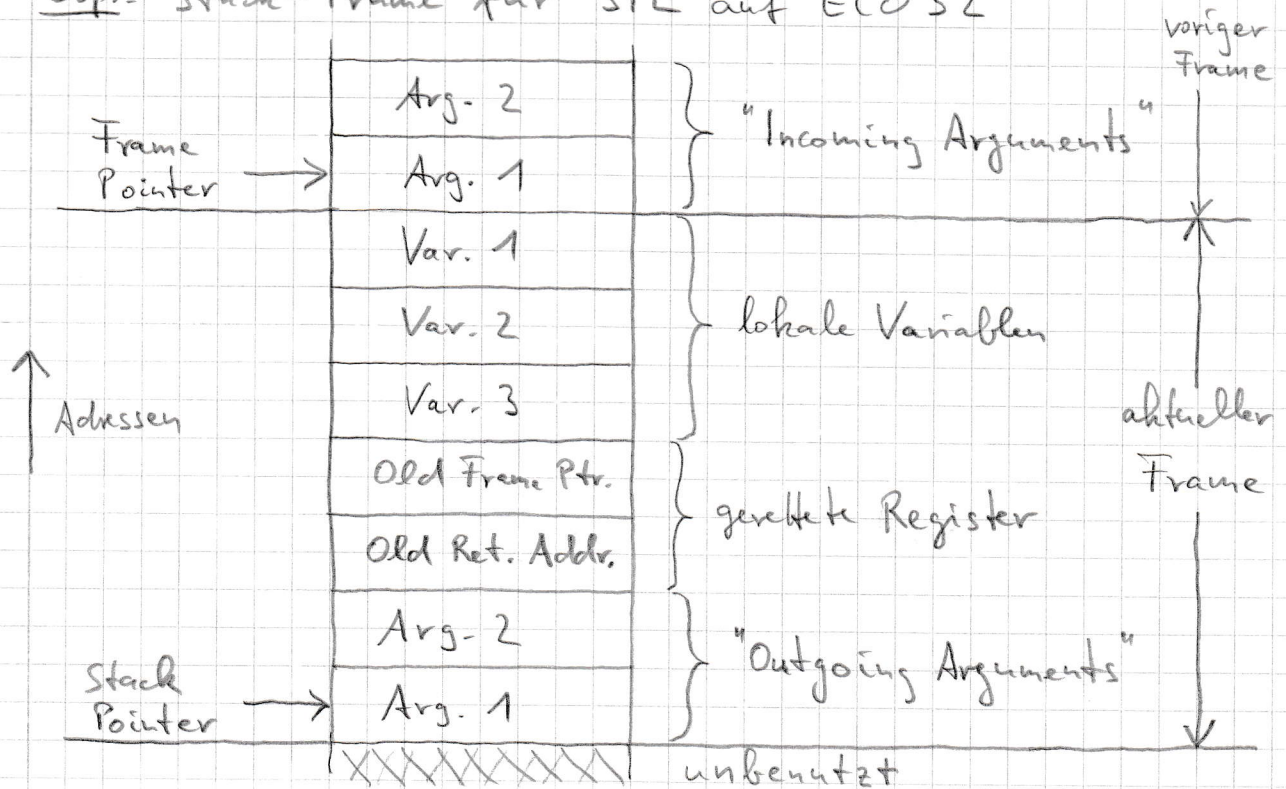
(48)

Variablen werden in größerer Zahl angelegt / freigegeben (nicht einzeln) \Rightarrow push/pop von ganzen Blöcken ("Stack-Frame", "Activation Record"). Die Struktur eines Stack-Frames wird typischerweise nicht von der Hardware vorgegeben, sondern durch Konventionen vereinbart (wichtig für das Binden von Programmen, die in verschiedenen Progr. Sprachen geschrieben sind).

