

# **Compilerbau SS2004**

Prof. Indermark  
GeTeXt von  
Michael Neuendorf

July 27, 2004

Dieses Skript ist eine Abschrift meiner Vorlesungsmitschriften der Vorlesung Compilerbau im Sommersemester 2004 und entbehrt sicherlich nicht einiger Tippfehler. Es schafft hoffentlich eine leserliche Übersicht über die behandelten Themen. Leider fehlen zur Zeit noch die Grafiken, ich hoffe sie aber noch bald ergänzen zu können. Da ich Fehler im Skript möglichst ausmerzen möchte, bitte ich darum, entdeckte Fehler per eMail an *michael.neuendorf@rwth-aachen.de* zu schicken.

Michael Neuendorf

# 1 Einleitung

**Compiler** = Programm zur Übersetzung von PS-Programm (Quellsprache) in äquivalente MS-Programme (Zielsprache)

**PS:** imperativ, deklarativ (funktional, logisch), nebenläufig, objekt-orientiert

**MS:** Maschinensprache, von-Neumann-Rechner

- RISC (reduced instruction set computer)
- CISC (complex instruction set computer)

## Aspekte einer PS

1. Syntax: formaler hierarchischer Aufbau eines Programms aus strukturellen Komponenten
2. Semantik: Bedeutung eines Programms, Zustandstransformation einer abstrakten Maschine
3. Pragmatik: benutzerfreundliche Formulierung, natürliche Sprache, Maschinenabhängigkeiten

Äquivalenz von Programmen: semantische Gleichheit

## 1.1 Struktur eines Compilers

logisch unabhängige Phasen, die aber verzahnt als Passes (Läufe) ablaufen:

- Analyse: Bestimmung der syntaktischen Struktur, Fehlerbestimmung
  - lexikalische Analyse:** Erkennung von Symbolen, Trennzeichen und Kommentaren (reguläre Ausdrücke, endliche Automaten)
  - syntaktische Analyse:** Erkennung des hierarchischen Programmaufbaus, Ableitungsbaum (CFG, Kellerautomaten)
  - semantische Analyse:** Kontextabhängigkeiten, statische Semantik, Typinformation, Attributierung des Ableitungsbaums mit semantisch relevanten Informationen (attributierte Grammatiken)
- Synthese: Erzeugung von MS-Code aus attributiertem Ableitungsbaum
  - Übersetzung in Zwischencode für eine abstrakte Maschine

- Optimierung: Verbesserung von Laufzeit und Speicherbedarf
- Codegenerierung: effiziente Verwendung von Registern und MS-Befehlssatz zur Erzeugung von MS-Code

**Frontend (MS-unabhängig)** Analyse, Zwischencodeerzeugung und Optimierung

**Backend (MS-abhängig)** MS-Code-Erzeugung und Optimierung

**One-Pass, n-Pass-Compiler** Zahl der Läufe durch das Quellprogramm

## 2 Lexikalische Analyse

**Ausgangspunkt** Quellprogramm  $P$  als Zeichenfolge

- $\Sigma_0$  Zeichensatz (z.B. ASCII),  $a \in \Sigma_0$  Zeichen, lexikalisches Atom (Lexem)
- $P \in \Sigma_0^*$  besitzt aufgrund der Pragmatik von PS eine lexikalische Struktur

**Pragmatischer Aspekt**

- natürliche Sprache für Bezeichner und Schlüsselwörter
- mathematische Formelsprache für Zahlen und Formeln
- Leerzeichen, Zeilenwechsel, Einrücken, ...
- Kommentare

↪ gute Lesbarkeit und Wartbarkeit

Die Semantik von  $P$  und damit die Übersetzung von  $P$  ist syntaxorientiert: Sie folgt dem hierarchischen Aufbau. Dabei sind pragmatische Aspekte irrelevant.

**1. Beobachtung** Syntaktische Atome (Symbole) werden dargestellt als Folgen lexikalischer Atome, sogenannter Lexeme.

**1. Aufgabe der lexikalischen Analyse** Zerlegung des Quellprogrammes  $P$  in eine Folge von Lexemen.

**2. Beobachtung** Für die syntaktische Analyse ist der Unterschied zwischen Lexem oft irrelevant, z.B. müssen Bezeichner nicht unterschieden werden.

**Symbolklassen** Lexeme werden zu Symbolklassen zusammengefaßt. Darstellung einer Symbolklasse: Token (z.B.  $(id, num)$ )

Die syntaktische Analyse bearbeitet eine Tokenfolge, Identifizierung eines Symbols durch zusätzliches Attribut für semantische Analyse und Codegenerierung.

$$Symbol = (Token, Attribut)$$

**2. Aufgabe der lexikalischen Analyse** Transformation einer Lexem-Folge in eine Symbolfolge.

**Lexikalische Analyse:** Zerlegung eines Quellprogrammes in eine Folge von Lexemen und deren Transformation in eine Folge von Symbolen.

**Scanner (Lexer):** Programm für die lexikalische Analyse

## 2.1 Wichtigste Symbolklassen (eigentlich Lexemklassen)

- Bezeichner
- Zahlwörter
- Schlüsselwörter
- einfache Symbole: ein Sonderzeichen, z.B. +, -, (, bildet jeweils eine Symbolklasse
- zusammengesetzte Symbole: Folgen von zwei oder mehr Zeichen
- Leerzeichen:  $\sqcup, \sqcup \dots \sqcup$
- Spezielle Symbole: Kommentare, Pragmas wie Compileroptionen

**Token:** id, const, divsym, semsym

**Attribute:** Zeiger in Symboltabelle, Binärdarstellung einer Zahl, leer bei Symbolklassen mit (nur) einem Symbol

**Feststellung:** Symbolklassen sind reguläre Mengen

$\rightsquigarrow$  Beschreibung durch reguläre Ausdrücke, Erkennung durch endliche Automaten

## 2.2 Scanner-Konstruktion

**Definition (reguläre Ausdrücke)** Für ein Alphabet  $\Sigma$  ist die Menge  $\text{RegE}(\Sigma)$  der regulären  $\Sigma$ -Ausdrücke definiert durch

- $\Lambda \in \text{RegE}(\Sigma), \Sigma \subseteq \text{RegE}(\Sigma)$
- $\alpha, \beta \in \text{RegE}(\Sigma) \rightsquigarrow (\alpha \vee \beta), (\alpha \wedge \beta), (\alpha^*) \in \text{RegE}(\Sigma)$

Ihre Semantik:

- $\llbracket \cdot \rrbracket : \text{RegE}(\Sigma) \rightarrow P(\Sigma^*)$
- $\llbracket \Lambda \rrbracket : \emptyset$
- $\llbracket a \rrbracket := \{a\}$
- $\llbracket (\alpha \vee \beta) \rrbracket := \llbracket \alpha \rrbracket \cup \llbracket \beta \rrbracket$
- $\llbracket (\alpha \cdot \beta) \rrbracket := \llbracket \alpha \rrbracket \cdot \llbracket \beta \rrbracket$
- $\llbracket (\alpha^*) \rrbracket := \llbracket \alpha \rrbracket^*$

## 2.3 Das einfache Matching-Problem

Entscheide für  $\alpha \in \text{RegE}(\Sigma)$  und  $w \in \Sigma^*$ , ob  $w \in \llbracket \alpha \rrbracket$  oder  $w \notin \llbracket \alpha \rrbracket$ .

### 2.3.1 Hilfsmittel: endliche Automaten

$\mathfrak{A} = \langle Q, \Sigma, \delta, q_0, F \rangle \in \text{NFA}(\Sigma)$  mit  $\delta : Q \times \Sigma_\varepsilon \rightarrow P(Q)$ ,  $\Sigma_\varepsilon := \Sigma \cup \{\varepsilon\}$  und  $q_0 \in Q, F \subseteq Q$ .

Für  $T \subseteq Q$  ist die  $\varepsilon$ -Hülle  $\varepsilon(T)$  definiert durch

- $T \in \varepsilon(T)$
- $q \in \varepsilon(T) \curvearrowright \delta(q, \varepsilon) \subseteq \varepsilon(T)$

### 2.3.2 Die erweiterte Transitionsfunktion

$\bar{\delta} : P(Q) \times \Sigma^* \rightarrow P(Q)$  ist definiert durch

- $\bar{\delta}(T, \varepsilon) := \varepsilon(T)$
- $\bar{\delta}(T, wa) := \varepsilon\left(\bigcup_{q \in \bar{\delta}(T, w)} \delta(q, a)\right)$

$\mathfrak{A}$  erkennt die Sprache  $L(\mathfrak{A}) := \{w \in \Sigma^* \mid \bar{\delta}(\{q_0\}, w) \cap F \neq \emptyset\}$ .

### 2.3.3 Die DFA-Methode zur Lösung des Matching-Problems

(bekannt aus ATFS)

$\alpha \in \text{RegE}(\Sigma) \xrightarrow{(1)} \mathfrak{A}(\alpha) \in \text{NFA}(\Sigma) \xrightarrow{(2)} \mathfrak{A}(\alpha)^P \in \text{DFA}(\Sigma)$  mit  $\llbracket \alpha \rrbracket = L(\mathfrak{A}(\alpha)) = L(\mathfrak{A}(\alpha)^P)$ .

- (1) Methode von Thompson (Satz von Kleene) (linearer Platz- und Zeitbedarf)
- (2) Potenzmengenkonstruktion (exponentieller Zeit- und Platzbedarf)

Aber  $\mathfrak{A}(\alpha)^P$  löst das Matching-Problem ( $w \in \llbracket \alpha \rrbracket$ ) in  $(|w| + 1)$  Schritten.

### 2.3.4 Die NFA-Methode zur Lösung des Matching-Problems

Verkleinerung des Platzbedarf, jedoch längere Laufzeit durch den Verzicht auf die Potenzmengenkonstruktion.

**Idee** Berücksichtigung der Eingabe  $w \in \Sigma^*$ , Potenzmengenkonstruktion nur für den "Lauf von  $w$  durch  $\mathfrak{A}(\alpha)$ ".

#### Komplexität (NFA-Methode)

- Platz:  $O(|\alpha| + |w|)$
- Zeit:  $O(|\alpha| \cdot |w|)$

**Kombination von NFA- und DFA-Methode** Zwischenspeichern von bereits berechneten Transitionen  $\bar{\delta}(T, a)$  im Cache.

**Das erweiterte Matching-Problem** Seien  $\alpha_1, \dots, \alpha_n \in \text{RegE}(\Sigma)$  und  $w \in \Sigma^*$ . Sei ferner  $\Delta := \{T_1, \dots, T_n\}$  ein Alphabet von Token. Wenn  $w = w_1 w_2 \dots w_k$  und  $w_j \in \llbracket \alpha_{i_j} \rrbracket$  mit  $1 \leq i_j \leq n$  für  $j = 1, \dots, k$ , dann heißt  $w_1, w_2, \dots, w_k$  eine Zerlegung von  $w$  bezüglich  $\alpha_1, \dots, \alpha_n$  und  $v = T_{i_1} \dots T_{i_k}$  eine Analyse von  $w$  bezüglich  $\alpha_1, \dots, \alpha_n$ .  
Hinweis: Mehrdeutigkeit

### Konventionen für Eindeutigkeit

1. Das Prinzip des längsten Match (max. munch) für die Eindeutigkeit der Zerlegung

Eine Zerlegung  $(w_1, \dots, w_k)$  von  $w$  bezüglich  $\alpha_1, \dots, \alpha_n$  heißt lm-Zerlegung, wenn für alle  $j = 1, \dots, k$ ,  $x, y \in \Sigma^*$  und  $p, q \in \{1, \dots, n\}$  gilt:

$$w = w_1 w_2 \dots w_j x y, \quad w_j \in \llbracket \alpha_p \rrbracket, \quad w_j x \in \llbracket \alpha_q \rrbracket \curvearrowright x = \varepsilon$$

Folgerung: Für  $w, \alpha_1, \dots, \alpha_n$  gibt es höchstens eine lm-Zerlegung.

2. Prinzip des ersten Matches (bezüglich  $\alpha_1, \dots, \alpha_n$  für die Eindeutigkeit der Analyse.

Trotz der Eindeutigkeit der lm-Zerlegung kann es mehrere zugehörige Analysen geben, weil  $\llbracket \alpha_p \rrbracket \cap \llbracket \alpha_q \rrbracket \neq \emptyset$  möglich ist.

Konvention: Erster Match in einer Reihe  $\alpha_1, \dots, \alpha_n$  zählt.

Sei  $(w_1, \dots, w_l)$  eine lm-Zerlegung und  $v = T_{i_1} \dots T_{i_k}$  eine (zugehörige) Analyse von  $w$  bezüglich  $\alpha_1 \dots \alpha_n$ . Dann heißt  $v$  eine flm-Analyse (first longest match) falls für alle  $j = 1, \dots, k$  und  $v = 1, \dots, n$  gilt:

$$w_j \in \llbracket \alpha_v \rrbracket \curvearrowright i, j \leq v$$

### Das erweiterte Matching-Problem

1. Existiert für  $w \in \Sigma^*$  bezüglich  $\{\alpha_1 \dots \alpha_n\} \subseteq \text{RegE}(\Sigma)$  eine lm-Zerlegung ?

2. Berechnung einer flm-Analyse

$\Delta = \{T_1, \dots, T_n\}$  Tokenalphabet

- a) Kontruieren für  $i = 1, \dots, n$

$$\mathfrak{A}_i = \langle Q_i, \Sigma, \delta_i, q_0^{(i)}, F_i \rangle \in \text{DFA}(\Sigma) \text{ mit } \llbracket \alpha_i \rrbracket = L(\mathfrak{A}_i).$$

- b) Bilde aus diesen den Produktautomaten

$$\mathfrak{A} = \langle Q, \Sigma, \delta, q_0, F \rangle \in \text{DFA}(\Sigma) \text{ mit}$$

- $Q = Q_1 \times \dots \times Q_n,$
- $q_0 := (q_0^{(1)}, \dots, q_0^{(n)}),$
- $\delta(q^{(1)}, \dots, q^{(n)}) := (\delta_1(q^{(1)}, a), \dots, \delta_n(q^{(n)}, a))$  und



- $(q^{(1)}, \dots, q^{(n)}) \in F \Leftrightarrow \exists i \in \{1, \dots, n\} q^{(i)} \in F$

Dann gilt  $L(\mathfrak{A}) = \bigcup_{i=1}^n \llbracket \alpha_i \rrbracket$ .

Zerlege  $F$  wegen "first match" in  $F = \bigcup_{i=1}^n f^{(i)}$ .  $(q^{(1)}, \dots, q^{(n)}) \in F^{(i)} \Leftrightarrow q^{(i)} \in F_i$  und  $q^{(j)} \notin F_j$  für alle  $1 \leq j \leq i$ . Dann gilt:  $\bar{\delta}(q_0, w) \in F^{(i)} \Leftrightarrow w \in \llbracket \alpha_i \rrbracket$  und  $w \notin \bigcup_{j=1}^{i-1} \llbracket \alpha_j \rrbracket$ .

$q \in Q$  heißt *produktiv*, falls es ein  $v \in \Sigma^*$  gibt, so daß  $\bar{\delta}(q, v) \in F$ .  $P$  ist die Menge der produktiven Zustände, also  $F \subseteq P$ .

- c) Erweitere  $\mathfrak{A} = \langle Q, \Sigma, \delta, q_0, F \rangle$  zu einem Backtrack-DFA  $\mathfrak{B}$  mit Ausgabe.

