



Theoretische Informatik

Kontextfreie Sprachen
und Parser



Inhalt

1. Grammatiken und Sprachen
 - .. Kontextfreie Grammatiken
 - .. Herleitungen, Linksherleitungen
 - .. Sprachen zu einer Grammatik
 - .. Äquivalenz
 - .. Chomsky-Normalform
 - .. Wortproblem, CYK
 - .. Pumping Lemma für CF-Sprachen

2. Stackautomaten
 - .. Definitionen und Beispiele
 - .. Konfigurationen, Läufe,
 - .. Sprache eines Stackautomaten
 - .. Parser

3. Parser
 - .. Mehrdeutigkeit
 - .. Bottom Up, Top Down
 - .. Recursive descent Parser
 - .. First, Follow
 - .. Semantische Aktionen
4. Shift-Reduce Parser
 - .. Konflikte
 - .. LR(0)-Zustände
 - .. Automat zur Grammatik
 - .. Shift-Reduce, Reduce-Reduce Konflikte
 - .. Präzedenz und Assoziativität
 - .. lex und yacc
 - .. Arbeitsweise eines Compilers



Der zweistufige Aufbau von Sprachen

Aus Alphabet der Token wird der syntaktische Anteil der Sprache aufgebaut,

Σ = Menge der Token

L = Menge aller syntaktisch korrekter Sätze.

Auf dieser Stufe sind Sprachen i.A. nicht regulär. (Geschachtelte Strukturen kann man nicht regulär beschreiben - Klammersprache).

Für die Definition muss was neues her:

Grammatiken



Syntax - Semantik

- n Grammatik beschreibt die **Syntax** (Struktur) von Sätzen, nicht deren **Semantik** (Bedeutung).
- n Semantik hängt von Syntax ab, aber nicht jeder syntaktisch korrekte Satz muss eine Semantik haben.

Syntaktisch falscher Satz:

gestohlen der hat Gans Fuchs die.

Syntaktisch richtiger, semantisch falscher Satz:

der Gans hat die Fuchs gestohlen.

Syntaktisch und semantisch richtiger Satz:

der Fuchs hat die Gans gestohlen.





Syntax - Semantik in PASCAL

Syntaktisch falsches Programm:

```
Test Program.  
  i for 1 do n ; to n+1 := n
```

Syntaktisch richtiges, semantisch falsches Programm:

```
Program Test ;  
Var i : Real;  
Begin  
  For i := 1 To n Do i := n-1  
End.
```

Syntaktisch und semantisch korrektes Programm:

```
Program Test ;  
Var i,n : Integer;  
Begin  
  For i := 1 To 50 Do n := n+i  
End.
```



Sprachen und Grammatiken

Eine **Grammatik** ist eine Menge von Regeln mit deren Hilfe „syntaktisch korrekte“ Worte (Sätze) der Sprache konstruiert werden können.

Erzeugen

Mit Hilfe der Grammatik kann man Sätze einer Sprache erzeugen. Die Menge dieser Sätze bilden die zur Grammatik G gehörende Sprache, $L(G)$.

Zerlegen

Ist ein Satz gegeben, so möchte man feststellen, ob er zur Sprache $L(G)$ gehört oder nicht. Dabei wird auch die Struktur des Satzes erkannt.



Beispiel einer Grammatik

Java

```
Stmt      :: while ( Expr ) Stmt | if (Expr) Stmt else Stmt
           | id = Expr ; | { Stmts } ...

Expr      :: Expr op Expr | ( Expr ) | num | id ...

Stmts     :: Stmt Stmts | ε

etc.
```

Bestandteile:

```
Token      : while, (, ), if, else, id, =, ;, {, }, op, num, id
Variablen  : Stmt, Expr, Stmts,
Satzformen : ( Expr ) , id = Expr ; , ε , ...
```

Die Zeichen

```
:: (definiert als)
| (Alternative)
```

sind Sonderzeichen für Grammatikdefinitionen



Kontextfreie Grammatik

Gegeben:

- T - ein Alphabet (Terminale, Token)
- V - eine Menge von Variablen (Nonterminale)
- $S \in V$ - Startsymbol

mit $T \cap V = \emptyset$

Eine **Kontextfreie Grammatik** ist ein Quadrupel $G = (V, T, P, S)$
mit $T \cap V = \emptyset$,

$S \in V$ und $P \subseteq V \times (T \cup V)^*$

Begriffe

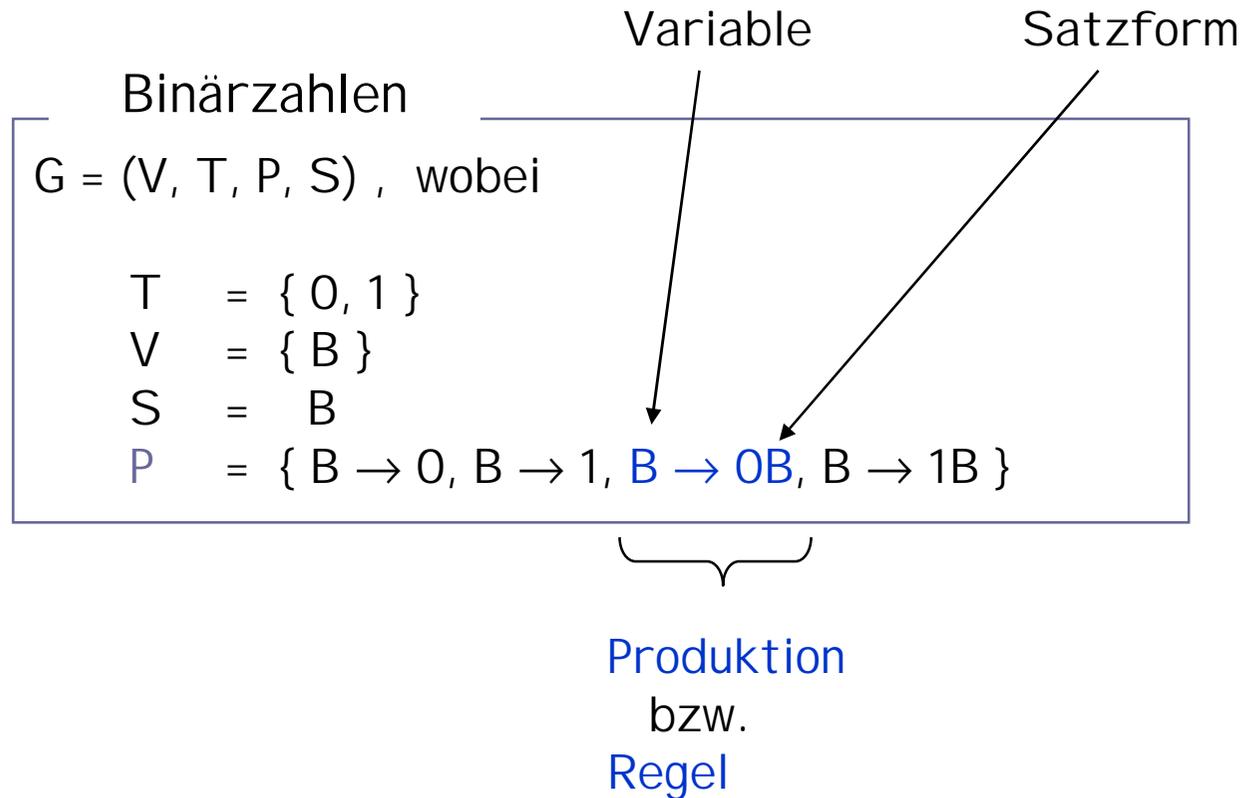
Satzform : Element α von $(V \cup T)^*$

Produktion : Jedes $(A, \alpha) \in P$. Notation: $A \rightarrow \alpha$

Startsymbol : S.



CF-Grammatik - Beispiele





CF-Grammatik - Beispiel

Ausdrücke

$G = (V, T, P, S)$, wobei

$T = \{ \text{id, num, +, *, }, (\}$

$V = \{ \text{Term, Factor, Expr} \}$

$S = \text{Expr}$

$P = \{ \text{Expr} \rightarrow \text{Term}, \text{Expr} \rightarrow \text{Expr} + \text{Term},$
 $\text{Term} \rightarrow \text{Factor}, \text{Term} \rightarrow \text{Term} * \text{Factor},$
 $\text{Factor} \rightarrow \text{id}, \text{Factor} \rightarrow \text{num},$
 $\text{Factor} \rightarrow (\text{Expr})$
 $\}$

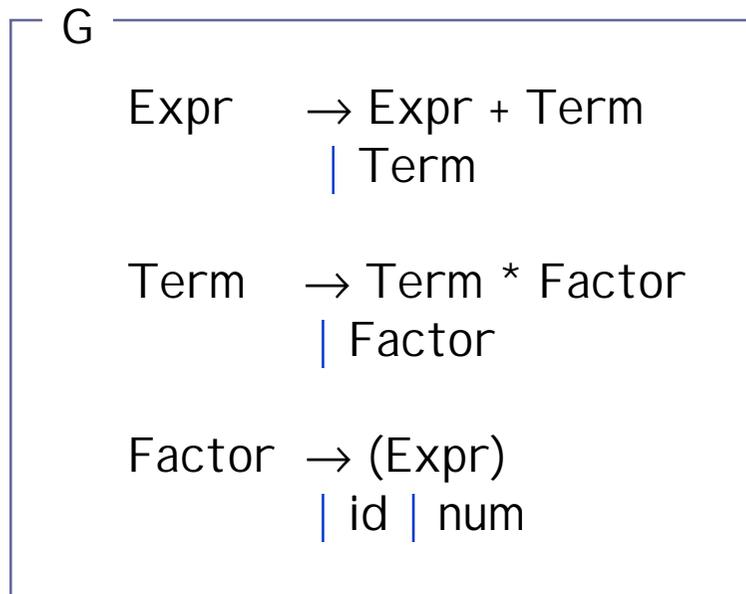


Konventionen

- n **Variablen** : Alle Symbole für die es mind. eine Produktion gibt
- n **Terminale** : alle anderen Symbole
- n **Startsymbol** :
 - Default: Variable der ersten Produktion.
- n Produktionen mit gleicher linker Seite zusammenfassen.
 - Rechte Seiten durch | trennen.



Das vorige Beispiele kann man dann schreiben als :





Ersetzung - Herleitung

Eine Produktion $A \rightarrow \beta$ erlaubt, in einer Satzform $\alpha A \gamma$ die Variable A durch β ersetzen. Resultat ist die Satzform $\alpha \beta \gamma$.

Eine Produktion:

$$A \rightarrow \beta$$

Eine Ersetzung: $\alpha A \gamma \Rightarrow \alpha \beta \gamma$

$\sigma \Rightarrow^* \tau$ falls es $\alpha_1, \dots, \alpha_n$ gibt mit
 $\sigma = \alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_n = \tau$

Herleitung



Beispiel einer Herleitung

Arithmetik

$$\text{Expr} \rightarrow \text{Expr} + \text{Term} \\ | \text{Term}$$
$$\text{Term} \rightarrow \text{Term} * \text{Factor} \\ | \text{Factor}$$
$$\text{Factor} \rightarrow (\text{Expr}) \\ | \text{id} | \text{num}$$

Es gilt also:

$$\text{Expr} \Rightarrow^* \text{id} * (\text{num} + \text{id})$$

Eine Herleitung:

$$\begin{aligned} \text{Expr} &\Rightarrow \text{Term} \\ &\Rightarrow \text{Term} * \text{Factor} \\ &\Rightarrow \text{Term} * (\text{Expr}) \\ &\Rightarrow \text{Term} * (\text{Expr} + \text{Term}) \\ &\Rightarrow \text{Term} * (\text{Term} + \text{Term}) \\ &\Rightarrow \text{Factor} * (\text{Term} + \text{Term}) \\ &\Rightarrow \text{id} * (\text{Term} + \text{Term}) \\ &\Rightarrow \text{id} * (\text{Term} + \text{Factor}) \\ &\Rightarrow \text{id} * (\text{Term} + \text{id}) \\ &\Rightarrow \text{id} * (\text{Factor} + \text{id}) \\ &\Rightarrow \text{id} * (\text{num} + \text{id}). \end{aligned}$$



Eine Links-Herleitung

Arithmetik

$$\text{Expr} \rightarrow \text{Expr} + \text{Term} \\ | \text{Term}$$
$$\text{Term} \rightarrow \text{Term} * \text{Factor} \\ | \text{Factor}$$
$$\text{Factor} \rightarrow (\text{Expr}) \\ | \text{id} | \text{num}$$

Jedes Wort $w \in L(G)$ kann man auch durch eine **Linksherleitung** gewinnen

Bei einer **Links-Herleitung** wird immer das linkeste Nonterminal ersetzt

Eine Links-Herleitung:

$$\begin{aligned} \text{Expr} &\Rightarrow \text{Term} \\ &\Rightarrow \text{Term} * \text{Factor} \\ &\Rightarrow \text{Factor} * \text{Factor} \\ &\Rightarrow \text{id} * \text{Factor} \\ &\Rightarrow \text{id} * (\text{Expr}) \\ &\Rightarrow \text{id} * (\text{Expr} + \text{Term}) \\ &\Rightarrow \text{id} * (\text{Term} + \text{Term}) \\ &\Rightarrow \text{id} * (\text{Factor} + \text{Term}) \\ &\Rightarrow \text{id} * (\text{num} + \text{Term}) \\ &\Rightarrow \text{id} * (\text{num} + \text{Factor}) \\ &\Rightarrow \text{id} * (\text{num} + \text{id}). \end{aligned}$$



Links-Herleitung eines anderen Wortes

Arithmetik

$$\text{Expr} \rightarrow \text{Expr} + \text{Term} \\ | \text{Term}$$
$$\text{Term} \rightarrow \text{Term} * \text{Factor} \\ | \text{Factor}$$
$$\text{Factor} \rightarrow (\text{Expr}) \\ | \text{id} | \text{num}$$

Auswahl anderer Regeln
erzeugt i.A. andere Worte
 $w \in L(G)$

Bei einer **Links-Herleitung** wird immer das linkeste Nonterminal ersetzt

Eine Links-Herleitung unter Verwendung anderer Regeln:

$\text{Expr} \Rightarrow \text{Expr} + \text{Term}$
 $\Rightarrow \text{Term} + \text{Term}$
 $\Rightarrow \text{Factor} + \text{Term}$
 $\Rightarrow \text{num} + \text{Term}$
 $\Rightarrow \text{num} + \text{Factor}$
 $\Rightarrow \text{num} + (\text{Expr})$
 $\Rightarrow \text{num} + (\text{Term})$
 $\Rightarrow \text{num} + (\text{Term} * \text{Factor})$
 $\Rightarrow \text{num} + (\text{Factor} * \text{Factor})$
 $\Rightarrow \text{num} + (\text{id} * \text{Factor})$
 $\Rightarrow \text{num} + (\text{id} * \text{num}).$



Herleitung als Baum

Arithmetik

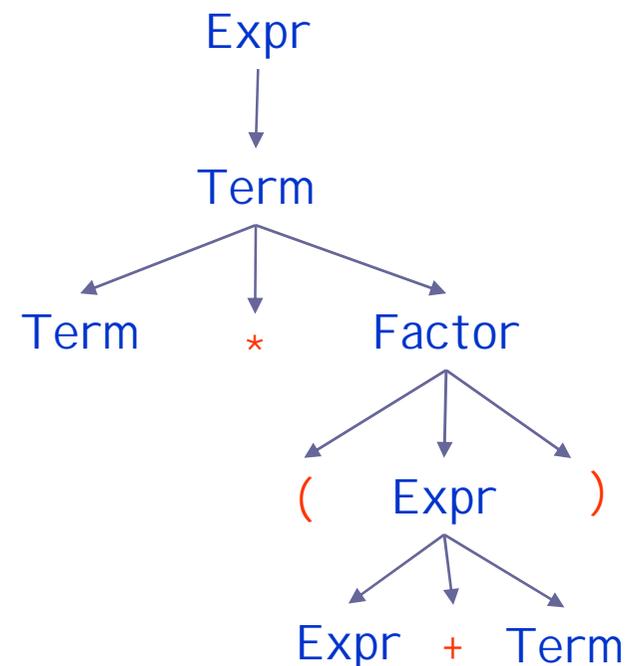
$$\text{Expr} \rightarrow \text{Expr} + \text{Term} \\ | \text{Term}$$
$$\text{Term} \rightarrow \text{Term} * \text{Factor} \\ | \text{Factor}$$
$$\text{Factor} \rightarrow (\text{Expr}) \\ | \text{id} | \text{num}$$

$\text{Expr} \Rightarrow \text{Term}$
 $\Rightarrow \text{Term} * \text{Factor}$
 $\Rightarrow \text{Term} * (\text{Expr})$
 $\Rightarrow \text{Term} * (\text{Expr} + \text{Term})$

angefangene Herleitung

Herleitungen in Baumnotation:

- n **Wurzel:** Startsymbol
- n Anwendung einer Regel $A \rightarrow \alpha$
 - Knotenwachstum.
- n Symbole aus α werden Kindknoten



entsprechender Herleitungsbaum

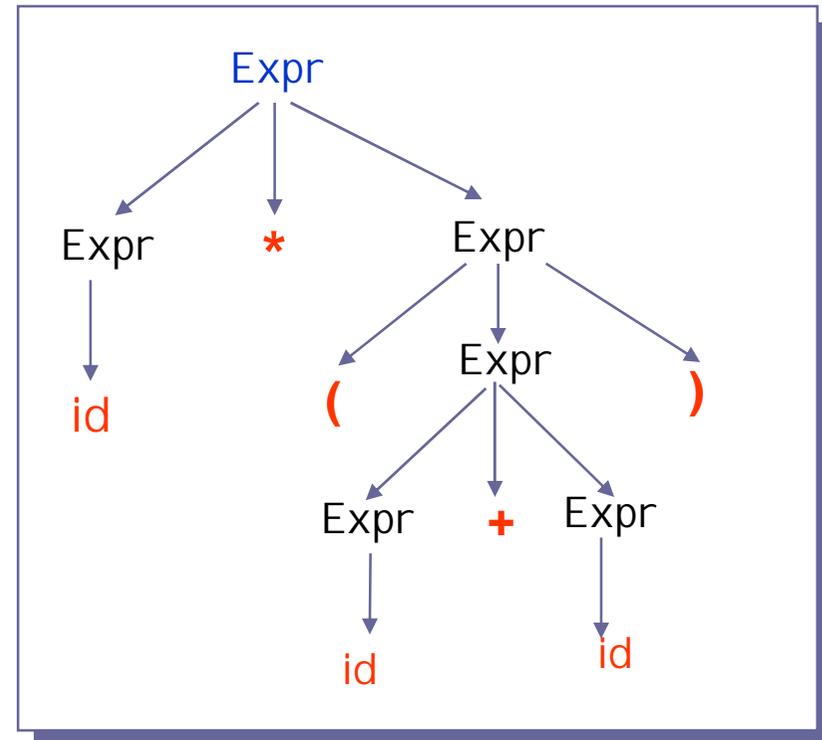


Fertiger Herleitungsbaum

G

```
Expr → Expr + Expr
      | Expr * Expr
      | ( Expr )
      | id
```

Wurzel: Startsymbol
innere Knoten: Variablen
Söhne von A : Symbole von α ,
für $A \rightarrow \alpha \in P$
Blätter: Terminalsymbole



repräsentiert Herleitung
 $\text{Expr} \Rightarrow^* \text{id} * (\text{id} + \text{id})$

Blätter des Baumes ergeben von links nach rechts das hergeleitete Wort.



Grammar-Editor (C. Burch)

The screenshot shows the Grammar Editor interface. On the left, the 'Grammar Editor' window displays the following grammar rules:

```
# Ausdrucks-Grammatik
EXPR -> TERM
      | EXPR plus TERM
TERM  -> FACTOR
      | TERM mal FACTOR
FACTOR -> num
        | id
        | auf EXPR zu
```

At the bottom, the 'Text:' field contains the input string: `auf num plus id zu mal auf id plus num plus id zu`. Below the text field are three buttons: 'Parse', 'Generate', and 'Clear'.

On the right, the 'Parse Tree' window displays a hierarchical tree structure for the input string. The root node is 'EXPR', which branches into 'TERM' and 'mal'. The 'TERM' node further branches into 'FACTOR', 'auf', and 'FACTOR'. The left 'FACTOR' branches into 'auf', 'EXPR', and 'zu'. The 'EXPR' node under this 'FACTOR' branches into 'EXPR', 'plus', and 'TERM'. The left 'EXPR' branches into 'EXPR', 'plus', and 'TERM'. The leftmost 'EXPR' branches into 'FACTOR' and 'id'. The 'FACTOR' node under this 'EXPR' branches into 'num'. The right 'TERM' node branches into 'TERM' and 'FACTOR'. The 'TERM' node under this 'TERM' branches into 'FACTOR' and 'id'. The 'FACTOR' node under this 'TERM' branches into 'num'. The right 'FACTOR' node branches into 'plus' and 'TERM'. The 'TERM' node under this 'FACTOR' branches into 'FACTOR' and 'id'. The 'FACTOR' node under this 'TERM' branches into 'id'.

nettes Spielzeug-Tool
erzeugt Worte einer Grammatik
generiert und zeigt Herleitungsbaum

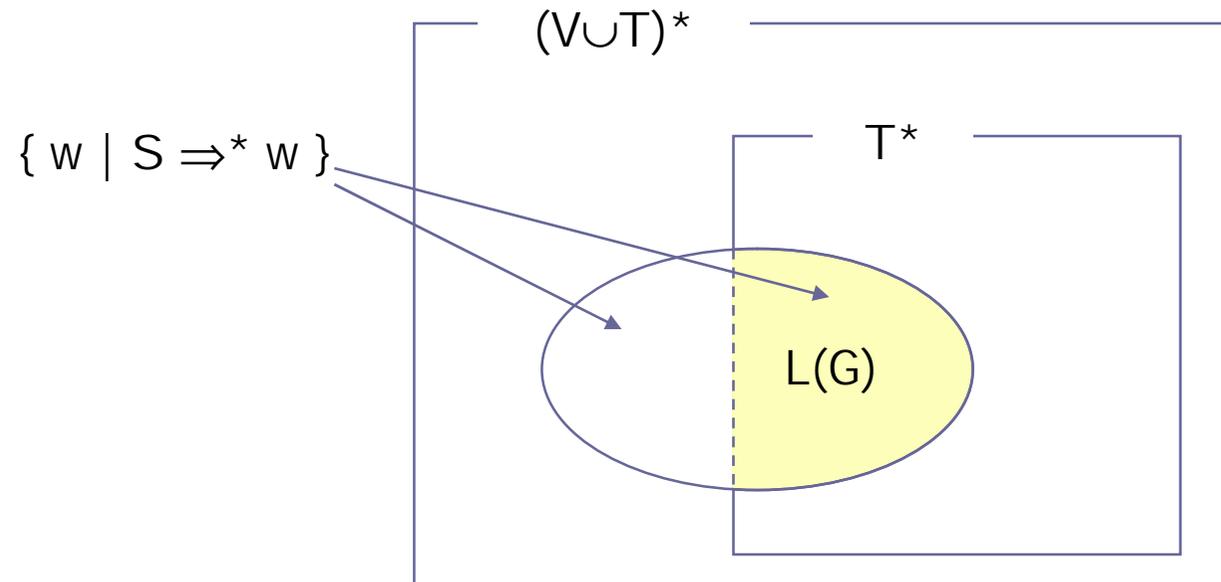


Die Sprache zu einer Grammatik

Ist $G = (V, T, P, S)$ eine Grammatik, so ist

$$L(G) = \{ w \in T^* \mid S \Rightarrow^* w \}$$

die von G erzeugte Sprache.





Kontextfreie Sprachen

n Eine Sprache $L \subseteq T^*$ heißt **kontextfrei**, falls es eine kontextfreie Grammatik G gibt, mit $L = L(G)$.

Beispiel:

$L_{anbn} = \{ a^n b^n \mid n \geq 0 \}$ ist kontextfrei:

$$\begin{array}{l} G \\ \hline S \rightarrow a S b \\ \quad | \varepsilon \end{array}$$

Muss man das beweisen? Kaum, denn es ist ziemlich offensichtlich. Aber wir können daran üben, wie man einen Beweis macht. Also:

1. $L_{anbn} \subseteq L(G)$:

Induktionsanfang $n=0$: klar, denn $a^0 b^0 = \varepsilon$ und $S \Rightarrow^* \varepsilon \in T^*$

Ind. schritt:

Sei $a^k b^k \in L(G)$ schon gezeigt, dann gibt es eine Herleitung $S \Rightarrow a^k b^k$.

... und jetzt ??? Tja - doch nicht so einfach ??



Kontextfreie Sprachen

n Eine Sprache $L \subseteq T^*$ heißt **kontextfrei**, falls es eine kontextfreie Grammatik G gibt, mit $L = L(G)$.

Beispiel:

$L_{anbn} = \{ a^n b^n \mid n \geq 0 \}$ ist kontextfrei:

$$\begin{array}{l} G \\ S \rightarrow a S b \\ \quad | \varepsilon \end{array}$$

$$L_{anbn} = L(G).$$

Die Behauptung folgt aus der stärkeren Aussage:

Für jede Satzform $\alpha \in \{S, a, b\}^*$: $S \Rightarrow^* \alpha$ gdw. $\alpha = a^n S b^n \vee \alpha = a^n b^n$

Beweis: Sei $F = \{a^n S b^n \mid n \geq 0\} \cup \{a^n b^n \mid n \geq 0\}$. Dann gilt

$S \in F$ und

$\sigma \in F$, $\sigma \Rightarrow \tau$ impliziert $\tau \in F$

Daher: $S \Rightarrow^* \alpha$ impliziert $\alpha \in F$.

Umgekehrt: Für jedes n gilt $S \Rightarrow^* a^n S b^n$ (Induktion). Folglich auch $S \Rightarrow^* a^n S b^n \Rightarrow a^n b^n$.



Kontextfreie Sprachen

n Eine Sprache $L \subseteq T^*$ heißt *kontextfrei*, falls es eine kontextfreie Grammatik G gibt, mit $L = L(G)$.

Die Klammersprache ist kontextfrei:

$$\begin{array}{l} \text{K} \\ \hline S \rightarrow (S) \\ \quad | S S \\ \quad | \varepsilon \end{array}$$

Beispiel einer Herleitung:

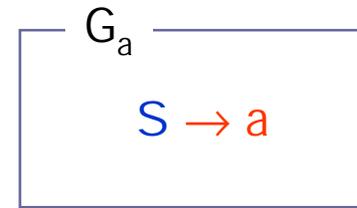
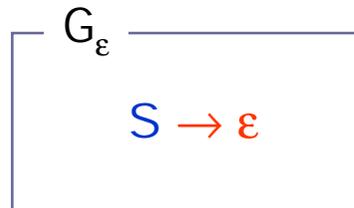
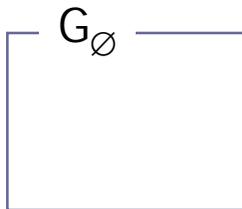
$$\begin{array}{l} S \\ \Rightarrow (S) \\ \Rightarrow (S S) \\ \Rightarrow ((S) S) \\ \Rightarrow (((S)) S) \\ \Rightarrow ((())(S)) \\ \Rightarrow ((())()) \end{array}$$

Hier ist nichts zu beweisen, weil wir die Klammersprache nie definiert haben.
Im Gegenteil, wir nehmen die obige Grammatik als [Definition für die Klammersprache](#).

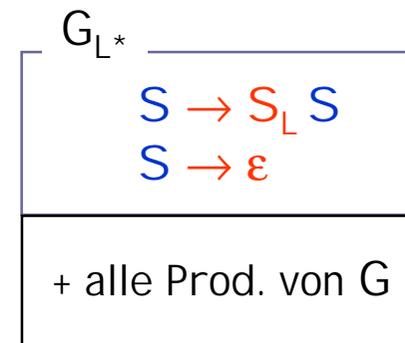
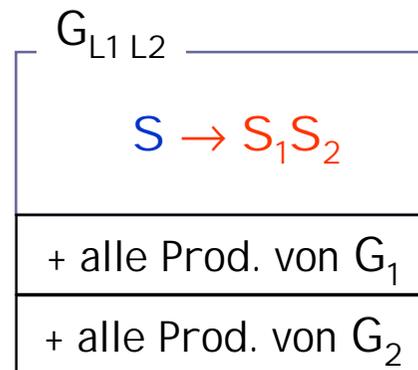
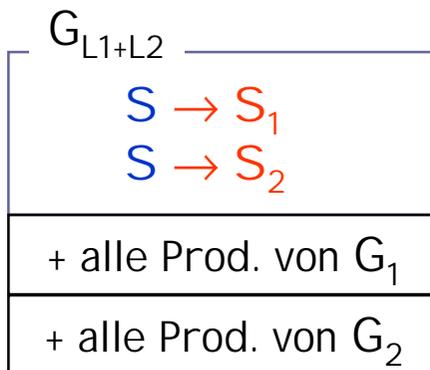


Jede reguläre Sprache ist kontextfrei

Sei Σ gegeben. Die Basissprachen :



Seien L, L_1, L_2 durch Grammatiken G, G_1, G_2 mit Startsymbolen S_L, S_1 , bzw. S_2 gegeben:





Sprachen zu einer Grammatik

Welche Sprache gehört zu folgender Grammatik ?

$$\begin{array}{l} G \\ S \rightarrow aB \mid bA \mid \varepsilon \\ A \rightarrow aS \mid bAA \\ B \rightarrow bS \mid aBB \end{array}$$

Vermutung: $L(G) = \{ w \in \{a,b\}^* \mid \#(a,w) = \#(b,w) \}$,
wobei

$\#(a,w)$ = Anzahl der a's in w.

Seien L_A , L_B , L_S die Sprachen die entstehen, wenn man in G das entsprechende Nonterminal als Startsymbol nimmt. Wir zeigen die stärkere Behauptung:

$$\begin{array}{l} L_S = \{ w \in \{a,b\}^* \mid \#(a,w) = \#(b,w) \} \\ L_A = \{ w \in \{a,b\}^* \mid \#(a,w) = \#(b,w) + 1 \} \\ L_B = \{ w \in \{a,b\}^* \mid \#(a,w) = \#(b,w) - 1 \} \end{array}$$



Simultaner Beweis-Teil 1

G

$$\begin{aligned} S &\rightarrow aB \mid bA \mid \varepsilon \\ A &\rightarrow aS \mid bAA \\ B &\rightarrow bS \mid aBB \end{aligned}$$

$$L_S \subseteq \{ w \in \{a,b\}^* \mid \#(a,w) = \#(b,w) \}$$

$$L_A \subseteq \{ w \in \{a,b\}^* \mid \#(a,w) = \#(b,w) + 1 \}$$

$$L_B \subseteq \{ w \in \{a,b\}^* \mid \#(a,w) = \#(b,w) - 1 \}$$

Behauptung :

$$S \Rightarrow^* w \text{ impliziert } \#(a,w) + \#(A,w) = \#(b,w) + \#(B,w).$$

Induktion über Länge der Herleitung $S \Rightarrow^* w$:

Gilt für $w=S$ (Herleitung der Länge 0).

Jede Regelanwendung erhält die Eigenschaft (einzeln nachprüfen)

exemplarisch: $A \rightarrow bAA$

$$\alpha A \beta \Rightarrow \alpha bAA \beta \quad \begin{array}{l} // \text{ links verschwindet ein } A \\ // \text{ rechts kommen 2 } A\text{'s und ein } b \text{ hinzu.} \end{array}$$

Für $w \in L_S$ folgt also:

$$S \Rightarrow^* w \in \{a,b\}^*, \text{ somit } \#(a,w) = \#(b,w),$$

$$\text{also } L_S \subseteq \{ w \in \{a,b\}^* \mid \#(a,w) = \#(b,w) \}.$$

q.e.d.

Analog zeigt man die Behauptungen für L_A und L_B .



Simultaner Beweis-Teil 2

$$\begin{array}{l}
G \\
S \rightarrow aB \mid bA \mid \varepsilon \\
A \rightarrow aS \mid bAA \\
B \rightarrow bS \mid aBB
\end{array}$$

$$\begin{array}{ll}
L_S \supseteq \{ w \in \{a,b\}^* \mid \#(a,w) = \#(b,w) \} & =: L_0 \\
L_A \supseteq \{ w \in \{a,b\}^* \mid \#(a,w) = \#(b,w) + 1 \} & =: L_{+a} \\
L_B \supseteq \{ w \in \{a,b\}^* \mid \#(a,w) = \#(b,w) - 1 \} & =: L_{+b}
\end{array}$$

Behauptung:

w aus einer der rechten Seiten, dann auch aus der entsprechenden linken.

Induktion über die Länge von w :

$w = \varepsilon$: Wenn w aus einer der rechten Seiten ist, dann muss es aus der ersten sein.

$w \neq \varepsilon$: Beispiel: $w \in L_{+a}$: ($w \in L_0, w \in L_{+b}$ genauso)

1. Fall: $w = a.v$:
 $v \in L_0$ // Def. L_{+a}, L_0
also $v \in L_S$ // Ind. hypothese (v kürzer als w)
also $S \Rightarrow^* v$ // Def L_B
 $A \Rightarrow aS \Rightarrow^* a.v$ // Erste Regel für A
 $A \Rightarrow^* w$ // \Rightarrow transitiv, $w = a.v$
 $w \in L_A$ // Def. L_A

2. Fall: $w = b.v$:
 $\#(a,v) = \#(b,v) + 2$ // klar
 $v = st$ mit $s, t \in L_{+a}$ // v kann so zerlegt werden
also $s, t \in L_A$ // Ind. hypothese (s, t kürzer als w)
also $A \Rightarrow^* s, A \Rightarrow^* t$ // Def L_A
 $A \Rightarrow bAA \Rightarrow^* bsA \Rightarrow^* bst$ // Zweite Regel für A
 $A \Rightarrow^* w$ // \Rightarrow transitiv, $w = a.st$
 $w \in L_A$ // Def. L_A



Äquivalenz

Grammatiken G_1 und G_2 heißen **äquivalent**, falls $L(G_1) = L(G_2)$.

$$\begin{array}{l} G_1 \\ \hline \text{Sum} \rightarrow \text{Sum} + \text{Sum} \\ \quad | \quad \text{id} \end{array}$$

$$\begin{array}{l} G_2 \\ \hline \text{Sum} \rightarrow \text{Sum} + \text{id} \\ \quad | \quad \text{id} \end{array}$$

G_1 und G_2 sind äquivalent:

$$L(G_1) = \text{id} ('+' \text{id})^* = L(G_2)$$



Chomsky-Normalform



<http://www.chomsky.info/>

Zu jeder Grammatik G gibt es eine äquivalente Grammatik G' , deren Regeln folgende einfache Bauart haben:

$$A \rightarrow a$$

$$A \rightarrow BC$$

$$S \rightarrow \varepsilon$$

a Terminal

A, B, C Nichtterminale.

S Startsymbol

Die Produktion $S \rightarrow \varepsilon$ ist genau dann vorhanden, wenn $\varepsilon \in L(G)$ ist.
In diesem Falle kommt S in keiner Regel rechts vor.

Der Satz von Chomsky ist ein theoretisches Ergebnis. In der Praxis ist es angenehmer, mit Regeln beliebiger Bauart spezifizieren zu dürfen.



Erreichbare Variablen

$A \in V$ heißt **erreichbar**, falls $S \Rightarrow^* \alpha A \beta$

Rekursive Charakterisierung

A **erreichbar** $\Leftrightarrow A = S$
oder $(E \rightarrow \alpha A \beta) \in P$ und E **erreichbar**.

Fast schon ein rekursives Programm. Nur für Terminierung sorgen.
Iterative Implementierung *markiert* alle erreichbaren Variablen:

MarkierungsAlgorithmus

Markiere S .

DO

FALLS $(E \rightarrow \alpha A \beta) \in P$ und E schon markiert, **markiere** A

BIS in einer Runde nichts neues markiert wurde

Nicht erreichbare Variablen kann man entfernen.



Produktive Variablen

$A \in V$ heißt **produktiv**, falls $A \Rightarrow^* w \in T^*$

Rekursive Charakterisierung

A **produktiv** $\Leftrightarrow (A \rightarrow w) \in P$ für ein $w \in T^*$
oder $(A \rightarrow \alpha) \in P$ für ein $\alpha \in (V - \{A\} \cup T)^*$
und alle Variablen in α **produktiv**.

MarkierungsAlgorithmus

Markiere alle $A \in V$ mit $(A \rightarrow w) \in P$ und $w \in T^*$.

DO

FALLS $(A \rightarrow \alpha) \in P$ und alle Nonterminale in α schon markiert,
markiere A

BIS in einer Runde nichts neues markiert wurde

Nicht produktive Variablen kann man entfernen.



Kollabierende Variablen

$A \in V$ heißt **kollabierend**, falls $A \Rightarrow^* \varepsilon$

Rekursive Charakterisierung

A **kollabierend** $\Leftrightarrow (A \rightarrow \varepsilon) \in P$
oder $(A \rightarrow B_1 B_2 \dots B_n) \in P$, $B_i \neq A$,
und alle B_i **kollabierend**).

MarkierungsAlgorithmus

Markiere alle $A \in V$ mit $(A \rightarrow \varepsilon) \in P$.

DO

FALLS $(A \rightarrow B_1 \dots B_n) \in P$ und alle B_i schon markiert, **markiere** A
BIS in einer Runde nichts neues markiert wurde



Entfernung von ε -Regeln

n Jede Regel der Form $A \rightarrow \varepsilon$ heißt ε -Regel

Für jede kollabierende Variable C

DO FÜR jede Regel der Form

$$A \rightarrow \alpha C \beta$$

füge eine neue Regel hinzu:

$$A \rightarrow \alpha \beta$$

BIS in einer Runde keine neue Regel hinzukam

Streiche alle ε -Regeln

Vorsicht: Falls $A \rightarrow \alpha C \beta C \gamma$ und C kollabierend, müssen sowohl $A \rightarrow \alpha \beta C \gamma$ als auch $A \rightarrow \alpha C \beta \gamma$ und $A \rightarrow \alpha \beta \gamma$ hinzukommen.

Aus Grammatik G entsteht eine neue Grammatik G^+ ohne ε -Regeln und es gilt

$$L(G^+) = L(G) - \{\varepsilon\}.$$



ϵ -Freiheit

G heißt ϵ -frei, falls

n keine Regel der Form $A \rightarrow \epsilon$ vorkommt,

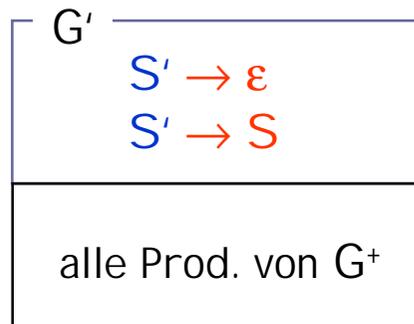
n oder ($S \rightarrow \epsilon$ ist einzige ϵ -Regel,

und S kommt in keiner Produktion rechts vor).

$L(G^+) = L(G) - \{\epsilon\}$ und G^+ ist ϵ -frei. Daher:

1. Fall: $\epsilon \notin L(G)$. Setze $G' = G^+$.

2. Fall: $\epsilon \in L(G)$:



Zu jeder Grammatik G gibt es eine äquivalente ϵ -freie Grammatik.



Variablen-Regeln

Alle Variablen-Regeln, das sind Regeln der Form $A \rightarrow B$ mit $A, B \in V$ kann man eliminieren

Algorithmus

DO

FÜR jede Variablenregel $A \rightarrow B$:

FÜR jede Regel $B \rightarrow \alpha$

füge neue Regel $A \rightarrow \alpha$ hinzu

Streiche $A \rightarrow B$

BIS nichts mehr passiert.



Chomsky - Normalform

n Satz: Jede kontextfreie Grammatik G lässt sich äquivalent umformen in eine Grammatik G' , deren Regeln von der Form sind:

$$\begin{aligned} A &\rightarrow a \\ A &\rightarrow BC \end{aligned}$$

Falls $\varepsilon \in L(G)$ gibt es noch die Produktion $S \rightarrow \varepsilon$ und S kommt in keiner Produktion rechts vor.

1. Entferne nicht terminierende und nicht erreichbare Variablen
2. Mache G ε -frei und entferne alle Variablen-Regeln.
3. In jeder Regel $A \rightarrow \alpha$ mit $|\alpha| > 1$ ersetze Terminalsymbol a durch neue Variable X_a und füge die Regel $X_a \rightarrow a$ hinzu.

4. Ersetze jede Regel der Form

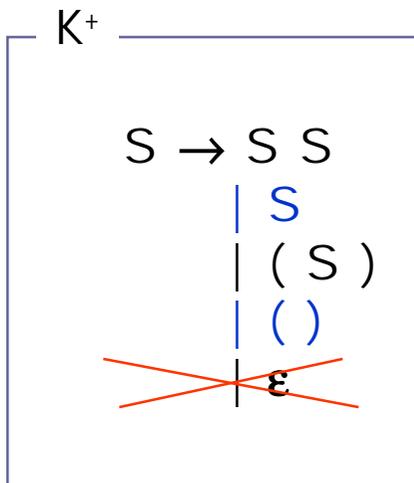
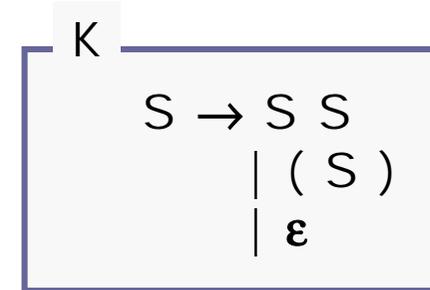
$A \rightarrow B_1 B_2 \dots B_n$ mit Hilfe neuer Variablen $C_1 \dots C_{n-1}$ durch

$$A \rightarrow B_1 C_1, \quad C_1 \rightarrow B_2 C_2 \quad \dots \quad C_{n-1} \rightarrow B_{n-1} B_n$$

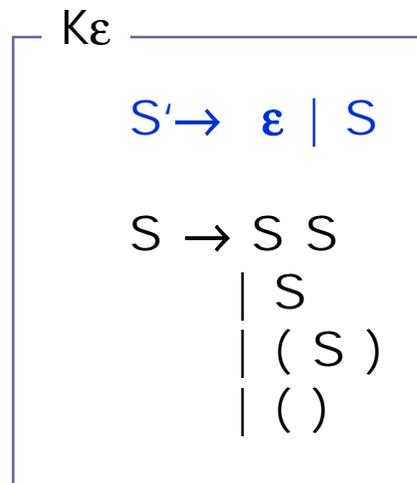


Beispiel: Gewinnung einer CNF – Anfang

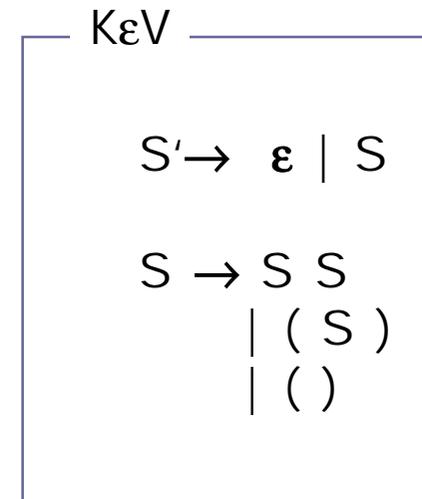
- n Wir wandeln die Klammersprache K in CNF um
- n Jede Variable ist erreichbar und terminierend
 - .. Das ist normal für realistische Grammatiken !
- n S ist kollabierend. Wir erhalten nacheinander
 - .. K^+ durch Entfernen der ϵ -Regeln
 - .. K_ϵ als zu K äquivalente ϵ -freie Grammatik
 - .. $K_{\epsilon V}$ durch Elimination von Variablenregeln



$$L(K^+) \neq L(K)$$



$$L(K_\epsilon) = L(K)$$



$$L(K_{\epsilon V}) = L(K)$$



Beispiel: Gewinnung einer CNF – Schluss

n Wir wandeln die Kammersprache K in CNF um

- Wir haben bereits $K_{\epsilon V}$ zu K äquivalent
 - n alle Variablen terminierend und ϵ -frei
 - n ohne Variablenregeln
- Wir ersetzen nun Terminale durch Variablen
- und verkürzen rechte Seiten
- K_{CNF} ist äquivalent zu K und in Chomsky-Normalform

$$\begin{array}{l}
 K \\
 \hline
 S \rightarrow S S \\
 \quad \quad | (S) \\
 \quad \quad | \epsilon
 \end{array}$$

$$\begin{array}{l}
 K_{\epsilon V} \\
 \hline
 S' \rightarrow \epsilon \\
 \quad \quad | S \\
 \\
 S \rightarrow S S \\
 \quad \quad | (S) \\
 \quad \quad | ()
 \end{array}$$

$$\begin{array}{l}
 K_{\epsilon VT} \\
 \hline
 S' \rightarrow \epsilon \mid S \\
 \\
 S \rightarrow S S \\
 \quad \quad | A S Z \\
 \quad \quad | A Z \\
 \\
 A \rightarrow (\\
 Z \rightarrow)
 \end{array}$$

$$\begin{array}{l}
 K_{CNF1} \\
 \hline
 S' \rightarrow \epsilon \mid S \\
 \\
 S \rightarrow S S \\
 \quad \quad | A C \\
 \quad \quad | A Z \\
 \\
 C \rightarrow S Z \\
 \\
 A \rightarrow (\\
 Z \rightarrow)
 \end{array}$$



Die Normalform

n Eine Variablenregel ist noch zu entfernen:

$$\begin{array}{l} K \\ S \rightarrow S S \\ \quad | (S) \\ \quad | \epsilon \end{array}$$

Ausgangs
grammatik

$$\begin{array}{l} K_{\text{CNF1}} \\ S' \rightarrow \epsilon \mid S \\ S \rightarrow S S \\ \quad | A C \\ \quad | A Z \\ C \rightarrow S Z \\ A \rightarrow (\\ Z \rightarrow) \end{array}$$

$$\begin{array}{l} K_{\text{CNF}} \\ S' \rightarrow \epsilon \\ \quad | S S \\ \quad | A C \\ \quad | A Z \\ S \rightarrow S S \\ \quad | A C \\ \quad | A Z \\ C \rightarrow S Z \\ A \rightarrow (\\ Z \rightarrow) \end{array}$$



Wortproblem für kontextfreier Sprachen

Sei G eine Grammatik, $A \in V$ und $w = a_1 a_2 \dots a_n \in T^*$.

Gilt

$$A \Rightarrow^* w \quad ?$$

Wir haben einen Algorithmus, um G in Chomsky-Normalform G_{CNF} zu überführen. Dann wenden wir den folgenden Algorithmus von Cocke, Younger und Kasami an:

CYK-Algorithmus - rekursiv

$$\begin{aligned} \underline{A \Rightarrow^* w} \quad \Leftrightarrow \quad & w = \varepsilon \text{ und } S \rightarrow \varepsilon \in G_{\text{CNF}} \\ & \text{oder } n=1 \text{ und } S \rightarrow w \in G_{\text{CNF}} \\ & \text{oder } n > 1 \text{ und es gibt ein } 1 \leq k < n \text{ mit} \\ & \quad A \rightarrow BC \in G_{\text{CNF}} \text{ und } B \Rightarrow^* a_1 \dots a_k \\ & \quad \text{und } C \Rightarrow^* a_{k+1} \dots a_n \end{aligned}$$

Jede kontextfreie Sprache ist entscheidbar.



CYK-Algorithmus - imperativ

Imperativ implementiert man den CYK durch **dynamisches Programmieren**:

Gegeben $w = a_1 \dots a_n \in T^*$. Für jedes $1 \leq i, k \leq n$ sei

$$M_{ik} = \{ A \in V \mid A \Rightarrow^* a_i \dots a_{i+k-1} \}$$

Dann gilt:

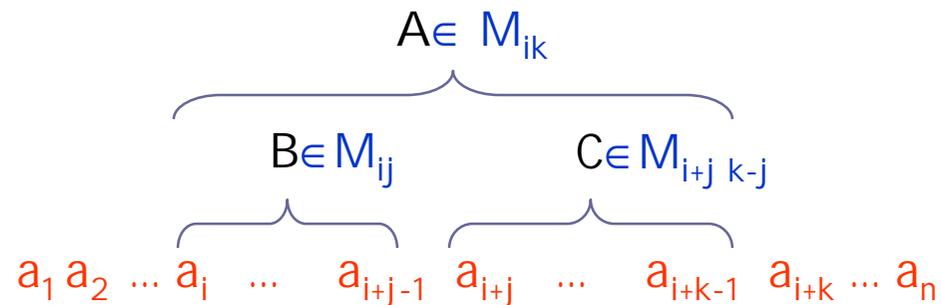
$$w \in L(G) \Leftrightarrow S \in M_{1n}$$

Ist G in Chomsky-Normalform, so kann man die M_{ik} in eine Tabelle eintragen.

$$M_{i1} = \{ A \in V \mid A \rightarrow a_i \in G \}$$

$$M_{ik} = \{ A \in V \mid \exists_{1 \leq j < k} : \exists A \rightarrow BC \in G : B \in M_{ij}, C \in M_{i+j, k-j} \}$$

5	M_{15}				
4	M_{14}	M_{24}			
3	M_{13}	M_{23}	M_{33}		
2	M_{12}	M_{22}	M_{32}	M_{42}	
1	M_{11}	M_{21}	M_{31}	M_{41}	M_{51}
M_{ik}	1	2	3	4	5





CYK- Algorithmus – ein Beispiel

- n Gegeben sei die Sprache der nichtleeren Klammerausdrücke
- n Gefragt: Ist $((())()) \in L(G)$?

8								
7								
6								
5								
4								
3								
2								
1	A	A	A	Z	Z	A	Z	Z
M_{ij}	((())	())

K_{CNF}^+
$S \rightarrow S S$
A C
A Z
$C \rightarrow S Z$
$A \rightarrow ($
$Z \rightarrow)$



CYK- Algorithmus – ein Beispiel

- n Gegeben sei die Sprache der nichtleeren Klammerausdrücke
- n Gefragt: Ist $((())()) \in L(G)$?

8									
7									
6									
5									
4									
3									
2	-	-	S	-	-	S	-		
1	A	A	A	Z	Z	A	Z	Z	
M_{ij}	((())	())	

K_{CNF}^+	
S	\rightarrow S S
	A C
	A Z
C	\rightarrow S Z
A	\rightarrow (
Z	\rightarrow)



CYK- Algorithmus – ein Beispiel

- n Gegeben sei die Sprache der nichtleeren Klammerausdrücke
- n Gefragt: Ist $((())()) \in L(G)$?

8									
7									
6									
5									
4									
3	-	-	C	-	-	C			
2	-	-	S	-	-	S	-		
1	A	A	A	Z	Z	A	Z	Z	
M_{ij}	((())	())	

K_{CNF}^+	
S	$\rightarrow S S$
	A C
	A Z
C	$\rightarrow S Z$
A	$\rightarrow ($
Z	$\rightarrow)$



CYK- Algorithmus – ein Beispiel

- n Gegeben sei die Sprache der nichtleeren Klammerausdrücke
- n Gefragt: Ist $((())()) \in L(G)$?

8									
7									
6									
5									
4	-	S	-	-	-				
3	-	-	C	-	-	C			
2	-	-	S	-	-	S	-		
1	A	A	A	Z	Z	A	Z	Z	
M_{ij}	((())	())	

K_{CNF}^+	
S	\rightarrow S S
	A C
	A Z
C	\rightarrow S Z
A	\rightarrow (
Z	\rightarrow)



CYK- Algorithmus – ein Beispiel

- n Gegeben sei die Sprache der nichtleeren Klammerausdrücke
- n Gefragt: Ist $((())()) \in L(G)$?

8								
7								
6								
5	-	-	-	-				
4	-	S	-	-	-			
3	-	-	C	-	-	C		
2	-	-	S	-	-	S	-	
1	A	A	A	Z	Z	A	Z	Z
M_{ij}	((())	())

K_{CNF}^+
$S \rightarrow S S$
A C
A Z
$C \rightarrow S Z$
$A \rightarrow ($
$Z \rightarrow)$



CYK- Algorithmus – ein Beispiel

- n Gegeben sei die Sprache der nichtleeren Klammerausdrücke
- n Gefragt: Ist $((())()) \in L(G)$?

8								
7								
6	-	S	-					
5	-	-	-	-				
4	-	S	-	-	-			
3	-	-	C	-	-	C		
2	-	-	S	-	-	S	-	
1	A	A	A	Z	Z	A	Z	Z
M_{ij}	((())	())

K_{CNF}^+
$S \rightarrow S S$
A C
A Z
$C \rightarrow S Z$
$A \rightarrow ($
$Z \rightarrow)$



CYK- Algorithmus – ein Beispiel

- n Gegeben sei die Sprache der nichtleeren Klammerausdrücke
- n Gefragt: Ist $((())()) \in L(G)$?

8								
7	-	C						
6	-	S	-					
5	-	-	-	-				
4	-	S	-	-	-			
3	-	-	C	-	-	C		
2	-	-	S	-	-	S	-	
1	A	A	A	Z	Z	A	Z	Z
M_{ij}	((())	())

K_{CNF}^+	
S	$\rightarrow S S$
	A C
	A Z
C	$\rightarrow S Z$
A	$\rightarrow ($
Z	$\rightarrow)$



CYK- Algorithmus – ein Beispiel

- n Gegeben sei die Sprache der nichtleeren Klammerausdrücke
- n Gefragt: Ist $((())()) \in L(G)$?

8	S								
7	-	C							
6	-	S	-						
5	-	-	-	-					
4	-	S	-	-	-				
3	-	-	C	-	-	C			
2	-	-	S	-	-	S	-		
1	A	A	A	Z	Z	A	Z	Z	
M_{ij}	((())	())	

$S \in M_{18} \iff w \in L(K_{\text{CNF}}^+)$

K_{CNF}^+	
S	\rightarrow S S
	A C
	A Z
C	\rightarrow S Z
A	\rightarrow (
Z	\rightarrow)



Grenzen von CF-Sprachen

n Gibt es Sprachen, die nicht kontextfrei sind ?

.. ja, z.B. $\{a^n b^n c^n \mid n \geq 0\}$

n Wie kann man das beweisen ?

.. Pumping Lemma für CF Sprachen – kommt gleich

n Sind CF-Sprachen unter \cap abgeschlossen

.. Nein, denn

$$\{a^n b^n c^n \mid n \geq 0\} = \{a^n b^n c^k \mid n \geq 0, k \geq 0\} \cap \{a^k b^n c^n \mid n \geq 0, k \geq 0\}$$

und letztere beiden sind offensichtlich CF.

n Unter Komplement ?

.. Nein, denn sie sind unter \cup abgeschlossen, und mit Komplement und Vereinigung kann man den Schnitt ausdrücken:

$$L \cap M = \Sigma^* - ((\Sigma^* - L) \cup (\Sigma^* - M))$$



Pumping Lemma für CF-Sprachen

Für jede CF-Sprache L gibt es eine Zahl k , so daß jedes Wort $w \in L$ mit $|w| \geq k$ sich zerlegen läßt als

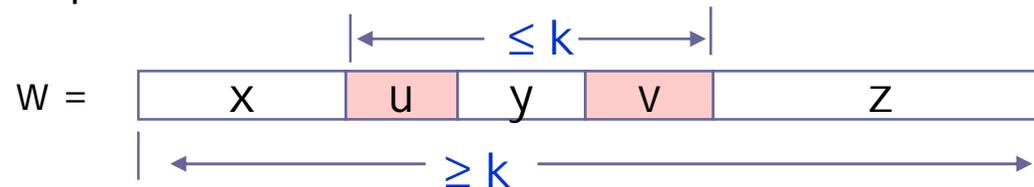
$$w = x u y v z$$

so dass

$$0 < |uv| \leq |uyv| \leq k$$

$$\forall n \in \text{Nat} : x u^n y v^n z \in L .$$

Also jedes große ($\geq k$) Wort hat ein kleines ($\leq k$) Teilwort uyv , das sich ringförmig aufpumpen läßt:



$$\boxed{x \quad y \quad z} \in L$$

$$\boxed{x \quad u \quad u \quad y \quad v \quad v \quad z} \in L$$

$$\boxed{x \quad u \quad u \quad u \quad y \quad v \quad v \quad v \quad z} \in L$$

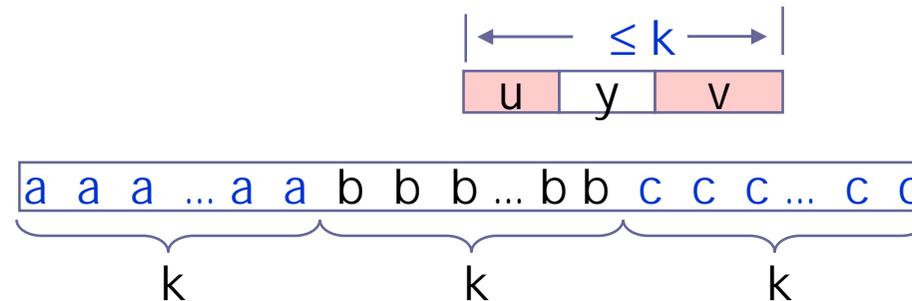


Anwendung des CF-Pumping Lemma

$L = \{a^n b^n c^n \mid n \geq 0\}$ ist nicht CF

- **Angenommen** L sei CF, dann gäbe es ein k wie im PL.
- Betrachte $w = a^k b^k c^k \in L$.
- Gemäß PL lässt sich dies zerlegen:

$$a^k b^k c^k = x u y v z \text{ mit } |u y v| \leq k$$

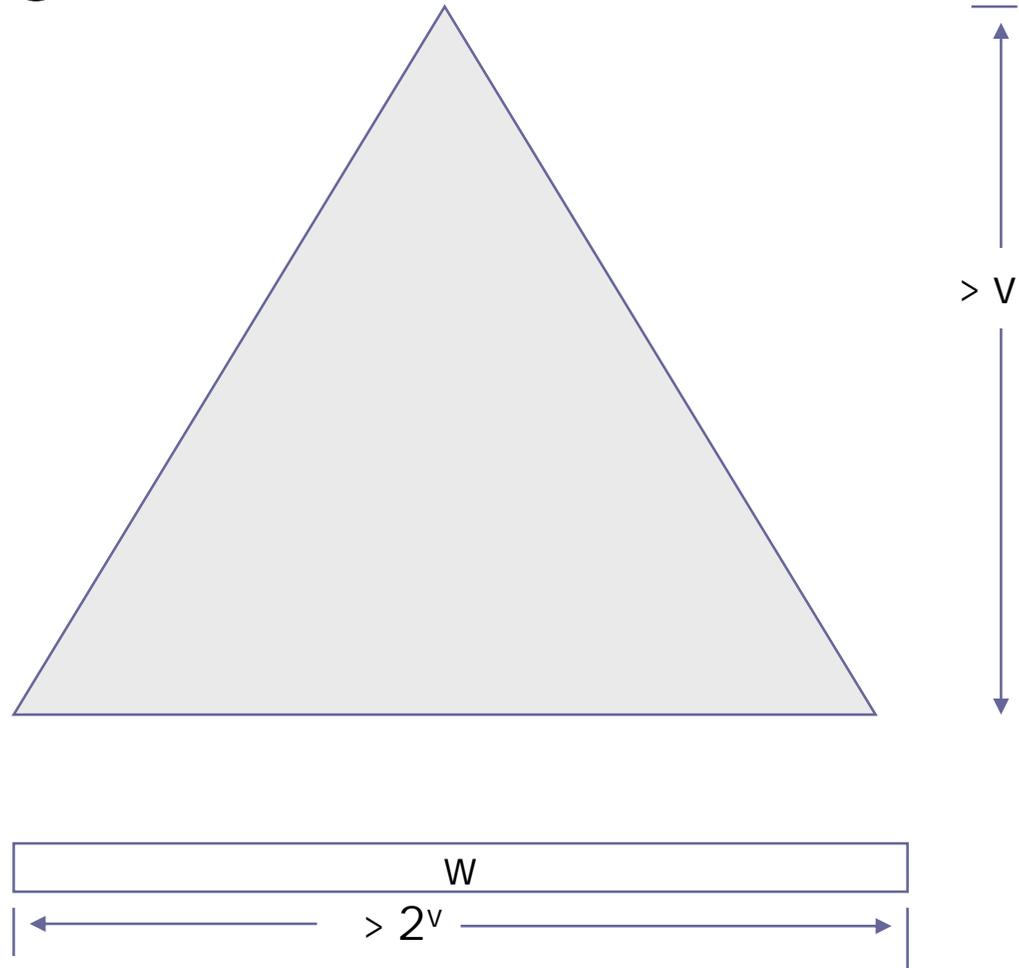


In $u y v$ kommt entweder kein a oder kein c vor. OBdA kein a .
Wegen $|u y v| > 0$ hat dann $x y z$ **zu viele** a 's um in L zu sein.
Widerspruch zum PL.



