

Formale Grundlagen der Informatik

10

Backus-Naur-Form
Kontextfreie Grammatiken
Compiler • Parsing

Recap: Hierarchiebeziehungen

$$\mathcal{L}(\text{REG}) \subset \mathcal{L}(\text{REC}) \subset \mathcal{L}(\text{RE})$$

echt wegen $\{ a^n b^n \mid n \geq 0 \}$

echt wegen L_u

Die Sprache $\{ a^n b^n \mid n \geq 0 \}$

- ist relevant für die Beschreibung syntaktischer korrekter Eigenschaften, wie sie vom Compiler geprüft werden
- ist nicht regulär \rightarrow kann nicht mit regulären Ausdrücken spezifiziert werden ☹
- kann von DTM akzeptiert werden, aber das Wortproblem und das Äquivalenzproblem sind für DTMs unentscheidbar ☹
- Benötigen Mechanismen, die stärker sind als reguläre Ausdrücke und schwächer als DTMs!

Syntax-Definition von HTML5

tag-open := '<' tag-name ws* attr-list? ws* '>'

tag-empty := '<' tag-name ws* attr-list? ws* '/>'

tag-close := '</' tag-name ws* '>'

attr-list := (ws+ attr)*

attr := attr-empty | attr-unquoted | attr-single-quoted | attr-double-quoted

attr-empty := attr-name

attr-unquoted := attr-name ws* = ws* attr-unquoted-value

attr-single-quoted := attr-name ws* = ws* ' attr-single-quoted-value '

attr-double-quoted := attr-name ws* = ws* " attr-double-quoted-value "

tag-name := (alphabets | digits)+

attr-name := /[^\s">/=\\p{Control}]+/

attr-unquoted-value := /[^\s"'=<>`]+/

attr-single-quoted-value := /[^\']*+/

attr-double-quoted-value := /[^\"]*/

alphabets := /[a-zA-Z]/

digits := /[0-9]/

ws := /\s/

*Definition von syntaktischen
Kategorien in (erweiterter)
Backus-Naur- Form*

Backus-Naur-Form (BNF)

- entwickelt von John Backus und Peter Naur
- Definition der Syntax durch einen Satz von Ersetzungsregeln der Form

Kategorie := **Ausdruck**

- wobei **Kategorie** durch einen Namen angegeben wird und
- **Ausdruck** eine durch | getrennte Folge von Wörtern ist, wobei jedes Wort aus Buchstaben, wie sie im Programm vorkommen, und Kategorien besteht

`attr := attr-empty | attr-unquoted | attr-single-quoted | attr-double-quoted`

`tag-open := '<' tag-name ws* attr-list? ws* '>'`

- In erweiterter BNF dürfen die Wörter reguläre Ausdrücke sein.
(*betrachten wir hier nicht weiter*)

Backus-Naur-Form (BNF)

- Beispiel aus der Definition von Programmiersprachen:

statement := assignment | loop | branch

assignment := identifier '=' arithmetic-expression

arithmetic-expression := ...

- Die Ersetzungsregel für statement kann in drei Regeln aufgeteilt werden.
- Dann hat jede Regel eine linke Seite (Kategorie) und eine rechte Seite, durch die die Kategorie ersetzt („verfeinert“) werden kann.
 - Sie sind also Paare der Form (linke-Seite, rechte-Seite).

Kontextfreie Grammatiken

- sind eine Abstraktion von BNF
- Komponenten:
 - Alphabet der nichtterminalen Symbole („Kategorien“)
 - müssen weiter ersetzt werden, damit ein syntaktisch korrektes Programm entsteht
 - Alphabet der terminalen Symbole (wie sie im Programmen vorkommen)
 - müssen also nicht weiter ersetzt werden
 - endliche Menge von Ersetzungsregeln
 - beschreiben, wie ein Nichtterminal ersetzt werden kann
 - ein Startsymbol aus der Menge der Nichtterminale
 - die („größte“/allgemeinste) syntaktische Kategorie, die als erstes zu ersetzen ist
 - z.B. Programm bzw. HTML-Dokument