Bestandteile:

- lokale Variablen
- Return-Adresse
- Hilfsvariablen
- geregelte Register
- Argumente f. aufzurufende Prozeduren

Bsp.: Stack-Frame für SPL auf ECO32



- Bem.: a) Zur Vereinfachung werden alle lokalen Variablen im Stack-Frame gespeichert (keine in Registern) und alle Hilfsvariablen in Registern (keine im Frame).
- b) Der Speicher unterhalb des Stack-Pointers darf nicht benutzt werden (Interrupts können dort hinein schreiben).
- c) Die Speicherplätze für "Old Ret. Addr." und "Outgoing Args" werden nur angelegt, falls die aktuelle Prozedur eine andere aufrufen kann (kann - denn ob sie es wirklich tut, ist zur Übersetzungszeit nicht berechenbar).

6.2. Der Frame-Pointer, Prozedur-Eintritt/Austritt

Der Frame-Pointer dient zum Adressieren von Argumenten und lokalen Variablen. Er muß bei Eintritt in eine Prozedur gesetzt und bei Austritt restauriert werden:

proc: $sp \leftarrow sp - framesize$

$mem[sp + offset(old fp)] \leftarrow fp$

$fp \leftarrow sp + framesize$

⋮

(Hier können Argumente mit $fp + offset(\text{Argument})$

und lokale Variablen mit $fp - offset(\text{Variable})$

adressiert werden. Argumente für gerufene

Prozeduren werden an $sp + offset(\text{Argument})$ gelegt.)

⋮

$fp \leftarrow mem[sp + offset(old fp)]$

$sp \leftarrow sp + framesize$

return

Zum Zeitpunkt der Codegenerierung werden also

- benötigt:
 - framesize
 - offset (old fp)
 - offset (Argument) f. alle eingehenden Argumente
 - offset (Variable) f. alle lokalen Variablen
 - offset (Argument) f. alle ausgehenden Argumente

Bem.: Der Frame-Pointer muß nicht tatsächlich in einem Register gespeichert werden: $fp \equiv sp + framesize$. Alle Zugriffe auf Variablen, Parameter etc. können über den sp erfolgen. Warum gibt's dann überhaupt den fp ? Der framesize ist normalerweise erst spät im Übersetzungsprozess bekannt - nachdem man weiß, wieviel Hilfsvariablen gebraucht werden. Deshalb in jedem Falle zunächst: Zugriff auf Variablen etc. über (evtl. fiktiven) fp .
 Bei SPL auf ECO32: realer fp in einem Register.

6.3. Berechnungen zur Laufzeitorganisation

1. Jeder Typ (nicht nur benannte!) speichert in seinem Typgraphen seine Größe, d.h. die Anzahl Bytes, die ein Objekt von diesem Typ im Speicher belegen wird. Das ist notwendig, um die Offsets (siehe 2.) berechnen zu können. Diese Information entsteht bei der Typgraphen-Konstruktion.

2. Jedes Argument, jeder Parameter und jede lokale Variable speichert den zugehörigen Offset relativ zum Frame-Pointer. Diese Information wird bei einem Durchgang durch die Prozedurdeklarationen gewonnen.

- 3. Jede Prozedur speichert die Größe von 3 Bereichen:
 - a) eingehende Argumente
 - b) lokale Variable
 - c) ausgehende Argumente

a) und b) werden bei der Berechnung oben (siehe 2.) mit ausgerechnet.

c) wird am besten in einem separaten Durchlauf ermittelt: Für jede Prozedur p bestimme das Maximum der eingehenden Argumente der Prozeduren, die p ruft. Dieser Wert ist die Größe des Bereichs der ausgehenden Argumente für p.

Bem.: Hierbei wird auch bemerkt, ob p überhaupt eine weitere Prozedur aufruft. Diese Information wird ebenfalls festgehalten (bei der Codeerzeugung: Sichern der Return-Adresse sollte ansonsten entfallen).

