

**Bachelor-Programm**

# **Compilerbau**

im SoSe 2014

Prof. Dr. Joachim Fischer  
Dr. Klaus Ahrens  
Dipl.-Inf. Ingmar Eveslage

[fischer@informatik.hu-berlin.de](mailto:fischer@informatik.hu-berlin.de)



# Position

- ② **Teil I**  
Die Programmiersprache

- ② **Teil II**  
Methodische Grundlagen

- ② **Teil III**  
Entwicklung der Compiler

- ② **Kapitel 1**  
Compilationsprozess

- ② **Kapitel 2**  
Formalismen zur Sprachbeschreibung

- ② **Kapitel 3**  
Lexikalische Analyse: die Scanner

- ② **Kapitel 4**  
Syntaktische Analyse: die Parser

- ② **Kapitel 5**  
Parser-Generatoren: Yacc

- ② **Kapitel 6**  
Statische Semantikanalyse

- ② **Kapitel 7**  
Laufzeitsysteme

- ② **Kapitel 8**  
Ausblick: Codegenerierung

- ② **6.1**  
Überblick: Grammatik-basierte Übersetzung

- ② **6.2**  
Attributgrammatiken

- ② **6.3**  
S-attributierte Syntaxdefinitionen

- ② **6.4**  
Attributierte Syntaxdefinitionen mit synthetisierten und ererbten Attributen

- ② **6.5**  
L-attributierte Syntaxdefinitionen

- ② **6.6**  
Verfahren syntaxgesteuerter Übersetzungen im Überblick

- ② **6.7**  
Entwurf syntaxgesteuerter Übersetzungen

- ② **6.8**  
Drei-Adress-Code-Generierung (einige Aspekte)

- ② **6.9**  
Symboltabelle

Eingabe ... while ( C ) S1

S → WHILE '(' M C ')' N S1

# Umsetzung mit Bison

kritischer Fall: \$0, \$-1, \$-2, ...

Type-Cast

```

M → ε { L1:= newLabel();
        L2:= newLabel();
        C.false:= val[tos-3].next;
        C.true:= L2;
        }

```

```

m: { ... c_false= $<typ>-2.next; ... }

```

```

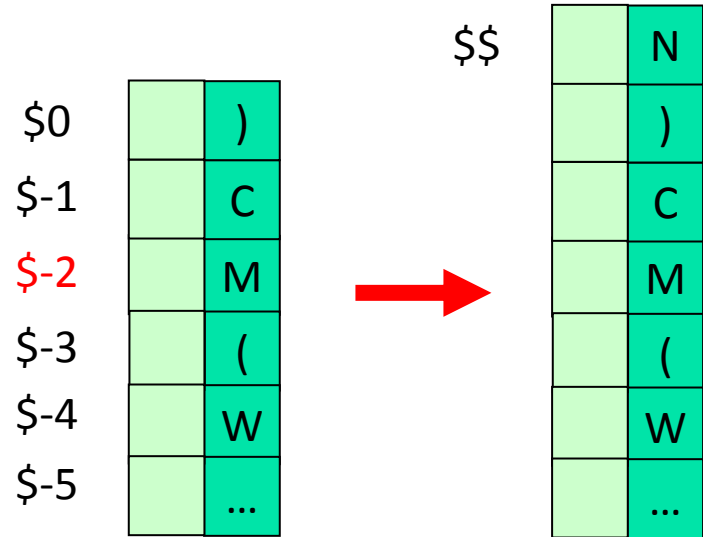
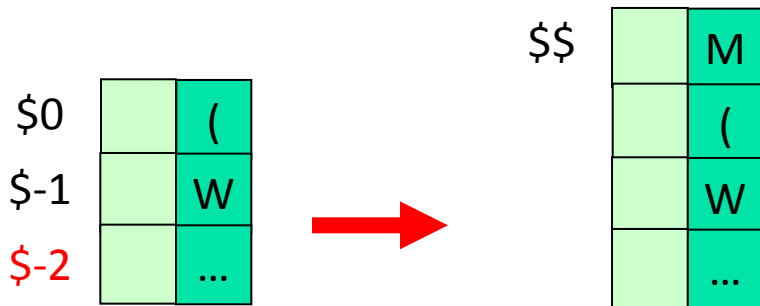
N → ε { S1.next= val[tos-3].L1;
        }

```

```

n: { $$S1_next= $<typ>-2.L1; }

```



**Achtung Fehlerquelle:**

- (1) allgemeine Kellerindizierung stimmt **nicht** konsequent mit Bison-Indizierung überein !
- (2) Zugriff auf Komponenten **außerhalb** der aktuellen Regel verlangt eine explizite Typwandlung (ohne Test)

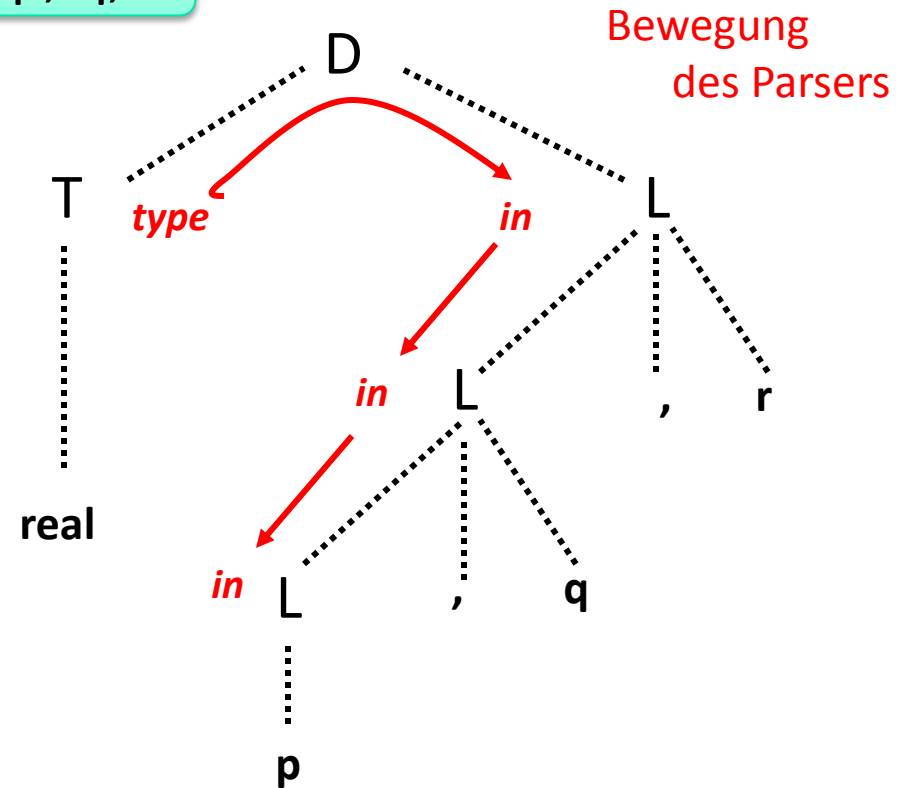
# Effizienzsteigerung von Attributberechnungen:

erneut: Variablendeklarations-Beispiel

Eingabe **real p, q, r**

Produktion	Semantische Regel
$D \rightarrow T L$	$L.in := T.type$
$T \rightarrow \text{int}$	$T.type := \text{integer}$
$T \rightarrow \text{real}$	$T.type := \text{real}$
$L \rightarrow L_1, \text{id}$	$L_1.in := L.in$ $\text{addtype}(\text{id.entry}, L.in)$
$L \rightarrow \text{id}$	$\text{addtype}(\text{id.entry}, L.in)$

Symboltabellen-Update



Welche Zustände(Symbole)  
kommen wann/wie auf den Ableitungsstack?

→ Kenntnis hilft, um auf zugriffs-synchrone Attributwert zugreifen zu können

# Auswertung ererbter Attribute

Eingabe	Zustand	benutzte Produktion, <b>semantische Regel</b>
real p, q, r	-	
p, q, r	real	
p, q, r	T	$T \rightarrow \text{real}$ <b>T.type := real</b>
, q, r	T p	
, q, r	T L	$L \rightarrow \text{id}$ <b>addtype(id.entry, L.in)</b>
q, r	T L,	
, r	T L, q	
, r	T L	$L \rightarrow L, \text{id}$ <b>L1.in = L.in; addtype(id.entry, L.in)</b>
r	T L,	
	T L, r	
	T L	$L \rightarrow L, \text{id}$ <b>L1.in = L.in; addtype(id.entry, L.in)</b>
	D	$D \rightarrow T L$ <b>L.in := T.type</b>

?