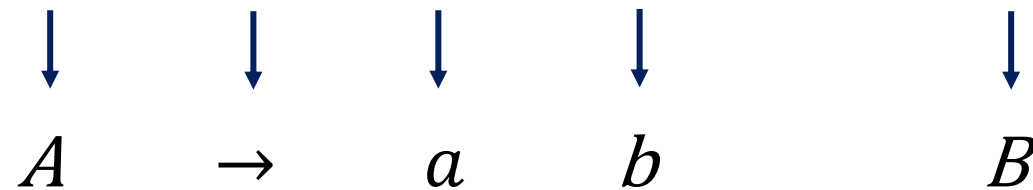
Kontextfreie Grammatiken - Definition

Definition: Eine **kontextfreie Grammatik (kfG)** ist ein Quadrupel $G = (N, T, P, S)$, wobei

- N und T Alphabete sind, $N \cap T = \emptyset$,
deren Symbole **Nichtterminale** bzw. **Terminale** heißen;
- $P \subseteq N \times (N \cup T)^*$ eine endliche Menge von **Ersetzungsregeln**
(Regeln/Produktionen) ist,
deren Elemente in der Form $A \rightarrow \alpha$ (statt (A, α)) geschrieben werden;
- $S \in N$ das **Startsymbol** (oder Axiom) ist.

Kontextfreie Regeln

- *Beispiel:* assignment := identifier '=' arithmetic-expression



mit $A, B \in N$ und $a, b \in T$

- *Konvention:*
 - Nichtterminale durch Großbuchstaben angeben,
 - Terminale durch Kleinbuchstaben,
 - Wörter, in denen Terminale und Nichtterminale vorkommen, durch griechische Buchstaben

Kontextfreie Grammatik - Beispiel

- einfache arithmetische Ausdrücke mit den Operatoren + und *:
 - Nichtterminal E für „expression“
 - Terminale + und *
 - Terminal **id** für „identifizier“ als elementarer Ausdruck
- $G = (\{E\}, \{+, *, \mathbf{id}\}, P, E)$, wobei P aus folgenden Regeln besteht:

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow (E)$$

$$E \rightarrow \mathbf{id}$$

oder: $E \rightarrow E + E \mid E * E \mid (E) \mid \mathbf{id}$

Ableitungen

- Sei $G = (N, T, P, S)$ eine kontextfreie Grammatik.
- Ein direkter Ableitungsschritt in G ist wie folgt definiert:
Für Wörter $\alpha, \beta \in (N \cup T)^*$ gilt $\alpha \xRightarrow{G} \beta$ falls

- $\alpha = \gamma_1 A \gamma_2$
- $\beta = \gamma_1 v \gamma_2$
- $A \rightarrow v \in P$

Der Index G kann weggelassen werden, wenn aus dem Kontext eindeutig klar.

- Es gilt $\alpha \xRightarrow{G}^m \beta$ falls Wörter $\alpha_0, \alpha_1, \dots, \alpha_m$ existieren, so dass

$$\alpha = \alpha_0 \xRightarrow{G} \alpha_1 \xRightarrow{G} \dots \xRightarrow{G} \alpha_m = \beta.$$

- Es gilt $\alpha \xRightarrow{G}^* \beta$ falls $\alpha \xRightarrow{G}^m \beta$ für ein $m \geq 0$.

Wir sagen dann, dass β aus α abgeleitet wird.

$\xRightarrow{}$ ist die reflexive und transitive Hülle von \Rightarrow .*

Ableitungen – Beispiele

$$E \rightarrow E + E \mid E * E \mid (E) \mid \text{id}$$

- $E \Rightarrow E + E$

- $E \Rightarrow E + E \Rightarrow E + E * E$

- $E \Rightarrow \text{id}$

- $E \Rightarrow E * E \Rightarrow (E) * E \Rightarrow (E + E) * E$
 $\Rightarrow (\text{id} + E) * E \Rightarrow (\text{id} + \text{id}) * E \Rightarrow (\text{id} + \text{id}) * \text{id}$

also

$$E \xRightarrow{*} E + E$$

$$E \xRightarrow{*} E + E * E$$

$$E \xRightarrow{*} \text{id}$$

$$E \xRightarrow{*} (\text{id} + \text{id}) * \text{id}$$

terminale
Wörter

Sprache einer kontextfreien Grammatik

- Falls $S \xRightarrow{*} \alpha$ für ein $\alpha \in (N \cup T)^*$, dann heißt α **Satzform** von G .
- Terminale Satzformen
 - erlauben keine weiteren Ableitungsschritte
 - enthalten nur noch Symbole, die in Programmen (als Token) so vorkommen
 - sind Elemente der durch die Grammatik definierten Sprache

- **Definition:** Sei $G = (N, T, P, S)$ eine kontextfreie Grammatik.