Idee: Einweg-Leseband mit 2 Köpfen

- ein Backtrack-Kopf  $b$  zur Markierung eines Matches
- ein Lookahead-Kopf  $l$  zur Bestimmung des längsten Matches

Konfigurationsmenge von  $\mathfrak{B}$ :  $\{N\} \cup \Delta \times \Sigma^* Q \Sigma^* \times \Delta^* \{\varepsilon, \text{lexerr}\}$

Ausgangskonfiguration für  $w \in \Sigma^*$ :  $(N, q_0 w, \varepsilon)$

Transitionen:  $q' := \delta(q, a)$

- i. normal mode: match suchen

$$(N, qaw, W) \vdash \begin{cases} (N, q'w, W) & \text{falls } q' \in P \setminus F \\ (T_i, q'w, W) & \text{falls } q' \in F^{(i)} \\ \text{Ausgabe: } W \text{ lexerr} & \text{falls } q' \notin P \end{cases}$$

- ii. backtrack mode: längsten match suchen

$$(T, vqaw, W) \vdash \begin{cases} (T, vaq'w, W) & \text{falls } q' \in P \setminus F \\ (T_i, q'w, W) & \text{falls } q' \in F^{(i)} \\ (N, q_0vaw, WT) & \text{falls } q' \notin P \end{cases}$$

- iii. Eingabeende

- $(N, q, W) \vdash$  Ausgabe:  $W \text{lexerr}$ , falls  $q \in P \setminus F$
- $(T, q, W) \vdash$  Ausgabe:  $WT$ , falls  $q \in F$
- $(T; vaq, W) \vdash (N, q_0va, WT)$ , falls  $q \in P \setminus F$

Dann gilt für  $w \in \Sigma^*$ :

- $(N, q_0w, \varepsilon) \vdash W \in \Delta^* \Leftrightarrow W$  ist flm-Analyse von  $w$
- $(N, q_0w, \varepsilon) \vdash W \cdot \text{lexerr} \Leftrightarrow$  es gibt keine flm-Analyse von  $w$

Zeitaufwand:  $O(|w|^2)$  im worst case, Verbesserung durch Tabular-Methode möglich zu Linearzeit. Literatur: Th. Reps "Maximal Munch Tokenization in Linear Time", ACM-TOPLAS 20(1998).

## 2.3.5 Praktische Aspekte der Scannerkonstruktion

- Unterschiede
  - automatische Scannererzeugung (lex u.a.) mit beliebigen  $\alpha_1 \dots \alpha_n$
  - Scanner-Konstruktion für eine Programmiersprachen mit speziellen  $\alpha_1 \dots \alpha_n$  (Effizienz wichtig, häufiger Aufruf durch nexttoken, direkte Programmierung [Assembler])
  - Spracherweiterung von  $\text{RegE}(\Sigma)$ 
    1. Vereinfachende Bezeichnungen
      - \* Präzedenzregeln zur Vermeidung von Klammern: \* bindet stärker als  $\cdot$  bindet stärker als  $\vee$ ,  $\cdot$  und  $\vee$  sind assoziativ (Linksklammerung),  $\cdot$  weglassen,  $|$  statt  $\vee$
      - \* Beispiel:  $a|b * c$  statt  $(a \vee ((b^*) \cdot c))$
      - \* Abkürzungen:

$a^+$	:=	$aa^*$	
$a?$	:=	$a \Lambda^*$	
$\cdot$	:=	$a \dots z$	( $\Sigma = \{a, \dots, z\}$ )
$[abc]$	:=	$a b c$	
$[0 - 9]$	:=	$0 1 \dots 9$	
    2. Reguläre Definitionen

Schrittweise Beschreibung von Symbolklassen (regulären Mengen) durch zusätzliche Bezeichner  $id_1 = \alpha_1, \dots, id_n = \alpha_n$  mit  $id_1, \dots, id_n \notin \Sigma$  und  $\alpha_i \in \text{RegE}(\Sigma \cup \{id_1, \dots, id_{i-1}\})$

Beachte: keine Rekursion, Entschachtelung möglich (Rekursion  $\leadsto$  EBNF, ECFG)
- Prinzip des längsten Matches
  - Allgemein: beliebig langer "look-ahead" (backtrack-Phase) erforderlich
  - PS(Pascal): look-ahead um 2 Zeichen
  - aber: Beschränkung kann die lexikalische Analyse verändern
  - Abweichung von lm-Prinzip möglich:
    - \* lex:  $\alpha/\beta$  Match von  $\alpha$  erfordert zusätzlich einen lookahead-String, der  $\beta$  matcht
  - Besondere Rolle von "blanks" (white space): Trennung von "eentlichen" Lexem, um mit 1-look-ahead auszukommen
- Attributberechnung

Nach Erkennen eines Lexems erneutes Lesen:

  - Dezimalzahl  $\rightsquigarrow$  Binärzahl (lex: `install_num`)
  - Bezeichner: Gültigkeit prüfen, Symboltabelle (lex: `install_id`)

- Schlüsselwörter (alternativ): zunächst als Bezeichner behandeln, Symbolta-  
belle mit Schlüsselwörtern initialisieren, Attributberechnung (`install_id`)  
liefert das Schlüsselworttoken
- Automatische Scannergenerierung mit Lex
  - Unixtool: `man flex` ("fast lexical analyzer") oder `info flex`
  - `scan.l`  $\xrightarrow{\text{lex}}$  `lex.yy.c`  $\xrightarrow{\text{C-Compiler}}$  `a.out`
  - Programm `a.out` Tokenfolge mit Attributbenennung
  - Aufbau eine Lex-Spezifikation
    - Definitionen (optional)
    - Regeln
    - C-Hilfsprozeduren (optional)

- Definitionen**
1. Direkter C-Code `%{... C-Code ... %}`
  2. Substitutionen (reguläre Definitionen) `NAME regexp...`
  3. Startzustände (zur Einleitung in besondere Phasen)

**Regeln** `muster {aktion}`

**muster** von der Form `regexp [/regexp]`

**aktion** C-Code zur Berechnung von `<Token, Attribut>`

**Finden der passenden Regel** \* longest match

\* first match

\* kein match: abfangen mit

`.\n`

↪ Fehlerbehandlung

**Attribute** Weitergabe über globale Variable `yylval`

# 3 Syntaktische Analyse

**Aufgabe:** Zerlegung der Symbolfolge, die der Scanner ausgibt, in syntaktische Einheiten oder Behandlung syntaktischer Fehler

## 3.1 Syntaktische Einheiten

Variablen, Ausdrücke, Anweisungen . . .

Beachte: Schachtelung syntaktischer Einheiten, Baumstruktur im Unterschied zur linearen Symbolfolge

**Beschreibung der syntaktischen Struktur, kontextfreie Grammatiken** (BNF, EBNF, Syntaxdiagramme)

**Problem:** deterministische Simulation

**Allgemeiner Fall:** beliebige CFG, Lösung durch das Tabularverfahren von Cocke, Younger, Kasami mit  $O(n^2)$  Platz- und  $O(n^3)$  Zeitbedarf

**Programmiersprachen:** spezielle CFG, Analyse durch deterministische Kellerautomaten mit "Input-Lookahead" bei linearem Platz- und Zeitbedarf

1. Top-Down-Analyse: Konstruktion des Ableitungsbaumes von der Wurzel zu den Blättern in Form einer Linksanalyse
2. Bottom-Up-Analyse: umgekehrt, gespiegelte Rechtsanalyse

**Kontextfreie Grammatiken**  $\mathcal{G} = \langle N, \Sigma, P, S \rangle \in \text{CFG}(\Sigma)$

$A, B, C, \dots$	$\in N$	Nichtterminalsymbole
$a, b, c, \dots$	$\in \Sigma$	Terminalsymbole
$u, v, w, \dots$	$\in \Sigma^*$	Terminalwörter
$\alpha, \beta, \gamma, \dots$	$\in \mathfrak{X}^*$	Satzformen
$A \rightarrow \alpha$	$\in P$	Produktion / Regel

Table 3.1: Bezeichnungskonventionen

Ableitungsrelation:  $\Rightarrow \subseteq (\mathcal{X}^*)^2$  definiert durch

$$\alpha \Rightarrow \beta \Leftrightarrow \begin{array}{l} \alpha = \alpha_1 A \alpha_2 \quad A \mapsto \gamma \in P \\ \beta = \alpha_1 \gamma \alpha_2 \end{array}$$

Gilt außerdem  $\alpha_1 \in \Sigma^*$  bzw.  $\alpha_2 \in \Sigma^*$ , so

$$\alpha \underset{l}{\Rightarrow} \text{ bzw. } \alpha \underset{r}{\Rightarrow} \beta$$

Sprechweise: Ableitungsschritt, Linksableitungsschritt, Rechtsableitungsschritt

### Erzeugte Sprachen

$$\begin{aligned} L(\mathcal{G}) &:= \{w \in \Sigma^* \mid S \overset{*}{\Rightarrow} w\} \\ &= \{w \in \Sigma^* \mid S \underset{l}{\overset{*}{\Rightarrow}} w\} \\ &= \{w \in \Sigma^* \mid S \underset{r}{\overset{*}{\Rightarrow}} w\} \end{aligned}$$

Beispiel:  $\mathcal{G} : S \mapsto aSb \mid \varepsilon$

$$S \Rightarrow aSb \Rightarrow aaSbb \overset{*}{\Rightarrow} a^n S b^n \Rightarrow a^n \cdot \varepsilon \cdot b^n \Rightarrow a^n b^n$$

**Ableitungsbläume** repräsentieren die syntaktische Struktur des abgeleiteten Wortes  $a^2 b^2$ .

**Definition** Eine Grammatik  $\mathcal{G}$  heißt eindeutig, wenn es für jedes  $w \in L(\mathcal{G})$  genau einen Ableitungsbaum gibt.

Folgerung:  $\mathcal{G}$  ist mehrdeutig, wenn es für jedes  $w \in L(\mathcal{G})$  genau eine Linksableitung (bzw. Rechtsableitung) gibt.

**Definition**  $\mathcal{G}$  heißt mehrdeutig, wenn  $\mathcal{G}$  nicht eindeutig ist.

**l-Analyse, r-Analyse** Darstellung von l/r-Ableitungen durch Nummernfolgen:

$$|P| = p \quad [p] := \{1, \dots, p\}$$

$$\pi : [p] \rightarrow P \text{ bijektiv}$$

Sei  $i \in [p]$ :

$$wA\alpha \underset{l}{\overset{i}{\Rightarrow}} w\gamma\alpha \text{ bzw. } \alpha Aw \underset{r}{\overset{i}{\Rightarrow}} \alpha\gamma W \text{ falls } \pi_i := \pi(i) = A \mapsto \gamma$$

Für  $z = i_1 \dots i_n \in [p]^+$  soll  $\alpha \overset{z}{\underset{l}{\Rightarrow}} \beta \Leftrightarrow \alpha \overset{i_1}{\underset{l}{\Rightarrow}} \alpha_1 \overset{i_2}{\underset{l}{\Rightarrow}} \alpha_2 \dots \overset{i_n}{\underset{l}{\Rightarrow}} \alpha_n = \beta$  und  $\alpha \overset{\varepsilon}{\underset{l}{\Rightarrow}} \alpha$  (ein "leerer" Ableitungsschritt) gelten.

**Definition** Sei  $z \in [p]^*$ .  $z$  heißt l-Analyse von  $\alpha : S \overset{z}{\underset{l}{\Rightarrow}} \alpha$ . Analog: r-Analyse.

**Syntaxanalyse** Gegeben:  $\mathcal{G} \in \text{CFG}(\Sigma)$  und  $w \in \Sigma^*$ . Berechnung einer Links- bzw. Rechtsanalyse, falls  $w \in L(\mathcal{G})$ . Andernfalls: Bestimmung syntaktischer Fehler.

**Generalvoraussetzung**  $G \in \text{CFG}$  ist reduziert, d.h. für jedes  $A \in N$  gibt es  $\alpha, \beta \in \mathfrak{X}^*$  und  $w \in \Sigma^*$ , sodaß

- $S \xrightarrow{*} \alpha A \beta$  ( $A$  erreichbar) und
- $A \xrightarrow{*} w$  ( $A$  produktiv).

### 3.1.1 Top-Down-Analyse, $LL(k)$ -Grammatiken

**Definition** Sei  $\mathcal{G} \in \text{CFG}(\Sigma)$ . Der (nicht-deterministische Top-Down-Analyseautomat von  $\mathcal{G}$  (Bezeichnung:  $\text{NTA}(\mathcal{G})$ )

- Eingabealphabet:  $\Sigma$
- Kelleralphabet:  $\mathfrak{X} = N \cup \Sigma$
- Ausgabealphabet:  $[p] = \{1, \dots, p\}$
- ("nur 1 Zustand")
- Bild (einfügen)
- Konfigurationsmenge:  $\Sigma^* \times \mathfrak{X}^* \times [p]^*$
- Transitionen:
  - Ableitungsschritt:  $(w, A\alpha, z) \vdash (w, \beta\alpha, z_i)$ , falls  $\pi_i = A \mapsto \beta$
  - Vergleichsschritt:  $(aw, a\alpha, z) \vdash (w, \alpha, z)$ , für  $a \in \Sigma$
- Anfangskonfiguration für  $w \in \Sigma$ :  $(w, S, \varepsilon)$
- Endkonfigurationen:  $(\varepsilon, \varepsilon, z)$
- Beachte: Nicht-Determinismus wegen  $A \mapsto \beta|\gamma$

**Satz** Der  $\text{NTA}(\mathcal{G})$  berechnet l-Analysen, d.h.  $(w, S, \varepsilon) \vdash^* (\varepsilon, \varepsilon, z) \leftrightarrow z$  ist l-Analyse von  $w$ .

Beweis:

- " $\curvearrowright$ " Automat arbeitet korrekt
- (\*)  $(w, \alpha, y) \vdash^* (\varepsilon, \varepsilon, yz) \curvearrowright \alpha \stackrel{\curvearrowright}{\downarrow} w$

Beweis von (\*) durch Induktion über  $k = |z|$ .

1.  $k = 0$ : Nur Vergleichsschritte,  $w = \alpha$ , also (\*) mit  $w \stackrel{\curvearrowright}{\downarrow} w$ .

2.  $k \rightarrow k + 1$ :  $z = iz', \alpha = wA\beta, w = uv, \pi_i = A \rightarrow \gamma$   
 $(w, \alpha, y) = (uv, uA\beta, y) \xrightarrow{*} (v, A\beta, y) \mapsto (v, \gamma\beta, yi) \xrightarrow{*} (\varepsilon, \varepsilon, yiz')$   
 Nach Induktionsvoraussetzung:  $\gamma\beta \xrightarrow{z'}_l v$ . Damit folgt (\*):

$$\alpha = uA\beta \xrightarrow{i}_l u\gamma\beta \xrightarrow{z'}_l uv = w$$

q.e.d.

**Vorüberlegung**  $3 : A \rightarrow a\beta \quad 4 : A \rightarrow b\beta$   
 $(aw, A\alpha, z) \mapsto (aw, a\beta\alpha, z3) \mapsto (w, \beta\alpha, z3)$   
 $A \rightarrow BC \quad B \rightarrow \varepsilon \quad C \rightarrow \varepsilon$

**Ziel:** Nicht-Determinismus des  $\text{NTA}(\mathcal{G})$  durch  $k$ -look-ahead auf der Eingabe beseitigen,  $k \in \mathbb{N}$ .

**Definition (first-Mengen):** Sei  $\mathcal{G} \in \text{CFG}, \alpha \in \Sigma^*$  und  $k \in \mathbb{N}$ . Wir definieren  $\text{first}_k(\alpha) \subseteq \Sigma^*$  durch

$$\begin{aligned} \text{first}_k(\alpha) := & \{v \in \Sigma^* \mid \exists w \in \Sigma^* : \alpha \xrightarrow{*} vw, |v| = k\} \\ & \cup \{v \in \Sigma^* \mid \alpha \xrightarrow{*} v, |v| < k\} \end{aligned}$$

### Folgerungen (Übung)

1.  $\text{first}_k(\alpha) \neq \emptyset$ , weil  $\mathcal{G}$  reduziert ist.
2.  $\varepsilon \in \text{first}_k(\alpha) \leftrightarrow k = 0$  oder  $\alpha \xrightarrow{*} \varepsilon$
3.  $\alpha \xrightarrow{*} \beta \curvearrowright \text{first}_k(\beta) \leq \text{first}_k(\alpha)$
4.  $v \in \text{first}_k(\alpha) \leftrightarrow \exists x \in \Sigma^* : \alpha \xrightarrow{*}_l x, \{v\} = \text{first}_k(x)$

$LL(k)$ : Lesen der Eingabe von links nach rechts mit  $k$ -look-ahead, Berechnung einer l-Analyse

**$LL(k)$ -Grammatik** Sei  $\mathcal{G} \in \text{CFG}$  und  $k \in \mathbb{N}$ .  
 $\mathcal{G} \in LL(k) : \leftrightarrow$  für alle Linksableitungen der Form

$$S \xrightarrow{*}_l wA\alpha \left\{ \begin{array}{l} \xrightarrow{*}_l w\beta\alpha \xrightarrow{*}_l wx \\ \xrightarrow{*}_l w\gamma\alpha \xrightarrow{*}_l wy \end{array} \right.$$

gilt:  $\text{first}_k(x) = \text{first}_k(y) \curvearrowright \beta = \gamma$ .  
 Bemerkung:

- Linksableitungsschritt für  $wA\alpha$  ist durch die nächsten  $k$  auf  $w$  folgenden Symbole bestimmt.
- Der NTA kann für  $\mathcal{G} \in LL(k)$  deterministisch mit  $k$ -look-ahead auf der Eingabe arbeiten.