# Auswertung ererbter Attribute

Eingabe	Zustand	benutzte Produktion, semantische Regel
real p, q, r	-	
p, q, r	real	
p, q, r	T	$T \rightarrow \text{real}$ $T.\text{type} := \text{real}$
, q, r	T p	
, q, r	T L	$L \rightarrow \text{id}$ $\text{addtype}(\text{id.entry}, L.\text{in})$
q, r	T L,	
, r	T L, q	
, r	T L	$L \rightarrow L, \text{id}$ $L1.\text{in} = L.\text{in}; \text{addtype}(\text{id.entry}, L.\text{in})$
r	T L,	
	T L, r	
	T L	$L \rightarrow L, \text{id}$ $L1.\text{in} = L.\text{in}; \text{addtype}(\text{id.entry}, L.\text{in})$
	D	$D \rightarrow T L$ $L.\text{in} := T.\text{type}$

immer  
wenn zu **L** reduziert wird,  
steht **T** an  
einer bestimmten Position

➔ gesicherter Zugriff auf  
**T.type** möglich

(also direkter Zugriff!)

# Simulation ererbter Attribute: Beispiel

Eingabe	Keller	Produktion
real p, q, r	-	
p, q, r	real	
p, q, r	T	$T \rightarrow \text{real}$
, q, r	$T \overleftarrow{p}$	
, q, r	T L	$L \rightarrow \text{id}$
q, r	T L ,	
, r	$T \overleftarrow{L, q}$	
, r	T L	$L \rightarrow L, \text{id}$
r	T L ,	
	$T \overleftarrow{L, r}$	
	T L	$L \rightarrow L, \text{id}$
	D	$D \rightarrow T L$

## Beobachtung

- **T.type** ist immer auf festem Platz im Keller (bezüglich **tos**) zu finden
  - bei " $L \rightarrow L, \text{id}$ " ist dies **tos - 3**
  - bei " $L \rightarrow \text{id}$ " ist dies **tos - 1**

## Folgerung

- deshalb braucht **T.type** nicht nach **L.in** kopiert zu werden, sondern kann **direkt** vom Keller benutzt werden

Simulation ererbter Attribute erhöht i.allg. das Laufzeitverhalten eines Compilers

# Auswertung ererbter Attribute

## Beobachtung

- **T.type** ist immer auf festem Platz im Keller (bezüglich **tos**) zu finden
  - bei "**L** → **L**, **id**" ist dies **tos - 3**
  - bei "**L** → **id**" ist dies **tos - 1**

## Folgerung

- deshalb braucht **T.type** nicht nach **L.in** kopiert zu werden, sondern kann **direkt** vom Keller benutzt werden

Eingabe	Zustand	benutzte Produktion, semantische Aktion
real p, q, r	-	
p, q, r	real	
p, q, r	T	T → real    T.type := real
, q, r	T p	
, q, r	T L	L → id    addtype(id.entry, L.in)
q, r	T L,	
, r	T L, q	
, r	T L	L → L, id    L1.in = L.in; addtype(id.entry, L.in)
r	T L,	
	T L, r	
	T L	L → L, id    L1.in = L.in; addtype(id.entry, L.in)
	D	D → T L    L.in := T.type

Simulation ererbter Attribute erhöht i.allg. das Laufzeitverhalten eines Compilers



# Simulation ererbter Attribute: Beispiel

## angepasste semantische Aktionen

Produktion	Semantische Regel - Codefragment	Bison
$D \rightarrow T L$		
$T \rightarrow \text{int}$	<code>val[ntos] = integer</code>	<code>\$\$= integer</code>
$T \rightarrow \text{real}$	<code>val[ntos] = real</code>	<code>\$\$= real</code>
$L \rightarrow L_1, \text{id}$	<code>addtype (val[tos], val[tos-3])</code>	<code>addtype (\$3, \$&lt;typ&gt;-2)</code>
$L \rightarrow \text{id}$	<code>addtype (val[tos], val[tos-1])</code>	<code>addtype (\$1, \$&lt;typ&gt;0)</code>

tos = momentaner »top-of-stack«  
ntos = »top-of-stack« nach der Reduktion

Positionen  
in der aktuellen Regel  
positiv

Positionen  
vor der aktuellen Regel  
0, negativ

**Achtung:** Spezieller Zugriff auf Keller aber nur dann möglich,  
wenn Position (unabhängig von anderen Regeln) immer vorher bekannt ist

# Übersetzung arithmetischer Ausdrücke

## Idee:

1. dynamische Bereitstellung von temporären Adressen zur Speicherung von Zwischenergebnissen der Verkettung

Befehl: `newTemp()` // liefert neue temporäre Adresse

2. Meta-Symbol **E**

**2 synthetische Attribute**: `E.code` // 3-Adresscode für E

`E.addr` // Adresse, für Resultat der Ausführung von E.code

**Fazit**: S-Attributierte Grammatik

Produktion	Semantische Regel
$S \rightarrow id = E$	<code>S.code = E.code    top.get(id.lexem)    '='    E.addr</code>
$E \rightarrow E_1 + E_2$	<code>E.addr = newTemp()</code> <code>E.code = E_1.code    E_2.code    E.addr    '='    E_1.addr    '+'    E_2.addr</code>
$E \rightarrow - E_1$	<code>E.addr = newTemp()</code> <code>E.code = E_1.code    E.addr    '=' '-'    E_1.addr</code>
$E \rightarrow ( E_1 )$	<code>E.addr = E_1.addr</code> <code>E.code = E_1.code</code>
$E \rightarrow id$	<code>E.addr = top.get(id.lexem)</code> <code>E.code = ""</code>

# Übersetzung arithmetischer Ausdrücke

Produktion	Semantische Regel
$S \rightarrow id = E$	$S.code = E.code \parallel top.get(id.lexem) \parallel '=' \parallel E.addr$
$E \rightarrow E_1 + E_2$	$E.addr = \mathbf{new} Temp()$ $E.code = E_1.code \parallel E_2.code \parallel E.addr \parallel '=' \parallel E_1.addr \parallel '+' \parallel E_2.addr$
$E \rightarrow - E_1$	$E.addr = \mathbf{new} Temp()$ $E.code = E_1.code \parallel E.addr \parallel '=' '-' \parallel E_1.addr$
$E \rightarrow ( E_1 )$	$E.addr = E_1.addr$ $E.code = E_1.code$
$E \rightarrow id$	$E.addr = top.get(id.lexem)$ $E.code = ""$

# Übersetzung arithmetischer Ausdrücke

Produktion	Semantische Regel
$S \rightarrow id = E$	<code>S.code= E.code    top.get(id.lexem)    '='    E.addr</code>
$E \rightarrow E_1 + E_2$	<code>E.addr= new Temp() E.code= E_1.code    E_2.code    E.addr    '='    E_1.addr    '+'    E_2.addr</code>
$E \rightarrow - E_1$	<code>E.addr= new Temp() E.code= E_1.code    E.addr    '=' '-'    E_1.addr</code>
$E \rightarrow ( E_1 )$	<code>E.addr= E_1.addr E.code= E_1.code</code>
$E \rightarrow id$	<code>E.addr= top.get(id.lexem) E.code= ""</code>

# Übersetzung arithmetischer Ausdrücke

Produktion	Semantische Regel
$S \rightarrow id = E$	<code>S.code= E.code    top.get(id.lexem)    '='    E.addr</code>
$E \rightarrow E_1 + E_2$	<code>E.addr= new Temp() E.code= E_1.code    E_2.code    E.addr    '='    E_1.addr    '+'    E_2.addr</code>
$E \rightarrow - E_1$	<code>E.addr= new Temp() E.code= E_1.code    E.addr    '=' '-'    E_1.addr</code>
$E \rightarrow ( E_1 )$	<code>E.addr= E_1.addr E.code= E_1.code</code>
$E \rightarrow id$	<code>E.addr= top.get(id.lexem) E.code= ""</code>

# Übersetzung arithmetischer Ausdrücke

Produktion	Semantische Regel
$S \rightarrow id = E$	<code>S.code= E.code    top.get(id.lexem)    '='    E.addr</code>
$E \rightarrow E_1 + E_2$	<code>E.addr= <b>new</b> Temp() E.code= E_1.code    E_2.code    E.addr    '='    E_1.addr    '+'    E_2.addr</code>
$E \rightarrow - E_1$	<code>E.addr= <b>new</b> Temp() E.code= E_1.code    E.addr    '=' '-'    E_1.addr</code>
$E \rightarrow ( E_1 )$	<code>E.addr= E_1.addr E.code= E_1.code</code>
$E \rightarrow id$	<code>E.addr= top.get(id.lexem) E.code= ""</code>