Die von G **erzeugte Sprache** ist definiert als

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

- Eine Sprache heißt kontextfrei, wenn sie von einer kfG erzeugt wird.

Kontextfreie Sprache - Beispiel

- Konstruieren eine kfG für $L = \{ a^n b^n \mid n \geq 0 \}$.
 - *Idee*: Erzeugen von einem Nichtterminal gleichzeitig ein a nach links und ein b nach rechts: $A \rightarrow aAb$.
 - Dies kann wiederholt werden, bis durch Anwendung von $A \rightarrow \varepsilon$ ein Terminalwort entsteht.
 - Da ε in der Sprache L ist, muss A nicht vom Axiom unterschieden werden.
- $G = (\{A\}, \{a, b\}, \{A \rightarrow aAb, A \rightarrow \varepsilon\}, A)$

$$\begin{array}{ccccccc}
 A & \Rightarrow & aAb & \Rightarrow & aaAbb & \Rightarrow & \dots \Rightarrow a^n Ab^n \Rightarrow \dots \\
 \Downarrow & & \Downarrow & & \Downarrow & & \Downarrow \\
 \varepsilon & & ab & & aabb & & a^n b^n \quad \dots
 \end{array}$$

Korrektheitsbeweis

- $L = \{ a^n b^n \mid n \geq 0 \}$, $G = (\{A\}, \{a, b\}, \{A \rightarrow aAb, A \rightarrow \varepsilon\}, A)$; *Beh.:* $L(G) = L$
- Zeigen durch Induktion: $A \xRightarrow{n} w$ **gdw.** $w \in \{a^n Ab^n, a^{n-1} b^{n-1}\}$ für alle $n \geq 1$.
- **I.A.:** $n = 1$. $A \Rightarrow w$ durch Anwendung einer der Regeln, die das Axiom A ersetzen.
Somit $A \Rightarrow aAb$, $A \Rightarrow \varepsilon$ und keine weiteren direkten Ableitungen.
Also $A \xRightarrow{1} w$ **gdw.** $w \in \{aAb, \varepsilon\} = \{a^1 Ab^1, a^0 b^0\}$.
- **I.V.:** $A \xRightarrow{n} w$ **gdw.** $w \in \{a^n Ab^n, a^{n-1} b^{n-1}\}$
- **I.S.:** Keine Ableitung von $a^{n-1} b^{n-1}$, also genau zwei Ableitungen der Länge $n + 1$:
 - $A \xRightarrow{n} a^n Ab^n \Rightarrow a^{n+1} Ab^{n+1}$ durch Anwendung von $A \rightarrow aAb$
 - $A \xRightarrow{n} a^n Ab^n \Rightarrow a^n b^n$ durch Anwendung von $A \rightarrow \varepsilon$
 Also gilt $A \xRightarrow{n+1} w$ **gdw.** $w \in \{a^{n+1} Ab^{n+1}, a^n b^n\}$.
- $A \xRightarrow{*} w$ und $w \in T^*$ (also $w \in L(G)$) **gdw.** $w \in \{ a^n b^n \mid n \geq 0 \} = L$. ■

Korrektheitsbeweise

Für den Nachweis von $L(G) = L$ ist es auch hier wieder wichtig, beide „Richtungen“ zu beachten:

1. Es können alle Wörter von L erzeugt werden, also $L \subseteq L(G)$.
2. Es kann kein Wort erzeugt werden, dass nicht aus L ist, also $L(G) \subseteq L$.

