

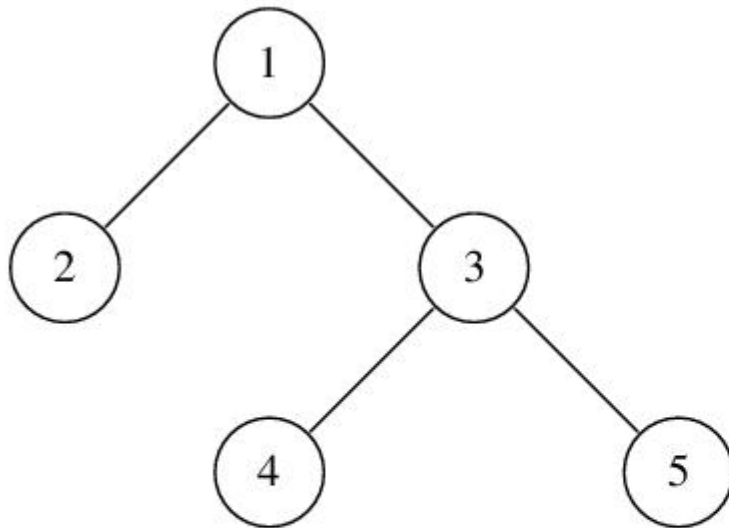
Syntaxanalyse (Parsing)

- **Eingabe:** Tokensequenz
- **Ausgabe:** Parsebaum (in der Regel)
 - innere Knoten: Nichtterminale
 - Blätter: Terminale (von links nach rechts: Eingabe)
- **analysierbare Grammatiken:**
 - kontextfrei (linke Seite: ein Nichtterminal)
 - eindeutig (nur ein Parsebaum pro Eingabe)
- **Erkenner:** Stackautomaten
 - Rekursion (Laufzeitstack)
 - Tabellen-gesteuert (separater Stack)

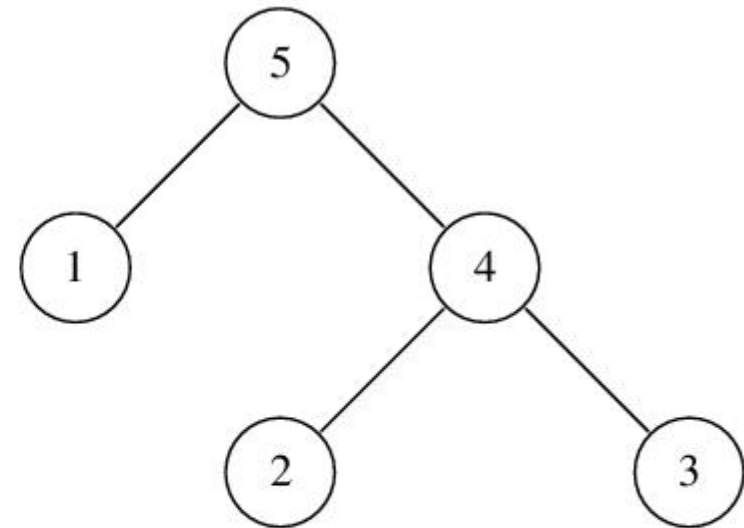
Zwei Strategien

	LL(k)	LR(k)
Eingabe lesen	links nach rechts (L_)	links nach rechts (L_)
Produktion generieren	links nach rechts (_L)	rechts nach links (_R)
Vorausschau	k Token	k Token
Zeitkomplexität	linear	linear
linksrekursive Produktionen	nicht möglich	möglich und effizient
Baumrekonstruktion	top-down	bottom-up
Bild	Abb. 2.51	Abb. 2.52

Preorder (LL), Postorder (LR)



(a) Pre-order



(b) Post-order

Figure 2.50 A tree with its nodes numbered in pre-order and post-order.