**Problem:** Bestimmung der  $A$ -Regel aus  $k$ -look-ahead.

**Lemma**  $\mathcal{G} \in LL(k) \leftrightarrow$  für alle Linksableitungen der Form  $S \xrightarrow{*} wA\alpha$   $\left\{ \begin{array}{l} \xrightarrow[l]{*} w\beta\alpha \\ \xrightarrow[l]{*} w\gamma\alpha \end{array} \right.$  gilt:

$$\beta \neq \gamma \curvearrowright first_k(\beta\alpha) \cap first_k(\gamma\alpha) = \emptyset$$

Beweis:

1. "Definition"  $\curvearrowright$  "Lemma"  
 $\beta \neq \gamma$ , aber  $v \in first_k(\beta\alpha)$  und  $v \in first_k(\gamma\alpha)$   
 $\curvearrowright$  es existiert  $\beta\alpha \xrightarrow[l]{*} x$  und  $\gamma\alpha \xrightarrow[l]{*} y$  mit  $\{v\} = first_k(x) = first_k(y)$   
 $\curvearrowright$  (nach "Definition")  $\beta = \gamma$ , Widerspruch!
2. "Lemma"  $\curvearrowright$  "Definition"  
 "Lemma" gilt,  $first_k(X) = first_k(y)$ , aber angenommen  $\beta \neq \gamma$ :

$$first_k(x) \subseteq first_k(\beta\alpha)$$

$$first_k(y) \subseteq first_k(\gamma\alpha)$$

Widerspruch zum Lemma!

q.e.d.

**Folgerung:** Bestimmung der  $A$ -Regel durch die look-ahead-Mengen

$$first_k(\beta\alpha), first_k(\gamma\alpha), \dots \text{ für } (A \rightarrow \beta|\gamma|\dots)$$

**Problem** Abhängigkeit der look-ahead-Mengen vom Rechtskontext  $\alpha$

**Ziel** Bestimmung der look-ahead-Menge aus der Regel allein

**Idee** Mögliche Rechtskontexte vereinigen

**Definition (follow-Menge)** Sei  $\mathcal{G} \in CFG, A \in N$  und  $k \in \mathbb{N}$ . Wir definieren  $follow_k(A) \subseteq \Sigma^*$  durch

$$follow_k(A) := \{v \in \Sigma^* \mid S \xrightarrow[l]{*} wA\alpha, v \in first_k(\alpha)\}$$

Bemerkung: In  $follow_k(A)$  werden alle möglichen Rechtskontexte berücksichtigt.



**Der Fall**  $k = 1$  Im allgemeinen genügt  $k = 1$ .  $k > 1$  aufwendiger (Hinweis: ANTLR)

**Abkürzungen**  $F_i := first_1$  und  $fo := follow_1$

**Satz** Für  $\mathcal{G} \in CFG$  ( $\mathcal{G}$  reduziert) gilt:  $\mathcal{G} \in LL(1) \leftrightarrow$  Für alle Regelpaare  $A \rightarrow \beta | \gamma$  folgt:  $fi(\beta fo(A)) \cap fi(\gamma fo(A)) = \emptyset$ .

**Definition**  $la(A \rightarrow \beta) := fi(\beta fo(A))$  heißt look-ahead-Menge von  $A \rightarrow \beta$  (la-Menge).

**Beachte:**

- $fi(\alpha) \subseteq \Sigma_\varepsilon$
- $fo(A) \subseteq \Sigma_\varepsilon$
- $\beta fo(A) \subseteq \mathfrak{X}^*$
- für  $\Gamma \subseteq \mathfrak{X}^*$  ist  $fi(\Gamma) := \bigcup_{\alpha \in \Gamma} fi(\alpha)$
- $la(A \rightarrow \beta) = fi(\beta fo(A)) \subseteq \Sigma_\varepsilon$
- $\varepsilon \in fi(\beta fo(A)) \leftrightarrow \beta \xrightarrow{*} \varepsilon$  und  $\varepsilon \in fo(A)$
- $a \in fi(\beta fo(A)) \leftrightarrow a \in fi(\beta)$  oder  $(\beta \xrightarrow{*} \varepsilon$  und  $a \in fo(A))$

**Beweis:**

" $\curvearrowright$ ": Angenommen, es gibt  $A \rightarrow \beta | \gamma$  mit  $\beta \neq \gamma$  und  $c \in fi(\beta fo(A)) \cap fi(\gamma fo(A))$  für ein  $c \in \Sigma_\varepsilon$ , dann treten folgende Fälle auf:

1.  $c = \varepsilon \curvearrowright \beta \xrightarrow{*} \varepsilon, \gamma \xrightarrow{*} \varepsilon, \varepsilon \in fo(A)$

$$S \xrightarrow[t]{*} wA\alpha \begin{cases} \xrightarrow[t]{*} w\beta\alpha \\ \xrightarrow[t]{*} w\gamma\alpha \end{cases}$$

mit  $\alpha \xrightarrow[t]{*} \varepsilon \curvearrowright \varepsilon \in fi(\beta\alpha) \cap fi(\gamma\alpha)$  Widerspruch!

2.  $c \in \Sigma$ . Es folgt:

- a)  $c \in fi(\beta) \cap fi(\gamma)$
- b)  $c \in fi(\beta), \gamma \xrightarrow{*} \varepsilon$  und  $c \in fo(A)$
- c)  $\beta \xrightarrow{*} \varepsilon, c \in fo(A), c \in fi(\gamma)$
- d)  $\beta \xrightarrow{*} \varepsilon, \gamma \xrightarrow{*} \varepsilon$  und  $c \in fo(A)$

$\mathcal{G}$  ist reduziert  $\curvearrowright$  es existieren Linksableitungen

$$S \xrightarrow[\imath]{*} wA\alpha \left\{ \begin{array}{l} \xrightarrow[\imath]{} w\beta\alpha \\ \xrightarrow[\imath]{} w\gamma\alpha \end{array} \right.$$

mit  $a \in \text{fi}(\beta\alpha) \cap \text{fi}(\gamma\alpha)$ .

” $\curvearrowright$ “: Sei  $S \xrightarrow[\imath]{*} wA\alpha \left\{ \begin{array}{l} \xrightarrow[\imath]{} w\beta\alpha \\ \xrightarrow[\imath]{} w\gamma\alpha \end{array} \right.$  mit  $\beta \neq \gamma$ . Nach Voraussetzung gilt  $\text{fi}(\beta\text{fo}(A)) \cap \text{fi}(\gamma\text{fo}(A)) = \emptyset$ . Da  $\text{fi}(\beta\alpha) \subseteq \text{fi}(\beta\text{fo}(A))$ , muß auch  $\text{fi}(\beta\alpha) \cap \text{fi}(\gamma\alpha) = \emptyset$  sein.

q.e.d

### Berechnung der la-Mengen

1.  $\text{fi}(X)$  für  $X \in \mathfrak{X}$ :

- $X = a \in \Sigma \curvearrowright \text{fi}(X) = \{X\}$
- $X \rightarrow a\alpha \curvearrowright a \in \text{fi}(X)$
- $X \rightarrow \varepsilon \curvearrowright \varepsilon \in \text{fi}(X)$
- $X \rightarrow A_1 \dots A_k Y \alpha, k \geq 0, Y \in \mathfrak{X}, \varepsilon \in \text{fi}(A_1) \cap \dots \cap \text{fi}(A_k), a \in \text{fi}(Y) \curvearrowright a \in \text{fi}(X)$
- $X \rightarrow A_1 \dots A_k, k \geq 1, \varepsilon \in \text{fi}(A_1) \cap \dots \cap \text{fi}(A_k) \curvearrowright \varepsilon \in \text{fi}(X)$

2.  $\text{fi}(X_1 \dots X_n)$  für  $X_i \in \mathfrak{X}, n \in \mathbb{N}$

- $\varepsilon \in \text{fi}(X_1) \cap \dots \cap \text{fi}(X_{i-1}), a \in \text{fi}(X_i) \curvearrowright a \in \text{fi}(X_1 \dots X_n)$
- $\varepsilon \in \text{fi}(X_1) \cap \dots \cap \text{fi}(X_n) \curvearrowright \varepsilon \in \text{fi}(X_1 \dots X_n)$
- $\text{fi}(\varepsilon) = \{\varepsilon\}$

3.  $\text{fo}(A)$

- $\varepsilon \in \text{fo}(A)$ , falls  $A = S$
- $B \rightarrow \alpha A \beta, a \in \text{fi}(\beta) \curvearrowright a \in \text{fo}(A)$
- $B \rightarrow \alpha A, x \in \text{fo}(B) \curvearrowright x \in \text{fo}(A)$
- $B \rightarrow \alpha A \beta, \varepsilon \in \text{fi}(\beta), x \in \text{fo}(B) \curvearrowright x \in \text{fo}(A)$

**Lemma** Die Mengen  $\text{fi}(\alpha)$  für  $\alpha \in \mathfrak{X}^*$  und  $\text{fo}(A)$  für  $A \in N$  sind die kleinsten unter den Regeln (1) - (3) abgeschlossenen Teilmengen von  $\Sigma_\varepsilon$ .

**Beispiel**

$$\begin{aligned} \mathcal{G}'_{AE}: \quad E &\rightarrow TE' & (1) \\ E' &\rightarrow +TE'|\varepsilon & (2,3) \\ T &\rightarrow FT' & (4) \\ T' &\rightarrow *FT'|\varepsilon & (5,6) \\ F &\rightarrow (E)|a & (7,8) \end{aligned}$$

$$la(A \rightarrow \alpha) := \text{fi}(\alpha\text{fo}(A))$$

Nichtterminalsymbole	E	E'	T	T'	F
fi	(, a	+, ε	(, a	*, ε	(, a
fo	ε, )	ε, )	+, ε	+, ε, )	*, +, ε, )

la-Mengen:

1	2	3	4	5	6	7	8
(, a	+	ε, )	(, a	*	+, ε, )	(	a

**Folgerung** Die la-Mengen der Alternative (2,3), (5,6) und (7,8) sind disjunkt. Also ist  $\mathcal{G}'_{AE}$  eine LL(1)-Grammatik.

**LL(1)-Test:** la-Mengen berechnen und Alternativen auf Disjunktheit prüfen.

**Deterministischer TD-Analyseautomat DTA( $\mathcal{G}$ ) für  $\mathcal{G} \in LL(1)$**

**Idee:** Die Zugehörigkeit des Eingangssymbols zu einer la-Menge steuert die Regelauswahl; 1-look-ahead auf der Eingabe

Modifikation des Ableitungsschrittes:

- $(aw, A\alpha, z) \vdash (aw, \beta\alpha, zi)$ , falls  $\pi_i = A \rightarrow \beta$  und  $a \in la(\pi_i)$
- $(\varepsilon, A\alpha, z) \vdash (\varepsilon, \beta\alpha, zi)$ , falls  $\pi_i = A \rightarrow \beta$  und  $\varepsilon \in la(\pi_i)$

**Folge:** Deterministische Arbeitsweise

**Beachte:** Das Eingangssymbol wird bei Ableitungsschritten nicht gelöscht

**Beispiel**

$$\begin{aligned} &\vdots \\ A &\rightarrow BC(1) \\ B &\rightarrow bB|\varepsilon \\ C &\rightarrow cC|\varepsilon \\ &\vdots \end{aligned}$$

$$\begin{aligned}
la(A \rightarrow BC) &= fi(BCfo(A)) \\
fi(B) &= \{b, \varepsilon\} \\
fi(C) &= \{c, \varepsilon\} \\
la(1) &= \{b, c\} \cup fo(A)
\end{aligned}$$

**action-Funktion** Darstellung des  $DTA(\mathcal{G})$  durch die Analysetabelle von  $\mathcal{G}$ , auch action-Funktion genannt.

$$act : \Sigma_\varepsilon \times (N \cup \Sigma_\varepsilon) \rightarrow \{\alpha | A \rightarrow \alpha \text{ in } \mathcal{G}\} \times [p] \cup \{pop, error, accept\}$$

$$\begin{aligned}
act(x, A) &:= (\alpha, i) \text{ falls } \pi_i : A \rightarrow \alpha \wedge x \in la(\pi_i) \\
act(a, a) &:= pop \\
act(\varepsilon, \varepsilon) &:= accept \\
act(x, X) &:= error \text{ sonst}
\end{aligned}$$

### Parser-Konstruktion nach TD-Methode

- $\mathcal{G} \in \text{CFG}$
- Berechnung der  $la$ -Mengen (fi- und fo-Mengen)
- Analysetabelle, Eindeutigkeit prüfen
- tabellengesteuerter Parser  $\rightsquigarrow$  ANTLR
- Problem: Mehrdeutige Analysetabelle ( $\mathcal{G} \notin LL(1)$ )

**Transformation nach  $LL(1)$**  Es gibt 2 Methoden,  $\mathcal{G}$  in eine äquivalente  $LL(1)$ -Grammatik zu transformieren; dies ist jedoch nicht immer möglich.

1. Beseitigung von Linksrekursion
2. Linksfaktorisieren

Dazu werden parsererzeugende Systeme verwendet. Es ist jedoch Vorsicht geboten, da die Transformationen zwar die Äquivalenz erhalten, i.a. aber nicht die syntaktische Struktur.

1. Beseitigung von Linksrekursion

**Definition:**  $\mathcal{G} \in \text{CFG}$  heißt linksrekursiv:  $\leftrightarrow \exists \forall A \in N, \alpha \in \mathfrak{X}^* : A \stackrel{\pm}{\Rightarrow} A\alpha$

**Folgerung:**  $\mathcal{G}$  linksrekursiv  $\rightsquigarrow \mathcal{G} \notin LL(k)$  für alle  $k \in \mathbb{N}$

**Grund:** Wenn ein TD-Parser  $A \stackrel{\pm}{\Rightarrow} A$  simuliert, so bleibt der Eingabekopf stehen,

gleicher look-ahead. Schleife. Keine Ableitung der Form  $S \xrightarrow{*} wA\beta \xrightarrow{\pm} wA\alpha\beta \xrightarrow{*} wv$  simulierbar.

**Beispiel:**

$$\begin{aligned} \mathcal{G}_{AE} : \quad E &\rightarrow E + T|T & (1, 2) \\ T &\rightarrow T * F|F & (3, 4) \\ F &\rightarrow (E)|a & (5, 6) \end{aligned}$$

$$\text{fi}(F) = \{(, a\} = \text{fi}(T) = \text{fi}(E)$$

$$la(1) = \{(, a\} = la(2) = la(3) = la(4)$$

Darauf folgt:  $\mathcal{G}_{AE} \notin LL(1)$ .

**Spezialfall:** direkte Linksrekursion und ihre Beseitigung

$A \rightarrow A\alpha|\beta$  mit  $\beta \neq \alpha \dots$  und  $\alpha \neq \varepsilon$  wird ersetzt durch  $A \rightarrow \beta A'$  mit neuem  $A'$  und  $A' \rightarrow \alpha A'|\varepsilon$ .

**Folgerung:**  $L(\mathcal{G})$  ist unverändert, bekommt jedoch eine neue syntaktische Struktur. Dies ist kein Problem bei assoziativen Operatoren wie  $*$  bzw.  $+$ .

Die Semantik ist invariant.

$$\begin{aligned} \mathcal{G}_{AE} : \quad E &\rightarrow TE' \\ E' &\rightarrow +TE'|\varepsilon \\ T &\rightarrow FT' \\ T' &\rightarrow *FT'|\varepsilon \\ F &\rightarrow (E)|a \end{aligned}$$

**Allgemeiner Fall:** indirekte Linksrekursion

$$\begin{aligned} A &\rightarrow A_1\alpha_1|\dots \\ A_1 &\rightarrow A_2\alpha_2|\dots \\ &\vdots \\ A_n &\rightarrow A\beta|\dots \end{aligned}$$

Beseitigung durch Transformation in GNF (Greibach-Normalform)

$$\begin{aligned} \text{GNF} : \quad A &\rightarrow aB_1 \dots B_n \quad (B_i \neq S) \\ S &\rightarrow \varepsilon \end{aligned}$$

**Beachte:** Die Beseitigung von Linksrekursion ergibt nicht notwendigerweise  $\mathcal{G}' \in LL(1)$ .

**Grund:** Jedes  $\mathcal{G} \in \text{CFG}$  ist äquivalent transformierbar in GNF, aber nicht jede  $L \in \text{CFL}$  ist durch  $\mathcal{G} \in LL(1)$  erzeugbar.