Beweis des Pumping Lemma

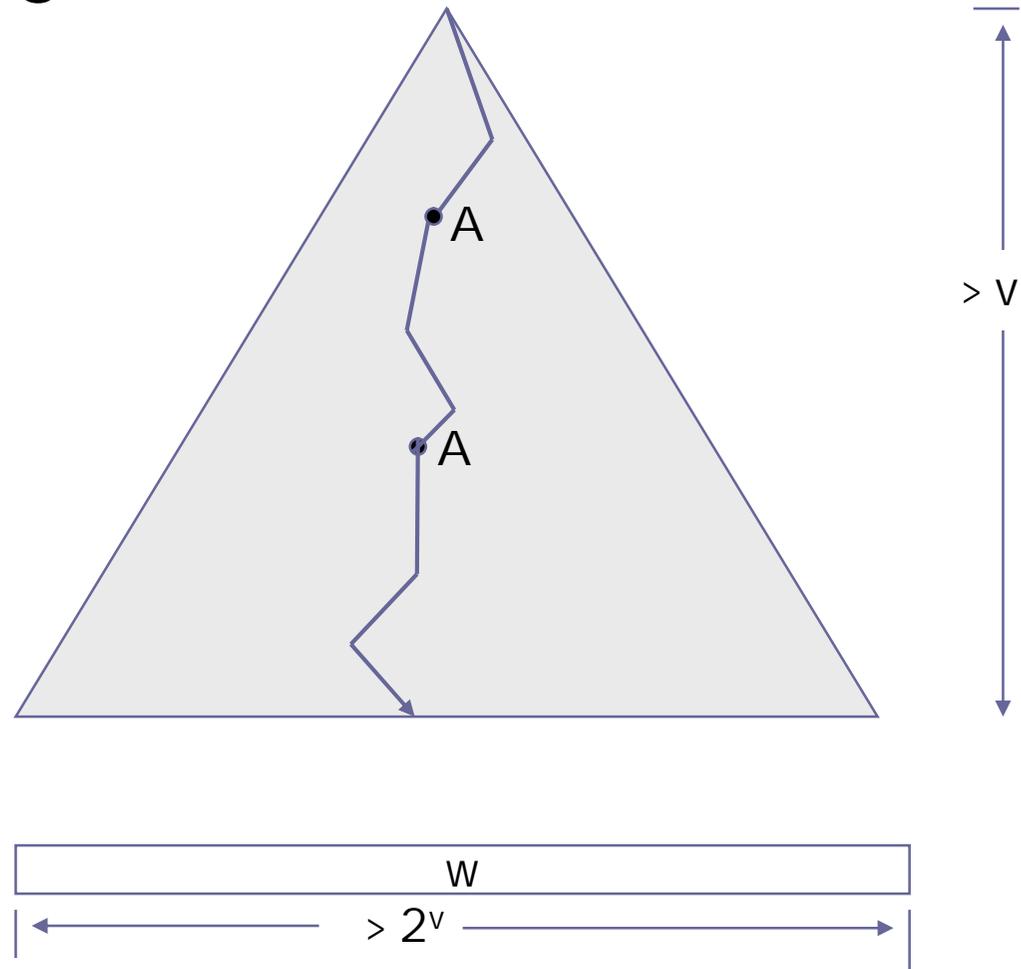
- n Sei G eine Grammatik in **Chomsky-Normalform** für L mit $v-1$ Variablen.
- n Setze $k=2^v$.
- n Ist $w \in L(G)$ mit $|w| \geq k$, dann hat der Herleitungsbaum für w Tiefe $> v$.
- n





Beweis des Pumping Lemma

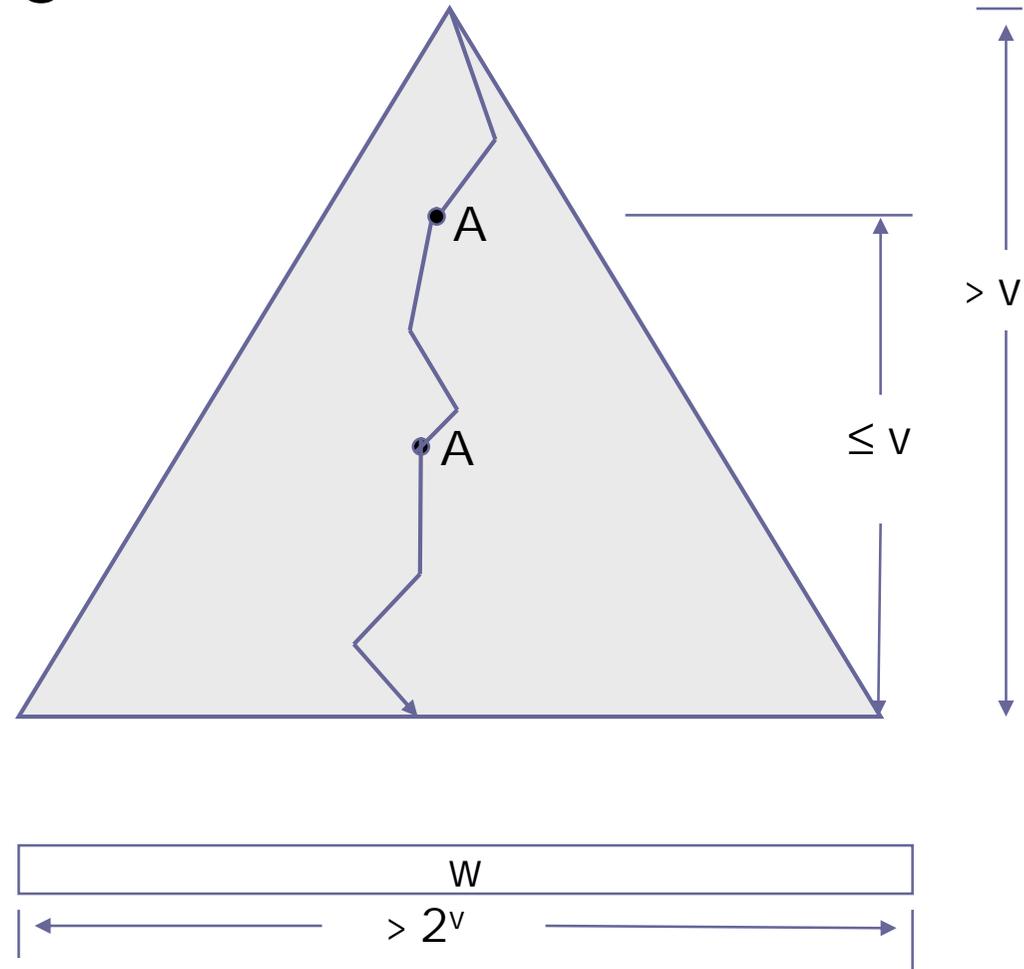
- n Sei G eine Grammatik in Chomsky-Normalform für L mit $v-1$ Variablen.
- n Setze $k=2^v$.
- n Ist $w \in L(G)$ mit $|w| \geq k$, dann hat der Herleitungsbaum für w Tiefe $> v$.
- n Es gibt also einen Pfad zur Wurzel, auf dem ein Nonterminal A zweimal vorkommt.





Beweis des Pumping Lemma

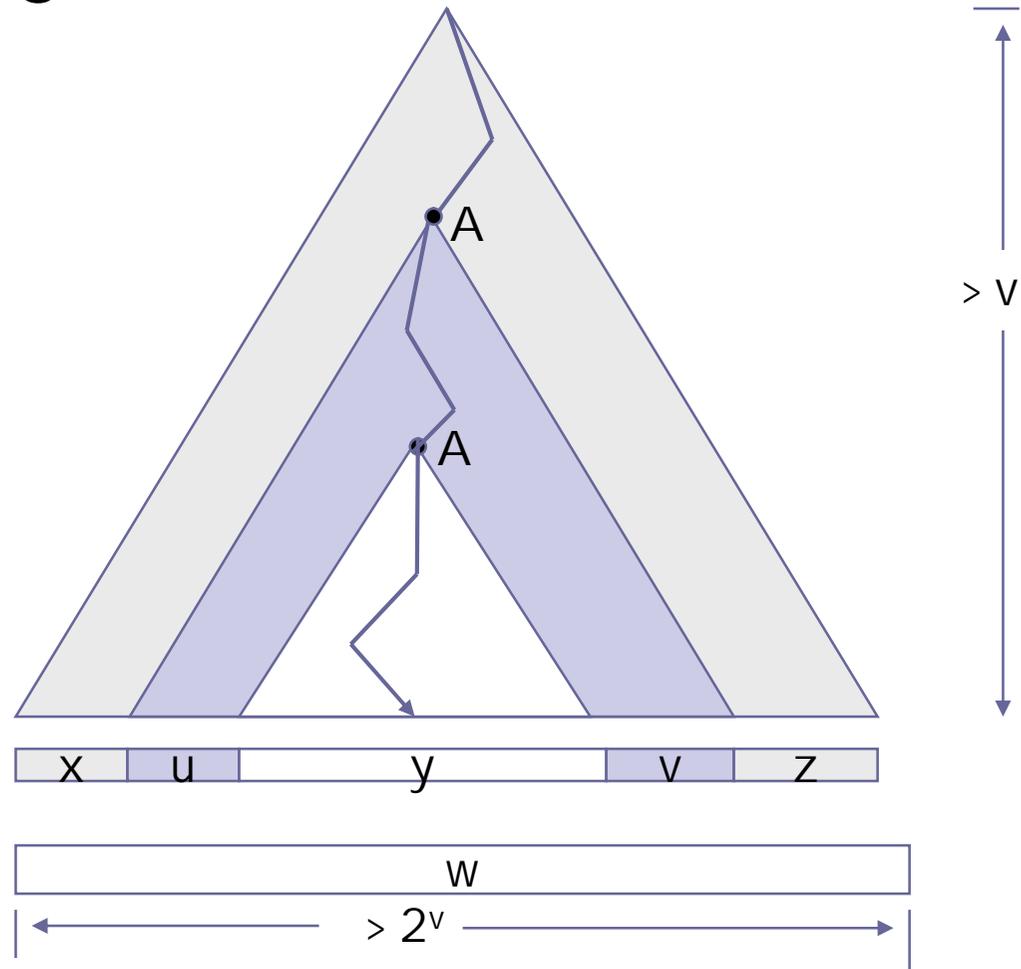
- n Sei G eine Grammatik in Chomsky-Normalform für L mit $v-1$ Variablen.
- n Setze $k=2^v$.
- n Ist $w \in L(G)$ mit $|w| \geq k$, dann hat der Herleitungsbaum für w Tiefe $> v$.
- n Es gibt also einen Pfad zur Wurzel, auf dem ein Nonterminal A zweimal vorkommt. Man kann die beiden A sogar von Tiefe $\leq v$ wählen.





Beweis des Pumping Lemma

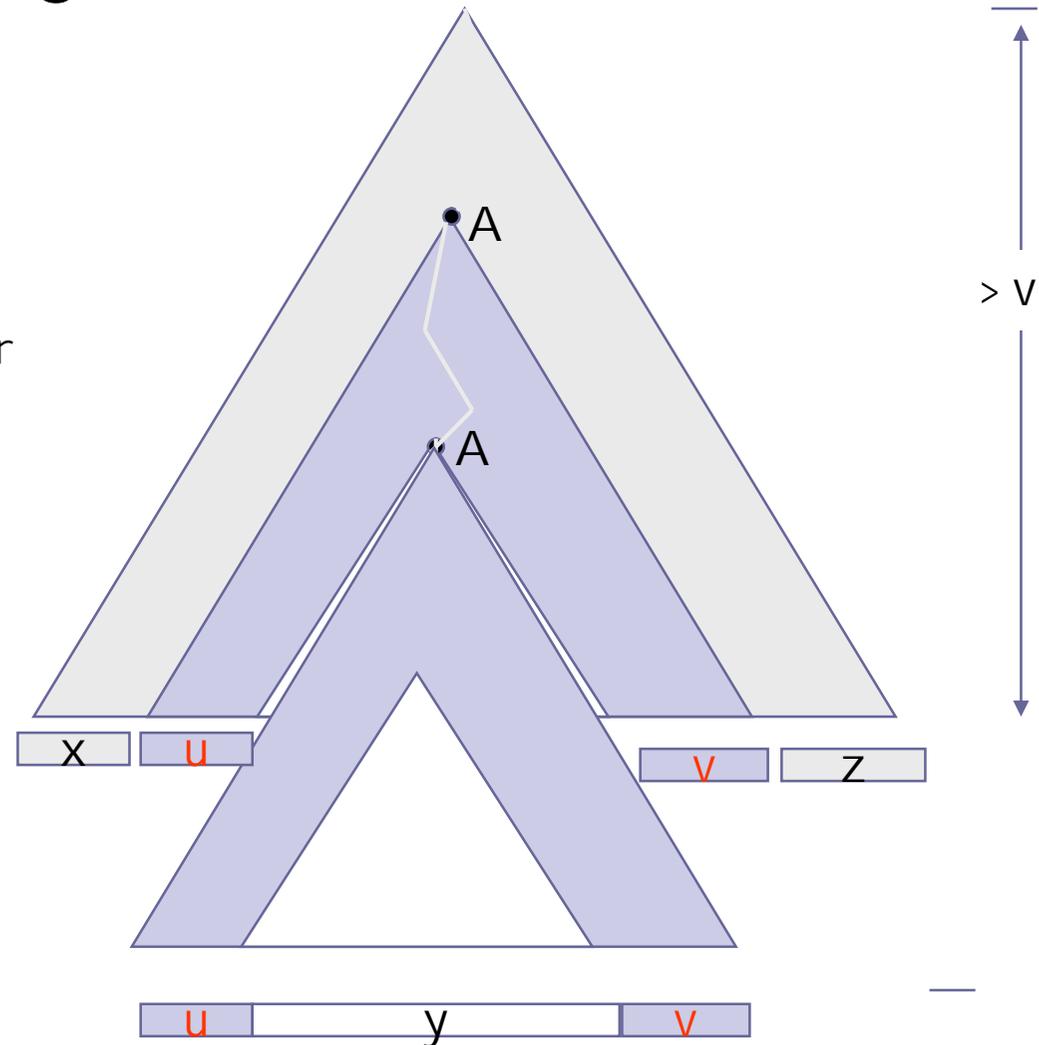
- n Sei G eine Grammatik in Chomsky-Normalform für L mit $v-1$ Variablen.
- n Setze $k=2^v$.
- n Ist $w \in L(G)$ mit $|w| \geq k$, dann hat der Herleitungsbaum für w Tiefe $> v$.
- n Es gibt also einen Pfad zur Wurzel, auf dem ein Nonterminal A zweimal vorkommt. Man kann die beiden A sogar von Tiefe $\leq v$ wählen.
- n w zerlegt sich also als $w = x u y v z$ wie in der Figur





Beweis des Pumping Lemma

- n Sei G eine Grammatik in Chomsky-Normalform für L mit $v-1$ Variablen.
- n Setze $k=2^v$.
- n Ist $w \in L(G)$ mit $|w| \geq k$, dann hat der Herleitungsbaum für w Tiefe $> v$.
- n Es gibt also einen Pfad zur Wurzel, auf dem ein Nonterminal A zweimal vorkommt. Man kann die beiden A sogar von Tiefe $\leq v$ wählen.
- n w zerlegt sich also als $w = x u y v z$ wie in der Figur
- n Wir könnten den Baum unter dem ersten Vorkommen von A auch nochmal verwenden
- n und nochmal ... $xuuuyvvvz \in L$

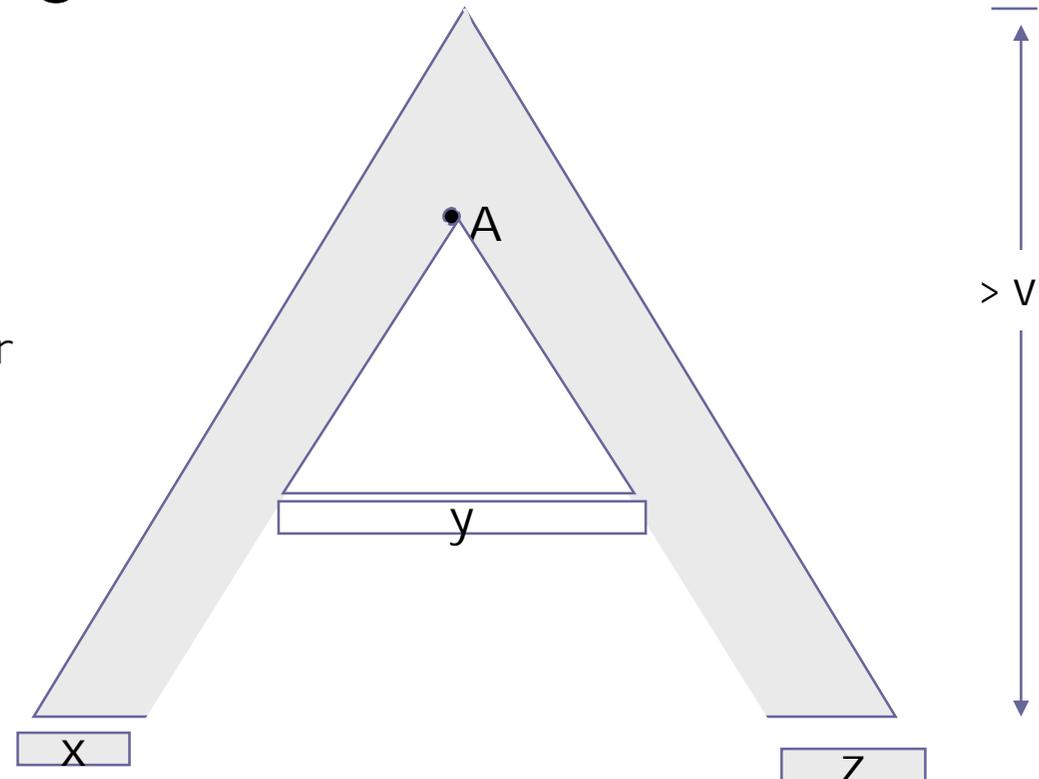


$$w = x u u y v v z \in L(G)$$



Beweis des Pumping Lemma

- n Sei G eine Grammatik in Chomsky-Normalform für L mit $v-1$ Variablen.
- n Setze $k=2^v$.
- n Ist $w \in L(G)$ mit $|w| \geq k$, dann hat der Herleitungsbaum für w Tiefe $> v$.
- n Es gibt also einen Pfad zur Wurzel, auf dem ein Nonterminal A zweimal vorkommt. Man kann die beiden A sogar von Tiefe $\leq v$ wählen.
- n w zerlegt sich also als $w = x u y v z$ wie in der Figur
- n Wir könnten den Baum unter dem ersten Vorkommen von A auch nochmal verwenden
- n und nochmal ... $xuuuyvvvz \in L$
- n oder auch gar nicht $xyz \in L$



$$w = x y z \in L(G)$$



Abschluss unter Operationen



- n Sind L_1, L_2, L kontextfrei, dann auch
 - $L_1 + L_2$ (Vereinigung)
 - $L_1 L_2$ (Verkettung)
 - L^* (Kleene-Stern)

- n aber: kontextfreie Sprachen sind **nicht** abgeschlossen unter
 - $\Sigma^* - L$ (Komplement)
 - $L_1 \cap L_2$ (Schnitt)

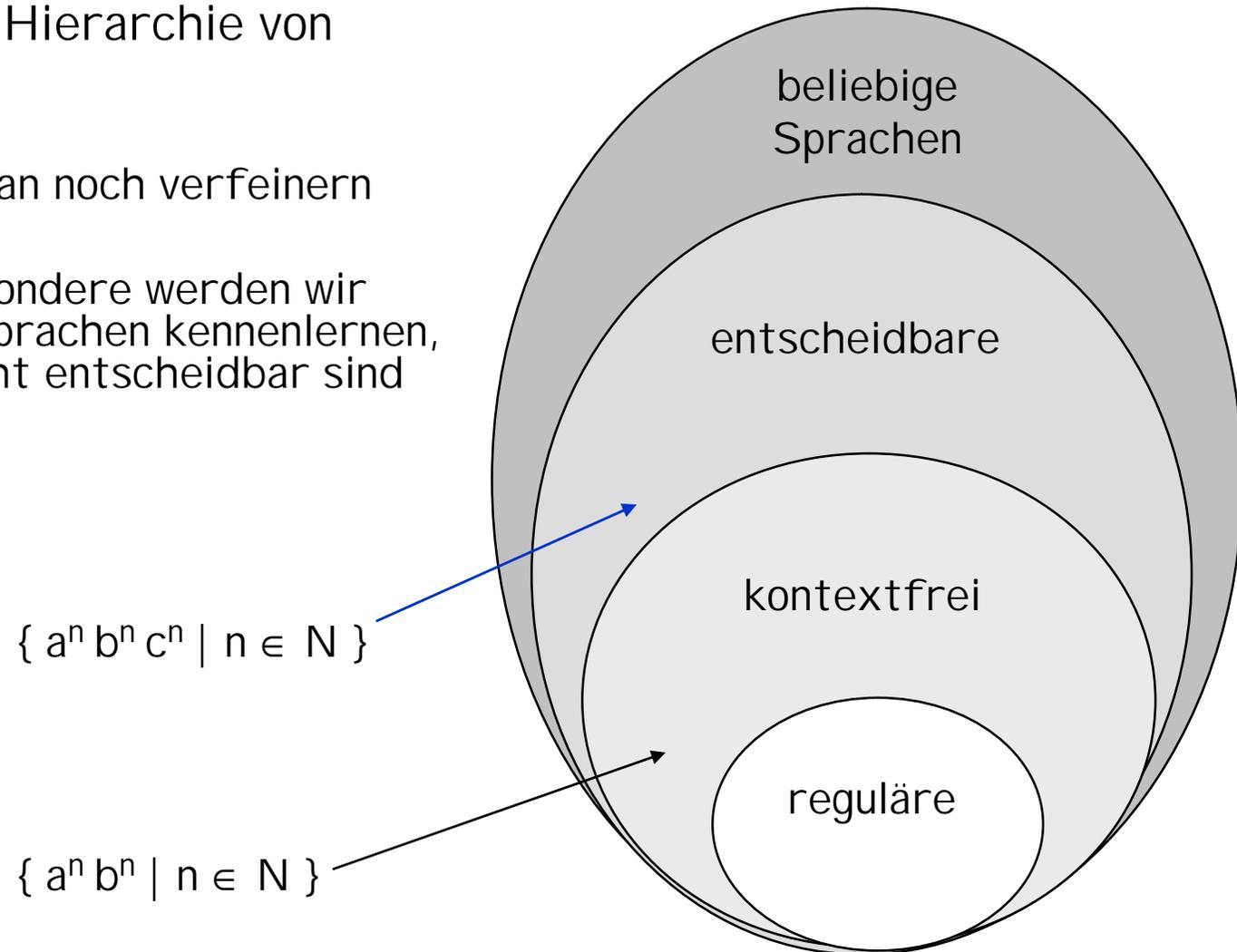
- n Gegenbeispiel :
 - $L_1 = \{ a^k b^m c^n \mid k=m \}$ **Kontextfrei**
 - $L_2 = \{ a^k b^m c^n \mid m=n \}$ **Kontextfrei**
 - $L_1 \cap L_2 = \{ a^k b^m c^n \mid k=m=n \}$ **nicht** Kontextfrei



Hierarchie von Sprachen

n Bisherige Hierarchie von Sprachen

- .. Kann man noch verfeinern
- .. Insbesondere werden wir noch Sprachen kennenlernen, die nicht entscheidbar sind





Inhalt

1. Grammatiken und Sprachen

- .. Kontextfreie Grammatiken
- .. Herleitungen, Linksherleitungen
- .. Sprachen zu einer Grammatik
- .. Äquivalenz
- .. Chomsky-Normalform
- .. Wortproblem, CYK
- .. Pumping Lemma für CF-Sprachen

2. Stackautomaten

- .. Definitionen und Beispiele
- .. Konfigurationen, Läufe,
- .. Sprache eines Stackautomaten
- .. Parser

3. Parser

- .. Mehrdeutigkeit
- .. Bottom Up, Top Down
- .. Recursive descent Parser
- .. First, Follow
- .. Semantische Aktionen

4. Shift-Reduce Parser

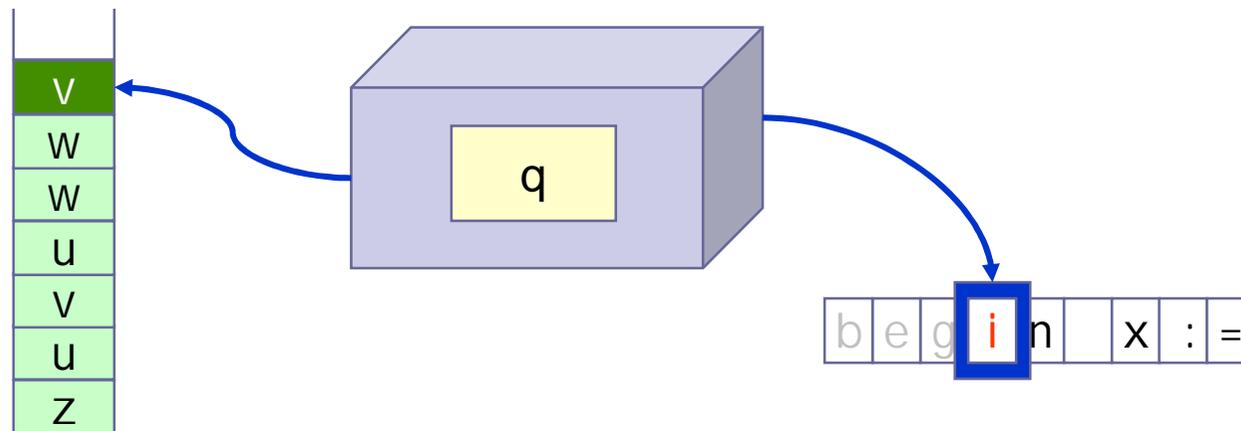
- .. Konflikte
- .. LR(0)-Zustände
- .. Automat zur Grammatik
- .. Shift-Reduce, Reduce-Reduce Konflikte
- .. Präzedenz und Assoziativität
- .. lex und yacc
- .. Arbeitsweise eines Compilers



Stackautomat (Kellerautomat)

- n DFA mit Stack. Dadurch kann sich der Automat beliebig viele Informationen "merken".

In jedem Schritt kann er neben seinem Zustand auch noch das oberste Element des Stacks zu Rate ziehen und endlich viele neue Elemente auf dem Stack ablegen.



Im deutschen Sprachbereich findet man häufig das Wort "Keller" statt des Englischen "Stack".



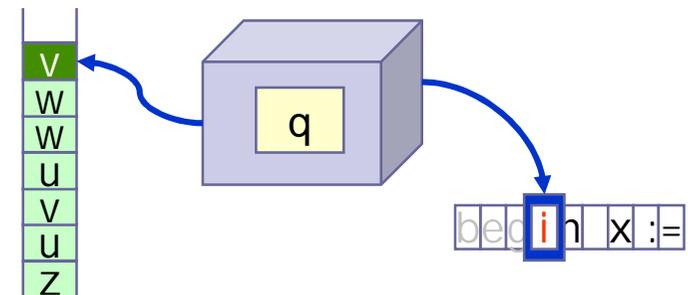
Nichtdeterministischer Stackautomat

Ein nichtdeterministischer Stackautomat (PDA) ist ein 6-Tupel

$$M = (Q, \Sigma, \Gamma, \delta, q_0, s_0)$$

bestehend aus

- n Q - endliche Menge (Zustände)
- n Σ - endliches Alphabet (Inputalphabet)
- n Γ - endliches Alphabet (Stackalphabet)
- n $\delta : Q \times (\Sigma \cup \{ \varepsilon \}) \times \Gamma \rightarrow \wp_{\text{finite}}(Q \times \Gamma^*)$ (Transitionsfunktion)
- n $q_0 \in Q$ (Startzustand)
- n $s_0 \in \Gamma$ (Stack-Startsymbol)





Nichtdeterministischer Stackautomat

Ein nichtdeterministischer Stackautomat (PDA) ist ein 6-Tupel

$$M = (Q, \Sigma, \Gamma, \delta, q_0, s_0)$$

bestehend aus

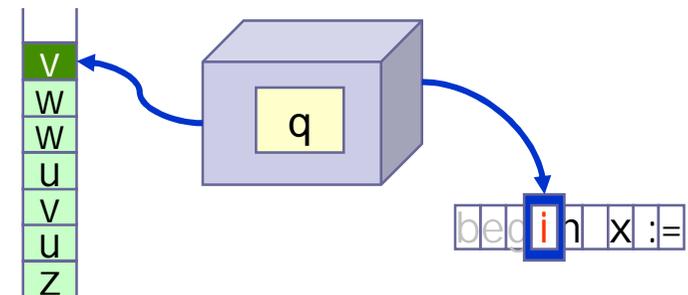
- n Q - endliche Menge (Zustände)
- n Σ - endliches Alphabet (Inputalphabet)
- n Γ - endliches Alphabet (Stackalphabet)
- n $\delta : Q \times (\Sigma \cup \{ \varepsilon \}) \times \Gamma \rightarrow \wp_{\text{finite}}(Q \times \Gamma^*)$ (Transitionsfunktion)
- n $q_0 \in Q$ (Startzustand)
- n $s_0 \in \Gamma$ (Stack-Startsymbol)

Interpretation von $(q', g_1 \dots g_k) \in \delta(q, a, g)$:

Falls

Zustand = q
Input = a
top(Stack) = g ,

dann





Nichtdeterministischer Stackautomat

Ein nichtdeterministischer Stackautomat (PDA) ist ein 6-Tupel

$$M = (Q, \Sigma, \Gamma, \delta, q_0, s_0)$$

bestehend aus

- n Q - endliche Menge (Zustände)
- n Σ - endliches Alphabet (Inputalphabet)
- n Γ - endliches Alphabet (Stackalphabet)
- n $\delta : Q \times (\Sigma \cup \{ \varepsilon \}) \times \Gamma \rightarrow \wp_{\text{finite}}(Q \times \Gamma^*)$ (Transitionsfunktion)
- n $q_0 \in Q$ (Startzustand)
- n $s_0 \in \Gamma$ (Stack-Startsymbol)

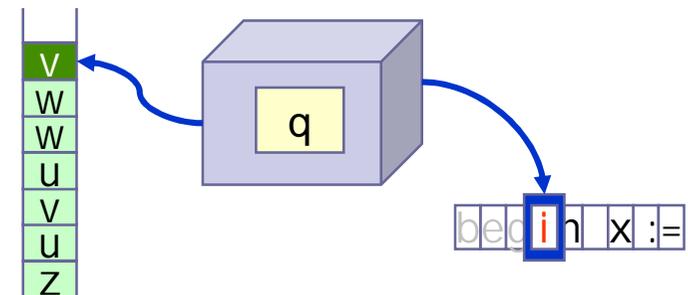
Interpretation von $(q', g_1 \dots g_k) \in \delta(q, a, g)$:

Falls

Zustand = q
Input = a
top(Stack) = g ,

dann

neuer Zustand = q'
pop, push $g_k, \dots, \text{push } g_1$





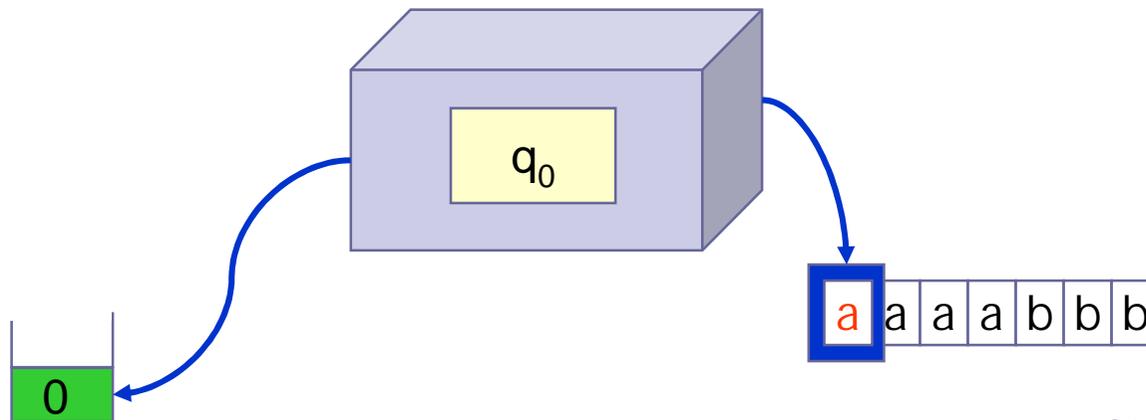
Beispiel einer Stackmaschine

$Q = \{ q_0, q_1, q_2 \}$
 $\Sigma = \{ a, b \}$
 $\Gamma = \{ 0, 1 \}$
 $s_0 = 0$

M

$\delta(q_0, a, 0) = \{ (q_1, 0), (q_1, 10) \}$
 $\delta(q_1, a, 0) = \{ (q_1, 10) \}$
 $\delta(q_1, a, 1) = \{ (q_1, 11) \}$
 $\delta(q_1, b, 1) = \{ (q_2, \epsilon) \}$
 $\delta(q_2, b, 1) = \{ (q_2, \epsilon) \}$
 $\delta(q_2, \epsilon, 0) = \{ (q_2, \epsilon) \}$
 $\delta(q_i, u, k) = \{ \}$ sonst

Stackmaschine beginnt immer mit s auf dem ansonsten leeren Stack.
Im Beispiel kann das Wort **aaaabbb** den Stack entleeren :





Beispiel einer Stackmaschine

$$Q = \{ q_0, q_1, q_2 \}$$

$$\Sigma = \{ a, b \}$$

$$\Gamma = \{ 0, 1 \}$$

$$s_0 = 0$$

M

$$\delta(q_0, a, 0) = \{ (q_1, 0), (q_1, 10) \}$$

$$\delta(q_1, a, 0) = \{ (q_1, 10) \}$$

$$\delta(q_1, a, 1) = \{ (q_1, 11) \}$$

$$\delta(q_1, b, 1) = \{ (q_2, \epsilon) \}$$

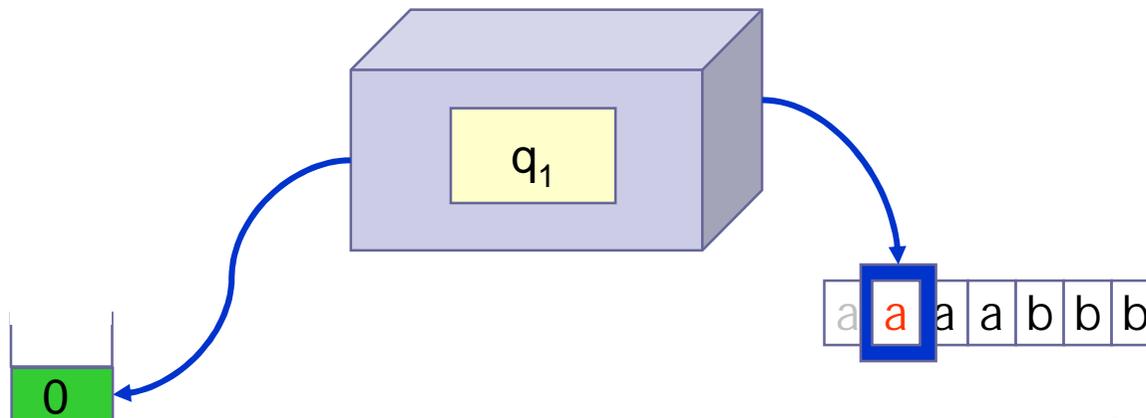
$$\delta(q_2, b, 1) = \{ (q_2, \epsilon) \}$$

$$\delta(q_2, \epsilon, 0) = \{ (q_2, \epsilon) \}$$

$$\delta(q_i, u, k) = \{ \} \quad \text{sonst}$$

Stackmaschine beginnt immer mit s auf dem ansonsten leeren Stack.

Im Beispiel kann das Wort **aaaabbb** den Stack entleeren :





Beispiel einer Stackmaschine

$$Q = \{ q_0, q_1, q_2 \}$$

$$\Sigma = \{ a, b \}$$

$$\Gamma = \{ 0, 1 \}$$

$$s_0 = 0$$

M

$$\delta(q_0, a, 0) = \{ (q_1, 0), (q_1, 10) \}$$

$$\delta(q_1, a, 0) = \{ (q_1, 10) \}$$

$$\delta(q_1, a, 1) = \{ (q_1, 11) \}$$

$$\delta(q_1, b, 1) = \{ (q_2, \epsilon) \}$$

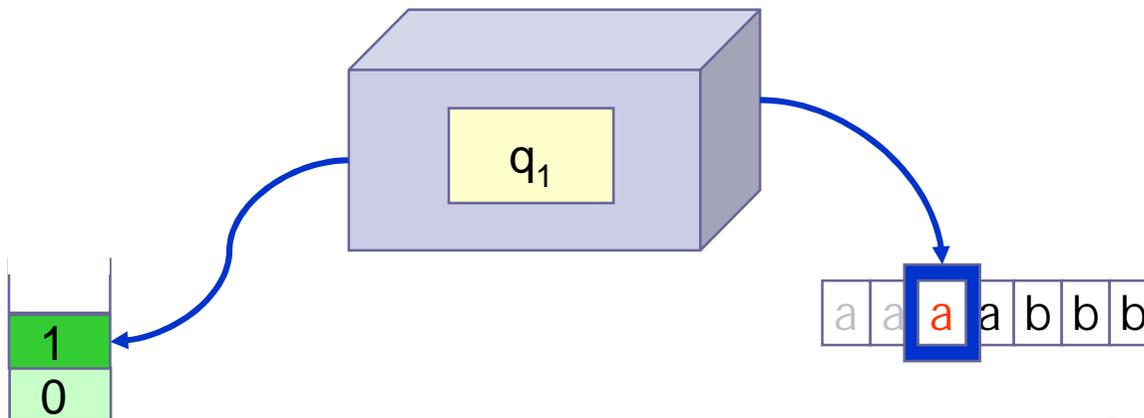
$$\delta(q_2, b, 1) = \{ (q_2, \epsilon) \}$$

$$\delta(q_2, \epsilon, 0) = \{ (q_2, \epsilon) \}$$

$$\delta(q_i, u, k) = \{ \} \quad \text{sonst}$$

Stackmaschine beginnt immer mit s auf dem ansonsten leeren Stack.

Im Beispiel kann das Wort **aaaabbb** den Stack entleeren :





Beispiel einer Stackmaschine

$$Q = \{ q_0, q_1, q_2 \}$$

$$\Sigma = \{ a, b \}$$

$$\Gamma = \{ 0, 1 \}$$

$$s_0 = 0$$

M

$$\delta(q_0, a, 0) = \{ (q_1, 0), (q_1, 10) \}$$

$$\delta(q_1, a, 0) = \{ (q_1, 10) \}$$

$$\delta(q_1, a, 1) = \{ (q_1, 11) \}$$

$$\delta(q_1, b, 1) = \{ (q_2, \epsilon) \}$$

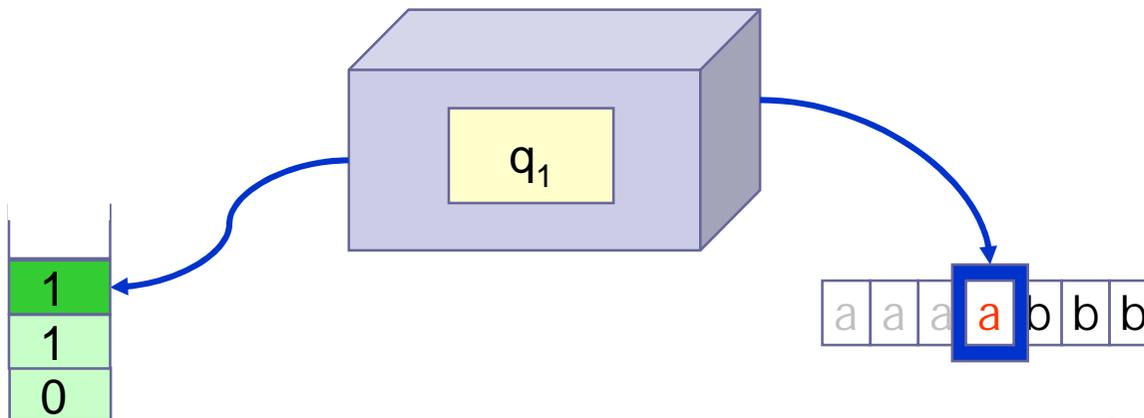
$$\delta(q_2, b, 1) = \{ (q_2, \epsilon) \}$$

$$\delta(q_2, \epsilon, 0) = \{ (q_2, \epsilon) \}$$

$$\delta(q_i, u, k) = \{ \} \quad \text{sonst}$$

Stackmaschine beginnt immer mit s auf dem ansonsten leeren Stack.

Im Beispiel kann das Wort **aaaabbb** den Stack entleeren :





Beispiel einer Stackmaschine

$$Q = \{ q_0, q_1, q_2 \}$$

$$\Sigma = \{ a, b \}$$

$$\Gamma = \{ 0, 1 \}$$

$$s_0 = 0$$

M

$$\delta(q_0, a, 0) = \{ (q_1, 0), (q_1, 10) \}$$

$$\delta(q_1, a, 0) = \{ (q_1, 10) \}$$

$$\delta(q_1, a, 1) = \{ (q_1, 11) \}$$

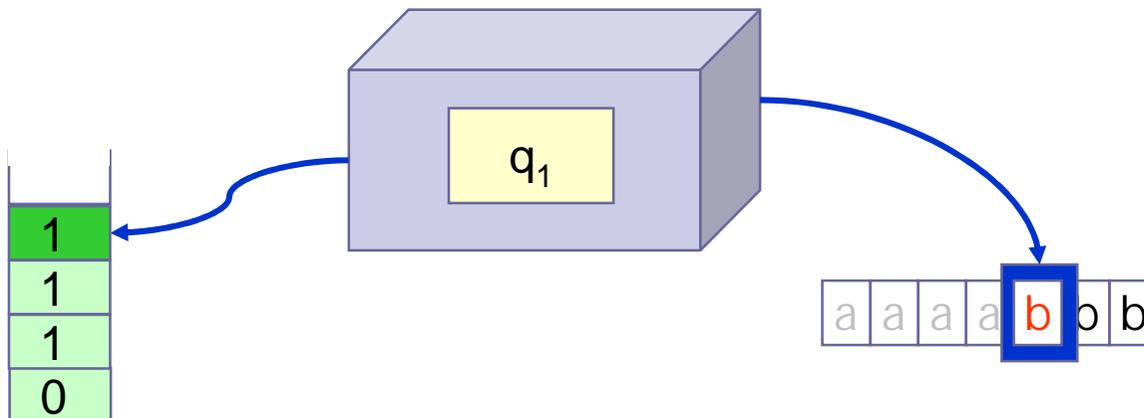
$$\delta(q_1, b, 1) = \{ (q_2, \varepsilon) \}$$

$$\delta(q_2, b, 1) = \{ (q_2, \varepsilon) \}$$

$$\delta(q_2, \varepsilon, 0) = \{ (q_2, \varepsilon) \}$$

$$\delta(q_i, u, k) = \{ \} \quad \text{sonst}$$

Stackmaschine beginnt immer mit s auf dem ansonsten leeren Stack.
Im Beispiel kann das Wort **aaaabbb** den Stack entleeren :





Beispiel einer Stackmaschine

$$Q = \{ q_0, q_1, q_2 \}$$

$$\Sigma = \{ a, b \}$$

$$\Gamma = \{ 0, 1 \}$$

$$s_0 = 0$$

M

$$\delta(q_0, a, 0) = \{ (q_1, 0), (q_1, 10) \}$$

$$\delta(q_1, a, 0) = \{ (q_1, 10) \}$$

$$\delta(q_1, a, 1) = \{ (q_1, 11) \}$$

$$\delta(q_1, b, 1) = \{ (q_2, \epsilon) \}$$

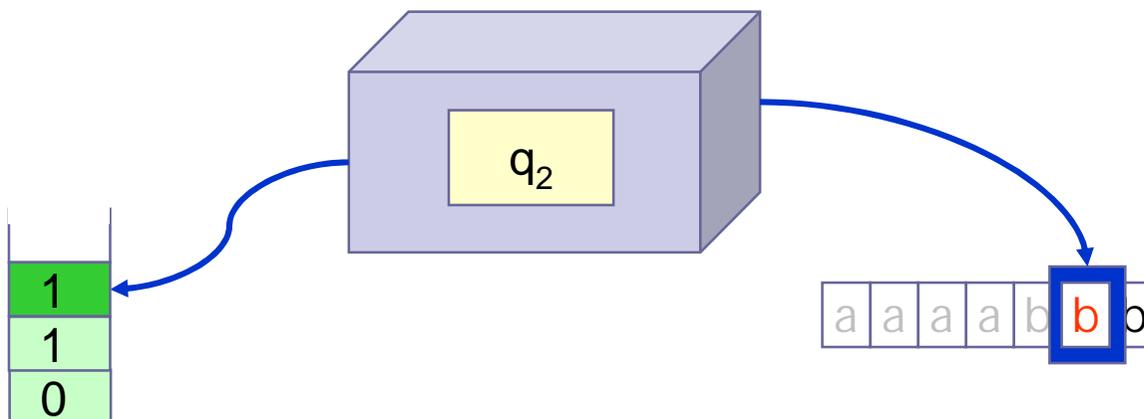
$$\delta(q_2, b, 1) = \{ (q_2, \epsilon) \}$$

$$\delta(q_2, \epsilon, 0) = \{ (q_2, \epsilon) \}$$

$$\delta(q_i, u, k) = \{ \} \quad \text{sonst}$$

Stackmaschine beginnt immer mit s auf dem ansonsten leeren Stack.

Im Beispiel kann das Wort **aaaabbb** den Stack entleeren :





Beispiel einer Stackmaschine

$$Q = \{ q_0, q_1, q_2 \}$$

$$\Sigma = \{ a, b \}$$

$$\Gamma = \{ 0, 1 \}$$

$$s_0 = 0$$

M

$$\delta(q_0, a, 0) = \{ (q_1, 0), (q_1, 10) \}$$

$$\delta(q_1, a, 0) = \{ (q_1, 10) \}$$

$$\delta(q_1, a, 1) = \{ (q_1, 11) \}$$

$$\delta(q_1, b, 1) = \{ (q_2, \epsilon) \}$$

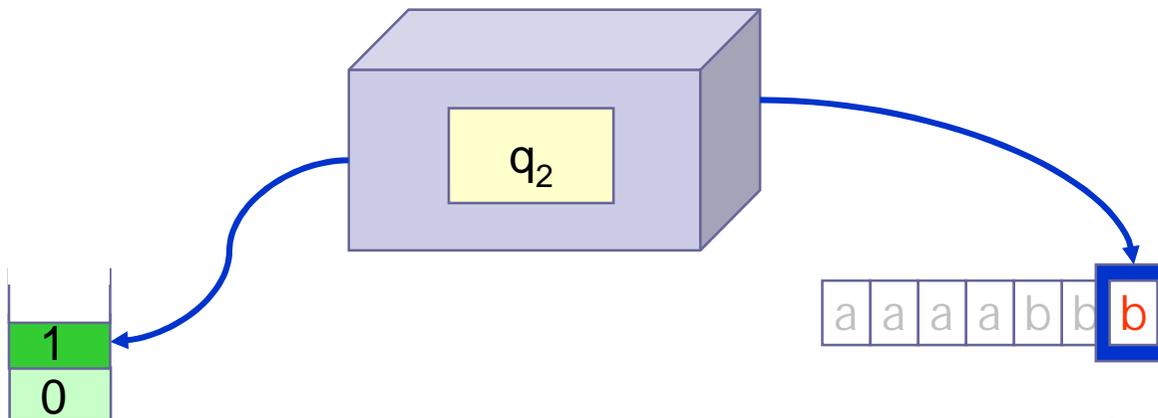
$$\delta(q_2, b, 1) = \{ (q_2, \epsilon) \}$$

$$\delta(q_2, \epsilon, 0) = \{ (q_2, \epsilon) \}$$

$$\delta(q_i, u, k) = \{ \} \quad \text{sonst}$$

Stackmaschine beginnt immer mit s auf dem ansonsten leeren Stack.

Im Beispiel kann das Wort **aaaabbb** den Stack entleeren :





Beispiel einer Stackmaschine

$$Q = \{ q_0, q_1, q_2 \}$$

$$\Sigma = \{ a, b \}$$

$$\Gamma = \{ 0, 1 \}$$

$$s_0 = 0$$

M

$$\delta(q_0, a, 0) = \{ (q_1, 0), (q_1, 10) \}$$

$$\delta(q_1, a, 0) = \{ (q_1, 10) \}$$

$$\delta(q_1, a, 1) = \{ (q_1, 11) \}$$

$$\delta(q_1, b, 1) = \{ (q_2, \epsilon) \}$$

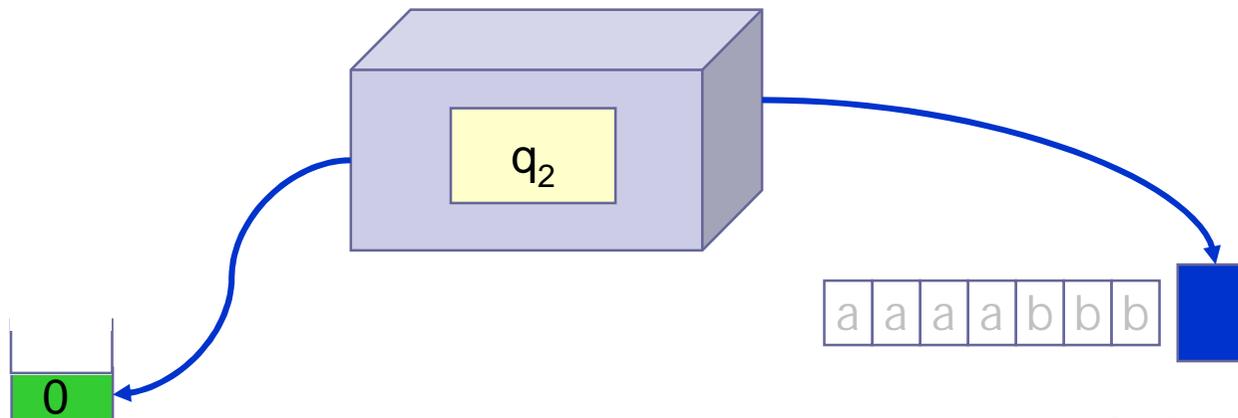
$$\delta(q_2, b, 1) = \{ (q_2, \epsilon) \}$$

$$\delta(q_2, \epsilon, 0) = \{ (q_2, \epsilon) \}$$

$$\delta(q_i, u, k) = \{ \} \quad \text{sonst}$$

Stackmaschine beginnt immer mit s auf dem ansonsten leeren Stack.

Im Beispiel kann das Wort **aaaabbb** den Stack entleeren :





Beispiel einer Stackmaschine

$$Q = \{ q_0, q_1, q_2 \}$$

$$\Sigma = \{ a, b \}$$

$$\Gamma = \{ 0, 1 \}$$

$$s_0 = 0$$

M

$$\delta(q_0, a, 0) = \{ (q_1, 0), (q_1, 10) \}$$

$$\delta(q_1, a, 0) = \{ (q_1, 10) \}$$

$$\delta(q_1, a, 1) = \{ (q_1, 11) \}$$

$$\delta(q_1, b, 1) = \{ (q_2, \varepsilon) \}$$

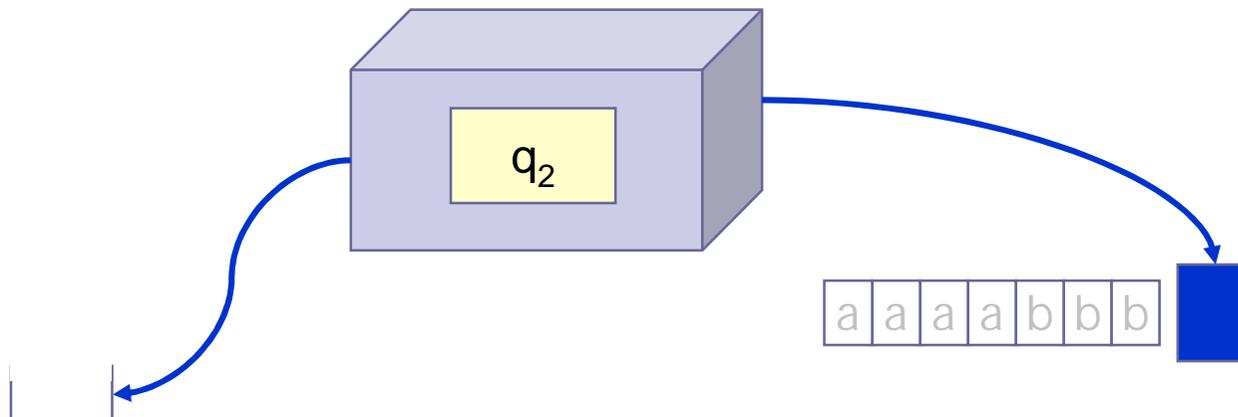
$$\delta(q_2, b, 1) = \{ (q_2, \varepsilon) \}$$

$$\delta(q_2, \varepsilon, 0) = \{ (q_2, \varepsilon) \}$$

$$\delta(q_i, u, k) = \{ \} \quad \text{sonst}$$

Stackmaschine beginnt immer mit s auf dem ansonsten leeren Stack.

Im Beispiel kann das Wort **aaaabbb** den Stack entleeren :

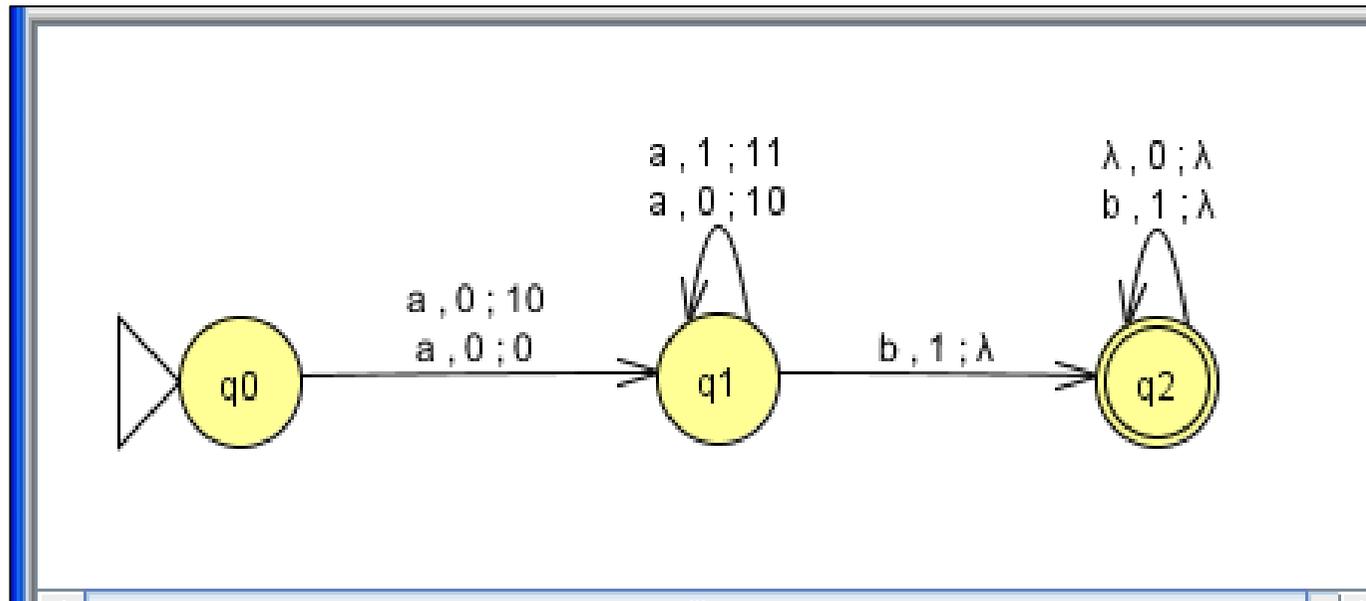




Die Stackmaschine in JFLAP

Es gibt Varianten in der Definition von Stackautomaten. Beispielsweise kann man Endzustände einführen und andere Terminierungsbedingungen.

Wir werden sehen, dass für die Spracherkennung Zustände unwichtig sind: man kommt mit einem Zustand aus, könnte daher auf Zustände gänzlich verzichten.



Die Stackmaschine von vorhin in JFLAP



Konfiguration und Übergang

Konfiguration $\kappa = (q, w, s)$ besteht aus :

- n q - gegenwärtigen Zustand,
- n w - Rest des noch zu lesenden Inputwortes,
- n s - gegenwärtigen Stack.