Rekursiver Abstieg bei LL

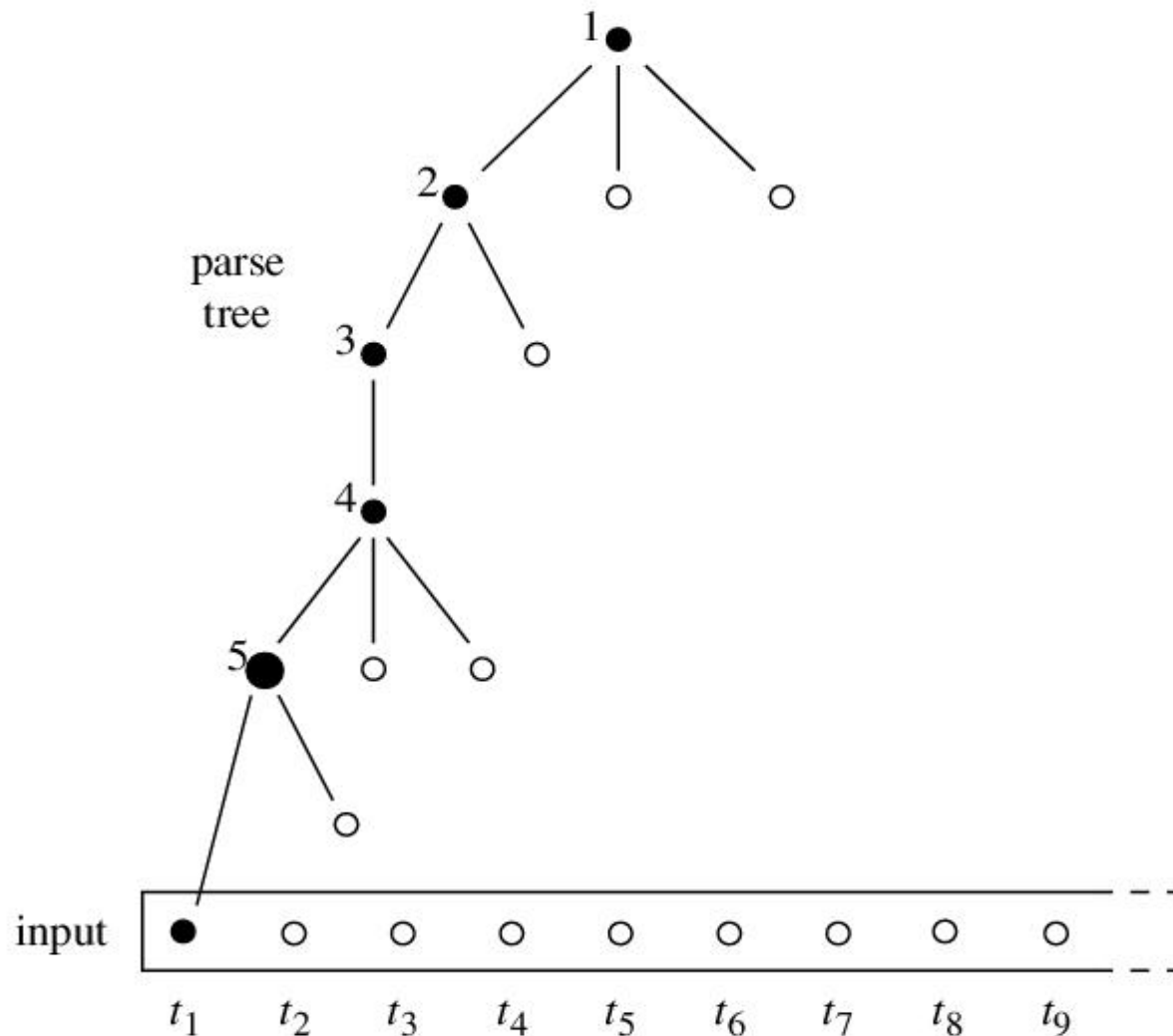


Figure 2.51 A top-down parser recognizing the first token in the input.