**Es gilt:**  $\mathcal{L}(LL(k)) \subsetneq \mathcal{L}(LL(k+1)) \subsetneq \mathcal{L}(\text{DPDA}) \subsetneq \text{CFL}$

**Komplexität der  $LL(1)$ -Analyse**  $\mathcal{G}LL(1)$ -Grammatik  $\curvearrowright \mathcal{G}$  nicht linksrekursiv  
DTA( $\mathcal{G}$ ) mit Eingabe  $w \in \Sigma^*$ :

- $|w|$  Vergleichsschritte und 1 accept-Schritt
  - Frage: Wieviele Schritte  $i$  kann der DTA( $\mathcal{G}$ ) durchführen, ohne den Eingabekopf zu verschieben?  $wA\alpha \xrightarrow{i} wB\beta$
- Aho/Ullman - The Theory Of Parsing, ..., Vol. I, Parsing, S.356:  $\exists c \in \mathbb{N}$  nicht linksrekursiv  $\curvearrowright A \xrightarrow{i} B\gamma \curvearrowright i < c = |N| \underset{\text{falsch}}{\curvearrowright} (S \xrightarrow{n} w \curvearrowright n \leq$

$$c(|w| + 1)$$

Besser: O. Mayer, Syntaxanalyse, S.48-55:

**Satz:**  $\mathcal{G}$  nicht linksrekursiv:  $\exists c \in \mathbb{N} \ S \xrightarrow[l]{n} w \rightsquigarrow n \leq c(|w| + 1)$

2. Links-Faktorisierung:

**Beispiel:**  $\text{statement} \rightarrow \text{if cond then stat else stat fi}$  Keine Regelentscheidung  
 $\text{stat} \rightarrow \text{if cond then stat fi}$   
 mit beschränktem look-ahead.

**Idee:** Verschieben der Entscheidung bis Alternative erkennbar. Links-Faktorisieren:

$A \rightarrow \alpha\beta|\alpha\gamma$  ersetzen durch

$A \rightarrow \alpha A'$  und

$A' \rightarrow \beta|\gamma$

**Beispiel:**  $\text{stat} \rightarrow \text{if cond then stat } S'$   
 $S' \rightarrow \text{else stat fi} \text{ — fi}$

### Top-Down-Analyse mit rekursiven Prozeduren

**Idee:** Keine explizite Kellerbenutzung wie im TDA( $\mathcal{G}$ ), sondern Verwendung des Laufzeitkellers durch Einsatz rekursiver Prozeduren

**Vorteil:** Leichte Programmierung

**Spezialfall:**  $\mathcal{G} \in LL(1)$  - Beispiel  $\mathcal{G}'_{AE}$   
 "Recursive Descent Parser" (Analyse durch rekursiven Abstieg)

**Methode:**  $A \in N \rightarrow A()$  Parameterlose Prozedur  $r$  (zur Simulation des Ableitungsschrittes)  
 Alternativen durch Eingabesymbol unterscheidbar, da  $\mathcal{G} \in LL(1)$ .

**Eingabe:**  $sym$  Variable für das Eingabesymbol  
 $nextsym$  Zum Lesen des nächsten Eingabesymbols

**Ausgabe:**  $print(i)$  zur Ausgabe einer Regelnummer

**Folie:** Beispiel zu  $\mathcal{G}'_{AE}$

## 3.2 Bottom-Up-Analyse, $LR(k)$ -Grammatiken

**Idee:** Bottom-Up-Berechnung des Ableitungsbaumes in Form einer gespiegelten Rechtsanalyse durch einen Kellerautomaten:

**Shift-Schritte:** Verschieben von Eingabesymbolen auf den Keller

**Reduce-Schritte:** Umkehrung von Ableitungs- zu Reduktionsschritten  $\rightarrow$  Shift-Reduce-Verfahren

**Definition:** Der nichtdeterministische Bottom-Up-Analyseautomat von  $\mathcal{G}$ , kurz:  $\text{NBA}(\mathcal{G})$ :

**Eingabealphabet:**  $\Sigma$

**Kelleralphabet:**  $\mathfrak{X}$

**Ausgabealphabet:**  $[p]$

**Konfigurationsmenge:**  $\mathfrak{X}^* \times \Sigma^* \times [p]^*$

**Transitionen:** • Shift-Schritt:  $(\alpha, aw, z) \vdash (\alpha a, w, z)$  für  $a \in \Sigma$   
 • Reduce-Schritt:  $(\beta\alpha, w, z) \vdash (\beta A, w, zi)$  falls  $\pi_i = A \rightarrow \alpha$

**Ausgangskonfiguration für  $w \in \Sigma^*$ :**  $(\varepsilon, w, \varepsilon)$

**Endkonfigurationen:**  $(S, \varepsilon, z)$

**Satz:** Der  $\text{NBA}(\mathcal{G})$  berechnet gespiegelte Rechtsanalysen, d.h. für  $w \in \Sigma^*$  und  $z \in [p]^*$  gilt:  $z$  ist eine r-Analyse von  $w \leftrightarrow (\varepsilon, w, \varepsilon) \xrightarrow{*} (s, w, \overleftarrow{z})$

**Nicht-Determinismus des  $\text{NBA}(\mathcal{G})$**

1. Shift- oder Reduce-Schritt
2. Reduce-Schritt: linker Henkelrand
3. Reduce-Schritt: linke Regelseite
4. Analyseende

**Ziel:** Nicht-Determinismus durch  $k$ -look-ahead auf der Eingabe beseitigen.  $LR(k)$ -Grammatiken

**Generelle Voraussetzung**  $\mathcal{G}$  startsepariert, d.h.  $S$  nur in  $S \rightarrow A$  mit  $A \neq S$ . Jedes  $\mathcal{G} \in \text{CFG}$  startsepariert machen durch Hinzufügen von  $S' \rightarrow S$  (neues  $S'$ ). Im folgenden nur startseparierte Grammatiken mit  $S' \rightarrow S(0)$ .

**Folgerung:**  $(S', \varepsilon, z)$  ist eine Endkonfiguration, die zwar Folgekonfiguration haben kann ( $\varepsilon$ -Regeln) ohne Übergang in weitere Endkonfigurationen (Kellerlänge  $\geq 2$ ). Beseitigung des restlichen Nicht-Determinismus durch

**Definition ( $LR(k)$ -Grammatiken:** Sei  $\mathcal{G} \in \text{CFG}$ , startsepariert durch  $S' \rightarrow S, k \in \mathbb{N}$ .  $\mathcal{G} \in LR(k) :\leftrightarrow$  Für alle Rechtsableitungen der Form

$$S' \Rightarrow S \left\{ \begin{array}{l} \xrightarrow[r]{*} \alpha A w \xrightarrow[r]{*} \alpha \beta w \\ \xrightarrow[r]{*} \alpha' A' v' \xrightarrow[r]{*} \alpha \beta v \end{array} \right.$$

mit  $\text{first}_k(w) = \text{first}_k(v)$  gilt:  $\alpha = \alpha', A' = A$  und  $v' = v$ .

**Folgerung:** Der  $\text{NBA}(\mathcal{G})$  kann mit  $k$ -look-ahead auf der Eingabe die Transition entscheiden.

**LR(0)-Grammatiken**  $k = 0 \curvearrowright$  Entscheidung ohne look-ahead allein durch Kellerinhalt. Abstraktion endlicher Information aus  $\alpha\beta$ , welche für die Entscheidung ausreicht.

**Definition (LR(0)-Auskünfte, LR(0)-Mengen:** Sei  $\mathcal{G} \in \text{CFG}$  mit  $S' \rightarrow S$  und  $S' \xrightarrow[r]{*} \alpha A w \Rightarrow \alpha \beta_1 \beta_2 w$ . Dann heißt  $[A \rightarrow \beta_1 \cdot \beta_2]$  eine LR(0)-Auskunft für  $\alpha\beta_1$ .

Für  $\gamma \in \mathfrak{X}^*$  bezeichnet  $LR(0)(\gamma)$  die Menge aller LR(0)-Auskünfte für  $\gamma$ , die sogenannte LR(0)-Menge von  $\gamma$ .

**Folgerung:**

1.  $LR(0)(\gamma)$  endlich
2.  $LR(0)(\mathcal{G}) =: \{LR(0)(\gamma) \mid \gamma \in \mathfrak{X}^*\}$  endlich
3.  $[A \rightarrow \beta_1 \cdot] \in LR(0)(\gamma)$  signalisiert Reduktionsmöglichkeit  $(\alpha\beta_1, w, z) \vdash (\alpha A, w, zi)$  für  $\pi_i = A \rightarrow \beta_1$  und  $\gamma = \alpha\beta_1$ .
4.  $[A \rightarrow \beta_1 \beta_2] \in LR(0)(\gamma)$  mit  $\beta_2 \neq \varepsilon$  bedeutet Shift-Möglichkeit wegen unvollständigem Henkel.
5.  $\mathcal{G} \in LR(0) \leftrightarrow$  Die LR(0)-Mengen enthalten keine widersprüchlichen Auskünfte.

**Berechnung der LR(0)-Mengen einer Grammatik**

**Satz**  $\mathcal{G} \in \text{CFG}$  mit  $S' \rightarrow S, \mathcal{G}$  reduziert. Dann gilt

1.  $LR(0)(\varepsilon)$  ist die kleinste Menge, welche
  - a)  $[S' \rightarrow \cdot S]$  enthält
  - b) mit  $[A \rightarrow \cdot B \delta]$  und  $B \rightarrow \beta$  in  $\mathcal{G}$  auch  $[B \rightarrow \cdot \beta]$  enthält.
2.  $LR(0)(\alpha X)$  mit  $X \in \mathfrak{X}$  ist die kleinste Menge, welche
  - a)  $[A \rightarrow \beta_1 X \cdot \beta_2]$  enthält, falls  $[A \rightarrow \beta_1 \cdot X \beta_2] \in LR(0)(\alpha)$
  - b) und mit  $[A \rightarrow \gamma \cdot B \delta]$  und  $B \rightarrow \beta$  in  $\mathcal{G}$  auf  $[B \rightarrow \cdot \beta]$  enthält.

**Wiederholung:**  $\mathcal{G} \in LR(0) \leftrightarrow$  Für alle Ableitungen  $S' \xrightarrow[r]{*} \alpha A w \Rightarrow \alpha \beta w$  und  $S' \xrightarrow[r]{*} \alpha' A' v' \Rightarrow \alpha \beta v$  gilt:

$$\alpha' = \alpha, \quad A' = A, \quad v' = v$$



## Beispiel

$$A \rightarrow \beta_1 B \beta_2, \quad B \rightarrow \gamma$$

$$S' \xrightarrow[r]{*} \alpha A w \xrightarrow[r]{(1)} \alpha \beta_1 B \beta_2 w \xrightarrow[r]{*} \alpha \beta_1 w' w \xrightarrow[r]{(2)} \alpha \beta_1 \gamma w' w$$

1. a)  $[A \rightarrow \beta_1 B \beta_2 \cdot] \in LR(0)(\alpha \beta_1 B \beta_2)$
- b)  $[A \rightarrow \cdot \beta_1 B \beta_2] \in LR(0)(\alpha)$
- c)  $[A \rightarrow \beta_1 \cdot B \beta_2] \in LR(0)(\alpha \beta_1)$
- d)  $[A \rightarrow \beta_1 B \cdot \beta_2] \in LR(0)(\alpha \beta_1 B)$

Folgerung: Auskünfte längerer Keller folgen durch Punktverschiebung aus Auskünften kürzerer Keller.

2. c)  $[B \rightarrow \cdot \gamma] \in LR(0)(\alpha \beta_1)$   
 Folgerung:  $[A \rightarrow \beta_1 \cdot B \beta_2] \in LR(0)(\alpha \beta_1), \quad B \rightarrow \gamma \quad \curvearrowright [B \rightarrow \cdot \gamma] \text{ in } LR(0)(\alpha \beta_1)$

**Die goto-Funktion**  $\mathcal{G} \in LR(0) \curvearrowright LR(0)(\gamma)$  liefert eine Shift/Reduce-Entscheidung für den Bottom-Up-Analyseautomaten mit Kellerinschrift  $\gamma$ .

- Neues Kellularphabet:  $LR(0)(\mathcal{G})$  statt  $\mathfrak{X}$
- Beachte:  $LR(0)(\gamma X)$  ist bestimmt durch  $LR(0)(\gamma)$  und  $X$ , aber unabhängig von  $\gamma$ .
- goto:  $LR(0)(\gamma) \times \mathfrak{X} \rightarrow LR(0)(\mathcal{G})$  ist definiert durch  $goto(I, X) := I' \leftrightarrow \exists \gamma \in \mathfrak{X}^*$  mit  $I = LR(0)(\gamma)$  und  $I' = LR(0)(\gamma X)$ .

**Berechnung der  $LR(0)$ -Mengen und goto-Funktion durch Potenzmengenkonstruktion eines NFA** Sei  $\mathcal{G} \in \text{CFG}$  startsepariert mit  $S' \rightarrow S$ . Konstruktion eines  $\mathfrak{A}(\mathcal{G}) \in \text{NFA}_\varepsilon$ .

- Zustandsmenge  $Q = \{[A \rightarrow \beta_1 \cdot \beta_2] \mid A \rightarrow \beta_1 \beta_2 \in \mathcal{G}\}$
- Eingabealphabet  $\mathfrak{X} = N \cup \Sigma$
- Anfangszustand  $q_0 := [S' \rightarrow \cdot S]$
- Endzustandsmenge  $F := Q$  (ohne Bedeutung)
- Transitionsfunktion:  $\delta : Q \times \mathfrak{X}_{\text{veps}} \rightarrow \mathcal{P}(Q)$

$$\delta([A \rightarrow \beta_1 \cdot X \beta_2], x) \ni [A \rightarrow \beta_1 X \cdot \beta_2]$$

$$\delta([A \rightarrow \beta_1 \cdot B \beta_2], \varepsilon) \ni [B \rightarrow \cdot \gamma], \text{ falls } B \rightarrow \gamma \in \mathcal{G}$$

**Potenzmengenkonstruktion (Thompson)** Konstruktion von  $\mathfrak{A}(\hat{\mathcal{G}}) = \langle \hat{Q}, \mathfrak{X}, \hat{\delta}, \hat{q}_0 \emptyset \rangle \in \text{DFA}$

Erweiterte Transitionsfunktion:

- $\bar{\delta} : \mathcal{P}(Q) \times \mathfrak{X}^* \rightarrow \mathcal{P}(Q)$
- $\bar{\delta}(T, \varepsilon) := \varepsilon(T)$
- $\bar{\delta}(T, wa) := \varepsilon\left(\bigcup_{q \in \bar{\delta}(T, w)} \delta(q, a)\right)$
- $\hat{Q} := \{\bar{\delta}(\{[S' \rightarrow \cdot S]\}, \alpha) \mid \alpha \in \mathfrak{X}^*\}$
- $\hat{q}_0 := \varepsilon(\{[S' \rightarrow \cdot S]\})$
- $\hat{\delta}(T, X) := \bar{\delta}(T, X)$

Dann gilt:  $\hat{Q} = LR(0)(\mathcal{G})$  und  $\hat{\delta} = goto$ .

**Konstruktion des deterministischen BU-Analyseautomaten für  $\mathcal{G} \in LR(0)$**  Hilfsmittel:  $LR(0)(\mathcal{G})$  und goto-Funktion

Die action-Funktion von  $\mathcal{G}$  gibt die Shift/Reduce-Entscheidung an:

$$act : LR(0)(\mathcal{G}) \rightarrow \{shift, red\ i, error, accept \mid i \in [p]\}$$

$$act(I) := \begin{cases} red\ i & \text{falls } \pi_i = A \rightarrow \alpha \text{ und } [A \rightarrow \alpha \cdot] \in I \\ shift & \text{falls } [A \rightarrow \alpha \cdot X\beta] \in I \\ accept & \text{falls } [S' \rightarrow S \cdot] \in I \\ error & \text{falls } I = \emptyset \end{cases}$$

**DBA für  $\mathcal{G} \in LR(0)$**

$$act : LR(0)(\mathcal{G}) \rightarrow \{shift, red\ i, error, accept \mid i \in [p]\}$$

Die act und goto-Funktionen bilden  $LR(0)$ -Analysetabelle von  $\mathcal{G}$ . Die Tabelle bestimmt  $LR(0)(\mathcal{G})$ -Analyseautomaten.

