

Compilerkonstruktion

Wintersemester 2015/16

Prof. Dr. R. Parchmann

17. November 2015

Zwischencode-Erzeugung

Warum verwendet man einen Zwischencode und erzeugt nicht direkt ein Programm in der Zielsprache?

Zwischencode-Erzeugung

Warum verwendet man einen Zwischencode und erzeugt nicht direkt ein Programm in der Zielsprache?

1. Leichteres Ändern des Compilers auf eine neue Zielsprache - im Idealfall müsste nur der Maschinencode-Erzeuger geändert werden (das sogenannte *Back-End* des Compilers).

Zwischencode-Erzeugung

Warum verwendet man einen Zwischencode und erzeugt nicht direkt ein Programm in der Zielsprache?

1. Leichteres Ändern des Compilers auf eine neue Zielsprache - im Idealfall müsste nur der Maschinencode-Erzeuger geändert werden (das sogenannte *Back-End* des Compilers).
2. Die Code-Optimierung ist weitestgehend unabhängig von der Zielsprache und kann somit für mehrere Compiler benutzt werden.

Zwischencode-Erzeugung

Warum verwendet man einen Zwischencode und erzeugt nicht direkt ein Programm in der Zielsprache?

1. Leichteres Ändern des Compilers auf eine neue Zielsprache - im Idealfall müsste nur der Maschinencode-Erzeuger geändert werden (das sogenannte *Back-End* des Compilers).
2. Die Code-Optimierung ist weitestgehend unabhängig von der Zielsprache und kann somit für mehrere Compiler benutzt werden.
3. Es ist einfacher, eine Übersetzung in zwei (oder mehr) kleineren Schritten als in einem großen Schritt durchzuführen.

Form des Zwischencodes

Man unterscheidet die folgenden Hauptkategorien von Formen des Zwischencodes:

Form des Zwischencodes

Man unterscheidet die folgenden Hauptkategorien von Formen des Zwischencodes:

- ▶ HIL - High-level Intermediate Languages werden nur in den ersten Stufen des Übersetzungsprozesses benutzt

Form des Zwischencodes

Man unterscheidet die folgenden Hauptkategorien von Formen des Zwischencodes:

- ▶ HIL - High-level Intermediate Languages werden nur in den ersten Stufen des Übersetzungsprozesses benutzt
- ▶ MIL - Medium-level Intermediate Languages spiegeln die Möglichkeiten der Quellsprache in sehr vereinfachter Form

Form des Zwischencodes

Man unterscheidet die folgenden Hauptkategorien von Formen des Zwischencodes:

- ▶ HIL - High-level Intermediate Languages werden nur in den ersten Stufen des Übersetzungsprozesses benutzt
- ▶ MIL - Medium-level Intermediate Languages spiegeln die Möglichkeiten der Quellsprache in sehr vereinfachter Form
- ▶ LIL - Low-level Intermediate Languages sind nahe an der Zielsprache.

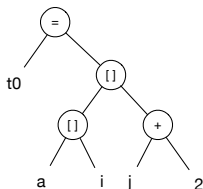
Beispiel

a sei deklariert als float a[20][10] mit unterer Feldgrenze 0
Der Ausdruck a[i][j+2] würde übersetzt

HIL: in einen Ausdruck wie etwa

$$t0 = a[i, j+2]$$

mit t0 als temporären Namen oder in einen Syntaxbaum wie



MIL:

```
t1 = j + 2
t2 = i * 10
t3 = t1 + t2
t4 = 4 * t3      /* Float-Zahl = 4 Bytes */
t0 = a[t4]       /* a[t4] entspricht */
                 /*Anfangsadresse a + t4 */
```

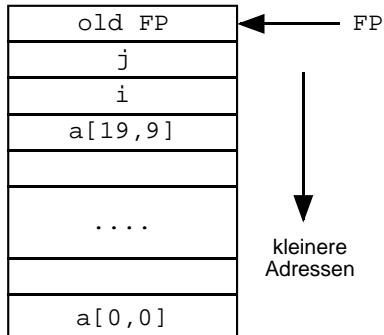
wobei die t_i temporäre Namen sind.