Bottom-Up Konstruktion bei LR

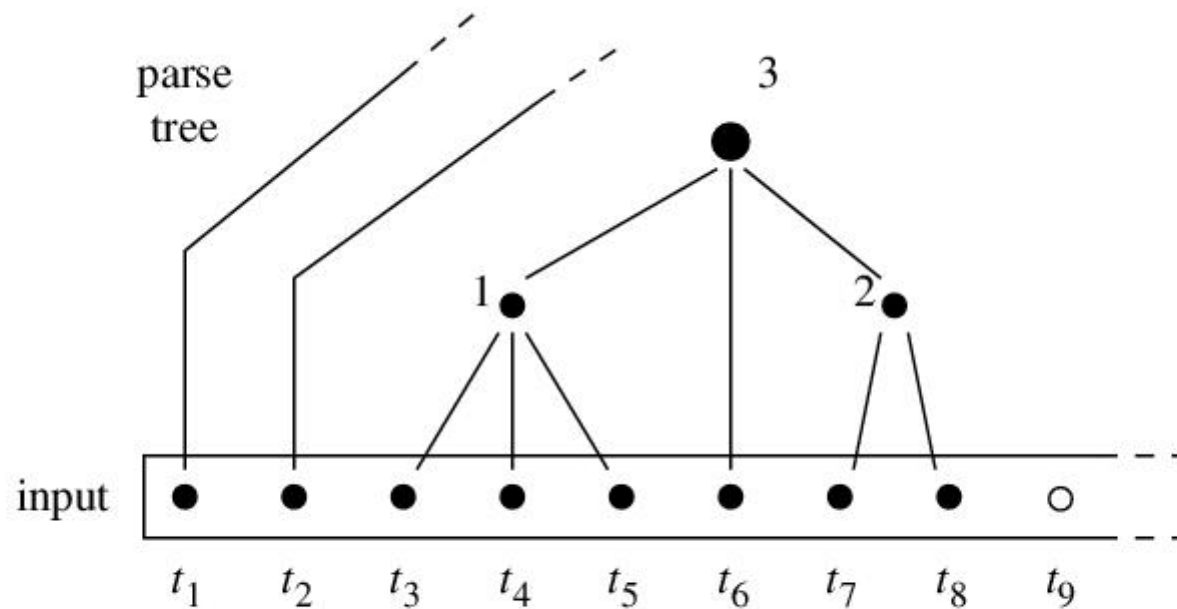


Figure 2.52 A bottom-up parser constructing its first, second, and third nodes.

LL(1) mit rekursivem Abstieg

- **Beispielgrammatik:**

input	→ expression EOF
expression	→ term rest_expression
term	→ IDENTIFIER parenthesized_expression
parenthesized_expression	→ '(' expression ')'
rest_expression	→ '+' expression ε

- **Beispiel-"Erkenner"** (erzeugt keinen Parsebaum) mit rekursivem Abstieg:

- **Beispielwort:**

IDENTIFIER + (IDENTIFIER + IDENTIFIER) EOF

- **Implementierung (Abb. 2.54-2.55)**

- pro Nichtterminal eine Funktion

- gdw. aus diesem Nichtterminal am Anfang das Eingabesymbol erzeugt werden kann (**Produktion muss eindeutig sein**)

(a) wird das Eingabetoken verbraucht

(b) gibt die Funktion den Wert true (**≠0 in C**) zurück

```

int input(void) {
    return expression() && require(token(EOF));
}

int expression(void) {
    return term() && require(rest_expression());
}

int term(void) {
    return token(IDENTIFIER) || parenthesized_expression();
}

int parenthesized_expression(void) {
    return token('(') && require(expression()) && require(token(')'));
}

int rest_expression(void) {
    return token('+') && require(expression()) || 1;
}

int token(int tk) {
    if (tk != Token.class) return 0;
    get_next_token(); return 1;
}

int require(int found) {
    if (!found) error();
    return 1;
}

```

Figure 2.54 A recursive descent parser/recognizer for the grammar of Figure 2.53.

```
int main(void) {
    start_lex(); get_next_token();
    require(input());
    return 0;
}

void error(void) {
    printf("Error in expression\n"); exit(1);
}
```

Figure 2.55 Driver for the recursive descent parser/recognizer.