# Übersetzung arithmetischer Ausdrücke

Produktion	Semantische Regel
$S \rightarrow id = E$	<code>S.code= E.code    top.get(id.lexem)    '='    E.addr</code>
$E \rightarrow E_1 + E_2$	<code>E.addr= new Temp() E.code= E_1.code    E_2.code    E.addr    '='    E_1.addr    '+'    E_2.addr</code>
$E \rightarrow - E_1$	<code>E.addr= new Temp() E.code= E_1.code    E.addr    '=' '-'    E_1.addr</code>
$E \rightarrow ( E_1 )$	<code>E.addr= E_1.addr E.code= E_1.code</code>
$E \rightarrow id$	<code>E.addr= top.get(id.lexem) E.code= ""</code>

# Übersetzung boolescher Ausdrücke

## Vorkommen:

1. Berechnung boolescher Werte (erfolgt analog zu arithmetischen Ausdrücken)
2. Steuerung des Kontrollflusses (als Bedingung in **if-then-else**, **while**, ...)

man erzeugt ebenfalls sog. **Sprung-Code**

**Metasymbol B** bekommt 3 Attribute

- **B.code** //synthetisch: 3-Adress-Code für B
- **B.l\_true** // ererbt: Marke, zu der B.code springt, wenn B-Code-Ausführung **true** liefert
- **B.l\_false** // ererbt: Marke, zu der B.code springt, wenn B-Code-Ausführung **false** liefert

**Frage:** Wo kommen **B.l\_true** und **B.l\_false** her?

- entweder von einem umfassenderen booleschen Ausdruck oder
- von der **Anweisung**, in der B als boolescher Ausdruck enthalten ist

**Anweisung S** hat zwei Attribute (code und Sprungmarke)



# Sprung-Code für boolsche Ausdrücke

Produktion	Semantische Regel
$B \rightarrow B_1 \parallel B_2$	$B_1.l\_true \parallel B.l\_true$ $B_1.l\_false = \mathbf{new} \text{ Label}()$ $B_2.l\_true = B.l\_true$ $B_2.l\_false = B.l\_false$ $B.code = B_1.code \parallel B_1.l\_false \parallel B_2.code$
$B \rightarrow B_1 \ \&\& \ B_2$	$B_1.l\_true = \mathbf{new} \text{ Label}()$ $B_1.l\_false = B.l\_false$ $B_2.l\_true = B.l\_true$ $B_2.l\_false = B.l\_false$ $B.code = B_1.code \parallel B_1.l\_true \parallel B_2.code$
$B \rightarrow ! B_1$	$B_1.l\_true = B.l\_false$ $B_1.l\_false = B.l\_true$ $B.code = B_1.code$
$B \rightarrow E_1 \ \mathbf{rel} \ E_2$	$B.code = E_1.code \parallel E_2.l\_true \parallel \mathit{gen}('if' \ E_1.addr \ \mathbf{rel}.op \ E_2.addr \ 'goto' \ B.l\_true)$ $\parallel \mathit{gen}('goto' \ B.l\_true)$
$B \rightarrow \mathbf{true}$	$B.code = \mathit{gen}('goto' \ B.l\_true)$
$B \rightarrow \mathbf{false}$	$B.code = \mathit{gen}('goto' \ B.l\_false)$

# Sprung-Code für boolesche Ausdrücke

Produktion	Semantische Regel
$B \rightarrow B_1 \parallel B_2$	$B_1.l\_true = B.l\_true$ $B_1.l\_false = \mathbf{new}$ Label() $B_2.l\_true = B.l\_true$ $B_2.l\_false = B.l\_false$ $B.code = B_1.code \parallel B_1.l\_false \parallel B_2.code$
$B \rightarrow B_1 \ \&\& \ B_2$	<p>Wenn <math>B_1 == \mathbf{true}</math>, dann <math>B = \mathbf{true}</math> (ohne dass <math>B_2</math> ausgewertet wird)  <math>B_1</math> und <math>B</math> haben in diesem Fall das gleiche Sprungziel</p> <p>Wenn <math>B_1 == \mathbf{false}</math>, muss auch <math>B_2</math> ausgewertet werden.  D.h. man muss an den Anfang von <math>B_2</math> springen.</p>
$B \rightarrow ! B_1$	<p>Wenn nun <math>B_2 == \mathbf{true}</math> (bzw. <math>B_2 == \mathbf{false}</math>), dann ist auch <math>B == \mathbf{true}</math> (bzw. <math>B == \mathbf{false}</math>) weil man <math>B_2</math> nur auswertet, wenn <math>B_1 == \mathbf{false}</math></p>
$B \rightarrow E_1 \ \mathbf{rel} \ E_2$	<p><u>Also</u>: beide Sprungziele werden einfach vererbt</p>
$B \rightarrow \mathbf{true}$	$B.l\_true = B.l\_true$
$B \rightarrow \mathbf{false}$	$B.l\_false = B.l\_false$

# Sprung-Code für boolsche Ausdrücke

Produktion	Semantische Regel
$B \rightarrow B_1 \parallel B_2$	$B_1.l\_true \parallel B.l\_true$ $B_1.l\_false = \mathbf{new} \text{ Label}()$ $B_2.l\_true = B.l\_true$ $B_2.l\_false = B.l\_false$ $B.code = B_1.code \parallel B_1.l\_false \parallel B_2.code$
$B \rightarrow B_1 \ \&\& \ B_2$	$B_1.l\_true = \mathbf{new} \text{ Label}()$ $B_1.l\_false = B.l\_false$ $B_2.l\_true = B.l\_true$ $B_2.l\_false = B.l\_false$ $B.code = B_1.code \parallel B_1.l\_true \parallel B_2.code$
$B \rightarrow ! B_1$	$B_1.l\_true = B.l\_false$ $B_1.l\_false = B.l\_true$ $B.code = B_1.code$
$B \rightarrow E_1 \ \mathbf{rel} \ E_2$	$B.code = E_1.code \parallel E_2.l\_true \parallel \mathit{gen}('if' \ E_1.addr \ \mathbf{rel.op} \ E_2.addr \ 'goto' \ B.l\_true)$ $\parallel \mathit{gen}('goto' \ B.l\_true)$
$B \rightarrow \mathbf{true}$	$B.code = \mathit{gen}('goto' \ B.l\_true)$
$B \rightarrow \mathbf{false}$	$B.code = \mathit{gen}('goto' \ B.l\_false)$

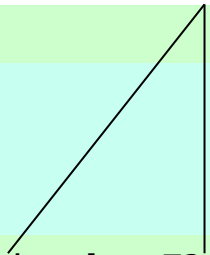
# Sprung-Code für boolsche Ausdrücke

Produktion	Semantische Regel
$B \rightarrow B_1 \parallel B_2$	$B_1.l\_true \parallel B.l\_true$ $B_1.l\_false = \mathbf{new} \text{Label}()$ $B_2.l\_true = B.l\_true$ $B_2.l\_false = B.l\_false$ $B.code = B_1.code \parallel B_1.l\_false \parallel B_2.code$
$B \rightarrow B_1 \&\& B_2$	$B_1.l\_true = \mathbf{new} \text{Label}()$ $B_1.l\_false = B.l\_false$ $B_2.l\_true = B.l\_true$ $B_2.l\_false = B.l\_false$ $B.code = B_1.code \parallel B_1.l\_true \parallel B_2.code$
$B \rightarrow ! B_1$	$B_1.l\_true = B.l\_false$ $B_1.l\_false = B.l\_true$ $B.code = B_1.code$
$B \rightarrow E_1 \text{ rel } E_2$	$B.code = E_1.code \parallel \text{gen}('goto' B.l\_true)$ <div style="border: 1px solid black; padding: 2px; display: inline-block;">B hat gleichen Code wie <math>B_1</math>, nur Sprungziele sind vertauscht</div>
$B \rightarrow \mathbf{true}$	$B.code = \text{gen}('goto' B.l\_true)$
$B \rightarrow \mathbf{false}$	$B.code = \text{gen}('goto' B.l\_false)$

# Sprung-Code für boolsche Ausdrücke

Produktion	Semantische Regel
$B \rightarrow B_1 \parallel B_2$	$B_1.l\_true \parallel B.l\_true$ $B_1.l\_false = \text{new Label}()$ $B_2.l\_true = B.l\_true$ $B_2.l\_false = B.l\_false$ $B.code = B_1.code \parallel B_1.l\_false \parallel B_2.code$
$B \rightarrow B_1 \&\& B_2$	$B_1.l\_true = \text{new Label}()$ $B_1.l\_false$ $B_2.l\_true =$ $B_2.l\_false$ $B.code = B_1.code \parallel B_1.l\_true \parallel B_2.code$
$B \rightarrow ! B_1$	$B_1.l\_true = B.l\_false$ $B_1.l\_false = B.l\_true$ $B.code = B_1.code$
$B \rightarrow E_1 \text{ rel } E_2$	$B.code = E_1.code \parallel E_2.l\_true \parallel \text{gen}('if' E_1.addr \text{ rel.op } E_2.addr 'goto' B.l\_true)$ $\parallel \text{gen}('goto' B.l\_true)$
$B \rightarrow \text{true}$	$B.code = \text{gen}('goto' B.l\_true)$
$B \rightarrow \text{false}$	$B.code = \text{gen}('goto' B.l\_false)$

