

# Reguläre Ausdrücke

## Bemerkung

- 1 *Reguläre Ausdrücke können geklammert werden, um die Reihenfolge der Operationen festzulegen. Um mit wenig Klammern auszukommen, vereinbart man Prioritäten, und zwar haben  $*$  und  $+$  die höchste Priorität, dann kommt das Produkt und schließlich die Vereinigung. Alle Operationen sind linksassoziativ.*
- 2 *Häufig werden zur Abkürzung regulären Ausdrücken Namen zugeordnet, die in anderen Ausdrücken benutzt werden können.*

# Reguläre Ausdrücke

## Satz

*Die durch reguläre Ausdrücke bezeichneten Mengen sind reguläre Sprachen, und jede reguläre Sprache lässt sich durch einen regulären Ausdruck bezeichnen.*

# Reguläre Ausdrücke

## Beispiel

$(a \mid b)^* abb$	Wörter über $\{a, b\}$ , die auf $abb$ enden
$(0 \mid 1)^* 0 (0 \mid 1)(0 \mid 1)$	Binär-Zahlen mit Achter-Resten 0,1,2,3
$0 (0 \mid 1)^* 0$	Binär-Zahlen, die mit 0 beginnen und enden
$(A \mid B \mid \dots \mid Z)(A \mid B \mid \dots \mid Z \mid 0 \mid \dots \mid 9)^*$	Bezeichner

Der reguläre Ausdruck  $(+ \mid - \mid \varepsilon) \langle \text{digit} \rangle^+ (\varepsilon \mid . \langle \text{digit} \rangle^+)$  beschreibt Festkommazahlen (ggf. mit Vorzeichen).

Dabei ist  $\langle \text{digit} \rangle = 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$  ein Name für einen regulären Ausdruck.

# Abkürzende Schreibweisen in regulären Ausdrücken

## Definition

- ① Null- oder einmaliges Vorkommen:  
Ein Ausdruck  $r \mid \varepsilon$  wird zu  $r?$  abgekürzt.
- ② Terminalzeichen  $c_i$ :  
Für den Ausdruck  $c_1 \mid c_2 \mid \dots \mid c_k$  schreibt man  $[c_1 c_2 \dots c_k]$ ,  
bei im Alphabet *aufeinanderfolgenden* Zeichen:  $[c_1 - c_k]$ .

# Reguläre Ausdrücke

## Beispiel

Der reguläre Ausdruck für Festkommazahlen wird jetzt kürzer formuliert:

Statt

$$\begin{aligned} & (+ \mid - \mid \varepsilon) \langle \text{digit} \rangle^+ (\varepsilon \mid . \langle \text{digit} \rangle^+) \\ & \langle \text{digit} \rangle = 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \end{aligned}$$

nun

$$\begin{aligned} & (+ \mid -)^? \langle \text{digit} \rangle^+ (. \langle \text{digit} \rangle^+)^? \\ & \langle \text{digit} \rangle = [0 - 9] \end{aligned}$$

oder

$$(+ \mid -)^? [0 - 9]^+ (. [0 - 9]^+)^?$$

# Von Zeichen zu Token

Für gängige Programmiersprachen nutzt man Scanner, um Zeichen zu Token zusammenzufassen, also zur Erkennung von

- reservierten Wörtern (**begin**, **for**, **if**, ...),
- Identifikatoren (Namen für z.B. Variable, Konstante, Typen),
- Konstanten (ganze Zahlen, Gleitpunktzahlen, Zeichenketten...),
- (zusammengesetzten) Sonderzeichen ( $:=$ ,  $<=$ , ...) und
- Zwischenraum (*white space*)  
(als Folge von Leerzeichen, Tabulatoren, Zeilenwechsel)

# Tokenklassen für eine einfache Programmiersprache

Beispiel (Tokenklassen für eine einfache Programmiersprache mit arithmetischen Ausdrücken, Wertzuweisungen, und if-Anweisungen)

Mit den Abkürzungen

`<letter>` = [A - Z] | [a - z]

`<digit>` = [0 - 9]

können folgende Tokenklassen definiert werden:

**ident** (Identifizier, Namen)

`<letter>` (`<letter>`|`<digit>`)\*

**number** (vorzeichenlose Zahlen)

`<digit>`<sup>+</sup> ( `.``<digit>`<sup>+</sup> )?

**relop** (Vergleichsoperatoren)

= | <= | <> | < | > | >=

**arithop** (arithmetische Op.)

+ | - | \* | /

Die Tokenklassen **if**, **then** und **else** enthalten jeweils das gleichnamige reservierte Wort.

# Scanner-Generatoren

Verwendet man Scanner-Generatoren wie **LEX**, **FLEX** oder **JFLEX**, so definiert man reguläre Ausdrücke für die Tokenklassen und geeignete Aktionen, die beim Erkennen eines Tokens ausgeführt werden sollen.

Wichtig ist, dass *alle* Eingabezeichen zu Token verarbeitet werden.

Daher sollte man eine weitere Tokenklasse **ws** (*white space*) schaffen, die das Überlesen von Zwischenraum steuert.

```
<delimiter> = <blank> | <tab> | <newline>  
und die Tokenklasse ws wird definiert durch  
<delimiter>+
```

Die zugeordnete Aktion ist in vielen Sprachen die Null-Operation!

Hat eine **Einrückung** aber eine Bedeutung in der Programmiersprache (z.B. **Haskell**, **Python**) (*off-side rule*), so erzeugt der Scanner bei verstärkter Einrückung ein **INDENT**-Token und bei verringerter Einrückung ein **DEDENT**-Token.



# Scanner-Generatoren

Scanner-Generatoren erzeugen Scanner automatisch;  
typisch ist folgendes Vorgehen:

- Erzeugung eines (nicht-deterministischen) Zustandsübergangsgraphen eines endlichen Automaten für jeden regulären Ausdruck einer Tokenklasse.
- Entfernung der  $\varepsilon$ -Übergänge im nicht-deterministischen Zustandsübergangsgraphen.
- Umwandlung des nicht-deterministischen Zustandsübergangsgraphen ohne  $\varepsilon$ -Übergänge in einen deterministischen Zustandsübergangsgraphen.
- Ableitung eines Scanners aus dem deterministischen Zustandsübergangsgraphen.

Oft sind aber **direkt erzeugte Scanner** effizienter.

Oder der Parser erledigt diese Aufgabe gleich mit.

# Das Parsing-Problem

Gegeben ist  
eine kontextfreie Grammatik  $G = (N, T, P, S)$  und  
ein Wort  $w \in T^*$ .

Frage: Ist  $w \in L(G)$ ? „**Wortproblem** für kontextfreie Sprachen“

- Das Problem ist entscheidbar,  
d.h. es gibt einen Algorithmus, der diese Frage beantwortet.
- Der Cocke-Younger-Kasami-Algorithmus läuft in einer Zeit  $O(|w|^3)$   
(für kontextfreie Grammatiken in Chomsky-Normalform).
- Wir suchen ein lineares Verfahren!