Probleme beim LL(1)-Parsing

... der Parser ist nicht immer korrekt ...

- **Ausdrucksprache mit Arrays:**

term \rightarrow IDENTIFIER | Indexed_element
| parenthesized_expression
indexed_element \rightarrow IDENTIFIER '[' expression ']'

nicht links-faktoriert, Auswahl der Produktion unklar

- **Sprache {ab, aab}: $S \rightarrow A 'a' 'b'$**

$A \rightarrow 'a' \mid \epsilon$ Auswahl auch unklar

Lösung des Problems: Substitution von A und Links-Faktorisierung:

$S \rightarrow 'a' R$
 $R \rightarrow 'a' 'b' \mid 'b'$

jetzt ist das zweite 'a' optional, nicht das erste

- **Ausdrucksprache mit links-assoziativer Subtraktion**

expression \rightarrow expression '-' term

links-rekursiv, Parser gerät in Endlosschleife

Konflikte bei der Auswahl

$$\begin{aligned} S &\rightarrow A \ 'a' \ 'b' \\ A &\rightarrow \ 'a' \ | \ \varepsilon \end{aligned}$$

Figure 2.56 A simple grammar with a FIRST/FOLLOW conflict.

```
int S(void) {
    return A() && require(token('a')) && require(token('b'));
}

int A(void) {
    return token('a') || 1; /* erkennt nicht "ab" */
       1 || token('a'); /* erkennt nicht "aab" */
}
```

Figure 2.57 A faulty recursive recognizer for grammar of Figure 2.56.

LL(1)-Parsergenerierung: FIRST-Mengen

*... geben an, mit welchen Zeichen eine Satzform beginnen kann ...
(falls Satzform das leere Wort herleiten kann, zusätzlich ε)*

Inferenzregeln eines Hüllenalgorithmus (vgl. Abb. 2.58):

• $N \rightarrow \alpha$

$$\text{FIRST}(N) \supseteq \text{FIRST}(\alpha)$$

“füge alle Zeichen aus $\text{FIRST}(\alpha)$ zu $\text{FIRST}(N)$ hinzu”

• $\alpha = A\beta$

♦ $\neg(A \rightarrow \varepsilon)$ “aus A kann nicht der leere String erzeugt werden”

$$\text{FIRST}(\alpha) \supseteq \text{FIRST}(A)$$

♦ $A \rightarrow \varepsilon$

$$\text{FIRST}(\alpha) \supseteq (\text{FIRST}(A) \setminus \{\varepsilon\}) \cup \text{FIRST}(\beta)$$

Data definitions:

1. Token sets called FIRST sets for all terminals, non-terminals and alternatives of non-terminals in G .
2. A token set called FIRST for each alternative tail in G ; an alternative tail is a sequence of zero or more grammar symbols α if $A\alpha$ is an alternative or alternative tail in G .

Initializations:

1. For all terminals T , set $\text{FIRST}(T)$ to $\{T\}$.
2. For all non-terminals N , set $\text{FIRST}(N)$ to the empty set.
3. For all non-empty alternatives and alternative tails α , set $\text{FIRST}(\alpha)$ to the empty set.
4. Set the FIRST set of all empty alternatives and alternative tails to $\{\epsilon\}$.

Inference rules:

1. For each rule $N \rightarrow \alpha$ in G , $\text{FIRST}(N)$ must contain all tokens in $\text{FIRST}(\alpha)$, including ϵ if $\text{FIRST}(\alpha)$ contains it.
2. For each alternative or alternative tail α of the form $A\beta$, $\text{FIRST}(\alpha)$ must contain all tokens in $\text{FIRST}(A)$, excluding ϵ , should $\text{FIRST}(A)$ contain it.
3. For each alternative or alternative tail α of the form $A\beta$ and $\text{FIRST}(A)$ contains ϵ , $\text{FIRST}(\alpha)$ must contain all tokens in $\text{FIRST}(\beta)$, including ϵ if $\text{FIRST}(\beta)$ contains it.

