

# Einschub: Reguläre Ausdrücke

(E1)

Reguläre Sprachen ("Typ 3 - Sprachen") werden durch reguläre Grammatiken beschrieben ("generiert"). Das ist i. A. ziemlich unständlich! Bsp.: Ganze Zahlen

$$S \rightarrow \emptyset$$

$$S \rightarrow 1$$

⋮

$$S \rightarrow 9$$

$$S \rightarrow \emptyset S$$

$$S \rightarrow 1 S$$

⋮

$$S \rightarrow 9 S$$

Deswegen benutzt man "reguläre Ausdrücke" zum Beschreiben von regulären Sprachen. Diese können vollständig durch reguläre Grammatiken ausgedrückt werden: keine Veränderung der prinzipiellen Ausdruckskraft! Aber es ist viel bequemer:

$$[\emptyset - 9]^+$$

Bestandteile von regulären Ausdrücken:

einzelne Zeichen

$a$

Alternative

$\beta | \gamma$

Konkatenation

$\beta \gamma$

Epsilon

$\epsilon$  (leerer String)

Kleene'scher Abschluss

$\beta^*$  ( $= \epsilon | \beta | \beta\beta | \dots$ )

Meist ebenfalls zugelassen:

Charakterklasse

$[abcd]$  ( $= a|b|c|d$ )

$[b-e]$  ( $= b|c|d|e$ )

Optionaler Ausdruck

$\beta?$  ( $= \beta | \epsilon$ )

Positiv Abschluss

$\beta^+$  ( $= \beta\beta^*$ )

Irreguläre Zeichen (außer Zeichenumlauf)

# Aufgaben

(E2)

1. Binärzahlen, die Vielfache von 2 sind:

$$(\emptyset | 1) * \emptyset$$

2. Strings aus a und b, die mind. einmal zwei a's hintereinander enthalten:

$$(a | b) * a a (a | b) *$$

3. Strings aus a und b, die keine zwei a's hintereinander enthalten:

$$b * (a | b * ) * (a | \epsilon)$$

4. Hexadezimalzahlen

$$\emptyset * [\emptyset - 9 a - f A - F] +$$

5. Ein einzelnes Plus-Zeichen:

$$\backslash +$$

6. Identifizier in der Prog.-Sprache C:

$$[a - z A - Z _] [a - z A - Z \emptyset - 9 _] *$$

Hier sieht man die Nützlichkeit von "regulären Abkürzungen"  
(= Namen für reguläre Ausdrücke):

$$\text{Definiere } \{L\} = [a - z A - Z _] \quad (\text{"Letter"})$$

$$\{D\} = [\emptyset - 9] \quad (\text{"Digit"})$$

Dann kann man den reg. Ausdruck für "Identifizier" so schreiben:

$$\{L\} (\{L\} | \{D\}) *$$



## Einschub: Kontextfreie Grammatiken

(E3)

Kontextfreie Sprachen ("Typ 2-Sprachen") werden durch kontextfreie Grammatiken beschrieben ("generiert").

Bestandteile:

- Menge von Terminalsymbolen  $T$  (die "Tokens")
- Menge von Nichtterminalsymbolen  $N$  (die "Variablen")
- Menge von Produktionen der Form  $l \rightarrow r$ , wobei  $l \in N$  und  $r \in (T \cup N)^*$  ("Stück v. Grammatiksyml.")
- Ein ausgezeichnetes  $S \in N$ , das "Startsymbol".

Bsp.: Korrekte Klammerung

$$T = \{a, l, r\}, N = \{S\}$$

$$P: S \rightarrow a$$

$$S \rightarrow l S r$$

Abkürzung f. die Produktionen:  $S \rightarrow a \mid l S r$

Bem.: Jede reguläre Sprache ist eine kontextfreie Sprache.

$$a \quad S \rightarrow a$$

$$\varepsilon \quad S \rightarrow \varepsilon$$

$$\beta \mid \gamma \quad S \rightarrow \beta \mid \gamma$$

$$\beta \gamma \quad S \rightarrow \beta \gamma$$

$$\beta^* \quad S \rightarrow \varepsilon \mid \beta S \quad (\text{"rechtsrekursiv"})$$

$$\text{alternativ: } S \rightarrow \varepsilon \mid S \beta \quad (\text{"linksrekursiv"})$$

Also: Reguläre Sprachen  $\subsetneq$  Kontextfreie Sprachen

Warum gibt es dann überhaupt Scanner?

Gründe: Effizienz, Einfachheit!