# Deterministische linksableitende Syntaxanalyse

Das Verfahren zur **linksableitenden Syntaxanalyse** (*top-down parsing*) soll

- 1 das gegebene Wort  $w$  nur einmal von links nach rechts lesen.
- 2 eine **Linksableitung** von  $w$  erzeugen.  
Das bedeutet, dass der Ableitungsbaum von oben nach unten und von links nach rechts konstruiert wird.

Außerdem soll das Verfahren **deterministisch** arbeiten:

- 3 Muss in der Linksableitung das nichtterminale Symbol  $A$  als nächstes verarbeitet werden, so ist die anzuwendende  $A$ -Produktion durch das nächste zu lesende Symbol von  $w$  (**Vorschau-Symbol**, *lookahead*) **eindeutig** bestimmt.

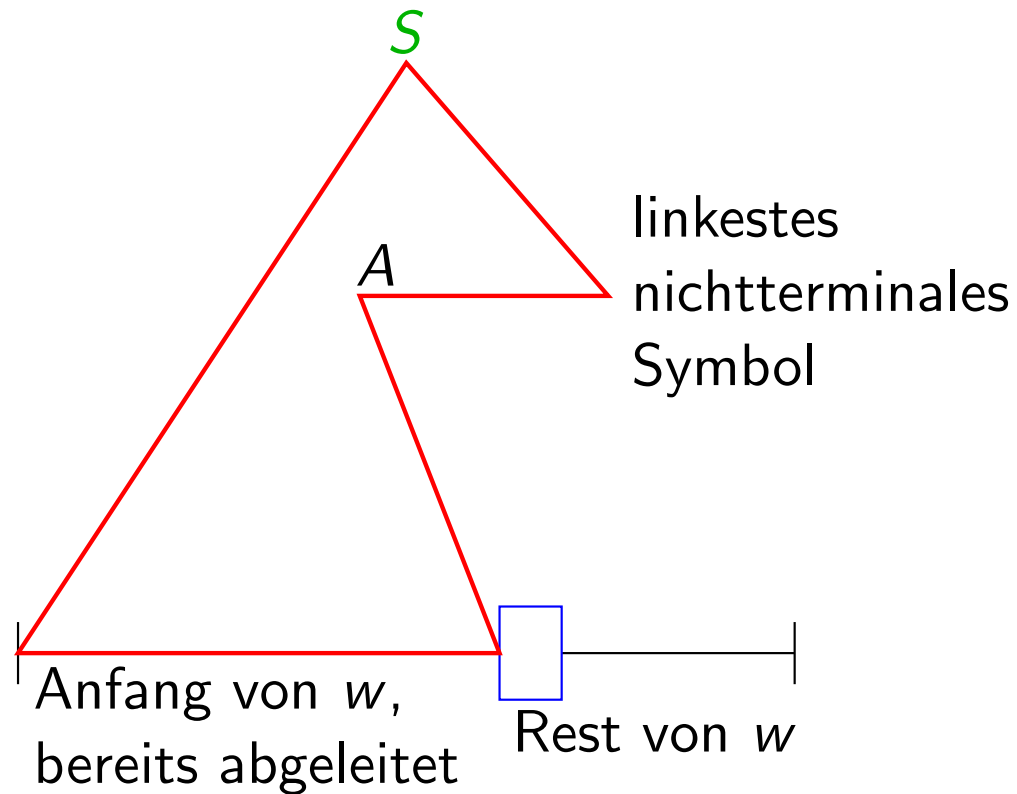
## Anfangssituation bei der linksableitenden Syntaxanalyse

$S$  Startsymbol der Grammatik



Lesefenster für Vorschau-Symbol,  
zeigt zu Anfang auf das erste Symbol der Eingabe

## Situation bei der linksableitenden Syntaxanalyse



## Bemerkung

- 1 *Nicht alle Grammatiken erlauben eine Syntaxanalyse dieser Art.*

## Beispiel

$$E \rightarrow E + E \mid E * E \mid a$$

$w = a + a$  Startsymbol  $E$  Vorschau  $a$  Welche Produktion anwenden?

- 2 *Es gibt sogar kontextfreie Sprachen, für die es keine Grammatik gibt, die eine deterministische linksableitende Syntaxanalyse ermöglicht:*

$$\{a^n 0 b^n \mid n \geq 0\} \cup \{a^n 1 b^{2n} \mid n \geq 0\}$$

Kfr. Produktionen:  $\{S \rightarrow A \mid B, \quad A \rightarrow aAb \mid 0, \quad B \rightarrow aBbb \mid 1\}$

Aber eine deterministische linksableitende Syntaxanalyse ist stets möglich, falls für jedes nichtterminale Zeichen  $A$  gilt:

Alle rechten Seiten der  $A$ -Produktionen beginnen mit *verschiedenen terminalen* Zeichen (**s-Grammatik**).

(Das leere Wort wird nur direkt aus dem Startsymbol abgeleitet, das dann nicht auf der rechten Seite einer Produktion vorkommen darf.)

# Hindernisse für deterministisches Top-Down Parsing

## Eigenschaften von Grammatiken, die eine deterministische linksableitende Syntaxanalyse offensichtlich verhindern

- 1 Es gibt verschiedene Produktionen mit der gleichen linken Seite, deren rechte Seiten ein gleiches Anfangsstück (*Präfix*) haben.
- 2 Es gibt *linksrekursive Produktionen*, d.h. Produktionen der Form  $A \rightarrow A\alpha$ ,  $\alpha \in (N \cup T)^+$ .

# Entfernen gemeinsamer Präfixe

$$A \rightarrow \alpha\beta_1 \mid \dots \mid \alpha\beta_r \mid \gamma_1 \mid \dots \mid \gamma_s$$

seien alle  $A$ -Produktionen der Grammatik,  
wobei  $\alpha \neq \varepsilon$  und kein  $\gamma_i$  das Präfix  $\alpha$  hat.

Ersetze diese Produktionen durch

$$\begin{aligned} A &\rightarrow \alpha A' \mid \gamma_1 \mid \dots \mid \gamma_s \quad \text{und} \\ A' &\rightarrow \beta_1 \mid \dots \mid \beta_r, \end{aligned}$$

wobei  $A'$  ein neues nichtterminales Symbol ist.

Man sieht leicht: Die umgeformte Grammatik erzeugt die gleiche Sprache.



# Entfernen gemeinsamer Präfixe

## Beispiel

$$S \rightarrow a + S \mid a * S \mid (S) * S \mid a$$

$$S \rightarrow aS' \mid (S) * S \quad \text{und}$$

$$S' \rightarrow +S \mid *S \mid \varepsilon$$

# Entfernen linksrekursiver Produktionen

$$A \rightarrow A\alpha_1 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \dots \mid \beta_n$$

seien alle  $A$ -Produktionen,  
wobei kein  $\beta_i$  mit  $A$  beginnt und alle  $\alpha_i \neq \varepsilon$  sind.