$\sigma \mid - \tau \quad :\Leftrightarrow$ Automat **kann** von σ in τ übergehen

Inputzeichen wird verbraucht:

$$(q, a.w, g.u) \mid - (p, w, vu) \Leftrightarrow (p, v) \in \delta(q, a, g)$$

Kein Inputzeichen verbraucht

$$(q, w, g.u) \mid - (p, w, vu) \Leftrightarrow (p, v) \in \delta(q, \varepsilon, g)$$



Lauf

Ein **Lauf** ist eine Folge von Übergängen

M

$$\delta(q_0, a, 0) = \{ (q_1, 0), (q_1, 10) \}$$

$$\delta(q_1, a, 0) = \{ (q_1, 10) \}$$

$$\delta(q_1, a, 1) = \{ (q_1, 11) \}$$

$$\delta(q_1, b, 1) = \{ (q_2, \epsilon) \}$$

$$\delta(q_2, b, 1) = \{ (q_2, \epsilon) \}$$

$$\delta(q_2, \epsilon, 0) = \{ (q_2, \epsilon) \}$$

$$\delta(q_i, u, k) = \{ \} \quad \text{sonst}$$

$$(q_0, aabb, 0) \mid - (q_1, abb, 10)$$

$$\mid - (q_1, bb, 110)$$

$$\mid - (q_2, b, 10)$$

$$\mid - (q_2, \epsilon, 0)$$

$$\mid - (q_2, \epsilon, \epsilon)$$

$$\text{d. h. } (q_0, aabb, 0) \mid -^* (q_2, \epsilon, \epsilon)$$

Zwei Läufe

$$(q_0, aaabb, 0) \mid - (q_1, aabb, 0)$$

$$\mid - (q_1, abb, 10)$$

$$\mid - (q_1, bb, 110)$$

$$\mid - (q_2, b, 10)$$

$$\mid - (q_2, \epsilon, 0)$$

$$\mid - (q_2, \epsilon, \epsilon)$$

$$\text{also : } (q_0, aaabb, 0) \mid -^* (q_2, \epsilon, \epsilon)$$



Sprache einer Stackmaschine

- n *Sprache einer Stackmaschine*: Menge aller Worte, die von Anfangs- in eine Endkonfiguration überführen können.

$$L(M) = \{ w \in \Sigma^* \mid (q_0, w, s_0) \vdash^* (q_t, \varepsilon, \varepsilon) \}$$

M

$$\delta(q_0, a, 0) = \{ (q_1, 0), (q_1, 10) \}$$

$$\delta(q_1, a, 0) = \{ (q_1, 10) \}$$

$$\delta(q_1, a, 1) = \{ (q_1, 11) \}$$

$$\delta(q_1, b, 1) = \{ (q_2, \varepsilon) \}$$

$$\delta(q_2, b, 1) = \{ (q_2, \varepsilon) \}$$

$$\delta(q_2, \varepsilon, 0) = \{ (q_2, \varepsilon) \}$$

$$\delta(q_i, u, k) = \{ \} \quad \text{sonst}$$

Sprache von M im Beispiel:

$$L(M) = \{ a^n b^m \mid n = m \vee n = m + 1 \}$$



Stackmaschinen für kontextfreie Sprachen

Satz : Für jede kontextfreie Sprache L gibt es eine **non-deterministische** Stackmaschine M , so daß $L = L(M)$.

Aus $G = (V, T, P, S)$ konstruiere $M(G) = (Q, \Sigma, \Gamma, \delta, q_0, s_0)$ mit

$Q = \{S\}$ - nur ein Zustand, also irrelevant

$\Sigma = T$ - Input besteht aus Token

$\Gamma = V \cup T$ - Stack für Terminale und Nonterminale

$q_0 = s_0 = S$ - das Startsymbol

$M(G)$

$\delta(S, a, a) = \{(S, \varepsilon)\}$ für alle $a \in T$

Zustand irrelevant

$\delta(S, \varepsilon, B) = \{(S, \beta) \mid B \rightarrow \beta \in P\}$



Stackmaschinen für kontextfreie Sprachen

Satz : Für jede kontextfreie Sprache L gibt es eine **non-deterministische** Stackmaschine M , so daß $L = L(M)$.

Aus $G = (V, T, P, S)$ konstruiere $M(G) = (Q, \Sigma, \Gamma, \delta, q_0, s_0)$ mit

$Q = \{S\}$ - nur ein Zustand, also irrelevant

$\Sigma = T$ - Input besteht aus Token

$\Gamma = V \cup T$ - Stack für Terminale und Nonterminale

$q_0 = s_0 = S$ - das Startsymbol

$M(G)$

$\delta(a, a) = \{ \varepsilon \}$ für alle $a \in T$

$\delta(\varepsilon, B) = \{ \beta \mid B \rightarrow \beta \in P \}$

Zustand irrelevant

$\delta : (\Sigma \cup \{ \varepsilon \}) \times \Gamma \rightarrow \wp_{\text{finite}}(\Gamma^*)$



Beispiel: Maschine für Klammersprache

$$M(K) = (\{S\}, \{ \}, \{ (\}, \{) \}, (, S \} , \delta, S, S)$$

M(G)

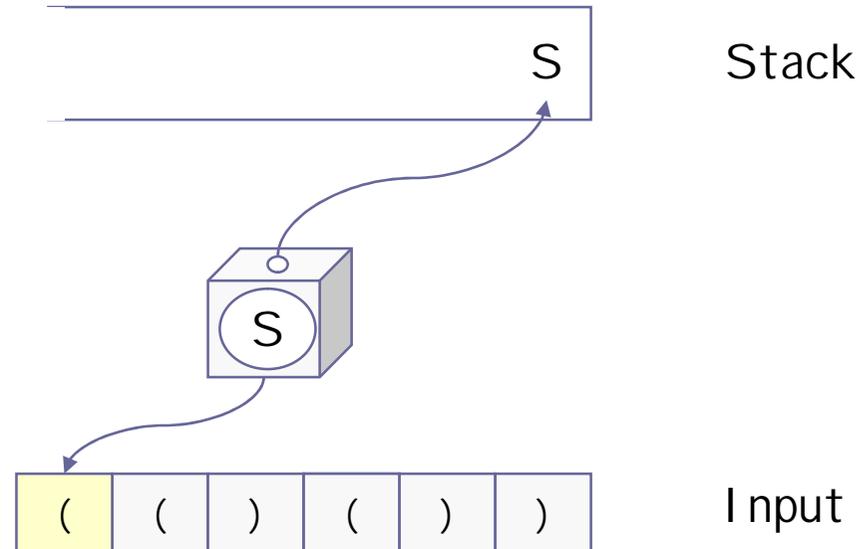
$$\delta(a, a) = \varepsilon \quad \text{für alle } a \in T$$

$$\delta(\varepsilon, B) = \{ \beta \mid B \rightarrow \beta \in P \}$$

Invariante:
Satzform auf Stack
 \approx Rest des Inputs

K

$$\begin{array}{l} S \rightarrow S S \\ \quad | (S) \\ \quad | \varepsilon \end{array}$$





Beispiel: Maschine für Klammersprache

$$M(K) = (\{S\}, \{ \}, \{ (\}, \{) \}, (, S \} , \delta, S, S)$$

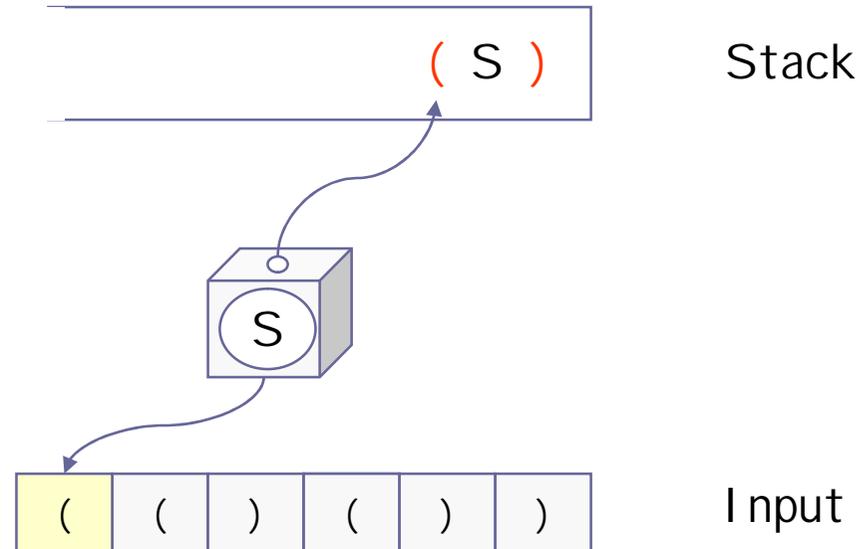
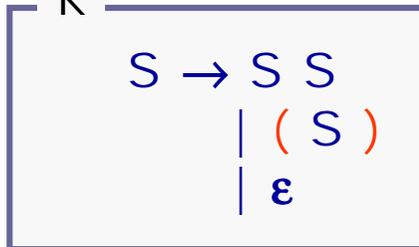
M(G)

$$\delta(S, a, a) = \varepsilon \quad \text{für alle } a \in T$$

$$\delta(S, \varepsilon, B) = \{ \beta \mid B \rightarrow \beta \in P \}$$

Invariante:
Satzform auf Stack
 \approx Rest des Inputs

K





Beispiel: Maschine für Klammersprache

$$M(K) = (\{S\}, \{), (\}, \{) \}, (, S \}, \delta, S, S)$$

M(G)

$$\delta(S, a, a) = \epsilon \quad \text{für alle } a \in T$$

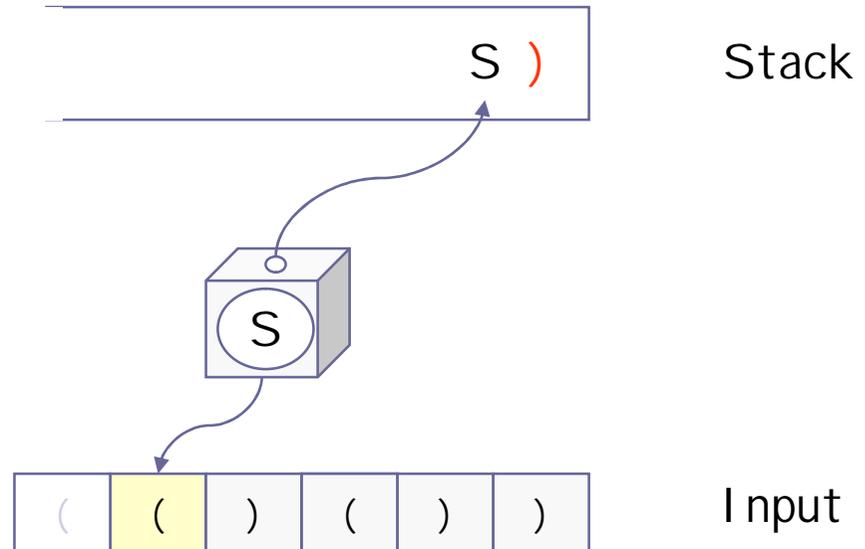
$$\delta(S, \epsilon, B) = \{ \beta \mid B \rightarrow \beta \in P \}$$

Invariante:

Satzform auf Stack
≈ Rest des Inputs

K

$$\begin{array}{l} S \rightarrow S S \\ \quad | \quad (S) \\ \quad | \quad \epsilon \end{array}$$





Beispiel: Maschine für Klammersprache

$$M(K) = (\{S\}, \{ \}, \{ (,) \}, \{ \}, (, S \}, \delta, S, S)$$

M(G)

$$\delta(S, a, a) = \varepsilon \quad \text{für alle } a \in T$$

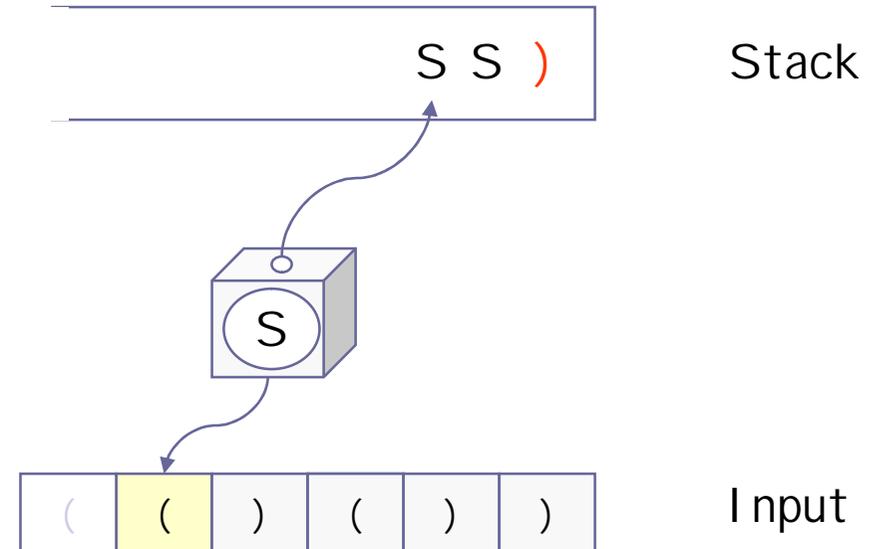
$$\delta(S, \varepsilon, B) = \{ \beta \mid B \rightarrow \beta \in P \}$$

Invariante:

Satzform auf Stack
 \approx Rest des Inputs

K

$$\begin{array}{l} S \rightarrow S S \\ \quad | \quad (S) \\ \quad | \quad \varepsilon \end{array}$$





Beispiel: Maschine für Klammersprache

$$M(K) = (\{S\}, \{ \}, \{ (\}, \{) \}, (, S \} , \delta, S, S)$$

M(G)

$$\delta(S, a, a) = \varepsilon \quad \text{für alle } a \in T$$

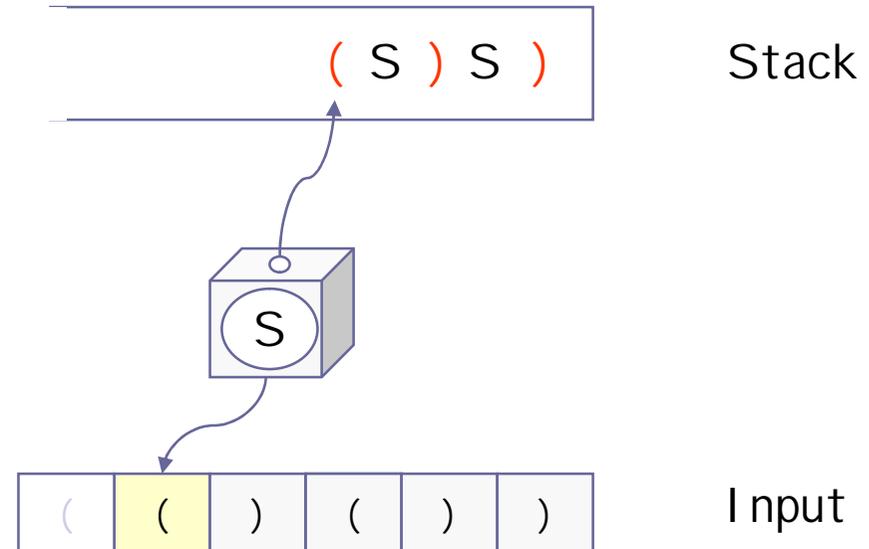
$$\delta(S, \varepsilon, B) = \{ \beta \mid B \rightarrow \beta \in P \}$$

Invariante:

Satzform auf Stack
 \approx Rest des Inputs

K

$$\begin{array}{l} S \rightarrow S S \\ \quad | \quad (S) \\ \quad | \quad \varepsilon \end{array}$$





Beispiel: Maschine für Klammersprache

$$M(K) = (\{S\}, \{ \}, \{ (\}, \{) \}, (, S \} , \delta, S, S)$$

M(G)

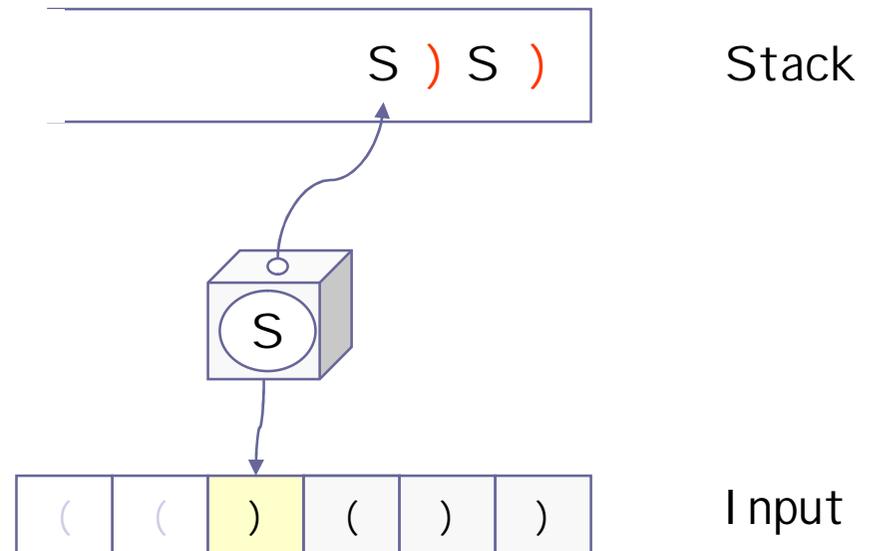
$$\delta(S, a, a) = \varepsilon \quad \text{für alle } a \in T$$

$$\delta(S, \varepsilon, B) = \{ \beta \mid B \rightarrow \beta \in P \}$$



K

$$\begin{array}{l} S \rightarrow S S \\ \quad | \quad (S) \\ \quad | \quad \varepsilon \end{array}$$





Beispiel: Maschine für Klammersprache

$$M(K) = (\{S\}, \{), (\}, \{) \}, (, S \}, \delta, S, S)$$

M(G)

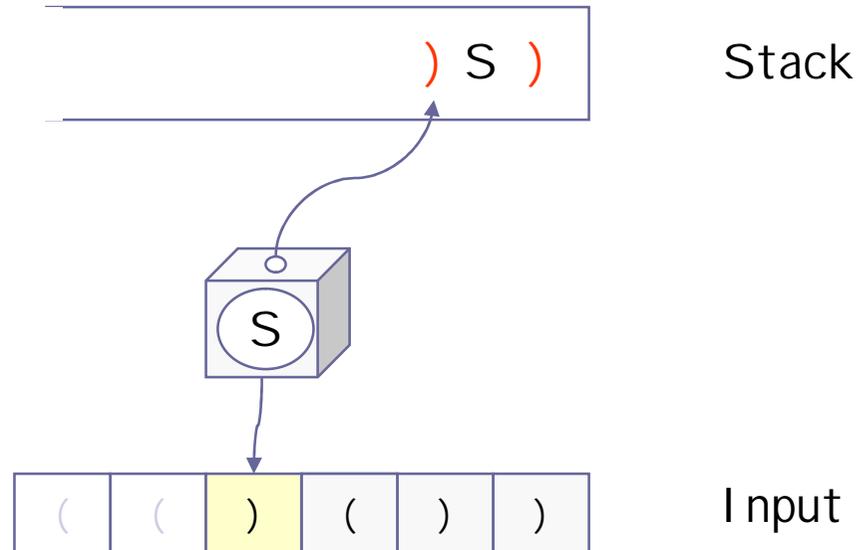
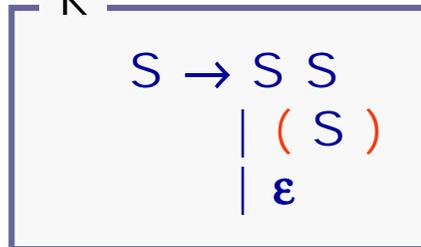
$$\delta(S, a, a) = \epsilon \quad \text{für alle } a \in T$$

$$\delta(S, \epsilon, B) = \{ \beta \mid B \rightarrow \beta \in P \}$$

Invariante:

Satzform auf Stack
≈ Rest des Inputs

K





Beispiel: Maschine für Klammersprache

$$M(K) = (\{S\}, \{), (\}, \{) \}, (, S \}, \delta, S, S)$$

M(G)

$$\delta(S, a, a) = \epsilon \quad \text{für alle } a \in T$$

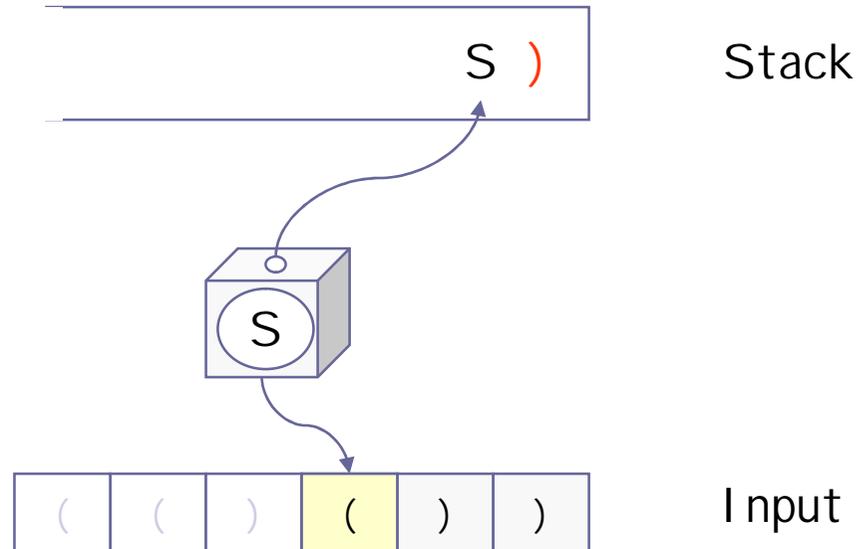
$$\delta(S, \epsilon, B) = \{ \beta \mid B \rightarrow \beta \in P \}$$

Invariante:

Satzform auf Stack
 \approx Rest des Inputs

K

$$\begin{array}{l} S \rightarrow S S \\ \quad | (S) \\ \quad | \epsilon \end{array}$$





Beispiel: Maschine für Klammersprache

$$M(K) = (\{S\}, \{), (\}, \{) \}, (, S \}, \delta, S, S)$$

M(G)

$$\delta(S, a, a) = \varepsilon \quad \text{für alle } a \in T$$

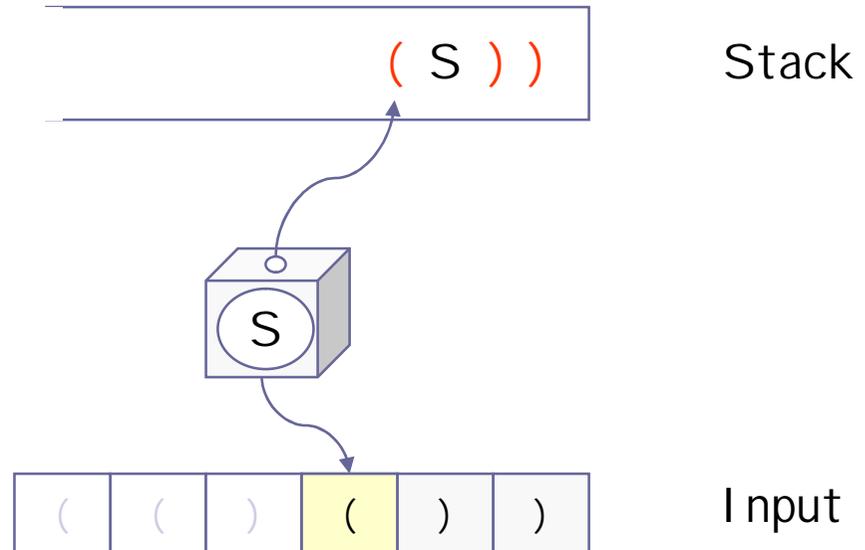
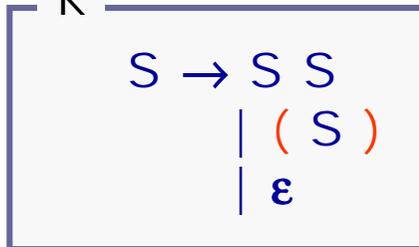
$$\delta(S, \varepsilon, B) = \{ \beta \mid B \rightarrow \beta \in P \}$$



Invariante:

Satzform auf Stack
 \approx Rest des Inputs

K





Beispiel: Maschine für Klammersprache

$$M(K) = (\{S\}, \{ \text{), (} \}, \{ \text{), (} \}, (, S \}, \delta, S, S)$$

M(G)

$$\delta(S, a, a) = \varepsilon \quad \text{für alle } a \in T$$

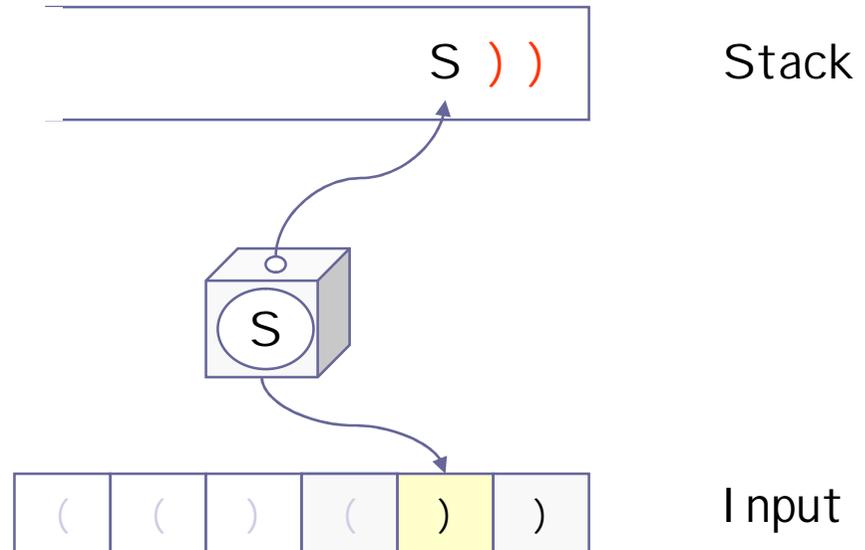
$$\delta(S, \varepsilon, B) = \{ \beta \mid B \rightarrow \beta \in P \}$$

Invariante:

Satzform auf Stack
 \approx Rest des Inputs

K

$$\begin{array}{l} S \rightarrow S S \\ \quad | \quad (S) \\ \quad | \quad \varepsilon \end{array}$$





Beispiel: Maschine für Klammersprache

$$M(K) = (\{S\}, \{ \}, \{ (\}, \{) \}, (, S \} , \delta, S, S)$$

M(G)

$$\delta(S, a, a) = \epsilon \quad \text{für alle } a \in T$$

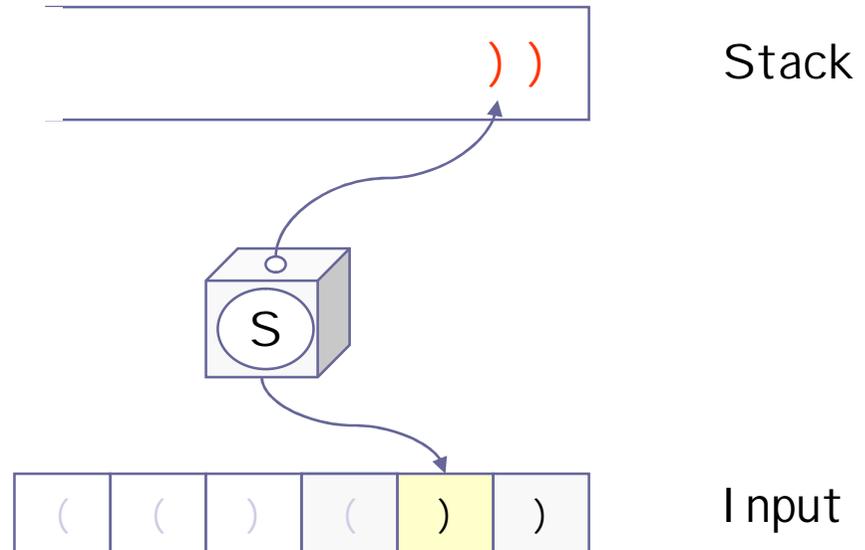
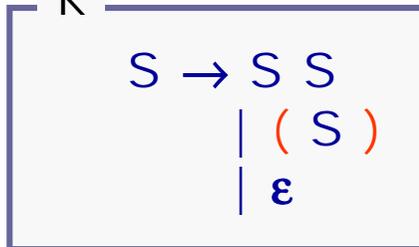
$$\delta(S, \epsilon, B) = \{ \beta \mid B \rightarrow \beta \in P \}$$



Invariante:

Satzform auf Stack
≈ Rest des Inputs

K





Beispiel: Maschine für Klammersprache

$$M(K) = (\{S\}, \{), (\}, \{) \}, (, S \}, \delta, S, S)$$

M(G)

$$\delta(S, a, a) = \varepsilon \quad \text{für alle } a \in T$$

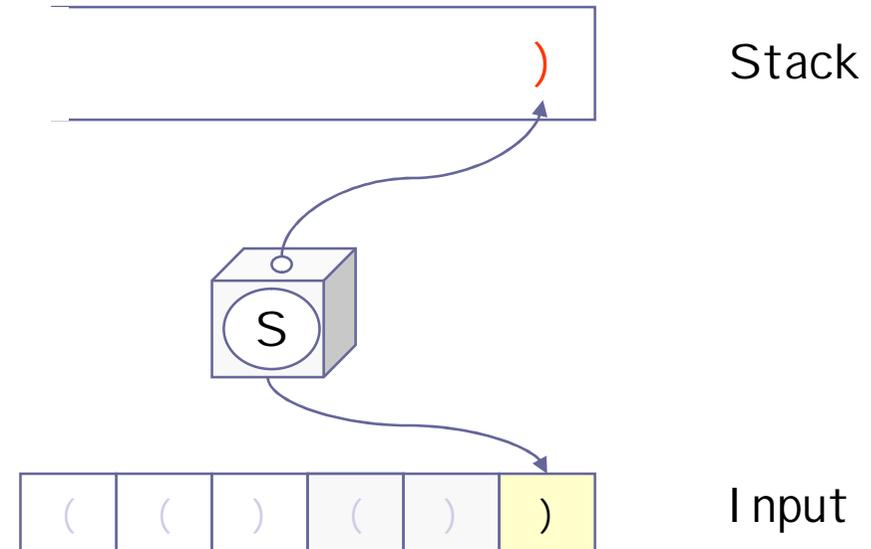
$$\delta(S, \varepsilon, B) = \{ \beta \mid B \rightarrow \beta \in P \}$$

Invariante:

Satzform auf Stack
 \approx Rest des Inputs

K

$$\begin{array}{l} S \rightarrow S S \\ \quad | \quad (S) \\ \quad | \quad \varepsilon \end{array}$$





Beispiel: Maschine für Klammersprache

$$M(K) = (\{S\}, \{ \}, \{ (\}, \{) \}, (, S \} , \delta, S, S)$$

M(G)

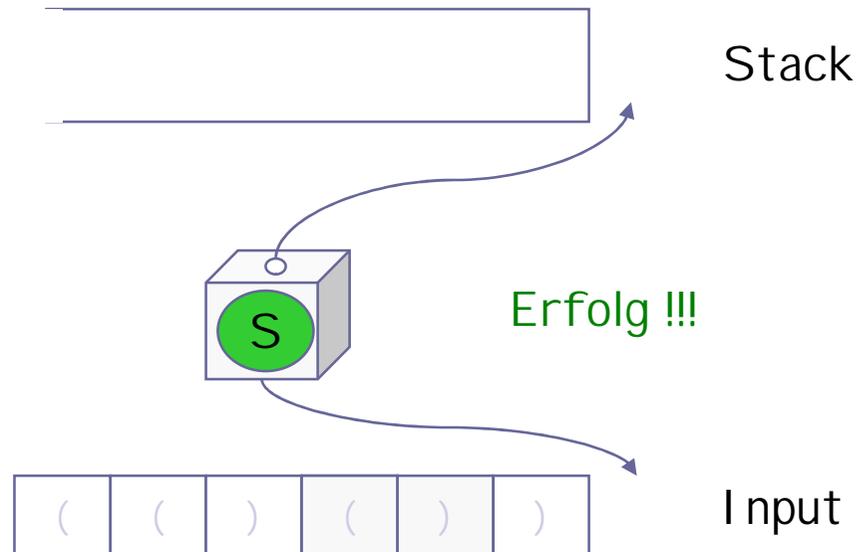
$$\delta(S, a, a) = \epsilon \quad \text{für alle } a \in T$$

$$\delta(S, \epsilon, B) = \{ \beta \mid B \rightarrow \beta \in P \}$$

Invariante:
Satzform auf Stack
 \approx Rest des Inputs

K

$$\begin{array}{l} S \rightarrow S S \\ \quad | (S) \\ \quad | \epsilon \end{array}$$





Stackmaschinen für kontextfreie Sprachen

Satz : Für jede kontextfreie Sprache L gibt es eine **nicht-deterministische** Stackmaschine M , so daß $L = L(M)$.

Beweis: Die Stackmaschine mache n Schritte und akzeptiere w . Im i -ten Schritt sei u_i der bereits verarbeitete Input und v_i der Rest (also immer $u_i v_i = w$) und α_i der Inhalt des Stacks.

Behauptung:

$S = u_0 \alpha_0 \Rightarrow u_1 \alpha_1 \Rightarrow \dots \Rightarrow u_n \alpha_n = u_n$ ist eine Linksableitung von u_n .

Klar, denn anfangs ist $u_0 = \varepsilon$ und $\alpha_0 = S$ und am Schluss gilt $u_n = w$.

Für den k -ten Schritt gilt

entweder

$$\alpha_k = a \cdot \alpha_{k+1} \text{ und } v_k = a \cdot v_{k+1} \text{ und } u_{k+1} = u_k a$$

// Regel 1

$$\text{also } u_k \alpha_k = u_{k+1} \alpha_{k+1}$$

$$\text{oder } \alpha_k = B \cdot \gamma, \alpha_{k+1} = \beta \cdot \gamma \text{ für ein } B \rightarrow \beta \in P$$

// Regel 2

$$\text{also } u_k \alpha_k = u_k B \cdot \gamma \Rightarrow u_k \beta \cdot \gamma = u_{k+1} \alpha_{k+1}$$

in jedem Fall ein (evtl. trivialer) Schritt einer Linksherleitung

Umgekehrt, kann man jede Linksherleitung simulieren.

$M(G)$

$$\delta(S, a, a) = \varepsilon \quad \text{für alle } a \in T$$

$$\delta(S, \varepsilon, B) = \{ \beta \mid B \rightarrow \beta \in P \}$$

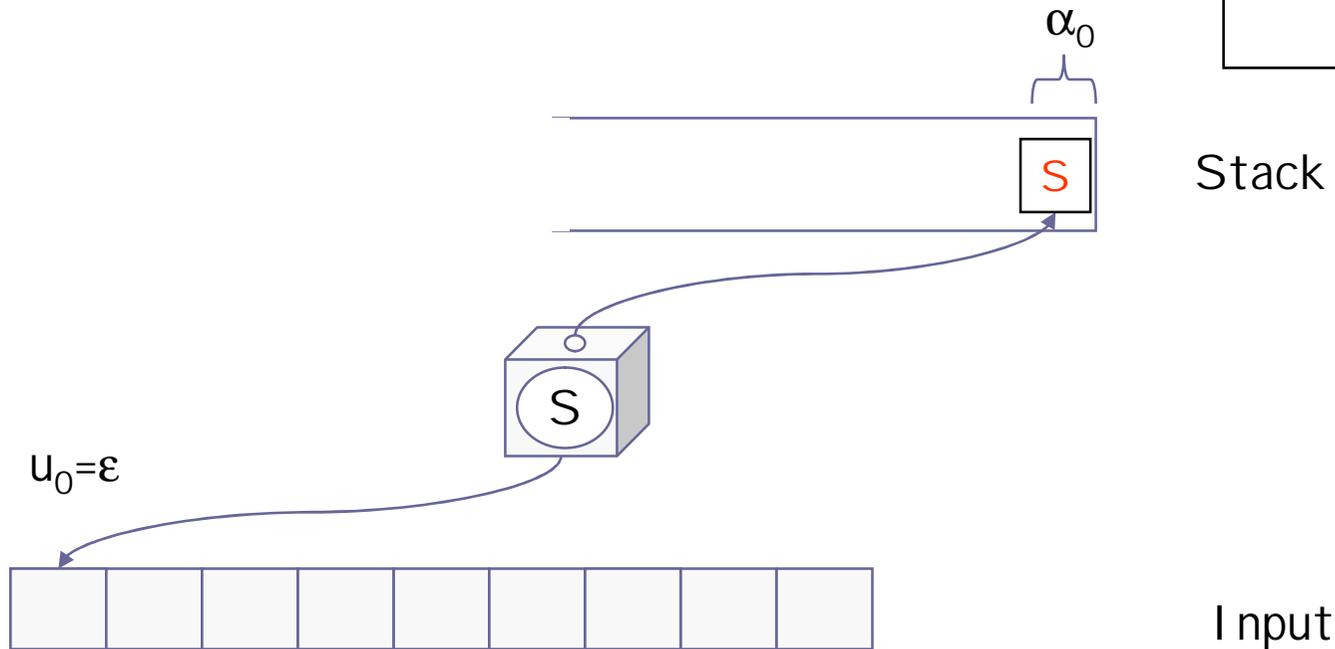


Anfangssituation

$M(G)$

$$\delta(S, a, a) = \varepsilon \quad \text{für alle } a \in T$$

$$\delta(S, \varepsilon, B) = \{ \beta \mid B \rightarrow \beta \in P \}$$



Linksherleitung:

$$S = u_0 \alpha_0$$

$$\Rightarrow \dots$$

$$\Rightarrow u_k \alpha_k$$

$$\Rightarrow$$

$$u_{k+1} \alpha_{k+1}$$

$$\Rightarrow \dots$$

$$\Rightarrow W$$

Stack

Input

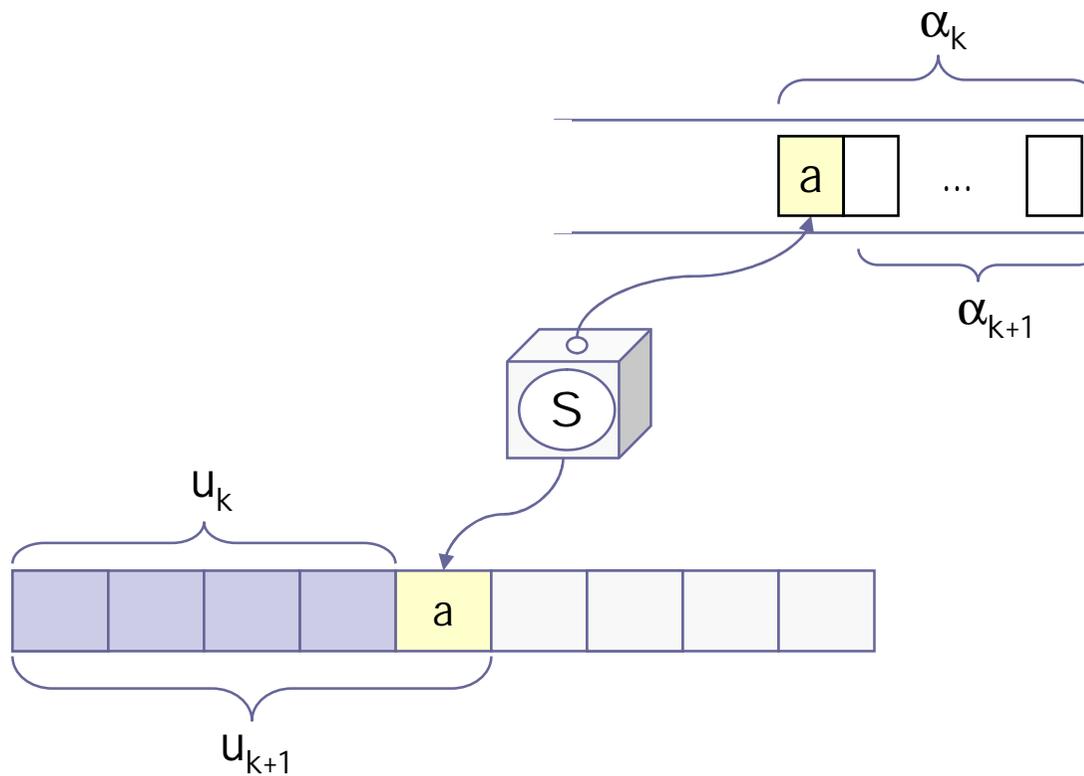


Anwendung der ersten Regel

M(G)

$$\delta(S, a, a) = \varepsilon \quad \text{für alle } a \in T$$

$$\delta(S, \varepsilon, B) = \{ \beta \mid B \rightarrow \beta \in P \}$$



Linksherleitung:

$$S = u_0 \alpha_0$$

$\Rightarrow \dots$

$$\Rightarrow u_k \alpha_k$$

$$= u_{k+1} \alpha_{k+1}$$

$\Rightarrow \dots$

$\Rightarrow W$

Stack

Input

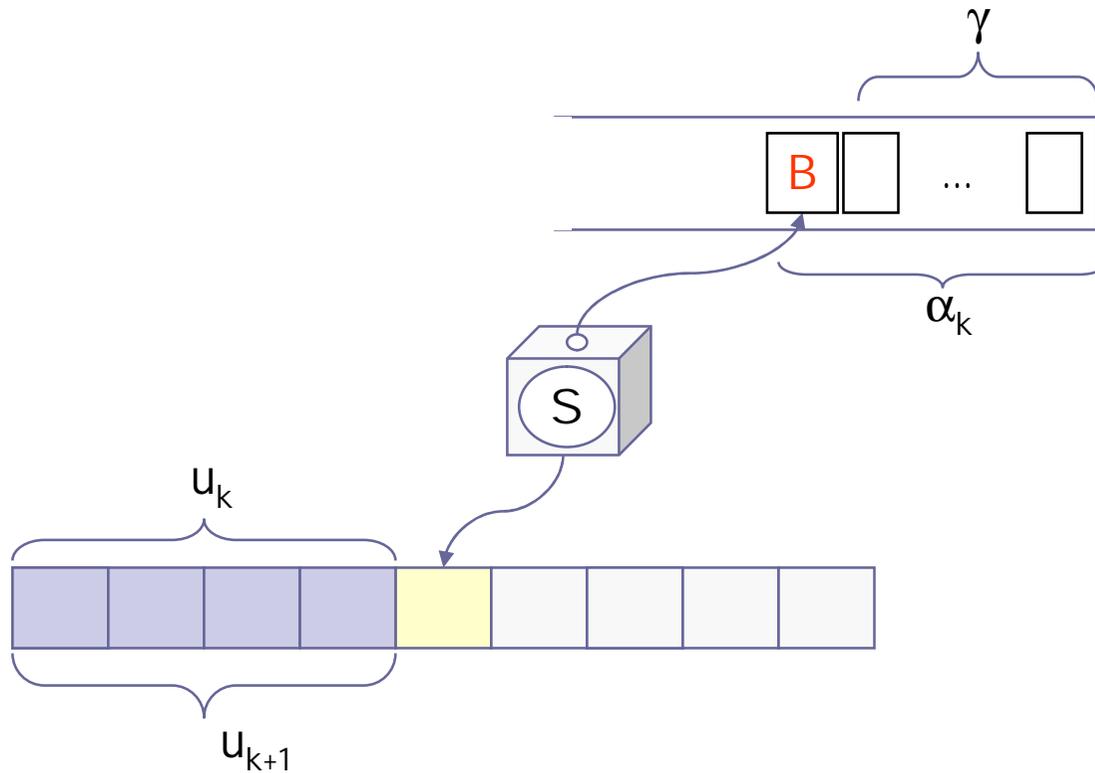


Anwendung der zweiten Regel

M(G)

$$\delta(S, a, a) = \varepsilon \quad \text{für alle } a \in T$$

$$\delta(S, \varepsilon, B) = \{ \beta \mid B \rightarrow \beta \in P \}$$



Linksherleitung:

$$S = u_0 \alpha_0$$

$$\Rightarrow \dots$$

$$\Rightarrow u_k \alpha_k$$

$$= u_k B \gamma$$

$$\Rightarrow u_k \beta \gamma$$

$$= u_{k+1} \alpha_{k+1}$$

$$\Rightarrow \dots$$

$$\Rightarrow W$$

Stack

Input

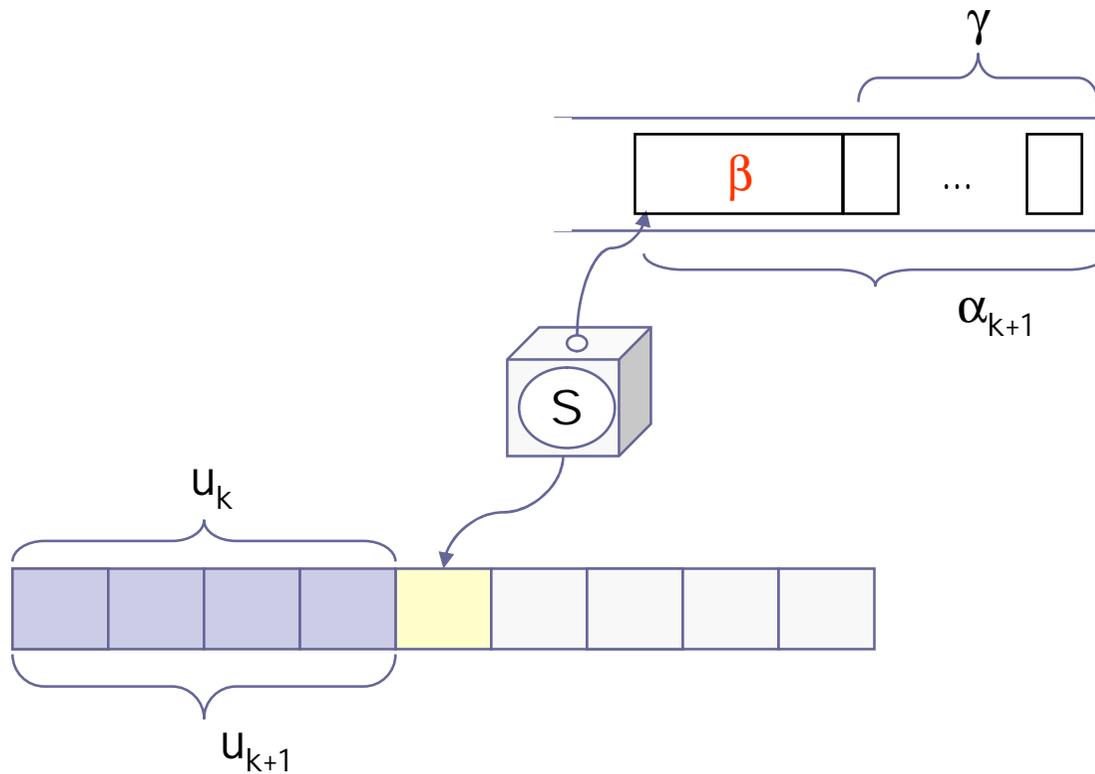


Anwendung der zweiten Regel

M(G)

$$\delta(S, a, a) = \varepsilon \quad \text{für alle } a \in T$$

$$\delta(S, \varepsilon, B) = \{ \beta \mid B \rightarrow \beta \in P \}$$



Linksherleitung:

$$S = u_0 \alpha_0$$

$\Rightarrow \dots$

$$\Rightarrow u_k \alpha_k$$

$$= u_k B \gamma$$

$$\Rightarrow u_k \beta \gamma$$

$$= u_{k+1} \alpha_{k+1}$$

$\Rightarrow \dots$

$\Rightarrow W$

Stack

Input



Endsituation

M(G)

$$\delta(S, a, a) = \varepsilon \quad \text{für alle } a \in T$$

$$\delta(S, \varepsilon, B) = \{ \beta \mid B \rightarrow \beta \in P \}$$

Linksherleitung:

$$S = u_0 \alpha_0$$

$$\Rightarrow \dots$$

$$\Rightarrow u_k \alpha_k$$

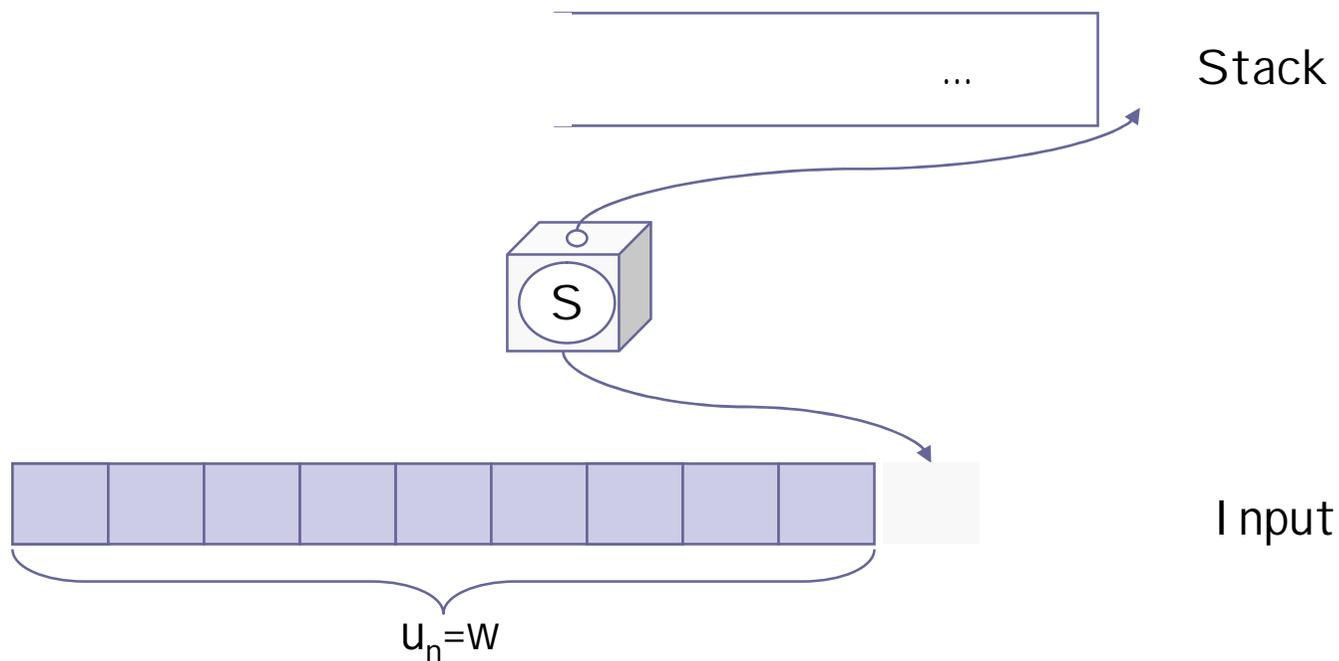
$$\Rightarrow$$

$$u_{k+1} \alpha_{k+1}$$

$$\Rightarrow \dots$$

$$\Rightarrow u_n = W$$

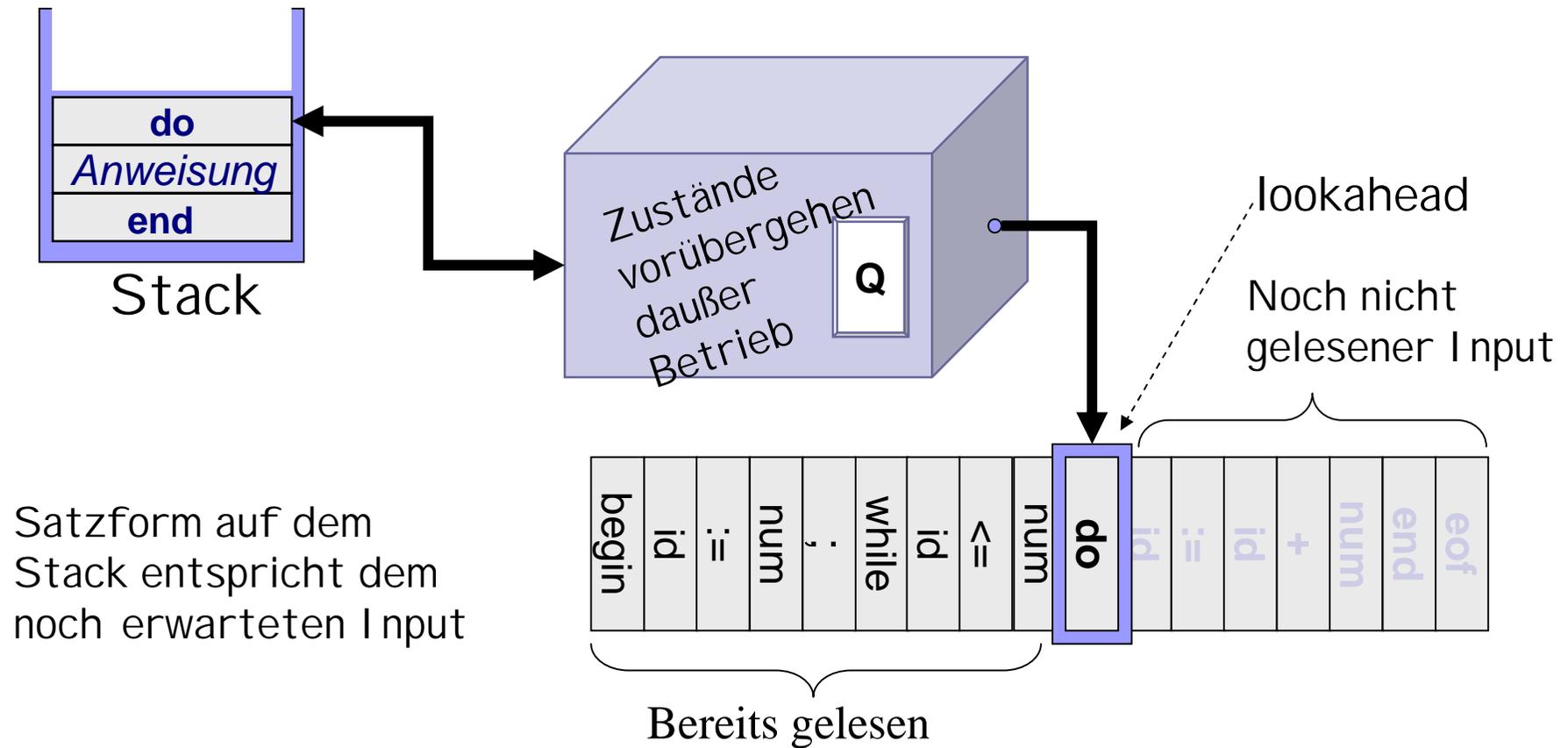
$$\alpha_n = \varepsilon$$





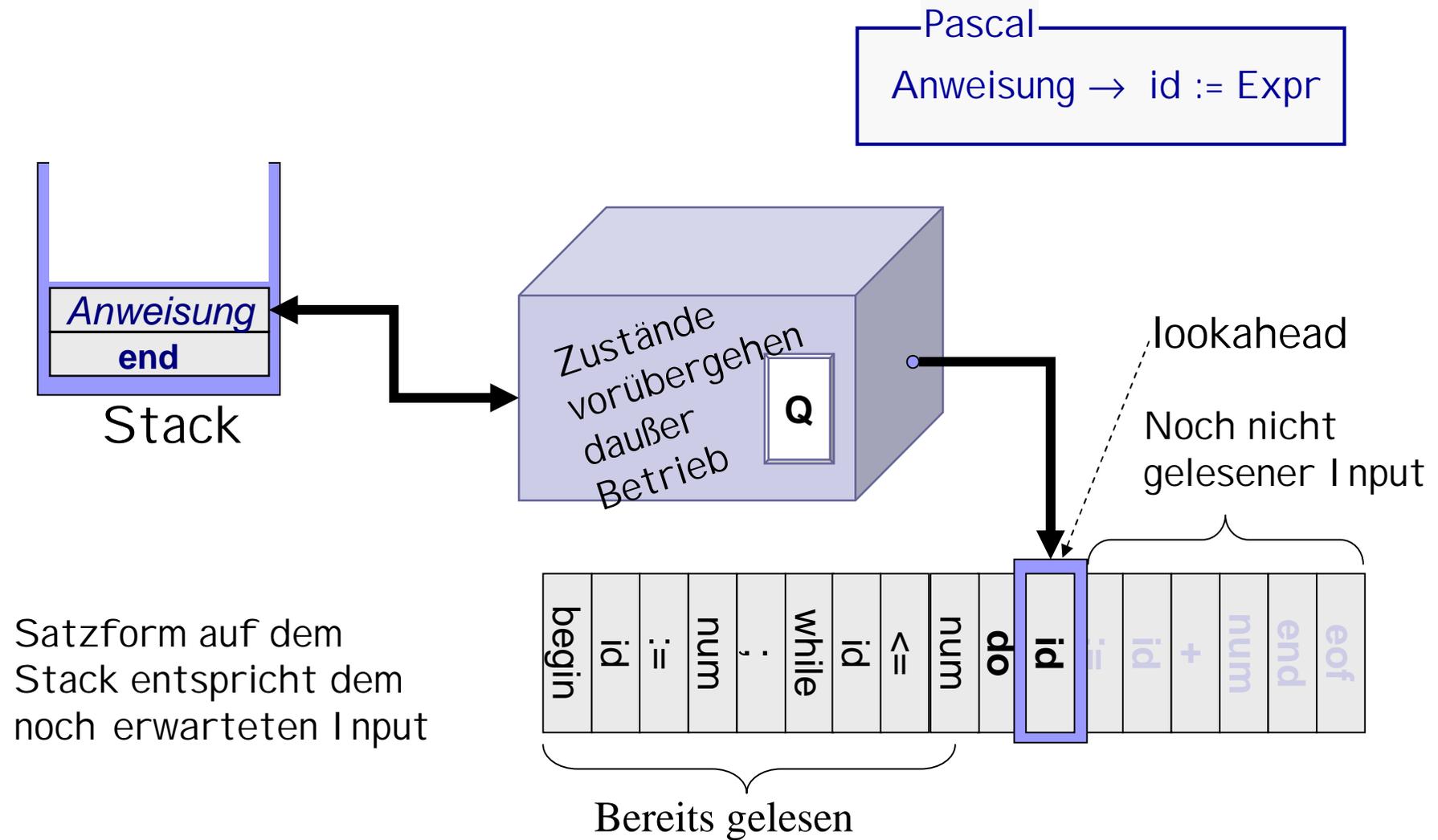
Nichtdeterministischer Stackautomat

Input hat das erwartete Token



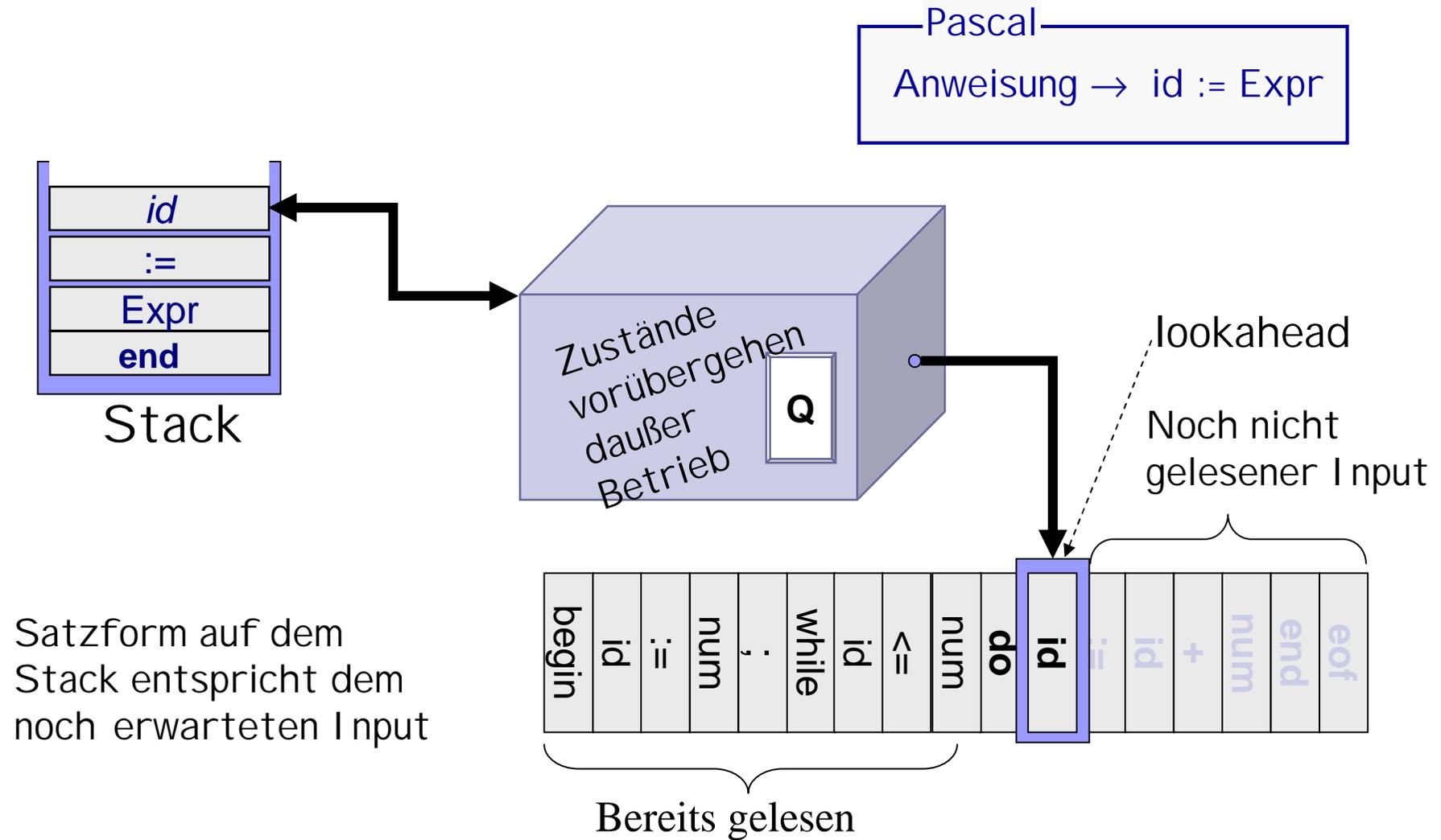


Nichtdeterministischer Stackautomat





Nichtdeterministischer Stackautomat





Stackmaschinen und kontextfreie Sprachen

Satz : Für jede kontextfreie Sprache L gibt es eine nicht-deterministische Stackmaschine M , so daß $L = L(M)$.