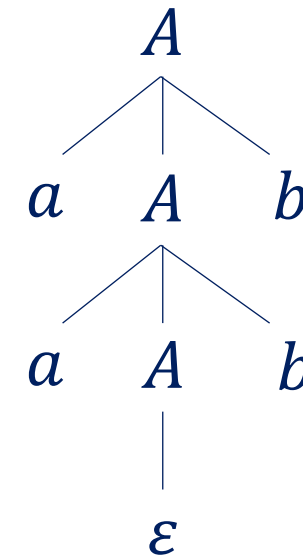
Sprachfamilie

- $\mathcal{L}(\mathbf{CF}) = \{ L \mid L = L(G) \text{ für eine kfG } G \}$
- $\{ a^n b^n \mid n \geq 0 \} \in \mathcal{L}(\mathbf{CF})$

Ableitungsbäume

- Wurzel markiert mit Axiom
- Innere Knoten mit Nichtterminalen
- Wird Regel $A \rightarrow \alpha$ angewendet, dann hat der mit A markierte Knoten $|\alpha| = k$ viele Kindknoten, die (v.l.n.r.) mit den Buchstaben von α markiert sind.
- Bei $A \rightarrow \varepsilon$ gibt es einen Kindknoten, der mit ε markiert ist.
- Blätter sind mit Terminalen oder ε markiert.
V.l.n.r. gelesen ergibt sich das abgeleitete Wort.

Bsp. $A \rightarrow aAb \mid \varepsilon$,
 $A \Rightarrow aAb \Rightarrow aaAbb \Rightarrow a^2b^2$:

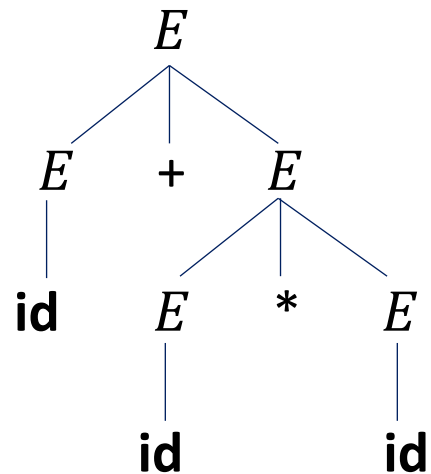


Struktur in den Wörtern

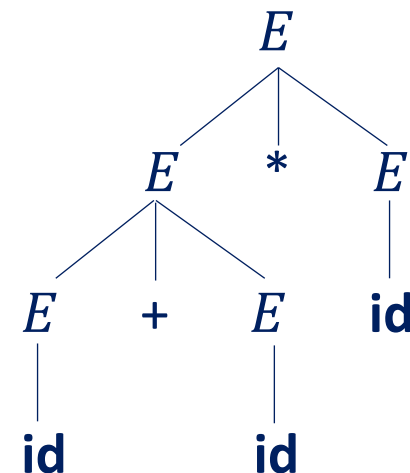
- Ableitungsbäume offenbaren die Struktur in den abgeleiteten Wörtern.
- Bedeutsam für korrekte Übersetzung durch Compiler.
- *am Beispiel $E \rightarrow E + E \mid E * E \mid (E) \mid \text{id}$*
 - Es gibt zwei strukturell verschiedene Ableitungen für das terminale Wort **id + id * id**:

$$E \Rightarrow E + E \Rightarrow E + E * E \xRightarrow{3} \text{id} + \text{id} * \text{id}$$

$$E \Rightarrow E * E \Rightarrow E + E * E \xRightarrow{3} \text{id} + \text{id} * \text{id}$$



*übliche
Vorrang-
regel*

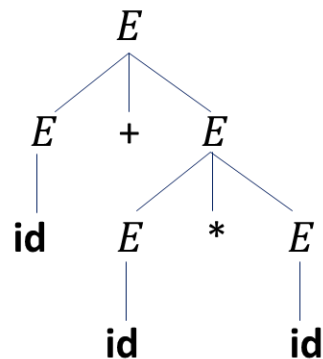


*verkehrte
Vorrang-
regel*

Eindeutigkeit der Zuordnung

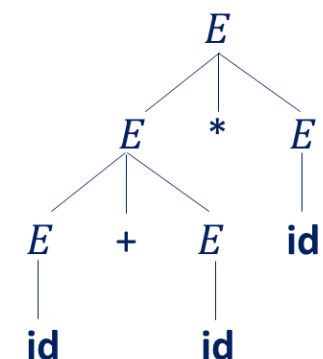
Zu jeder Ableitung gehört ein eindeutig bestimmter Ableitungsbaum.

$$E \Rightarrow E + E \Rightarrow E + E * E \xRightarrow{3} \text{id} + \text{id} * \text{id}$$