- Eingabealphabet:  $\Sigma$
- Kelleralphabet:  $\Gamma := LR(0)(\mathcal{G})$
- Ausgabealphabet:  $\Delta := [p] \cup \{error\} \cup \{0\}$
- Konfigurationsmenge:  $\Gamma^* \times \Sigma^* \times \Delta^*$
- Bezeichnung:  $wa_1 \dots a_n - (n) = w$
- Transitionen:
  - Shift-Schritt:  $(\alpha I, aw, z) \vdash (\alpha II', w, z)$  falls  $act(I) = shift$  und  $goto(I) = I'$

- Reduce-Schritt:  $(\alpha I, w, z) \vdash (\tilde{\alpha} \tilde{I} \tilde{I}', w, zi)$  falls  $act(I) = red$   $i$ ,  $\pi_i = A \rightarrow X_1, \dots, X_n, X_i \in \mathfrak{X}$  und  $\alpha I - (n) = \tilde{\alpha} \tilde{I}$  und  $goto(\tilde{I}, A) = I'$
- Accept-Schritt:  $(I_0 I, \varepsilon, z) \vdash (\varepsilon, \varepsilon, z0)$  falls  $act(I) = accept$
- Fehlererkennung:  $(\alpha I, w, u) \vdash (\varepsilon, \varepsilon, z \text{ error})$  in allen anderen Fällen

- Anfangskonfiguration für  $w \in \Sigma^*$ :  $(I_0, w, \varepsilon)$ ,  $I_0 = LR(0)(\varepsilon)$
- Folgerung: Wenn  $LR(0)$ -Mengen konfliktfrei sind, also  $act$  eindeutig, so arbeitet der  $LR(0)$ -Analyseautomat deterministisch, und es gilt: Für  $w \in \Sigma^*$  und  $z \in ([p] \cup \{0\})^*$

- $(I_0, w, \varepsilon) \stackrel{A}{\vdash} (\varepsilon, \varepsilon, z) \leftrightarrow \overleftarrow{z}$  r-Analyse von  $w$
- $(I_0, w, \varepsilon) \stackrel{*}{\vdash} (\varepsilon, \varepsilon, z \cdot error) \leftrightarrow w \notin L(\mathcal{G})$

- Beispiel:  $w = aac$

$$\begin{array}{lll}
(I_0, aac, \varepsilon) & \vdash & (I_0 I_4, ac, \varepsilon) \quad \vdash & (I_0 I_4 I_4, c, \varepsilon) \\
\vdash & (I_0 I_4 I_4 I_6, \varepsilon, \varepsilon) & \vdash & (I_0 I_4 I_4 I_8, \varepsilon, 6) \quad \vdash & (I_0 I_4 I_8, \varepsilon, 65) \\
\vdash & (I_0 I_3, \varepsilon, 655) & \vdash & (I_0 I_1, \varepsilon, 6552) \quad \vdash & (\varepsilon, \varepsilon, 65520)
\end{array}$$

$$\overleftarrow{z} = 02556$$

$$S' \xrightarrow[r]{0} S \xrightarrow[r]{2} C \xrightarrow[r]{5} aC \xrightarrow[r]{5} aaC \xrightarrow[r]{6} aac = w$$

- Praxis: Widersprüchliche Informationen /  $LR(0)$ -Mengen  $\rightsquigarrow \mathcal{G} \notin LR(0) \rightsquigarrow SLR(1)$ -Analyse
- Beachte: Shift/Reduce-Konflikte in  $I_1, I_2, I_9$  (CB 2-16)  
"Trick": Beseitigung des Konfliktes durch Eingabesymbol

### Beobachtung:

1.  $[A \rightarrow \beta_1 a \beta_2] \in LR(0)(\alpha \beta_1) \rightsquigarrow S \xrightarrow[r]{*} \alpha A w \Rightarrow \alpha \beta_1 a \beta_2 w$   
Also: Shift nur bei Eingabe von  $a$
2.  $[A \rightarrow \beta \cdot] \in LR(0)(\alpha \beta) \rightsquigarrow S \xrightarrow[r]{*} \alpha A a w \Rightarrow \alpha \beta a w \rightsquigarrow a \in \text{fo}(A)$   
Reduktion mit  $A \rightarrow \beta$  nur falls  $a \in \text{fo}(A)$ .

Für obiges Beispiel (CB 2-16):

- $I_1$ : shift bei Eingabe von  $+$
- $I_1$ : accept bei Eingabe von  $\$, \text{fo}(E') = \{\varepsilon\}$
- $I_2$ : shift bei Eingabe von  $*$

- $I_2$ : reduce 2 bei Eingabe von +, ), \$,  $\text{fo}(E) = \{+, ), \varepsilon\}$
- $I_9$ : shift bei Eingabe von \*
- $I_9$ : reduce 1 bei Eingabe von +, ), \$,  $\text{fo}(E) = \{+, ), \varepsilon\}$

Konflikte beseitigt. Weiterer Vorteil: Frühere Fehlererkennung (genauere Kontrolle der Aktionen).

### SLR(1)-action-Funktion (simple)

$$act : LR(0)(\mathcal{G}) \times (\Sigma \cup \{\$\}) \rightarrow \{shift, red\ i, accept, error | i \in [p]\}$$

sei definiert durch:

$$act(I, a) := \begin{cases} shift & \text{falls } [A \rightarrow \alpha \cdot a\beta] \in I \\ red\ i & \text{falls } \pi_i : A \rightarrow \alpha, [A \rightarrow \alpha \cdot] \in I, a \in \text{fo}(A) \\ accept & \text{falls } [S' \rightarrow S \cdot] \in I \text{ und } a = \$ \\ error & \text{sonst} \end{cases}$$

**Definition:**  $\mathcal{G} \in SLR(1) :\leftrightarrow act(I, a)$  ist stets eindeutig. action und goto bilden  $SLR(1)$ -Analysetabelle von  $\mathcal{G}$ . (CB 2-17)

Bemerkung:

- Tabellenkompression
- Bessere Konfliktbeseitigung durch Verwendung des lookahead-Symbols in der  $[A \rightarrow \beta_1 \cdot \beta_2 a] \rightsquigarrow LR(1), LALR(1)$

(CB 2-18) Konflikt in  $I_2$  nicht nach der SLR-Methode lösbar:  $= \in \text{fo}(R)$ , aber nicht jedes Element von  $\text{fo}(R)$  in beliebiger r-Ableitung hinter  $R$  möglich.  $\curvearrowright$  Verfeinerung der  $LR(0)$ -Auskünfte durch Mitführen der möglichen la-Symbole.

### Definition: ( $LR(1)$ -Auskünfte und $LR(1)$ -Mengen für $\mathcal{G} \in \text{CFG}$ )

1. Wenn  $S \xrightarrow[r]{*} \alpha A a w \Rightarrow \alpha \beta_1 \beta_2 a w$ , so  $[A \rightarrow \beta_1 \cdot \beta_2, a] \in LR(1)(\alpha \beta_1)$
2. Wenn  $S \xrightarrow[r]{*} \alpha A \Rightarrow \alpha \beta_1 \beta_2$ , so  $[A \rightarrow \beta_1 \cdot \beta_2, \$] \in LR(1)(\alpha \beta_1)$

$$LR(1)(\mathcal{G}) := \{LR(1)(\gamma) | \gamma \in \mathfrak{X}^*\}$$

Beobachtung:  $LR(1)(\gamma)$  "enthält in der ersten Komponente" gerade  $LR(0)(\gamma)$ .

**Berechnung der  $LR(1)$ -Mengen** Modifikation der Berechnung von  $LR(0)(\mathcal{G})$  unter Berücksichtigung des r-Kontextes.

- $LR(1)(\varepsilon)$ :
  - $[S' \rightarrow \cdot S, \$]$
  - Wenn  $[A \rightarrow \cdot B\delta, x] \in LR(1)(\varepsilon)$ ,  $B \rightarrow \beta$  in  $\mathcal{G}$  und  $y \in \text{fi}(\delta x)$   $[B \rightarrow \cdot \beta, y]$
- $LR(1)(\alpha X)$ 
  - Wenn  $[A \rightarrow \beta_1 \cdot X\beta, x] \in LR(1)(\alpha)$ , so  $[A \rightarrow \beta_1 X \cdot \beta_2, x] \in LR(1)(\alpha X)$
  - Wenn  $[A \rightarrow \gamma \cdot B\delta, x] \in LR(1)(\alpha X)$ ,  $B \rightarrow \beta$ ,  $y \in \text{fi}(\delta x)$ , so  $[B \rightarrow \cdot \beta, y] \in LR(1)(\alpha X)$

**Die  $LR(1)$ -action-Funktion von  $\mathcal{G}$**

$$act : LR(1)(\mathcal{G}) \times (\Sigma \cup \{\$\}) \rightarrow \{shift, red\ i, accept, error \mid 1 \leq i \leq r\}$$

sei definiert durch

$$act(I, x) := \begin{cases} red\ i & \text{falls } \pi_i = A \rightarrow \alpha \text{ und } [A \rightarrow \alpha \cdot, x] \in I \\ accept & \text{falls } x = \$ \text{ und } [S' \rightarrow S \cdot, \$] \in I \\ shift & \text{falls } x \neq \$ \text{ und } [A \rightarrow \alpha_1 \cdot x\alpha_2, y] \in I \\ error & \text{sonst} \end{cases}$$

Dann gilt:  $\mathcal{G} \text{ in } LR(1) \leftrightarrow act$  eindeutig (keine Konflikte)

**$LALR(1)$ -Analyse** Beseitigung von Entscheidungskonflikten nach  $LR(1)$  zu aufwendig (s. CB 2-19)

Beobachtung: Mögliche Informationsredundanz bei  $LR(1)(\mathcal{G})$ . Jede  $LR(1)$ -Menge "enthält" eine  $LR(0)$ -Menge.

**Definition:**  $I_1, I_2 \in LR(1)(\mathcal{G})$  heißen  $LR(0)$ -äquivalent,  $\sim_0$ :  $\leftrightarrow$  Die  $LR(0)$ -Anteile von  $I_1$  und  $I_2$  sind gleich.

Beispiel:  $I_4 \sim_0 I_{11}$ ,  $I_5 \sim_0 I_{12}$ ,  $I_7 \sim_0 I_{13}$ ,  $I_8 \sim_0 I_{10}$  (CB2 – 19)

**Folgerung:**  $|LR(1)(\mathcal{G}) / \sim_0| = |LR(0)(\mathcal{G})|$

Oft können  $LR(0)$ -äquivalente  $LR(1)$ -Informationen vereinigt werden, ohne daß die Konfliktlösbarkeit verloren geht. Insbesondere können nur RR-Konflikte durch Vereinigung entstehen! Eventuell entsteht dadurch eine spätere Fehlererkennung.

**Definition:**  $I \in LR(1)(\mathcal{G})$  bestimmt eine  $LALR(1)$ -Menge  $\bigcup \{I' \in LR(1)(\mathcal{G}) \mid I' \sim_0 I\}$ .

$LALR(1)(\mathcal{G})$ : Menge der  $LALR(1)$ -Mengen von  $\mathcal{G}$ .

Es gilt:  $|LALR(1)(\mathcal{G})| = |LR(0)(\mathcal{G})|$

Die  $LALR(1)$ -action-Funktion von  $\mathcal{G}$  ist analog zu  $LR(1)$  definiert.

**Definition:**  $\mathcal{G} \in LALR(1) \Leftrightarrow LALR(1)$ -action-Funktion eindeutig.

Die  $LR(1)$ -goto-Funktion überträgt sich auf  $LALR(1)(\mathcal{G})$ , weil für  $LR(1)$ -Mengen  $I_1$  und  $I_2$  gilt:  $I_1 \underset{0}{\sim} I_2 \leadsto goto(I_1, X) \underset{0}{\sim} goto(I_2, X)$

Grund: Der  $LR(0)$ -Kern von  $LR(1)(\alpha X)$  ist durch den  $LR(0)$ -Kern von  $LR(1)(\alpha X)$  und  $X$  vollständig bestimmt.

### 3.2.1 Bottom-Up-Analyse mehrdeutiger Grammatiken

Es gilt für  $\mathcal{G} \in CFG$ :  $\mathcal{G}$  mehrdeutig  $\leadsto \mathcal{G} \notin LR = \bigcup_{k \in \mathbb{N}} LR(k)$

Mehrdeutigkeit ist natürliches Beschreibungsmittel von Programmiersprachen zur Vermeidung aufwendiger Klammerung. Auflösung durch Regeln für Präzedenz und Assoziativität von Op-Symbolen.

1. Beispiel:

$$\mathcal{G}_{AE}^m : E \rightarrow E + E | E * E | (E) | a$$

Präzedenz:                   \* vor +

Assoziativität:           links

Konflikte (CB 2-25):  $I_1$   $SLR(1)$ -lösbar;  $I_7, I_8$  nicht auflösbar wegen Mehrdeutigkeit

Beispielrechnung:  $a + a * a$

$I_0$	$a + a * a$
$I_0 I_3$	$+ a * a$
$I_0 I_1$	$+ a * a$
$I_0 I_1 I_4$	$a * a$
$I_0 I_1 I_4 I_3$	$* a$
$I_0 I_1 I_4 I_7$	$* a$

$\{+, *\} \subseteq \text{fo}(E) \leadsto$  keine Konfliktlösung; Präzedenz:

a)  $* \succ + \leadsto \text{act}(I_7, *) = \text{shift}$

b)  $+ \succ * \leadsto \text{act}(I_7, *) = \text{red } 1$

2. Beispiel (CB 2-27): Mehrdeutigkeit bei Verzweigungen ("dangling else")

$$S \rightarrow iSeS | iS | a$$

Konflikt bei  $I_4 : e \in \text{fo}(S) \leftrightarrow$  nicht  $LR$ -lösbar: 2 Zerlegungen

a) if b then (if b then a else a) [korrekt]

b) if b then (if b then a) else a [nicht!]

K. John Gough: Syntax Analysis and Software Tools, Addison-Wesley, 1988

## 4 Semantische Analyse, Attributgrammatiken

- Ergebnis der syntaktischen Analyse: Ableitungsbaum
- Kontextabhängige Grammatiken (nicht durch CFG beschreibbar):
  - Deklariertheit von Bezeichnern (vgl. ATFS:  $\{ww|w \in \Sigma^* \setminus CFL\}$ )
  - Typinformation
- Festlegung dieser Eigenschaften:
  - Gültigkeitsregeln: Gültigkeitsbereich einer Deklaration
  - Sichtbarkeitsregeln: Sichtbar im Gültigkeitsbereich (Überdeckung globaler durch lokale Deklaration)
  - Typvorschriften: Typkonsistenz

**Statische Grammatik:** kontextabhängige, laufzeitunabhängige Eigenschaften eines Programms

- Formale Beschreibung durch Attributgrammatiken
- Idee (Knuth): CFG + semantische Regeln  $\rightsquigarrow$  Zusatzinformationen für Ableitungsbaum
- Semantische Analyse = Attributberechnung  
ihr Ergebnis: Attributierter Ableitungsbaum
- Grundlage für die anschließende Synthesephase

### Attributgrammatiken

- Idee: Attribute für  $A \in N$ , semantische Regeln für ihre Berechnung
- synthetische Attribute: BU-Berechnung
- inherite Attribute: TD-Berechnung  $\curvearrowright$  beliebiger Informationstransfer im Ableitungsbaum
- Attributwerte: Symboltabellen, Typen, Code, Fehler, ...; breite Anwendbarkeit von Attributgrammatiken (AG), syntaxgerichtete Programmierung, automatische Attributauswertung in Compilergeneratoren

CB:

- statische Semantik
- Programmanalyse für Optimierung
- Codegenerierung
- Fehlerbehandlung

**Historisch:** Knuth [68]: Semantics of CFLs

**Beispiel (Binärzahlen)  $\mathcal{G}_B$ :**

Bits	$B \rightarrow 0$	$v.0 = 0$
	$B \rightarrow 1$	$v.0 = 1$
Lists	$L \rightarrow B$	$v.0 = v.1$
		$l.0 = 1$
	$L \rightarrow LB$	$v.0 = w * v.1 + v.2$
		$l.0 = l.1 + 1$
Numbers	$N \rightarrow L$	$v.0 = v.1$
	$N \rightarrow L.L$	$v.0 = v.1 + \frac{v.2}{2^{l.2}}$

$\mathcal{G}_B$  erzeugt Binärzahlen mit und ohne Punkt,

- Synthetische Attribute:  
 $B, N : v(\text{value})$   
 $L : v.l(\text{length})$
- Semantische Regeln: Attributgleichungen mit Attributvariablen
- Index  $i$ :  $i$ -tes (Nichtterminal-) Symbol
- Ziel: Bestimmung des Zahlenwerts
- Attribute von
  - $v$  : rationale Zahlen ( $A^v = \mathbb{Q}$ )
  - $l$  : natürliche Zahlen ( $A^l = \mathbb{N}$ )

Attributierung von  $\mathcal{G}_B$  mit synthetischen und inheriten Attributen, zusätzliches inherites Attribut für Bits und Listen:  $p$  (position)

$B \rightarrow 0$	$v.0$			
$B \rightarrow 1$	$v.0 = 2^{p.0}$			
$L \rightarrow B$	$v.0 = v.1$	$l.0 = 1$	$p.1 = p.0$	
$L \rightarrow LB$	$v.0 = v.1 + v.2$	$l.0 = l.1 + 1$	$p.1 = p.0 + 1$	$p.2 = p.0$
$N \rightarrow L$	$v.0 = v.1$		$p.1 = 0$	
$N \rightarrow L.L$	$v.0 = v.1 + v.2$	$p.1 = 0$	$p.2 = -l.2$	



Berechnung des Wurzelattributs:

1. Längen  $\uparrow$  BU
2. Positionen  $\downarrow$  TD
3. Werte  $\uparrow$  BU

**Definition (Attributgrammatik):** Sei  $\mathcal{G} = \langle N, \Sigma, P, S \rangle \in \text{CFG}$ . Sei  $\text{Att}$  eine Menge von Attributen,  $A = (A^\alpha | \alpha \in \text{Att})$  eine Familie von Attribut-Wertmengen und  $\text{att} : \mathfrak{X} \rightarrow P(\text{Att})$  eine Attributzuordnung.