Ein Algorithmus, der überprüft, ob ein Wort w in der Sprache ist, insbesondere auch wie es hergeleitet werden kann, nennt man einen

Parser

Der obige Satz hat auch eine Umkehrung:

Ztas : Für Stackmaschine M gibt es eine CF-Grammatik G so dass $L(M)=L(G)$..

Da dieser aber für die Praxis (bisher) ohne Relevanz ist, nehmen wir ihn ohne Beweis zur Kenntnis.



Inhalt

1. Grammatiken und Sprachen
 - .. Kontextfreie Grammatiken
 - .. Herleitungen, Linksherleitungen
 - .. Sprachen zu einer Grammatik
 - .. Äquivalenz
 - .. Chomsky-Normalform
 - .. Wortproblem, CYK
 - .. Pumping Lemma für CF-Sprachen
2. Stackautomaten
 - .. Definitionen und Beispiele
 - .. Konfigurationen, Läufe,
 - .. Sprache eines Stackautomaten
 - .. Parser

3. Parser
 - .. Mehrdeutigkeit
 - .. Bottom Up, Top Down
 - .. Recursive descent Parser
 - .. First, Follow
 - .. Semantische Aktionen
4. Shift-Reduce Parser
 - .. Konflikte
 - .. LR(0)-Zustände
 - .. Automat zur Grammatik
 - .. Shift-Reduce, Reduce-Reduce Konflikte
 - .. Präzedenz und Assoziativität
 - .. lex und yacc
 - .. Arbeitsweise eines Compilers



Mehrdeutigkeit

Besitzt $L(G)$ ein Wort w , zu dem es mehrere Herleitungsbäume gibt, so heißt G **mehrdeutig**.

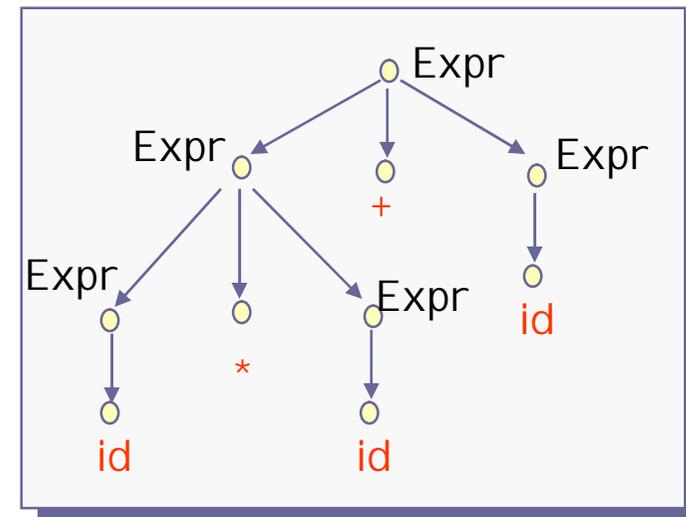
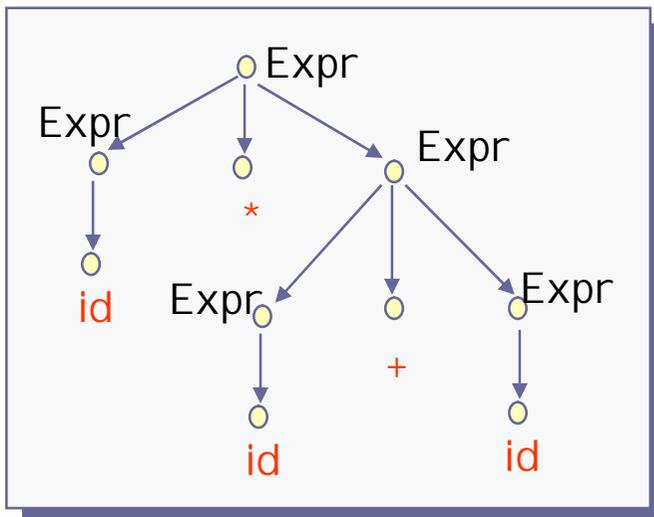
E

$$\begin{array}{l} \text{Expr} \rightarrow \text{Expr} + \text{Expr} \\ \quad \quad | \text{Expr} * \text{Expr} \\ \quad \quad | (\text{Expr}) \\ \quad \quad | \text{id} \end{array}$$

Diese Grammatik ist mehrdeutig, denn das Wort

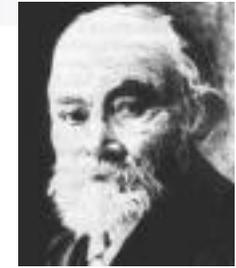
$\text{id} * \text{id} + \text{id}$

hat zwei Herleitungsbäume (äqu.: zwei Linksherleitungen)





Bedeutung - Semantik

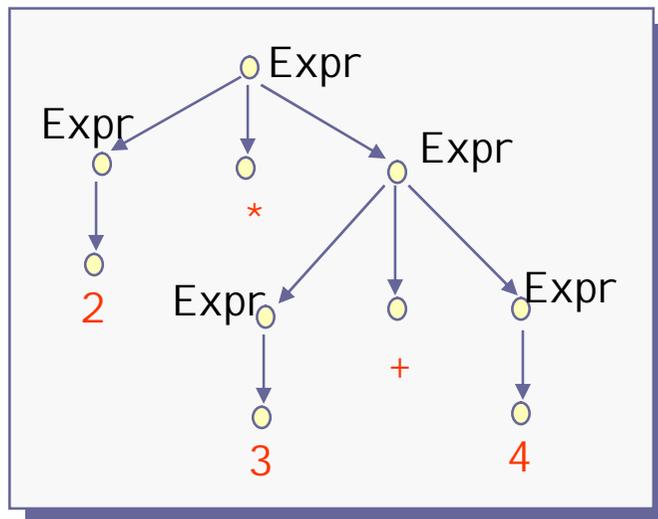


G. Frege

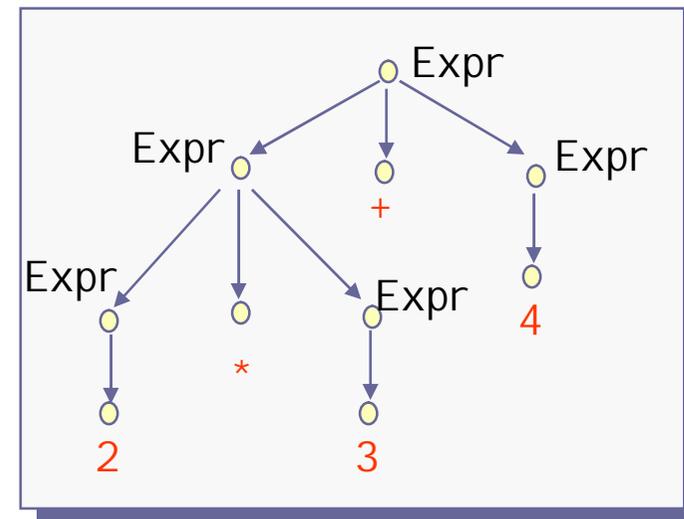
Denotationelles Prinzip (Frege): Die Bedeutung des Ganzen ist eine Funktion der Bedeutung der Teile.

Herleitungsbaum gibt Aufbau und damit die intendierte Bedeutung des Satzes eindeutig wieder. Lineare Schreibweise als String ist oft nicht eindeutig:

$$2 * 3 + 4 = ?$$



14



10



Strategien zur Mehrdeutigkeit

Im Sinne einer einfachen Sprachbeschreibung nimmt man Mehrdeutigkeiten einer Grammatik oft in Kauf. Beim Parsen gibt es verschiedene Strategien mit der Mehrdeutigkeit umzugehen :

$$\begin{array}{l} G_1 \\ \text{Expr} \rightarrow \text{Expr} + \text{Expr} \\ \quad | \text{ num} \end{array}$$

Mehrdeutigkeit unerheblich,
da Addition assoziativ.

$$\begin{array}{l} G_2 \\ \text{Expr} \rightarrow \text{Expr} - \text{Expr} \\ \quad | \text{ num} \end{array}$$

Mehrdeutigkeit wesentlich,
da Subtraktion nicht assoziativ.


$$\begin{array}{l} G'_2 \\ \text{Expr} \rightarrow \text{Expr} - \text{num} \\ \quad | \text{ num} \end{array}$$

Umgewandelt in
eindeutige Grammatik.



Strategien zur Mehrdeutigkeit

G_3

```
Expr → Expr + Expr
      | Expr * Expr
      | ( Expr )
      | num
```

Mehrdeutigkeit erheblich,
z.B. $7 + 3 * 4$ hat zwei Bäume.

Neue syntaktische Kategorien
(Nonterminale) für jede
Präzedenzstufe

G'_3

```
Expr → Expr + Term
      | Term
Term  → Term * Factor
      | Factor
Factor → ( Expr )
      | num
```

Äquivalente eindeutige Grammatik .

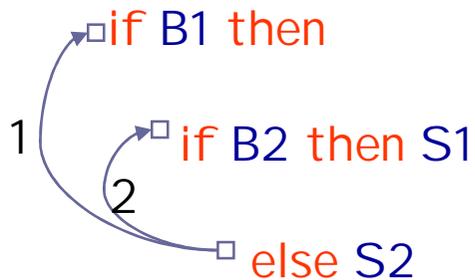


Strategien zur Mehrdeutigkeit

G₄

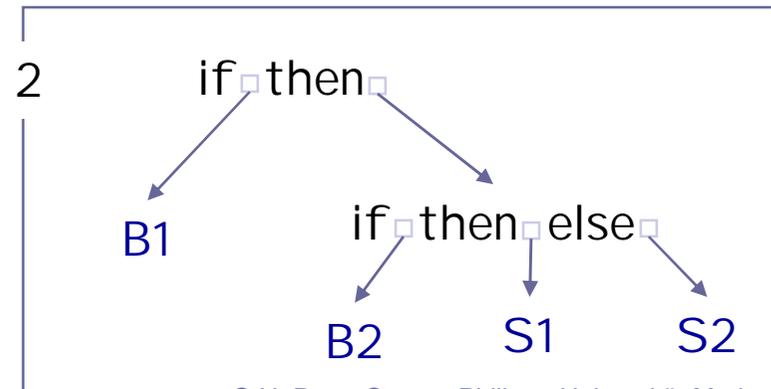
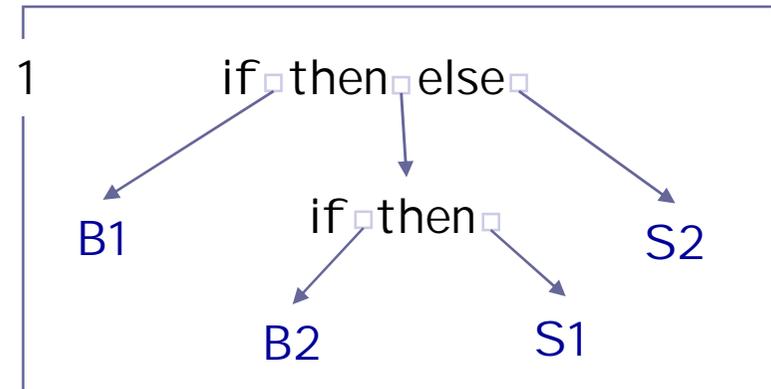
Stmt → if Bexpr then Stmt else Stmt
| if Bexpr then Stmt
| ...

Mehrdeutig:



Verabredung :

Jedes **else** ergänzt das letzte unergänzte **then**

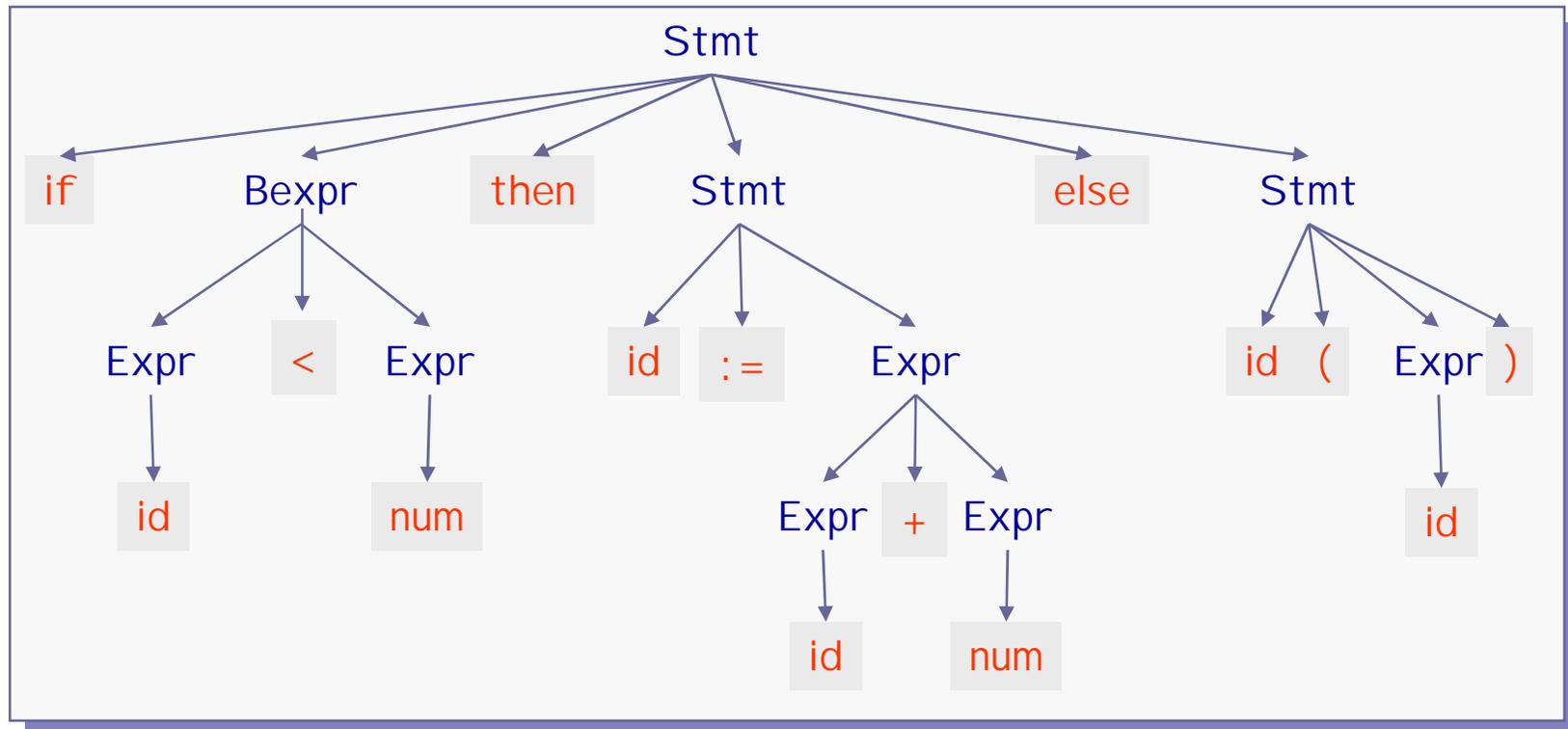




Parsen

Parsen bedeutet, zu entscheiden, ob $w \in L(G)$ und ggf. einen Herleitungsbaum zu finden.

End-
Ergebnis:



Input:

`if x < 0 then y := y + 1 else Inc (z)`



Parserstrategien



- n Bottom Up
 - .. Baum wächst von den Blättern zur Wurzel
- n Top Down
 - .. Baum wächst von der Wurzel zu den Blättern

Beispiel anhand
nebenstehender
Grammatik :

G

```
Stmt → if Bexpr then Stmt else Stmt  
      | id := Expr  
      | id (Expr)
```

```
Expr → Expr + Expr  
      | id  
      | num
```

```
Bexpr → Expr < Expr  
       | true
```



Bottom Up Strategie

n Bottom Up

- .. Baum wächst von den Blättern zur Wurzel
- .. jederzeit ein Wald (mehrere Bäume) vorhanden

n Zwei Aktionen

.. Shift

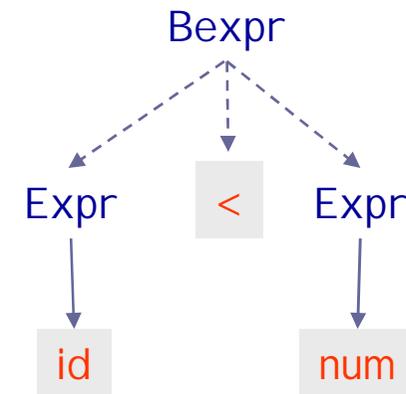
- n Lese nächstes Token im Input
- n mache es zur Wurzel eines neuen Baumes

.. Reduce

- n Die Wurzeln **der letzten** Bäume bilden die rechte Seite **einer** Regel $A \rightarrow w_1 w_2 \dots w_n$
- n Füge sie unter dem neuen Knoten A zusammen

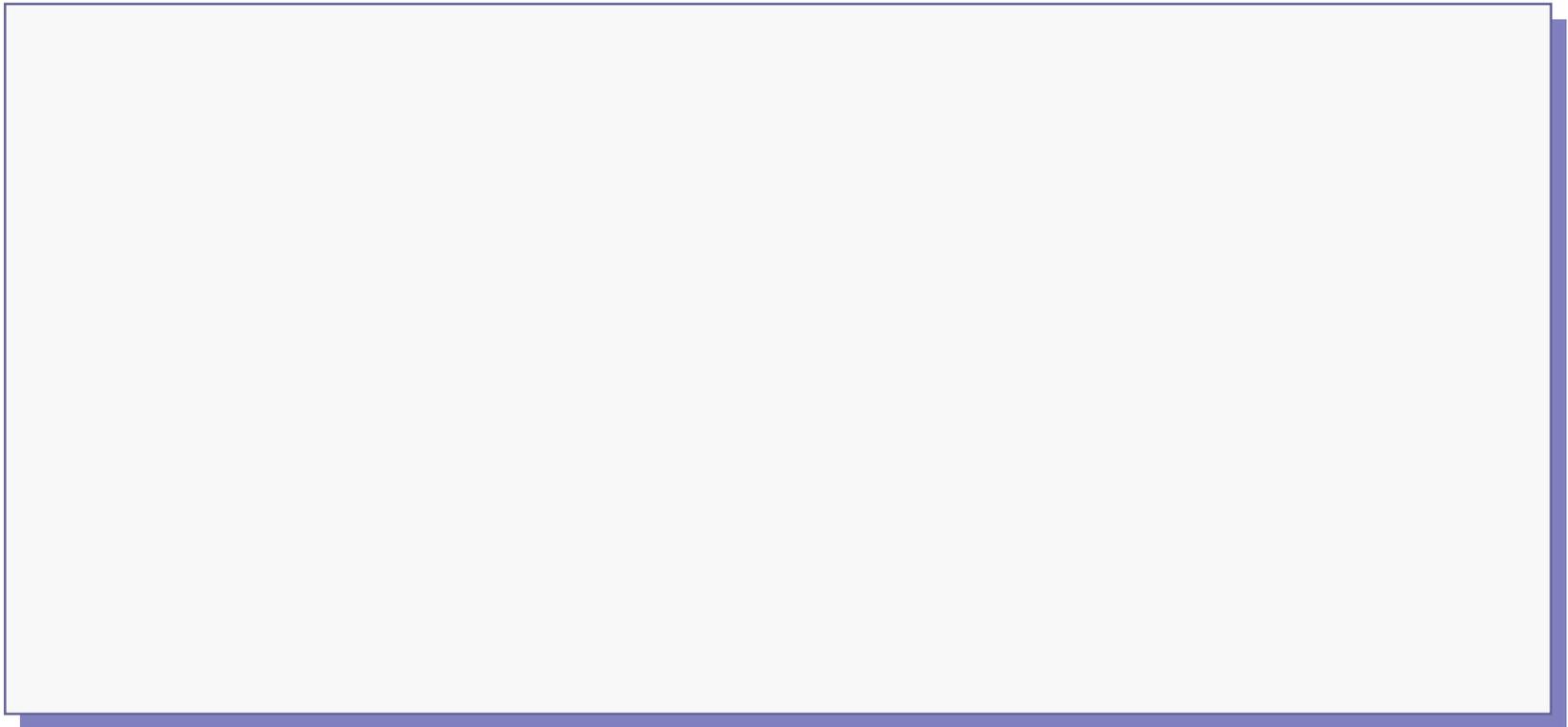


$Bexpr \rightarrow Expr < Expr$
| true





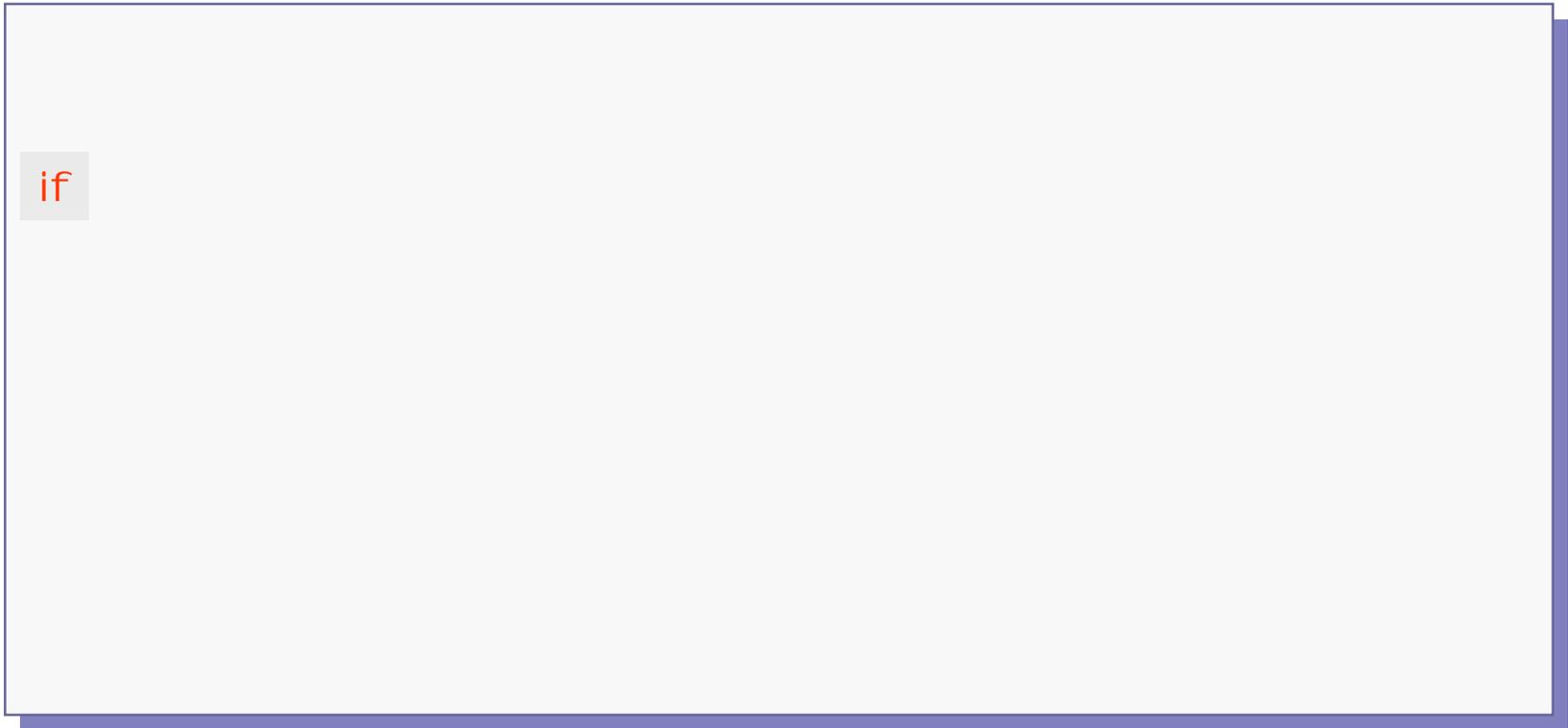
Bottom Up



<code>if</code>	
-----------------	--



Bottom Up



if





Bottom Up



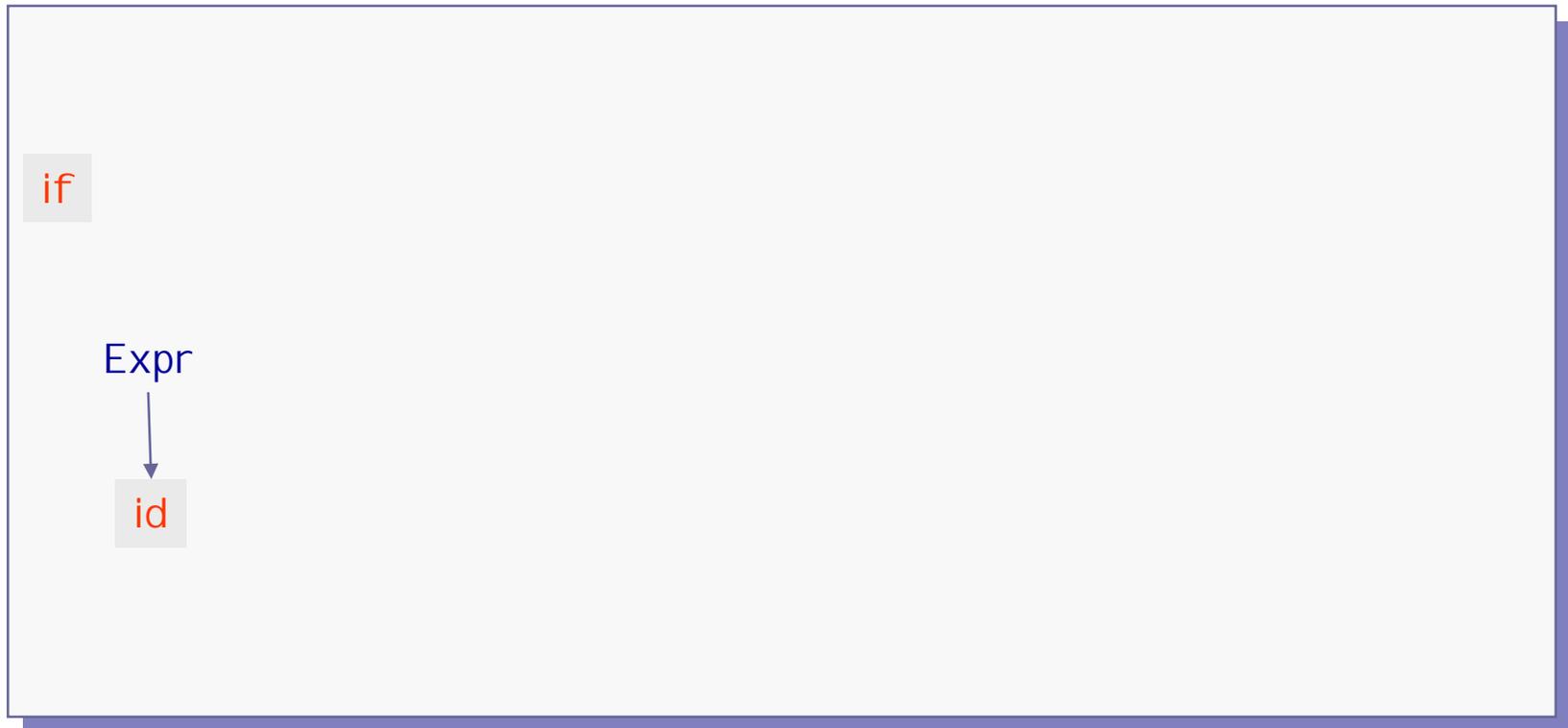
Expr \rightarrow Expr + Expr
Expr \rightarrow id
Expr \rightarrow num



if x <



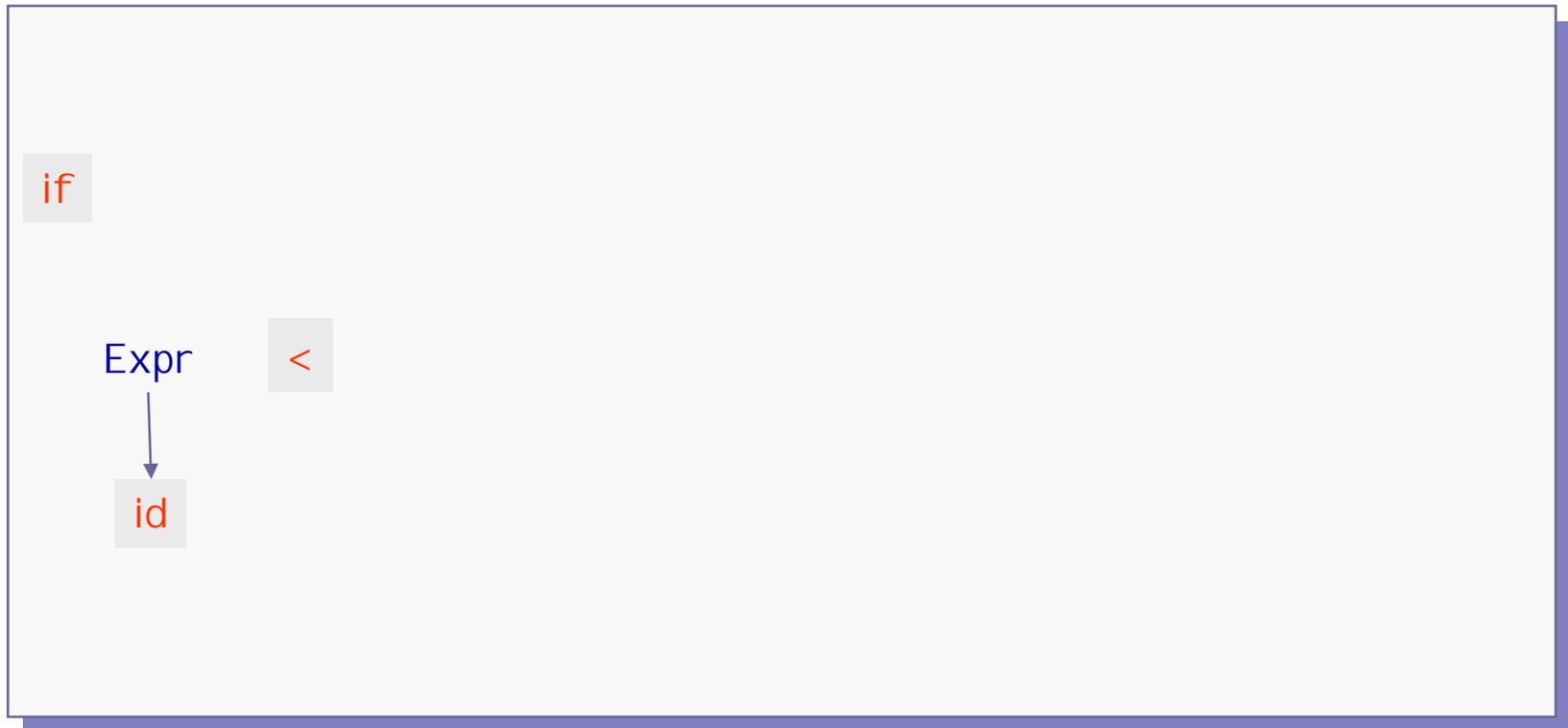
Bottom Up



`if` `x` `<`



Bottom Up



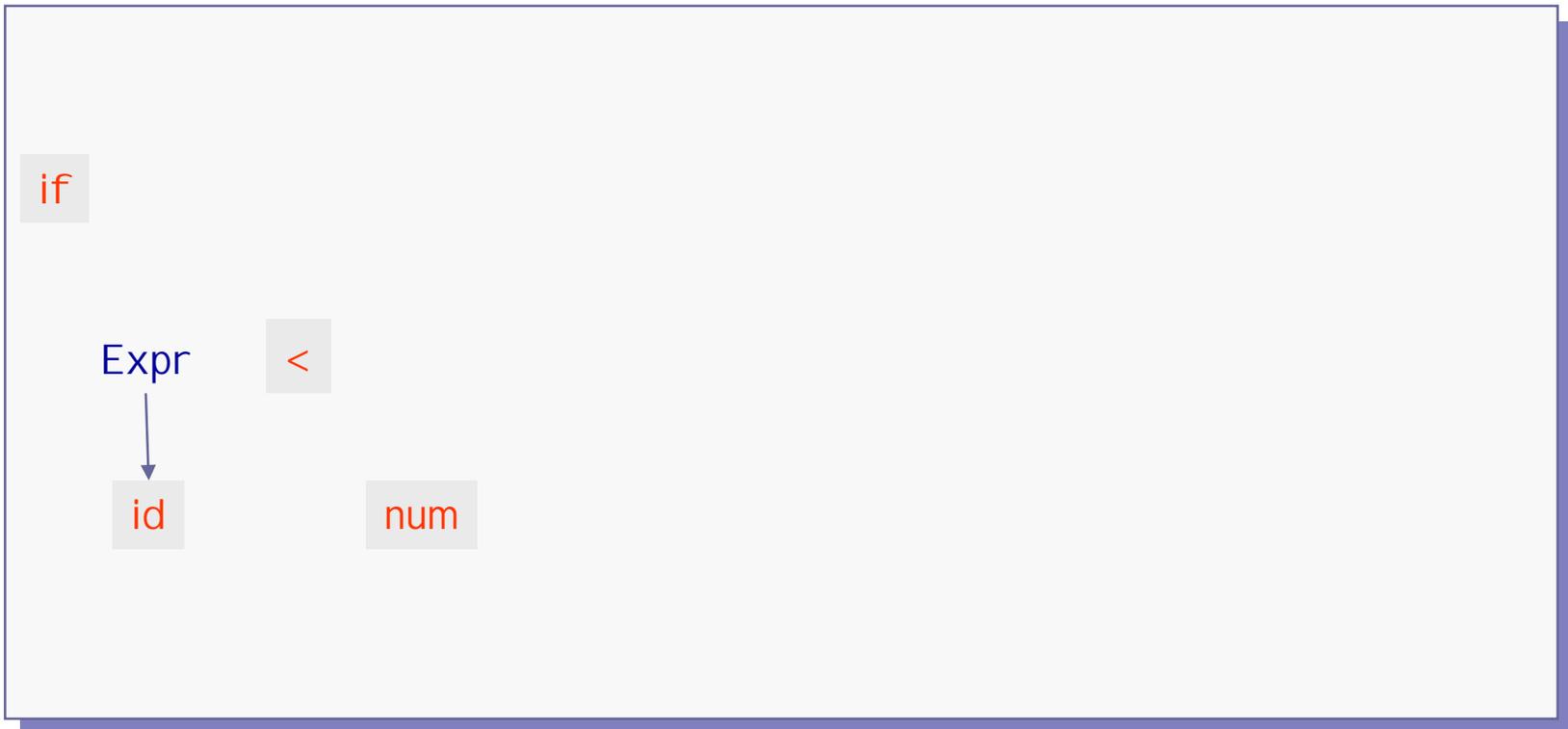
if	x	<	0	
----	---	---	---	--



Bottom Up



Expr \rightarrow Expr + Expr
Expr \rightarrow id
Expr \rightarrow num



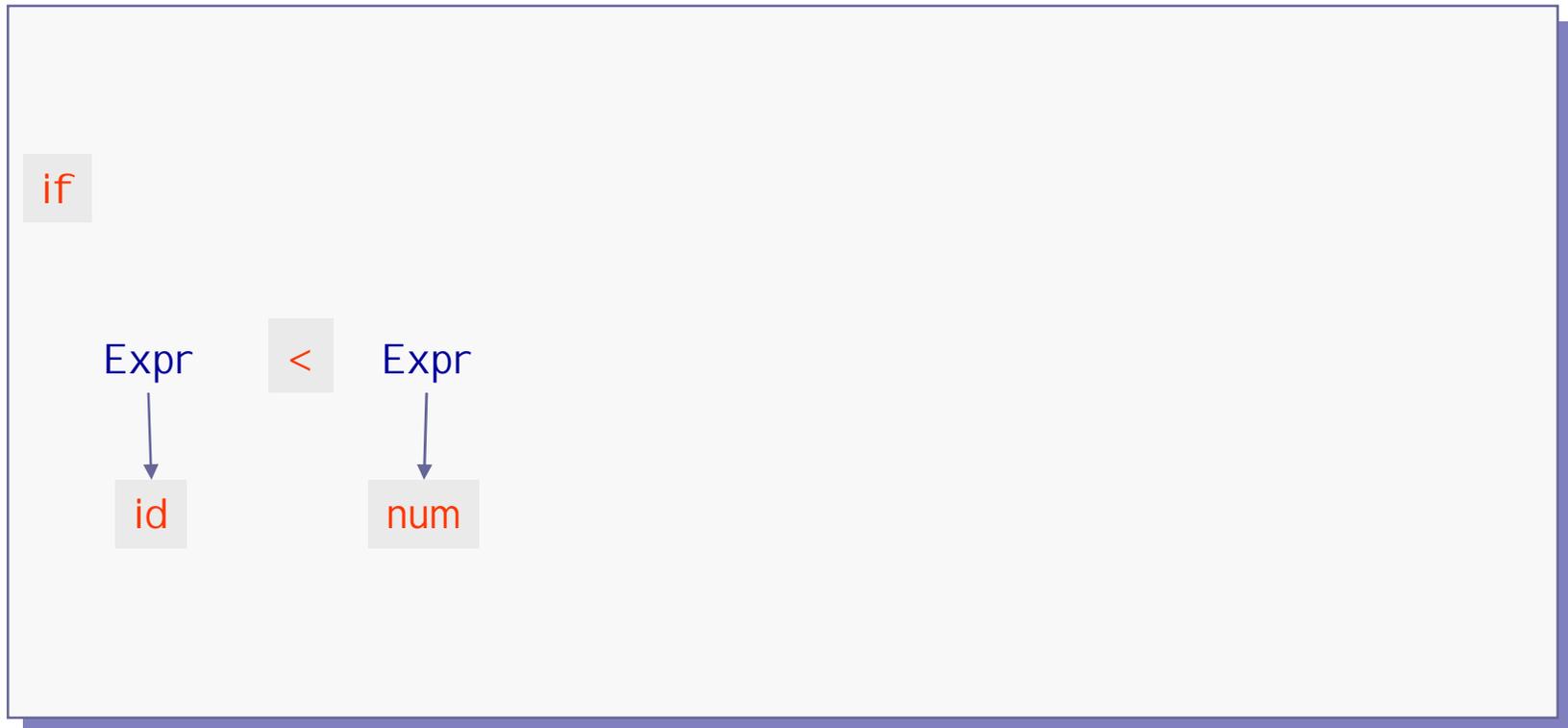
if	x	<	0	then	
----	---	---	---	------	--



Bottom Up



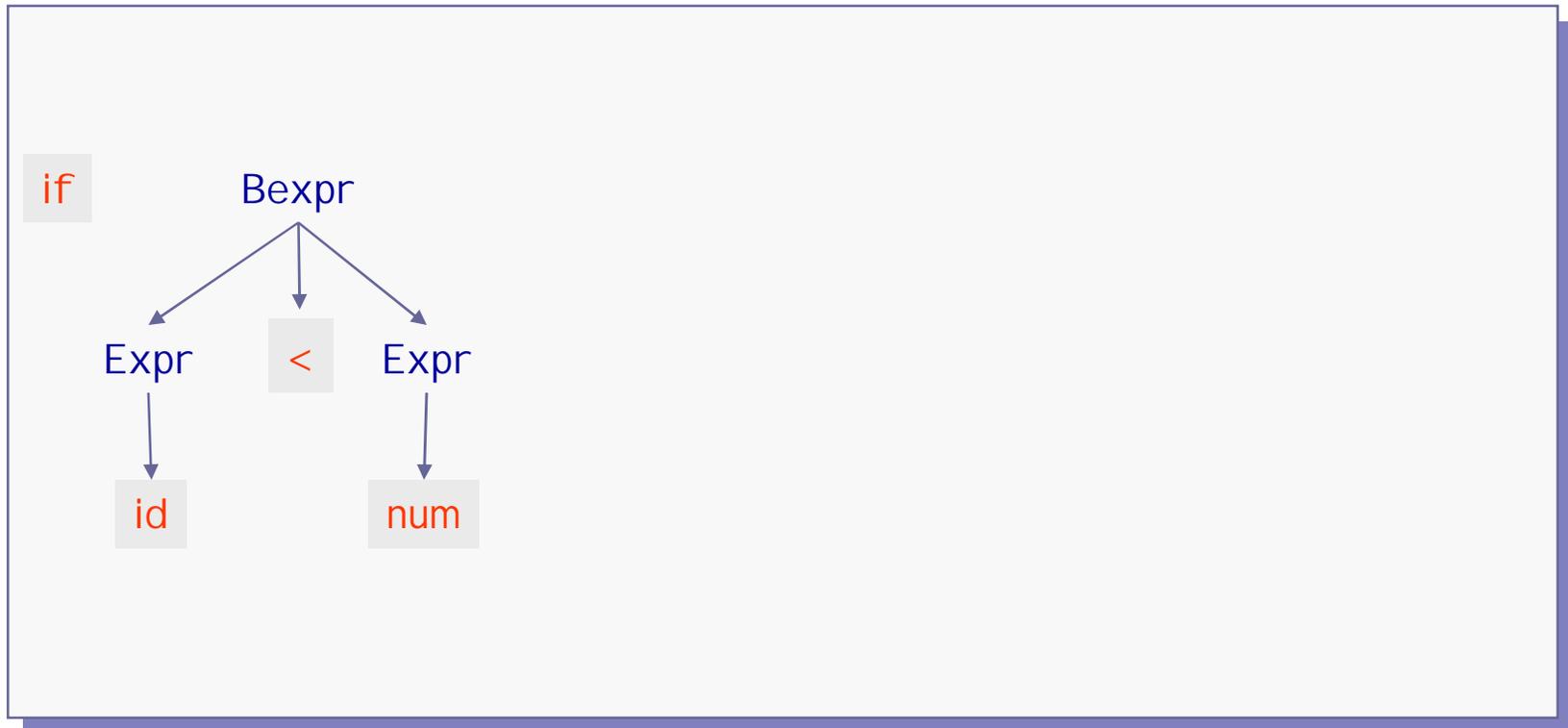
$Bexpr \rightarrow Expr < Expr$
 $Bexpr \rightarrow true$



if x < 0 then



Bottom Up



<code>if</code>	<code>x</code>	<code><</code>	<code>0</code>	<code>then</code>	
-----------------	----------------	-------------------	----------------	-------------------	--



Bottom Up



if	x	<	0	then	y	
----	---	---	---	------	---	--



Bottom Up



<code>if</code> <code>x</code> <code><</code> <code>0</code> <code>then</code> <code>y</code> <code>:=</code>	
--	--



Bottom Up



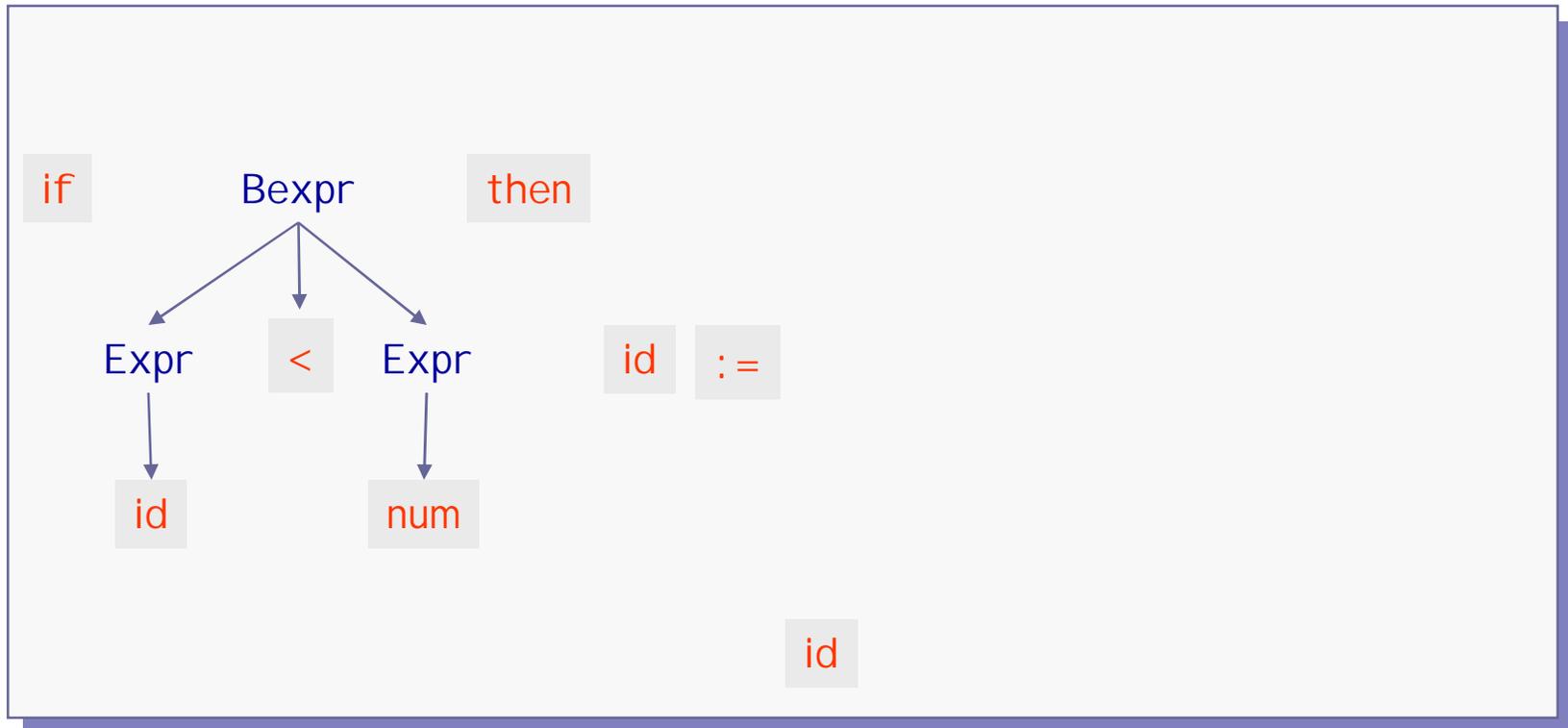
if x < 0 then y := y



Bottom Up



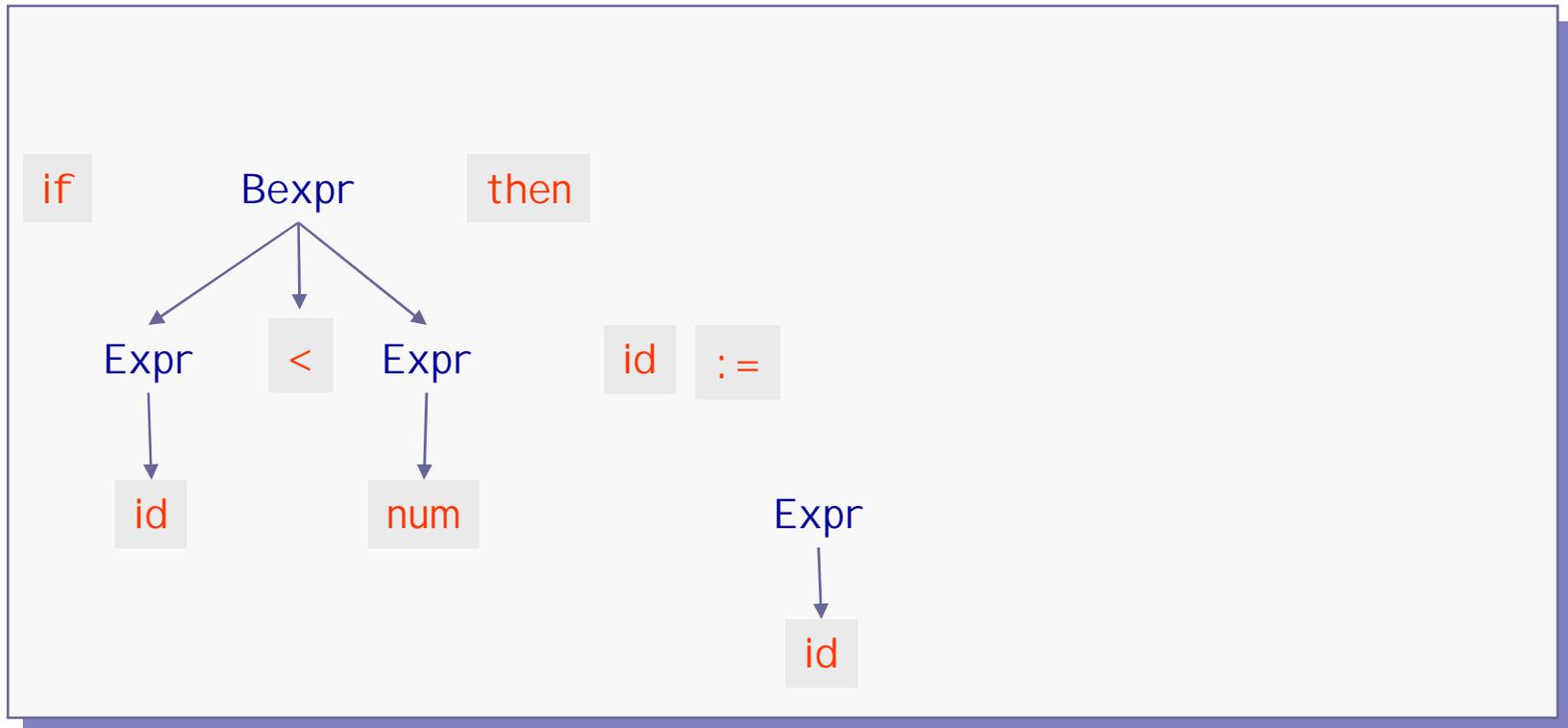
Expr \rightarrow Expr + Expr
Expr \rightarrow id
Expr \rightarrow num



if x < 0 then y := y +



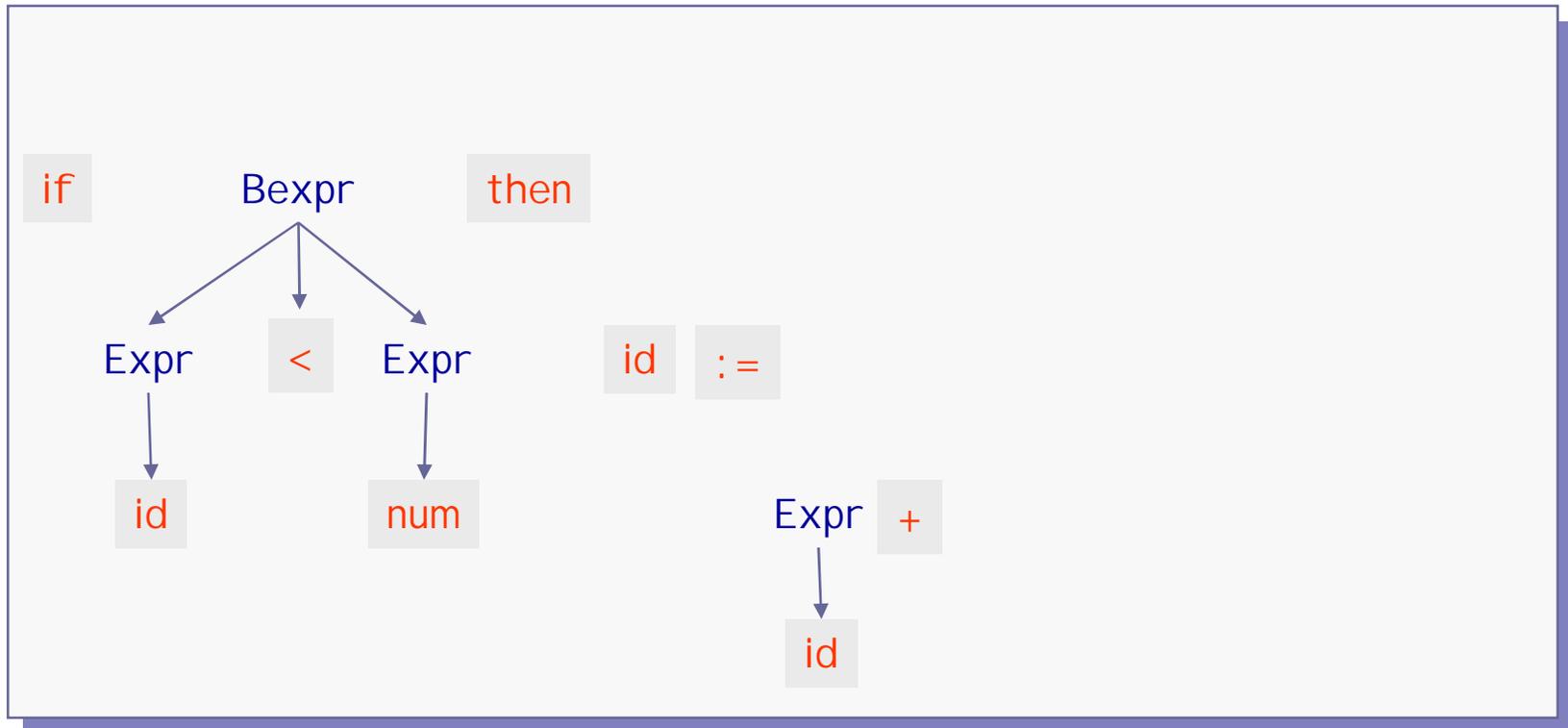
Bottom Up



if x < 0 then y := y +



Bottom Up



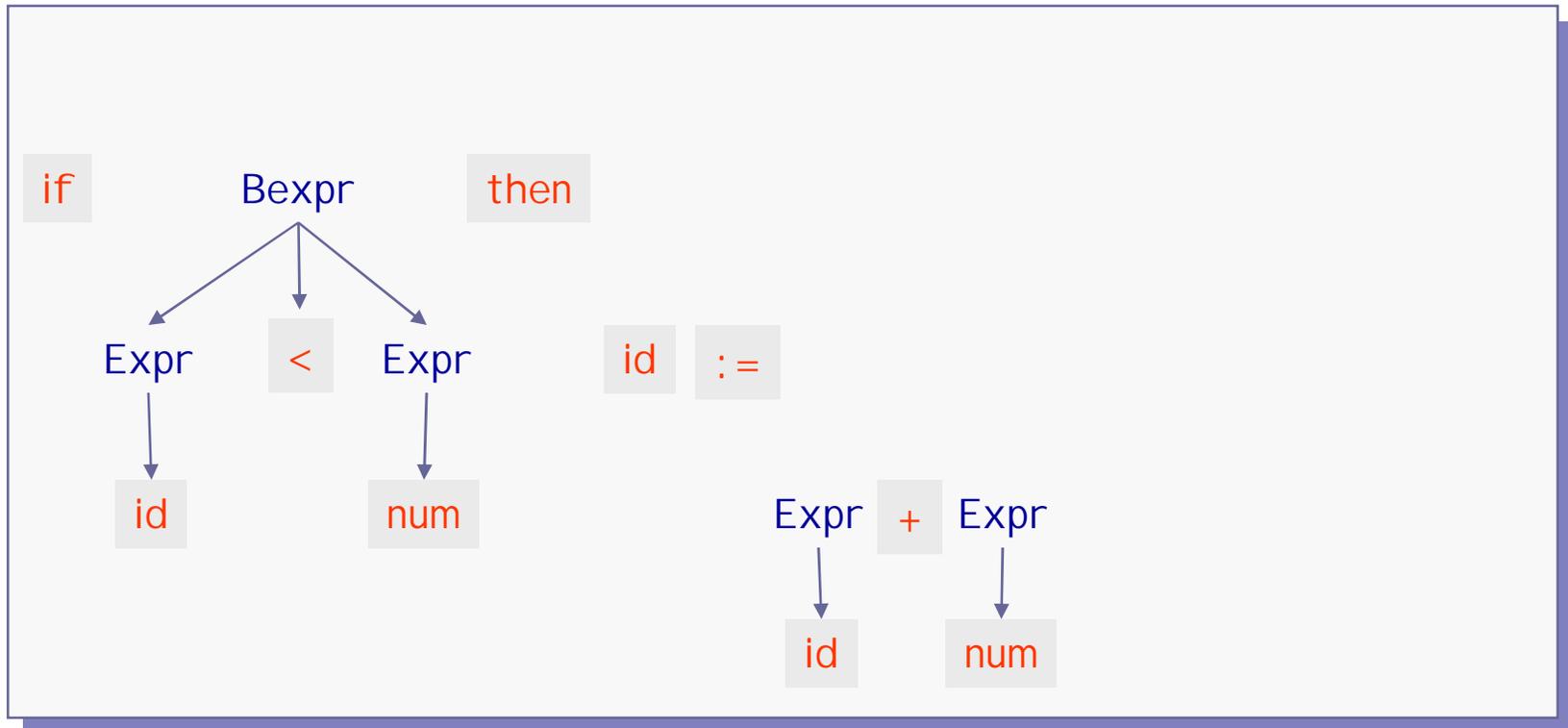
if x < 0 then y := y + 1



Bottom Up



Expr \rightarrow Expr + Expr
Expr \rightarrow id
Expr \rightarrow num



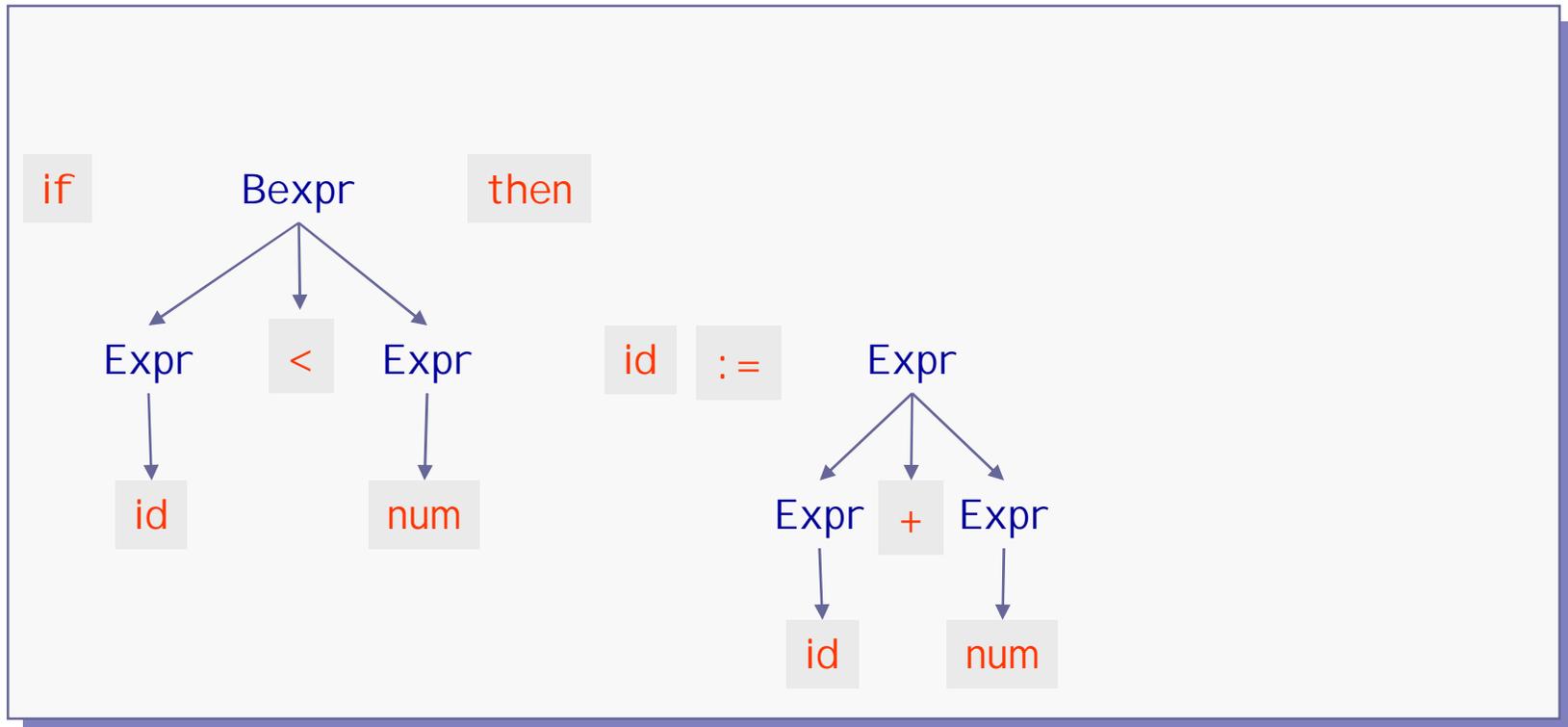
if	x	<	0	then	y	:=	y	+	1	else	
----	---	---	---	------	---	----	---	---	---	------	--