$E \Rightarrow E + E \Rightarrow \text{id} + E \Rightarrow \text{id} + E * E \Rightarrow \text{id} + \text{id} * E \Rightarrow \text{id} + \text{id} * \text{id}$ und
 $E \Rightarrow E + E \Rightarrow \text{id} + E \Rightarrow \text{id} + E * E \Rightarrow \text{id} + E * \text{id} \Rightarrow \text{id} + \text{id} * \text{id}$ und
 $E \Rightarrow E + E \Rightarrow E + E * E \Rightarrow E + E * \text{id} \Rightarrow E + \text{id} * \text{id} \Rightarrow \text{id} + \text{id} * \text{id}$
 werden z.B. durch denselben Ableitungsbaum beschrieben

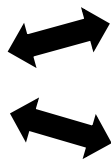
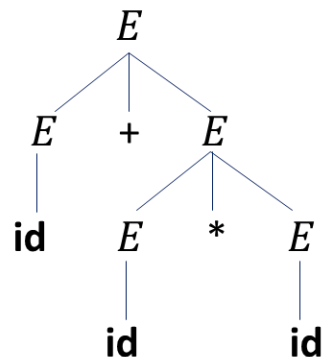
$$E \Rightarrow E * E \Rightarrow E + E * E \xRightarrow{3} \text{id} + \text{id} * \text{id}$$



Umgekehrt ist die Reihenfolge der Ableitungsschritte i.A. nicht aus dem Baum ersichtlich!

Links- und Rechtableitungen

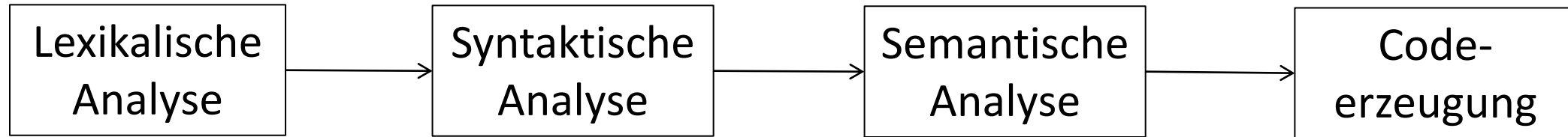
- In einer **Linksableitung** wird stets das Nichtterminal ersetzt, das in der Satzform am weitesten links auftritt; für jeden direkten Ableitungsschritt gilt:
 $uA\gamma \Rightarrow_L uv\gamma$ falls $uA\gamma \Rightarrow uv\gamma$ und $u \in T^*$.
- In einer **Rechtableitung** wird stets das Nichtterminal ersetzt, das in der Satzform am weitesten rechts auftritt; für jeden direkten Ableitungsschritt gilt:
 $\gamma Au \Rightarrow_R \gamma vu$ falls $\gamma Au \Rightarrow \gamma vu$ und $u \in T^*$.
- Zu jedem Ableitungsbaum gehört eine eindeutig bestimmte Links- und eine eindeutig bestimmte Rechtableitung und umgekehrt.



$$\begin{aligned} E &\Rightarrow_L E + E \Rightarrow_L \text{id} + E \Rightarrow_L \text{id} + E * E \\ &\Rightarrow_L \text{id} + \text{id} * E \Rightarrow_L \text{id} + \text{id} * \text{id} \end{aligned}$$

$$\begin{aligned} E &\Rightarrow_R E + E \Rightarrow_R E + E * E \Rightarrow_R E + E * \text{id} \\ &\Rightarrow_R E + \text{id} * \text{id} \Rightarrow_R \text{id} + \text{id} * \text{id} \end{aligned}$$

Compilerphasen (grob)



Scanner (Lexer)

Quellcode
(Char-Stream)
zu Folge von
Token (z.B. **id**,
Schlüsselwörter,
Kommentar, ...)