Sei  $\text{Att} = \text{Syn} \dot{\cup} \text{Inh}$  eine Zerlegung in Teilmengen synthetischer und inheriter Attribute, sodaß  $\text{att}$  zerfällt in

$$\text{syn} : \mathfrak{X} \rightarrow P(\text{Syn}) \text{ mit } \text{syn}(X) := \text{att}(X) \cap \text{Syn}$$

und

$$\text{inh} : \mathfrak{X} \rightarrow P(\text{Inh}) \text{ mit } \text{inh}(X) := \text{att}(X) \cap \text{Inh}.$$

Eine Regel  $\pi = X_0 \rightarrow X_1 \dots X_r \in P$  bestimmt die Menge  $\text{Var}_\pi := \{\alpha.i | \alpha \in \text{att}(X_i), 0 \leq i \leq r\}$  der formalen Attributvariablen von  $\pi$  mit den Teilmengen

$$\text{IVar}_\pi := \{\alpha.i | (i = 0 \wedge \alpha \in \text{syn}(X_0)) \vee (1 \leq i \leq r \wedge \alpha \in \text{inh}(X_i))\}$$

und

$$\text{OVar}_\pi := \text{Var}_\pi \setminus \text{IVar}_\pi$$

der Innen- bzw. Außenvariablen.

Eine Attributgleichung von  $\pi$  (semantische Regel) hat die Form

$$\alpha.i = f(\alpha_1 \cdot i_1, \dots, \alpha_n \cdot i_n)$$

mit  $\alpha.i \in \text{IVar}_\pi$ ,  $\alpha_1 \cdot i_1, \dots, \alpha_n \cdot i_n \in \text{OVar}_\pi$ ,  $f : A^{\alpha_1} \times \dots \times A^{\alpha_n} \rightarrow A^\alpha$  und  $n \in \mathbb{N}$ .

Sei  $E_\pi$  eine Menge von Attributgleichungen, in der jede Innenvariable genau einmal vorkommt und  $E := (E_\pi | \pi \in P)$ . Dann heißt  $\mathfrak{A} = (\mathcal{G}, E)$  eine Attributgrammatik  $\mathfrak{A} \in \text{AG}$ .

**Definition (Das Attributgleichungssystem eines Ableitungsbaumes)** Sei  $\mathfrak{A} = (\mathcal{G}, E) \in \text{AG}$ .  $\mathfrak{A}$  induziert für jeden Ableitungsbaum  $t$  von  $\mathcal{G}$  ein Attributgleichungssystem  $E_t$ . Sei  $\text{Kn}(t)$  die Menge der Knoten von  $t$ . Sie bestimmt die Menge

$$\text{Var}_t := \{\alpha.k | k \in \text{Kn}(t) \text{ markiert durch } X \in \mathfrak{X}, \alpha \in \text{att}(X)\}$$

der aktuellen Attributvariablen von  $t$ .

Wird auf einen inneren Knoten (kein Blattknoten)  $k_0 \in \text{Kn}(t)$  die Regel  $\pi = X_0 \rightarrow X_1 \dots X_r$  angewandt und sind  $k_1, \dots, k_r \in \text{Kn}(t)$  die entsprechenden Nachfolgerknoten, so erhält man das Attributgleichungssystem  $E_{k_0}$  von  $k_0$  aus  $E_\pi$  durch die Indexsubstitution ( $i \rightarrow k_i | 0 \leq i \leq r$ ) bei den Attributvariablen. Dann ist  $E_t := \bigcup \{E_k | k \text{ innerer Knoten von } t\}$ .

- Beachte: Zu jeder aktuellen Attributvariablen  $\alpha.k$ , ausgenommen die inheriten der Wurzeln und die synthetischen der Blätter, gibt es genau eine Gleichung  $\alpha.k = \dots$
- Annahme:
  - keine inheriten Attribute des Startsymbols
  - synthetische Attribute der Terminalsymbole vom Scanner

**Lösbarkeit von  $E_t$ :**  $E_t$  kann keine, genau eine oder mehrere Lösungen besitzen.

**Beispiel:**  $\mathcal{G}$  enthalte die Regeln  $A \rightarrow aBc$  und  $B \rightarrow w$ . Ferner  $\alpha \in \text{syn}(B)$  und  $\beta \in \text{inh}(B)$ .

- Attributgleichungen: ...

$$\beta.2 = f(\alpha, 2)$$

$$\alpha.0 = g(\beta.0)$$

- $E_t$ : zirkuläre Abhängigkeit

$$\beta.k = f(\alpha.k)$$

$$\alpha.k = g(\beta.k)$$

- Für  $A^\alpha = A^\beta = \mathbb{N}, g = id_{\mathbb{N}}$  und

1.  $f(x) = x + 1$  ergibt keine Lösung
2.  $f(x) = 2x$  ergibt genau eine Lösung
3.  $f(x) = x$  ergibt unendlich viele Lösungen

**Folgerung:** Zirkularitäten sind unerwünscht. Sie entstehen erst in  $E_t = \bigcup \{E_k | k \text{ innerer Knoten von } t\}$ .

**Definition:**  $\mathfrak{A} \in \text{AG}$  heißt zirkulär:  $\leftrightarrow$  Es gibt einen Ableitungsbaum  $t$  mit zirkulärem Attributgleichungssystem  $E_t$ , d.h. eine aktuelle Attributvariable hängt von sich selbst ab.

**Ein Zirkularitätstest für Attributgrammatiken** Sei  $\mathfrak{A} = \langle \mathcal{G}, E \rangle \in \text{AG}$ . Eine Regel  $\pi = X_0 \rightarrow X_1 \dots X_r$  bestimmt ihren Abhängigkeitsgraphen  $\text{DG}_\pi$  mit der Knotenmenge  $\text{Kno}(\text{DG}_\pi) := \text{OKno}(\text{DG}_\pi) \cup \text{UKno}(\text{DG}_\pi)$

$$\text{OKno}(\text{DG}_\pi) := \{\alpha.0 | \alpha \in \text{att}(X_0)\} \quad \text{Oberknoten}$$

$$\text{UKno}(\text{DG}_\pi) := \{\alpha.i | \alpha \in \text{att}(X_i), 1 \leq i \leq r\} \quad \text{Unterknoten}$$

und der Kantenmenge  $\text{Kan}(\text{DG}_\pi)$ :

$$(x_1, x_2) \in \text{Kan}(\text{DG}_\pi) : \leftrightarrow x_2 = f(\dots x_1 \dots) \in E_\pi$$

**Beispiel:****Beachte:**  $\text{Var}_\pi = \text{UKno}(\text{DG}_\pi) \cup \text{OKno}(\text{DG}_\pi) = \text{OVar}_\pi \cup \text{IVar}_\pi$ **Problem:** Beim Verkleben der  $\text{DG}_\pi$  zu  $\text{DG}_t$  für einen Ableitungsbaum  $A$  können Schleifen auftreten.**Folgerung:**  $\mathfrak{A}$  ist zirkulär  $\leftrightarrow$  Es gibt einen Ableitungsbaum  $t$  von  $\mathfrak{A}$ , so daß  $\text{DG}_t$  eine Schleife enthält.**Bemerkung:**  $\text{DG}_t$  entsteht aus  $(\text{DG}_\pi | \pi \in P)$  wie  $E_t$  aus  $(E_\pi | \pi \in P)$  und  $t$  aus  $(\pi | \pi \in P)$  entstehen.**Beobachtung:** Eine Schleife in  $\text{DG}_t$  bestimmt einen Teilgraphen der Form, also eine Regel  $\pi$  mit verbindbaren Unterknoten von  $\text{DG}_\pi$ . Eine Verbindung gehört zu einem Nichtterminalsymbol und verläuft von einer inheriten zu einer synthetischen Attributvariablen.**Definition (Attributabhängigkeit):** Sei  $A \in N, \alpha \in \text{inh}(A), \alpha' \in \text{syn}(A)$ . Dann heißt  $\alpha'$  von  $\alpha$  unterhalb  $A$  abhängig, Bezeichnung:  $\alpha \underset{\sim}{A} \alpha'$ , wenn ein Ableitungsbaum  $t$  mit Wurzel  $A$  und Wurzelknoten  $k$  existiert, sodaß in  $\text{DG}_t$  ein Pfad von  $\alpha.k$  nach  $\alpha'.k$  führt:

$$(\alpha.k, \alpha'.k) \in \text{trans}(\text{Kan}(\text{DG}_t))$$

**Definition (Attributabhängigkeitsmengen):** Sei  $t$  ein Ableitungsbaum mit Wurzel  $A$ .

$$D(A, t) := \{(\alpha, \alpha') | \alpha \underset{\sim}{A} \alpha' \text{ in } t\}$$

$$\mathcal{D}(A) := \{D(A, t) | t \text{ Ableitungsbaum mit Wurzel } A\}$$

Bei der induktiven Bestimmung benutzen wir folgende Bezeichnung:

**Definition:** Für  $\pi = A_0 \rightarrow w_0 A_1 w_1 \dots A_r w_r$  und  $D_i \subseteq \text{inh}(A_i) \times \text{syn}(A_i)$  definieren wir

$$D[\pi_i; D_1, \dots, D_r] \subseteq \text{inh}(A_0) \times \text{syn}(A_0)$$

durch

$$\{(\alpha, \alpha') | (\alpha.0, \alpha'.0) \in \text{trans}(\text{Kan}(\text{DG}_{p_i}) \cup \bigcup_{i=1}^r \{(\beta.i, \beta'.i) | (\beta, \beta') \in D_i\})\}.$$

**Lemma:** Die Attributabhängigkeitssysteme  $\mathcal{D}(A)$  mit  $A \in N$  sind induktiv bestimmt durch:

1.  $\pi = A \rightarrow w \curvearrowright D[\pi_i] \in \mathcal{D}(A)$
2.  $\pi = A \rightarrow w_0 A_1 w_1 \dots A_r w_r, \quad D_i \in \mathcal{D}(A_i), \quad 1 \leq i \leq r \leftrightarrow D(\pi; D_1, \dots, D_r] \in \mathcal{D}(A)$

Beweis: Induktion über die Struktur der Ableitungsbäume.

**Beachte:** Da für  $D \in \mathcal{D}(A)$  gilt:  $D \subseteq \text{inh}(A) \times \text{syn}(A)$ , bricht der Berechnungsprozeß nach endlichen vielen Schritten ab.

**Folgerung (Zirkulartätstest):**  $\mathfrak{A} \in \text{AG}$  ist zirkulär  $\leftrightarrow$  es gibt  $\pi = A_0 \rightarrow w_0 A_1 w_1 \dots A_r w_r, \alpha.k \in \text{UKno}(\text{DG}_\pi)$  und  $D_i \in \mathcal{D}(A_i) (1 \leq i \leq r)$ , so daß  $(\alpha.k, \alpha.k) \in \text{trans}(\text{Kan}(\text{DG}_\pi) \cup \bigcup_{i=1}^r \{\beta.i, \beta'.i \mid (\beta, \beta') \in D_i\})$

**Komplexität**  $n = |\Omega|, T(n)$  Zeit zur Entscheidung der Zirkularität

$$2^{\frac{cn}{\log n}} \leq T(n) \leq 2^{dn^2}$$

**Stark nicht-zirkuläre Grammatiken** Keine Trennung der Abhängigkeitsmengen verschiedener Ableitungsbäume

$$D(A) := \{(\alpha, \alpha') \mid \alpha \underset{\sim}{A} \alpha'\}$$

Hinreichendes Kriterium für die Nichtzirkularität. Test in Polynomzeit.

**Beachte:** Es gibt  $\mathfrak{A} \in \text{AG}$ , sodaß  $\Omega$  nicht zirkulär ist, ohne stark nicht-zirkulär zu sein.

#### Attributberechnung:

1.  $E_t$  als Termersetzungssystem, TD-Berechnung  
keine Zirkularität  $\curvearrowright$  Termination mit eindeutiger Lösung
2. BU-Auflösung von  $E_t$ : Variablen durch Werte ersetzen und Terme ausrechnen
3. Uniforme Berechnung, unabhängig von Ableitungsbaum  
Rekursive Prozeduren / Funktionen für  $(E_\pi \mid \pi \in P)$   
Idee: Jedem synthetischen Attribut einer Variablen wird eine Prozedur / Funktion zugeordnet mit inheriten Attributen als Parameter
4. Spezialfälle: SAG, LAG mit Attributberechnung während der Syntaxanalyse

#### S-Attributgrammatiken

**Definition:**  $\mathfrak{A} \in \text{AG}$  heißt S-Attributgrammatik ( $\mathfrak{A} \in \text{SAG}$ ), wenn  $\text{Inh} = \emptyset$ , also  $\text{Att} = \text{Syn}$ .

- Folgerung: BU-Berechnung der Attributwerte
- Spezialfall: Bestimmung der Wurzelattribute
- Idee: Durchführung der Attributberechnung während der BU-Syntaxanalyse, Zwischenspeichern von Attributwerten im Analysekeiler (YACC).  
Beim Reduktionsschritt werden simultan die Werte der synthetischen Attribute von  $A$  aus Werten unter  $v_1, v_2, v_3$  berechnet (records); Acc-Schritt: Ausgabe der Wurzelattribute

**Beispiel 1: Stackcode für arithmetische Ausdrücke**

$$E \rightarrow (E + E)|(E * E)|id|num \quad (1, 2, 3, 4)$$

1.  $c.0 = c.2; c.4; ADD$
2.  $c.0 = c.2; c.3; MULT$
3.  $c.0 = LOADa.1$ , a: lexical address
4.  $c.0 = LITn.1$ , n: lexical value

Lexikalische Attribute, die mit dem Token vom Lexer übergeben werden:

$$(id, X) \quad (num, 42)$$

$$((3 + X) * (Y + 5)) \rightsquigarrow (((num, 3) + (id, X)) * ((id, X) + (num, 5)))$$

Codegenerierung durch Attributauswertung während der Syntaxanalyse:

1. shift: Aufrufe des Lexers durch "nextsym"
  - a) Token  $\mapsto LR(0)$ -Menge (mit goto)
  - b) lexikalisches Attribut

Beide Informationen werden auf den Analysekeiler abgelegt.

2. reduce:
  - a) Reduktionen auf Analysekeiler
  - b) simultane Attributberechnung
3. accept: Wurzelattribut

$$LIT3; LOADX; ADD; LOAY; LIT5; ADD; MULT$$

## Beispiel 2: Berechnung des abstrakten Syntaxbaumes

$$\pi = A_0 \rightarrow w_0 A_1 w_1 \dots A_r w_r$$

repräsentiert Op-Symbol  $F_\pi$  vom Typ  $A_1 \times \dots \times A_0$  ( $A_i$ : Sorte, Typ).