Bottom Up



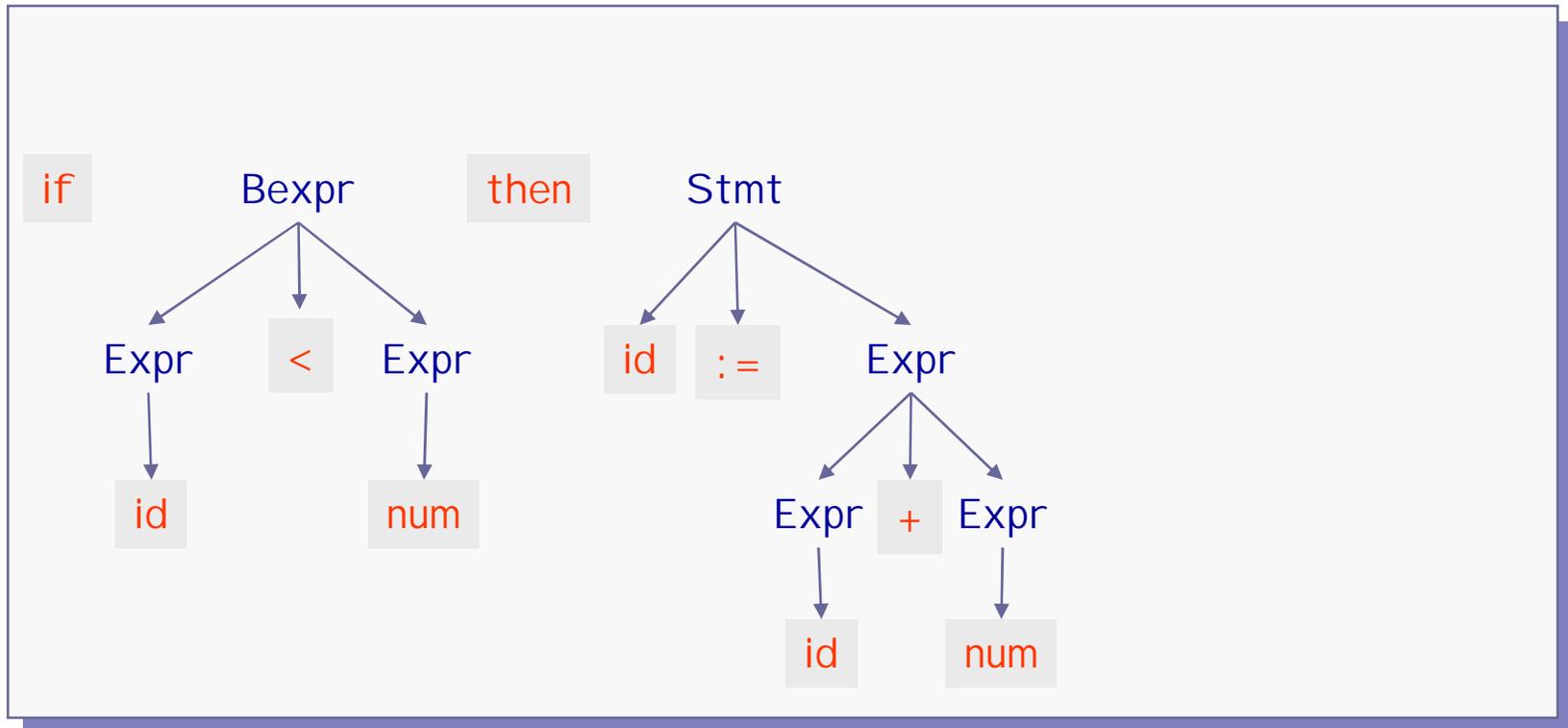
Stmt \rightarrow if Bexpr then Stmt else Stmt
Stmt \rightarrow id := Expr
Stmt \rightarrow id (Expr)



if	x	<	0	then	y	:=	y	+	1	else	
----	---	---	---	------	---	----	---	---	---	------	--



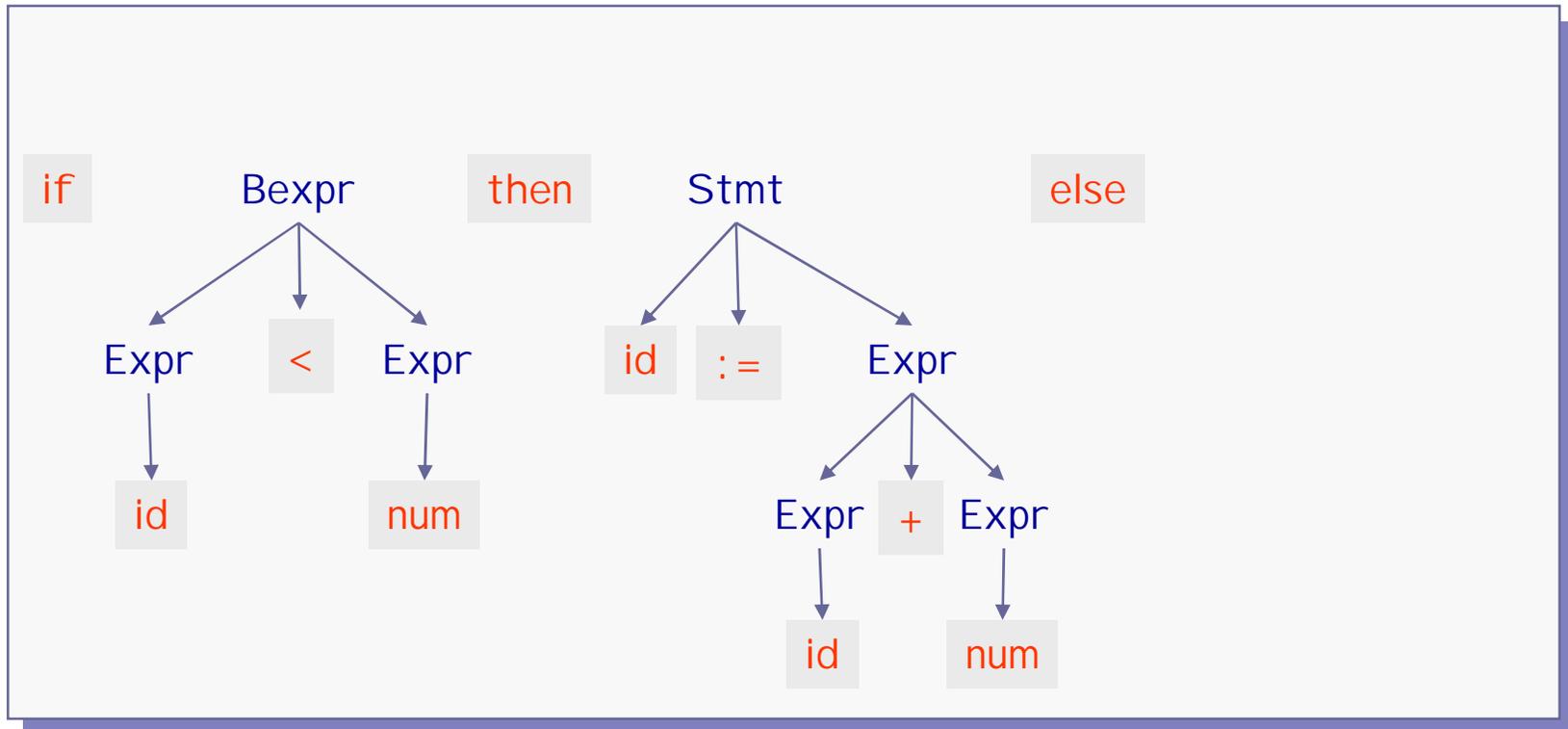
Bottom Up



if	x	<	0	then	y	:=	y	+	1	else	
-----------	----------	-------------	----------	-------------	----------	-----------	----------	----------	----------	-------------	--



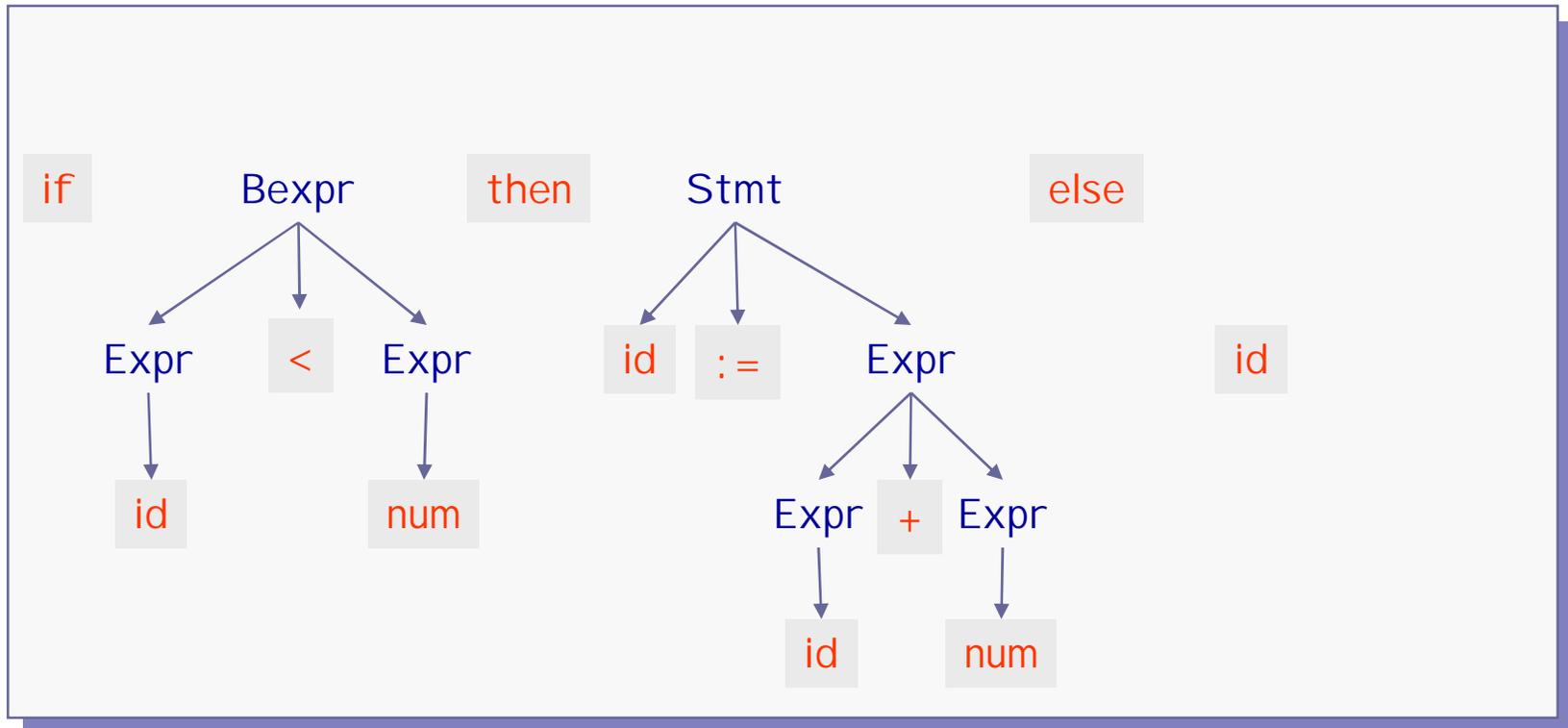
Bottom Up



if	x	<	0	then	y	:=	y	+	1	else	Inc	
----	---	---	---	------	---	----	---	---	---	------	-----	--



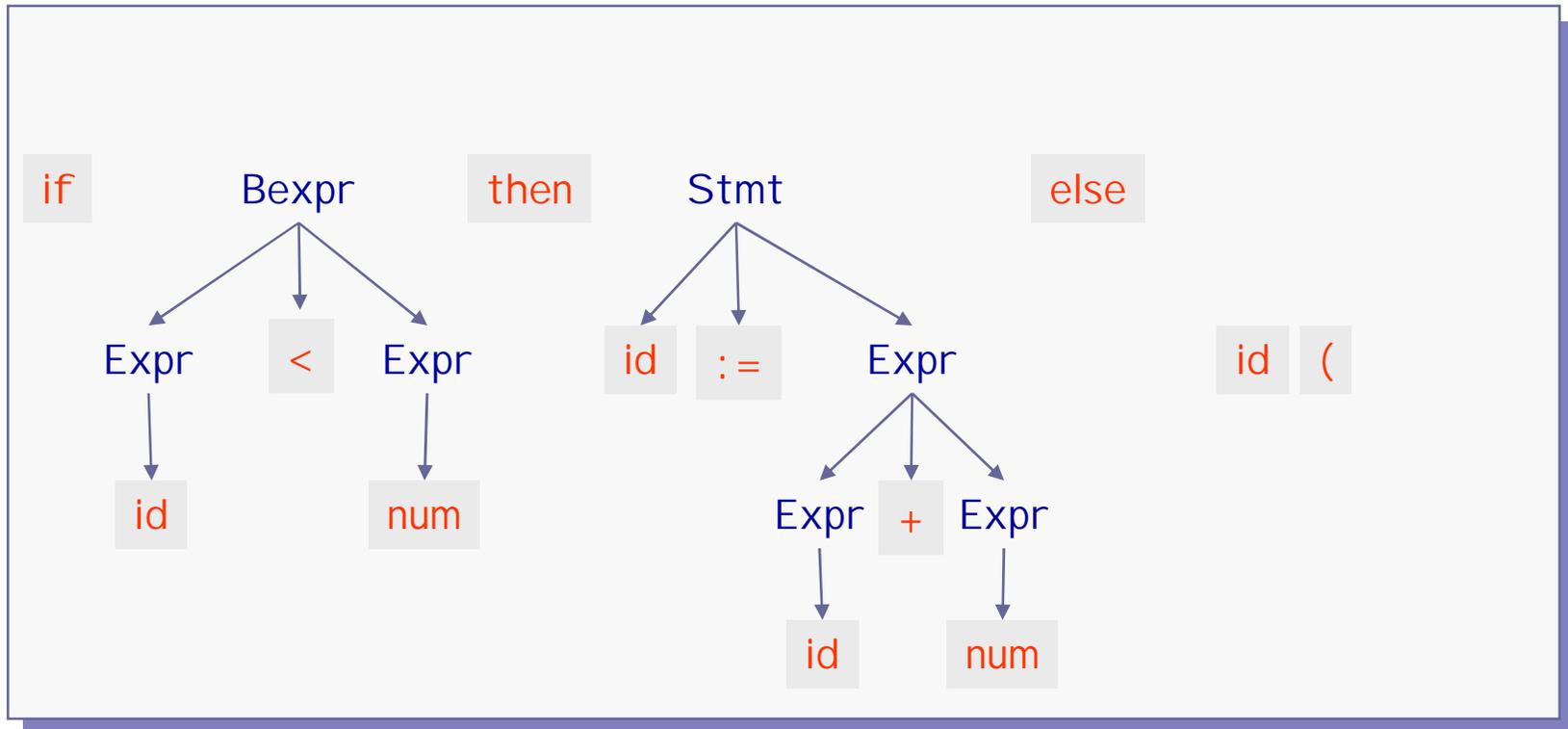
Bottom Up



<code>if</code>	<code>x</code>	<code><</code>	<code>0</code>	<code>then</code>	<code>y</code>	<code>:=</code>	<code>y</code>	<code>+</code>	<code>1</code>	<code>else</code>	<code>Inc</code>	<code>(</code>	
-----------------	----------------	-------------------	----------------	-------------------	----------------	-----------------	----------------	----------------	----------------	-------------------	------------------	----------------	--



Bottom Up



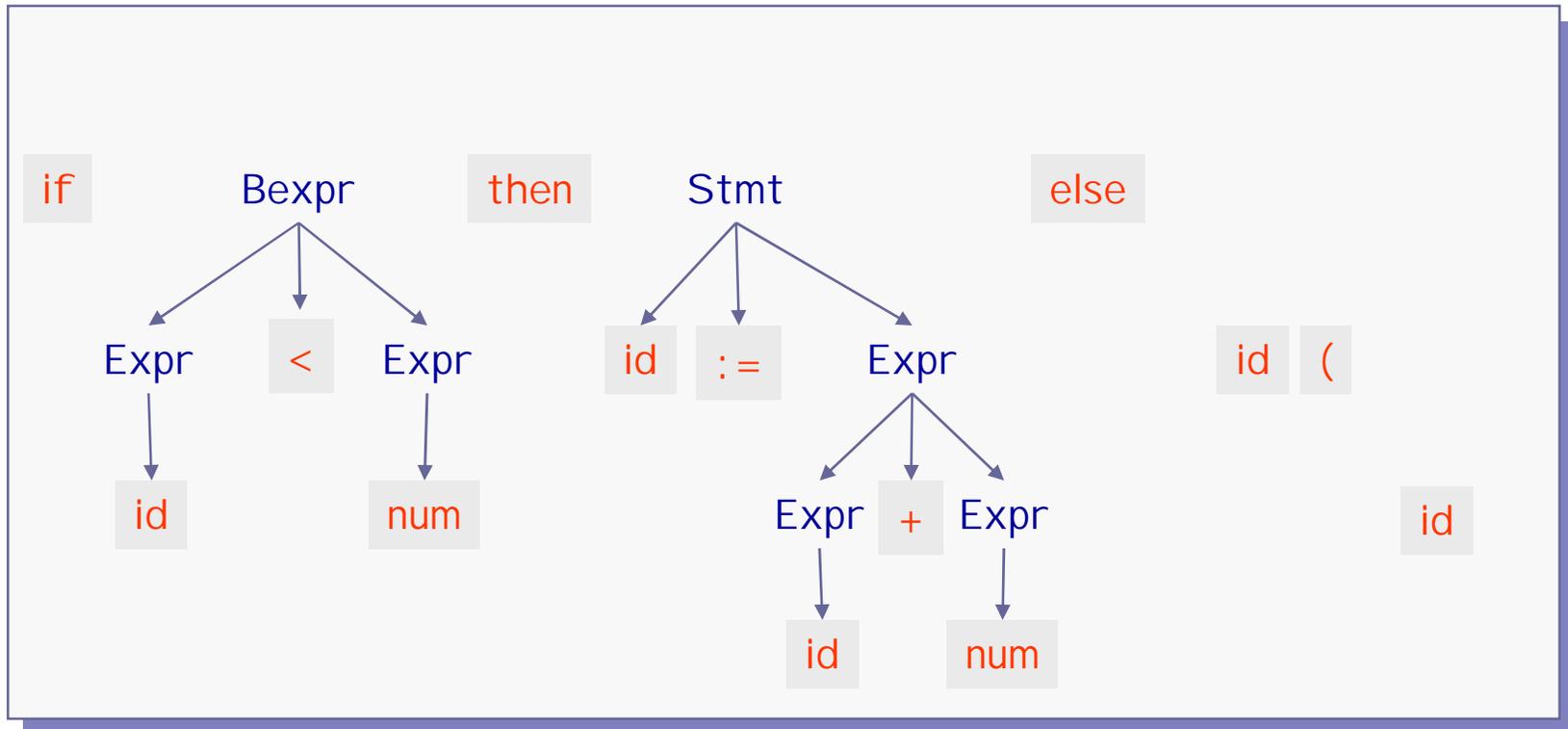
if	x	<	0	then	y	:=	y	+	1	else	Inc	(z	
----	---	---	---	------	---	----	---	---	---	------	-----	---	---	--



Bottom Up



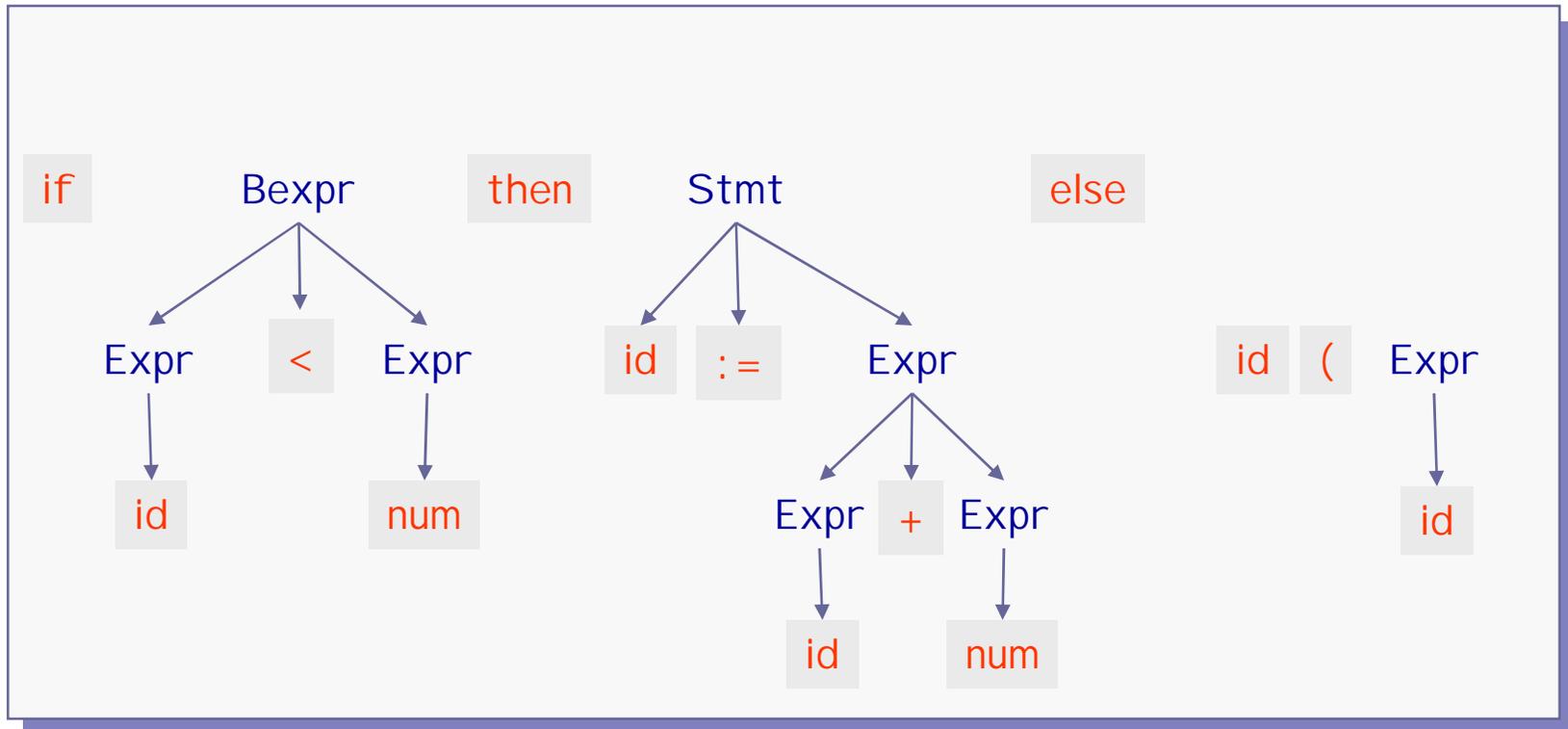
Expr \rightarrow Expr + Expr
Expr \rightarrow id
Expr \rightarrow num



if x < 0 then y := y + 1 else Inc (z)



Bottom Up



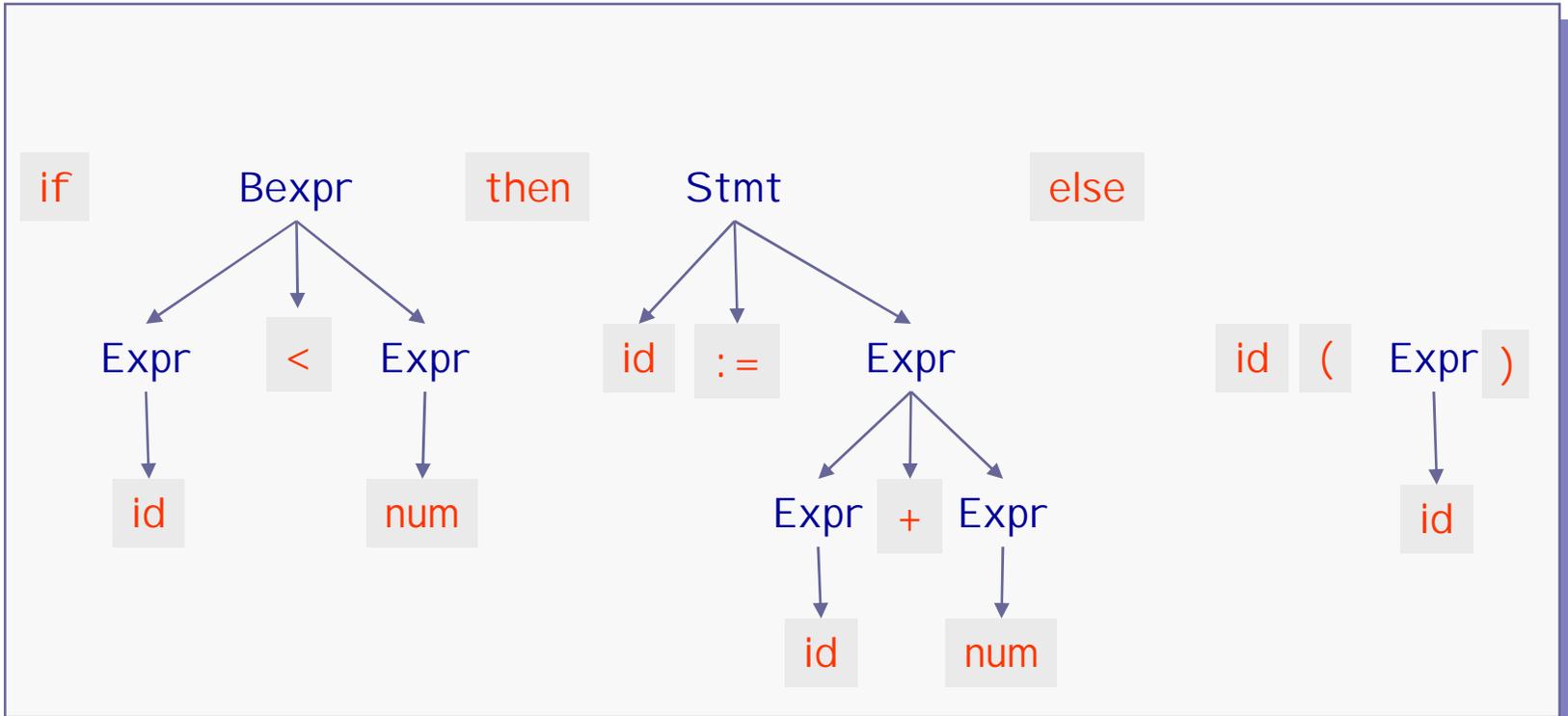
if x < 0 then y := y + 1 else Inc (z)



Bottom Up



Stmt \rightarrow if Bexpr then Stmt else Stmt
Stmt \rightarrow id := Expr
Stmt \rightarrow id (Expr)



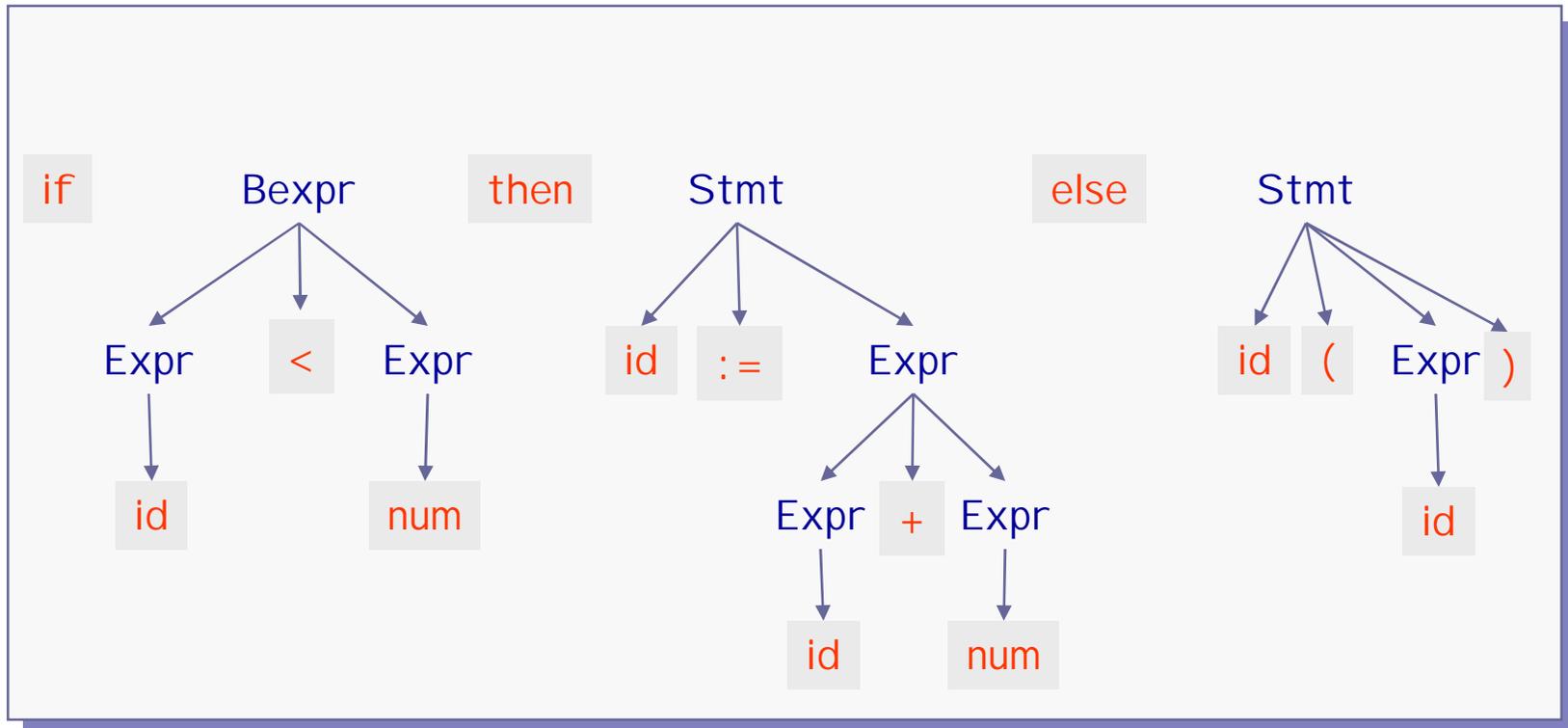
if x < 0 then y := y + 1 else Inc (z)



Bottom Up



Stmt \rightarrow if Bexpr then Stmt else Stmt
Stmt \rightarrow id := Expr
Stmt \rightarrow id (Expr)



if x < 0 then y := y + 1 else Inc (z)

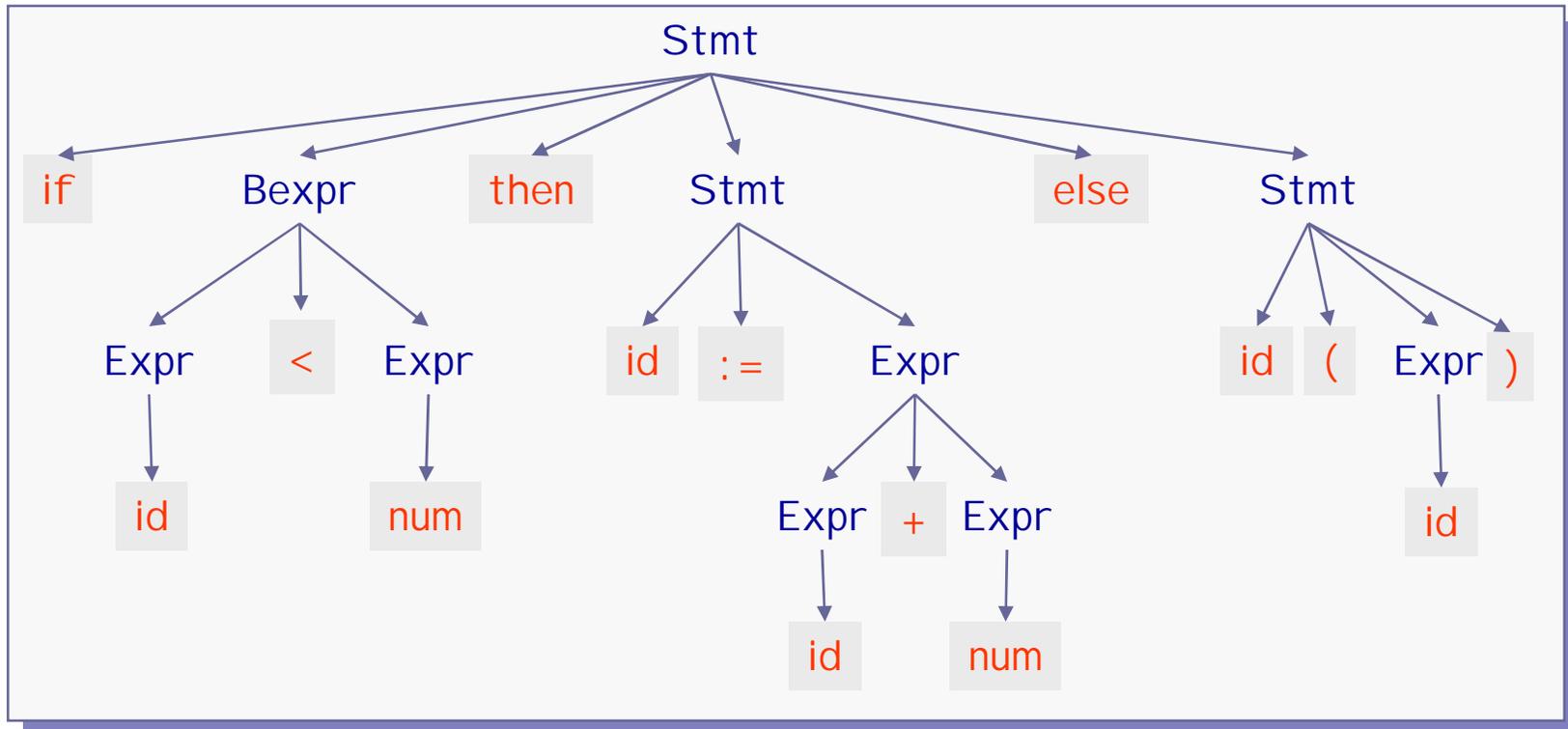


Bottom Up



Stmt \rightarrow if Bexpr then Stmt else Stmt
Stmt \rightarrow id := Expr
Stmt \rightarrow id (Expr)

Fertig



if x < 0 then y := y + 1 else Inc (z)



Top Down Strategie

n Top Down

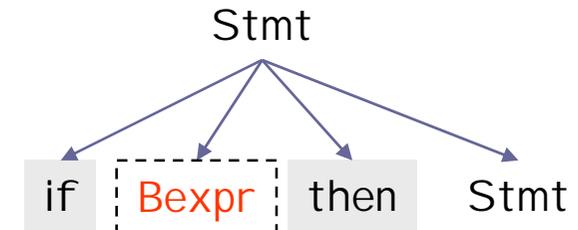
- .. Baum wächst von der Wurzel zu den Blättern
- .. jederzeit ein Baum
- .. Blätter: Terminale und Nonterminale

n Aktuelles Blatt

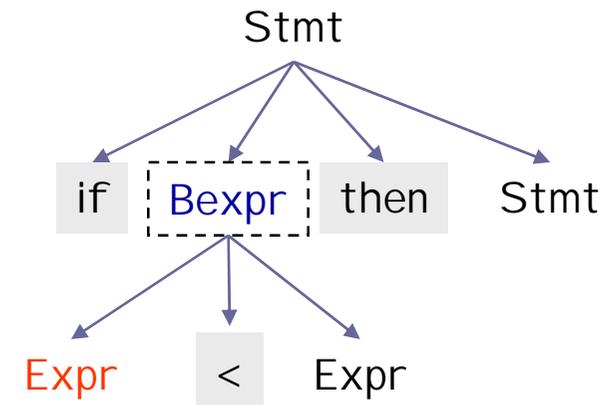
- .. gibt an, was im Input erwartet wird
- .. alle Blätter links davon sind schon terminal

n Falls aktuelles Blatt B

- .. terminal
 - n Lese nächstes Token im Input
 - n prüfe, dass es mit B übereinstimmt
 - n wenn nicht: „Error: B expected“
 - n rechter Nachbar wird aktuelles Blatt
- .. nonterminal
 - n Wähle Grammatikregel $B \rightarrow w_1 \dots w_n$
 - n Hänge neue Söhne $w_1 \dots w_n$ an
 - n w_1 wird neues aktuelles Blatt



$Bexpr \rightarrow Bexpr < Bexpr$
| true





Top Down

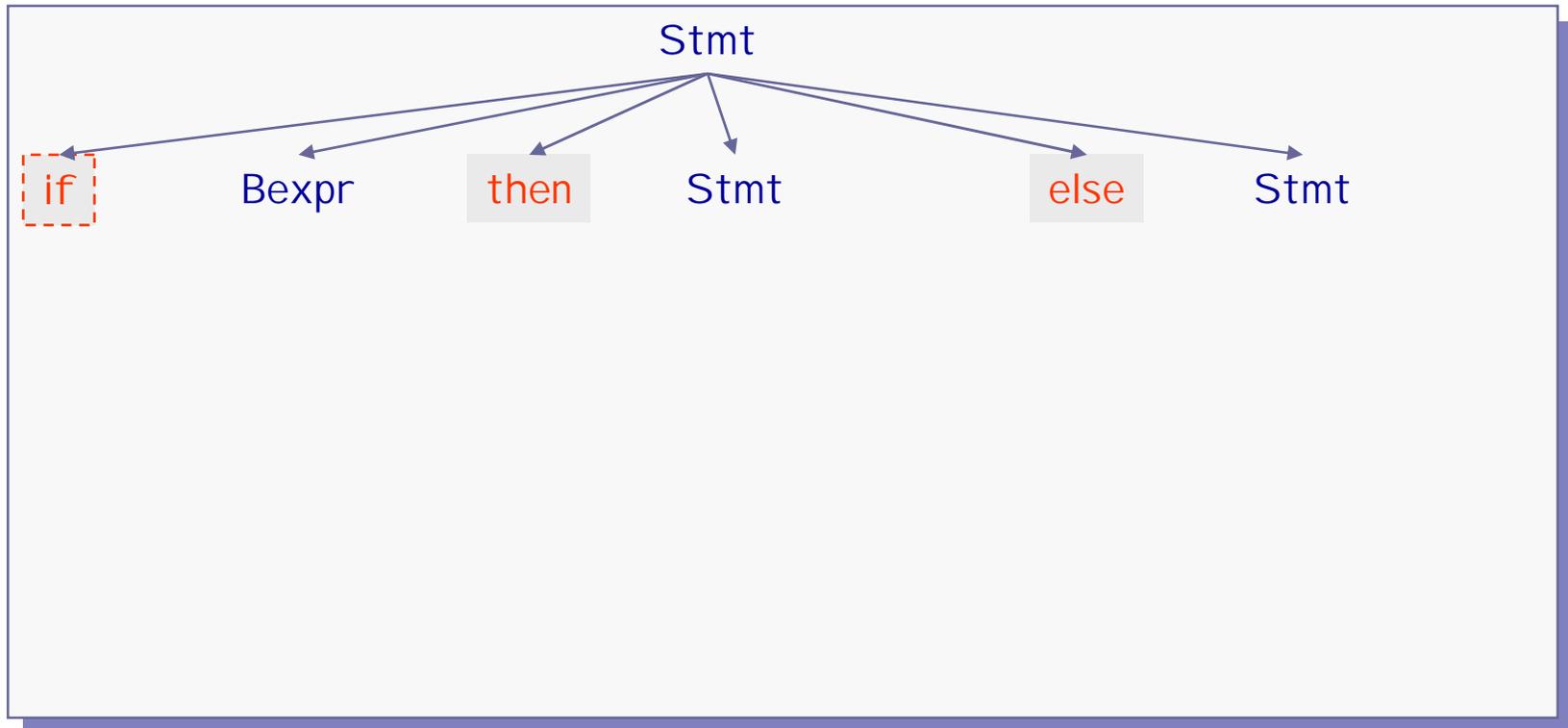
Stmt \rightarrow if Bexpr then Stmt else Stmt
Stmt \rightarrow id := Expr
Stmt \rightarrow id (Expr)