$E_1.code$  und  $E_2.code$  liefern ihre Resultate in  $E_1.addr$  und  $E_2.addr$  ab.  
 Man muss die Inhalte beider Adressen vergleichen (mit  $rel.op$ )



# Position

- Teil I  
Die Programmiersprache

- Teil II  
Methodische Grundlagen

- Teil III  
Entwicklung der Compiler

- Kapitel 1  
Compilationsprozess

- Kapitel 2  
Formalismen zur Sprachbeschreibung

- Kapitel 3  
Lexikalische Analyse: die Scanner

- Kapitel 4  
Syntaktische Analyse: die Parser

- Kapitel 5  
Parser-Generatoren: Yacc

- Kapitel 6  
Statische Semantikanalyse

- Kapitel 7  
Laufzeitsysteme

- Kapitel 8  
Ausblick: Codegenerierung

- 6.1  
Überblick: Grammatik-basierte Übersetzung

- 6.2  
Attributgrammatiken

- 6.3  
S-attributierte Syntaxdefinitionen

- 6.4  
Attributierte Syntaxdefinitionen mit synthetisierten und ererbten Attributen

- 6.5  
L-attributierte Syntaxdefinitionen

- 6.6  
Verfahren syntaxgesteuerter Übersetzungen im Überblick

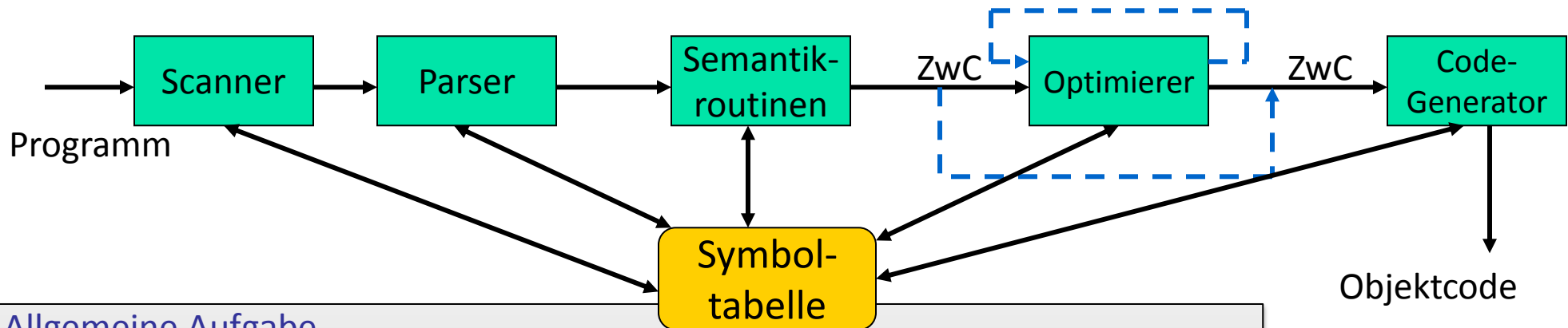
- 6.7  
Entwurf syntaxgesteuerter Übersetzungen

- 6.8  
Drei-Adress-Code-Generierung (einige Aspekte)

- 6.9  
Symboltabelle

- 6.10  
Typprüfung

# Einsatz der Symboltabelle



## Allgemeine Aufgabe

- Verbindung lexikalischer Namen (Symbole) mit ihren Attributen
- zentrales Hilfsmittel für Vereinbarung und Nutzung von Symbolen

## Was soll in der Symboltabelle gespeichert werden?

- Variablennamen
- ...
- Prozedur- und Funktionsnamen
- ...
- Sprungadressen (Labels)
- ...

**Achtung:** eine Symboltabelle ist eine Datenstruktur, die **nur** zu Übersetzungszeit existiert

# Rolle der Symboltabelle

Welche Art von Informationen braucht ein Compiler?

- Namen (Zeichenkette)
- Typ (Speicherplatzbedarf, Operatoranwendbarkeit)
- Dimensionsinformation (Felder)
- Prozedur-/ Funktionsdeklaration
- Anzahl und Typ der Argumente von Prozeduren/Funktionen
- **Gültigkeitsbereich von Deklarationen**  
(Blockstruktur → verschachtelte Gültigkeitsbereiche)
- Sichtbarkeitsbereich von Symbolen (Modul-, Klassenkonzept)
- Überladungsinfos für Operatoren
- Speicherklasse (static, extern, global, ...)
- Offset im Speicher
- falls Name eines Records: dann Strukturbeschreibung
- falls Parameter: call-by-value, call-by-reference, ...

## Beschränkung (Vorlesung):

blockstrukturierte  
Symboltabellen  
verwendet von  
Compilern für

- Pascal,
- C, C++,
- Modula,
- Java, ...



# Inhalt der Symboltabelle

---

Symboltabelle assoziiert Namen/Symbole  
mit Attributen (komplexe Werten)

- **Attribute** beschreiben zur Übersetzungszeit Eigenschaften einer Deklaration
- **Attribute** unterscheiden sich je nach der Bedeutung des Namens
  - **Variablennamen:** Typ, Prozedurebene, Speicherinfo (Rahmen)
  - **Typen:** Beschreibung des Typs, Größe und Speicheranforderungen: Alignment
  - **Konstanten:** Typ, Wert
  - **Prozeduren/Funktionen:** Formale Namen und Typen, Ergebnistyp, Speicherinfos, Rahmengröße

# Alternative Implementierung von Symboltabellen

## ■ lineare Liste

- Komplexität:  **$O(n)$**  Tests pro Suchoperation
- einfach zu erweitern, keine feste Größe
- ein neuer Eintrag pro neu eingefügtes Element

## ■ geordnetes lineares Feld

- Komplexität:  **$O(\log_2 n)$**  Tests pro Suchoperation mit binärer Suche
- Einfügen ist aufwendig, da Reihenfolge erhalten bleiben muss

## ■ Binärer Baum

- Komplexität:  **$O(n)$**  Tests pro Suchoperation bei nicht-balancierten Bäumen
- Komplexität:  **$O(\log_2 n)$**  Tests pro Suchoperation bei balancierten Bäumen
- einfach zu erweitern, keine feste Größe
- ein neuer Eintrag pro eingefügtes Element

## ■ Hash-Tabelle

- Komplexität:  **$O(1)$**  Tests pro Suchoperation (abhängig vom Füllstand)
- Einfügen ist unterschiedlich schwierig (abhängig vom gewählten Verfahren)

## Fall: Einpass-Compiler

### • Variante 1

Knoten, die Blockstrukturen eröffnen, erhalten individuelle Symboltabellen ihrer Blöcke (als Liste von Bezeichnern in Form von Hash-Tabellen)

### • Variante 2

Kellerspeicher mit Blockmarkierungen als eleganter Ansatz

# Kellerimplementation

---

## ▪ Aufbau

- Tabelle besteht aus **Blöcken**, die auf dem Stack liegen
- beim Betreten eines Blockes im Quell-Programm wird ein Block in der Symboltabelle angelegt, d.h. auf dem Keller angelegt
- beim Verlassen eines Blockes im Programm wird der letzte Block gelöscht
- Bezeichner werden nur im aktuellen Block eingetragen (auf bereits deklarierte Bezeichner wird nur im aktuellen Block geachtet)
- beim Suchen eines Bezeichners, werden die Blöcke ausgehend vom letzten zum ersten Eintrag hin untersucht, der erste Treffer wird gewählt
- Sonderfall: Prozedur-/Funktionsparameter, Rekords

## ▪ Feststellung

Keller ist am Ende des Syntaxanalyse leer  
(Kellerimplementierung ist ungeeignet für Mehrpass-Compiler)

# Blockstrukturierte Symboltabelle

## Verwaltung von **Gültigkeitsbereichen**!!

- wodurch entstehen Gültigkeitsbereiche?  
Definition eines Moduls/Programms, einer Funktion/Prozedur/ eines Blocks
- Anfragen zur Analysezeit nach einem Namen,
  - a) muss die in diesem Gültigkeitsbereich gültige Deklaration mit ihren Informationen zurückgegeben werden, **oder**
  - b) Deklaration aus einem "äußeren" Gültigkeitsbereich  
(innerer Gültigkeitsbereich überschreibt möglicherweise Deklarationen in äußeren Bereichen)

### Block

- Gültigkeitsbereich von Bezeichnern
- verschachtelt
- Zur Strukturierung von
  - a) Namensraum eines Programms
  - b) Symboltabelle