Es wird also der Ausdruck $(\text{addr } a) + 4 * (i * 10 + j + 2)$ berechnet. Über die Speicherung des Feldes a muss nur bekannt sein, dass sie zeilenweise erfolgt und dass jedes Feldelement 4 Bytes belegt.

LIL: Annahmen:

- ▶ das Feld `a` ist lokal
- ▶ Speicherplatz für `a` und für `i` und `j` ist im Aktivierungs-Record
- ▶ der Zugriff auf lokale Größen erfolgt über einen Framepointer `FP`
- ▶ die Objekte sind in der Reihenfolge `j`, `i` und dann `a` im Aktivierungs-Record
- ▶ jede Integergröße ist 2 Bytes und jedes Feldelement 4 Bytes lang.
- ▶ der Laufzeitstack wächst zu kleineren Adressen hin
- ▶ Die `ri` sind (fiktive) Integer-Register, die `fi` (fiktive) spezielle Floating-Point Register.

Aufbau des Aktivierungs-Records

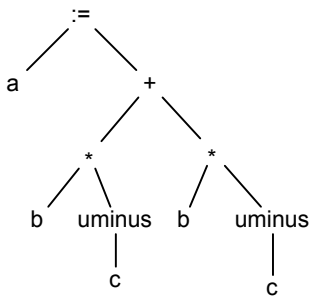


Äquivalentes Programm in einer LIL-Zwischensprache:

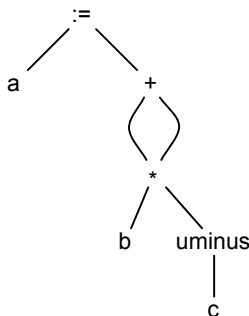
```
r1 <- [FP-2]      /* r1 enthält j */
r2 <- r1 + 2
r3 <- [FP-4]      /* r3 enthält i */
r4 <- r3 * 10
r5 <- r4 + r2
r6 <- 4 * r5
r7 <- FP - 804    /* Startadresse */
                  /* des Feldes a */
f1 <- [r7+r6]
```

Übersetzung von Wertzuweisungen in Syntaxbäume bzw. Syntaxgraphen

Gegeben sei der Ausdruck $a := b * -c + b * -c$, dann wäre der zugehörige Syntaxbaum bzw. Syntaxgraph:



Syntaxbaum



Syntaxgraph

Gegeben seien zwei Funktionen:

Gegeben seien zwei Funktionen:

- ▶ `mkleaf(id, id.place)` erzeugt ein Blatt mit Markierung `id` und einem Zeiger in die Symboltabelle auf den korrespondierenden Eintrag.

Gegeben seien zwei Funktionen:

- ▶ `mkleaf(id, id.place)` erzeugt ein Blatt mit Markierung `id` und einem Zeiger in die Symboltabelle auf den korrespondierenden Eintrag.
- ▶ `mknode(op, left, right)` erzeugt einen internen Knoten mit Markierung `op` und zwei Zeigern mit Werten `left` und `right`.

Gegeben seien zwei Funktionen:

- ▶ `mkleaf(id, id.place)` erzeugt ein Blatt mit Markierung `id` und einem Zeiger in die Symboltabelle auf den korrespondierenden Eintrag.
- ▶ `mknode(op, left, right)` erzeugt einen internen Knoten mit Markierung `op` und zwei Zeigern mit Werten `left` und `right`.
- ▶ Die Variante `mknode(op, link)` erzeugt einen internen Knoten mit nur einem Zeiger.

Alle Funktionen geben einen Zeiger auf den erzeugten Knoten zurück.