Parser

- Syntaxprüfung
- Tokenfolge zu
Parsebaum
(Ableitungsbaum
durch Finden
einer Links- bzw.
Rechtsableitung)

- Prüfung von
Constraints
- Symboltabelle/
Scoping
- Typprüfung

- Parsebaum zu
Zielcode
- ggf. Code-
Optimierung

Scanner und Parser

- werden oft automatisch von Compilergeneratoren erzeugt
- **Scanner:**
 - **Token** durch reguläre Ausdrücke spezifiziert
 - aus regulären Ausdrücken werden die DEAs generiert
 - etwas Software „drumherum“
- **Parser:**
 - **Syntax** durch kfG (meist durch erweiterte BNF gegeben) spezifiziert
 - aus der kfG wird der Parser generiert
 - falls Syntaxfehler: Fehlermeldungen erzeugen
 - falls Programm syntaktisch korrekt: Ableitungsbaum ausgeben

Ein einfacher Parser – Prinzip

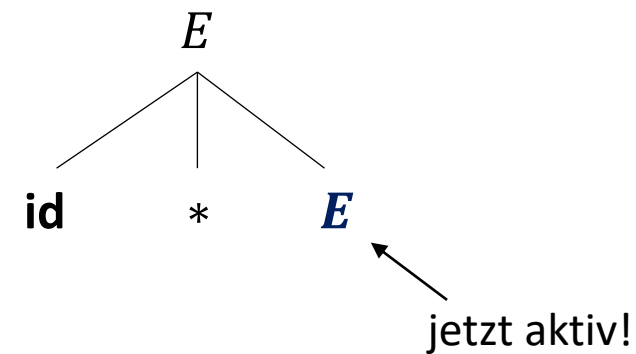
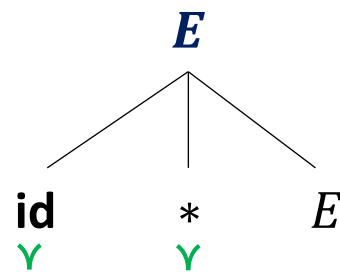
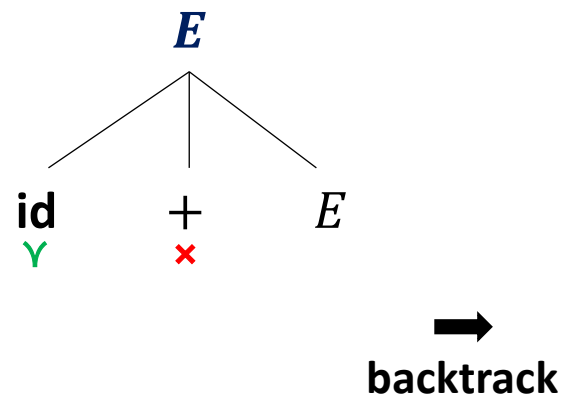
- **Eingabe:** kfG $G = (N, T, P, S)$, Wort $w \in T^*$
- **Ausgabe:** Ableitungsbaum für w , falls existent (*sonst:* Fehlermeldung)
- Für jedes Nichtterminal A : Festlegen einer Anordnung aller Regeln für A
- Konstruiere den Ableitungsbaum top-down
 - Füge S als Wurzel hinzu, markiere S als **aktiv**
 - Wenn A aktiv, dann wende die nächste A -Regel an (führt zu Kindern des Knotens mit aktivem Nichtterminal A)
 - Neu eingeführte Terminale von links nach rechts mit den nächsten noch nicht verglichenen Buchstaben von w vergleichen („Abhaken“)
 - am weitesten links stehendes Nichtterminal als aktiv markieren
 - im Fehlerfall (nicht passende Terminale) zurück zum Elternknoten mit A , nächste Regel