Ersetze diese Produktionen durch:

$$\begin{aligned} A &\rightarrow \beta_1 A' \mid \dots \mid \beta_n A' \quad \text{und} \\ A' &\rightarrow \alpha_1 A' \mid \dots \mid \alpha_m A' \mid \varepsilon, \end{aligned}$$

wobei  $A'$  ein neues nichtterminales Symbol ist.

## Bemerkung

*Diese Umformungen müssen nicht unbedingt zu einer Grammatik führen, die deterministisches Parsing erlaubt.*

# Entfernen linksrekursiver Produktionen

Beispiel (durch Komma getrennte Liste  $L$  von Werten  $v$ )

$$L \rightarrow v \mid L, v$$

$$L \rightarrow vR \quad \text{und}$$

$$R \rightarrow ,vR \mid \varepsilon$$

Beispiel (gültige Klammerungen)

$$S \rightarrow SS \mid (S) \mid [S] \mid \varepsilon$$

$$S \rightarrow (S)S' \mid [S]S' \mid \varepsilon S' \quad \text{und}$$

$$S' \rightarrow SS' \mid \varepsilon$$

# Definition der Funktionen **First** und **Follow**

Wann erlaubt eine kontextfreie Grammatik eine deterministische linksableitende Syntaxanalyse?

## Definition

Sei  $\alpha \in (N \cup T)^*$ . Dann ist

$$\mathbf{First}(\alpha) = \{a \in T \mid \alpha \xRightarrow{*} a\beta, \beta \in (N \cup T)^*\} \cup \{\varepsilon \mid \alpha \xRightarrow{*} \varepsilon\}$$

## Bemerkung

**First**( $\alpha$ ) enthält diejenigen terminalen Zeichen, die bei Linksableitungen von  $\alpha$  als erste auftreten.

Kann man aus  $\alpha$  das leere Wort  $\varepsilon$  herleiten, so ist  $\varepsilon$  ebenfalls in **First**( $\alpha$ ).

## Beispiel

$$S \rightarrow aS' \mid (S) * S \quad S' \rightarrow +S \mid *S \mid \varepsilon$$

$$\mathbf{First}(S) = \{a, ( \quad \mathbf{First}(S') = \{+, *, \varepsilon\}$$

# Definition der Funktionen **First** und **Follow**

## Definition

Sei  $A \in N$  und  $\$$  ein neues Symbol, das das Ende der Eingabe markiert.

$$\begin{aligned} \mathbf{Follow}(A) = & \{a \in T \mid S \xRightarrow{*} \alpha A a \beta, \alpha, \beta \in (N \cup T)^*\} \\ & \cup \{\$ \mid S \xRightarrow{*} \alpha A, \alpha \in (N \cup T)^*\} \end{aligned}$$

## Bemerkung

**Follow**( $A$ ) enthält die terminalen Zeichen, die bei einer Ableitung vom Startsymbol  $S$  direkt auf  $A$  folgen können.

Ist  $A$  das letzte Zeichen, so ist die Endmarke  $\$$  in **Follow**( $A$ ).

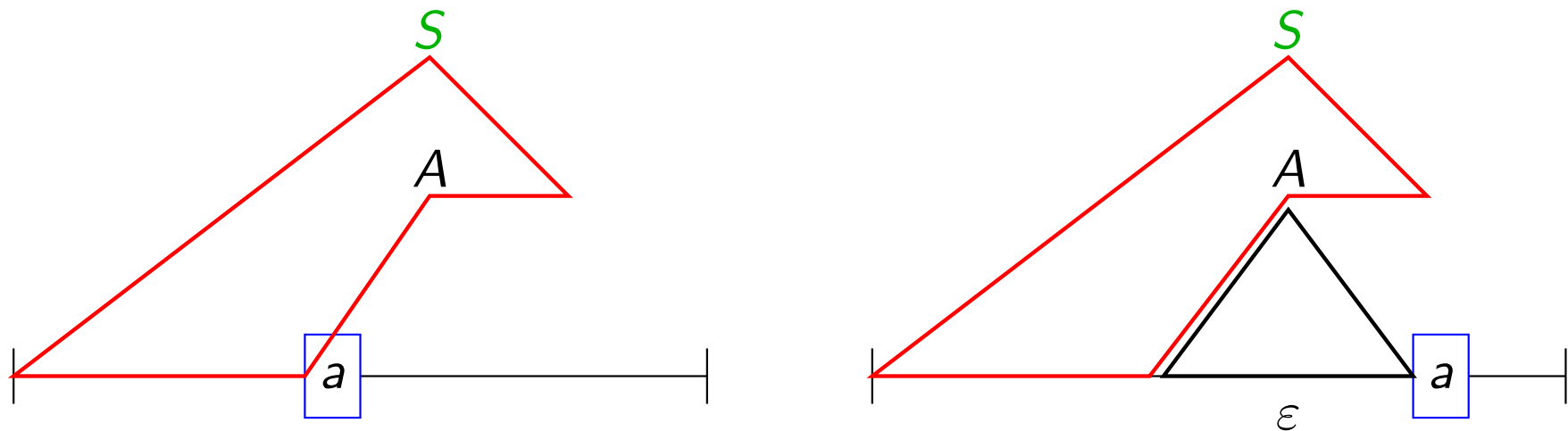
## Beispiel

$$S \rightarrow aS' \mid (S) * S \quad S' \rightarrow +S \mid *S \mid \varepsilon$$

$$\mathbf{Follow}(S) = \{\$, )\} \quad \mathbf{Follow}(S') = \{\$, )\}$$

# Linksableitung eines Wortes

Betrachten wir eine Linksableitung eines Wortes der Sprache und sei  $A$  das linkeste nichtterminale Zeichen, es ist also als nächstes eine  $A$ -Produktion anzuwenden:



# erwartete Vorschau-Symbole einer Produktion

Wir bestimmen zunächst mit Hilfe der Grammatik die Menge der **erwarteten Symbole**, die in einer Linksableitung bei Anwendung einer Produktion  $A \rightarrow \alpha$  als Vorschau-Symbol vorkommen können:

$$\mathbf{Erwartet}(A \rightarrow \alpha) = \begin{cases} \mathbf{First}(\alpha) & , \text{ falls } \varepsilon \notin \mathbf{First}(\alpha) \\ (\mathbf{First}(\alpha) \setminus \{\varepsilon\}) \cup \mathbf{Follow}(A) & , \text{ falls } \varepsilon \in \mathbf{First}(\alpha) \end{cases}$$

Das Vorschau-Symbol des zu testenden Wortes soll die anzuwendende  $A$ -Produktion *eindeutig* bestimmen, dazu müssen die Vorschau-Mengen aller  $A$ -Produktionen paarweise disjunkt sein.

