

Reguläre und kontextfreie Sprachen

Grundzüge digitaler Systeme (192.134)

Vortrag von: Gernot Salzer

Formale Sprachen – Übersicht

- 1 Operationen auf formalen Sprachen
- 2 Reguläre Sprachen
- 3 Reguläre Ausdrücke
- 4 Eigenschaften regulärer Sprachen
- 5 Vom regulären Ausdruck zum Automaten
- 6 Vom Automaten zum regulären Ausdruck
- 7 Kontextfreie Sprachen
- 8 Beispiele kontextfreier Grammatiken
- 9 Andere Notationen
- 10 Style Guide für Grammatiken
- 11 Grammatiken in freier Wildbahn
- 12 Jenseits der Kontextfreiheit

Operationen auf formalen Sprachen

- Σ Alphabet, d.h., endliche, nicht-leere Menge atomarer Symbole
 w Wort über Σ , (endliche) Folge von Zeichen aus dem Alphabet Σ
 ε Leerwort
 Σ^* Menge aller endlichen Wörter über Σ (inklusive Leerwort)
 $w \cdot w' = ww'$ Verkettung der Wörter $w, w' \in \Sigma^*$

Seien $L, L' \subseteq \Sigma^*$ zwei Sprachen.

$$L \cup L' = \{ w \mid w \in L \text{ oder } w \in L' \} \quad \text{Vereinigung}$$

$$L \cdot L' = \{ w \cdot w' \mid w \in L, w' \in L' \} \quad \text{Verkettung}$$

$$\begin{aligned} L^0 &= \{\varepsilon\} \\ L^{n+1} &= L \cdot L^n \quad (n \geq 0) \end{aligned} \quad \text{Potenzen}$$

$$\begin{aligned} L^+ &= \bigcup_{n \geq 1} L^n \\ L^* &= \bigcup_{n \geq 0} L^n = L^0 \cup L^+ = \{\varepsilon\} \cup L^+ \quad \text{Kleene-Stern} \end{aligned}$$

Verkettung

$$\{a, b\} \cdot \{b, c, d\} = \{ab, ac, ad, bb, bc, bd\}$$

$$\{b, c, d\} \cdot \{a, b\} = \{ba, bb, ca, cb, da, db\}$$

$$\begin{aligned} (\{a, b\} \cdot \{1, 2\}) \cdot \{\#, \$\} &= \{a1\#, a1\$, a2\#, a2\$, b1\#, b1\$, b2\#, b2\$\} \\ &= \{a, b\} \cdot (\{1, 2\} \cdot \{\#, \$\}) \end{aligned}$$

$$\{a, b\} \cdot \{\varepsilon\} = \{a \cdot \varepsilon, b \cdot \varepsilon\} = \{a, b\} = \{\varepsilon \cdot a, \varepsilon \cdot b\} = \{\varepsilon\} \cdot \{a, b\}$$

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

$$\{\varepsilon\} \cdot \{\varepsilon\} = \{\varepsilon\}$$

$$\{\} \cdot \{\} = \{\varepsilon\} \cdot \{\} = \{\} \cdot \{\varepsilon\} = \{\}$$

Beobachtungen:

- Sprachverkettung ist nicht kommutativ.
- Sprachverkettung ist assoziativ.
- $\{\varepsilon\}$ ist neutrales Element bzgl. Sprachverkettung.
- $\{\}$ ist Nullelement bzgl. Sprachverkettung.

Potenzen von $\{a, 42\}$

$$L = \{a, 42\}$$

$$L^0 = \{\varepsilon\}$$

$$L^1 = L \cdot L^0 = L \cdot \{\varepsilon\} = L = \{a, 42\}$$

$$L^2 = L \cdot L^1 = L \cdot L = \{aa, a42, 42a, 4242\}$$

$$L^3 = L \cdot L^2 = \{aaa, aa42, a42a, a4242, 42aa, 42a42, 4242a, 424242\}$$

\vdots

$$L^+ = \bigcup_{n \geq 1} L^n = L^1 \cup L^2 \cup L^3 \cup \dots$$

$$= \{a, 42, aa, a42, 42a, 4242, aaa, aa42, a42a, a4242, 42aa, \dots\}$$

$$L^* = \bigcup_{n \geq 0} L^n = L^0 \cup L^1 \cup L^2 \cup L^3 \cup \dots$$

$$= \{\varepsilon, a, 42, aa, a42, 42a, 4242, aaa, aa42, a42a, a4242, 42aa, \dots\}$$