# Aufgaben

(E4)

1. Eine Liste mit bel. vielen a's, einschll. keinem a:

linksrekursiv:  $S \rightarrow \varepsilon$

$S \rightarrow Sa$

rechtsrekursiv:  $S \rightarrow \varepsilon$

$S \rightarrow aS$

2. Eine Liste mit bel. vielen a's, aber mind. einem a:

linksrekursiv:  $S \rightarrow a$

$S \rightarrow Sa$

rechtsrekursiv:  $S \rightarrow a$

$S \rightarrow aS$

3. Eine Liste mit bel. vielen a's (einschließlich keinem),  
getrennt durch jeweils ein b ("Parameterliste"):

linksrekursiv:  $S \rightarrow \varepsilon \mid T$

$T \rightarrow a \mid Tba$

rechtsrekursiv:  $S \rightarrow \varepsilon \mid T$

$T \rightarrow a \mid a b T$

4. Allgemeiner Tipp: Trennung von Wiederholung und  
der Spezifikation, was wiederholt wird!

Bsp.: Eine bel. Folge von "Ablenkungen" d  
und "Statements" e:

$S \rightarrow \varepsilon \mid ST$

$T \rightarrow d \mid e$



## Einschub: Ableitungen, Parse-Bäume, Mehrdeutigkeit (ES)

Grammatik f. arithmet. Ausdrücke:

$$E \rightarrow E + E \mid E - E \mid E * E \mid E / E \mid (E) \mid -E \mid id$$

Ist  $-(id + id)$  ein Satz der Grammatik?

Ja, denn es gibt eine "Ableitung":  $E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E + E) \Rightarrow -(id + E) \Rightarrow -(id + id)$

In jedem Herleitungsschritt zwei Entscheidungen:

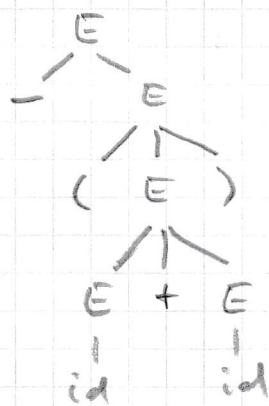
- Welches Nichtterminal wird ersetzt?
- Welche Alternative für dieses Nichtterminal wird gewählt?

Links ableitungen: das am weitesten links stehende N.T.

Rechts ableitungen: " " " rechts " N.T.

Parse-Baum: grafische Darstellung einer Ableitung ohne Berücksichtigung der Reihenfolge.

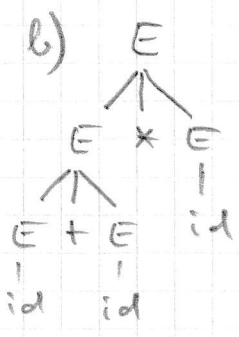
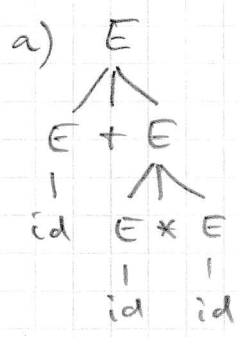
Bsp:  $-(id + id)$



Syntaxanalyse = Finden einer Ableitung = Konstr. des Parse-Baums bei gegeb. Satz

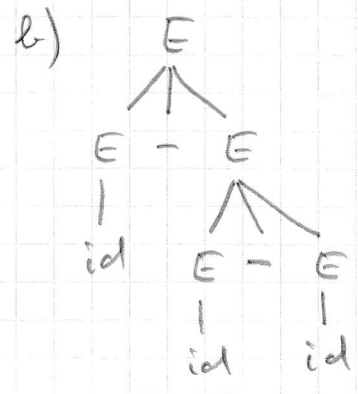
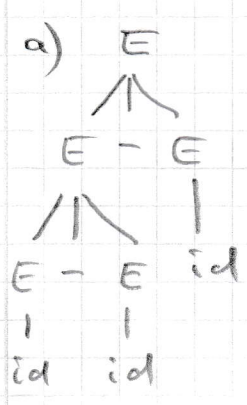
Eine Grammatik, bei der es mind. einen Satz gibt, der durch versch. Parse-Bäume repräsentiert wird, heißt "mehrdeutig".

Bsp: Grammatik von oben, Satz  $id + id * id$



Wir wollen a) haben, wg. "Priorität" der Operatoren

Bsp.: Grammatik von oben, Satz id-id-id



Wir wollen a) haben, wg. "Assoziativität" des Operators

→ Grammatik ist mehrdeutig! Wie wird sie eindeutig, bei gegeb. Prioritäten und Assoziativitäten?

- Neue Nichtterminale E ("Expression")
- T ("Term")
- F ("Faktor")

$$E \rightarrow E + T \mid E - T \mid T$$

$$T \rightarrow T * F \mid T / F \mid F$$

$$F \rightarrow -F \mid id \mid (E)$$

Die Priorität der Operatoren wird durch die "Staffelung" E/T/F erreicht, die (Links-) Assoziativität durch die (Links-) Rekursion der Produktionen.