# LL(1) Grammatik

## Definition

Eine kontextfreie Grammatik  $G = (N, T, P, S)$  heißt **LL(1)-Grammatik**  
(*Lesen der Eingabe von **L**inks nach rechts,  
Erzeugen einer **L**inksableitung und  
nur **1** Zeichen als Vorschau*),

falls für alle  $A \in N$  gilt:

Seien  $A \rightarrow \alpha_1 \mid \dots \mid \alpha_n$  alle  $A$ -Produktionen in  $P$ .

- ① **First**( $\alpha_1$ ), ..., **First**( $\alpha_n$ ) sind paarweise disjunkt,  
d.h. **First**( $\alpha_i$ )  $\cap$  **First**( $\alpha_j$ ) =  $\emptyset$  falls  $i \neq j$ .
- ② Ist  $\varepsilon \in$  **First**( $\alpha_j$ ), dann ist **Follow**( $A$ )  $\cap$  **First**( $\alpha_i$ ) =  $\emptyset$   
für alle  $1 \leq i \leq n, i \neq j$ .



# LL(1) Grammatik

## Beispiel

$S \rightarrow aAaa \mid bAba \quad A \rightarrow b \mid \varepsilon$

$S$ : **First**( $aAaa$ ) =  $\{a\}$ , **First**( $bAba$ ) =  $\{b\}$

Bedingung 1 erfüllt, Bedingung 2 trivial!

$A$ : **First**( $b$ ) =  $\{b\}$ , **First**( $\varepsilon$ ) =  $\{\varepsilon\}$

Bedingung 1 erfüllt, Bedingung 2 nicht trivial!

**Follow**( $A$ ) =  $\{a, b\}$

Bedingung 2: **Follow**( $A$ )  $\cap$  **First**( $b$ ) =  $\{a, b\} \cap \{b\} = \{b\} \neq \emptyset$  nicht erfüllt!

keine LL(1)-Grammatik, nicht zum det. top-down-Parsen geeignet

# Berechnung der **First**- und **Follow**-Funktion

Gegeben sei eine kontextfreie Grammatik  $G = (N, T, P, S)$ .

Wir bestimmen zuerst, aus welchen nichtterminalen Symbolen das leere Wort abgeleitet werden kann, d.h. die Menge  $N_\varepsilon = \{A \in N \mid A \xRightarrow{*} \varepsilon\}$ .

## Algorithmus zur Berechnung von $N_\varepsilon$ :

- 1 Bestimme  $N_0 = \{A \mid A \rightarrow \varepsilon \text{ ist in } P\}$
- 2 Berechne  $N_{i+1} = N_i \cup \{A \mid A \rightarrow \alpha \text{ ist in } P \text{ und } \alpha \in N_i^*\}$
- 3 Ist  $N_{i+1} = N_i$ , dann setze  $N_\varepsilon = N_i$ .

# Berechnung der **First**- und **Follow**-Funktion

Nun berechnen wir die First-Funktion für alle nichtterminalen Symbole der Grammatik.

## Algorithmus zur Berechnung der First-Funktion:

Berechnung von **First**( $A$ ) für alle  $A \in N$ .

Man konstruiert den **Parchmann-Graphen**  $\Gamma_{first}(G)$  folgendermaßen:

- ① Jedes Symbol aus  $N \cup T$  wird durch einen Knoten dargestellt.
- ② Für jede Produktion  $A \rightarrow X_1 \dots X_n$  mit  $n \geq 1$ , fügt man eine Kante von  $A$  nach  $X_i$  hinzu, falls alle davorstehenden Symbole  $X_1, \dots, X_{i-1} \in N_\epsilon$ .
- ③ Setze **First**( $A$ ) =
 
$$\begin{cases} \{a \in T \mid \text{ex. Weg in } \Gamma_{first}(G) \text{ von } A \text{ nach } a\} & , \text{ falls } A \notin N_\epsilon \\ \{\epsilon\} \cup \{a \in T \mid \text{ex. Weg in } \Gamma_{first}(G) \text{ von } A \text{ nach } a\} & , \text{ falls } A \in N_\epsilon \end{cases}$$

# Berechnung der **First**- und **Follow**-Funktion

Damit ist die **First**-Funktion auch für alle Wörter  $\alpha \in (N \cup T)^*$  bestimmt.

Wir unterscheiden, ob  $\alpha$  das leere Wort ist bzw. mit einem terminalen oder mit einem nichtterminalen Zeichen beginnt:

- 1 **First** $(\varepsilon) = \{\varepsilon\}$
- 2 **First** $(a\beta) = \{a\}$ , falls  $a \in T$ .
- 3 **First** $(A\beta) = \begin{cases} \mathbf{First}(A) & , \text{ falls } A \in N, A \notin N_\varepsilon \\ (\mathbf{First}(A) - \{\varepsilon\}) \cup \mathbf{First}(\beta) & , \text{ falls } A \in N_\varepsilon \end{cases}$

# Berechnung der First-Funktion

## Beispiel (1)

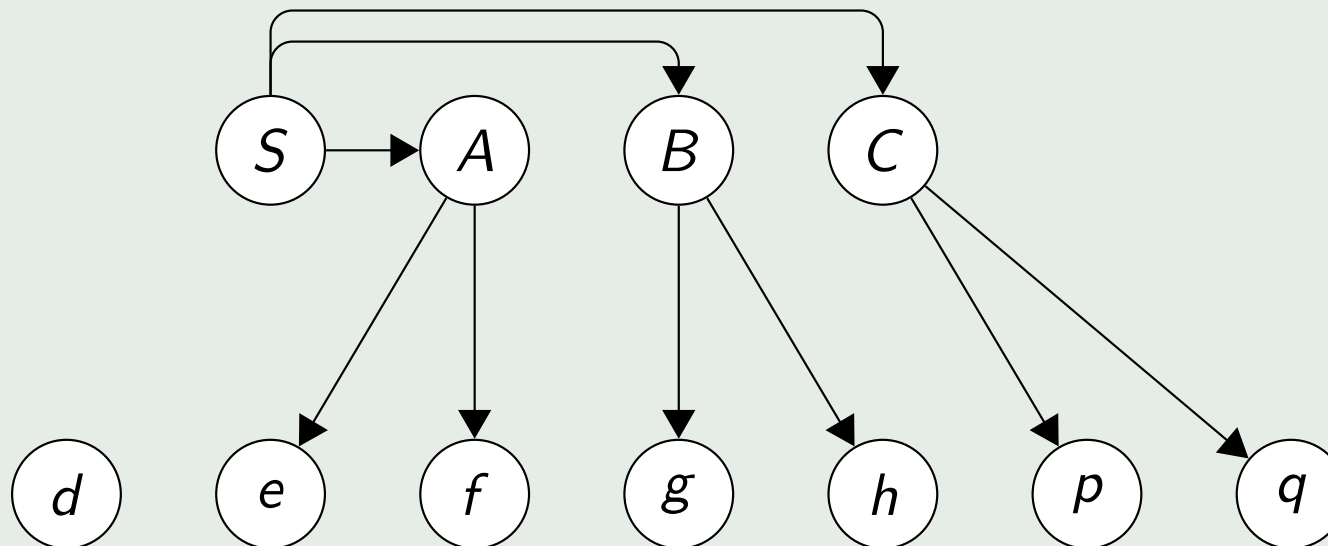
Betrachte die Grammatik  $G = (N, T, P, S)$