Potenzen eines Alphabets Σ

$$\Sigma^0 = \{\varepsilon\}$$

$$\Sigma^1 = \Sigma$$

$$\Sigma^n$$

alle Σ -Wörter der Länge n (d.h., mit n Symbolen)

$$\Sigma^+ = \bigcup_{n \geq 1} \Sigma^n$$

alle Σ -Wörter ohne Leerwort

$$\Sigma^* = \bigcup_{n \geq 0} \Sigma^n$$

alle Σ -Wörter mit Leerwort

Rechengesetze für formale Sprachen

$A, B, C \dots$ formale Sprachen

Vereinigung

$$(A \cup B) \cup C = A \cup (B \cup C)$$

$$\{\} \cup A = A$$

$$A \cup \{\} = A$$

$$A \cup B = B \cup A$$

$$A \cup A = A$$

Verkettung

$$(A \cdot B) \cdot C = A \cdot (B \cdot C)$$

$$\{\varepsilon\} \cdot A = A$$

$$A \cdot \{\varepsilon\} = A$$

$$\{\} \cdot A = \{\}$$

$$A \cdot \{\} = \{\}$$

Distributivität

$$A \cdot (B \cup C) = (A \cdot B) \cup (A \cdot C)$$

$$(B \cup C) \cdot A = (B \cdot A) \cup (C \cdot A)$$

Kleene-Stern und -Plus

$$(A^*)^* = A^*$$

$$A^* \cdot A = A^+$$

$$A \cdot A^* = A^+$$

$$(A \cup \{\varepsilon\})^* = A^*$$

$$A^+ \cup \{\varepsilon\} = A^*$$

Formale Sprachen – Übersicht

- 1 Operationen auf formalen Sprachen
- 2 Reguläre Sprachen**
- 3 Reguläre Ausdrücke
- 4 Eigenschaften regulärer Sprachen
- 5 Vom regulären Ausdruck zum Automaten
- 6 Vom Automaten zum regulären Ausdruck
- 7 Kontextfreie Sprachen
- 8 Beispiele kontextfreier Grammatiken
- 9 Andere Notationen
- 10 Style Guide für Grammatiken
- 11 Grammatiken in freier Wildbahn
- 12 Jenseits der Kontextfreiheit

Reguläre Sprachen

Alle Sprachen, die aus einem Alphabet mit Hilfe von Vereinigung, Verkettung und Stern gebildet werden können.

Anwendungen:

- Betriebssystem-Shells: DOS („Wildcards“), UNIX-Shells (wie `sh`, `csch`, `ash`, `bash`, `zsh`), ...
- UNIX Kommandozeilenprogramme: `grep`, `awk`, `ed`, `sed`, ...
- Editoren: `vi`, `emacs`, ...
- Compilerbau: *Tokens* bilden reguläre Sprache, die durch sog. Scanner (Lexer) wie `lex` oder `flex` verarbeitet werden.
- Programmiersprachen: PERL, TCL, PHP, PYTHON, RUBY, R, JAVA, JAVASCRIPT, .NET-Sprachen, ...
- Websprachen: XML Schema, XQuery, XPath, DTDs, ...
- Datenbanken: MySQL, Oracle, PostgreSQL, ...
- ...

Regulären Sprachen über einem Alphabet

Die Menge der regulären Sprachen über Σ , $\mathcal{L}_{\text{reg}}(\Sigma)$, ist die kleinste Menge, sodass gilt:

- $\{\}$, $\{\epsilon\}$ und $\{s\}$ sind reguläre Sprachen (für alle $s \in \Sigma$).
- Wenn L und L' reguläre Sprachen sind, dann auch $L \cup L'$, $L \cdot L'$ und L^* .

Reellen Numerale: reguläre Sprache über $\Sigma = \{0, \dots, 9, ., E, +, -\}$

$real = digit \cdot digit^* \cdot \{.\} \cdot digit^* \cdot (\{\epsilon\} \cup scale)$

$scale = \{E\} \cdot \{+, -, \epsilon\} \cdot digit \cdot digit^*$

$digit = \{0, \dots, 9\} = \{0\} \cup \dots \cup \{9\}$

Wichtig: Unterscheide Symbole des Alphabets von Meta-Symbolen!