Figure 2.58 Closure algorithm for computing the FIRST sets in a grammar G .

Rule/alternative (tail)	FIRST set
input	{ }
expression EOF	{ }
EOF	{ EOF }
expression	{ }
term rest_expression	{ }
rest_expression	{ }
term	{ }
IDENTIFIER	{ IDENTIFIER }
parenthesized_expression	{ }
parenthesized_expression	{ }
'(' expression ')'	{ '(' }
expression ')'	{ }
')'	{ ')' }
rest_expression	{ }
'+' expression	{ '+' }
expression	{ }
ϵ	{ ϵ }

Figure 2.59 The initial FIRST sets.

Rule/alternative (tail)	FIRST set
input	{ IDENTIFIER '(' }
expression EOF	{ IDENTIFIER '(' }
EOF	{ EOF }
expression	{ IDENTIFIER '(' }
term rest_expression	{ IDENTIFIER '(' }
rest_expression	{ '+' ϵ }
term	{ IDENTIFIER '(' }
IDENTIFIER	{ IDENTIFIER }
parenthesized_expression	{ '(' }
parenthesized_expression	{ '(' }
'(' expression ')'	{ '(' }
expression ')'	{ IDENTIFIER '(' }
')'	{ ')'
rest_expression	{ '+' ϵ }
'+' expression	{ '+' }
expression	{ IDENTIFIER '(' }
ϵ	{ ϵ }

Figure 2.60 The final FIRST sets.

LL(1)-Parsergenerierung: FOLLOW-Mengen

... geben an, welche Zeichen auf ein Nichtterminal folgen können ...

Inferenzregeln eines Hüllalgorithmus (vgl. Abb. 2.62):

• $M \rightarrow \alpha N \beta$

♦ $\neg(\beta \rightarrow \varepsilon)$

$\text{FOLLOW}(N) \supseteq \text{FIRST}(\beta)$

♦ $\beta \rightarrow \varepsilon$ (beinhaltet, dass β nicht existent ist)

$\text{FOLLOW}(N) \supseteq (\text{FIRST}(\beta) \setminus \{\varepsilon\}) \cup \text{FOLLOW}(M)$

LL(1)-Parsergenerierung: FOLLOW-Mengen

Data definitions:

1. Token sets called FOLLOW sets for all non-terminals in G .
2. Token sets called FIRST sets for all alternatives and alternative tails in G .

Initializations:

1. For all non-terminals N , set FOLLOW(N) to the empty set.
2. Set all FIRST sets to the values determined by the algorithm for FIRST sets.

Inference rules:

1. For each rule of the form $M \rightarrow \alpha N \beta$ in G , FOLLOW(N) must contain all tokens in FIRST(β), excluding ϵ , should FIRST(β) contain it.
2. For each rule of the form $M \rightarrow \alpha N \beta$ in G where FIRST(β) contains ϵ , FOLLOW(N) must contain all tokens in FOLLOW(M).

Figure 2.62 Closure algorithm for computing the FOLLOW sets in grammar G .

Beispiel: FIRST/FOLLOW-Konstruktion

	FIRST	FOLLOW
input → exp EOF	{ ID, '(' }	{ }
exp → term rest_exp	{ ID, '(' }	{ EOF, ')' }
term → ID par_exp	{ ID, '(' }	{ '+', EOF, ')' }
par_exp → '(' exp ')'	{ '(' }	{ '+', EOF, ')' }
rest_exp → '+' exp ε	{ '+', ε }	{ EOF, ')' }

Prädiktiver LL(1)-Parser

zusätzliche Herausforderungen:

- Behandlung des Kleene-Sterns
(verkompliziert FIRST/FOLLOW-Berechnungen)
- aussagekräftige Fehlermeldungen
- Erzeugung des Syntaxbaums
- Behandlung von Syntaxfehlern
- Optimierungen des Parsens

FIRST/FIRST-Konflikte

Beispielgrammatik:

```
term          → IDENTIFIER | Indexed_element  
              | parenthesized_expression  
indexed_element → IDENTIFIER '[' expression ']'
```

$\text{FIRST}(\text{indexed_element}) = \{ \text{IDENTIFIER} \}$

Parser:

```
void term(void) {  
    switch (Token.class) {  
    case IDENTIFIER: token(IDENTIFIER); break;  
    case IDENTIFIER: indexed_element(); break;  
    case '(':        parenthesized_expression();  
                    break;  
    default:        error(); }  
}
```

bei Links-Rekursion: auch FIRST/FIRST-Konflikte!