mit  $N = \{S, A, B, C\}$ ,  $T = \{d, e, f, g, h, p, q\}$  und den

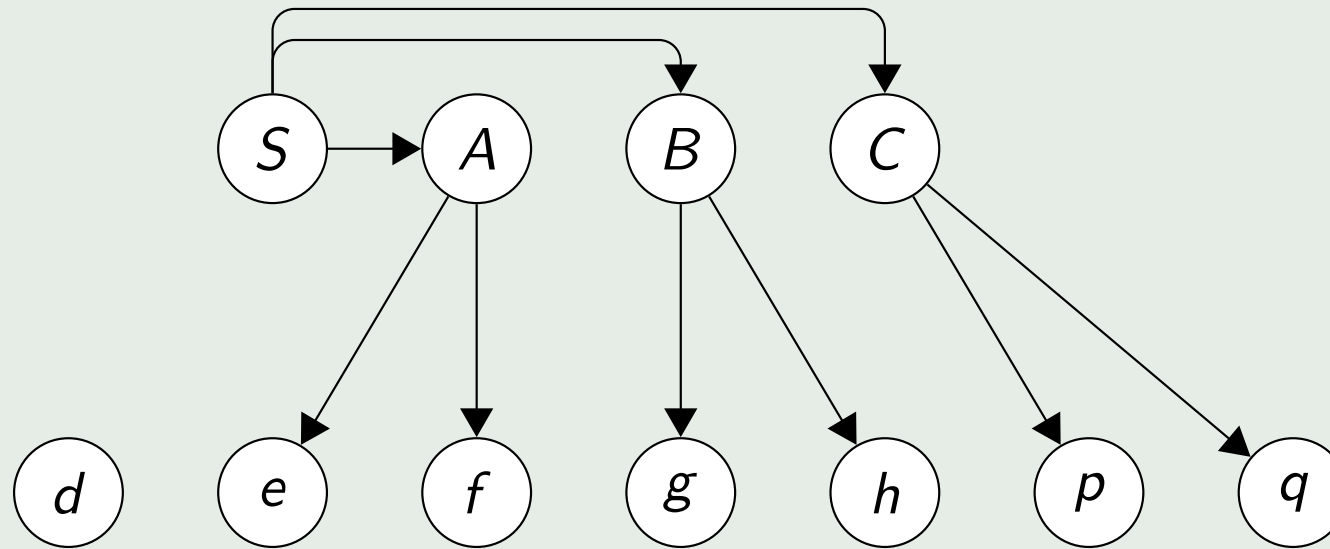
Produktionen  $S \rightarrow ABCd$ ,  $A \rightarrow e \mid f \mid \varepsilon$ ,  $B \rightarrow g \mid h \mid \varepsilon$  und  $C \rightarrow p \mid q$ .

Es ist  $N_\varepsilon = \{A, B\}$  und

der Parchmann-Graph  $\Gamma_{first}(G)$  ist:



## Beispiel (1)



Also gilt:

$$\mathbf{First}(S) = \{e, f, g, h, p, q\},$$

$$\mathbf{First}(A) = \{\varepsilon, e, f\},$$

$$\mathbf{First}(B) = \{\varepsilon, g, h\} \text{ und}$$

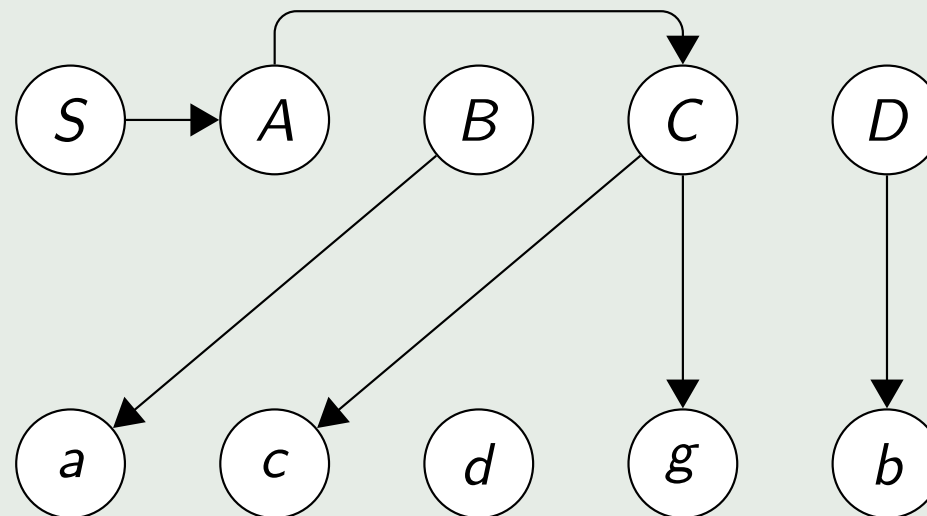
$$\mathbf{First}(C) = \{p, q\}.$$

## Beispiel (2)

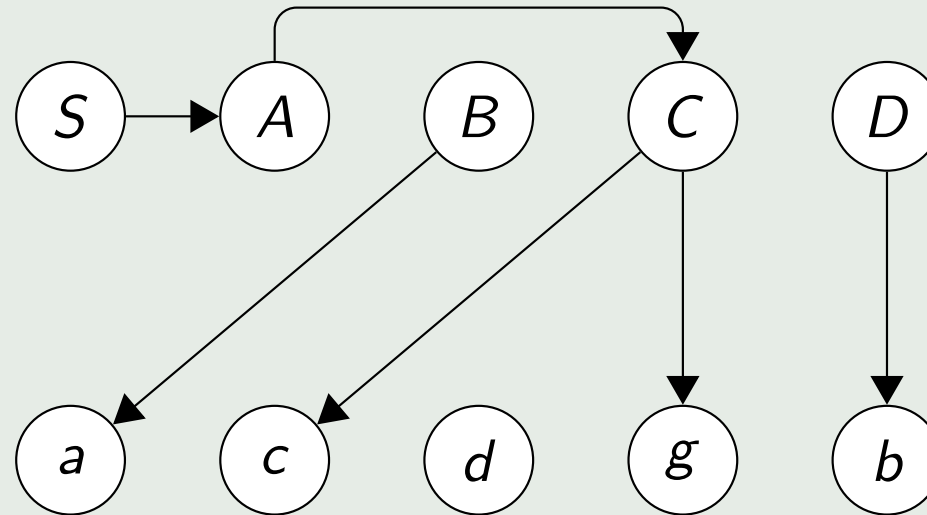
$G = (N, T, P, S)$  mit  $N = \{S, A, B, C, D\}$  und  $T = \{a, b, c, d, g\}$ ,  
 $P = \{S \rightarrow AB, B \rightarrow aAB \mid \varepsilon, A \rightarrow CD, D \rightarrow bCD \mid \varepsilon, C \rightarrow cSd \mid g\}$

Offenbar gilt:  $N_\varepsilon = \{B, D\}$

Der Parchmann-Graph  $\Gamma_{first}(G)$  ist:



## Beispiel (2)



Es gilt also ( $N_\varepsilon = \{B, D\}$ ):

**First**( $S$ ) = **First**( $A$ ) = **First**( $C$ ) =  $\{c, g\}$ ,

**First**( $B$ ) =  $\{\varepsilon, a\}$ , und

**First**( $D$ ) =  $\{\varepsilon, b\}$ .

und z.B. **First**( $BAD$ ) =  $\{a, c, g\}$ .