$0, \dots, 9, ., E, +, - \dots$ Symbole des Alphabets

$\epsilon, real, scale, digit \dots$ Meta-Symbole, Abkürzungen

Formale Sprachen – Übersicht

- 1 Operationen auf formalen Sprachen
- 2 Reguläre Sprachen
- 3 Reguläre Ausdrücke**
- 4 Eigenschaften regulärer Sprachen
- 5 Vom regulären Ausdruck zum Automaten
- 6 Vom Automaten zum regulären Ausdruck
- 7 Kontextfreie Sprachen
- 8 Beispiele kontextfreier Grammatiken
- 9 Andere Notationen
- 10 Style Guide für Grammatiken
- 11 Grammatiken in freier Wildbahn
- 12 Jenseits der Kontextfreiheit

Reguläre Ausdrücke

Ausdrücke wie $digit \cdot digit^*$ und $digit^* \cdot digit$ sind ununterscheidbar: Beides sind semantische Beschreibungen der Menge aller Ziffernfolgen. Um Aussagen über die Form der Ausdrücke treffen zu können, benötigen wir eine formale Sprache.

Reguläre Ausdrücke (algebraische Notation)

Die regulären Ausdrücke über Σ sind die kleinste Menge, für die gilt:

- \emptyset , ε und s sind reguläre Ausdrücke (für alle Symbole $s \in \Sigma$).
- Sind X und Y reguläre Ausdrücke, dann auch $(X + Y)$, (XY) und X^* .

Vereinfachte Klammerung: $+$ bindet am schwächsten, $*$ am stärksten. Keine Klammern bei gleichartigen Operatoren (wegen Assoziativität).

Die Sprache $\mathcal{L}(X)$ zu einem regulären Ausdruck X ist definiert durch:

$$\mathcal{L}(\emptyset) = \{\}$$

$$\mathcal{L}(\varepsilon) = \{\varepsilon\}$$

$$\mathcal{L}(s) = \{s\} \quad \text{für } s \in \Sigma$$

$$\mathcal{L}(X + Y) = \mathcal{L}(X) \cup \mathcal{L}(Y)$$

$$\mathcal{L}(XY) = \mathcal{L}(X) \cdot \mathcal{L}(Y)$$

$$\mathcal{L}(X^*) = (\mathcal{L}(X))^*$$

Regulärer Ausdruck für die reellen Numerale

$$R = DD^* \cdot D^*(\varepsilon + S)$$

$$S = E(+ + - + \varepsilon)DD^*$$

$$D = 0 + 1 + \dots + 9$$

(R , S und D sind Abkürzungen für die jeweiligen regulären Ausdrücke.)

Die zugehörigen Sprachen:

$$\begin{aligned}\mathcal{L}(D) &= \mathcal{L}(0 + 1 + \dots + 9) \\ &= \mathcal{L}(0) \cup \mathcal{L}(1) \cup \dots \cup \mathcal{L}(9) \\ &= \{0\} \cup \{1\} \cup \dots \cup \{9\} \\ &= \textit{digit}\end{aligned}$$

$$\begin{aligned}\mathcal{L}(S) &= \mathcal{L}(E(+ + - + \varepsilon)DD^*) \\ &= \mathcal{L}(E) \cdot \mathcal{L}(+ + - + \varepsilon) \cdot \mathcal{L}(D) \cdot \mathcal{L}(D^*) \\ &= \{E\} \cdot (\mathcal{L}(+) \cup \mathcal{L}(-) \cup \mathcal{L}(\varepsilon)) \cdot \textit{digit} \cdot \mathcal{L}(D)^* \\ &= \{E\} \cdot (\{+\} \cup \{-\} \cup \{\varepsilon\}) \cdot \textit{digit} \cdot \textit{digit}^* \\ &= \textit{scale}\end{aligned}$$

$$\mathcal{L}(R) = \dots = \textit{real}$$

Zwei reguläre Ausdrücke X und Y heißen äquivalent, geschrieben $X = Y$, wenn $\mathcal{L}(X) = \mathcal{L}(Y)$ gilt.