FIRST/FOLLOW-Konflikte

	FIRST	FOLLOW
$S \rightarrow A \text{ 'a' 'b'}$	{ 'a' }	{ }
$A \rightarrow \text{'a'} \mid \varepsilon$	{ 'a', ε }	{ 'a' }

Parser:

```
void S(void) {
    switch(Token.class) {
        case 'a': A(); token('a'); token('b'); break;
        default: error(); }
    }
void A(void) {
    switch(Token.class) {
        case 'a': token('a'); break;
        case 'a': break;
        default: error(); }
    }
```

LL(1)-Anforderungen

... eindeutige Auswahl jeder Alternative ...

- Voraussetzung für Auswahl von Alternative $N \rightarrow A_k$ für Nichtterminal N
 - ◆ $\neg(A_k \xrightarrow{\epsilon})$: Lookahead-Token $\in \text{FIRST}(A_k)$
 - ◆ $A_k \xrightarrow{\epsilon}$: Lookahead-Token $\in \text{FIRST}(A_k) \cup \text{FOLLOW}(N)$
- Bedingungen für eine eindeutige Bestimmung von A_k :
 - keine FIRST/FIRST-Konflikte
 - keine FIRST/FOLLOW-Konflikte
 - höchstens eine Alternative A_k , für die gilt: $A_k \xrightarrow{\epsilon}$

Grammatikanpassung an LL(1)

Problem

- **gegeben:** Grammatik $G \notin \text{LL}(1)$
- **gesucht:** LL(1)-Grammatik G' , so dass $L(G')=L(G)$
- **Vorteil von G' :** automatische Parsererstellung möglich
- **Nachteil von G' :** Struktur des Parsebaum passt nicht zu G

Tricks und Guidelines

- Substitution von Nichtterminalen durch ihre rechten Seiten
- Linksfaktorisierung
- Elimination von Linksrekursion

Tricks und Guidelines

	Beispiel 1	Beispiel 2
	$S \rightarrow A \text{ 'a' 'b'}$ $A \rightarrow \text{'a' } \ \epsilon$ <p>FIRST/FOLLOW-Konflikt</p>	$\text{term} \rightarrow \text{ID} \text{par_exp}$ $ \text{ind_el}$ $\text{ind_el} \rightarrow \text{ID '[' exp ']'}$ <p>indirekter FIRST/FIRST-Konflikt</p>
1. Substitution ↓	$S \rightarrow \text{'a' 'a' 'b'}$ $ \text{'a' 'b'}$ <p>FIRST/FIRST-Konflikt</p>	$\text{term} \rightarrow \text{ID} \text{par_exp}$ $ \text{ID '[' exp ']'}$ <p>FIRST/FIRST-Konflikt</p>
2. Links-faktorisierung ↓	$S \rightarrow \text{'a' } R$ $R \rightarrow \text{'a' 'b' } \text{'b'}$ <p>LL(1)</p>	$\text{term} \rightarrow \text{ID } R \text{par_exp}$ $R \rightarrow \text{'[' exp ']' } \ \epsilon$ <p>LL(1)</p>