- Folge: Ableitungsbaum vereinfacht sich zu abstraktem Syntaxbaum (AST); nur dieser ist für die Übersetzung notwendig
- Konkrete Syntax: Benutzerfreundliche "Mixfix"-Notation unter Verwendung natürlicher Sprache für Terminalsymbol von CFG
- Abstrakte Syntax: Darstellungsabhängige algebraische Struktur ( $\rightarrow$  ADT, Konstruktoren) (CB 3-6):  
`begin id:=num+id; if id<num then id:=num+id; end`
- Regeln selbst als Op-Symbole möglich: Regelbaum
- Aufgabe: S-Attributgrammatik
- Darstellung des Aufbaus von Syntaxbäumen durch Konstruktion von Graphen

1. Definiere

$a \in \text{Adr}$	unendliche Menge von Adressen für Speicherplätze (Heap)
$\Omega$	Op-Alphabet mit Stelligkeit, $f^{(n)} \in \Omega = \{assign^{(1)}, seq^{(2)}, cond^{(2)}, less^{(2)}, plus^{(2)}, id^{(0)}, num^{(0)}\}$
Kno	Menge der $\Omega$ -Knoten
$k \in \text{Kno}$	$\{(a, f, a_1, \dots, a_n) \mid a, a_i \in \text{Adr}, A \in \Omega^{(n)}\}$
Graph:	Menge von Knoten mit einer Wurzel
$G \in \text{Graph}$	$\{(a, k) \mid a \in \text{Adr}, k \subseteq \text{Kno}\}$

(CB 3-7)

2. Konstruktorfunktionen für Graphen:

$$f^{(n)} \in \Omega \mapsto mk - f : \text{Graph}^n \rightarrow \text{Graph}$$

$$mk - f((a_1, k_1), \dots, (a_n, k_n)) := (a, k)$$

mit neuer Adresse  $a$  und

$$k := \{(a, f, a_1, \dots, a_n)\} \cup \bigcup_{i=1}^n k_i$$

( $k_i$  disjunkt, erfolgt automatisch)

- S-Attributierung von  $\mathcal{G}_s$ :  
Wurzeladressen als S-Attribute im Analysekernel; Heap statt Ausgabeband  
`/* printf("%d", i) -> mk-f_i */`  
Attributwerte als Graphen; Analyseergebnis: Graphdarstellung des AST auf dem Heap  $\rightsquigarrow$  Basis für weitere Attributierung
- Beachte: Unterschied zwischen Wurzelattributberechnung (Beispiel 1) und voller Attributierung (Erweiterung von Beispiel 2)

**Beispiel 3: Typenberechnung für arithmetische Ausdrücke**  $Typ : \{int, real\}$  mit  $\max Typ^2 \mapsto Typ$  bzgl.  $int < real$ .  $nextsym()$  (Lexer) bestimmt Typattribut von  $num$ , Typattribut von  $id$  über lookup-table  $\rightsquigarrow$  Syntaxbaum mit Typinformationen (weiteres Knotenfeld) annotiert.

### L-Attributgrammatiken

**Definition:**  $\mathfrak{A} \in \langle \mathcal{G}, E \rangle \in AG$  ist eine L-Attributgrammatik ( $\mathfrak{A} \in LAG$ )  $\leftrightarrow$  Für jede Attributgleichung  $\alpha.i = f(\dots \beta.j \dots)$  gilt:  $\alpha \in \text{Inh}, \beta \in \text{Syn} \curvearrowright j < i$

**Beispiel (DG $_{\pi}$ ):** Bei L-Attributgrammatiken fallen die gestrichenen Kanten weg.

**Folgerung:**  $\mathcal{G} \in LAG$  nicht zirkulär. Die Attributberechnung erfolgt in "depth-first, left-to-right order". Die Baumreise geschieht dabei mit 2 Kantenbesuchen:

1. top-down: inherite Attribute
2. bottom-up: synthetische Attribute

Auswertungsreihenfolge:

**Syntaxanalyse mit L-Attributauswertung** Sei  $\mathfrak{A} = \langle \mathcal{G}, E \rangle \in LAG$  mit  $\mathcal{G} \in LL(1)$ .

- Ziel: Erweiterung der TD-Analyse zur Berechnung der synthetischen Wurzelattribute
- Methode: Expansion von  $A \in N$  auf dem Analysekeiler so gestalten, daß spätere Reduktion möglich: TD + BU, TD: inherite Attribute, BU: synthetische Attribute
- Kellularphabet:  $LR(0)_{\pi}(\mathcal{G}) := \{[A \rightarrow \alpha.\beta] | \pi = A \rightarrow \alpha.\beta\}$

$$\text{Val}_{\pi} = \{\text{val}_{\pi} | \text{val}_{\pi} : \text{Var}_{\pi} \rightarrow A\}$$

$$\Gamma := \bigcup_{\pi \in P} (LR(0)_{\pi}(\mathcal{G}) \times \text{Val}_{\pi} \cup \{\rightarrow \cdot S, \emptyset\}, (\rightarrow S \cdot, \text{val}))$$

$$\text{val} : \text{syn}(S) \rightarrow A$$

- Konfigurationen:  $(w, \gamma)$  ( $w$ : Eingabeband,  $\gamma \in \Gamma$ )
- Arbeitsweise (3 Fälle):

1. "expand"-Schritt mit Berechnung der inheriten Attribute

$$\begin{array}{|c|c|} \hline \vdots & \vdots \\ \hline [A \rightarrow X_1 \dots X_{i-1} \cdot B\gamma] & \text{val} \\ \hline \end{array} \Rightarrow \begin{array}{|c|c|} \hline \vdots & \vdots \\ \hline [A \rightarrow X_1 \dots X_{i-1} \cdot B\gamma] & \text{val} \\ \hline [B \rightarrow \cdot\beta] & \text{val}' \\ \hline \end{array}$$

(Kellerspitze unten) Dabei ist  $B \rightarrow \beta$  wegen  $\mathcal{G} \in LL(1)$  durch Eingabe-  
Lookahead bestimmt.  $\text{val}'$  belegt die inheriten Attributvariablen von  $B$  nach  
Attributgleichungen für  $A \rightarrow X_1 \dots X_{i-1} B\gamma$ : Sei  $\alpha \in \text{inh}(B)$  und  $\alpha.i = t \in$   
 $E_{A \rightarrow X_1 \dots \gamma}$ , so  $\text{val}'(\alpha.0) := \text{val}(t) \in A$ . (Die Funktion  $\text{val}$  übernimmt das  
Einsetzen und ausrechnen der Attributwerte.)

2. "match"-Schritt (ohne Attributberechnung)

$$\begin{array}{|c|c|} \hline \vdots & \vdots \\ \hline [A \rightarrow \gamma \dots a\gamma'] & \text{val} \\ \hline \end{array} \Rightarrow \begin{array}{|c|c|} \hline \vdots & \vdots \\ \hline [A \rightarrow \gamma a \cdot \gamma'] & \text{val} \\ \hline \end{array}$$

Falls  $a$  nächstes Eingabesymbol ist; Eingabekopf rückt vor.

3. "reduce"-Schritt mit Berechnung der synthetischen Attribute

$$\begin{array}{|c|c|} \hline \vdots & \vdots \\ \hline [A \rightarrow X_1 \dots X_{i-1} \cdot B\gamma] & \text{val} \\ \hline [B \rightarrow \beta \cdot] & \text{val}' \\ \hline \end{array} \Rightarrow \begin{array}{|c|c|} \hline \vdots & \vdots \\ \hline [A \rightarrow X_1 \dots X_{i-1} B \cdot \gamma] & \text{val}'' \\ \hline \end{array}$$

$\text{val}''$  erweitert  $\text{val}$  um synthetische Attribute von  $B$ :

$$\text{val}''(\alpha.i) = \begin{cases} \text{val}(t), & \text{falls } \alpha.0 = t \in E_{B \rightarrow \beta} \\ \text{val}(\alpha.i) & \text{sonst} \end{cases}$$

**Anwendung von LAG** Überprüfung der Deklariertheit von Bezeichnern:  $\mathcal{G}$  sei gegeben  
durch

$P \rightarrow DL; SL$	Programm
$DL \rightarrow V N; DL$	Deklarationsliste
$SL \rightarrow S S; SL$	Statementliste      Beispiel: $a; c; a; b; a := c; c := d$
$V \rightarrow a b  \dots$	Variablen
$S \rightarrow V := V$	Statements

**Attribute** Die Attribute werden wie folgt deklariert:

$\text{syn } v$	deklarierte oder benutze Variable für $V \in \{a, b, \dots\}$
$\text{syn } dv$	Menge von deklarierten Variablen für $DL \subseteq \{a, b, \dots\}$
$\text{inh } v$	Umgebung für $S, SL \subseteq \{a, b, \dots\}$
$\text{syn } decl$	deklariert für $S, SL, P \in \{true, false\}$



### Attributgleichungen

$$\begin{array}{ll} P \rightarrow DL; SL & \text{decl.0} = \text{decl.3} \\ & \text{env.3} = \text{dv.1}(\text{LAG!}) \\ DL \rightarrow V & \text{dv.0} = \{v.1\} \\ DL \rightarrow V; DL & \text{dv.0} = \{v.1\} \cup \text{dv.3} \\ V \rightarrow a | \dots & v.0 = a \\ SL \rightarrow S & \text{env.1} = \text{env.0} \\ & \text{decl.0} = \text{decl.1} \\ SL \rightarrow S; SL & \text{env.1} = \text{env.0} \\ & \text{env.3} = \text{env.0} \\ & \text{decl.0} = \text{decl.1} \text{ AND } \text{decl.3} \\ S \rightarrow V := V & \text{decl.0} = v.1 \in \text{env.1} \text{ AND } v.0 \in \text{env.0} \end{array}$$

### Attributierter Ableitungsbaum

# 5 Übersetzung in Zwischencode

Aufteilung der Codeerzeugung:  $PS \xrightarrow{trans} Z \xrightarrow{code} MC$

- front-end: *trans* erzeugt maschinenunabhängigen Zwischencode für eine abstrakte Stackmaschine
- back-end: *code* erzeugt Maschinencode

Vorteil dieser Zerlegung:  $Z$  ist maschinenunabhängig  $\rightarrow$  Portabilität, Transparenz, Codeoptimierung

Beispiel:

- Java Virtual Machine (JVM)
- P-Code von Pascal

## 5.1 Übersetzung von Ausdrücken, Anweisungen, Blöcken und Prozeduren

Beispiel-Programmiersprache: BPS ("Pascal ohne Datenstrukturen und ohne Prozedurparameter")

- arithmetische und boolesche Ausdrücke mit strikter bzw. nicht-strikter Semantik
- Kontrollstrukturen: Sequenz, Verzweigung, Iteration
- Blöcke und Prozeduren: Lokale und globale Variablen, dynamische Speicherverwaltung mit Laufzeitkeller;  
unterschiedliche Verwendung in Programmiersprachen:
  - FORTRAN: Unterprogramme nicht geschachtelt, keine Rekursion  $\rightsquigarrow$  statische Speicherverwaltung, Speicherbedarf zur Übersetzungszeit bekannt
  - C: rekursive Prozeduren ohne Schachtelung  $\rightsquigarrow$  dynamische Speicherverwaltung, keine statischen Verweise
  - ALGOL-Familie: ALGOL 60, Pascal, Modula: Geschachtelte rekursive Prozedurdeklaration  $\rightsquigarrow$  dynamische Speicherverwaltung mit statischen Verweisen
- keine Datenstrukturen, keine Prozedurparameter
- Beschränkung auf Datentyp int
- Syntax: CB 4-6

## Semantik

- Bezeichner einer Deklaration  $\Delta$  müssen paarweise verschieden sein.
- Ein im Anweisungsteil  $\Gamma$  eines Blocks  $\Delta\Gamma$  auftretender Bezeichner muß deklariert sein, und zwar in  $\Delta$  oder in der Deklarationsliste eines  $\Delta\Gamma$  umfassenden Blocks.
- Mehrfach-Deklaration eines Bezeichners auf verschiedenen Niveaus möglich: Die "innerste" Deklaration ist für ein Auftreten gültig.
- Static scope: Beim Aufruf einer Prozedur ist ihre Deklarationsumgebung gültig und nicht ihre Aufrufumgebung.

## Zwischencode für BPS AM abstrakte Maschine mit Datenkeller und Prozedurkeller

- Zustandsraum  $ZR := BZ \times DK \times PK$  mit
  - Befehlszähler  $BZ := \mathbb{N}$
  - Datenkeller  $DK := \mathbb{Z}^*$  (Spitze rechts)
  - Prozedurkeller  $PK := \mathbb{Z}^*$  (Spitze links)
- Zustand  $s = (m, d, p) \in ZR$  mit
  - Befehlsmarke  $m$
  - DK-Zustand  $d = d.r : \dots : d.1$
  - PK-Zustand  $p = p.1 : \dots : p.t$

## AM-Befehle

- arithmetische: ADD, ...
- logische: LT, NOT, AND, OR, ...
- Sprungbefehle: JMP  $n$ , JFALSE  $n$  ( $n \in \mathbb{N}$ )
- Prozedurbefehle: CALL( $ca, dif, loc$ ) ( $ca, dif, loc \in \mathbb{N}$ ), RET
- Transportbefehle: LOAD( $dif, off$ ), STORE( $dif, off$ ) ( $dif, off \in \mathbb{N}$ )

## Befehlsemantik

- $\llbracket B \rrbracket : ZR \rightarrow ZR$  für jedem AM-Befehl
- $\llbracket ADD \rrbracket(m, d : z_1 : z_2, p) := (m + 1, d : (z_1 + z_2), p)$
- $\llbracket LT \rrbracket(m, d : z_1 : z_2, p) := (m + 1, d : b, p)$  mit  $b = \begin{cases} 1 & \text{falls } z_1 < z_2 \\ 0 & \text{falls } z_1 \geq z_2 \end{cases}$
- $\llbracket AND \rrbracket(m, d : b_1 : b_2, p) := (m + 1, d : b_1 \wedge b_2, p)$ , entsprechend für OR und NOT

- $\llbracket JFALSEn \rrbracket(m, d : b, p) := \begin{cases} (n, d, p) & \text{falls } b = 0 \\ (m + 1, d, p) & \text{falls } b = 1 \end{cases}$
- $\llbracket JMPn \rrbracket(m, d, p) := (n, d, p)$

Für die Semantik der Prozedur- und Transportbefehle: Besondere Struktur des Prozedurkellers  $p$ :  $p$  zerfällt in Aktivierungsblöcke (Frames) der Form  $sv : dv : ra : l_1 : \dots : l_k$

- $sv$ : statischer Verweis, zeigt auf Aktivierungsblock der Deklarationsumgebung
- $dv$ : dynamischer Verweis, zeigt auf den letzten Aktivierungsblock
- $ra$ : Rücksprungadresse, ist die Codeadresse nach Beendigung des Prozeduraufrufs
- $l_i$ : lokale Variablen des Prozeduraufrufs

**Berechnung des statischen Verweises**  $sv$  liefert die Differenz zwischen Aufruf- und Deklarationsniveau und damit die Länge der Verweiskette.

### Abstrakte Maschine

$$s = (m, d, p)$$

$p$  zerfällt in Aktivierungsblöcke  $sv : dv : ra : l_1 : \dots : l_k$

**Berechnung des statischen Verweises**  $sv$  liefert die Differenz zwischen Aufruf- und Deklarationsniveau und damit die Länge der Verweiskette.