7. Codegenerierung

(52)

7.1. Die Zielmaschine ECO32

a) Allgemeines:

- 32-Bit RISC-Architektur, big-endian
- 32 32-Bit Register $\$0 \dots \31 , $\$0 \equiv \phi$ durch Hardware
- 64 Instruktionen (wir brauchen ca. 20)
- Kernel/User-Mode (wir benutzen nur Kernel-Mode)
- Speicher: Wort-organisiert, Byte-adressiert, max. 1 GByte, Zugriffe müssen ausgerichtet sein
- I/O: Timer, Terminal, Disk, Grafik-Controller, Memory-mapped, Interrupt- und DMA-fähig
(Bedienung wird von den Bibliotheksprozeduren geleistet)
- Virtuelle Adressierung, 4 KByte Seitengröße, max. 4 GB, TLB-Unterstützung (das benutzen wir nicht: alle Adressen ab $\phi \times c \phi \phi \phi \phi \phi \phi$ werden durch die Hardware auf physikal. Adressen ab ϕ abgebildet)

b) Instruktionen:

add $\$4, \$17, \$8$; $\$4 \leftarrow \$17 + \$8$

add $\$4, \$17, 8$; $\$4 \leftarrow \$17 + 8$

entspr. auch für sub, mul, div

ldw $\$4, \$17, 8$; $\$4 \leftarrow \text{mem}[\$17 + 8]$

stw $\$4, \$17, 8$; $\text{mem}[\$17 + 8] \leftarrow \4

; $\$4 \rightarrow \text{mem}[\$17 + 8]$

hier: ; das ist ein Label (53)

beq \$4, \$17, hier ; if \$4 = \$17: weiter bei "hier"
entspr. auch für bne, bgt, bge, blt, ble

(sowie für vorzeichenlose Größen: bgtu, bgeu, bltu, bleu)

j hier ; unbedingter Sprung zu "hier"

jal hier ; Prozeduraufruf der Proz. "hier",
Sichern der Ret. Addr. in \$31

jr \$31 ; weiter bei Inhalt v. \$31, "return"

c) Pseudo-Instruktionen (= Assembler-Direktiven):

.export hier ; lokales Label "hier" wird außerhalb
des Moduls verfügbar

.import dort ; externes Label "dort" wird innerhalb
des Moduls verfügbar

.code ; die im folgenden erzeugten Bytes
kommen ins Code-Segment

.align 4 ; Einstellen der Adresse auf einen
durch 4 teilbaren Wert ("Alignment")

d) Register-Konventionen

\$0 immer 0 (Hardware)

\$31 Return-Adresse (Hardware)

\$25 Frame-Pointer

\$29 Stack-Pointer

\$8..\$23 Hilfsvariable

Alle anderen

Register dürfen

nicht benutzt

werden!

7.2. Code für Stack-Maschinen

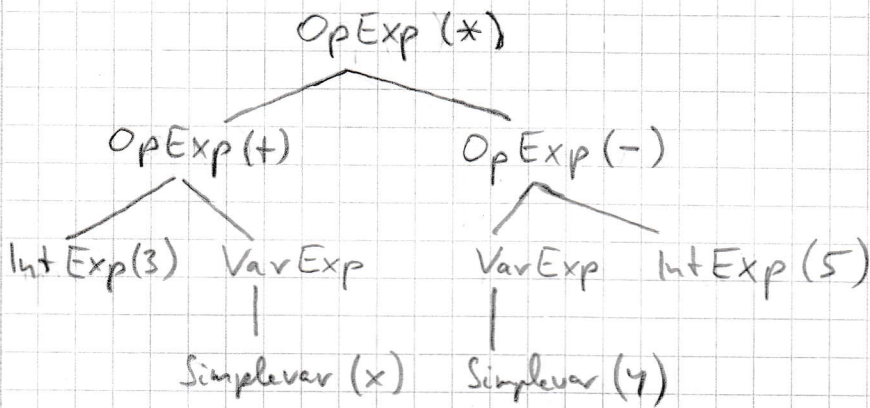
54

Vorteil: einfach zu erzeugen

Nachteil: ineffizienter Code

Methode: Post-Order-Durchgang durch Abstrakte Syntax

Bsp.: $(3+x) * (y-5)$ ← Depth-First-



push (3)

push (Adr(x))

load

add

push (Adr(y))

load

push (5)

sub

mul

7.3. Code für Register-Maschinen

Wie oben, aber mit Registern anstelle des Stacks:

die Stack-Positionen werden jetzt zu Register-Nummern.