Tricks und Guidelines

3. Elimination von Linksrekursion

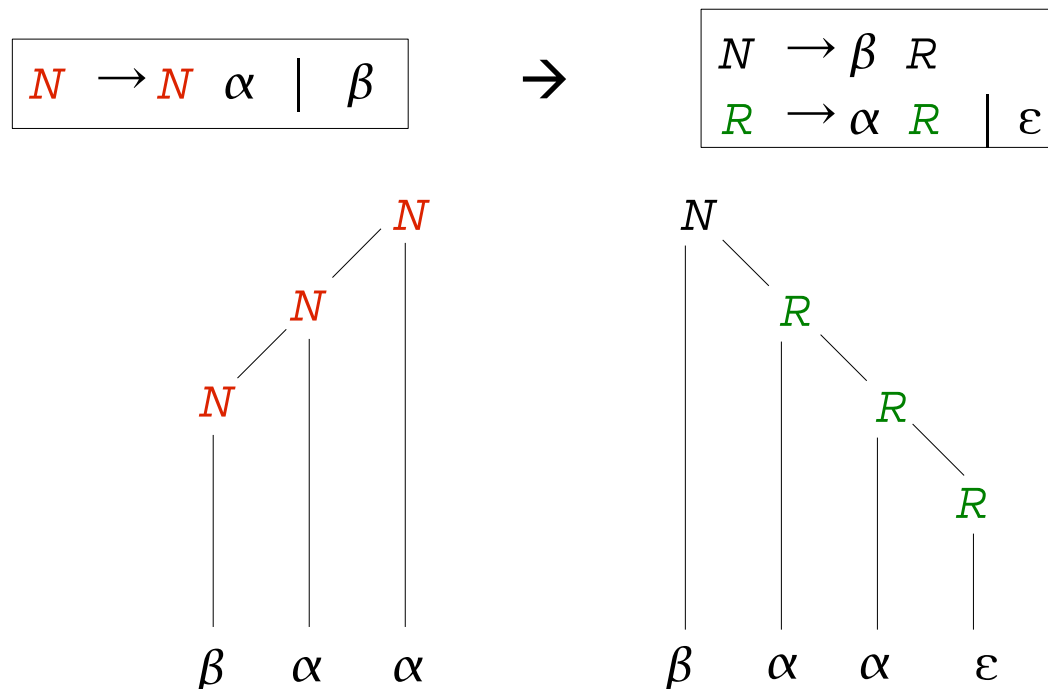
Beispiel:

$exp \rightarrow exp \text{ '-' } term \mid term$

↓

$exp \rightarrow term \text{ rest}$
$rest \rightarrow \text{'-' } term \text{ rest} \mid \epsilon$

Regel zur Elimination *direkter* Linksrekursion:



Tricks und Guidelines

Arten von Linksrekursion

- **direkt:** $N \rightarrow N \alpha \mid \beta$
- **indirekt:** $N \rightarrow M \alpha \mid \beta$
 $M \rightarrow N \gamma \mid \delta$
- **versteckt:** $N \rightarrow \alpha N \mid \beta$, wobei $\alpha \xrightarrow{*} \varepsilon$

Semantische Annotationen

... mit Markierungsregeln ...

- **Frage:** wie die semantischen Ergebnisse verwerten?
- **Idee:** füge der Grammatik Produktionen hinzu, die
 - alle nur ε generieren
 - die an der Stelle geforderten semantischen Aktionen ausführen
- **Beispiel in C** (Referenzen (`int *`) für Ein-/Ausgabeparameter)

◆ linksrekursiv

$$\begin{array}{l} \text{exp}(\text{int } *e) \rightarrow \text{exp}(\text{int } *e) \text{ '-' term}(\text{int } *t) \quad \{ *e \text{ -= } *t; \} \\ \quad \quad \quad | \text{ term}(\text{int } *t) \quad \quad \quad \{ *e \text{ = } *t; \} \end{array}$$

◆ nach Elimination von Linksrekursion

$$\begin{array}{l} \text{exp}(\text{int } *e) \rightarrow \text{term}(\text{int } *t) \quad \{ *e \text{ = } *t; \} \text{ rest}(\text{int } *e) \\ \text{rest}(\text{int } *e) \rightarrow \text{'-' term}(\text{int } *t) \quad \{ *e \text{ -= } *t; \} \text{ rest}(\text{int } *e) \\ \quad \quad \quad | \varepsilon \end{array}$$