**Hilfsfunktion:**  $base : PK \times \mathbb{N} \rightarrow \mathbb{N}$

bestimmt für einen Prozedurkeller bzgl. einer Niveaudifferenz  $dif$  den Beginn der Deklarationsumgebung (als absolute Adresse des aktuellen Prozedurkellers).

$$base(p, 0) := 1$$

$$base(p, dif + 1) := base(p, dif) + p.base(p, dif)$$

**Beispiel:** 2. Aufruf von A (s. Folie CB 4-7)

$$\begin{aligned} dif &= 2 \\ base(p, 0) &= 1 \\ base(p, 1) &= 1 + p.1 = 6 \\ base(p, 2) &= 6 + p.6 = 11 \end{aligned}$$

Es folgt:  $sv = base(p, 2) + 2 + 2$

- $CALL(ca, dif, loc)$  mit
  - $ca \in \mathbb{N}$  Codeadresse

- $diff \in \mathbb{N}$  Niveaudifferenz
- $loc \in \mathbb{N}$  Zahl der lokalen Variablen

erzeugt neuen Aktivierungsblock und springt zum Code des Prozedurrumpfes

$$\begin{aligned} \llbracket CALL(ca, diff, loc) \rrbracket(m, d, p) := \\ (ca, d, (base(p, diff) + 2 + loc) : (loc + 2) : (m + 1) : 0 : \dots : 0 : p) \end{aligned}$$

- *RET* löscht den letzten Aktivierungsblock und kehrt zur Aufruftabelle zurück

$$\begin{aligned} \llbracket RET \rrbracket(m, d, p.1 : p.2 : \dots : p.t) := \\ \text{if } t \geq 2 + p.2 \text{ then } (p.3, d, p.(2 + p.2) : \dots : p.t) \end{aligned}$$

Der Prozedurkeller wird stets mindestens einen Aktivierungsblock enthalten.

- *LOAD(diff, off)* und *STORE(diff, off)* laden und speichern Variablen zwischen Datenkeller und Prozedurkeller. Relative Adressierung mit
  - Niveaudifferenz  $diff$  zwischen Auftreten und Deklaration
  - Offset  $off$  als relative Adresse im Aktivierungsblock

Die Kette der statischen Verweise bestimmt die sichtbare Umgebung auf dem Prozedurkeller.

$$\begin{aligned} \llbracket LOAD(diff, off) \rrbracket(m, d, p) &:= (m + 1, d : p.[base(p, diff) + 2 + off], p) \\ \llbracket STORE(diff, off) \rrbracket(m, d : z, p) &:= (m + 1, d, p.[base(p, diff) + 2 + off/z]) \\ \llbracket LITz \rrbracket(m, d, p) &:= (m + 1, d : z, p) \end{aligned}$$

**AM-Code** Befehlsfolgen mit aufsteigenden Befehlsmarken

$P \in \text{AM-Code}$ :  $\leftrightarrow P = a_1 : B_1; \dots; a_p : B_p$  mit  $a_i \in \text{Adr} := \mathbb{N}$ ,  $a_i = a_1 + i - 1$  und  $B_i$  ein AM-Befehl ( $1 \leq i \leq p$ ).

Semantik von  $P$ : Iteration der Befehle gemäß Befehlszähler

$$\mathcal{J} : \text{AM-Code} \times \text{ZR} \rightarrow \text{ZR}$$

$$\mathcal{J}\llbracket P \rrbracket(m, d, p) := \text{if } a_1 \leq m \leq a_p \text{ and } m = a \text{ then } \mathcal{J}\llbracket P \rrbracket(\llbracket B_i \rrbracket(m, d, p)) \text{ else } (m, d, p)$$

## Übersetzung von BPS-Programmen in AM-Code

- trans: BPS-Prog  $\rightarrow$  AM-Code (als Zwischencode  $Z$ )
- Hilfsmittel: Symboltabelle

$$\begin{aligned} \text{Tab} := \{st \mid st : IDE \rightarrow \\ (\{const\} \times \mathbb{Z}) \vee (\{var\} \times Lev \times Off) \vee (\{proc\} \times Adr \times Lev \times Size)\} \end{aligned}$$

- Variablen-Deklaration:

- \* Deklarationsniveau:  $dl \in Lev$
- \* Offset:  $off \in Off := \mathbb{N}$
- Prozedurdeklaration:
  - \* Startadresse:  $ca \in Adr$  des Prozedurcodes
  - \* Deklarationsniveau:  $dl \in Lev$
  - \* Anzahl der lokalen Variablen:  $loc \in Size := \mathbb{N}$

### Aufbau der Symboltabelle

$$up : Decl \times Tab \times Adr \times Lev \dashrightarrow Tab$$

$up(\Delta, st, a, l)$  beschreibt den Update einer Symboltabelle  $st$  bzgl. einer Deklaration  $\Delta$  bei freier Adresse  $a$  und aktuellem Niveau  $l$  (Blockschachtelungstiefe).

$$\begin{aligned}
up(\Delta_c \Delta_v \Delta_p, st, a, l) &:= \text{if } diff - id(\Delta_c \Delta_v \Delta_p) \\
&\quad \text{then } up(\Delta_p, up(\Delta_v, up(\Delta_c, st, a, l), a, l), a, l) \\
up(\varepsilon, st, a, l) &:= st \\
up(const I_1 = Z_1, \dots, I_n = Z_n; st, a, l) &:= st[I_1/(const, Z_1), \dots, I_n/(const, Z_n)] \\
up(var I_1, \dots, I_n; st, a, l) &:= st[I_1/(var, l, 1), \dots, I_n/(var, l, n)] \\
up(proc I_1; B_1; \dots; proc I_n; B_n; st, a, l) &:= st[I_1/(proc, a_1, l, size(B_1)), \dots, I_n/(proc, a_n, l, size(B_n))] \\
size : Block \rightarrow \mathbb{N} & \\
size(\Delta_c var I_1, \dots, I_n; \Delta_p \Gamma) &:= n
\end{aligned}$$

**Anfangstabelle**  $P = in/out I_1, \dots, I_n; B.$  hat die Semantik  $\mathcal{M}[[P]] : \mathbb{Z}^n \rightarrow \mathbb{Z}^n$ . Für  $(z_1, \dots, z_n) \in \mathbb{Z}^n$  wählen wir den Anfangszustand  $(1, \varepsilon, 0 : 0 : 0 : z_1 : \dots : z_n)$ . Die entsprechende Anfangstabelle hat daher  $n$  Einträge:

$$st_{I/O}(I_j) = (var, 0, j) \text{ für } 1 \leq j \leq n$$

$$\begin{aligned}
bt : Block \times Tab \times Adr \times Lex &\dashrightarrow AM - Code \\
dt : Decl \times Tab \times Adr \times Lex &\dashrightarrow AM - Code \\
ct : Cmd \times Tab \times Adr \times Lex &\dashrightarrow AM - Code \\
et : AExp \times Tab \times Adr \times Lex &\dashrightarrow AM - Code \\
sbt : BExp \times Tab \times Adr \times Lex &\dashrightarrow AM - Code
\end{aligned}$$

$$trans : Prog \rightarrow AM - Code$$

$$\begin{aligned}
trans(in/out I_1, \dots, I_n; B.) &:= 1 : CALL(a_\Gamma, 0, size(B)); \\
&2 : JMP0; \\
&bt(B, st_{I/O}, a_\Gamma, 1)
\end{aligned}$$

$a_\Gamma$ : Startadresse des Anweisungscode von  $\Gamma$  in  $B = \Delta\Gamma$ .

$$\begin{aligned}
bt(\Delta\Gamma, st, a, l) &:= dt(\Delta, up(\Delta, st, a1, l), a1, l) \\
&ct(\Gamma, up(\Delta, st, a1, l), a, l) \\
&a' : RET;
\end{aligned}$$

Dabei erzeugt  $dt$  Cpde für die Prozedurrümpfe von  $\Delta$ ; in der zugehörigen Symbolta-  
 belle sind nicht nur Konstanten und Variablen von  $\Delta$ , sondern auch die Prozedur-  
 Bezeichner eingetragen (Rekursion!).  $ct$  enthält die übergebene Startadresse und  
 erzeugt Code, der mit einem Rücksprung endet.

$$\begin{aligned}
 dt(\Delta_c \Delta_v \Delta_p, st, a, l) &:= dt(\Delta_p, st, a, l) \\
 dt(\varepsilon, st, a, l) &:= \varepsilon \\
 dt(procI_1; B_1; \dots; procI_n; B_n; st, a, l) &:= bt(B_1, st, a1, l + 1) \\
 &\vdots \\
 &bt(B_n, st, an, l + 1)
 \end{aligned}$$

**Beachte:**

1.  $st(I_j) = (proc, a_j, l + 1, off)$ , weil  $bt$  die Funktionen  $bt$  und  $up$  mit dem gleichen  
 Anfangsparameter aufruft und beide Funktionen aus diesen in gleicher Weise die  
 Adressen für die Prozedurrümpfe erzeugen.
2.  $bt$  wird mit dem Level  $l + 1$  aufgerufen.

$$\begin{aligned}
 ct(I := E, st, a, l) &:= \text{if } st(I) = (var, dl, off) \text{ then} \\
 &\quad ct(E, st, a, l) \\
 &\quad a' : STO(l - dl, off)
 \end{aligned}$$

$$\begin{aligned}
 ct(I(), st, a, l) &:= \text{if } st(I) = (proc, ca, dl, loc) \text{ then} \\
 &\quad a : CALL(ca, l - dl, loc)
 \end{aligned}$$

$$\begin{aligned}
 ct(\Gamma_1; \Gamma_2, st, a, l) &:= \\
 &\quad ct(\Gamma_1, st, a, l) \\
 &\quad ct(\Gamma_2, st, a, l)
 \end{aligned}$$

$$\begin{aligned}
 ct(\text{if } BE \text{ then } \Gamma_1 \text{ else } \Gamma_2, st, a, l) &:= \\
 &\quad sbt(BE, st, a, l) \\
 &\quad a' : JFALSEa'' \\
 &\quad ct(\Gamma_1, st, a' + 1, l) \\
 &\quad a'' - 1 : JMPa''' \\
 &\quad ct(\Gamma_2, st, a'', l) \\
 &\quad a''' : \dots
 \end{aligned}$$

$$\begin{aligned}
 ct(\text{while } BE \text{ do } \Gamma, st, a, l) &:= \\
 &\quad sbt(B, st, a, l) \\
 &\quad a' : JFALSEa'' + 1 \\
 &\quad ct(\Gamma, st, a' + 1, l) \\
 &\quad a'' : JMPa
 \end{aligned}$$

$$et(Z, st, a, l) := a : LITZ$$

$et(I, st, a, l) :=$   
 if  $st(I) = (const, Z)$  then  $a : LITZ$   
 if  $st(I) = (var, dl, off)$  then  $a : LOAD(l - dl, off)$

$et(E_1 + E_2, st, a, l) :=$   
 $et(E_1, st, a, l)$   
 $et(E_2, st, a', l)$   
 $a'' : ADD$

$sbt(E_1 < E_2, st, a, l) :=$   
 $et(E_1, st, a, l)$   
 $et(E_2, st, a', l)$   
 $a'' : LT$

$sbt(notBE, st, a', l) :=$   
 $sbt(BE, st, a, l)$   
 $a' : not$

$sbt(BE_1 and BE_2, st, a, l) :=$   
 $sbt(BE_1, st, a, l)$   
 $sbt(BE_2, st, a', l)$   
 $a'' : AND$

**Satz (Korrektheit der Übersetzung)** Für jedes  $P \in Prog^{(n)}$  und  $(z_1, \dots, z_n), (z'_1, \dots, z'_n) \in \mathbb{Z}$  gilt:  $\mathcal{M}[[P]](z_1, \dots, z_n) = (z'_1, \dots, z'_n) \leftrightarrow \exists \mathcal{J}[[trans(P)]](1, \varepsilon, 0 : 0 : 0 : z_1 : \dots : z_n) = (0, \varepsilon, 0 : 0 : 0 : z'_1 : \dots : z'_n)$

**Beweis:** M. Mohnen: A Compiler Correctness Proof for the Static Link Technique, Fund. Inf. 29 (1997)

### Beispiel (Fakultät)

```

in/out X;
var E;
proc F;
  if (1<x) then
    begin
      E:=E*X;
      X:=X-1;
      F()
    end
begin
  E:=1;
  F();
  X:=E;
end.

```



$$\begin{aligned} \text{trans}(in/outX; \Delta\Gamma) = \\ 1 : CALL(a_\Gamma, 0, 1) \\ 2 : JMP0; \\ bt(\Delta\Gamma, st_{I/O}, a_\Gamma, 1) \end{aligned}$$

mit  $st_{I/O} = (var, 0, 1)$ .

**Jumping-Code für boolesche Ausdrücke** Übersetzung boolescher Ausdrücke mit nicht-strikter Semantik durch Sprungbefehle anstelle logischer Befehle.

- Idee: Vererbung von Sprungzeilen für boolesche Ergebnisse
- $nbt : BExp \times Tab \times Adr^3 \dashrightarrow AM - Code$ 
  1. freie Anfangsadresse
  2. true-Adresse
  3. false-Adresse

mit geeigneter Modifikation der Übersetzung von Verzweigung und Iteration:

$$\begin{aligned} nct(\text{if } BE \text{ then } \Gamma_1 \text{ else } \Gamma_2, st, a, l) := \\ nbt(BE, st, a, a_t, a_f, l) \\ nct(\Gamma_1, st, a_t, l) \\ a' : JMPA''; \\ nct(\Gamma_2, st, a_f, l) \end{aligned}$$

$$\begin{aligned} nct(\text{while } BE \text{ do } \Gamma, st, a, l) := \\ nbt(BE, st, a, a_t, a_f, l) \\ nct(\Gamma, st, a_t, l) \\ a_t - 1 : JMPa \end{aligned}$$

$$\begin{aligned} nbt(E_1 < E_2, st, a, a_t, a_f, l) := \\ et(E_1, st, a, l) \\ et(E_2, st, a', l) \\ a'' : LT \\ a'' + 1 : JFALSEa_f \\ a'' + 2 : JMPa_t \end{aligned}$$

$$\begin{aligned} nbt(\text{not } BE, st, a, a_t, a_f, l) := \\ nbt(BE, st, a, a_f, a_t, l) \end{aligned}$$

$$\begin{aligned} nbt(BE_1 \text{ and } BE_2, st, a, a_t, a_f, l) := \\ nbt(BE_1, st, a, a', a_f, l) \\ nbt(BE_2, st, a', a_t, a_f, l) \end{aligned}$$

$$\begin{aligned} nbt(BE_1 \text{ or } BE_2, st, a, a_t, a_f, l) := \\ nbt(BE_1, st, a, a_t, a', l) \\ nbt(BE_2, st, a', a_t, a_f, l) \end{aligned}$$

**Prozeduren mit Parametern** Erweiterung von BPS um Wert- und Variablenparameter für Prozeduren.

### Sukzessiver Aufbau eines Frames

1. Berechnung der aktuellen Parameter
2. Berechnen des statischen Verweises (Indexregister  $R$ )
3. Sprung zur aufgerufenen Prozedur mit Eintrag der Rücksprungadresse
4. Alten FP (Framepointer) als dl speichern
5. Speicherplatz für lokale Variablen bereitstellen

**Befehlssatz (ZPP)** Arithmetische, logische und Sprungbefehle wie bisher.

CALL $ca$	$SP \leftarrow SP - 1; \langle SP \rangle \leftarrow IC + 1; IC \leftarrow ca$	
RET $k$	$IC \leftarrow \langle SP \rangle; SP \leftarrow SP + k + 1$	
PUSH $z$	$SP \leftarrow SP - 1; \langle SP \rangle \leftarrow z$	
PUSH $FP$	analog	
PUSH $\langle FP \rangle$	analog	
PUSH $\langle R + 2 \rangle$	analog	
POP $FP$	$FP \leftarrow \langle SP \rangle; SP \leftarrow SP + 1$	
POP $\langle n \rangle$	$S_n \leftarrow \langle SP \rangle; SP \leftarrow SP + 1$	$S_n :=$ Stackzelle mit Adresse $n$
SUB $SP, n$	$SP \leftarrow SP - n$	
LOAD $FP, SP$	$FP \leftarrow SP$	
LOAD $R, \langle n \rangle$	analog	
LOAD $R, \langle R + 2 \rangle$	analog	

**Ein- / Ausgabe**  $P = in/outI_1, \dots, I_n; B.$

**Anfangstabelle**  $st_{I/O}(I_j) := (var, 0, 3 + n - j)$

Zur weiteren Erläuterung siehe Folien (CB 4-19) bis (CB 4-23).

## 5.2 Übersetzung von Datenstrukturen

- Datenstrukturen  $\rightsquigarrow$  Variablen mit Komponenten, strukturierter Zustandsraum
- Abstrakte Maschine: Weiterhin lineare Speicherstruktur, Speicherzellen für atomare Daten
- Übersetzungsaufgabe: Abbildung des strukturierten Zustandsraumes auf linearen Speicherbereich
  - statische Datenstrukturen: Speicherbedarf zur Übersetzungszeit bekannt

- dynamische Datenstrukturen: Speicherbedarf laufzeitabhängig
- Heap / Halde als zusätzliche Maschinenkomponente
- Zeigertypen, garbage collection / Speicherbereinigung