Attributierte Grammatik für die Übersetzung

nptr sei ein synthetisches Attribut eines jeden nichtterminalen Symbols, das auf den Wurzelknoten des zugehörigen Syntaxbaums verweist.

Attributierte Grammatik für die Übersetzung

nptr sei ein synthetisches Attribut eines jeden nichtterminalen Symbols, das auf den Wurzelknoten des zugehörigen Syntaxbaums verweist.

Ein SDTS für die Übersetzung wäre:

<i>Produktion</i>	<i>Semantische Regel</i>
$S \rightarrow id := E$	$S.nptr := \text{mknode}('assign', \text{mkleaf}(id, id.entry), E.nptr);$
$E \rightarrow E_1 + E_2$	$E.nptr := \text{mknode}('+', E_1.nptr, E_2.nptr);$
$E \rightarrow E_1 * E_2$	$E.nptr := \text{mknode}('*', E_1.nptr, E_2.nptr);$
$E \rightarrow - E_1$	$E.nptr := \text{mknode}('uminus', E_1.nptr);$
$E \rightarrow (E_1)$	$E.nptr := E_1.nptr;$
$E \rightarrow id$	$E.nptr := \text{mkleaf}(id, id.entry);$

Drei-Adress-Befehle

Der Drei-Adress-Code ist eine MIL-Zwischensprache.

Drei-Adress-Befehle:

- ▶ Wertzuweisungen der Form: $x := y \text{ op } z$,
- ▶ Wertzuweisungen der Form: $x := \text{op } y$,
- ▶ Wertzuweisungen der Form: $x := y$
- ▶ Wertzuweisungen der Form $x := y[i]$ und $x[i] := y$, dabei bedeutet $y[i]$: Adresse von $y + i$ Speichereinheiten
- ▶ Wertzuweisungen der Form: $x := \&y$ (Adresse von y), $x := *y$ (Dereferenzierung) oder $*x := y$
- ▶ Sprünge der Form: `goto L`, dabei ist L eine Marke. (Marken können vor jedem Drei-Adress-Befehl auftreten.)
- ▶ Bedingte Sprünge der Form: `if x relop y goto L`,

Zwei Arten von Sprungzielen: explizite Marken oder Indizes

Beispiel

Für die Anweisung `do i = i+1; while (a[i] < v);` sind die folgenden Übersetzungen in Drei-Adress-Code möglich:

```
L:   t1 = i + 1  
      i = t1  
      t2 = i * 8  
      t3 = a [ t2 ]  
      if t3 < v goto L
```

```
100: t1 = i + 1  
101: i = t1  
102: t2 = i * 8  
103: t3 = a [ t2 ]  
104: if t3 < v goto 100
```

Übersetzung von Wertzuweisungen mit einfachen arithmetischen Ausdrücken in Drei-Adress-Code.

Der Variablen E, die für einen arithmetischen Ausdruck steht, sind zwei Attribute zugeordnet:

1. E.place: Namen einer Variablen im Drei-Adress-Code, die den Wert des zu E korrespondierenden Ausdrucks enthält
 2. E.code: Folge von Drei-Adress-Befehlen, die den zu E korrespondierenden Ausdruck auswertet.
- ▶ `newtemp()` erzeugt bei jedem Aufruf einen neuen, temporären Variablennamen
 - ▶ `gen(...)` erzeugt einen Drei-Adress-Befehl.

<i>Produktion</i>	<i>Semantische Regel</i>
$S \rightarrow id := E$	$S.code := E.code \parallel gen(id.lexeme ':=' E.place);$
$E \rightarrow E_1 + E_2$	$E.place := newtemp();$ $E.code := E_1.code \parallel E_2.code \parallel$ $gen(E.place ':=' E_1.place + E_2.place);$
$E \rightarrow E_1 * E_2$	$E.place := newtemp();$ $E.code := E_1.code \parallel E_2.code \parallel$ $gen(E.place ':=' E_1.place * E_2.place);$
$E \rightarrow - E_1$	$E.place := newtemp();$ $E.code := E_1.code \parallel$ $gen(E.place ':=' 'uminus' E_1.place);$
$E \rightarrow (E_1)$	$E.place := E_1.place;$ $E.code := E_1.code;$
$E \rightarrow id$	$E.place := id.lexeme;$ $E.code := '';$