$$((a + b)^* + \varepsilon)^* = (a + b)^*$$

$$\begin{aligned}\mathcal{L}(((a + b)^* + \varepsilon)^*) &= \dots \\ &= (\{a, b\}^* \cup \{\varepsilon\})^* \\ &= (\{a, b\}^*)^* \quad \text{da } \varepsilon \in L^* \text{ für alle } L \\ &= (\{a, b\}^*)^0 \cup (\{a, b\}^*)^1 \cup (\{a, b\}^*)^2 \cup \dots \\ &= \{a, b\}^* \cup (\{a, b\}^*)^0 \cup (\{a, b\}^*)^2 \cup \dots \\ &= \{a, b\}^* \quad \text{da } L^* \text{ alle Wörter über } L \text{ enthält} \\ &= \dots \\ &= \mathcal{L}((a + b)^*)\end{aligned}$$

Reguläre Ausdrücke in EBNF-Notation

EBNF ... Erweiterte Backus-Naur-Form (Formalismus zur Beschreibung der Syntax von Programmiersprachen, der reguläre Ausdrücke zulässt)



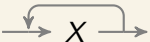
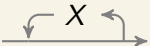

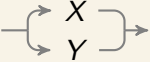

XY	$X \cdot Y$	Aufeinanderfolge
$X Y$	$X \cup Y$	Alternativen
$[X]$	$\{\epsilon\} \cup X$	Option
$\{X\}$	X^*	Wiederholung
(X)	(X)	Gruppierung
$"s"$	$\{s\}$	Symbol

Reelle Numerale in EBNF-Notation

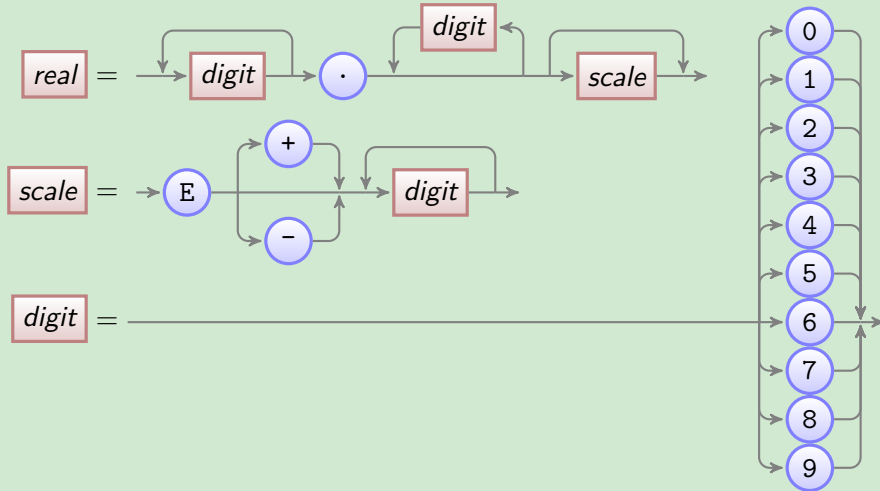
```
real = digit {digit} "." {digit} [scale]
scale = "E" ["+" | "-"] digit {digit}
digit = "0" | "1" | "2" | ... | "9"
```

Reguläre Ausdrücke als Syntaxdiagramme

Syntaxdiagramme ... graphische Form der EBNF

	A	Abkürzung
	$\{s\}$	Symbol
	X^+	Wiederholung ≥ 1
	X^*	Wiederholung ≥ 0
	$X \cdot Y$	Aufeinanderfolge
	$X \cup Y$	Alternativen
	$\{\epsilon\} \cup X$	Option

Reelle Numerale als Syntaxdiagramm



Reguläre Ausdrücke im informatischen Alltag

UNIX := „30 definitions of regular expressions living under one roof“ [Donald E. Knuth]

```
grep -E -e "regex" file
```

liefert alle Zeilen der Datei *file*, die eine Zeichenkette enthalten, die dem regulären Ausdruck *regex* (in POSIX ERE Syntax) entspricht.

Verschiedene „Standards“:

- POSIX Basic Regular Expressions
- POSIX Extended Regular Expressions
- PERL Regular Expressions
- ...