## Algorithmus zur Berechnung der Follow-Funktion

Der **Parchmann-Graph**  $\Gamma_{follow}(G)$  hat einen Knoten für jedes Symbol in  $N \cup T \cup \{\$\}$ .

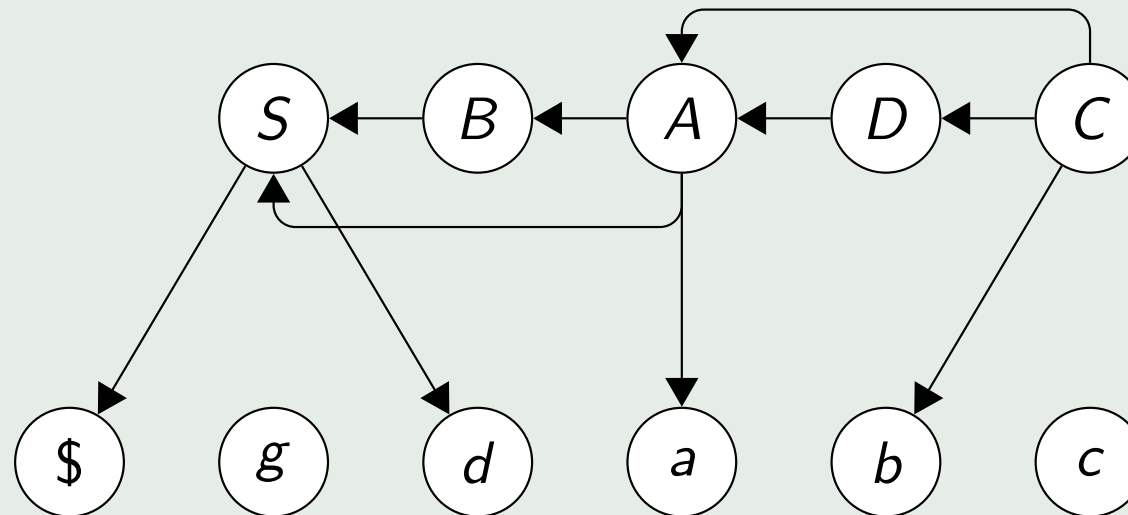
- ① Füge eine **Kante vom Startsymbol  $S$  nach  $\$$**  hinzu.
- ② Für jedes Vorkommen eines Nonterminals  $B$  auf der rechten Seite einer Produktion  $A \rightarrow \alpha B \beta$  füge **Kanten von  $B$  nach jedem terminalen Symbol aus  $\mathbf{First}(\beta)$**  hinzu.  
Ist  $\varepsilon \in \mathbf{First}(\beta)$  und  $A \neq B$ , so füge eine **Kante von  $B$  nach  $A$**  hinzu.
- ③ **Follow** $(A) := \{a \in T \cup \{\$\} \mid \text{ex. Weg von } A \text{ nach } a \text{ in } \Gamma_{follow}(G)\}$ .

# Berechnung der Follow-Funktion

## Beispiel (2)

$P = \{S \rightarrow AB, B \rightarrow aAB \mid \varepsilon, A \rightarrow CD, D \rightarrow bCD \mid \varepsilon, C \rightarrow cSd \mid g\}$

Der Parchmann-Graph  $\Gamma_{follow}(G)$  ist:



## Beispiel (2)

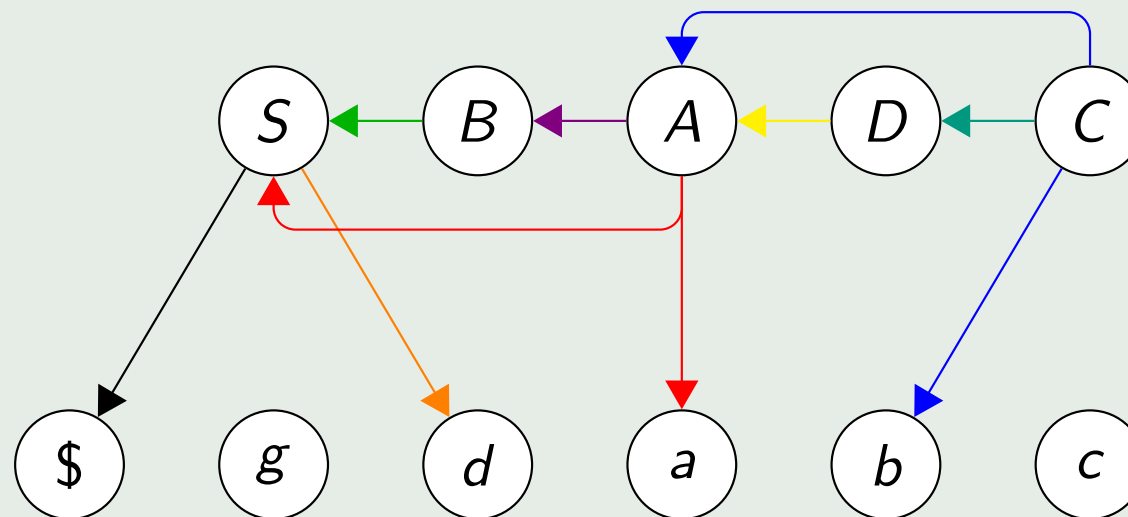
$$P = \{S \rightarrow AB, B \rightarrow aAB \mid \varepsilon, A \rightarrow CD, D \rightarrow bCD \mid \varepsilon, C \rightarrow cSd \mid g\}$$