Bemerkung

*Man nennt Regeln, in denen die Zeichenkettenattribute der einzelnen nichtterminalen Symbole in der Reihenfolge ihres Auftretens in der rechten Seite der Produktion mit eventuell weiteren konstanten Zeichenketten konkateniert werden **einfach**. Man betrachte die Produktion $A \rightarrow BC$ mit zugeordneter Regel*

$A.code := string_B || B.code || string_C || C.code || string_D$

Diese Regel kann ersetzt werden durch

$A \rightarrow \{emit(string_B)\} B \{emit(string_C)\} C \{emit(string_D)\}$

wobei $emit(\dots)$ die als Parameter übergebene Zeichenkette ausgibt.

Die Funktion $emit$ ist also eine Funktion mit Seiteneffekten, d.h. man muss auf die Reihenfolge der Auswertung achten!

Damit lässt sich das SDTS umschreiben zu:

<i>Produktion</i>	<i>Semantische Regel</i>
$S \rightarrow id := E$	<code>emit(id.lexeme ':= ' E.place);</code>
$E \rightarrow E_1 + E_2$	<code>E.place := newtemp(); emit(E.place ':= ' E₁.place + E₂.place);</code>
$E \rightarrow E_1 * E_2$	<code>E.place := newtemp(); emit(E.place ':= ' E₁.place * E₂.place);</code>
$E \rightarrow - E_1$	<code>E.place := newtemp(); emit(E.place ':= ' 'uminus' E₁.place);</code>
$E \rightarrow (E_1)$	<code>E.place := E₁.place;</code>
$E \rightarrow id$	<code>E.place := id.lexeme;</code>

Darstellung von Drei-Adress-Befehlen

Quadrupel-Darstellung:

Ein Befehl wird durch einen Record mit 4 Feldern dargestellt, etwa

Operator	1. Operand	2. Operand	Ergebnis
----------	------------	------------	----------

also z.B. für den Befehl $x := y * z$ den Record

*	y	z	x
---	---	---	---

,

wobei x, y und z Zeiger zum korrespondierenden Eintrag in der Symboltabelle sind. Temporäre Namen müssen also bei dieser Methode in die Symboltabelle eingetragen werden.

Tripel-Darstellung:

Idee: das Ergebnis wird durch den Befehl selbst repräsentiert.

Ein Befehl wird durch einen Record mit 3 Feldern dargestellt, etwa

Operator	1. Operand	2. Operand
----------	------------	------------

also z.B. für den Befehl $x := y * z$ den Record

*	y	z
---	---	---

,

wobei y und z Zeiger zum korrespondierenden Eintrag in der Symboltabelle oder aber auf einen Drei-Adress-Befehl sind.

Beispiel

Betrachten wir wieder den Ausdruck $a := b * -c + b * -c$.
Äquivalente Folgen von Drei-Adress-Befehlen wären:

Quadrupel-Darstellung

Nr.	Op	1. Op	2. Op	Erg.
0	uminus	c		t1
1	*	b	t1	t2
2	uminus	c		t3
3	*	b	t3	t4
4	+	t2	t4	t5
5	:=	t5		a

Tripel-Darstellung

Nr.	Op	1. Op	2. Op
0	uminus	c	
1	*	b	(0)
2	uminus	c	
3	*	b	(2)
4	+	(1)	(3)
5	:=	a	(4)

wobei (i) ein Zeiger auf den Drei-Adress-Befehl darstellt, mit dem der entsprechende Operand berechnet wurde.