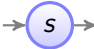

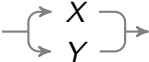
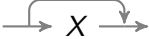
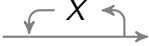
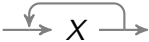
POSIX Extended Regular Expressions (ERE)

<i>regex</i>	trifft zu auf	<i>regex</i>	trifft zu auf
$\backslash s$	Zeichen s	XY	X gefolgt von Y
s	s , falls kein Sonderzeichen	$X Y$	X oder Y
$.$	alle Zeichen	X^*	≥ 0 Mal X
$^$	Zeilenanfang	X^+	≥ 1 Mal X
$\$$	Zeilenende	$X?$	≤ 1 Mal X
$[s_1 \cdots s_n]$	ein Zeichen aus $\{s_1, \dots, s_n\}$	$X\{i\}$	i Mal X
$[^s_1 \cdots s_n]$	alle Zeichen außer s_1, \dots, s_n	$X\{i, \}$	$\geq i$ Mal X
(X)	X	$X\{i, j\}$	i bis j Mal X

Reelle Numerale als ERE

$digit = \{0, \dots, 9\}$ [0-9]
 $scale = \{E\} \cdot \{+, -, \varepsilon\} \cdot digit \cdot digit^*$ $E[+-]?[0-9]^+$
 $real = digit \cdot digit^* \cdot \{.\} \cdot digit^* \cdot (\{\varepsilon\} \cup scale)$ $\wedge [0-9]^+ \cdot \backslash . [0-9]^* (E[+-]?[0-9]^+)? \$$
 [0-9] ... Kurzform von [0123456789]; analog [a-zA-Z] für Buchstaben

Zusammenfassung der Notationen

Reg. Sprache	Algebra	EBNF	Syntaxdiagramm	POSIX ERE	
A	A	A			Abkürzung
$\{\}$	\emptyset				Leersprache
$\{\varepsilon\}$	ε				Leerwortsprache
$\{s\}$	s	"s"		$\backslash s$	Terminalsymbol
$X \cdot Y$	XY	XY		XY	Aufeinanderfolge
$X \cup Y$	$X + Y$	$X Y$		$X Y$	Alternativen
$\{\varepsilon\} \cup X$	$\varepsilon + X$	$[X]$		$X?$	Option
X^*	X^*	$\{X\}$		X^*	Wiederholung ≥ 0
X^+	XX^*, X^+	$X\{X\}$		X^+	Wiederholung ≥ 1
(X)	(X)	(X)		(X)	Gruppierung

Weitere POSIX-Notationen:

POSIX ERE	Bedeutung	Reguläre Sprache
s	s , falls kein Sonderzeichen	$\{s\}$
$.$	alle Zeichen	Σ
\wedge	Zeilenanfang	
$\$$	Zeilenende	
$[s_1 \cdots s_n]$	eines der Zeichen s_i	$\{s_1, \dots, s_n\}$
$[\wedge s_1 \cdots s_n]$	alle Zeichen außer s_1, \dots, s_n	$\Sigma - \{s_1, \dots, s_n\}$
$X\{i\}$	i Mal X	X^i
$X\{i, \}$	$\geq i$ Mal X	$X^i \cdot X^*$
$X\{i, j\}$	i bis j Mal X	$X^i \cup X^{i+1} \cup \dots \cup X^j$

Formale Sprachen – Übersicht

- 1 Operationen auf formalen Sprachen
- 2 Reguläre Sprachen
- 3 Reguläre Ausdrücke
- 4 Eigenschaften regulärer Sprachen**
- 5 Vom regulären Ausdruck zum Automaten
- 6 Vom Automaten zum regulären Ausdruck
- 7 Kontextfreie Sprachen
- 8 Beispiele kontextfreier Grammatiken
- 9 Andere Notationen
- 10 Style Guide für Grammatiken
- 11 Grammatiken in freier Wildbahn
- 12 Jenseits der Kontextfreiheit

Eigenschaften regulärer Sprachen

Abschlusseigenschaften reguläre Sprachen

Reguläre Sprachen sind abgeschlossen gegenüber

- Vereinigung, Verkettung und Stern (warum?)
- Durchschnitt, Komplement, Differenz, Homomorphismen, Quotientenbildung, ...

Entscheidbarkeit eines Problems: Es gibt ein Verfahren (einen Algorithmus), das für jede Eingabe die richtige Antwort ja/nein liefert.

Nicht entscheidbar: Halteproblem Ihrer Lieblingsprogrammiersprache

- Gegeben ein Programm mit einer Eingabe, wird es anhalten?