$$\mathbf{First}(S) = \mathbf{First}(A) = \mathbf{First}(C) = \{c, g\},$$

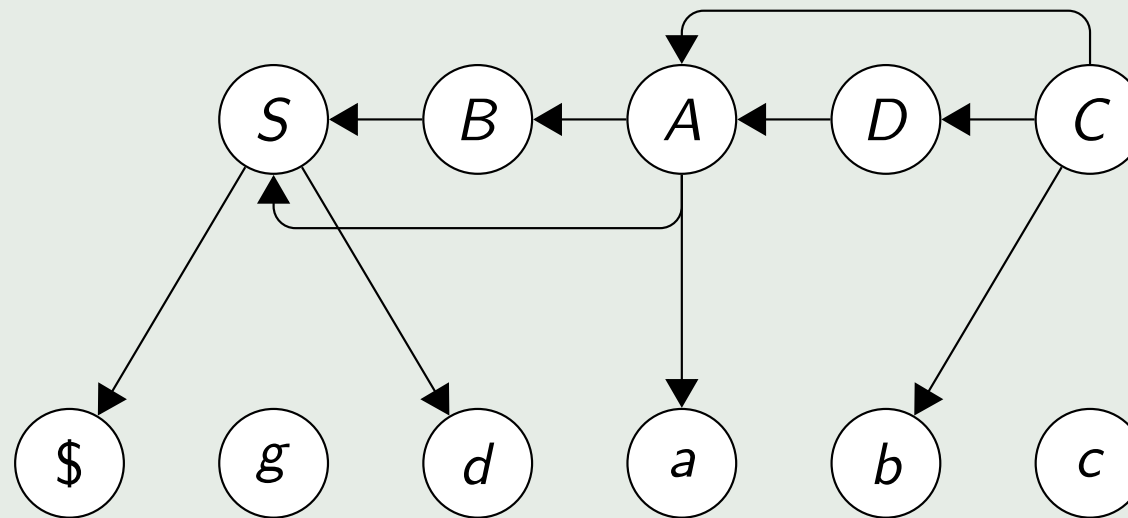
$$\mathbf{First}(B) = \{\varepsilon, a\} \text{ und}$$

$$\mathbf{First}(D) = \{\varepsilon, b\}.$$

Der Parchmann-Graph  $\Gamma_{follow}(G)$  ist:



## Beispiel (2)



$\text{Follow}(S) = \text{Follow}(B) = \{d, \$\},$   
 $\text{Follow}(A) = \text{Follow}(D) = \{a, d, \$\}$  und  
 $\text{Follow}(C) = \{a, b, d, \$\}.$

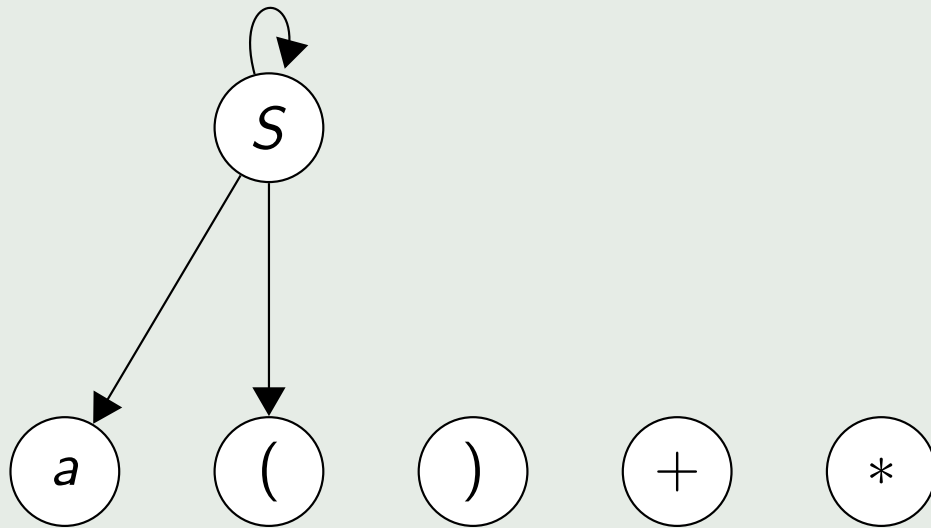
# Überprüfung einer Grammatik

## Beispiel (Grammatik für arithmetische Ausdrücke)

Produktionen:  $S \rightarrow S + S \mid S * S \mid (S) \mid a$

Berechnung der **First**-Mengen:

$$N_\epsilon = \emptyset$$



**First**( $S + S$ ) =  $\{ (, a \}$ ,

**First**( $S * S$ ) =  $\{ (, a \}$ ,

**First**( $(S)$ ) =  $\{ ( \}$  und

**First**( $a$ ) =  $\{ a \}$ .

**First**-Mengen nicht disjunkt!

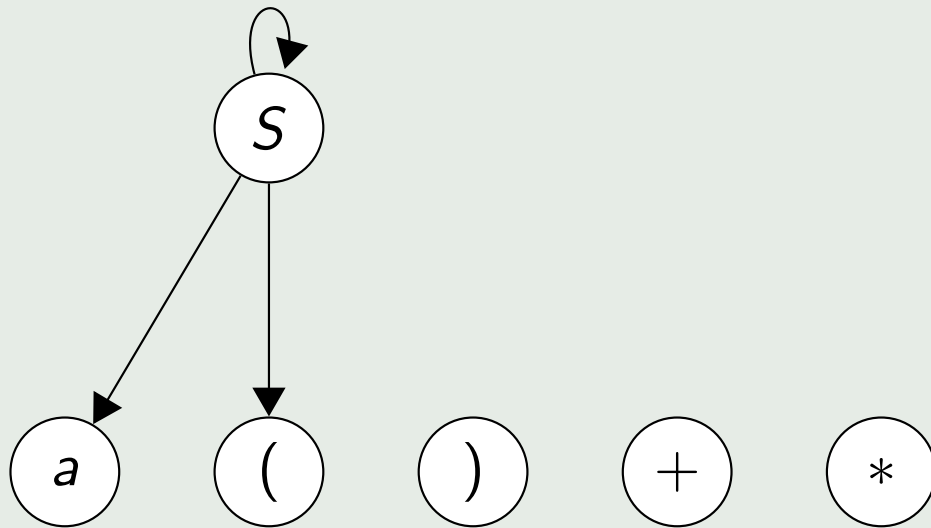
Keine LL(1)-Grammatik!

## Beispiel (Grammatik für arithmetische Ausdrücke)

Produktionen:  $S \rightarrow S + S \mid S * S \mid (S) \mid a$

Berechnung der **First**-Mengen:

$$N_\epsilon = \emptyset$$



**First**( $S + S$ ) =  $\{ (, a \}$ ,

**First**( $S * S$ ) =  $\{ (, a \}$ ,

**First**( $(S)$ ) =  $\{ ( \}$  und

**First**( $a$ ) =  $\{ a \}$ .

**First**-Mengen nicht disjunkt!

Keine LL(1)-Grammatik!

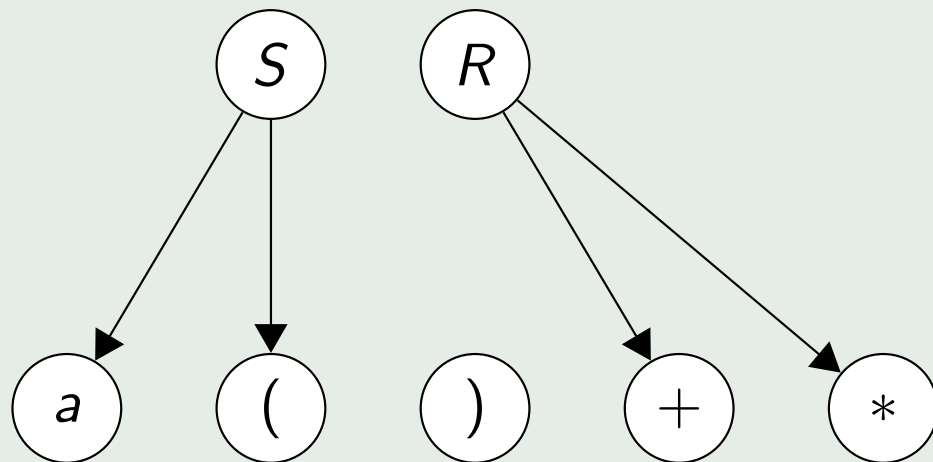
Nicht überraschend, da Produktionen linksrekursiv!