### Aufbau der Symboltabelle

- Folge von Bezeichnerblöcken
- Bezeichnerblock ::= alle Bezeichner eines Blocks mit Attributen
  - Sprache C: 2-stufig
  - Pascal: unbeschränkt (real 10)

Deklarationsniveaus:  
beliebig

Prozeduren/  
Funktionen

# Blockstrukturierte Symboltabelle: Beispiel

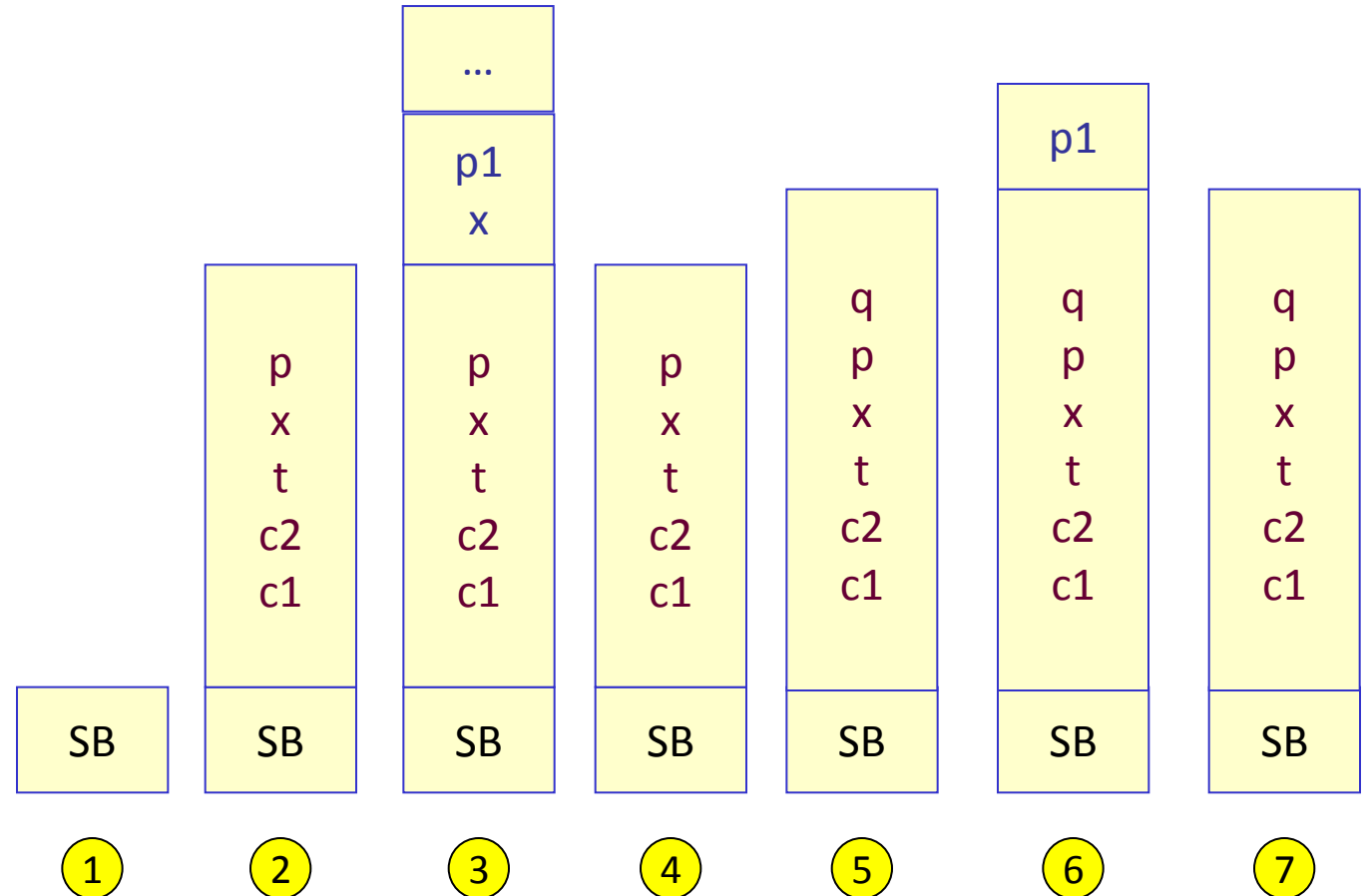
```

1 program Prog;
  const c1= 100;
        c2= 3.14;
  type  t= boolean;
  var   x: integer;

2  procedure p ();
3    var x : integer;
4    procedure p1
5    begin ...
6    end { p1 };
7  begin
8    x := c1;
  end { p };
  procedure q ();
    const p1= 20;
  begin ...
  end { q };
begin
...
end { Prog }

```

dynamische Tabellenverwaltung



# Operationales Interface: Symboltabelle

Realisierung als Hash-Tabelle

## Eigenschaften von Gültigkeitsbereich (GB)

- neue Deklarationen können nur im aktuellen GB definiert werden

## Welche Operationen?

- **void put (Symbol key, Object value)**  
Binden des Schlüssels (Namen) an einen (komplexen) Wert
- **Object get (Symbol key)**  
Auffinden des (komplexen) Wertes für einen gegebenen Schlüssel
- **void beginScope ()**  
Erzeuge neuen GB
- **void endScope ()**  
Schließe (und lösche) momentanen GB und setze den nächst äußeren GB als den jetzt gültigen

# Operationales Interface: Symboltabelle

- void initSymbolTable()
- void beginScope()
- void endScope()
- idPtr putSymbol(string s, idClass cl)
- idPtr getSymbol(string s)
- int currentScope()

- ... name(int entry)
- ... class(int entry)
- ... type(int entry)
- ... scope(int entry)
- ...
- Boolean isStandard(int entry)

} Attribut-Zugriff

Achtung Interface  
der Symboltabelle mit kompletten Signaturen  
wird in den Praktikumsübungen  
vorgestellt !!!

# Symboltabelle mit Codegenerierungsinfos

Attribute

Name	Klasse	Niveau	Typ	Spezielles	...
p	procedures	1	-	ParamListPtr	
x	variables	1	integers	Rel.Adr.: 12	
t	types	1	booleans	-	
c2	consts	1	reals	Wert: 3.14	
c1	consts	1	ints	Wert: 100	
program	keywords	1	-	-	
...	...	...	...	...	

p  
x  
t  
c2  
c1  
SB

weitere Attribute, z.B. zielcodeabhängige Infos



# Erweiterung der blockstrukturierten Symboltabelle

- Zuordnung von Speichergrößen für Namen mit einem Typ
- **relative Adressen** eines Namens  
lassen sich dann als **Offsets** in Bezug zum Beginn des Datenbereiches ermitteln

## Symboltabelle

Namen mit Typinfos und relativen Adressen  
bei Speicherbedarfsermittlung von Größen, die der  
Compiler für die Codegenerierung benötigt



**bislang für Codegenerierung noch ungelöst:**  
endgültiges Speicherlayout wird durch  
Adressierungsconstraints der jeweiligen Zielmaschine eingeschränkt

# Position

- Teil I  
Die Programmiersprache

- Teil II  
Methodische Grundlagen

- Teil III  
Entwicklung der Compiler

- Kapitel 1  
Compilationsprozess

- Kapitel 2  
Formalismen zur Sprachbeschreibung

- Kapitel 3  
Lexikalische Analyse: die Scanner

- Kapitel 4  
Syntaktische Analyse: die Parser

- Kapitel 5  
Parser-Generatoren: Yacc

- Kapitel 6  
Statische Semantikanalyse

- Kapitel 7  
Laufzeitsysteme

- Kapitel 8  
Ausblick: Codegenerierung

- 6.1  
Überblick

- 6.2  
Attribute

- 6.3  
S-Attribute

- 6.4  
Attribute und ereignisgesteuerte Semantik

- 6.5  
L-Attribute

- 6.6  
Verfahren zur semantischen Analyse  
Überblick

- 6.7  
Entwurf

- 6.8  
Drei-Adress-Code-Generierung (einige Aspekte)

- 6.9  
Symboltabelle

- 6.10  
Typprüfung

- 6.10.1  
Typsysteme, Typchecker

- 6.10.2  
Typchecker für eine einfache Sprache

- 6.10.3  
Wandlung kompatibler Typen

- 6.10.4  
Überladung von Operatoren und Funktionen

- 6.10.5  
Typcodierung

# Statische und dynamische Typprüfung

---

zur **Ausführungszeit**: dynamische Überprüfung

- prinzipiell kann jede Überprüfung dynamisch erfolgen, wenn Zielcode
  - den Typ des Elements zusammen mit
  - dem aktuellen Wert verwaltet(aber: schlechte Laufzeit- und Speichereigenschaften)