Entscheidbare Probleme regulärer Sprachen

- Gegeben ein Wort w und einen regulären Ausdruck X , gilt $w \in \mathcal{L}(X)$? (Wortproblem)
- Gegeben zwei reguläre Ausdrücke X und Y , gilt $\mathcal{L}(X) = \mathcal{L}(Y)$? (Äquivalenzproblem)
- Gegeben einen regulären Ausdruck X , ist $\mathcal{L}(X)$ leer/endlich/unendlich?

Ausdrucks kraft regulärer Sprachen

Die regulären Sprachen sind genau jene, die von endlichen Automaten akzeptiert werden, d.h.:

- Zu jedem regulären Ausdruck X gibt es einen endlichen Automaten \mathcal{A} , sodass $\mathcal{L}(\mathcal{A}) = \mathcal{L}(X)$ gilt.
- Zu jedem endlichen Automaten \mathcal{A} gibt es einen regulären Ausdruck X , sodass $\mathcal{L}(X) = \mathcal{L}(\mathcal{A})$ gilt.

Nicht regulär sind Sprachen, deren Analyse ein unbegrenztes Gedächtnis erfordert:

- Klammerausdrücke: $\{(), (()), ()(), ((())), (())(), ()()(), \dots\}$
- $\{a^n b^n \mid n \geq 0\} = \{\varepsilon, ab, aabb, aaabbb, \dots\}$
- $\{a^n b^n c^n \mid n \geq 0\} = \{\varepsilon, abc, aabbcc, aaabbbccc, \dots\}$
- Palindrome: Wörter, die identisch mit ihrem Spiegelbild sind.
 $\{\text{otto}, \text{anna}, \text{reliefpfeiler}, \text{ogenie der herr ehre dein ego}, \dots\}$
- Doppelwörter: $\{ww \mid w \in \Sigma^*\}$ (falls $|\Sigma| > 1$)

Formale Sprachen – Übersicht

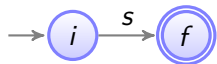
- 1 Operationen auf formalen Sprachen
- 2 Reguläre Sprachen
- 3 Reguläre Ausdrücke
- 4 Eigenschaften regulärer Sprachen
- 5 Vom regulären Ausdruck zum Automaten**
- 6 Vom Automaten zum regulären Ausdruck
- 7 Kontextfreie Sprachen
- 8 Beispiele kontextfreier Grammatiken
- 9 Andere Notationen
- 10 Style Guide für Grammatiken
- 11 Grammatiken in freier Wildbahn
- 12 Jenseits der Kontextfreiheit

Vom regulären Ausdruck zum Automaten

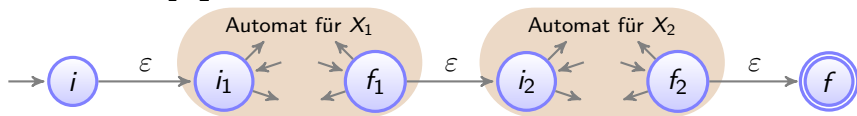
Automat für \emptyset :



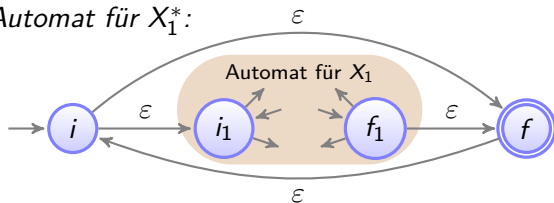
Automat für $s \in \Sigma \cup \{\varepsilon\}$:



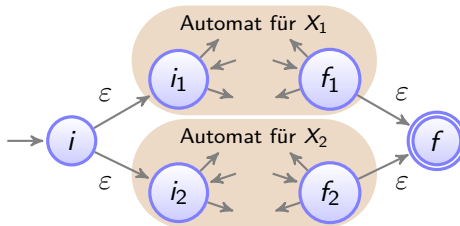
Automat für X_1X_2 :



Automat für X_1^* :

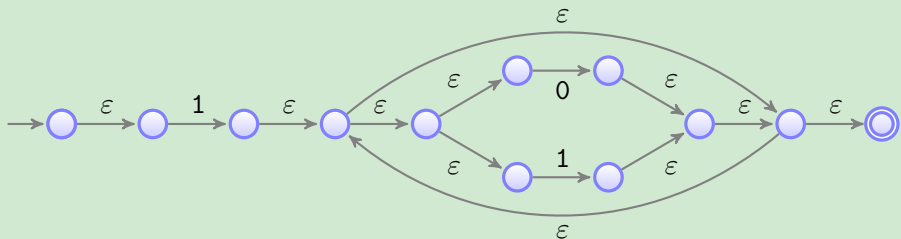


Automat für $X_1 + X_2$:

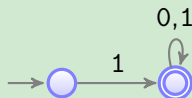


Automat für Binärnumerae ohne führende Null

Regulärer Ausdruck: $1(0 + 1)^*$



Minimaler deterministischer Automat:



Manuelle Konstruktion von Automaten:

- 1 Konstruiere die Automaten zu einfachen Teilsprachen „durch Hinschauen“.
- 2 Verwende die allgemeine Konstruktion mit ϵ -Übergängen für undurchsichtige Situationen.

Formale Sprachen – Übersicht

- 1 Operationen auf formalen Sprachen
- 2 Reguläre Sprachen
- 3 Reguläre Ausdrücke
- 4 Eigenschaften regulärer Sprachen
- 5 Vom regulären Ausdruck zum Automaten
- 6 Vom Automaten zum regulären Ausdruck**
- 7 Kontextfreie Sprachen
- 8 Beispiele kontextfreier Grammatiken
- 9 Andere Notationen
- 10 Style Guide für Grammatiken
- 11 Grammatiken in freier Wildbahn
- 12 Jenseits der Kontextfreiheit

Vom Automaten zum regulären Ausdruck

$R \dots$ Menge der regulären Ausdrücke über Σ

Verallgemeinerter endlicher Automat

\dots wird beschrieben durch ein 5-Tupel $\mathcal{A} = \langle Q, \Sigma, \delta, i, f \rangle$, wobei

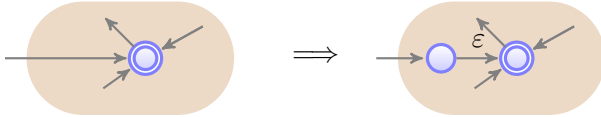
- $Q \dots$ endliche Zustandsmenge
- $\Sigma \dots$ Eingabealphabet
- $\delta: (Q - \{f\}) \times (\Sigma - \{\epsilon\}) \rightarrow R \dots$ Übergangsfunktion
- $i \in Q \dots$ Anfangszustand
- $f \in Q, f \neq i \dots$ Endzustand

Unterschiede zu „normalen“ Automaten:

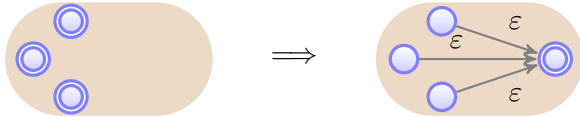
- keine Übergänge in den Anfangszustand;
- nur ein Endzustand, der nicht Anfangszustand ist;
- keine Übergänge weg vom Endzustand;
- nur ein Übergang zwischen je zwei Zuständen;
- Übergänge beschriftet mit regulären Ausdrücken.

Umwandlung eines endlichen Automaten in einen verallgemeinerten

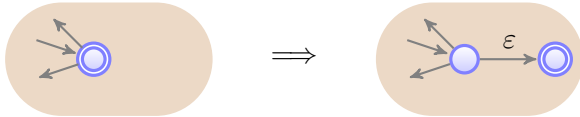
- Übergänge in den Anfangszustand oder Anfangszustand ist Endzustand:



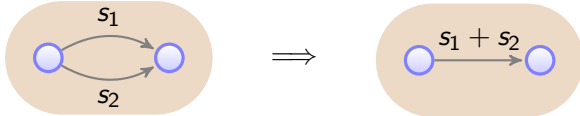
- Mehrere Endzustände:



- Übergänge weg vom Endzustand:



- Mehrere Übergänge zwischen zwei Zuständen:

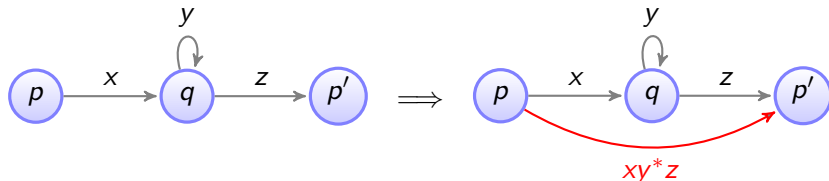


Vom verallgemeinerten Automaten zum regulären Ausdruck

Gegeben: Verallgemeinerter Automat $\mathcal{A} = \langle Q, \Sigma, \delta, i, f \rangle$

Für jeden Zustand $q \in Q - \{i, f\}$ führe folgende Schritte durch:

- 1 Füge zwischen allen Nachbarn p, p' von q neue Übergänge hinzu:



(x , y und z bezeichnen reguläre Ausdrücke.)