Bsp.: Codeerzeugung für binäre Operationen

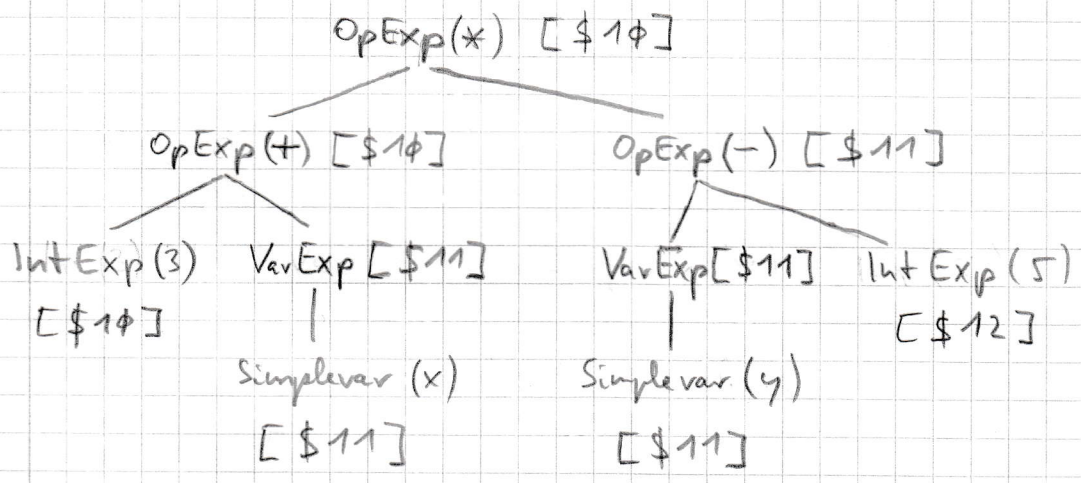
/ erzeuge Code, so daß Wert von "node" in "target" erscheint */*

*void genCodeOp (Absyn *node, int target) {*

```

genCode (node → left, target);
temp = target + 1;
if (temp > maxreg) error ("Ausdruck zu kompliziert");
genCode (node → right, temp);
output ("add %d, %d, %d \n",
        target, target, temp);
}

```



```

add $10, $0, 3
add $11, $25, offset(x)
ldw $11, $11, 0
add $10, $10, $11
add $11, $25, offset(y)
ldw $11, $11, 0
add $12, $0, 5
sub $11, $11, $12
mul $10, $10, $11

```

} kann zusammengefasst werden:
 ldw \$11, \$25, offset(x)

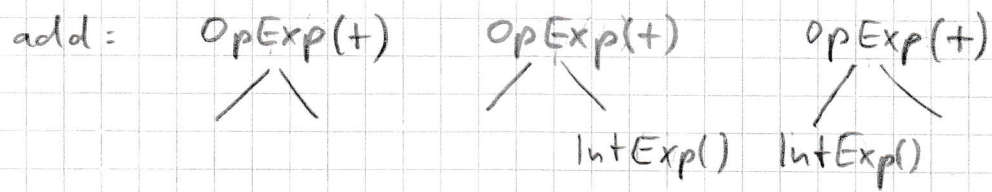
} entspr.

} sub \$11, \$11, 5

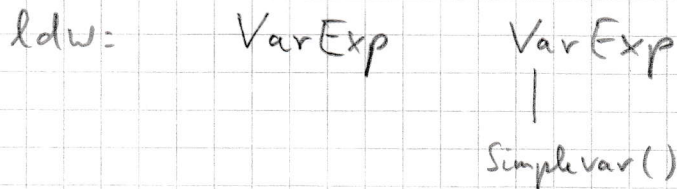
Bem.:

a) Das "Zusammenfassen" von Instruktionen kann folgendermaßen erreicht werden (Algorithmus heißt "Maximal Munch"):