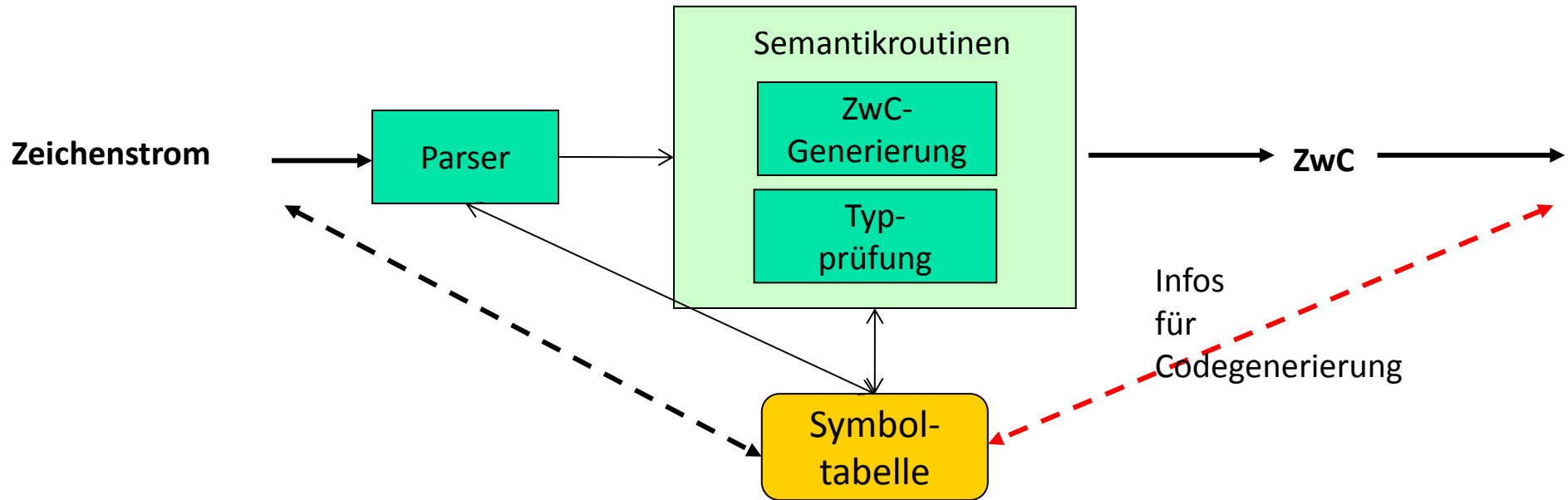
zur **Übersetzungszeit**: statische Überprüfung

- **streng-getypte** Sprachen erlauben (nahezu) komplette Typprüfung bereits durch den Compiler

**notwendige dynamische Überprüfungen**

- aber auch bei streng-getypten Sprachen
  - z.B. Indexverletzung
  - (Grenzüberprüfung von Feldern bzw. Feldzugriffen)

# Syntaxbaumaufbau und Typprüfung



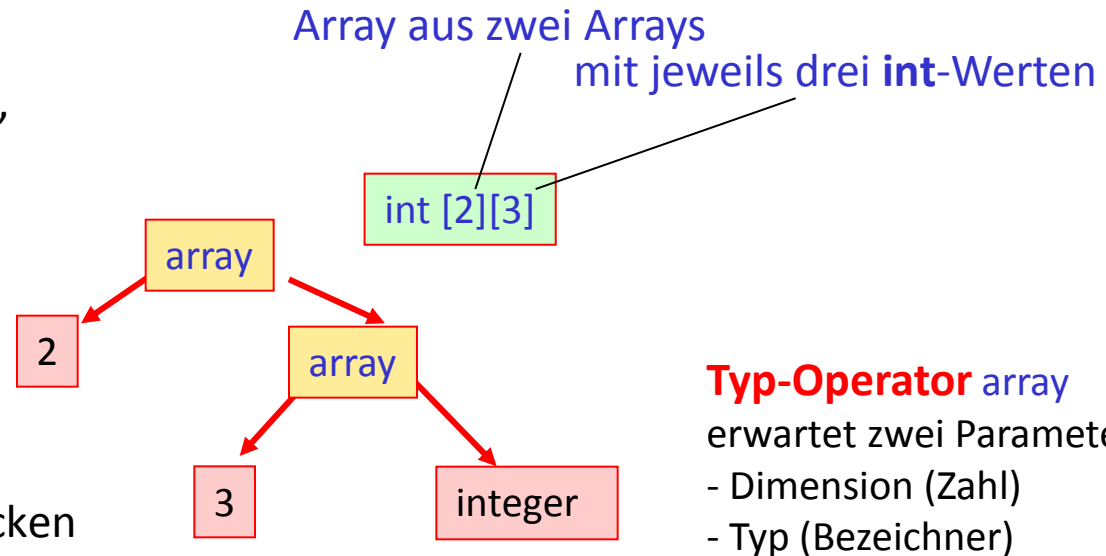
## Typchecker

- **allgemein** ein **Übersetzungsschema**, das den Typ jedes Ausdrucks aus den Typen seiner Teilausdrücke bestimmt
- **speziell** ein Analyseprogramm zur Behandlung von Arrays, Zeigern, Anweisungen und Funktionen

- bei komplexen Sprachen separater Pass für Typprüfung (z.B. Ada, SDL)
- bei einfachen Sprachen Kombination von Typprüfung mit Aufbau des Zwischencodes (z.B. Pascal)

# Typbeschreibungen

- Typen haben eine Struktur, die durch **Typausdrücke** formal beschrieben sind



- zwei Arten von Typausdrücken

## Basistypen

boolean, char, **integer**,  
real, ..., void

## konstruierte Typen

gebildet per Typoperator (Typkonstruktor)

Typausdrücke können mit Namen versehen werden (Typbezeichner sind Typausdrücke)

... durch Anwendung von Typ-Konstruktoren auf Typausdrücke (z.B. **array**, pointer, record, class, ...)

# Typkonstruktoren (1)

---

## ■ Felder

Ann.:  $I$  und  $T$  sind Typausdrücke,  
so beschreibt  $\text{array}(I, T)$  ein Feld von  $T$ , indiziert über  $I$

meist gibt es Einschränkung für den Indextyp

## ■ Produkte

Ann.:  $T1$  und  $T2$  sind Typausdrücke,  
so beschreibt  $T1 \times T2$  ein kartesisches Produkt über die Typausdrücke  $T1$  und  $T2$   
(Anwendung für Parameterlisten)

$\times$  ist linksassoziativ

## ■ Records

Felder (Elemente) haben - im Unterschied zu Produkten - Namen

z.B.:  $\text{record}((a \times \text{integer}) \times (b \times \text{real}))$

# Beispiele für Typkonstruktoren (2)

## ■ Zeiger

ist  $T$  ein Typausdruck,  
dann beschreibt `pointer (T)`  
den Typ »Zeiger auf ein Objekt vom Typ  $T$ «

z.B. Pascal: `var p: ^row`

## ■ Funktionen

$D \rightarrow R$  beschreibt den Typ einer Funktion,  
die Werte des Typs  $D$  auf Werte des Typs  $R$  abbildet,

z.B. `integer × char → pointer(integer)`

$\rightarrow$  ist der Funktionskonstruktor

Eine praktische Art der Darstellung von Typausdrücken  
ist die Verwendung von Graphen (ZwC-Repräsentation kommt später)

## Besonderheiten

### ■ Rekursive Typen

Beispiel: `class Link {  
    int info;  
    Link next;  
    ...  
}`

### ■ Typschablonen (Typausdrücke mit Variablen)

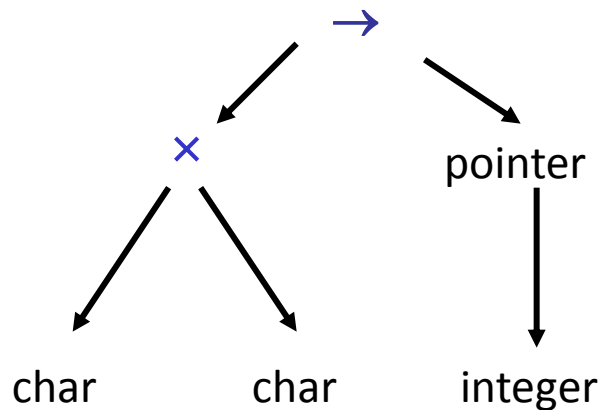
z.B. Templates in C++, Typen mit  
Kontextparametern in SDL