Syntax-gesteuerte Übersetzung von Deklarationen

SDTS zur Berechnung des Typs und des Offsets im Aktivierungs-Record für lokale Variablen.

<i>Produktion</i>	<i>Semantische Regel</i>
$P \rightarrow$	<code>offset := 0;</code>
$\quad D$	
$D \rightarrow D ; D$	
$D \rightarrow \text{id} : T$	<code>enter (id.name, T.type, offset);</code> <code>offset := offset + T.width;</code>
$T \rightarrow \text{integer}$	<code>T.type := integer;</code> <code>T.width := 4;</code>
$T \rightarrow \text{real}$	<code>T.type := real;</code> <code>T.width := 8;</code>
$T \rightarrow \text{array } [\text{num}] \text{ of } T_1$	<code>T.type := array (num.val, T₁.type);</code> <code>T.width := num.val * T₁.width;</code>
$T \rightarrow \uparrow T_1$	<code>T.type := pointer (T₁.type);</code> <code>T.width := 4;</code>

Bemerkung

T hat zwei synthetische Attribute:

- 1. type enthält den korrespondierenden Typausdruck*
 - 2. width enthält die zur Speicherung eines korrespondierenden Objektes benötigte Anzahl Bytes.*
- ▶ Es wird angenommen, dass der untere Index eines Feldes immer 0 ist.
 - ▶ offset ist eine globale Größe.

Übersetzung von Feldzugriffen

Annahme: die Feldelemente werden fortlaufend abgespeichert und jedes Element hat eine feste Größe von w Bytes.

Eindimensionale Felder

Die Deklaration des Feldes sei `a: array[low..high] of`
Dann ergibt sich die Adresse des Feldelements `a[i]` zu

$$\text{base} + (i - \text{low}) * w,$$

wobei `base` die Anfangsadresse des Speicherbereichs für das Feld `a` und `low` der Index des ersten Elementes in `a` ist, d.h. `a[low]` beginnt auf Adresse `base`.

Bemerkung

Durch Umformung erhält man

$$\text{base} + (i - \text{low}) * w = i * w + \text{base} - \text{low} * w,$$

*wobei der zweite Summand $\text{base} - \text{low} * w$ eventuell bereits zur Übersetzungszeit bestimmt werden kann.*

Zweidimensionale Felder

Die Deklaration des Feldes sei

`a: array[low1..high1, low2..high2] of ...`

Es gibt zwei prinzipielle Möglichkeiten:

- ▶ zeilenweise Speicherung (row-major) oder
- ▶ spaltenweise Speicherung (column-major)

Annahme: zeilenweise Speicherung.

Dann ergibt sich die Adresse des Feldelements `a[i1, i2]` zu

$$\text{base} + ((i_1 - \text{low}_1) * n_2 + i_2 - \text{low}_2) * w$$

wobei $n_2 = \text{high}_2 - \text{low}_2 + 1$ (Anzahl der Spalten) ist.

Bemerkung

Auch hier kann man die Formel umformen in einen Summanden, der zur Laufzeit berechnet werden muss und einen, der eventuell bereits zur Übersetzungszeit bekannt ist.

Man erhält:

$$\begin{aligned} & \text{base} + ((i_1 - \text{low}_1) * n_2 + i_2 - \text{low}_2) * w = \\ & ((i_1 * n_2) + i_2) * w + (\text{base} - ((\text{low}_1 * n_2) + \text{low}_2) * w). \end{aligned}$$

Mehrdimensionale Felder

Die Deklaration des Feldes sei

`a: array[low1..high1, low2..high2, ..., lowk..highk] of`

Verallgemeinert man die zeilenweise Speicherung auf mehr als zwei Dimensionen, dann erhält man eine sequentielle Speicherung der Feldelemente in einer Form, dass der letzte Index am schnellsten variiert.