Ein einfacher Parser – Beispiel

$$E \rightarrow \text{id} + E \mid \text{id} * E \mid \text{id}$$

$$w = \text{id} * \text{id}$$

Y Y



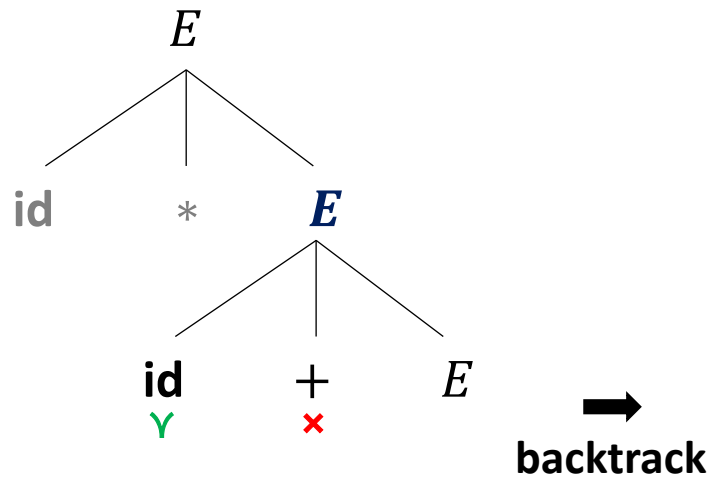
$$E \Rightarrow \text{id} + E$$

$$E \Rightarrow \text{id} * E$$

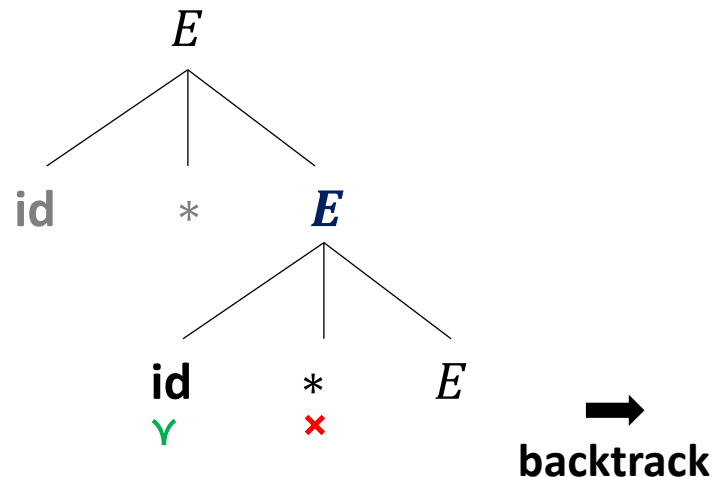
Ein einfacher Parser – Beispiel

$E \rightarrow \text{id} + E \mid \text{id} * E \mid \text{id}$

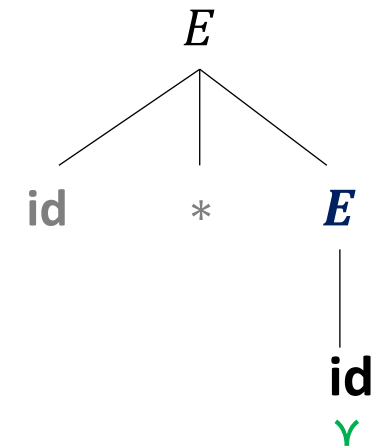
$w = \text{id} * \text{id}$
✓



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



$E \Rightarrow \text{id} * E \Rightarrow \text{id} * \text{id} * E$



$E \Rightarrow \text{id} * E \Rightarrow \text{id} * \text{id}$

Ein einfacher Parser – Diskussion

- Aufbau eines Ableitungsbaums von der Wurzel zum Eingabewort („Top-Down-Parsing“)
- Arbeitet stets am linken Ende der bisher gefundenen Satzform
 - Terminal: Abhaken
 - Nichtterminal: Aktivieren (also als nächstes Ersetzen)

➤ Aufbau einer **Linksableitung**

$$E \Rightarrow_L \text{id} * E \Rightarrow_L \text{id} * \text{id}$$

- **Backtracking** bei Fehlern
- Implementierung mit einem **Stack (linkes Symbol am Top)**

Ein einfacher Parser – Implementierung

$$E \rightarrow \text{id} + E \mid \text{id} * E \mid \text{id}$$

Stack	Eingabe	Aktion
E	$\text{id} * \text{id}$	$E \rightarrow \text{id} + E$
$\text{id} + E$	$\text{id} * \text{id}$	Terminal
$+ E$	$* \text{id}$	Backtrack
E	$\text{id} * \text{id}$	$E \rightarrow \text{id} * E$
$\text{id} * E$	$\text{id} * \text{id}$	Terminal
$* E$	$* \text{id}$	Terminal
E	id	$E \rightarrow \text{id} + E$
$\text{id} + E$	id	Terminal
$+ E$	ε	Backtrack

Ein einfacher Parser – Implementierung

$$E \rightarrow \text{id} + E \mid \text{id} * E \mid \text{id}$$

Stack	Eingabe	Aktion
E	id	$E \rightarrow \text{id} * E$
id * E	id	Terminal
* E	ε	Backtrack
E	id	$E \rightarrow \text{id}$
id	id	Terminal
ε	ε	ACCEPT