# Repräsentation von Typausdrücken

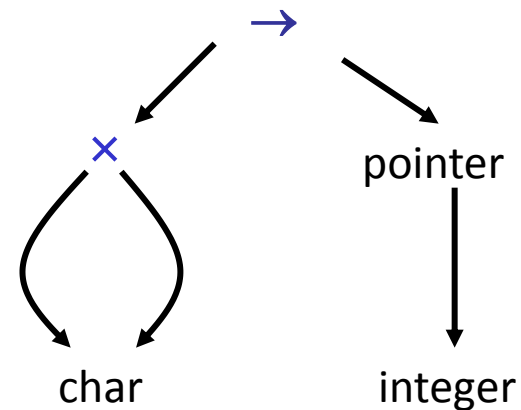
gerichteter azyklischer Graph (DAG),  
der durch Syntaxanalyse entsteht:

- innere Knoten sind Typkonstruktoren,
- Blätter sind einfache Typen, Typbezeichner, Typvariablen

Beispiel:  $\text{char} \times \text{char} \rightarrow \text{pointer}(\text{integer})$



oder





# Deklaration und Speicherlayout

für Namen von *Basis-* und *Array-Typen*

$D \rightarrow T \text{ id} ; D \mid \varepsilon$

$T \rightarrow B C \mid \text{record} \{ D \}$

$B \rightarrow \text{int} \mid \text{float}$

Basistypen

$C \rightarrow \varepsilon \mid [ \text{num} ] C$

Array-Typen

Einführung synthetisierter Attribute

für jedes Nichtterminalsymbol:

*type* [Typinfo]

*width* [Speicher in Bytes]

und globale Größen *t*, *w* als Ersatz ererbter Attribute

## Produktionen mit semantischen Regeln

$T \rightarrow B$  {*t*= B.type; *w*:= B.width}

$C$

$B \rightarrow \text{int}$  {B.type:= integer; B.width:=4;}

$B \rightarrow \text{float}$  {B.type:= float; B.width:=8;}

$C \rightarrow \varepsilon$  {C.type:= *t*; C.width:=*w*;}

$C \rightarrow [ \text{num} ] C1$  {C.type= array(num.value, C1.type);  
C.width:= num.value × C1.width; }

### **Achtung:**

*semantische Operationen,  
die leider **nicht** alle am Ende  
auftauchen*

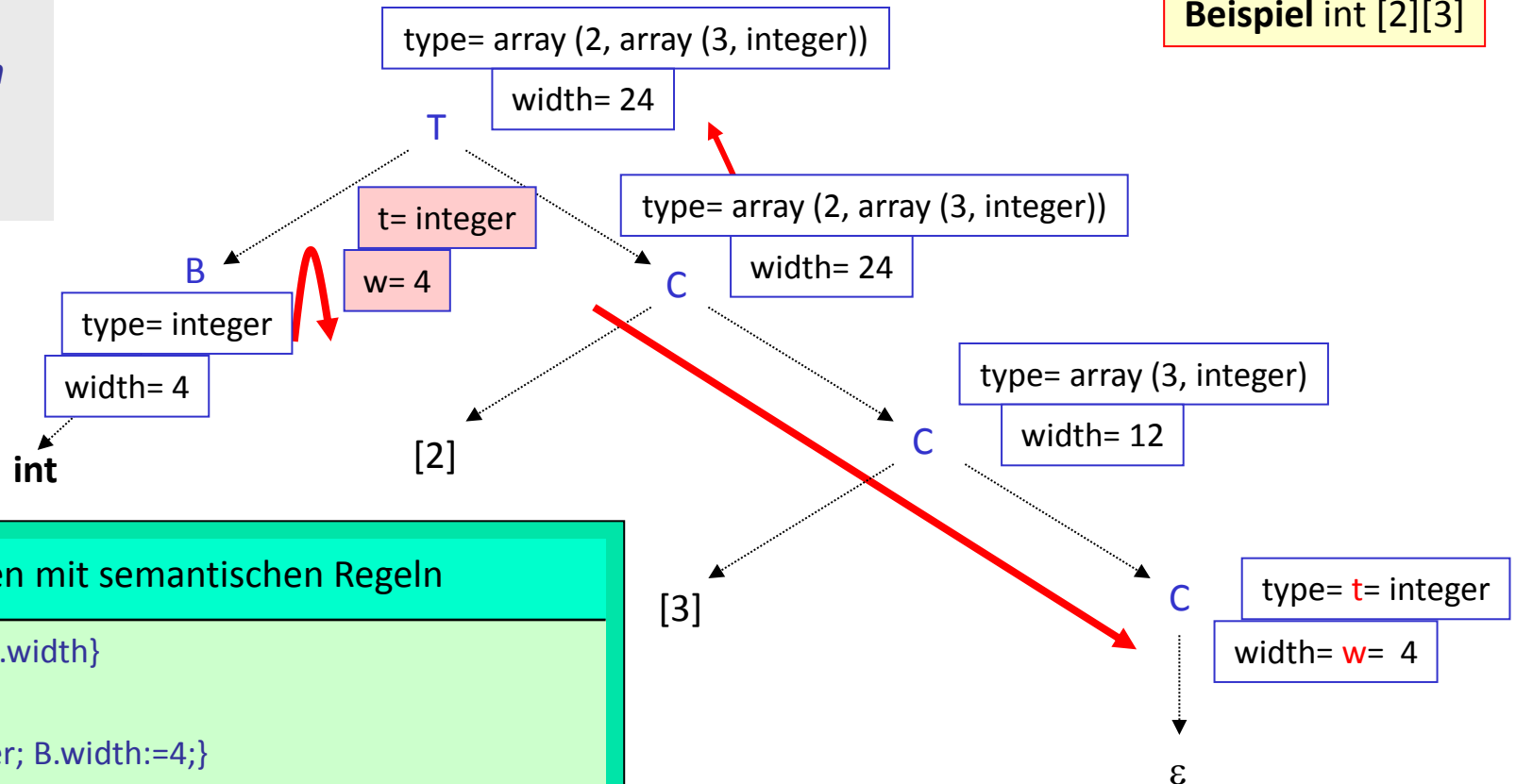
Ignorieren hierbei noch evtl. **Maschinenabhängigkeiten:**

Ausrichtung an Wortgrenzen entfällt (Beachtung in späterem Optimierungspass)

# Deklaration und Speicherlayout

Einsatz  
globaler Variablen  
anstatt ererbter  
Attribute

Beispiel int [2][3]



## Produktionen mit semantischen Regeln

```

T → B {t= B.type; w:= B.width}
    C
B → int {B.type:= integer; B.width:=4;}
B → float {B.type:= float; B.width:=8;}
C → ε {C.type:= t; C.width:=w;}
C → [ num ] C1 {C.type:= array(num.value, C1.type);
                C.width:= num.value × C1.width;}
    
```

# Erweiterung des Übersetzungsschemas

## Speicherlayout für Namen

betrachten jetzt  
Deklarationsliste

T als Typ,  
definiert wie vorher

Offset-Bestimmung  
bei Symboltabelleneintrag

### Produktionen mit semantischen Regeln

```
P → {offset:= 0}
    D
D → T id ; {top.put (id.lexeme, T.type, offset);
            offset:= offset + T.width; }
    D
D → ε
```

**offset**= relative Adresse  
des Speichers des bezeichneten  
Objektes  
im aktuellen Gültigkeitsbereich  
( top= aktuelle Symboltabelle)

### Produktionen mit semantischen Regeln

```
P → M D
M → ε {offset:= 0}
D → T id ; {top.put (id.lexeme, T.type, offset);
            offset:= offset + T.width; }
    D
D → e
```

**Achtung:**  
semantische Operationen,  
die leider auch  
**nicht** am Ende  
auftauchen

# Erweiterung des Übersetzungsschemas:

## Speicherlayout für Records (1)

Realisierung erfolgt über eigene  
(partielle) Symboltabelle je Record (!)