Jede Instruktion der Zielmaschine erzeugt eine 56 oder mehrere "Kacheln", je nach möglichen Argumenten:



IntExp() Simplevar()



Dann wird der ganze Baum der Abstrakten Syntax mit Kacheln belegt, ausgehend von der Wurzel, mit möglichst großen Kacheln. Das ist nicht immer das Optimum!

b) Eine Verbesserung läßt sich mit "Dynamischer Programmierung" erzielen.^(*) Weiterer Vorteil: Algorithmus kann in ein Tool integriert werden ("Codegenerator-Generator"). Das diese Optimierung beim Bau des Compilers durchgeführt

c) Der o.a. Algorithmus zur Allokation von Registern ist nicht immer optimal (\rightarrow "Lathi-Ullman-Algorithmus"; behandelt auch das Auslagern / Wiedereinlagern von Werten für temporäre Variable in den Stackframe, falls nicht genügend Register verfügbar sind).
, das sog. "Spilling"

(*) Das ist ein Algorithmus, der die optimale Lösung eines Problems aus den optimalen Lösungen der Subprobleme berechnet.

7.4. Einige Details zur Codeerzeugung

(57)

a) Array-Variable

- rekursive Berechnung der (Adresse der) Variablen

- rekursive Berechnung des Index

- Überprüfung, ob $0 \leq \text{Index} < \text{Feldgröße}$

Hinweis: bei Zweierkomplement-Darstellung kommt man mit einem Vergleich vorzeichenloser Größen aus!

- Adresse der Array-Variablen = Adresse der Variablen
+ Index * Datentypgröße des Basistyps

b) While-Statement

Bsp.: while ($e_1 < e_2$) stm

Zwei mögliche Übersetzungen:

$\alpha)$ L123:

<code für e_1 in r_1 >

<code für e_2 in r_2 >

bge $r_1, r_2, L124$

<code für stm>

j L123

L124:

} 2 Sprünge /
Durchlauf

$\beta)$

j L123

L124:

<code für stm>

L123:

<code für e_1 in r_1 >

<code für e_2 in r_2 >

b.lt $r_1, r_2, L124$

} 1 Sprung /
Durchlauf

Bem. L123 und L124 sind "generierte Marken".

58

Markengenerator:

```
int newLabel(void) {  
    static int numLabels = 0;  
    return numLabels++;  
}
```

c) If-Statement

d) einarmig, Bsp.: if ($e_1 < e_2$) stm

<code für e_1 in r_1 >

<code für e_2 in r_2 >

bge $r_1, r_2, L123$

<code für stm>

L123:

β) zweiarmig, Bsp.: if ($e_1 < e_2$) stm₁ else stm₂

<code für e_1 in r_1 >

<code für e_2 in r_2 >

bge $r_1, r_2, L123$

<code für stm₁>

̄ L124

L123:

<code für stm₂>

L124:

d) Prozedur-Aufruf

- Argumente rekursiv berechnen und abspeichern
- Prozedur aufrufen

e) Prozedur-Deklaration

(59)

- Framegröße berechnen:

ruft diese Prozedur andere Prozeduren?

nein: Framegröße = Platz f. lokale Variable + 4

(wegen Framepointer - Speicherplatz)

ja: Framegröße = Platz f. lokale Variable + 8 +

Platz f. ausgehende Argumente

(wegen Framepointer - und Returnadressen - Speicherplatz)

- Prozedur-Prolog ausgeben; siehe 6.2

(Label, Stackframe allozieren, fp und ggf. Ret. Addr. sichern, neuen fp etablieren)

- Code für Prozedurkörper erzeugen

- Prozedur-Epilog ausgeben; siehe 6.2

(fp und ggf. Ret. Addr. wiederherstellen, Stackframe freigeben)

- Prozedur verlassen

f) Referenzvariable

Gibt's bei uns nur in Form von Referenz-Parametern

Beim Prozeduraufruf: Übermittlung der Adresse
anstatt des Wertes

Bei Verwendung: automatische Dereferenzierung