## Beispiel (Grammatik nach Entfernung der Linksrekursion)

Produktionen:  $S \rightarrow (S)R \mid aR$   
 $R \rightarrow +SR \mid *SR \mid \varepsilon$

Berechnung der **First**-Mengen:

$$N_\varepsilon = \{R\}$$



Noch zu prüfen:  
und

$$\mathbf{First}(S) = \{a, (,$$

$$\mathbf{First}(R) = \{\varepsilon, +, *\},$$

$$\mathbf{First}((S)R) = \{ (,$$

$$\mathbf{First}(aR) = \{a\}. \quad \text{disjunkt!}$$

$$\mathbf{First}(+SR) = \{+\},$$

$$\mathbf{First}(*SR) = \{*\} \text{ und}$$

$$\mathbf{First}(\varepsilon) = \{\varepsilon\}. \quad \text{disjunkt!}$$

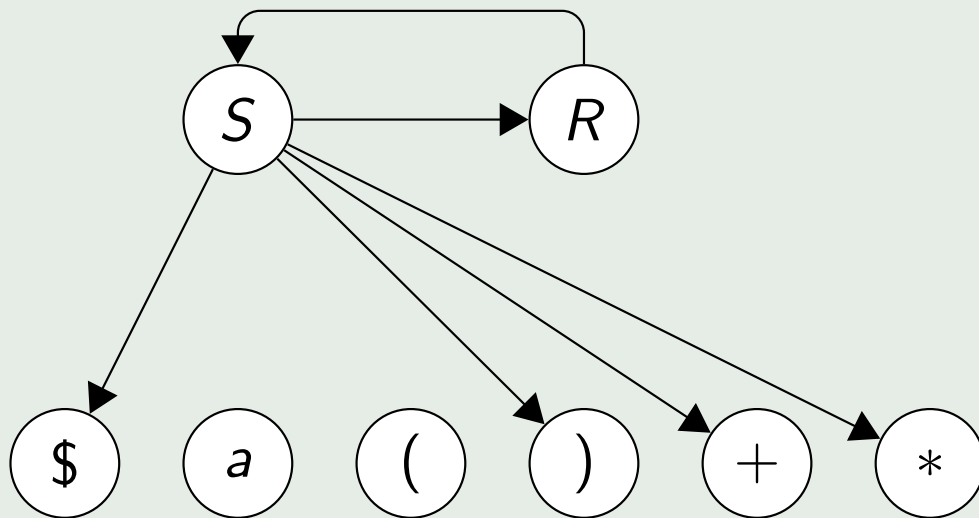
$$\mathbf{Follow}(R) \cap \mathbf{First}(+SR) = \emptyset$$

$$\mathbf{Follow}(R) \cap \mathbf{First}(*SR) = \emptyset$$

## Beispiel (Grammatik nach Entfernung der Linksrekursion)

Produktionen:  $S \rightarrow (S)R \mid aR$   
 $R \rightarrow +SR \mid *SR \mid \varepsilon$

Berechnung der **Follow**-Mengen:



**Follow**( $S$ ) =  $\{\$, ), +, *\}$ ,  
**Follow**( $R$ ) =  $\{\$, ), +, *\}$

**Follow**( $R$ )  $\cap$  **First**( $+SR$ ) =  
 $\{\$, ), +, *\} \cap \{+\} \neq \emptyset$

**Follow**( $R$ )  $\cap$  **First**( $*SR$ ) =  
 $\{\$, ), +, *\} \cap \{*\} \neq \emptyset$

Beispiel für  $+ \in \mathbf{Follow}(R)$ :  
 $S \Rightarrow aR \Rightarrow a+SR \Rightarrow a+aRR \Rightarrow$   
 $a+a\textcolor{red}{R}+SR$

**Keine LL(1)-Grammatik!**

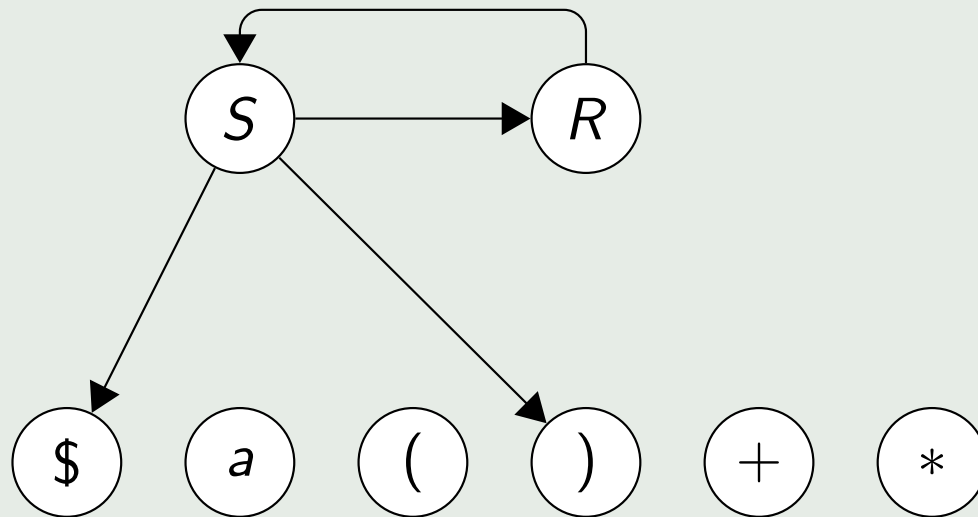


## Beispiel (modifizierte Grammatik)

Um  $+$  und  $*$  aus **Follow**( $R$ ) zu entfernen, ändern wir die Grammatik (aber nicht die Sprache!):

Produktionen:  $S \rightarrow (S)R \mid aR$   
 $R \rightarrow +S \mid *S \mid \varepsilon$

Berechnung der **Follow**-Mengen:



**Follow**( $S$ ) =  $\{\$, )\}$ ,  
**Follow**( $R$ ) =  $\{\$, )\}$

**Follow**( $R$ )  $\cap$  **First**( $+S$ ) =  
 $\{\$, )\} \cap \{+\} = \emptyset$

**Follow**( $R$ )  $\cap$  **First**( $*S$ ) =  
 $\{\$, )\} \cap \{*\} = \emptyset$

**LL(1)-Grammatik!**