Stmt

if



Top Down

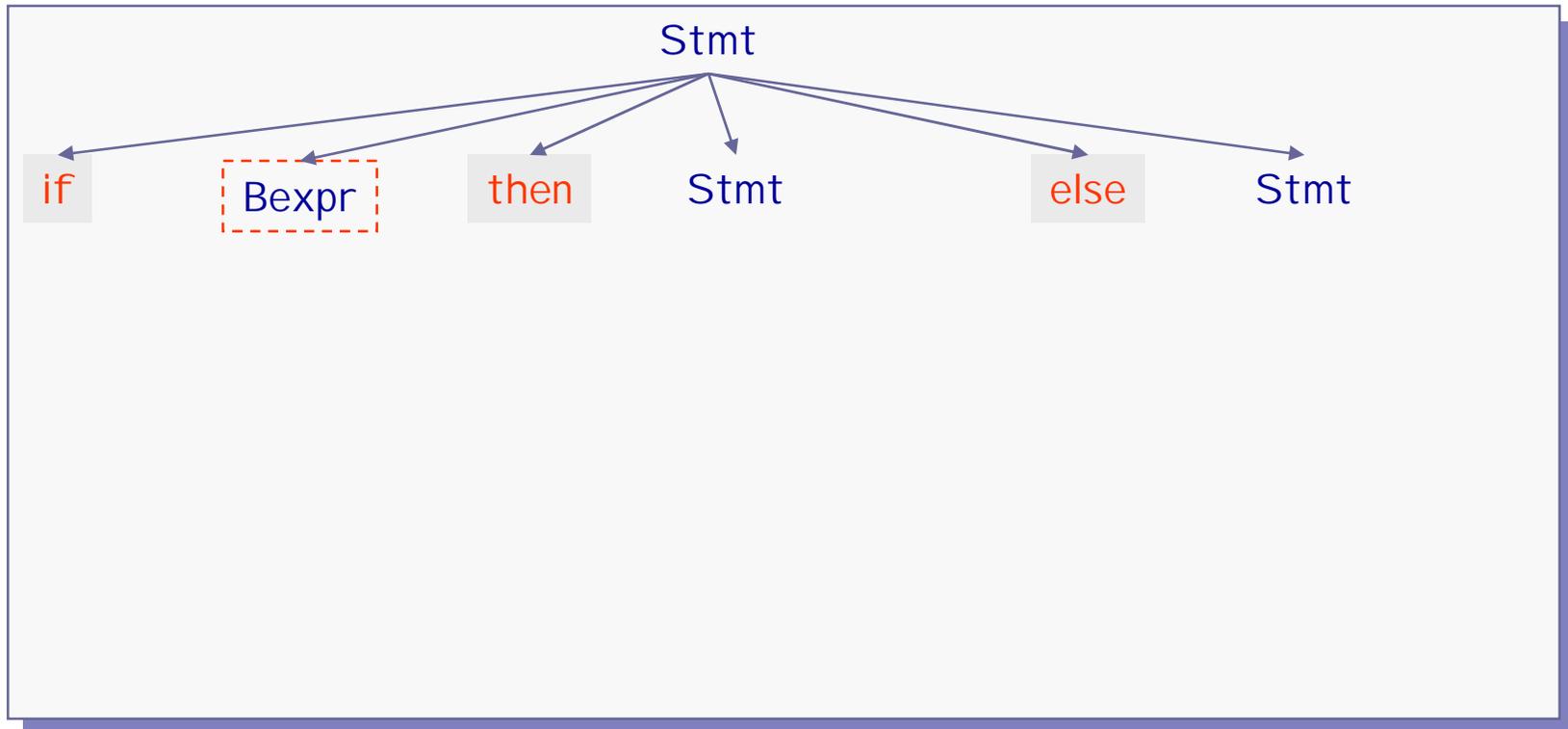


<code>if</code>	
-----------------	--



Top Down

$Bexpr \rightarrow Expr < Expr$
 $Bexpr \rightarrow true$



<code>if</code>	<code>x</code>	
-----------------	----------------	--

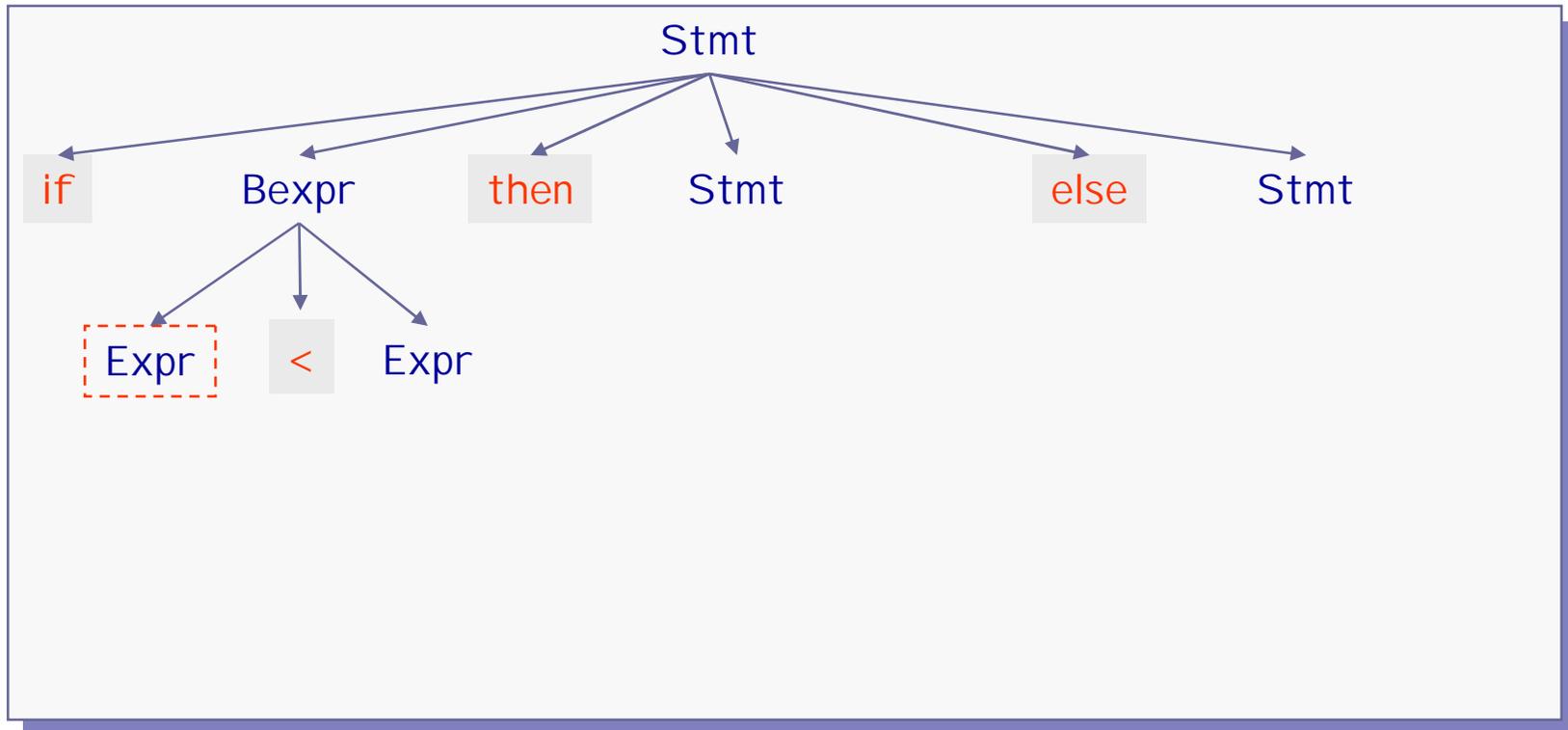


Top Down

Expr \rightarrow Expr + Expr

Expr \rightarrow id

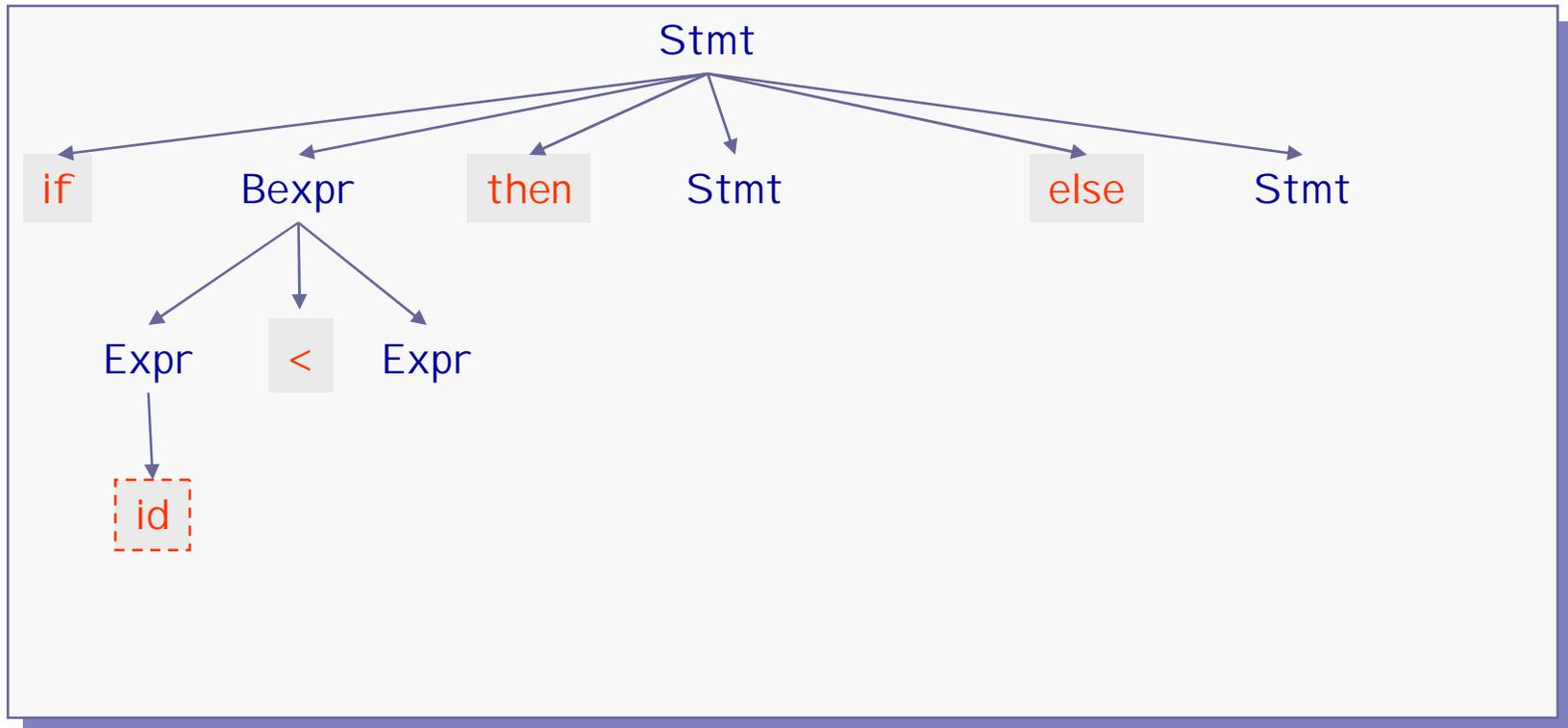
Expr \rightarrow num



if	x	
----	---	--



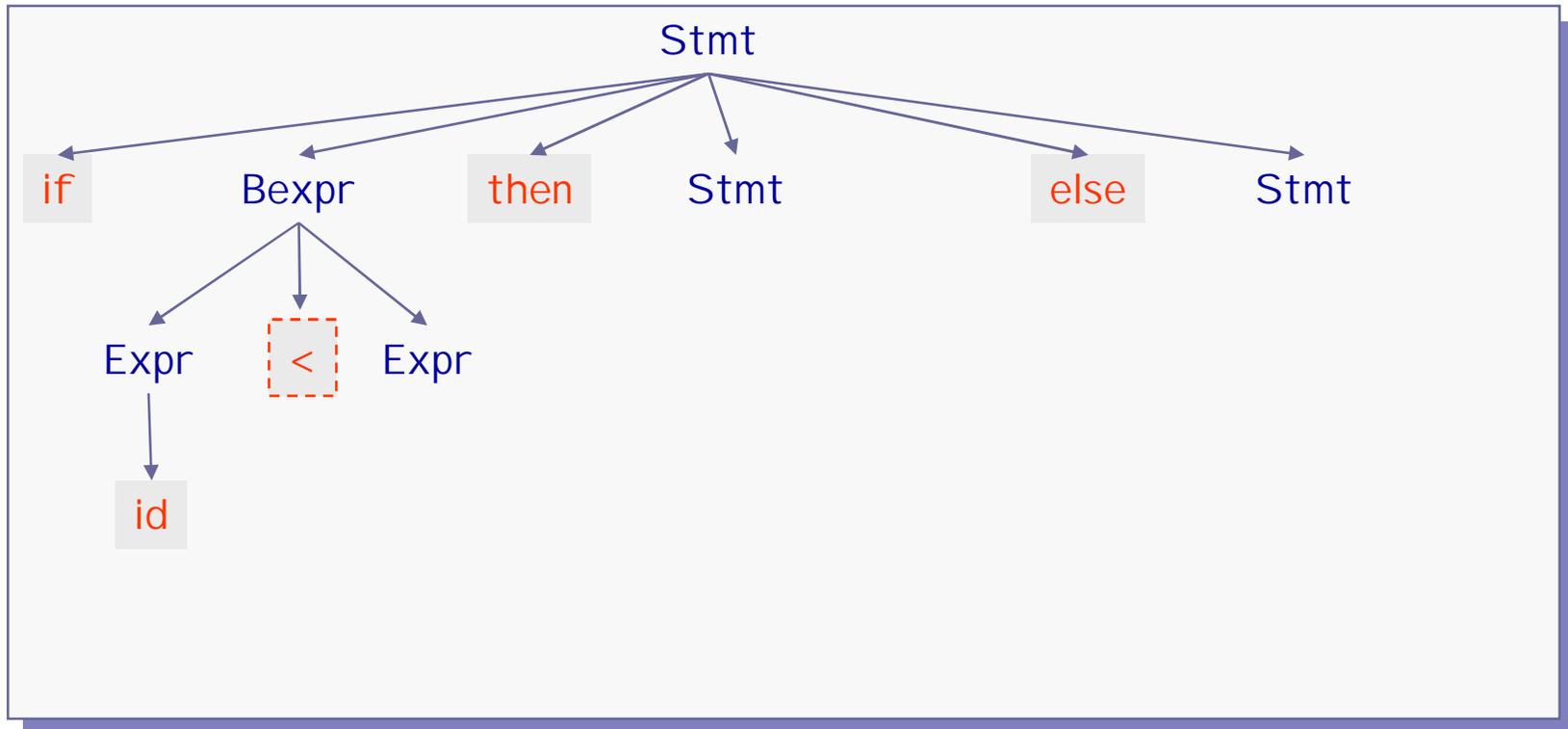
Top Down



if	x	
-----------	----------	--



Top Down

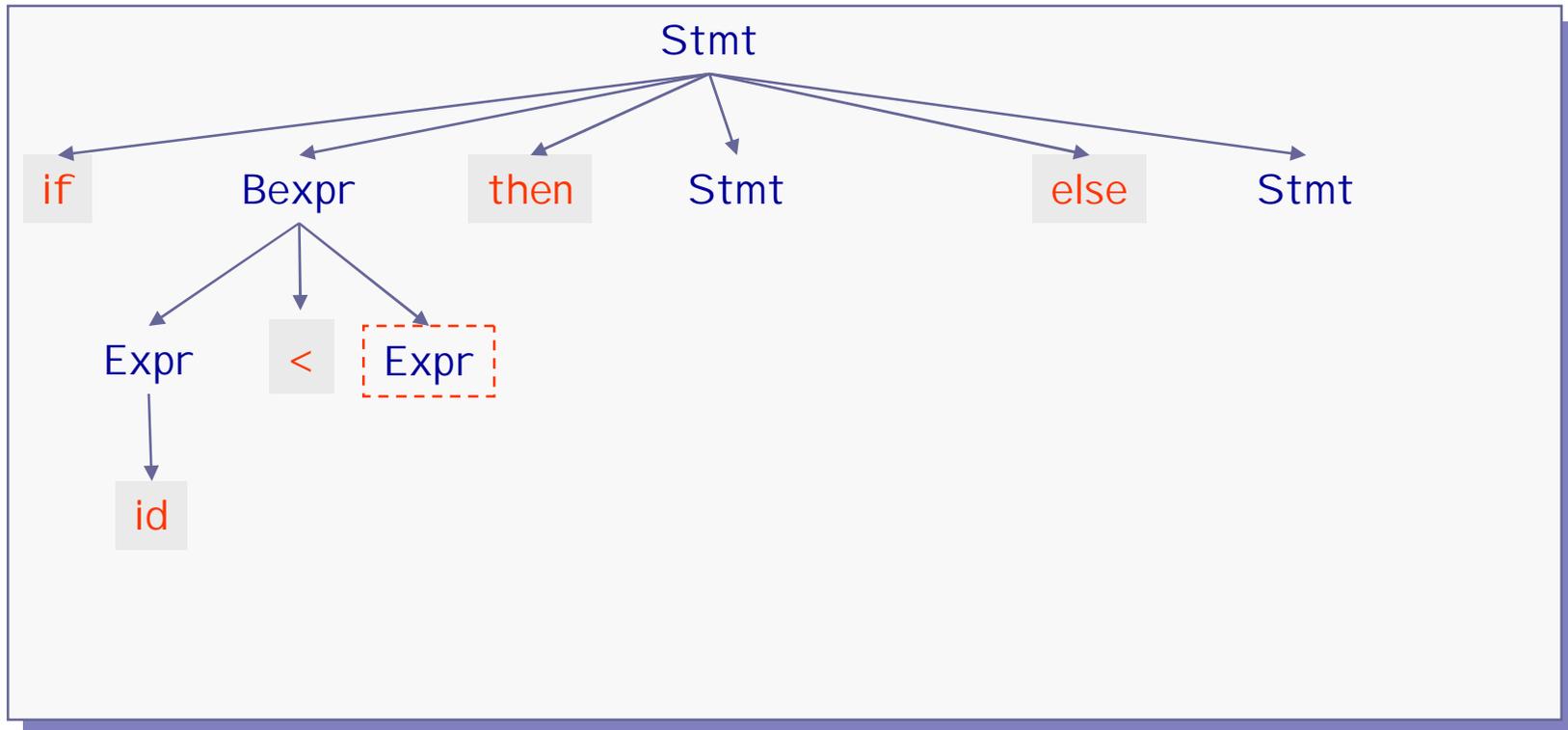


if	x	<	
----	---	---	--



Top Down

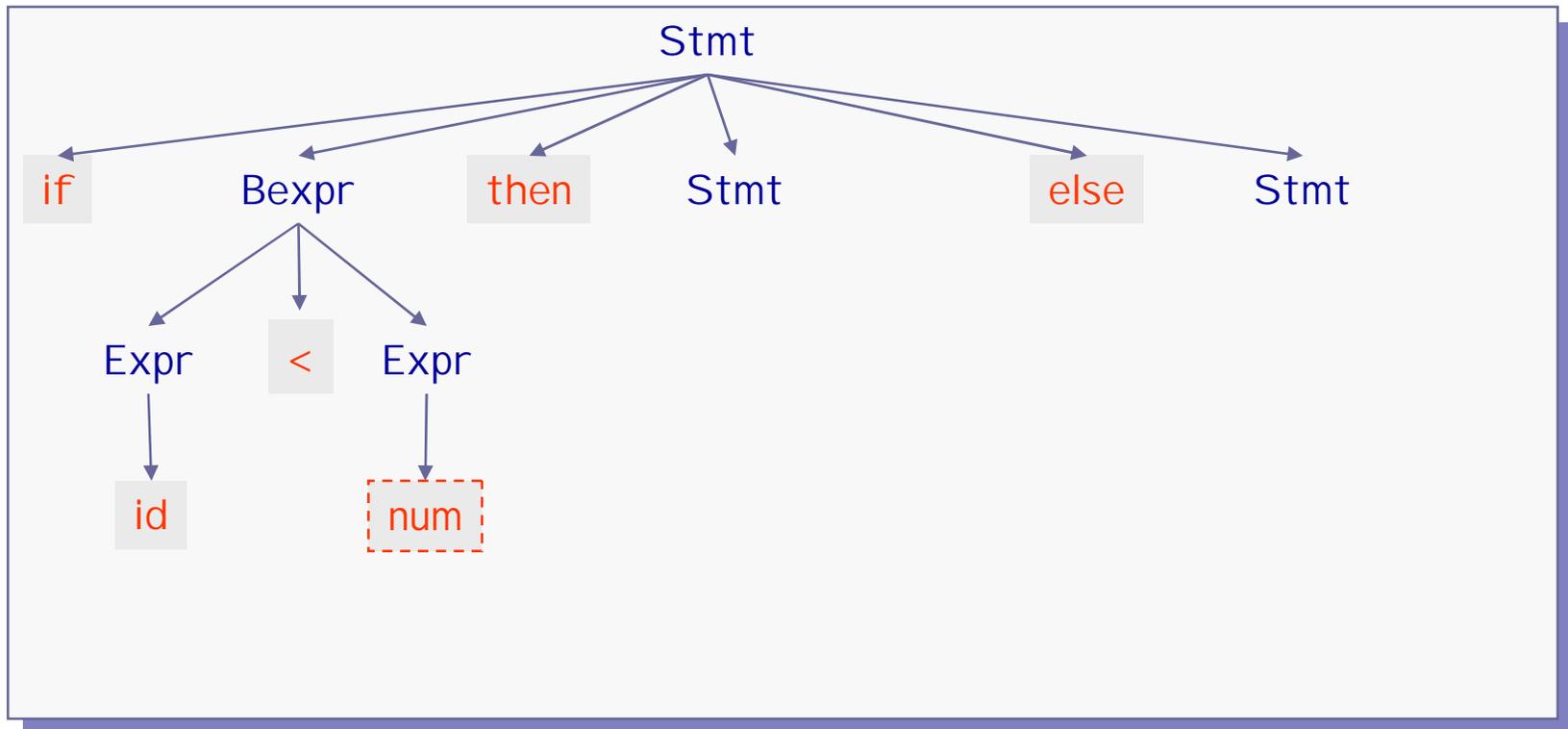
Expr \rightarrow Expr + Expr
Expr \rightarrow id
Expr \rightarrow num



if	x	<	0	
----	---	---	---	--



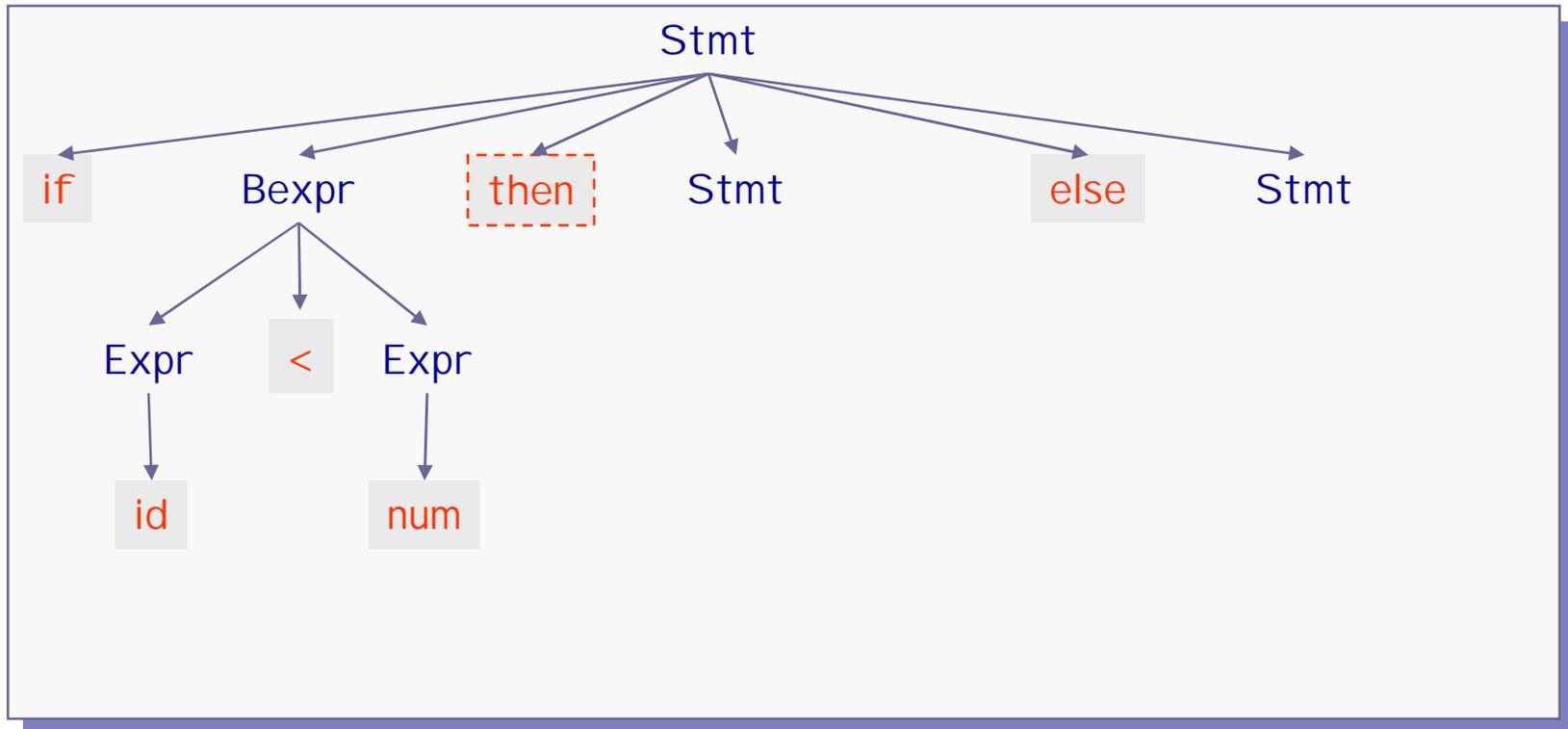
Top Down



<code>if</code>	<code>x</code>	<code><</code>	<code>0</code>	
-----------------	----------------	-------------------	----------------	--



Top Down

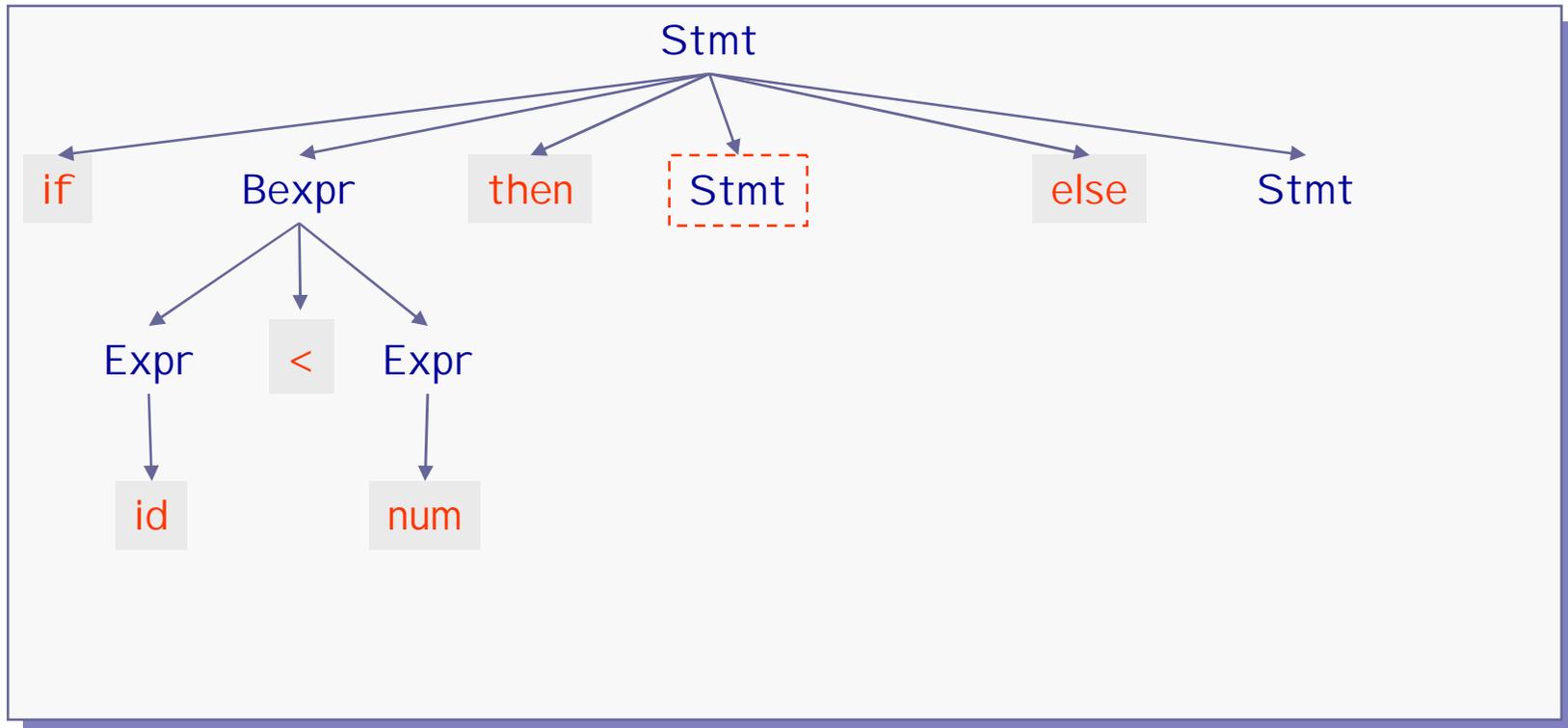


if	x	<	0	then	
----	---	---	---	------	--



Top Down

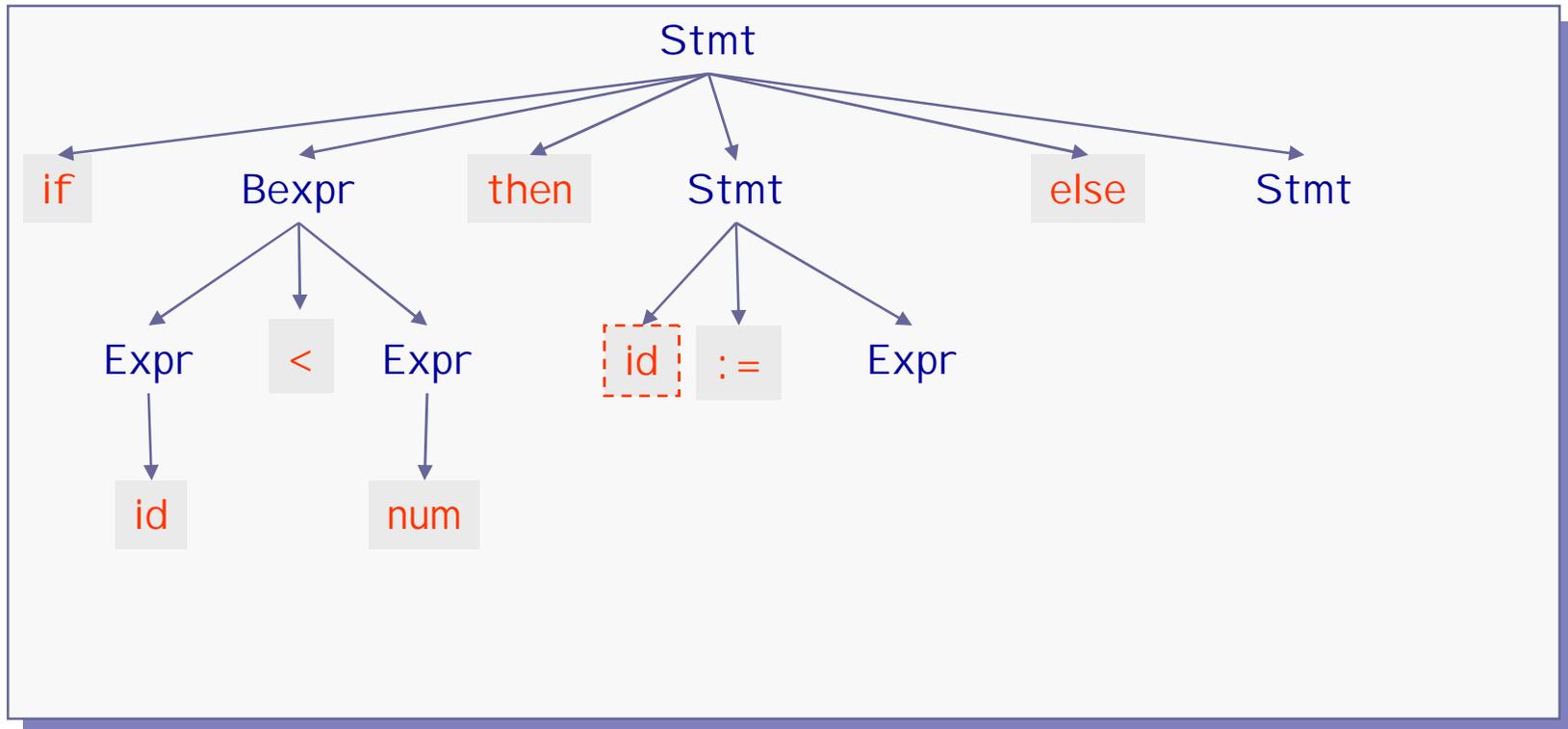
Stmt \rightarrow if Bexpr then Stmt else Stmt
Stmt \rightarrow id := Expr
Stmt \rightarrow id (Expr)



if	x	<	0	then	y	
----	---	---	---	------	---	--



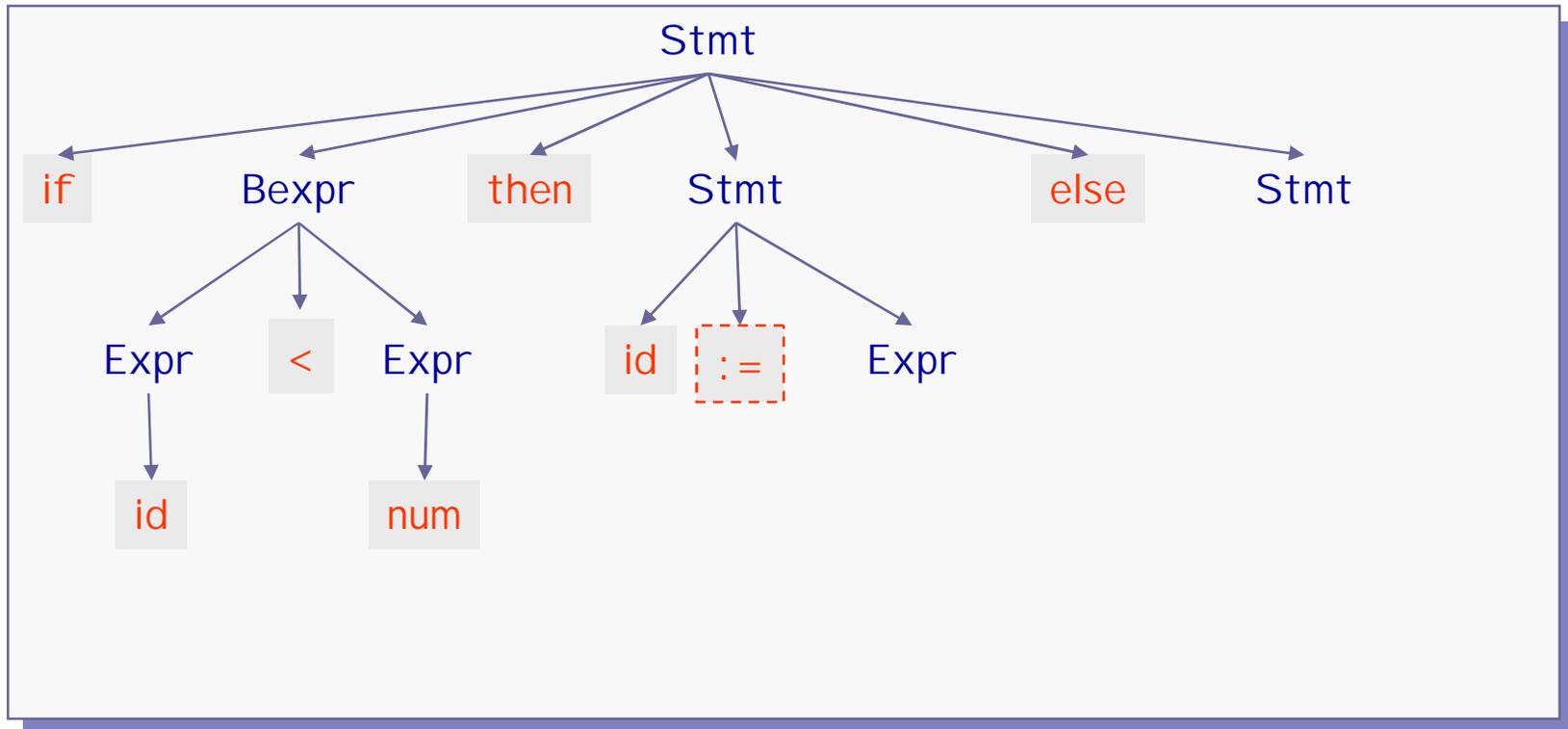
Top Down



<code>if</code>	<code>x</code>	<code><</code>	<code>0</code>	<code>then</code>	<code>y</code>	
-----------------	----------------	-------------------	----------------	-------------------	----------------	--



Top Down



if x < 0 then y :=

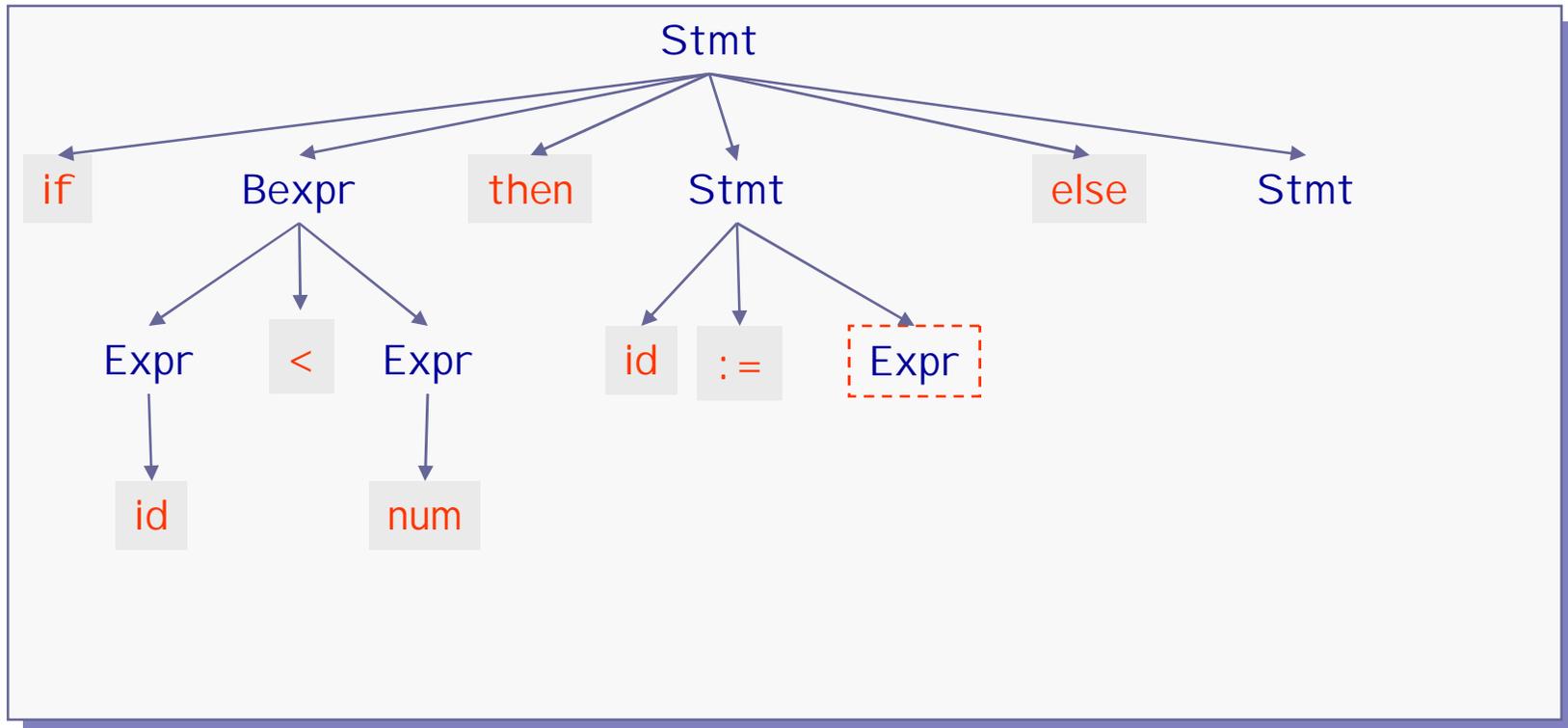


Top Down

Expr \rightarrow Expr + Expr

Expr \rightarrow id

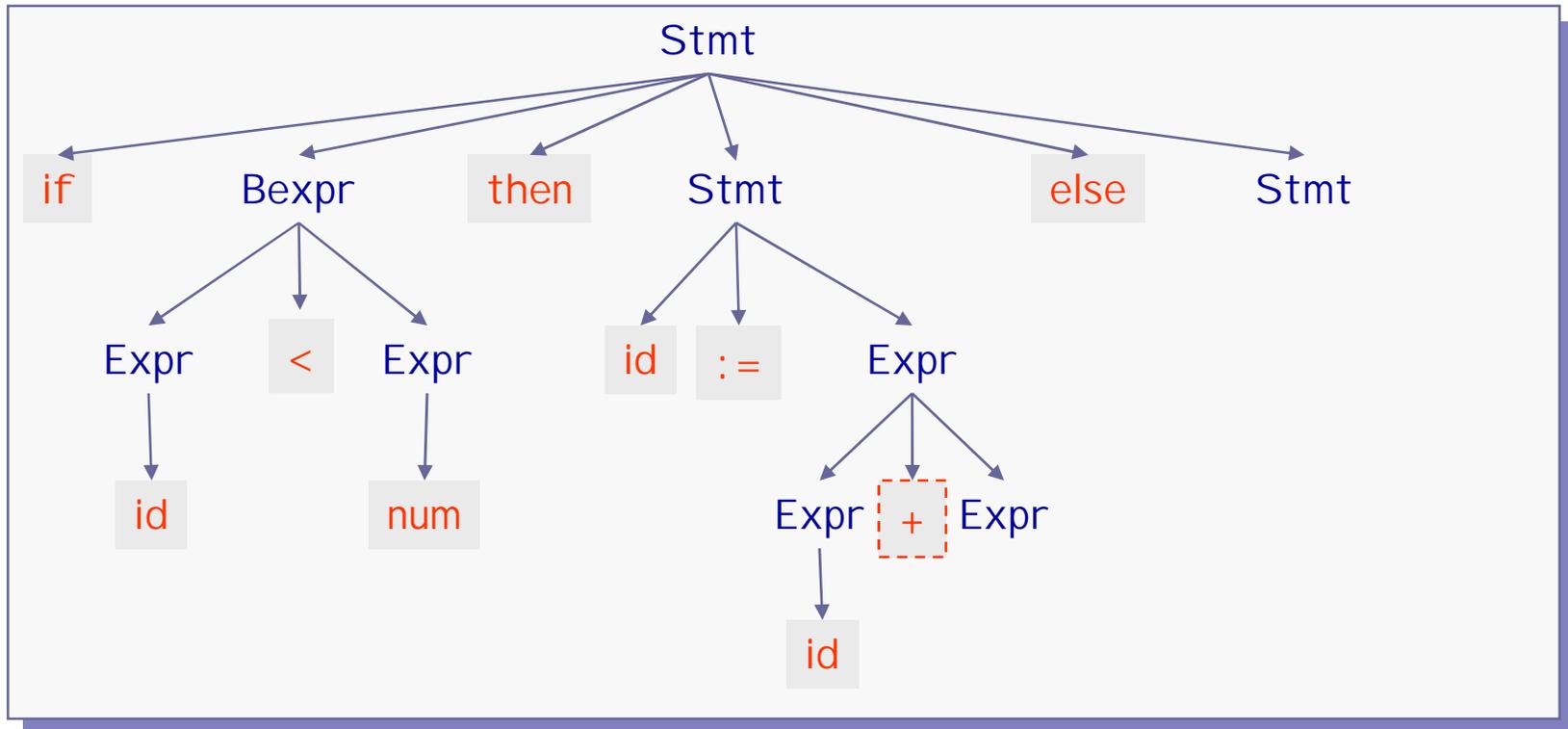
Expr \rightarrow num



if x < 0 then y := y



Top Down

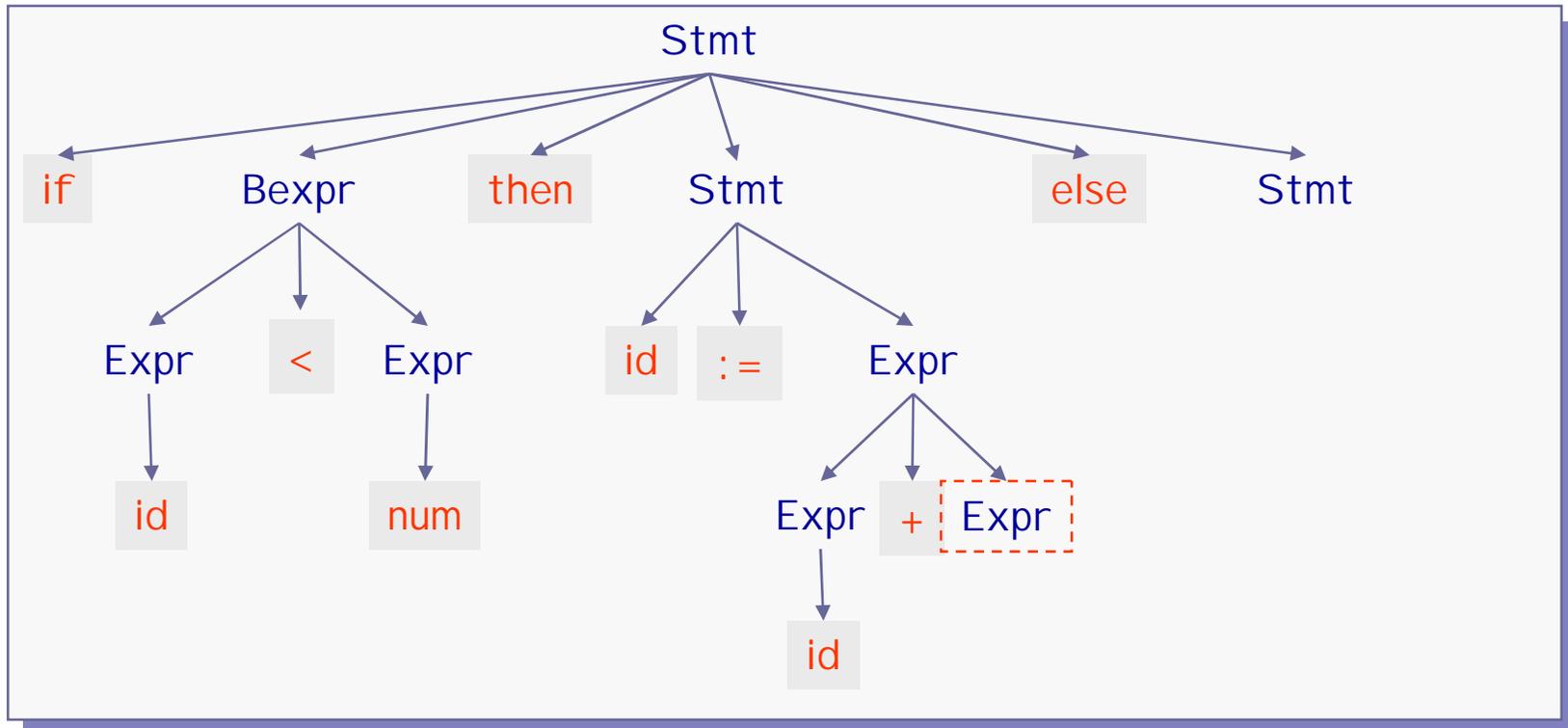


if x < 0 then y := y +



Top Down

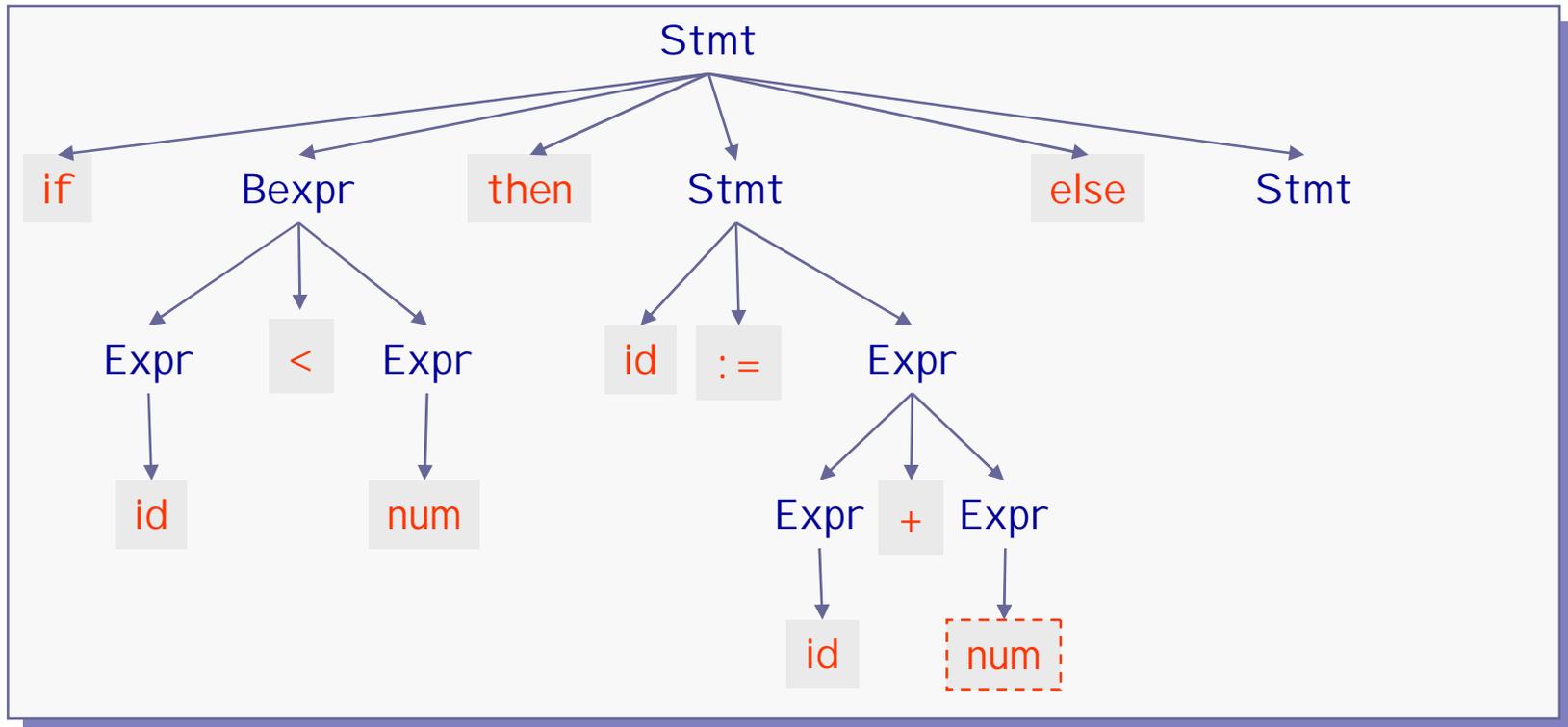
Expr \rightarrow Expr + Expr
Expr \rightarrow id
Expr \rightarrow num



if x < 0 then y := y + 1



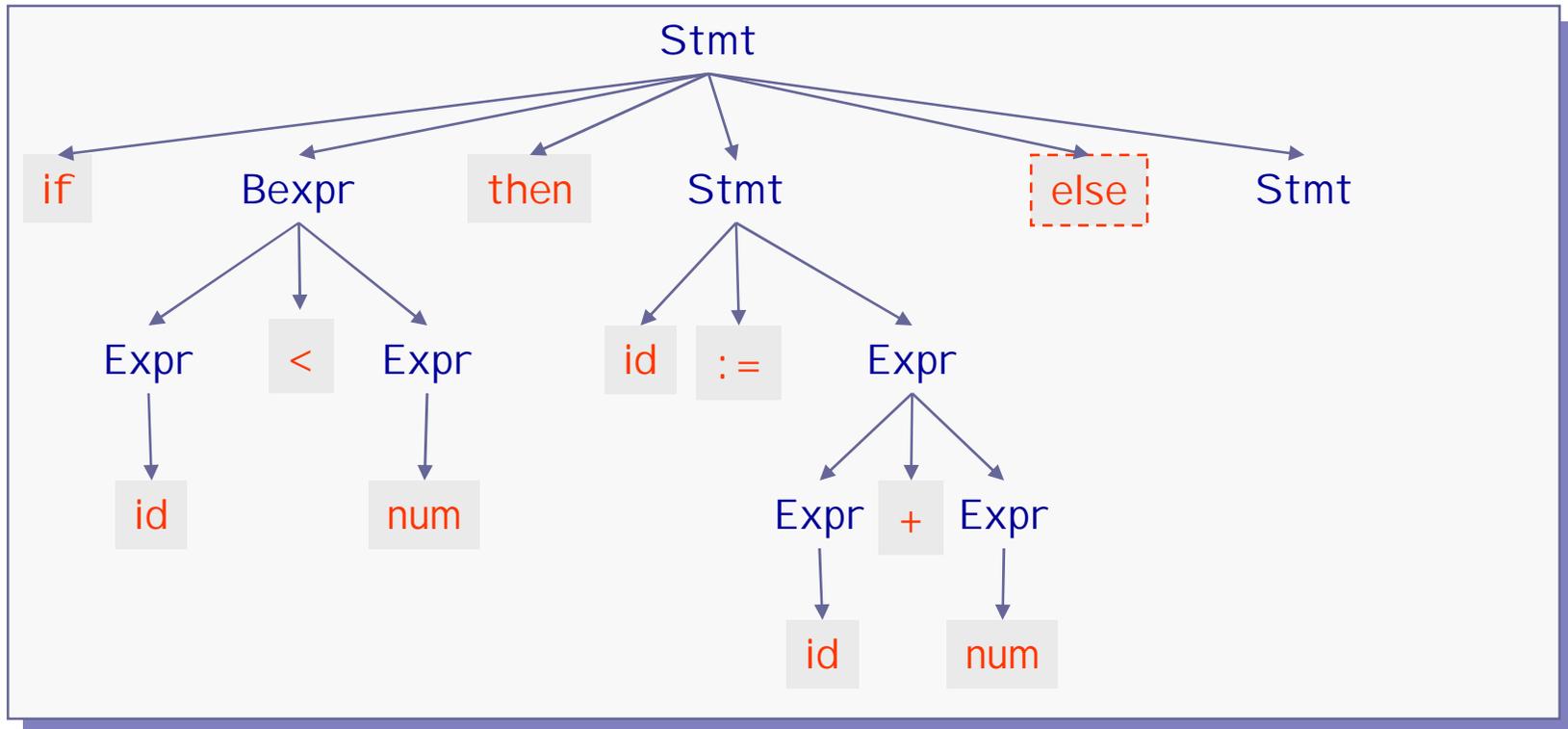
Top Down



if x < 0 then y := y + 1



Top Down

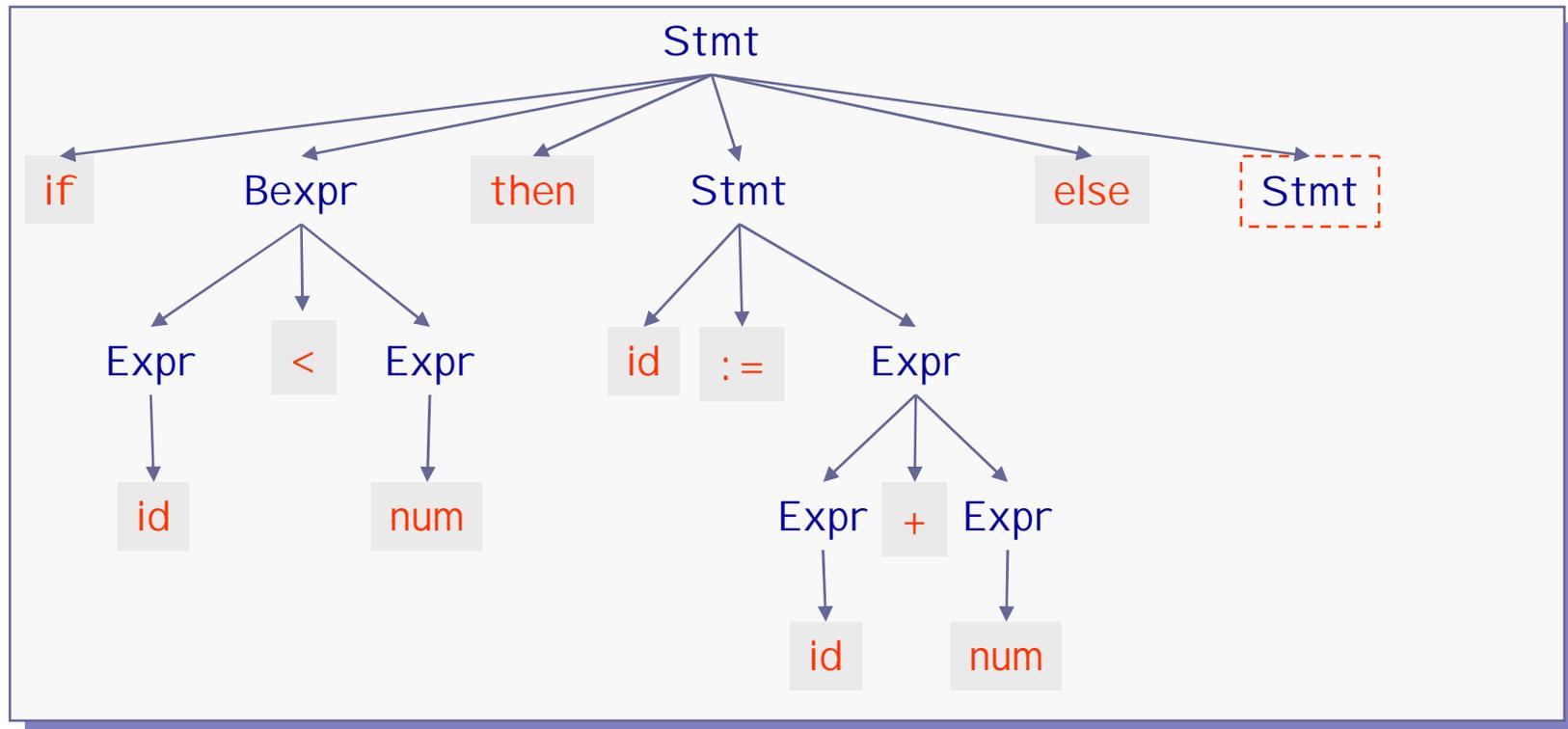


<code>if</code>	<code>x</code>	<code><</code>	<code>0</code>	<code>then</code>	<code>y</code>	<code>:=</code>	<code>y</code>	<code>+</code>	<code>1</code>	<code>else</code>	
-----------------	----------------	-------------------	----------------	-------------------	----------------	-----------------	----------------	----------------	----------------	-------------------	--



Top Down

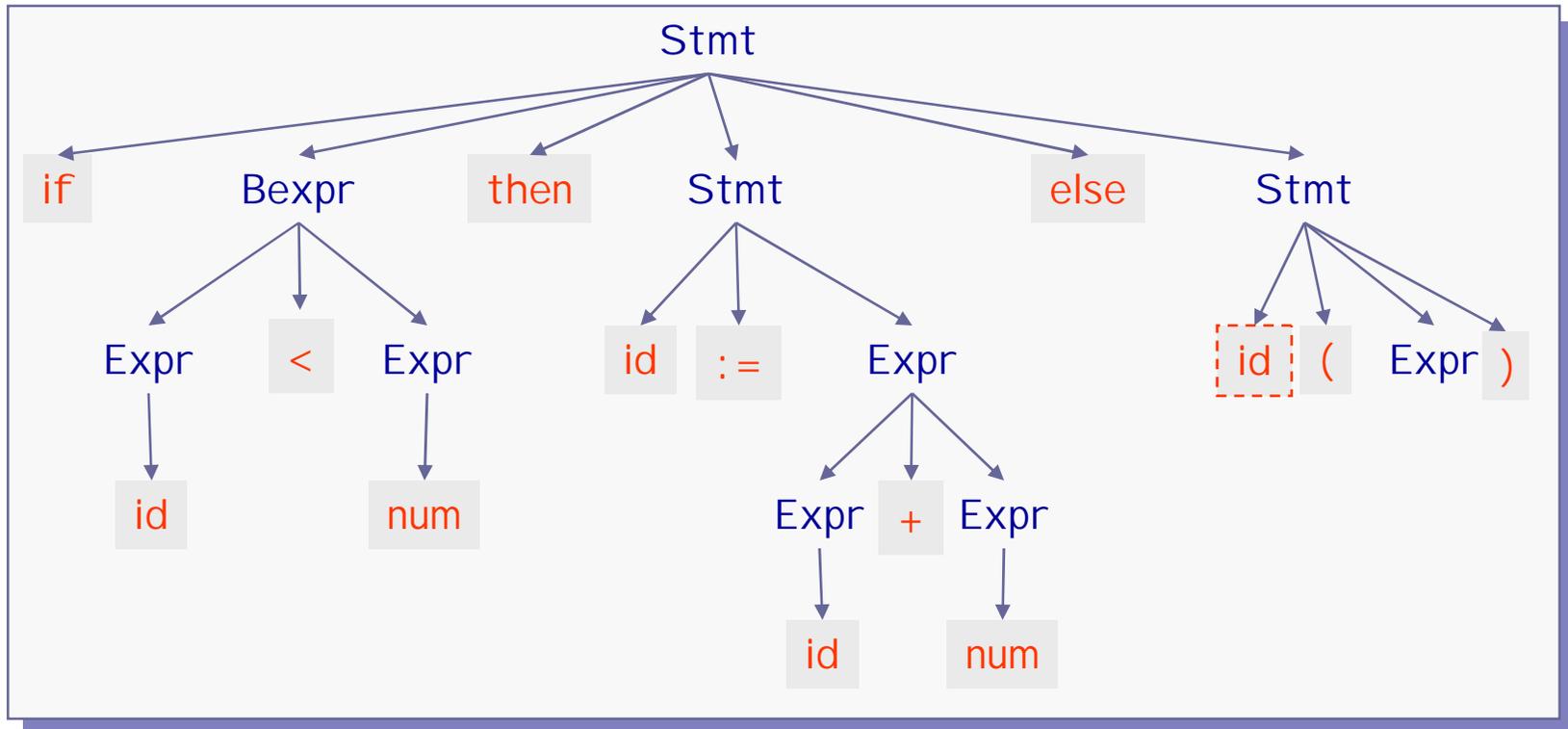
Stmt \rightarrow if Bexpr then Stmt else Stmt
Stmt \rightarrow id := Expr
Stmt \rightarrow id (Expr)



if x < 0 then y := y + 1 else Inc



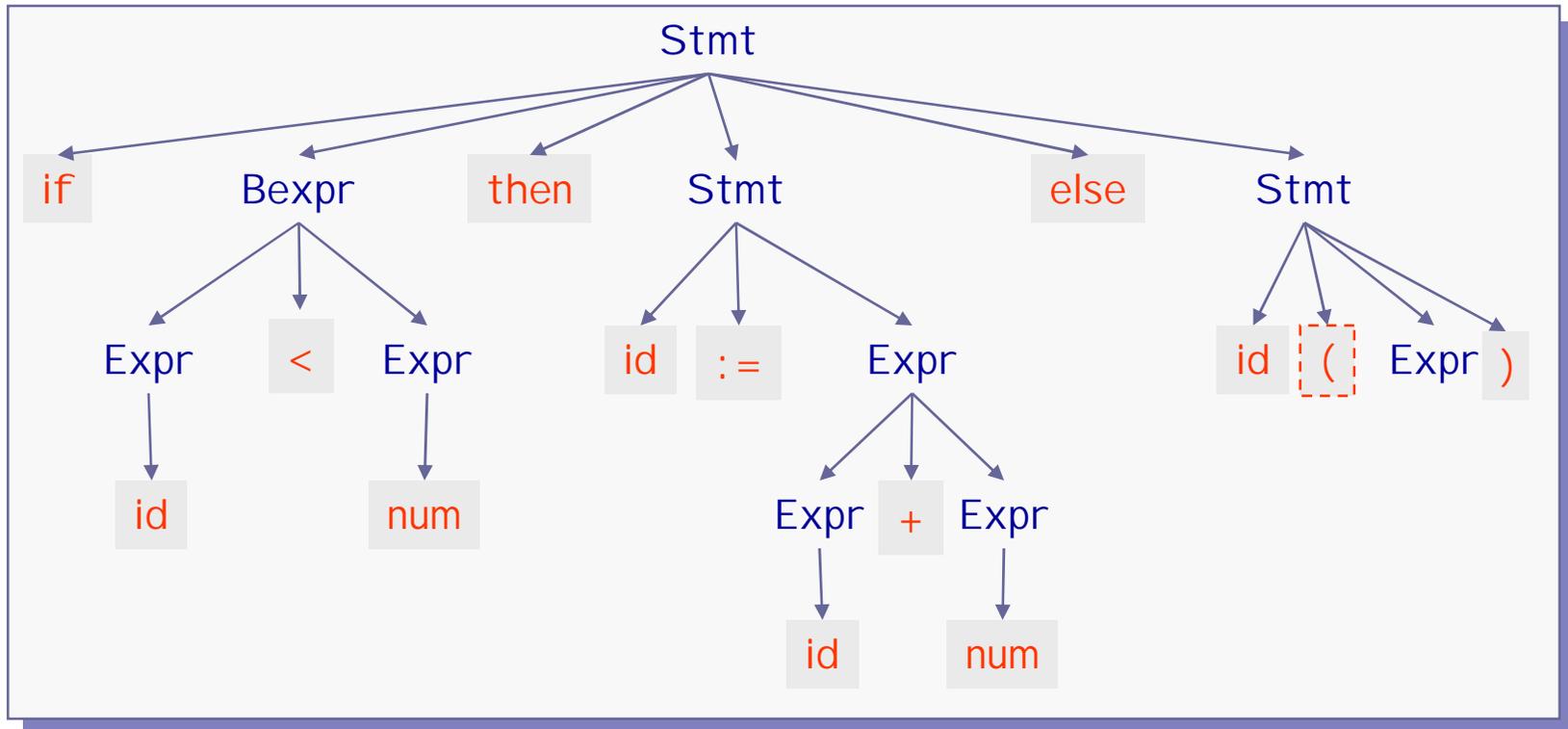
Top Down



<code>if</code>	<code>x</code>	<code><</code>	<code>0</code>	<code>then</code>	<code>y</code>	<code>:=</code>	<code>y</code>	<code>+</code>	<code>1</code>	<code>else</code>	<code>Inc</code>	
-----------------	----------------	-------------------	----------------	-------------------	----------------	-----------------	----------------	----------------	----------------	-------------------	------------------	--



Top Down



if x < 0 then y := y + 1 else Inc (

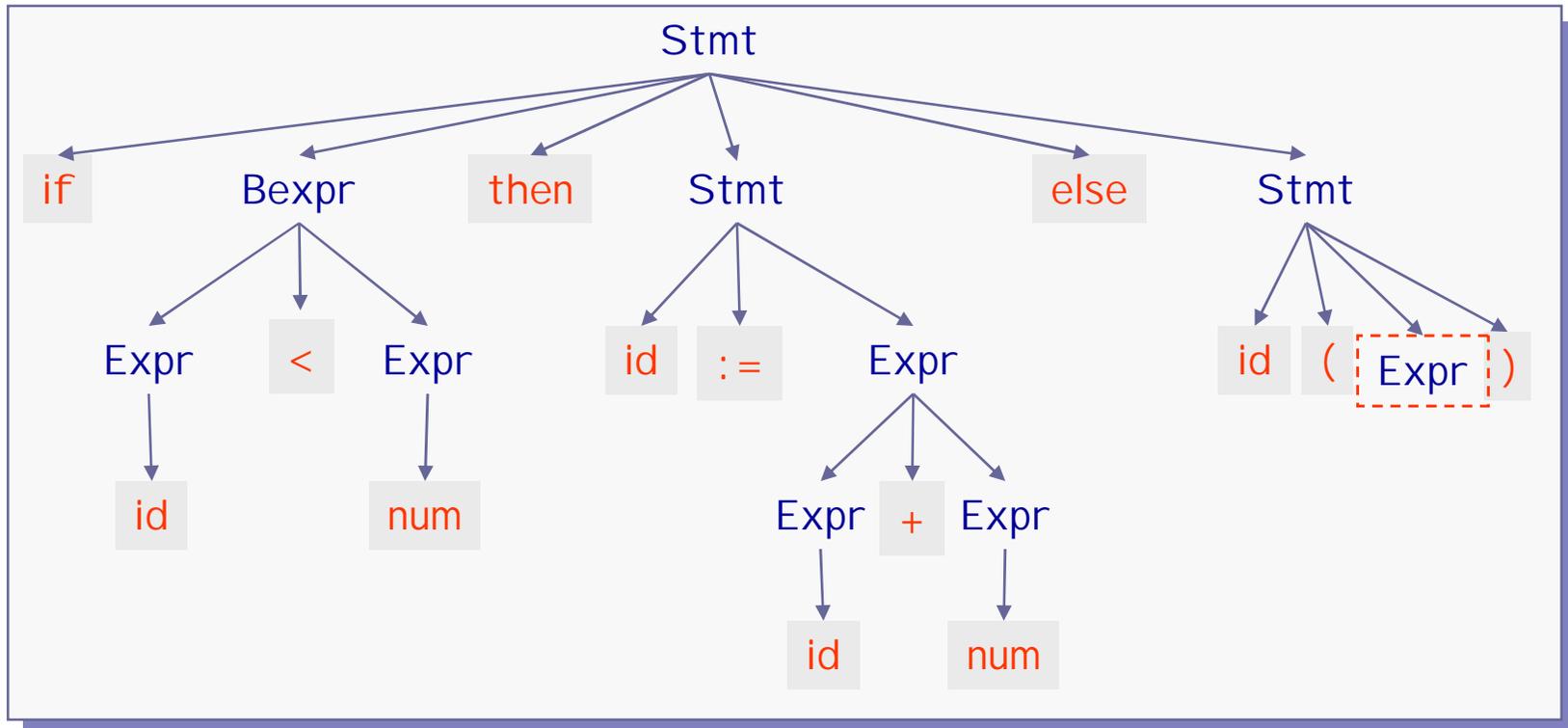


Top Down

Expr \rightarrow Expr + Expr

Expr \rightarrow id

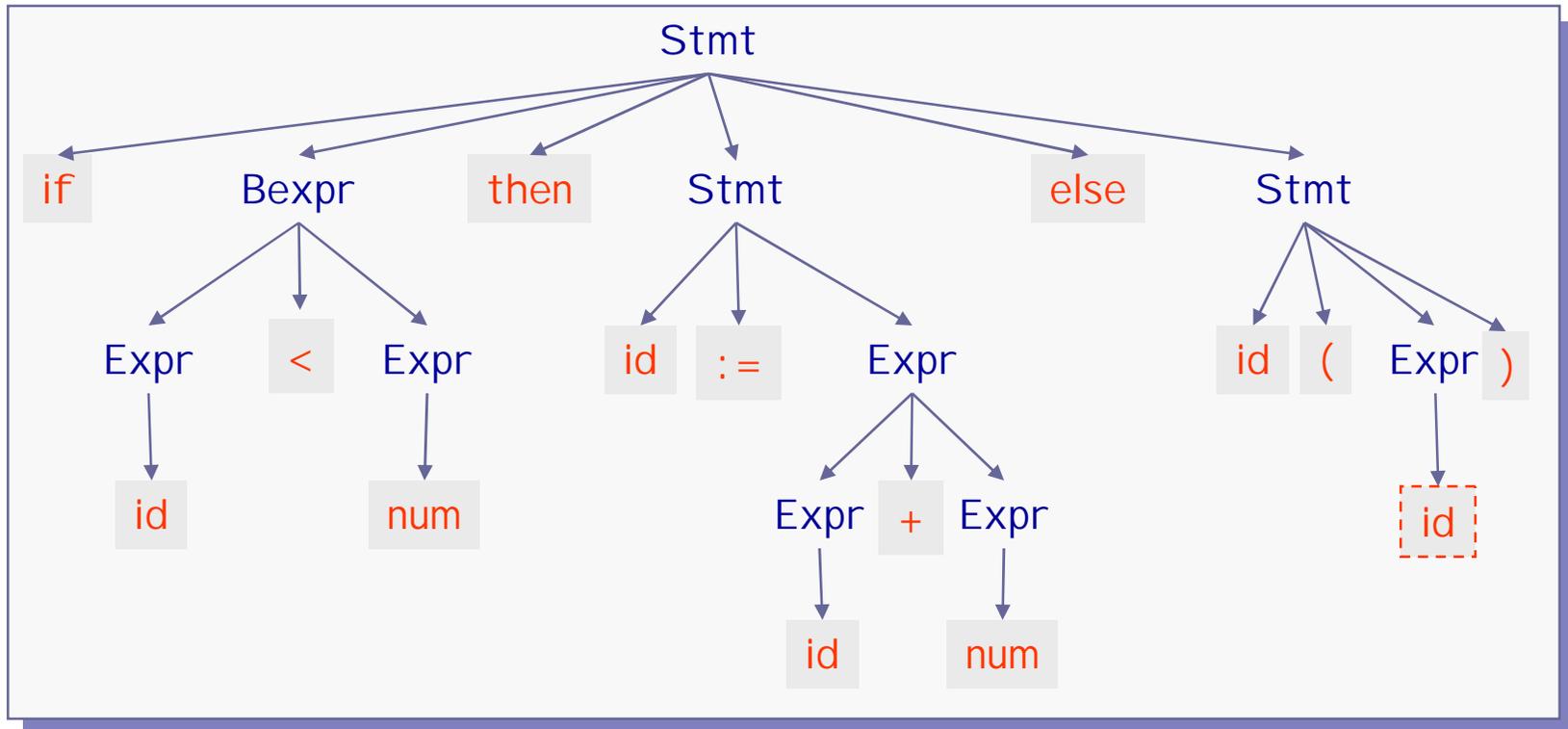
Expr \rightarrow num



if	x	<	0	then	y	:=	y	+	1	else	Inc	(z	
----	---	---	---	------	---	----	---	---	---	------	-----	---	---	--