Die Adresse des Feldelements `a[i1, i2, ..., ik]` ist dann

$$\text{base} + ((\dots (i_1 - \text{low}_1) * n_2 + i_2 - \text{low}_2) * n_3 + \dots + i_k - \text{low}_k) * w$$

wobei $n_i = \text{high}_i - \text{low}_i + 1$ für $1 < i \leq k$ gilt.

Bemerkung

Durch Umformung erhält man einen Summanden, der zur Laufzeit berechnet werden muss und einen, der eventuell bereits zur Übersetzungszeit bekannt ist:

$$\begin{aligned} & \text{base} + ((\dots (i_1 - \text{low}_1) * n_2 + i_2 - \text{low}_2) * n_3 \\ & \quad + \dots + i_k - \text{low}_k) * w = \\ & ((\dots (i_1 * n_2) + i_2) * n_3 + \dots) * n_k + i_k) * w \\ & \quad + \text{base} - ((\dots (\text{low}_1 * n_2) + \text{low}_2) * n_3 + \dots) * n_k + \text{low}_k) * w. \end{aligned}$$

Bemerkung

Der zur Laufzeit zu berechnende erste Summand der Formel lässt sich wie folgt rekursiv bestimmen:

$$e_1 := i_1 \text{ und}$$

$$e_j := e_{j-1} * n_j + i_j \text{ für } 1 < j \leq k$$

$e_k * w$ ergibt dann den Wert des ersten Summanden.

Übersetzung von Feldzugriffen

Als erster Ansatz könnte folgende Grammatik dienen:

$$F \rightarrow \text{id } [L]$$
$$L \rightarrow L, E \mid E$$

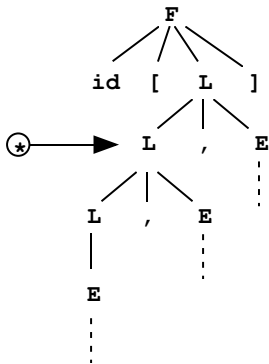
F steht dabei für ein Feldelement, L steht für eine Liste von Ausdrücken und E für einen arithmetischen Ausdruck..

F und L haben je zwei synthetische Attribute

- ▶ F.place enthält den Namen einer Variablen, die den Wert des zweiten, hier als zur Übersetzungszeit bekannt angenommenen Summanden der obigen Zugriffsformel enthält,
- ▶ F.offset enthält den Namen einer Variablen, die den Wert des ersten Summanden der obigen Zugriffsformel enthält.
- ▶ L.ndim enthält die Anzahl der arithmetischen Ausdrücke in L
- ▶ L.place gibt den Namen der Variablen an, die den Wert der Zugriffsformel für die Ausdrücke in L, also den Wert e_j enthält, wobei $j = L.ndim$ ist.
- ▶ `limit(array,j)` gibt für das Feld array den Wert $n_j = high_j - low_j + 1$ zurück.

Beispiel

Der Ableitungsbaum für den Zugriff auf $a[1,2,3]$. Probleme ergeben sich beim Auswerten der Attribute an der markierten Stelle.



Es gibt zwei Möglichkeiten, dieses Problem zu umgehen:

Es gibt zwei Möglichkeiten, dieses Problem zu umgehen:

1. Einführung eines inheriten Attributs für L und Herunterreichen des Feldnamens im Ableitungsbaum

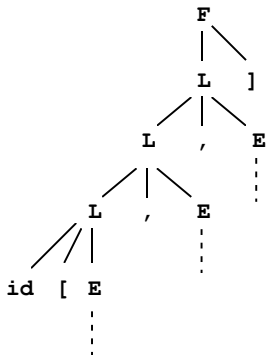
Es gibt zwei Möglichkeiten, dieses Problem zu umgehen:

1. Einführung eines inheriten Attributs für L und Herunterreichen des Feldnamens im Ableitungsbaum
2. Umformung der Grammatik, um bei einer S-attributierten Grammatik bleiben zu können.