### Produktionen mit semantischen Regeln

```
T → record { { Env.push(top);  
               top:=new Env();  
               Stack.push(offset);  
               offset:=0; }  
  
D } { T.type:= record(top);  
      T.width:= offset;  
      top= Env.pop();  
      offset:= Stack.pop(); }
```

**Achtung:** Records können insbesondere verschachtelt sein

- Feldbezeichner eines Records müssen 1-deutig sein
- relative Adresse (offset) eines Feldbezeichners wird relativ zum Datenbereich des Records angegeben

→ strukturierte Namensräume

# Erweiterung des Übersetzungsschemas:

## Speicherlayout für Records (2)

### Produktionen mit semantischen Regeln

```
T → record { { Env.push(top);  
                top:=new Env();  
                Stack.push(offset);  
                offset:=0; }  
  
D } { T.type:= record(top);  
      T.width:= offset;  
      top:= Env.pop();  
      offset:= Stack.pop(); }
```

#### Aktionen vor D

- **Retten** der existierenden Symboltabelle **top** auf einem **SymboltabellenStack** der Umgebung
- Anlegen einer neuen Symboltabelle (**new Env**)
- **Retten** des aktuellen **Offsets** auf einem (synchronen) Offset-Stack
- Intialisierung eines **neuen Offsets** für das aktuelle Record
- **Deklarationen** liefern Typen und Offsets, die in aktueller Symboltabelle vermerkt werden

#### Aktionen nach D

- **Erzeugung** eines Records unter Verwendung von **der Symboltabelle top**
- **Wiederherstellung** von Symboltabelle und Offsets

# Typsysteme (Definition)

... sind Sammlungen von Regeln zur Zuweisung von Typausdrücken

**Bem.:**

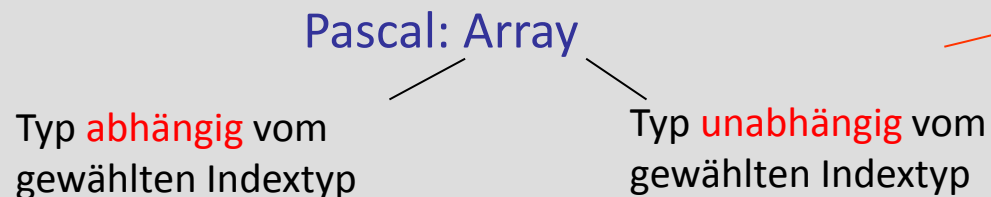
in verschiedenen Compilern der **gleichen** Sprache können durchaus (in Nuancen) unterschiedliche Typsysteme zum Einsatz kommen

## Spezifikation von Typsystemen

- syntaxgesteuerte Regeldefinition,  
so dass bekannte Techniken zur Implementation semantischer Aktionen eingesetzt werden können  
(attributierte Grammatiken)
- semantische Aktionen übernehmen **Typcheck**  
Sicherung einer partiellen Typäquivalenz



Einschränkung  
hinsichtlich  
Kompatibilität



# Aufgaben eines Typ-Checkers

Typ-Checker muss Typgleichheit / Typäquivalenz feststellen

Algorithmische Bestimmung der Typgleichheit hängt von Repräsentation der Typen ab

## Typäquivalenz

1. Behandlung impliziter (Namensgebung) und expliziter Typnamen
2. Behandlung rekursiver Typen
3. Behandlung expliziter und impliziter Typwandlungen

zwei Typen sind strukturell äquivalent

gdw. eine der folg. Bedingungen erfüllt ist:

1. sie sind dieselben Basistypen
2. sie wurden durch Anwendung desselben Konstruktors auf **strukturell** äquivalente Typen gebildet

zwei Typen sind namens äquivalent, wenn sie **namensgleich** sind

# Allgemeines Problem der Typäquivalenz

```
type link = ↑cell; {link ist der Typname für ↑cell }  
var next  : link;  
    last  : link;  
    p     : ↑cell;  
    q, r  : ↑cell;
```

**Beispiel:**  
Pascal-ähnliche  
Definition

**Problem:** sind next und last vom selben Typ wie p, q, r ?  
haben p und q überhaupt denselben Typ ?

## bei Namensäquivalenz

- next und last sind vom selben Typ
- p, q und r sind vom selben Typ
- p und next haben unterschiedliche Typen

bei **struktureller Typäquivalenz** müssten eigentlich alle Variablen jeweils denselben Typ haben  
(aber Praxis: Ada, Pascal und Modula-2

- sehen unterschiedliche Typdefinitionen als unterschiedliche Typen an,  
somit hat hier trotz geforderter struktureller Typäquivalenz p einen anderen Typ als q und r

**FRAGE: WARUM?**



# Berücksichtigung impliziter Typnamen

Für Pascal gilt:

enthält eine Deklaration einen Typausdruck, der kein Name ist, wird ein **impliziter** Name erzeugt:

```
type link = ↑cell;  
  
var next : link;  
    last : link;  
    p    : ↑cell;  
    q, r : ↑cell;
```



```
type link = ↑cell;  
    np1 = ↑cell;  
    np2 = ↑cell;  
  
var next : link;  
    last : link;  
    p    : np1;  
    q, r : np2;
```

p und q sind nicht vom selben Typ

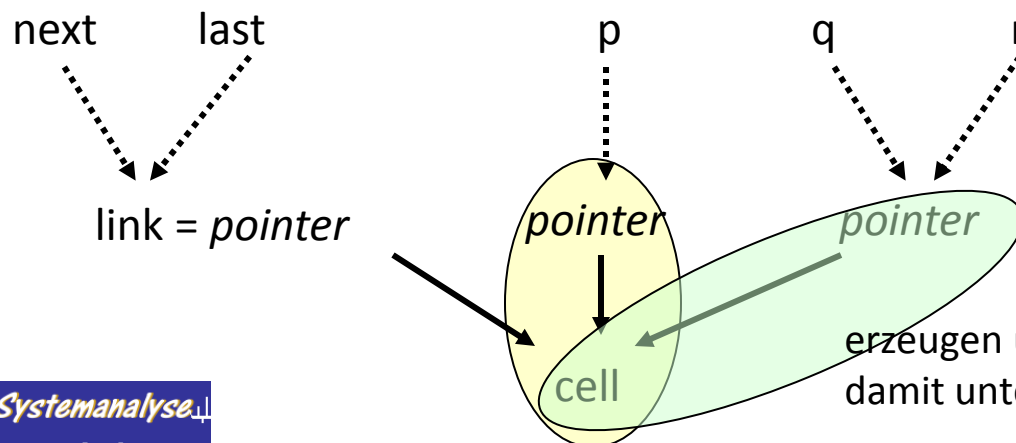
# Typäquivalenz in Pascal

Typ-Ausdrücke sind äquivalent, wenn sie durch denselben Knoten im Graph repräsentiert werden

## Typ-Graph: Struktur zur Übersetzungszeit

- jeder Konstruktor oder Basistyp erzeugt einen Knoten
- jeder Name erzeugt ein Blattknoten (der mit dem Typdeskriptor assoziiert wird)

in Pascal/Ada/ Modula: Deklaration von Bezeichnern mit **impliziten** Typnamen (immer dann, wenn ein Typausdruck verwendet wird)



```
type link = ↑cell;
var next : link;
    last : link;
    p    : ↑cell;
    q, r : ↑cell;
```

erzeugen unterschiedliche implizite Typnamen und damit unterschiedliche Knoten

# Fazit: Typausdrücke, Typprüfung und Typsystem

## Basis

### Compiler (Typ-Checker)

- weist jeder Komponente im Quellprogramm einen Typausdruck zu
- stellt fest, ob die Typausdrücke konform zu einer Menge von Regeln (Typsystem der Sprache) sind
- hat Potential für Fehlerbestimmung

**Typ-Checker** ist allgemein ein Übersetzungsschema,  
das den Typ jedes Ausdrucks aus den Typen seiner Teilausdrücke bestimmt  
und über semantische Regeln die Typtests realisiert

## Typausdrücke

- **array [256] of char**  
führt zum Typausdruck: *array ((1..256), char)*
- **↑int**

führt zum Typausdruck: *pointer (integer)*

# Typäquivalenz

Typ-Checker muss

(1) **Namensäquivalenz**

zwei Typen sind äquivalent, falls sie den gleichen Namen haben  
jeder neue Typname stellt einen neuen (anderen) Typ dar

(2) **Strukturelle Äquivalenz**

zwei Typen sind äquivalent, falls sie die gleiche Struktur haben (nachdem sämtliche Typnamen durch Typausdrücke ersetzt worden sind)

- $s \equiv t$ , falls  $s$  und  $t$  vom gleichen Basistyp sind
- $\text{array}(s_1, s_2) \equiv \text{array}(t_1, t_2)$ , falls  $s_1 \equiv t_1$  und  $s_2 \equiv t_2$
- $s_1 \times s_2 \equiv t_1 \times t_2$ , falls  $s_1 \equiv t_1$  und  $s_2 \equiv t_2$
- $\text{pointer}(s) \equiv \text{pointer}(t)$ , falls  $s \equiv t$
- $s_1 \rightarrow s_2 \equiv t_1 \rightarrow t_2$ , falls  $s_1 \equiv t_1$  und  $s_2 \equiv t_2$

feststellen

# Fazit: Typausdrücke, Typprüfung und Typsystem

## Basis

### Compiler (Typ-Checker)

- weist jeder Komponente im Quellprogramm einen Typausdruck zu
- stellt fest, ob die Typausdrücke konform zu einer Menge von Regeln (Typsystem der Sprache) sind
- hat Potential für Fehlerbestimmung

**Typ-Checker** ist allgemein ein Übersetzungsschema,  
das den Typ jedes Ausdrucks aus den Typen seiner Teilausdrücke bestimmt  
und über semantische Regeln die Typtests realisiert

## Typausdrücke

- **array [256] of char**  
führt zum Typausdruck: *array ((1..256), char)*
- **↑int**

führt zum Typausdruck: *pointer (integer)*