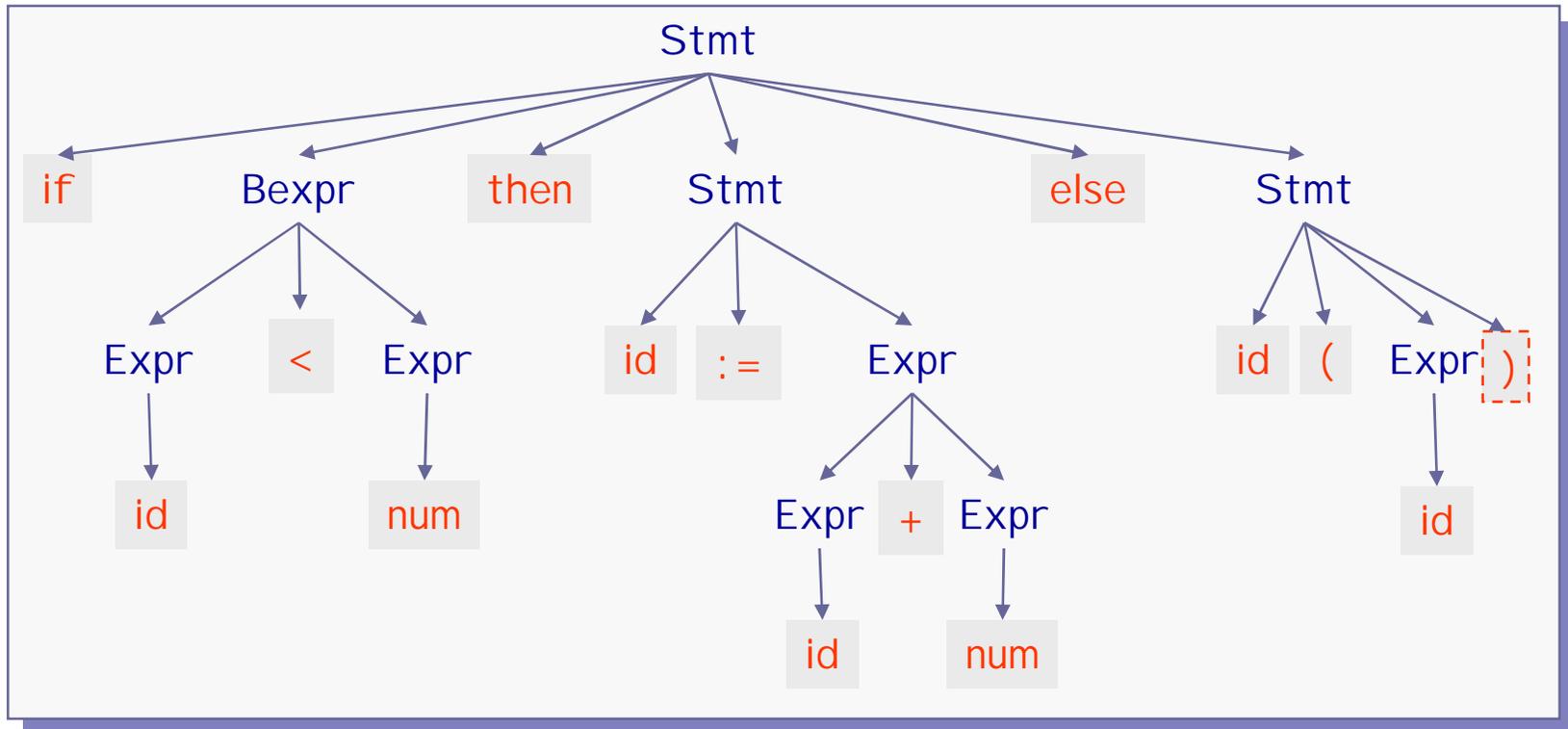
Top Down



if	x	<	0	then	y	:=	y	+	1	else	Inc	(z	
----	---	---	---	------	---	----	---	---	---	------	-----	---	---	--



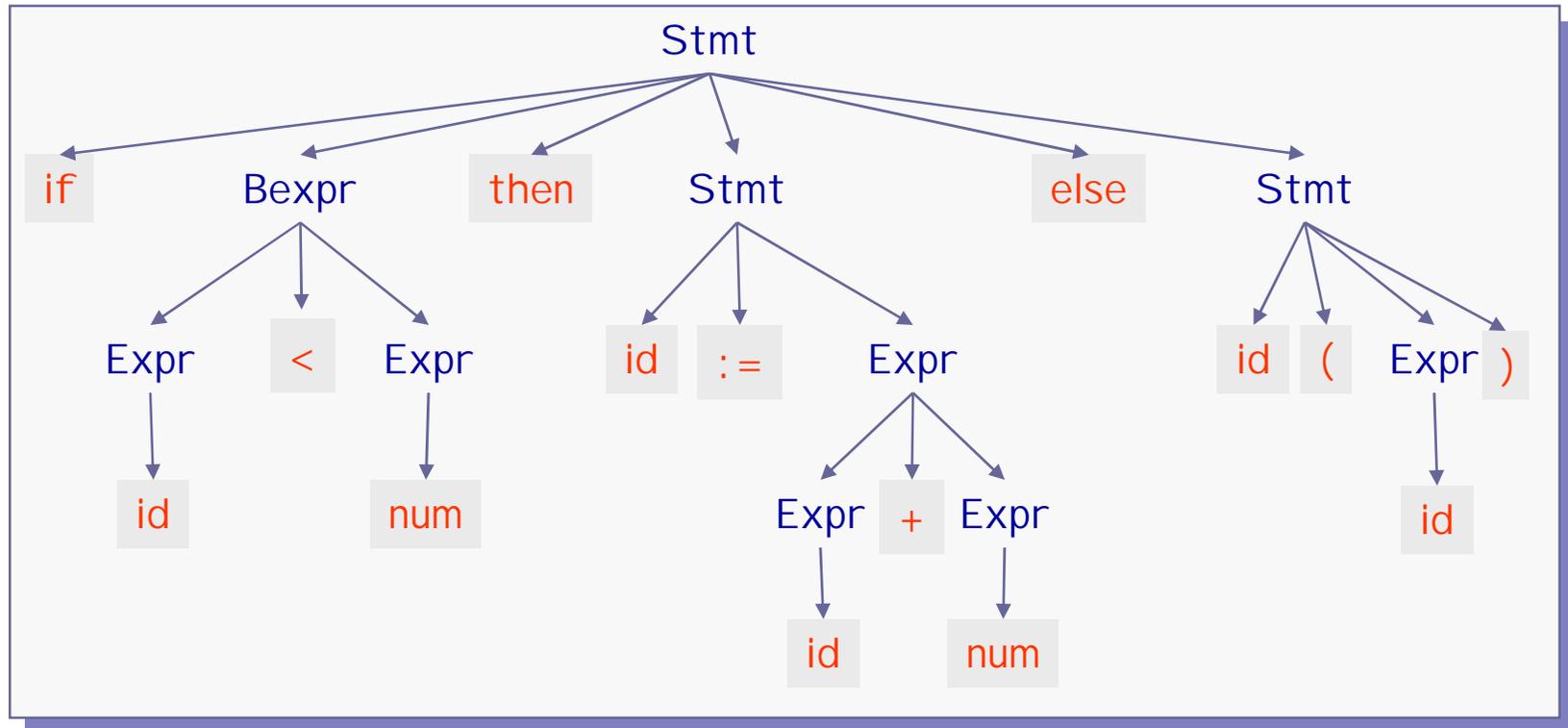
Top Down



if x < 0 then y := y + 1 else Inc (z)



Top Down



Fertig

`if x < 0 then y := y + 1 else Inc (z)`



Recursive Descent Parsing

Eine Produktion $A \rightarrow \beta$ kann man als *Bedingung* verstehen, unter welchen Umständen ein A akzeptiert werden soll.

```
G
Factor → id
       | ( Expr )
```

```
TO parse( Factor )
    consume( id )
OR
    consume( ( ) ; parse( Expr ) ; consume( ) ).
```



Die Grammatik *ist* (fast) schon ein Programm



Recursive Descent Parsing

Ein „Recursive Descent Parser“ (RD-Parser) entsteht aus einer Grammatik auf folgende Weise:

1. Eine Methode „consume“ für die Terminale:

```
void consume(Token t){  
    Token in = nextToken(); // ScannerAufruf  
    if (in==t) return;  
    else fehler(t + "expected");  
}
```

2. Für jedes Nonterminal **A** eine Methode `parseA`:

```
void parseA(){  
  
}
```



Grammatik zum RD-Parser

Aus der Grammatik ...

```
Expr → Term Rest
```

```
Rest → + Term Rest  
      | - Term Rest  
      | ε
```

Nondeterminismus muss
noch entfernt werden

... werden (fast) Methoden

```
void parseExpr() {  
    parseTerm();  
    parseRest();  
}
```

```
void parseRest() {  
    ENTWEDER  
        consume(+);  
        parseTerm();  
        parseRest();  
    ODER  
        consume(-);  
        parseTerm();  
        parseRest();  
    ODER  
        return;  
}
```



Auflösung des Nondeterminismus.

- n Entscheidung zwischen den Alternativen durch Vorausschau auf das nächste Token (**lookahead**).
- n **G** heißt LL(1), falls dieses die richtige Regel eindeutig bestimmt.

G

```
Rest → + Term Rest
      | - Term Rest
      | ε
```

```
void parseRest(){
    switch( lookahead )
        case + : consume(+);
                parseTerm();
                parseRest();
                return;
        case - : consume(-);
                parseTerm();
                parseRest();
                return;
        default : return;
    }
```



First

G enthalte

$$A \rightarrow \sigma$$

$$A \rightarrow \tau.$$

In `parseA` ist Entscheidung anhand **eines** lookahead nur möglich, falls aus σ bzw. τ hergeleitete Worte nicht mit dem gleichen Terminal beginnen können,

Def.

$$\text{First}(\alpha) := \{ t \in T \mid \alpha \Rightarrow^* tw \}$$

zusätzlich

$$\varepsilon \in \text{First}(\alpha) \Leftrightarrow \alpha \Rightarrow^* \varepsilon$$

Notwendige Bedingung für LL(1): Für je zwei $A \rightarrow \sigma, A \rightarrow \tau \in P$ mit $\sigma, \tau \neq \varepsilon$ ist

$$\text{First}(\sigma) \cap \text{First}(\tau) = \emptyset$$



Berechnung von $\text{First}(\alpha)$

$$\text{First}(\epsilon) = \{\epsilon\}$$

$$\text{First}(t.\tau) = \{t\} \quad // \text{ t terminal}$$

$$\begin{aligned} \text{First}(A.\tau) &= \cup \{ \text{First}(\alpha) \mid A \rightarrow \alpha \in G \} \\ &\cup \{ \text{First}(\tau) \mid \text{falls } A \text{ kollabierend} \} \end{aligned}$$

Fast wieder ein rekursives Programm. Nur für Terminierung sorgen !!





ϵ -Produktionen

Zwei kollabierende Variablen reklamieren das leere Wort:

$$\begin{array}{l} \text{String} \rightarrow \text{Expr} \\ \quad | \text{String} + \text{String} \\ \quad | \epsilon \end{array}$$

$$\begin{array}{l} \text{Expr} \rightarrow \text{num Rest}_E \\ \text{Rest}_E \rightarrow + \text{num Rest}_E \\ \quad | \epsilon \end{array}$$

Angenommen, wir parsen einen String (`parseString`)

3 + '=' + 1 + 2

An dieser Stelle kann sich `parseRestE` nicht allein anhand des nächsten Tokens `+` entscheiden! Hängt davon ab, was `folgen` kann.

3 + 4 '=' 7

Grammatik nicht LL(1) !!



Follow



- n RestE kann mit einem + beginnen,
" spricht für $\text{RestE} \rightarrow + \text{num RestE}$

- n auf ein RestE kann aber auch ein + folgen
" spricht für $\text{RestE} \rightarrow \varepsilon$

String \rightarrow Expr
| String + String
| ε

Expr \rightarrow num Rest_E

Rest_E \rightarrow + num Rest_E
| ε

$$\text{Follow}(A) = \{ t \mid t \in T, S \Rightarrow^* \alpha A t \beta \}$$

Problem im Beispiel: $\text{First}(+ \text{num Rest}_E) \cap \text{Follow}(\text{Rest}_E) = \{ + \} \neq \emptyset$



LL(1) Charakterisierung

- n Eine Grammatik G ist LL(1), falls sie für jedes Nonterminal A
- n die First-Bedingung erfüllt
 - $A \rightarrow \sigma$,
 - $A \rightarrow \tau \in P$mit $\sigma, \tau \neq \varepsilon$ impliziert
 - $\text{First}(\sigma) \cap \text{First}(\tau) = \emptyset$
- n die Follow-Bedingung erfüllt
 - $A \rightarrow \sigma$,
 - $A \rightarrow \varepsilon$mit $\sigma \neq \varepsilon$ impliziert
 - $\text{First}(\sigma) \cap \text{Follow}(A) = \emptyset$
- n keine Linksrekursionen enthält
 - keine Regel der Form $A \rightarrow A\sigma$





Die Berechnung der Follow-Mengen

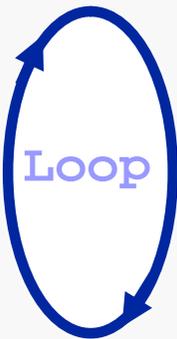
Algorithmus für Follow

Initialisierung:

$$\text{Follow}(S) = \{ \text{eof} \}$$

$$\text{Follow}(A) = \{ \} \text{ für alle } A \neq S.$$

DO



Für jede Produktion $A \rightarrow \alpha B \beta$

$$\text{Follow}(B) \cup = \text{First}(\beta) - \{\epsilon\}$$

Für jede Produktion $A \rightarrow \alpha B \beta$ mit $\epsilon \in \text{First}(\beta)$

$$\text{Follow}(B) \cup = \text{Follow}(A)$$

BIS sich nichts mehr ändert

Beachte

$\text{First}(\alpha)$ ist für jede Satzform α definiert,
 $\text{Follow}(A)$ nur für Variablen.



Beispiel

$$\begin{array}{l} \text{String} \rightarrow \text{Expr} \\ \quad | \text{String} + \text{String} \\ \quad | \varepsilon \end{array}$$
$$\text{Expr} \rightarrow \text{num Rest}_E$$
$$\begin{array}{l} \text{Rest}_E \rightarrow + \text{num Rest}_E \\ \quad | \varepsilon \end{array}$$

n First-Mengen der rechten Seiten

- $\text{First}(\text{Expr}) = \{\text{num}\}$
- $\text{First}(\text{String} + \text{String}) = \{\text{num}, +\}$
- $\text{First}(\text{num Rest}_E) = \{\text{num}\}$
- $\text{First}(\varepsilon) = \{\varepsilon\}$

n Follow-Mengen der Nonterminale

- $\text{Follow}(\text{String}) = \{\text{eof}, +, \}$
- $\text{Follow}(\text{Expr}) = \{\text{eof}, +\}$
- $\text{Follow}(\text{Rest}_E) = \{\text{eof}, +\}$

n $\text{First}(\text{Expr}) \cap \text{First}(\text{String} + \text{String}) \neq \emptyset$

n $\text{First}(\text{String} + \text{String}) \cap \text{Follow}(\text{String}) \neq \emptyset$

Grammatik
nicht LL(1)





Linksrekursion

Direkte Übersetzung einer Grammatik in (rekursiven) RD-Parser kann zu Endlosschleifen führen:

```
Expr → Expr + Expr  
Expr → id  
Expr → num
```

Hier würde `parseExpr` sofort wieder `parseExpr` aufrufen ohne einen Fortschritt zu machen.

Problem lösbar durch äquivalente Umformungen der Grammatik :

```
Expr → T Rest  
Rest → + T Rest  
      | ε  
T     → id | num
```



Elimination von Linksrekursionen

Angenommen die Produktionen für Nonterminal A sind

$$\begin{array}{l} A \rightarrow A \alpha_1 \\ \dots \\ A \rightarrow A \alpha_n \\ \\ A \rightarrow \beta_1 \\ \dots \\ A \rightarrow \beta_n \end{array}$$

wobei die α_i, β_i nicht mit A beginnen. Man kann sie ersetzen durch:

$$\begin{array}{l} A \rightarrow \beta_1 \text{Rest}_A \\ \dots \\ A \rightarrow \beta_n \text{Rest}_A \end{array}$$

$$\begin{array}{l} \text{Rest}_A \rightarrow \alpha_1 \text{Rest}_A \\ \dots \\ \text{Rest}_A \rightarrow \alpha_n \text{Rest}_A \\ \\ \text{Rest}_A \rightarrow \varepsilon \end{array}$$



LL(k) Sprachen

- n Verallgemeinerung: k token lookahead.
- n Sprache L ist LL(k), falls LL(k)-Grammatik für L existiert.
- n Falls Grammatik nicht LL(k), versuche in äquiv. LL(k)-Grammatik zu transformieren.

Grammatik ist **nicht LL(1)**,

Sprache **ist LL(1) !!**

G_1

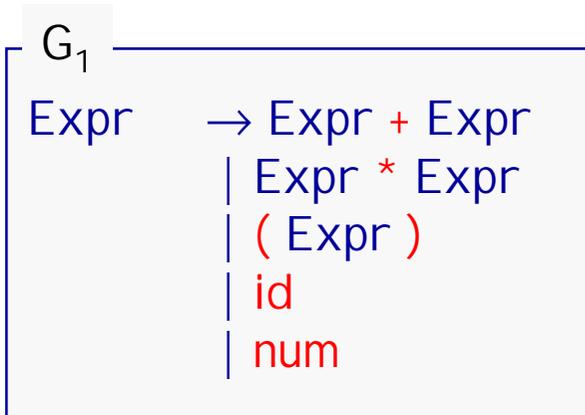
Expr	→	Expr + Expr
		Expr * Expr
		(Expr)
		id
		num

In einer Reihe von Schritten werden wir diese Grammatik in eine äquivalente LL(1) Grammatik umformen.

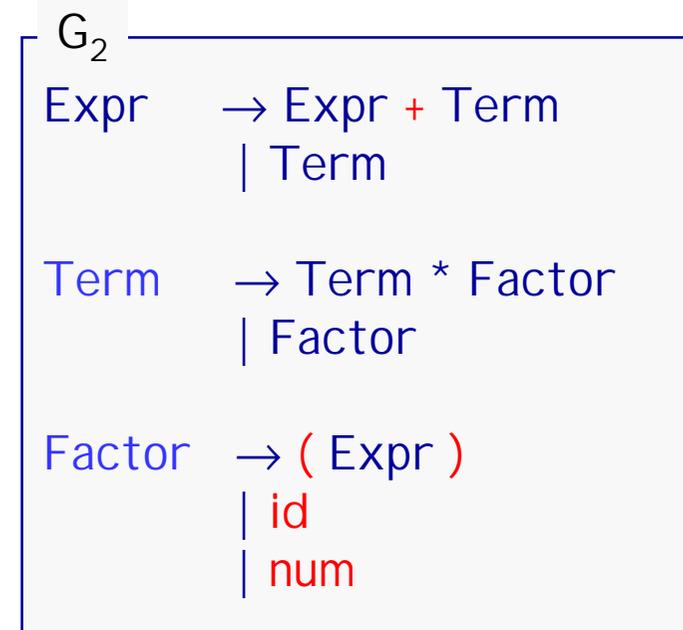
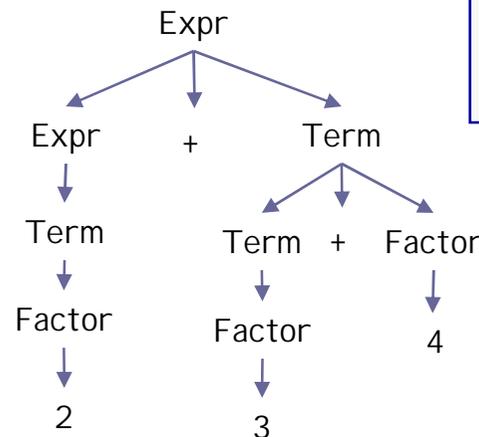


Von der Grammatik zum RD-Parser

Entferne Mehrdeutigkeit durch Einführung von getrennten Nonterminalen für jede Präzedenzstufe



2 + 3 * 4 hat zwei Herleitungsbäume



2 + 3 * 4 hat eindeutigen Herleitungsbaum



Von der Grammatik zum RD-Parser

n Linksrekursionen entfernen

G_2

$$\begin{aligned} \text{Expr} &\rightarrow \text{Expr} + \text{Term} \\ &\quad | \text{Term} \\ \\ \text{Term} &\rightarrow \text{Term} * \text{Factor} \\ &\quad | \text{Factor} \\ \\ \text{Factor} &\rightarrow (\text{Expr}) \\ &\quad | \text{id} \\ &\quad | \text{num} \end{aligned}$$

G_3

$$\begin{aligned} \text{Expr} &\rightarrow \text{Term Rest}_E \\ \\ \text{Rest}_E &\rightarrow + \text{Term Rest}_E \\ &\quad | \varepsilon \\ \\ \text{Term} &\rightarrow \text{Factor Rest}_T \\ \\ \text{Rest}_T &\rightarrow * \text{Factor Rest}_T \\ &\quad | \varepsilon \\ \\ \text{Factor} &\rightarrow (\text{Expr}) \\ &\quad | \text{id} \\ &\quad | \text{num} \end{aligned}$$



Von der Grammatik zum RD-Parser

G_3

$\text{Expr} \rightarrow \text{Term Rest}_E$

$\text{Rest}_E \rightarrow + \text{Term Rest}_E$
 $\quad \quad \quad | \varepsilon$

$\text{Term} \rightarrow \text{Factor Rest}_T$

$\text{Rest}_T \rightarrow * \text{Factor Rest}_T$
 $\quad \quad \quad | \varepsilon$

$\text{Factor} \rightarrow (\text{Expr})$
 $\quad \quad \quad | \text{id}$
 $\quad \quad \quad | \text{num}$

First-Mengen berechnen:

$\text{First}(\text{Term Rest}_E) = \{ (, \text{id}, \text{num} \}$

$\text{First}(+ \text{Term Rest}_E) = \{ + \}$

$\text{First}(\text{Rest}_E) = \{ +, \varepsilon \}$

$\text{First}(\text{Factor Rest}_T) = \{ (, \text{id}, \text{num} \}$

$\text{First}(* \text{Factor Rest}_T) = \{ * \}$

$\text{First}(\text{Rest}_T) = \{ *, \varepsilon \}$

$\text{First}((\text{Expr})) = \{ (\}$

$\text{First}(\text{id}) = \{ \text{id} \}$

$\text{First}(\text{num}) = \{ \text{num} \}$



Von der Grammatik zum RD-Parser

Follow-Mengen berechnen:

G_3

$\text{Expr} \rightarrow \text{Term Rest}_E$

$\text{Rest}_E \rightarrow + \text{Term Rest}_E$
 $\quad \quad \quad | \varepsilon$

$\text{Term} \rightarrow \text{Factor Rest}_T$

$\text{Rest}_T \rightarrow * \text{Factor Rest}_T$
 $\quad \quad \quad | \varepsilon$

$\text{Factor} \rightarrow (\text{Expr})$
 $\quad \quad \quad | \text{id}$
 $\quad \quad \quad | \text{num}$

$\text{Follow}(\text{Expr}) = \{), \text{eof} \}$

$\text{Follow}(\text{Rest}_E) = \{), \text{eof} \}$

$\text{Follow}(\text{Term}) = \{ +,), \text{eof} \}$

$\text{Follow}(\text{Rest}_T) = \{ +,), \text{eof} \}$

$\text{Follow}(\text{Factor}) = \{ *, +,), \text{eof} \}$



Von der Grammatik zum RD-Parser

G_3

Expr \rightarrow Term Rest_E

Rest_E \rightarrow + Term Rest_E
| ϵ

Term \rightarrow Factor Rest_T

Rest_T \rightarrow * Factor Rest_T
| ϵ

Factor \rightarrow (Expr)
| id
| num

Überprüfen der First Bedingungen:

$$\text{First}(+ \text{ Term Rest}_E) \cap \text{First}(\epsilon) = \emptyset$$

$$\text{First}(* \text{ Factor Rest}_T) \cap \text{First}(\epsilon) = \emptyset$$

$$\text{First}(\text{ (Expr) }) \cap \text{First}(\text{ id }) = \emptyset$$

$$\text{First}(\text{ (Expr) }) \cap \text{First}(\text{ num }) = \emptyset$$

$$\text{First}(\text{ id }) \cap \text{First}(\text{ num }) = \emptyset$$

Überprüfen der Follow Bedingungen:

$$\text{First}(+ \text{ Term Rest}_E) \cap \text{Follow}(\text{ Rest}_E) = \emptyset$$

$$\text{First}(* \text{ Factor Rest}_T) \cap \text{Follow}(\text{ Rest}_T) = \emptyset$$



Von der Grammatik zum RD-Parser

- n Da die LL(1) Bedingungen erfüllt sind, ist die Transformation in einen RD-Parser praktisch automatisch:

G_3

Expr	→	Term Rest _E
Rest _E	→	+ Term Rest _E ε
Term	→	Factor Rest _T
Rest _T	→	* Factor Rest _T ε
Factor	→	(Expr) id num

```
/* parser */  
void parseExpr() {  
    parseTerm();  
    parseRestE();  
}  
  
void parseRestE() {  
    switch(lookahead)  
    case PLUS :  
        consume(PLUS);  
        parseTerm();  
        parseRestE();  
        return;  
    default: return;  
}  
... etc.
```



Semantische Aktionen

Parser soll Sprache nicht nur entscheiden,
sondern z.B. auch

- Wert des Ausdrucks berechnen
- Ausdruck "pretty-print"en
- Ausdruck compilieren.



Wichtige Voraussetzung

- Resultat für Gesamtausdruck ist Kombination von Resultaten für Teilausdrücke
- denotationelles Prinzip

Vorgehensweise

- Führe an passender Stelle Aktionen ein
- z.B. Programmcode
- in der Ausgangsgrammatik
- Transformiere Grammatik
- Handle Aktionen in der Transformation wie Terminale
- ignoriere sie bei First- und Follow-Bestimmung



Beispiel: Infix-Postfix Übersetzer

n Aufgabe:

- Infix Ausdrücke lesen,
- in Postfix übersetzen.

G_1

Expr	→	Expr + Expr
		Expr * Expr
		(Expr)
		id
		num

Gehe davon aus, dass jedes Nonterminal auf der rechten Seite bereits seine Aktion ausführt

G_{postfix}

Expr	→	Expr + Expr	write("+");
		Expr * Expr	write("*");
		(Expr)	// Nix
		id	write("id");
		num	write("num");

- Für (Expr) keine Aktion notwendig:
- Expr führt seine Aktion selber aus
 - Postfix braucht keine Klammern



Ein Infix-Postfix Übersetzer

Mehrdeutigkeit entfernen durch Präzedenzstufen

G_{postfix}

Expr	→ Expr + Expr	write("+");
	Expr * Expr	write("*");
	(Expr)	
	id	write("id");
	num	write("num");

Man hätte auch gleich
hiermit beginnen können:

G_2

Expr	→ Expr + Term	write("+");
	Term	
Term	→ Term * Factor	write("*");
	Factor	
Factor	→ (Expr)	
	id	write("id");
	num	write("num");



Ein Infix-Postfix Übersetzer

Entferne Linksrekursion. Handle dabei Aktionen wie Terminale.

G_2

```
Expr  → Expr + Term write("+");  
      | Term  
Term   → Term * Factor write("*");  
      | Factor  
Factor → ( Expr )  
      | id  write("id");  
      | num write("num");
```

G_3

```
Expr  → Term RestE  
RestE → + Term { write("+"); } RestE  
      | ε  
Term   → Factor RestT  
RestT → * Factor { write("*"); } RestT  
      | ε  
Factor → ( Expr )  
      | id  { write("id"); }  
      | num { write("num"); }
```



Ein Infix-Postfix Übersetzer

Prüfe LL(1) Bedingungen. Ignoriere dabei die Aktionen
Transformiere in Programm:

G_3

```
Expr  → Term RestE
RestE → + Term { write("+"); } RestE
        | ε
Term   → Factor RestT
RestT  → * Factor { write("*"); } RestT
        | ε
Factor → ( Expr )
        | id   { write("id"); }
        | num  { write("num"); }
```

```
/* Postfix--Übersetzer */
void parseExpr() {
    parseTerm();
    parseRestE();
}

void parseRestE() {
    switch(lookahead)
    case PLUS :
        consume(PLUS);
        parseTerm();
        write("+");
        parseRestE();
        return;
    default: return;
}
... etc.
```



Inhalt

1. Grammatiken und Sprachen
 - .. Kontextfreie Grammatiken
 - .. Herleitungen, Linksherleitungen
 - .. Sprachen zu einer Grammatik
 - .. Äquivalenz
 - .. Chomsky-Normalform
 - .. Wortproblem, CYK
 - .. Pumping Lemma für CF-Sprachen
2. Stackautomaten
 - .. Definitionen und Beispiele
 - .. Konfigurationen, Läufe,
 - .. Sprache eines Stackautomaten
 - .. Parser

3. Parser
 - .. Mehrdeutigkeit
 - .. Bottom Up, Top Down
 - .. Recursive descent Parser
 - .. First, Follow
 - .. Semantische Aktionen
4. Shift-Reduce Parser
 - .. Konflikte
 - .. LR(0)-Zustände
 - .. Automat zur Grammatik
 - .. Shift-Reduce, Reduce-Reduce Konflikte
 - .. Präzedenz und Assoziativität
 - .. lex und yacc
 - .. Arbeitsweise eines Compilers



Shift-Reduce Parser

Shift-Reduce Parser sind "bottom-up parser", d.h. sie bauen den Syntaxbaum von den Blättern her auf. Aktionen sind :



Lese ein weiteres Zeichen aus dem Input



Die rechte Seite einer Produktion wurde im Input erkannt. Ersetze sie durch die linke Seite.

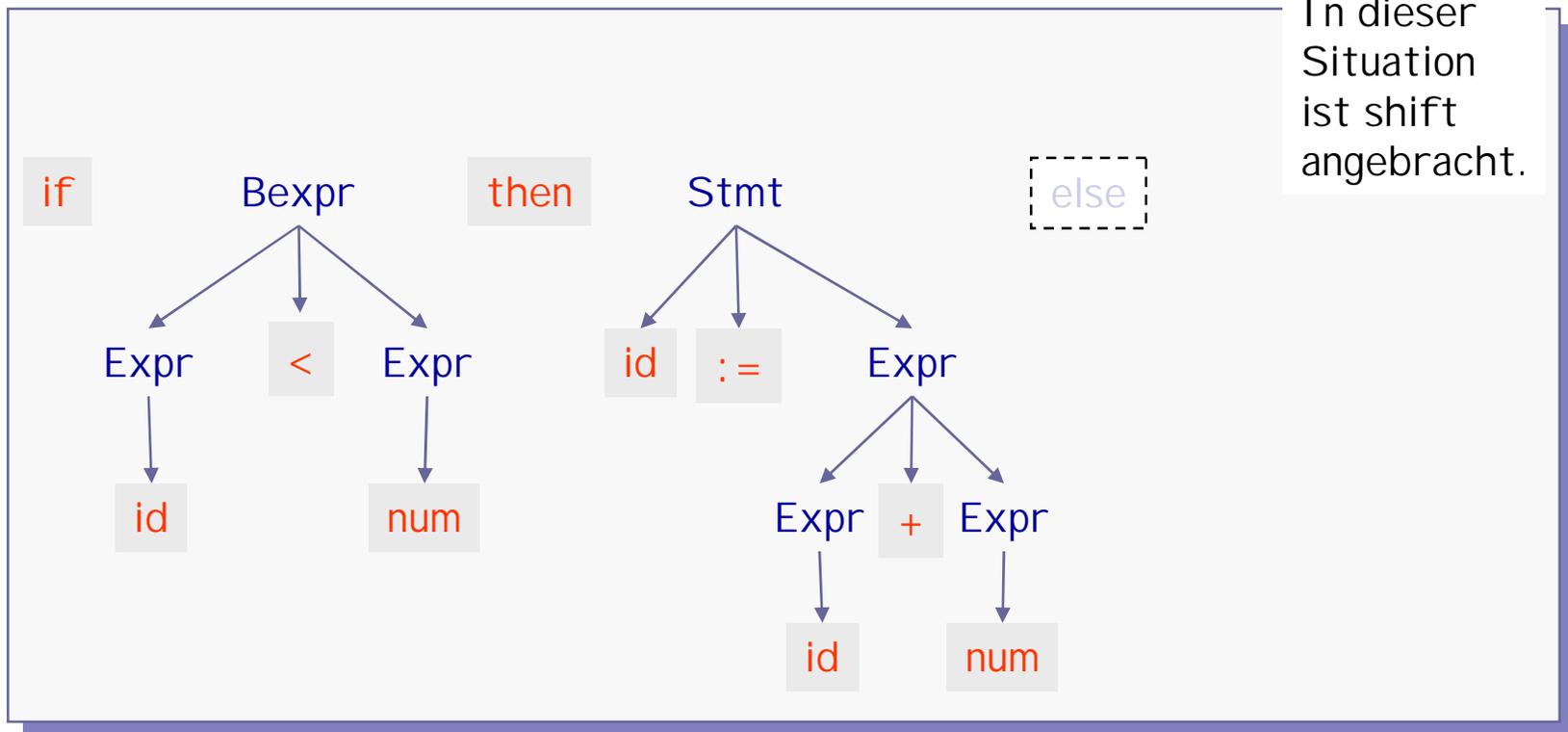
Input kann dabei eine beliebige Satzform sein.



Shift-Schritt

lookahead einlesen
Scanner legt nächstes
Token in lookahead

Stmt \rightarrow if Bexpr then Stmt else Stmt
Stmt \rightarrow if Bexpr then Stmt
Stmt \rightarrow id (Expr)



In dieser Situation ist shift angebracht.



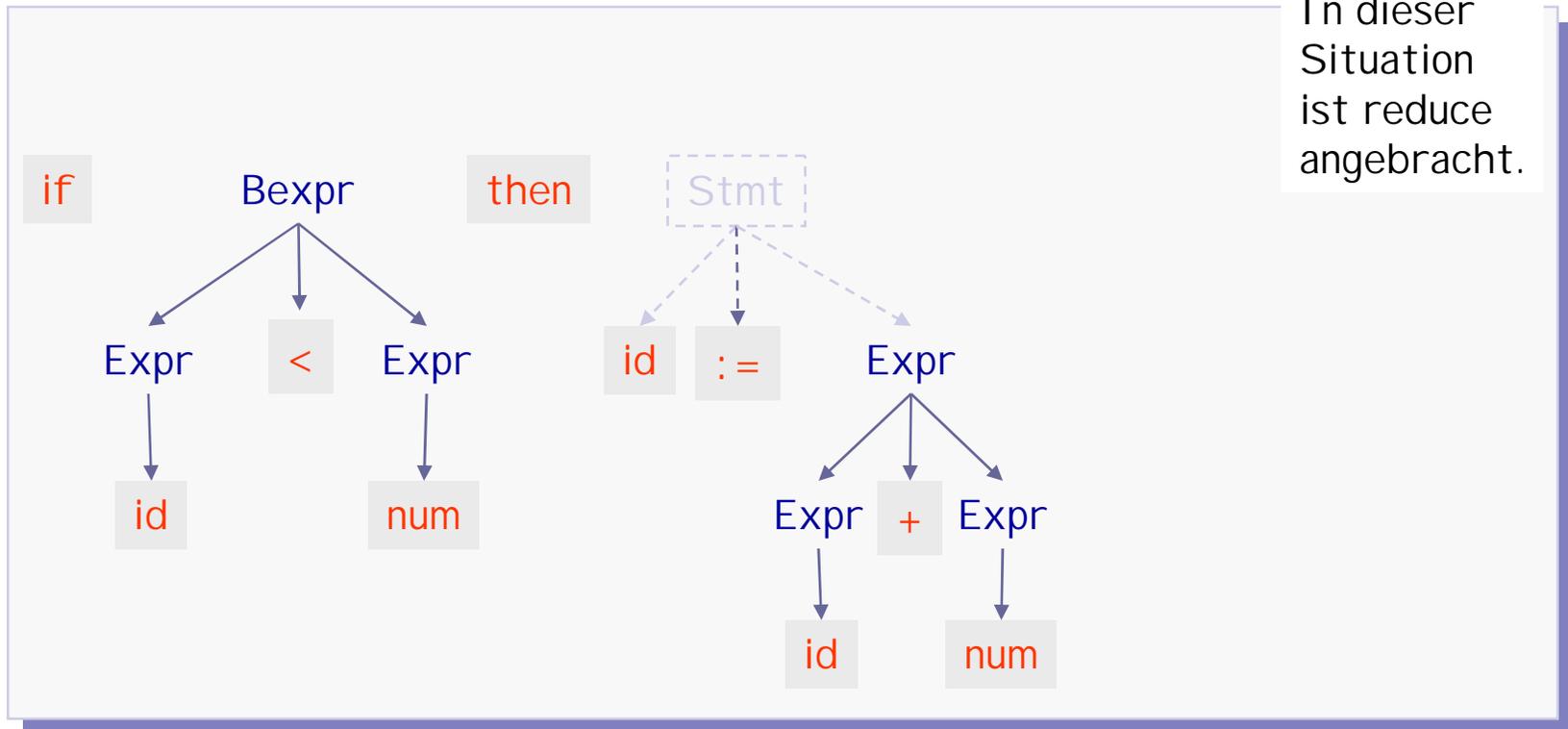


Reduce-Schritt

Regel auswählen, deren rechte Seite auf die letzten Bäume im Wald passt



Stmt \rightarrow if Bexpr then Stmt else Stmt
Stmt \rightarrow id := Expr
Stmt \rightarrow id (Expr)

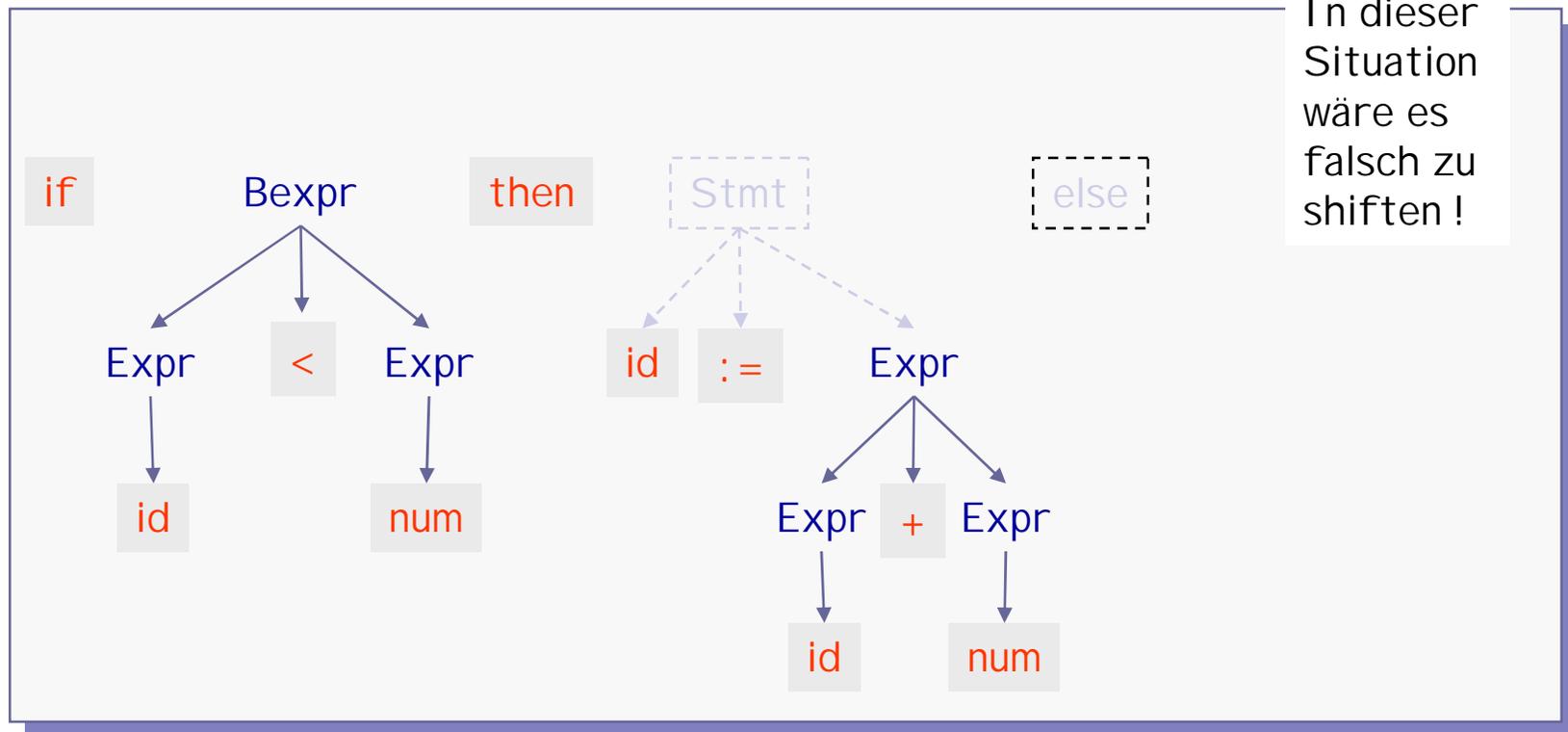


if x < 0 then y := y + 1 else lookahead



Shift oder reduce ?

Stmt \rightarrow if Bexpr then Stmt else Stmt
Stmt \rightarrow id := Expr
Stmt \rightarrow id (Expr)



In dieser Situation wäre es falsch zu shiften !



if x < 0 then y := y + 1 **else**





Reduce-reduce-Konflikt



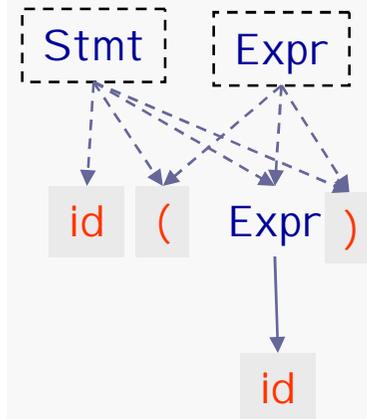
Stmt \rightarrow id (Expr)

Expr \rightarrow (Expr)

- Es gibt zwei Regeln, mit denen reduziert werden könnte

- Das lookahead bringt keine Entscheidung

Reduce-Reduce-Konflikte sind unangenehm. Sie werden oft durch Einführung von Präzedenzen entschieden



```
if x < 0 then y := y + 1 else Inc ( z )
```



Shift-Reduce-Konflikt



Stmt \rightarrow if Bexpr then Stmt else Stmt

Stmt \rightarrow if Bexpr then Stmt

if

Bexpr

then

Stmt

Es könnte mit der zweiten Regel reduziert werden

Die Existenz der ersten Regel lässt auch ein shift zu, da noch ein Stmt folgen kann.

Das lookahead bringt keine Entscheidung

if x < 0 then y := y + 1 else

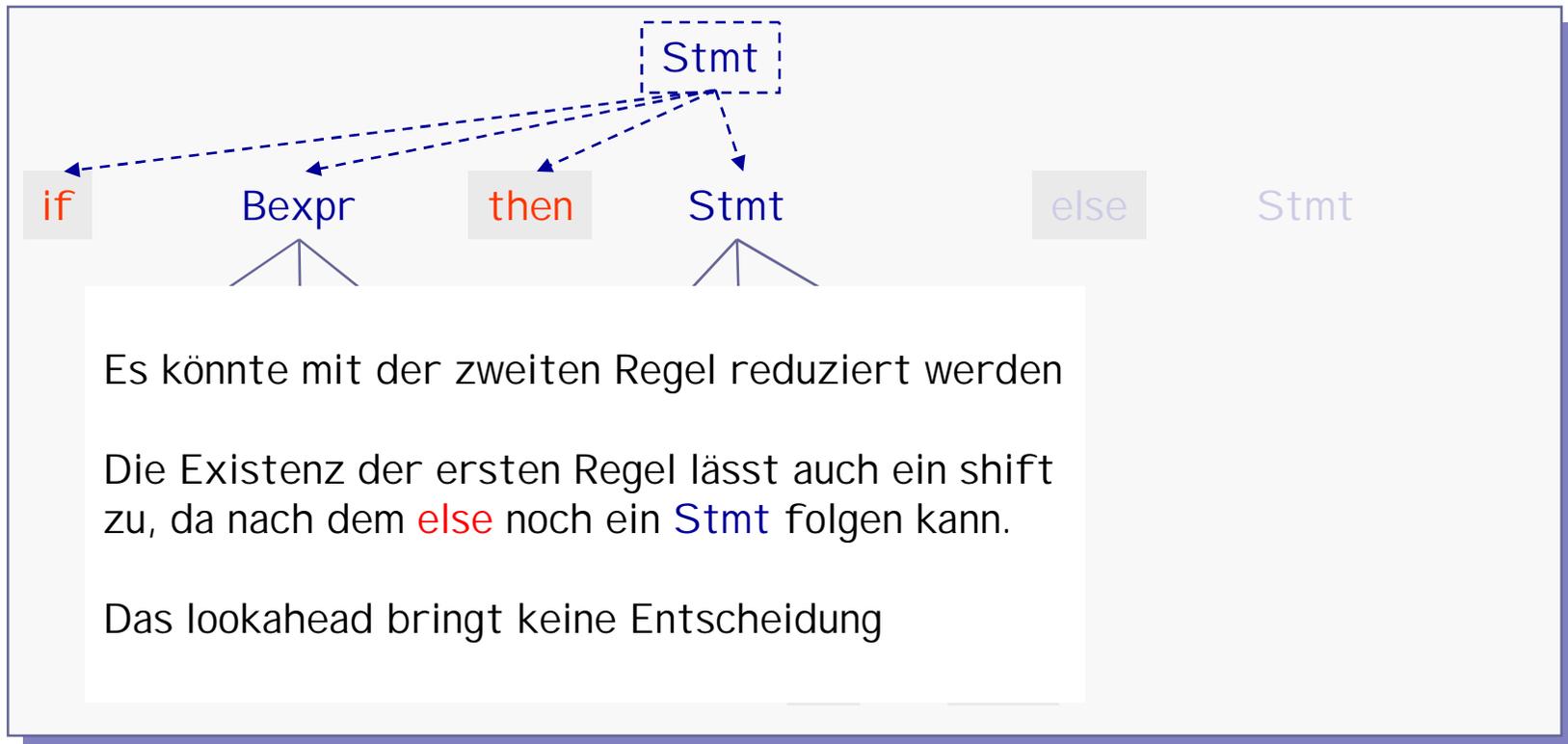


Shift-Reduce-Konflikt

Stmt \rightarrow if Bexpr then Stmt else Stmt



Stmt \rightarrow if Bexpr then Stmt



if x < 0 then y := y + 1 **else**



Zustände

- n Konflikte können nicht allein aufgrund des lookahead entschieden werden.
- n Parser merkt sich bisher nicht, was er gesehen hat
- n Lösung: Zustände einführen, die besagen, was der Parser gerade zu erkennen versucht.

G			
Expr	→ Expr + Term	// 1	
	Term	// 2	

Bisher erkannt:	lookahead	Aktion:
... THEN id := Term	+	red(2)
... THEN id := Expr + Term	+	red(1)
... THEN id := Expr	+	shift

Merke die lokale Situation. Notation :

Expr → Term • : wollte Expr sehen, habe Term erkannt

Expr → Expr • + Term : wollte Expr sehen, habe Expr erkannt
jetzt muss + kommen, dann ein Term

Zustand:	lookahead	Aktion:
Expr → Term •	+	red (2)
Expr → Expr + Term •	+	red (1)
Expr → Expr • + Term	+	shift



LR(0) - items

Ein item ist eine Produktion $A \rightarrow \alpha$ mit einer markierten Position in α .

Zustand beinhaltet also

- Nonterminal, das der Parser gerade erkennen will
- mit welcher Regel er dies versucht und
- wieviel von der rechten Seite der Regel er bereits erkannt hat.

Beispiel: Die Produktion $\text{Expr} \rightarrow \text{Expr} + \text{Term}$ liefert genau 4 items mit folgenden Bedeutungen:

$\text{Expr} \rightarrow \bullet \text{Expr} + \text{Term}$
 $\text{Expr} \rightarrow \text{Expr} \bullet + \text{Term}$
 $\text{Expr} \rightarrow \text{Expr} + \bullet \text{Term}$
 $\text{Expr} \rightarrow \text{Expr} + \text{Term} \bullet$

Bei dem Versuch ein Expr zu finden ...

... erwarte $\text{Expr} + \text{Term}$
... Expr gesehen, erwarte $+ \text{Term}$
... $\text{Expr} +$ gesehen, erwarte Term
... $\text{Expr} + \text{Term}$ gesehen

LR - steht für Left-Right, gemeint ist die Abarbeitung des Inputs von Links nach rechts



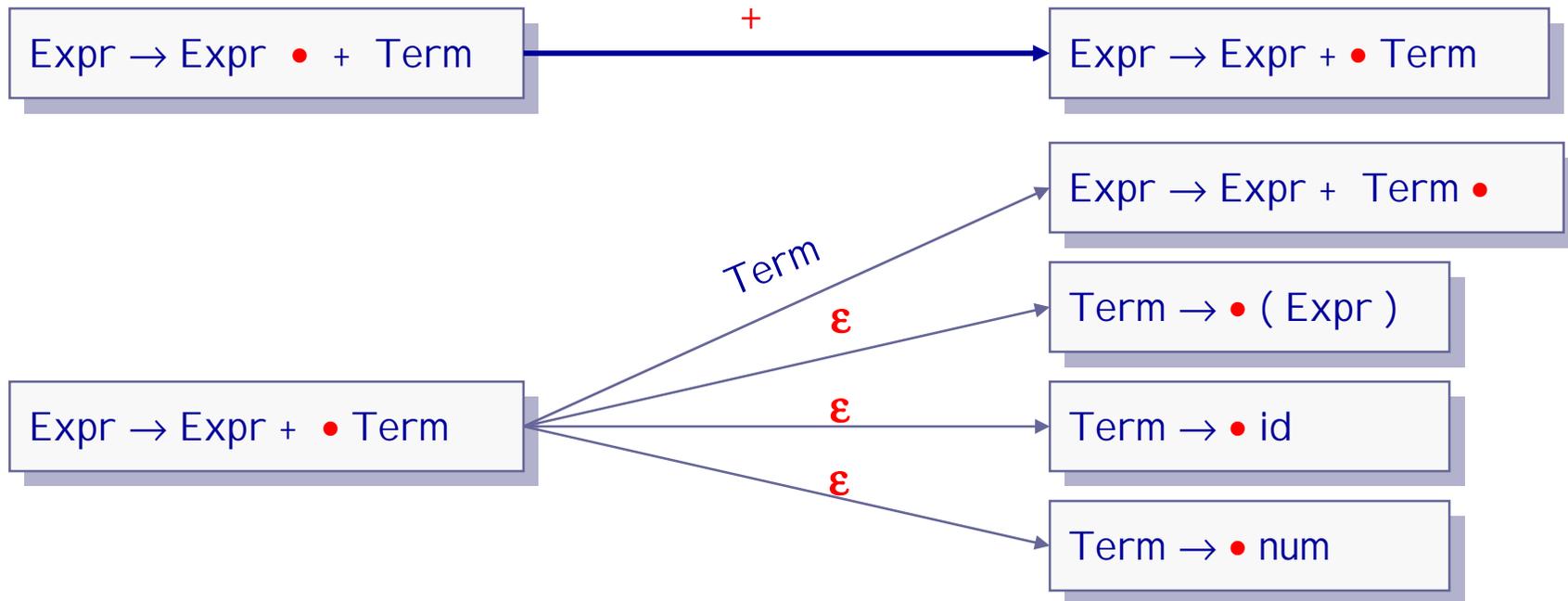
Automat einer Grammatik

- n Ein **item** drückt den jeweiligen Zustand des Parsers beim Erkennen eines Nonterminalaus.
- n Mit dem Lesen des nächsten Tokens geht der Parser in den nächsten Zustand über. Es gibt zwei mögliche Übergänge :

G_{simpl}

$\text{Expr} \rightarrow \text{Expr} + \text{Term}$
 $\quad \quad \quad | \text{Term}$

$\text{Term} \rightarrow (\text{Expr})$
 $\quad \quad \quad | \text{id}$
 $\quad \quad \quad | \text{num}$

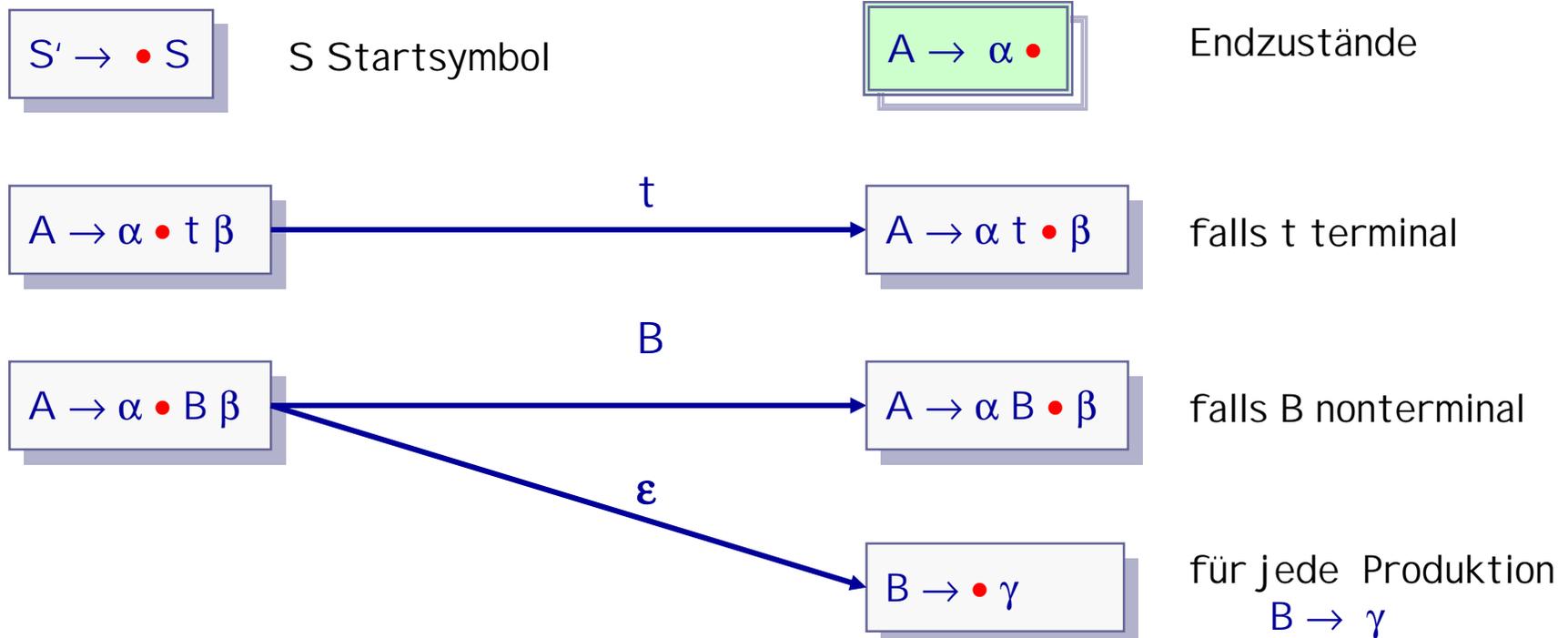


Dies definiert einen nichtdeterministischen Automaten: $N(G)$.