Eine Möglichkeit wäre:

$$F \rightarrow L]$$
$$L \rightarrow L, E \mid \text{id} [E$$

Damit erhielte man folgenden Ableitungsbaum für das Beispiel:



L benötigt noch ein weiteres synthetische Attribut `L.array`, das den Namen des Feldes bzw. einen Zeiger auf den entsprechenden Eintrag in der Symboltabelle enthält.

SDTS zur Übersetzung von Feldzugriffen

<i>Produktion</i>	<i>Semantische Regel</i>
$S \rightarrow id := E$	<code>emit (id.lexeme ':' E.place);</code>
$S \rightarrow F := E$	<code>emit (F.place '[' F.offset ']' ':' E.place);</code>
$E \rightarrow E_1 + E_2$	<code>E.place := newtemp();</code> <code>emit (E.place ':' E_1.place + E_2.place);</code>
$E \rightarrow (E_1)$	<code>E.place := E_1.place;</code>
$E \rightarrow id$	<code>E.place := id.lexeme;</code>
$E \rightarrow F$	<code>E.place := newtemp();</code> <code>emit (E.place ':' F.place '[' F.offset ']);</code>

<i>Produktion</i>	<i>Semantische Regel</i>
$F \rightarrow L \]$	<pre> F.place := newtemp(); emit (F.place ':=' L.array '-' σ(L.array)); F.offset := newtemp(); emit (F.offset ':=' w '*' L.place); </pre>
$L \rightarrow L_1 \ , \ E$	<pre> t := newtemp(); j := L₁.ndim + 1; emit (t ':=' L₁.place '*' limit (L₁.array, j)); emit (t ':=' t '+' E.place); L.array := L₁.array; L.place := t; L.ndim := j; </pre>
$L \rightarrow id \ [\ E$	<pre> L.place := E.place; L.ndim := 1; L.array := id.lexeme; </pre>

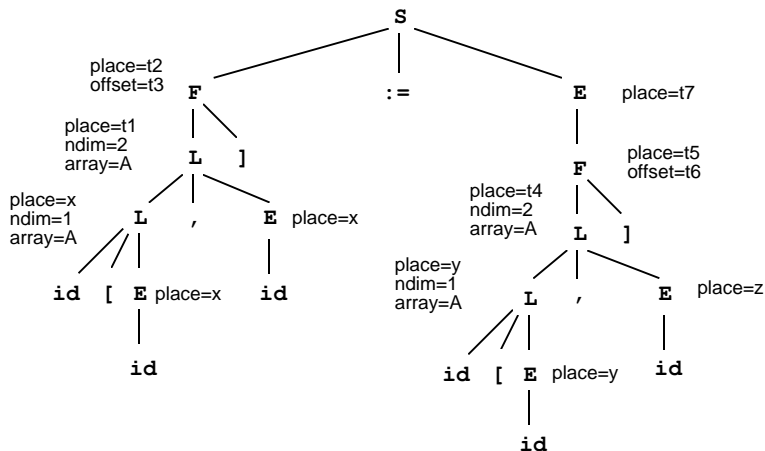
Beispiel

Sei A ein 10×20 -Feld, d.h. $A : \text{array}[1..10, 1..20] \text{ of } \dots$

Es ist also $n_1 = 10$ und $n_2 = 20$. und es gelte $w = 4$.

Übersetzt werden soll die Wertzuweisung $A[x, x] := A[y, z]$. Es ergibt sich folgender attributierter Ableitungsbaum:

Ableitungsbaum für das Beispiel:



Die erzeugten Drei-Adress-Befehle:

```
t1 := x * 20
t1 := t1 + x
t2 := A - 84
t3 := 4 * t1
t4 := y * 20
t4 := t4 + z
t5 := A -  $\sigma(A)$ 
t6 := 4 * t4
t7 := t5[t6]
t2[t3] := t7
```

wobei $\sigma(A) = ((low_1 * n_2) + low_2) * w$ den Wert 84 ergibt.