Automat einer Grammatik

n Allgemeine Definition





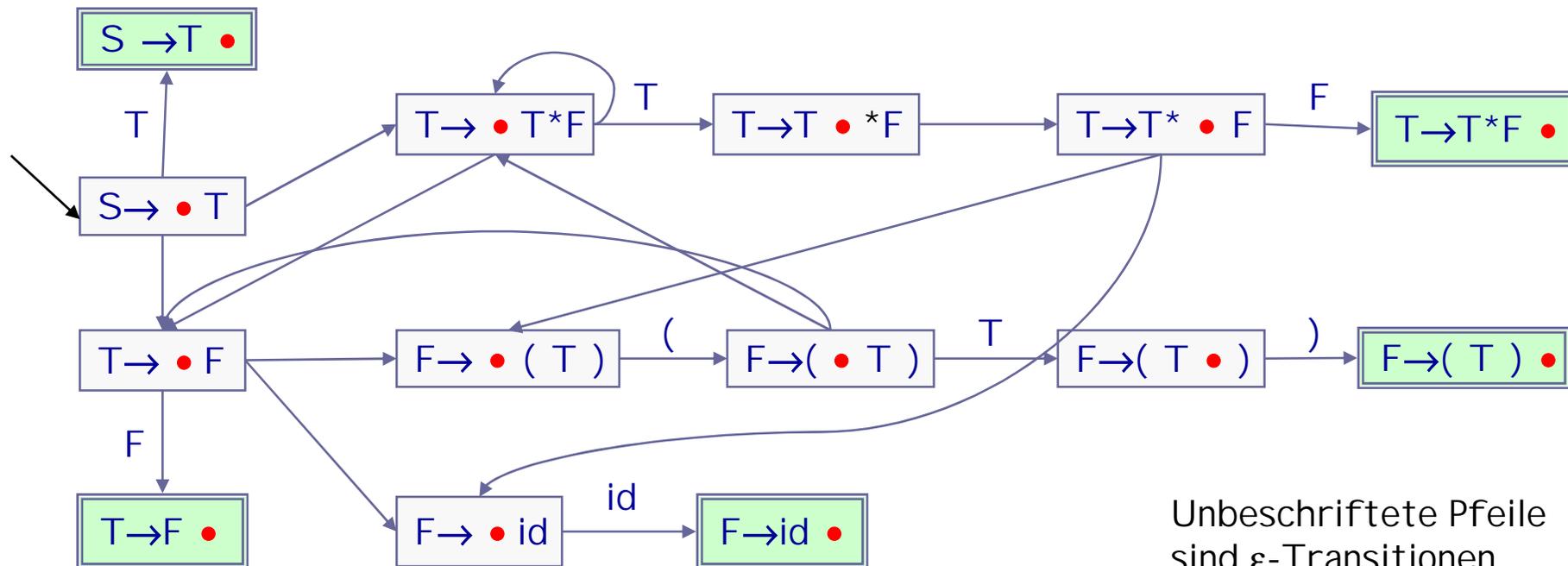
N(G) : Automat zu einer Grammatik G

Definiere NFA mit den items als Zuständen, den Terminalen und Nonterminalen der Grammatik als Alphabet :

- n S = Menge aller items
- n $\Sigma = V \cup T$
- n $S_0 = \bullet S$
- n $T = \{ X \rightarrow \beta \bullet \mid X \rightarrow \beta \in P \}$

G

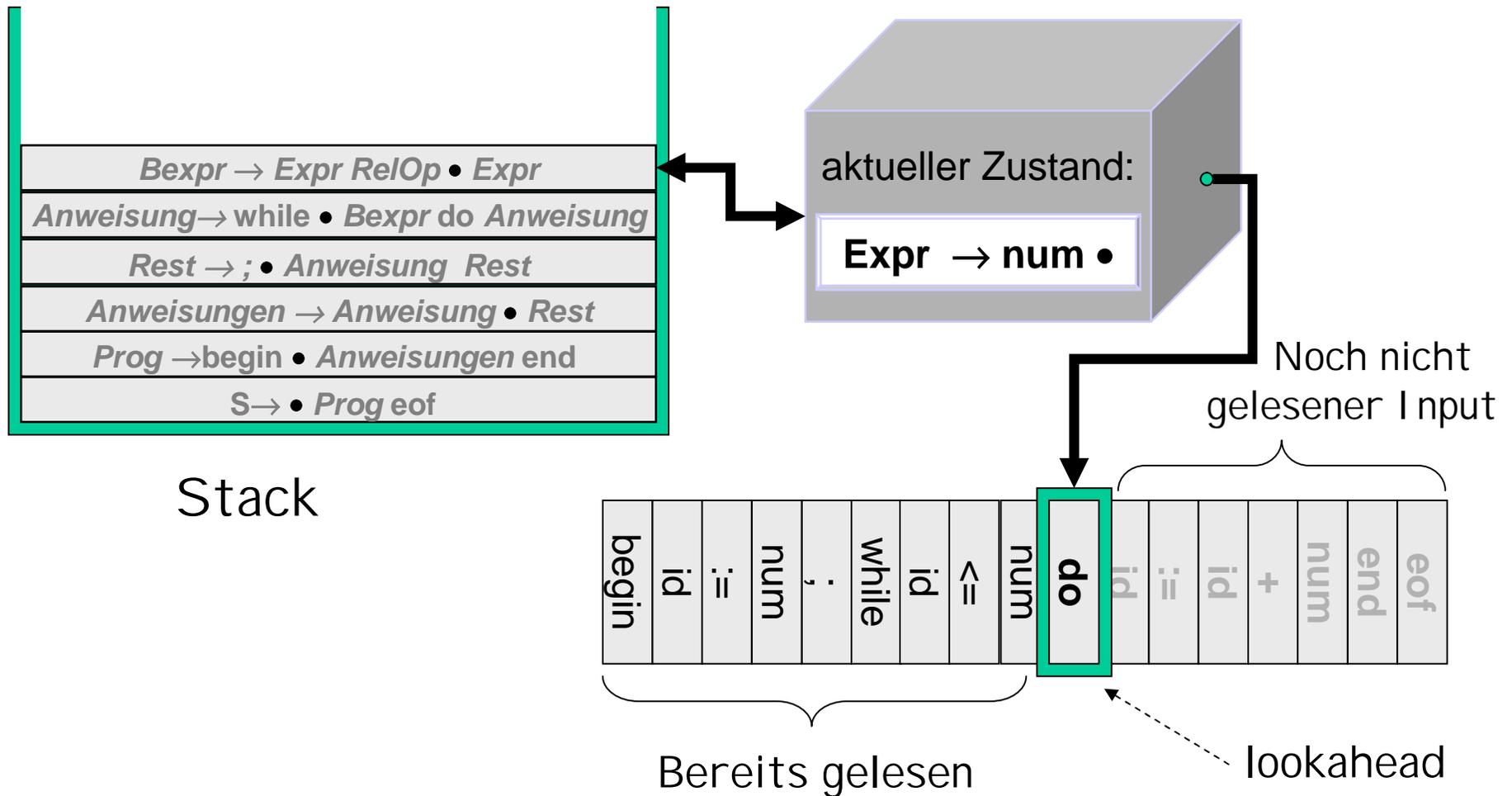
```
S → T
T → T * F
   | F
F → ( T )
   | id
```



Unbeschriftete Pfeile sind ϵ -Transitionen

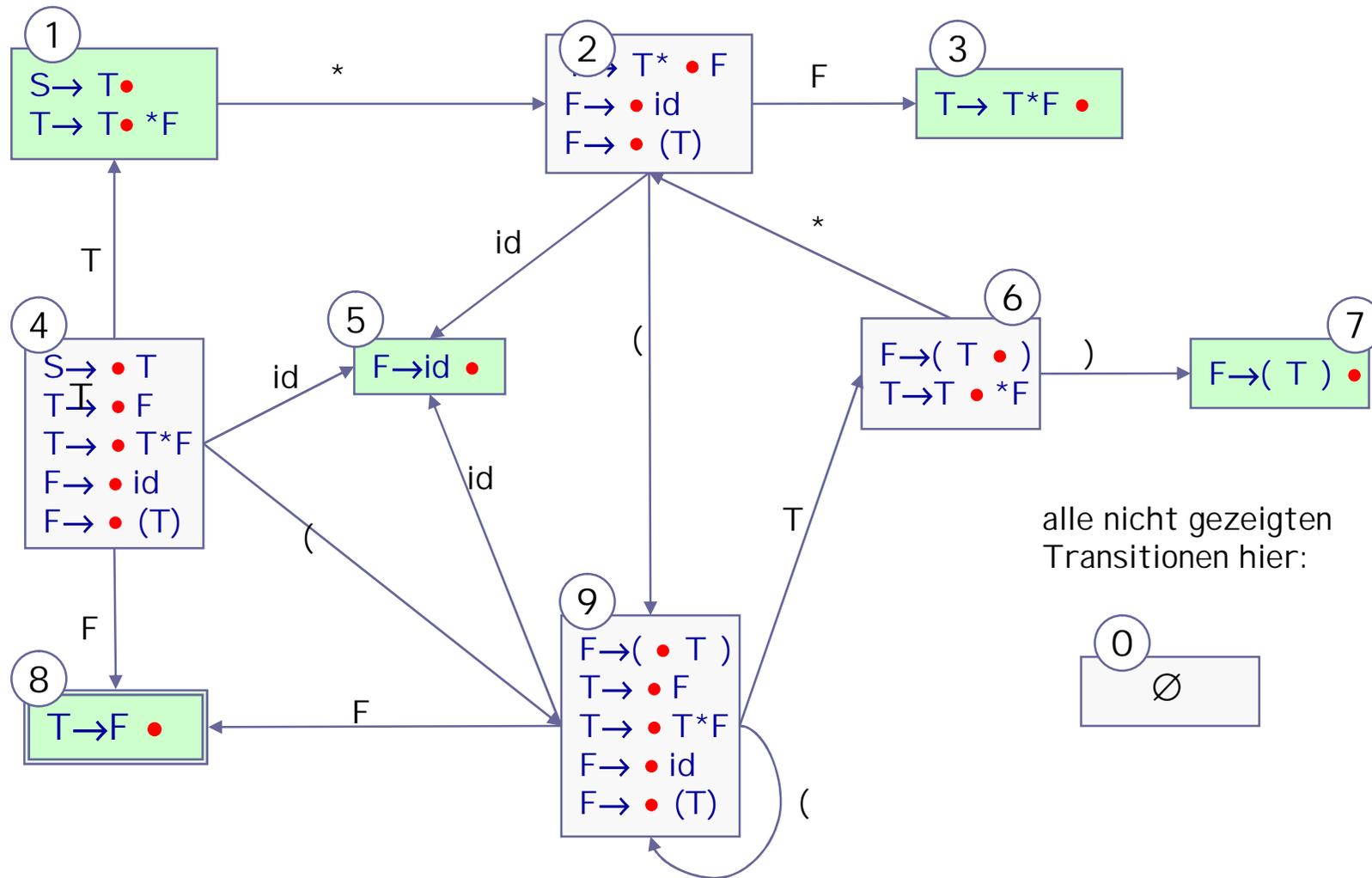


Nichtdeterministischer Stackautomat



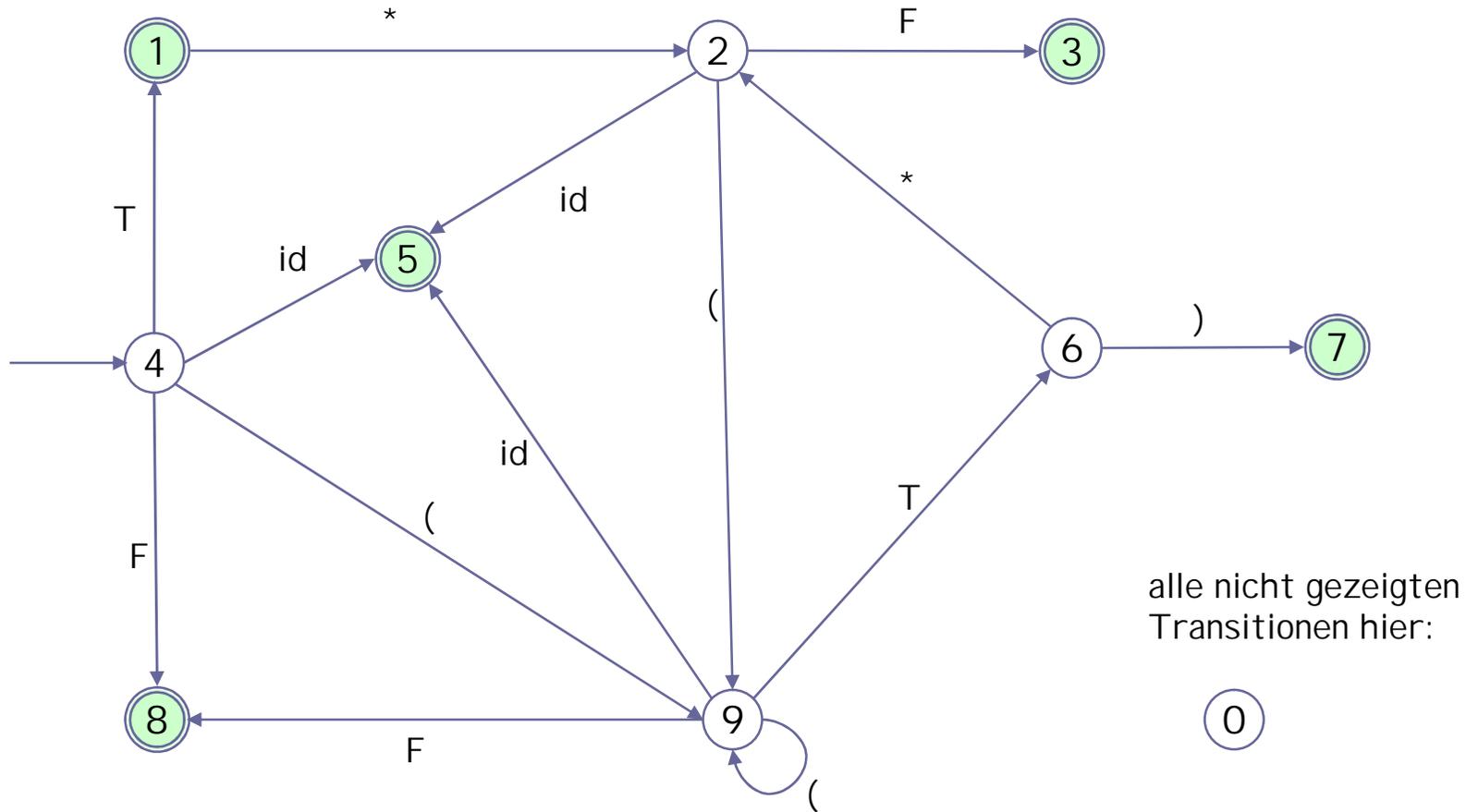


Deterministischer Automat



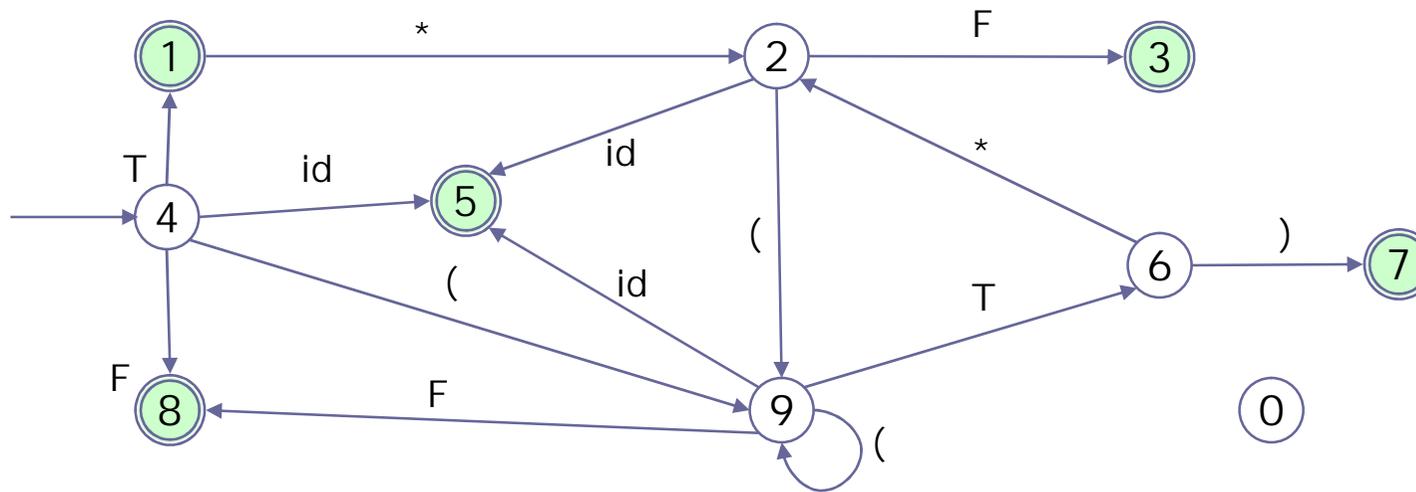


Deterministischer Automat





Deterministischer Automat



GoTo	T	F	id	*	()
1				2		
2		3	5			
4	1	8	5		9	
6				2		7
8						
9	6	8	5		9	

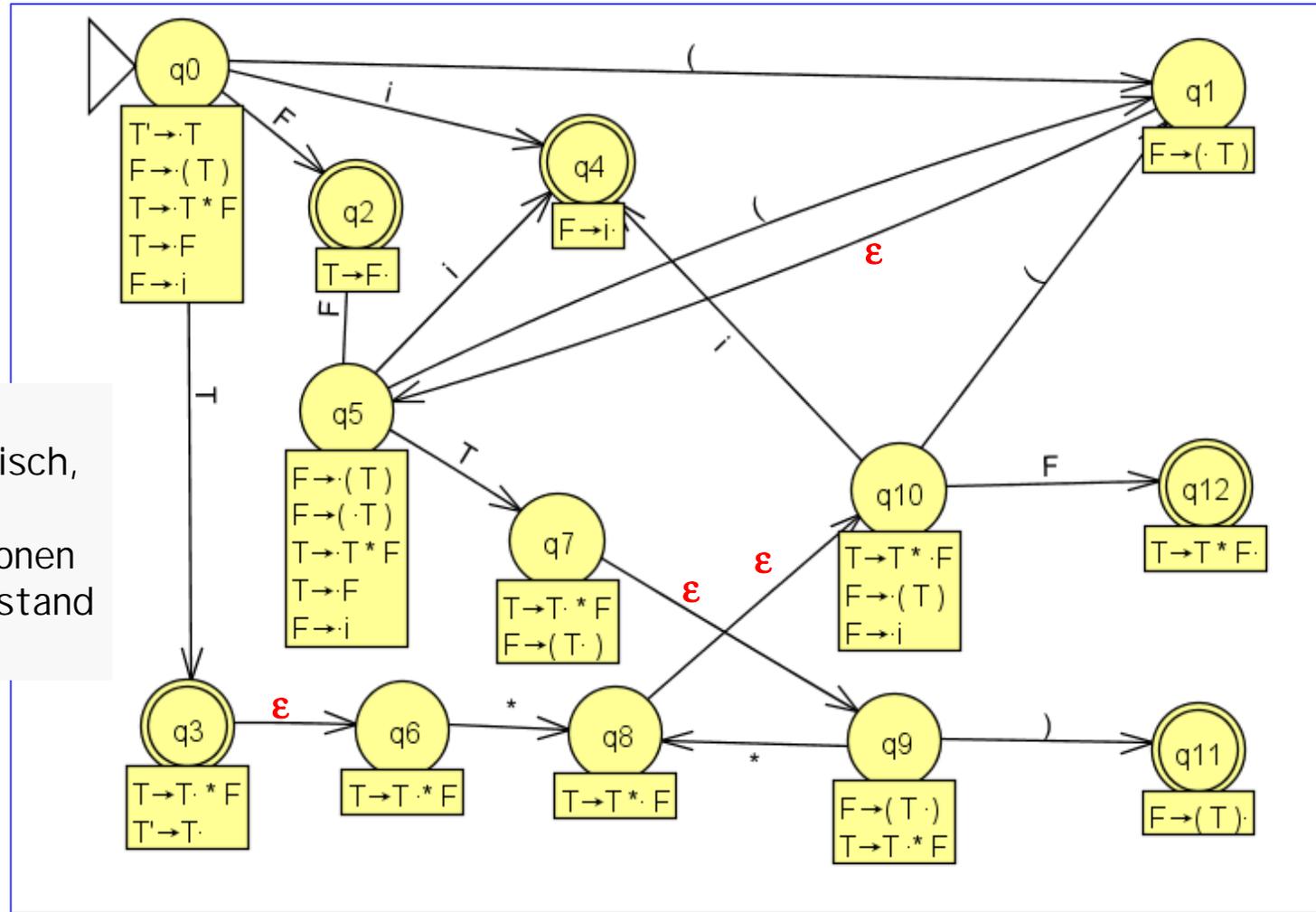


SLR-Automat in JFLAP erzeugt

SLR steht für Simple-LR

Der Automat ist nichtdeterministisch, weil er

- noch ϵ -Transitionen
- keinen Fehlerzustand enthält:



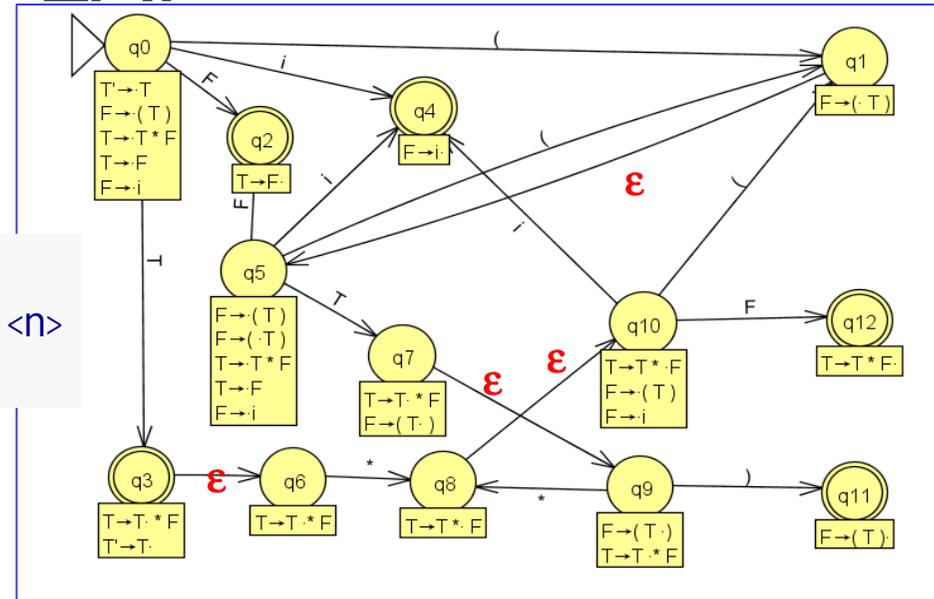


Parsertabelle in JFLAP

Einträge bedeuten:

$s\langle n \rangle$: Shift und gehe in Zustand $\langle n \rangle$

$r\langle k \rangle$: Reduce mit Regel Nr. $\langle k \rangle$

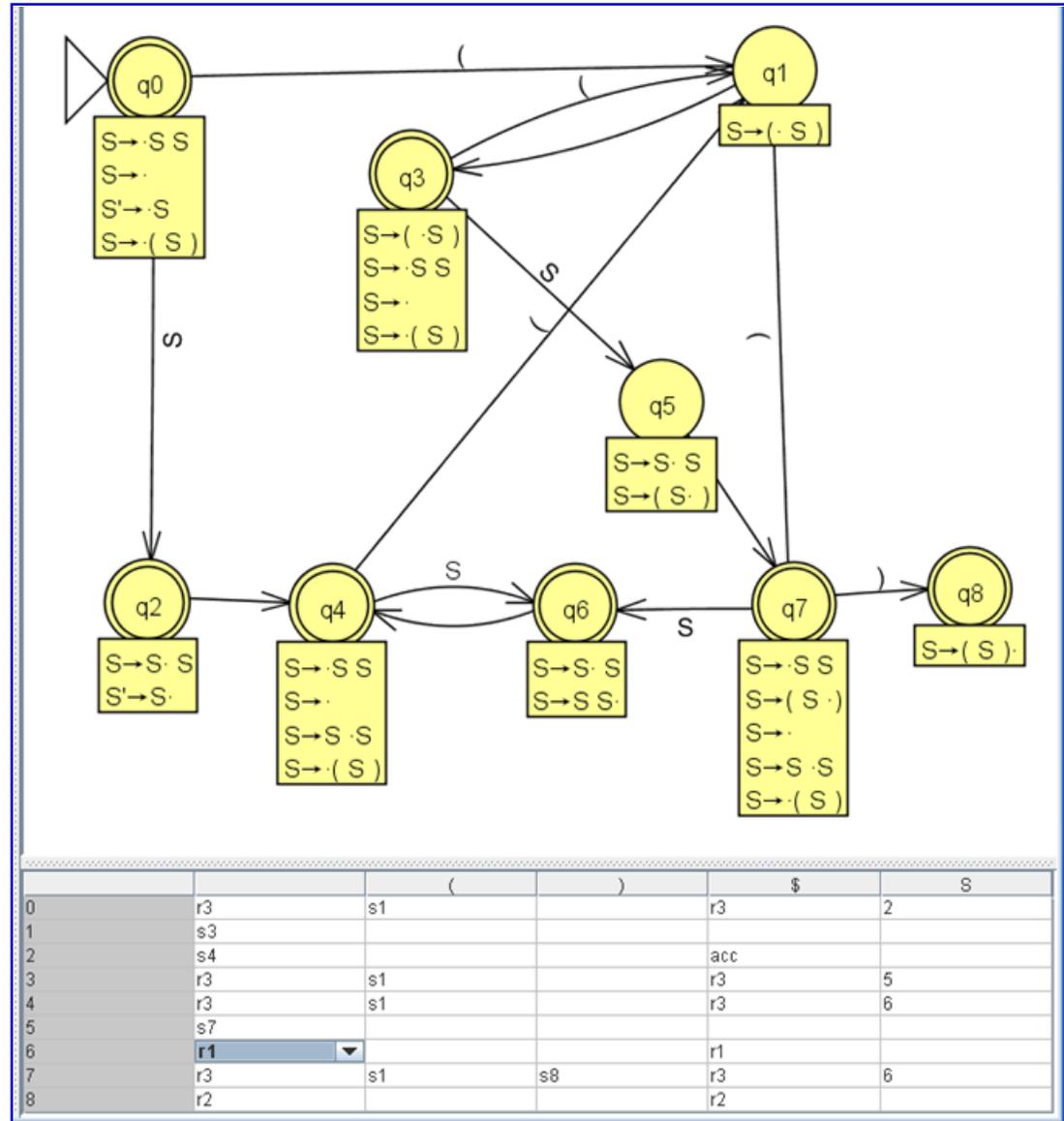


		()	*	i	\$	F	T
0		s1			s4		2	3
1	s5							
2	r2					r2		
3	s6					acc		
4	r3					r3		
5		s1			s4		2	7
6				s8				
7	s9							
8	s10							
9			s11	s8				
10		s1			s4		12	
11	r4					r4		
12	r1					r1		



SLR-Automat für Klammergrammatik

G_{klamm}
 $S' \rightarrow S$
 $S \rightarrow S S$
 $\quad | (S)$
 $\quad | \epsilon$





Die Notwendigkeit eines Stacks

Angenommen der Parser ist in einem Zustand (mit dem item)

$$\text{Term} \rightarrow (\bullet \text{Expr})$$

er muß nun also in einen Zustand springen, der das item

$$\text{Expr} \rightarrow \bullet \text{Expr} + \text{Term}$$

enthält. Nach dessen Abarbeitung kommt er in den Zustand :

$$\text{Expr} \rightarrow \text{Expr} + \text{Term} \bullet$$

Nun muß er aber zurückspringen in den Zustand

$$\text{Term} \rightarrow (\text{Expr} \bullet)$$

Offensichtlich muß bei jedem Sprung der alte Zustand auf einem Stack aufbewahrt werden und bei der Beendigung des items zurückgesprungen werden, verbunden mit einem POP des Stacks.



shift items / reduce items

Ein shift item hat die Form
 $A \rightarrow \alpha \bullet u \gamma$, mit $u \in T$.

shift \approx push

Drückt aus, daß ein u erwartet wird.

Im Parserzustand $Y = \{ \dots, A \rightarrow \alpha \bullet u \gamma \dots \}$ mit lookahead u wird dieses konsumiert und der Folgezustand $Z = \text{Goto}(Y, u) = \{ \dots, A \rightarrow \alpha u \bullet \gamma, \}$ auf den Stack gepusht.

Ein reduce item hat die Form $A \rightarrow \alpha \bullet$

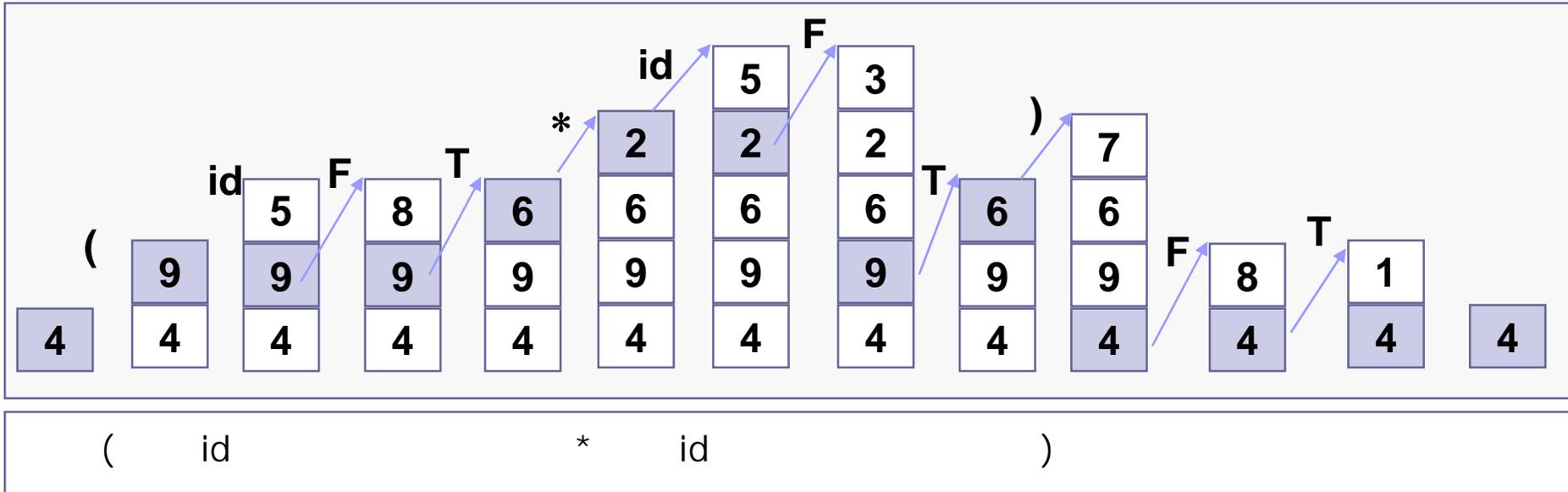
reduce \approx pop

Drückt aus, daß erfolgreich ein A erkannt wurde.

Auf dem Stack sollte ein Zustand $Y = \{ \dots, X \rightarrow \beta \bullet A \gamma, \dots \}$ liegen, darüber $|\alpha|$ weitere Zustände. Wir poppen also $|\alpha|$ viele Zustände und ersetzen sie durch $\text{Goto}(Y, A) = \{ \dots, X \rightarrow \beta \bullet A \gamma, \dots \}$.



Der Stack beim Parsen von $(id * id)$



GoTo	T	F	id	*	()
1				2		
2		3	5			
4	1	8	5		9	
6				2		7
8						
9	6	8	5		9	

1
 $S \rightarrow T \cdot$
 $T \rightarrow T \cdot * F$

2
 $T \rightarrow T * \cdot F$
 $F \rightarrow \cdot id$
 $F \rightarrow \cdot (T)$

3
 $T \rightarrow T * F \cdot$

4
 $S \rightarrow \cdot T$
 $T \rightarrow \cdot F$
 $T \rightarrow \cdot T * F$
 $F \rightarrow \cdot id$
 $F \rightarrow \cdot (T)$

5
 $F \rightarrow id \cdot$

6
 $F \rightarrow (T \cdot)$
 $T \rightarrow T \cdot * F$

7
 $F \rightarrow (T) \cdot$

8
 $T \rightarrow F \cdot$

9
 $F \rightarrow (\cdot T)$
 $T \rightarrow \cdot F$
 $T \rightarrow \cdot T * F$
 $F \rightarrow \cdot id$
 $F \rightarrow \cdot (T)$



Shift-Zustände - Reduce Zustände

Ein reduce-Zustand ist ein Zustand, der ein reduce item enthält.

Ein shift-Zustand ist ein Zustand, der ein shift item enthält.

Ein accept-Zustand ist ein reduce-Zustand der Form $\text{Start} \rightarrow \gamma \bullet$

Shift Zustände :

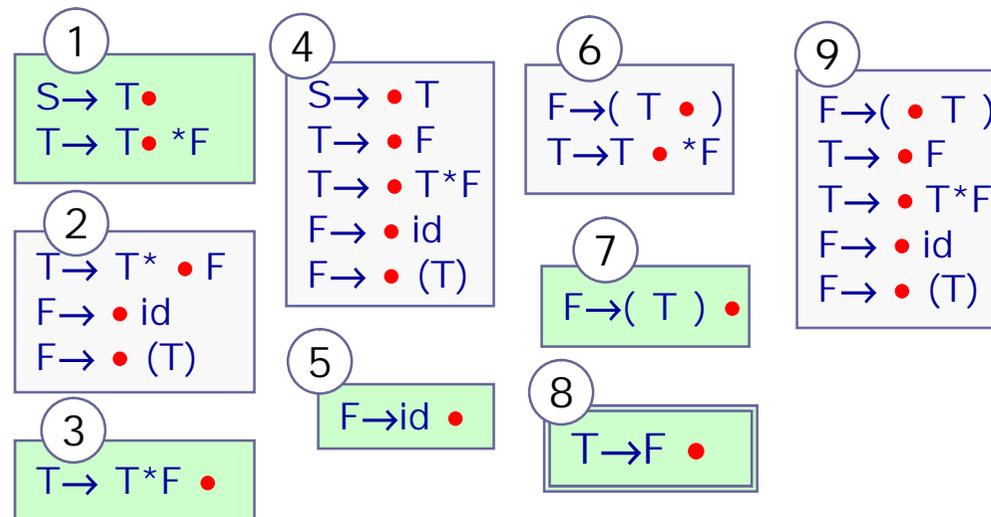
1, 2, 4, 6, 9

Reduce Zustände :

1, 3, 5, 7, 8

Accept Zustand:

8



Falls $L(X) \neq \emptyset$ für alle $X \in V$ dann ist jeder nichtleere Zustand ein Shift-Zustand oder ein Reduce-Zustand.

Der Zustand, der der leeren Menge von items entspricht kann immer als Fehlerzustand interpretiert werden.



Shift-Reduce Konflikte

Enthält ein Zustand Z sowohl ein shift item als auch ein reduce item, so muß das lookahead entscheiden, was zu tun ist

$$Z = \{ \dots, A \rightarrow \alpha \bullet u \gamma, \dots, B \rightarrow \beta \bullet, \dots \}$$

lookahead = u	\Rightarrow	shift,
Lookahead \in Follow(B)	\Rightarrow	reduce

Falls $u \in \text{Follow}(B)$ gibt es einen *Shift-Reduce-Konflikt* für u .
Pragmatisch entscheidet man sich meist für Shift



Reduce-Reduce Konflikte

Enthält ein Zustand Z mehrere reduce items so müssen die Follow-Mengen entscheiden, was zu tun ist

$$Z = \{ \dots, A \rightarrow \alpha \bullet, \dots, B \rightarrow \beta \bullet, \dots \}$$

lookahead \in Follow(A)	\Rightarrow	reduce mit $A \rightarrow \alpha$
Lookahead \in Follow(B)	\Rightarrow	reduce mit $B \rightarrow \beta$

Falls $u \in \text{Follow}(A) \cap \text{Follow}(B)$ gibt es einen *Reduce-Reduce-Konflikt* für u .

Solche Konflikte sollte man versuchen zu vermeiden.
Notfalls die Grammatik ändern !



Syntaxfehler

Ein Syntaxfehler im Input ist beim LR-Parsen immer zum frühestmöglichen Zeitpunkt erkennbar :

Ein Syntaxfehler liegt vor, wenn in Zustand Z mit lookahead u

Z kein shift-item $A \rightarrow \alpha \bullet u \gamma$ enthält und
für alle Reduce items $B \rightarrow \beta \bullet$ in Z gilt : $u \notin \text{Follow}(B)$

4

$S \rightarrow \bullet T$
 $T \rightarrow \bullet F$
 $T \rightarrow \bullet T * F$
 $F \rightarrow \bullet \text{id}$
 $F \rightarrow \bullet (T)$

id or (expected

1

$S \rightarrow T \bullet$
 $T \rightarrow T \bullet * F$

* or end-of-file expected

Gute Fehlermeldungen in Recursive Descent Parser einzubauen ist viel schwieriger!



Präzedenz und Assoziativität

Sowohl Präzedenz, als auch durch Assoziativität bedingte Mehrdeutigkeit der Grammatik führt zu Shift-Reduce-Konflikten.

G

S	→	E
E	→	E + E
		E * E
		(E)
		id

erzeugt u.a. die Zustände

Z₁

E	→	E + E •
E	→	E • + E
E	→	E • * E

shift-reduce Konflikt für + und für *.

Z₂

E	→	E * E •
E	→	E • + E
E	→	E • * E

shift-reduce Konflikt für + und für *.

Wie sind diese Konflikte zu lösen ?

Z₁

E	→	E + E •
E	→	E • + E
E	→	E • * E

Für * ⇒ shift (bewirkt Präzedenz von * über.)
für + ⇒ reduce (bewirkt Links-Klammerung von +)

Z₂

E	→	E * E •
E	→	E • + E
E	→	E • * E

Für * ⇒ reduce (bewirkt Links-Klammerung von.)
für + ⇒ reduce (garantiert Präzedenz von * über.)



LR(1) - items

LR Parsen ist mächtiger als Recursive Descent-Parsing.

Eine Grammatik heißt LR(0), falls es einen konfliktfreien SLR-Parser gibt.

Eine Sprache heißt LR(0), falls es eine LR(0) Grammatik für sie gibt.

Mächtigere Parsing-Methoden berücksichtigen noch lookaheads bei der Berechnung der Items:

Ein LR(1) item ist ein Paar $[A \rightarrow \alpha \bullet X \beta, a]$ aus einem LR(0) item und einem Terminalsymbol. Es gilt immer $a \in \text{Follow}(A)$

Anfangen mit

$[\text{Start} \rightarrow \bullet S, \text{eof}]$

konstruieren wir den NFA zur Grammatik, mit den Transitionen

$[A \rightarrow \alpha \bullet u \beta, a]$

\xrightarrow{u}

$[A \rightarrow \alpha u \bullet \beta, a]$

falls u ein Terminal

$[A \rightarrow \alpha \bullet B \beta, a]$

$\xrightarrow{\varepsilon}$

$[B \rightarrow \bullet \gamma, b]$

für jede Produktion

$B \rightarrow \bullet \gamma$

und jedes $b \in \text{First}(\beta a)$.



LALR(1)-Parsing

n Tatsachen

- .. LR(1) Parsing ist mächtiger als SLR-Parsing
- .. LR(1) Parsing hat weit mehr items, also umfangreichere Parsertabelle

n Kompromiss besteht aus

- .. Berechnung der LR(1) Zustände, danach
- .. Verschmelzung zweier Zustände, wenn sie
 - n die gleichen LR(0)-items, aber
 - n evtl. verschiedene lookahead Mengen
 - n besitzen

n Methode ist für praktische Zwecke mächtig genug

- .. implementiert in **yacc**
 - n **y**et **a**nother **c**ompiler-**c**ompiler



yacc

- n Yacc : Yet another Compiler Compiler (Parser Generator). Aus kontextfreier Grammatik erzeuge LALR(1) Parser.
- n C-Funktion `yyparse()`
- n yacc und lex arbeiten als Team. lex erkennt die Token, während yacc für die darauf aufbauende Grammatik verantwortlich ist.

```
# yacc Input File für Baby Deutsch

%start Satz

%Token  ARTIKEL NAME HAUPTWORT AUX VERB PUNKT
%%
Satz    :      Subjekt Praedikat Objekt PUNKT
        |      Subjekt Praedikat PUNKT
        ;

Subjekt :      ARTIKEL HAUPTWORT
        |      NAME
        ;

Objekt  :      Subjekt ;

Praedikat :      AUX VERB ;
          |      VERB
          |      error{"Verb erwartet"};
          ;
```

Startsymbol

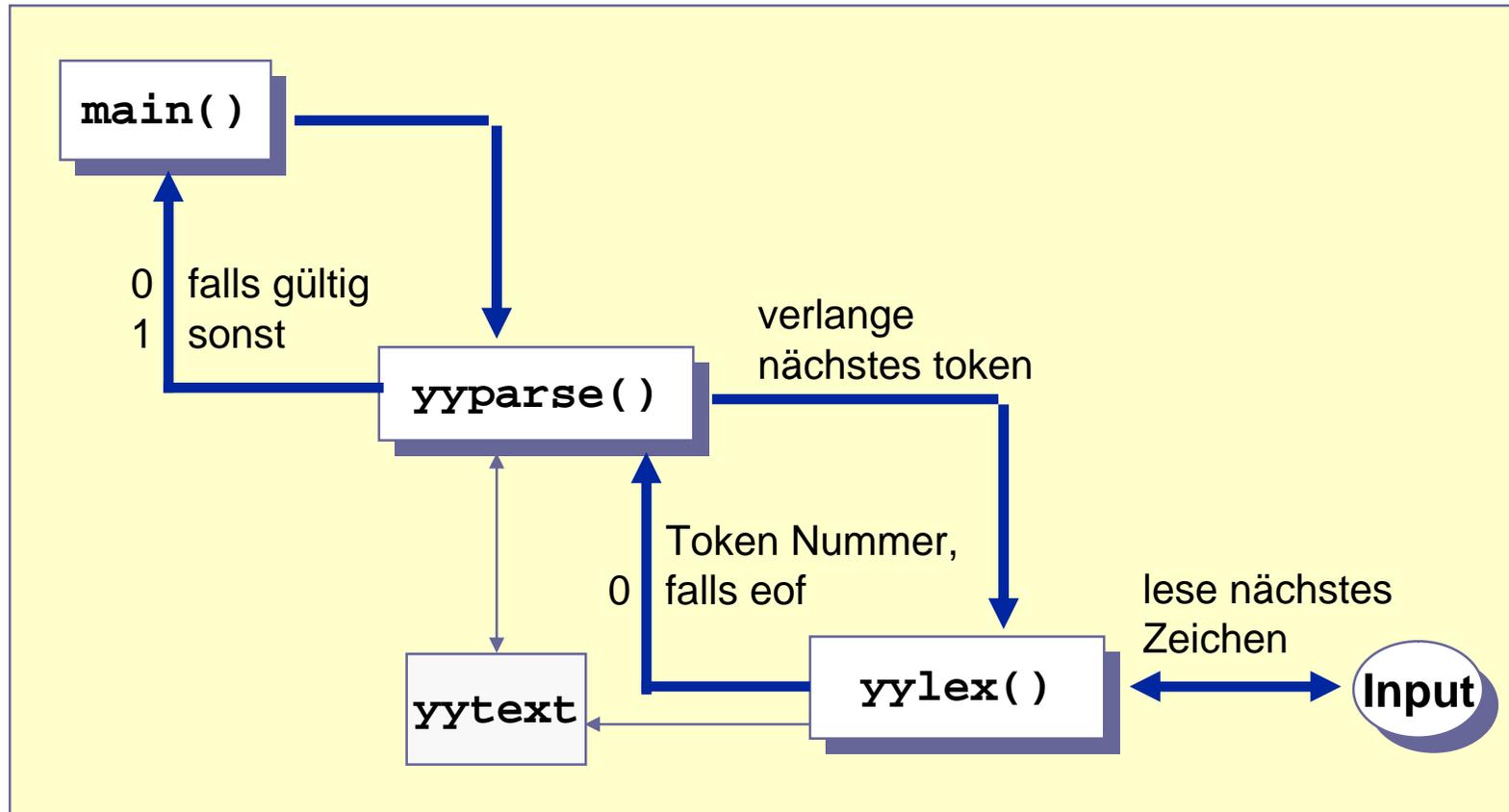
Für diese Token ist lex verantwortlich

Aktionen wie Terminal-symbole behandeln



Das Duo : lex und yacc

- n yacc erzeugt die Funktion `yyparse()`, lex die Funktion `yylex()`.
- n `yylex()` liefert die Nummer des Tokens. Der String aus dem dieses besteht ist immer in der (globalen) Variablen `yytext` vorhanden.





Beispiel: Ein Compiler für Expressions

Aufgabe: Compiler soll algebraische Ausdrücke in Postfix Notation verwandeln. Starte mit der Grammatik

G

Expr	→	Expr + Expr
		Expr - Expr
		Expr * Expr
		Expr / Expr
		(Expr)
		id
		num

Beseitigung der Linksrekursion : **nicht nötig**
Präzedenzstufen (Term, Faktor) : **nicht nötig**



Das lex file

- n C-Funktion `sscanf` wandelt eine Ziffernfolge (ASCII-String) in Zahlenwert (int).
- n Die int-Konstanten `PLUS`, `MINUS`, `TIMES`, `QUOT` werden später in dem zugehörigen yacc-file spezifiziert.

```
letter [a-zA-Z]
digit  [0-9]
%%
[ \t]+          ;
"+"           { return(PLUS) }
"-"           { return(MINUS) }
"*"           { return(TIMES) }
"/"           { return(QUOT) }
{digit}+       { sscanf(yytext,"%d", &yyval);
                return(NUM) }
{letter}({letter}|{digit})* {return(ID)}
%%
```



Das yacc file

- n Deklariere token `ID`, `NUM`, `PLUS`, `MINUS`, `TIMES`, `QUOT` als Operatoren
 - Wahlweise links-/rechts- oder nonassoziativ.
 - Bewirkt richtige Auflösung von shift-reduce Konflikten.
- n Reihenfolge wichtig (`PLUS`, `MINUS`, `TIMES`, `QUOT`)
 - bewirkt gewünschte Präzedenz.
- n `yylval` und `yytext` speichern Attribute der Token (hier `NUM` und `ID`)

```
%token ID, NUM
%left PLUS, MINUS
%left TIMES, QUOT
%start expr
%%
expr      : expr PLUS expr      { printf("add "); }
          | expr MINUS expr     { printf("sub "); }
          | expr TIMES expr     { printf("mult "); }
          | expr QUOT expr      { printf("div "); }
          | NUM                 { printf("%d ",yylval);}
          | ID                  { printf("%s ",yytext);}
          ;
%%
#include "lex.yy.c"
int main(){
    printf("Bitte geben Sie einen Ausdruck ein :\n");
    yyparse(); }

```

← Das von lex erzeugte C-Programm



Das fertige C-Programm

```
> lex in2post.l
```

```
  Erzeugt aus in2post.l ein C-File lex.yy.c.
```

```
> yacc in2post.y
```

```
  Erzeugt aus in2post.y ein C-File y.tab.c.
```

```
> cc -o myProg y.tab.c -ll -ly
```

```
  Compiliere mit Bibliotheksroutinen
```

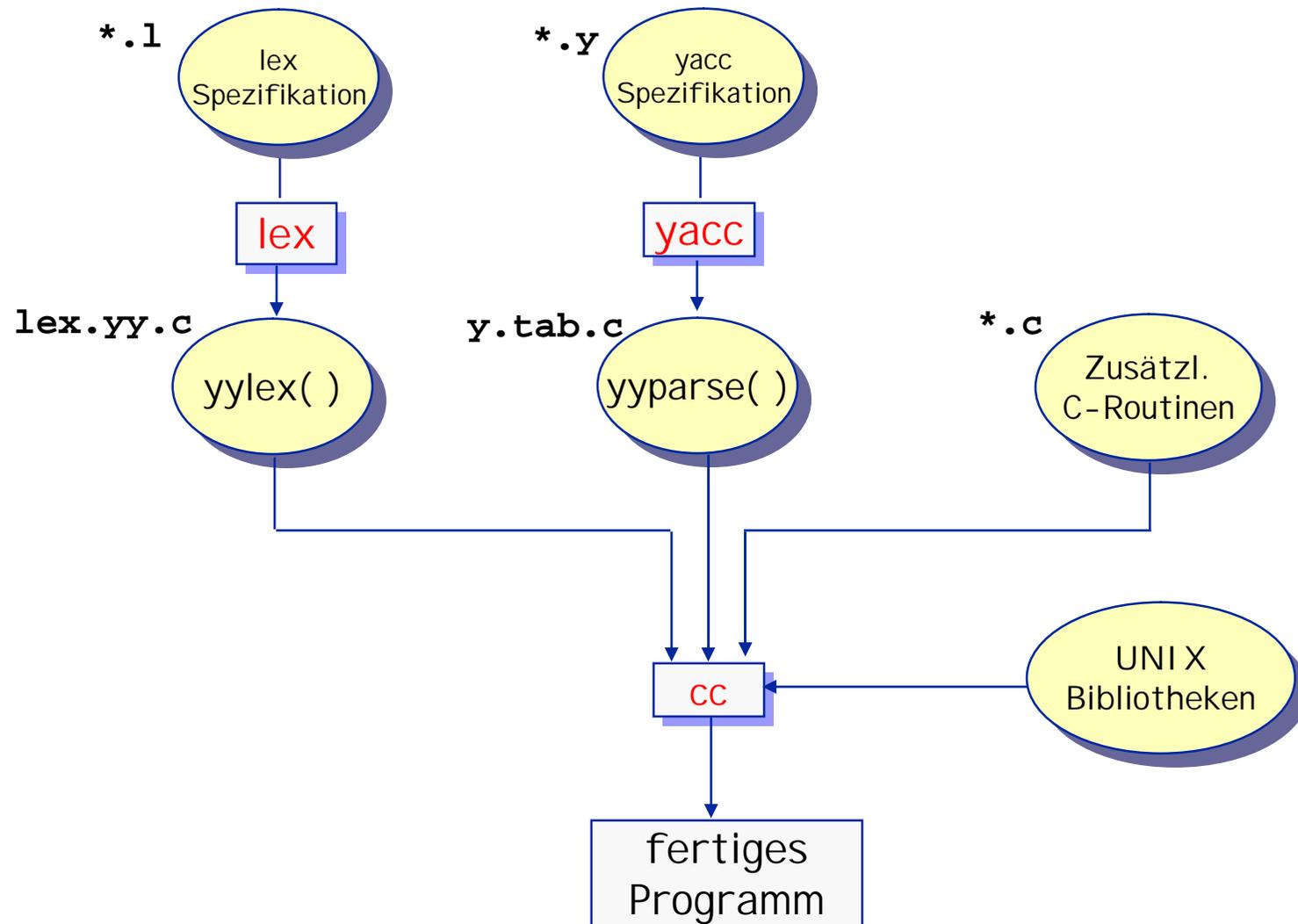
```
  n für lex -ll und
```

```
  n für yacc -ly
```

Endungen .l für lex/flex-Datei und .y für yacc/bison-Datei üblich



Das Zusammenspiel von lex & yacc





Der Aufbau eines Syntaxprüfers

Die Arbeitsweise eines Syntaxprüfers geschieht in 2 Phasen.
Diese können zeitlich verschachtelt ablaufen.

Phase 1: Scannen

Scanner zerlegt inputfile anhand regulärer Definitionen in Liste von Token.

File of char:

b	e	t	r	a	g		:	=		b	e	t	r	a	g		*		(1	+		z	i	n	s)	
---	---	---	---	---	---	--	---	---	--	---	---	---	---	---	---	--	---	--	---	---	---	--	---	---	---	---	---	--



File of token:

id	assign	id	*	(intConst	+	id)
----	--------	----	---	---	----------	---	----	---



Der Aufbau eines Syntaxprüfers

Phase 2: Parzen

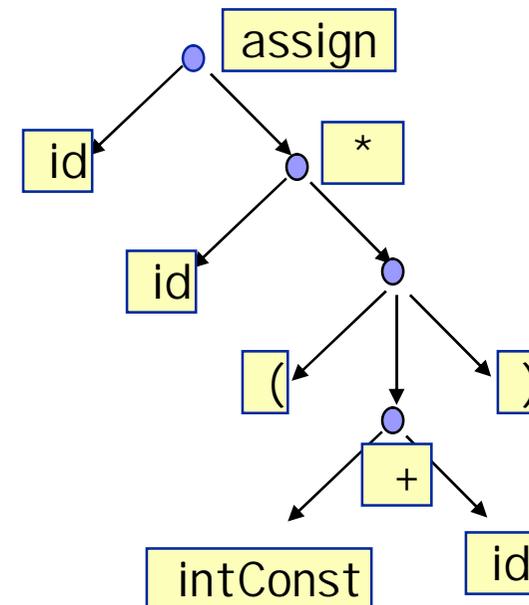
Parser versucht Tokenliste zu einem Herleitungsbaum zu gruppieren

File of token:

id
assign
id
*
(
intConst
+
id
)



Parse Tree :

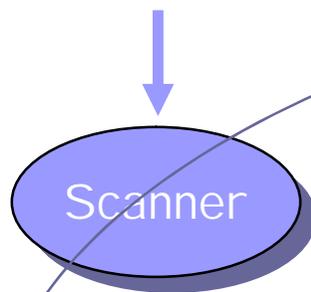




Symboltabelle

- n Um aus einem Syntaxbaum Code zu erzeugen muß die Identität einiger Token bekannt sein.
 - .. welche Bezeichner sind identisch
 - .. welchen Wert eine intConst , etc.
- n Scanner legt diese Information in Symboltabelle ab und liefert dem Parser die Token samt Zeiger in diese Tabelle.

```
b e t r a g   : =   n l
b e t r a g   *   (   1   n l
+ z i n s   )
```



Name	Typ	Sp. Platz
betrag	Real	17F4
zins	Real	17F8
x	Integer	201C
test	Boolean	C011

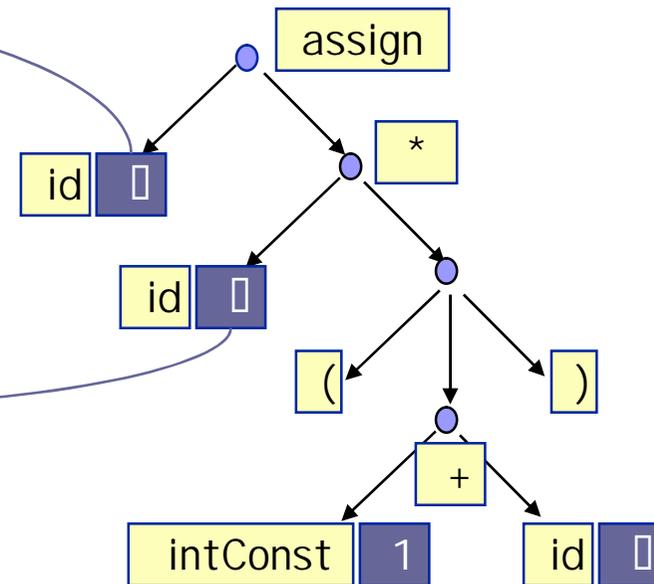
```
id  [ ]  assign  id  [ ]  *  (  intConst  1  +  id  [ ]  )
```



Symboltabelle

Name	Typ	Sp. Platz
betrag	Real	17F4
zins	Real	17F8
x	Integer	201C
test	Boolean	C011

Parse Tree :



Gewisse Token im Syntaxbaum haben einen Link in die Symboltabelle



Stackprozessor

- n Aus Programmtext soll Code für eine Maschine erzeugt werden.
 - .. Einfaches Maschinenmodell: Stackprozessor.

Typische Befehle eines Stackprozessors :

PUSH <num>

Lege den Zahlenwert <num> auf dem Stack ab

LOAD <hex>

Lege den Inhalt von Adresse <hex> auf dem Stack ab

STORE <hex>

Speichere den Top des Stacks an Adresse <hex>. Der Stack wird dabei gepopped.

ADD

Ersetze das oberste Element des Stacks durch die Summe der beiden obersten

MULT

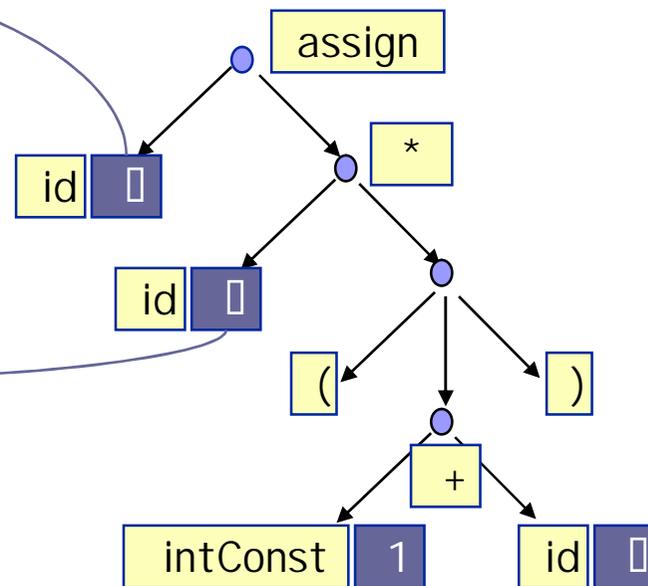
Ersetze das oberste Element des Stacks durch das Produkt der beiden obersten



Code-Erzeugung

Name	Typ	Sp. Platz
betrag	Real	17F4
zins	Real	17F8
x	Integer	201C
test	Boolean	C011

Parse Tree :



Code für einen Stackprozessor

```
LOAD 17F4
PUSH 1
LOAD 17f8
ADD
MULT
STORE 17F4
```



Code Erzeugung aus dem Parser

- n Code-Erzeugung durch semantische Aktionen:

```
%{ int* loc }%

expr      : expr + expr      { printf("ADD \n"); }
          | expr * expr      { printf("MULT\n"); }
          | intConst         { printf("PUSH ");
                              printf("%d\n",yylval);}
          | ID                { printf("LOAD ");
                              printf(lookup(yttext));}
          ;

stmt      : ID { loc=lookup(yttext);}
          ASSIGN expr { printf("STORE ");
                       printf(&loc); }
          | ... etc. ...

%%
#include lex.yy.c
...
```