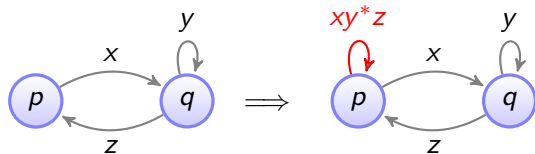
- 2 Entferne q und alle Kanten von und nach q aus dem Automaten.



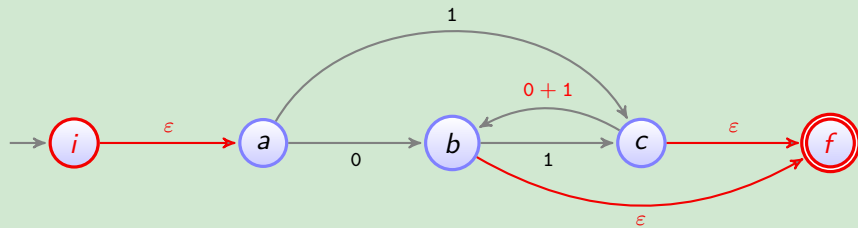
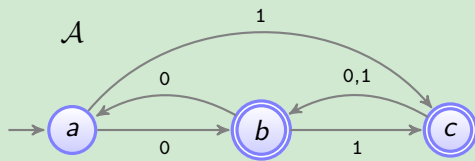
Ergebnis: r ist ein regulärer Ausdruck mit $\mathcal{L}(r) = \mathcal{L}(\mathcal{A})$.

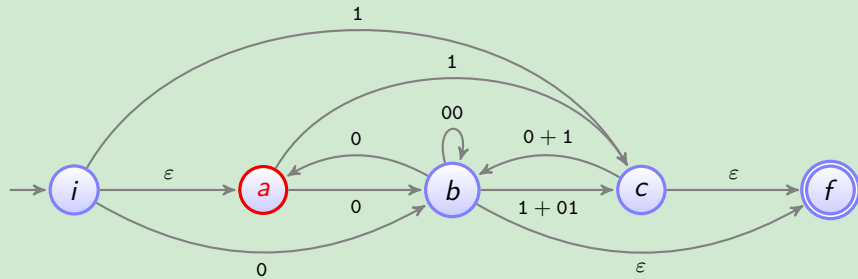
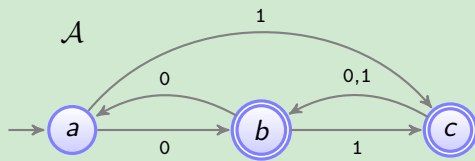
Anmerkungen:

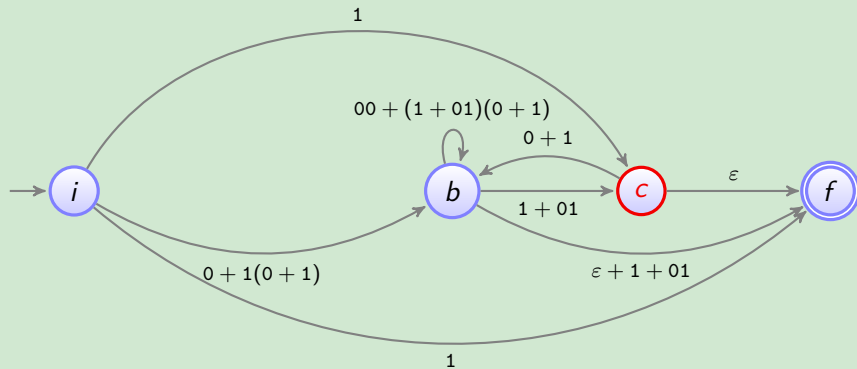
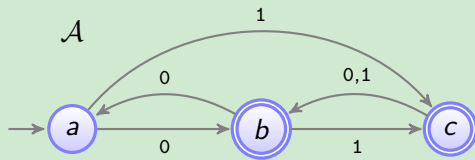
- Falls p und p' derselbe Knoten sind, erhält man:

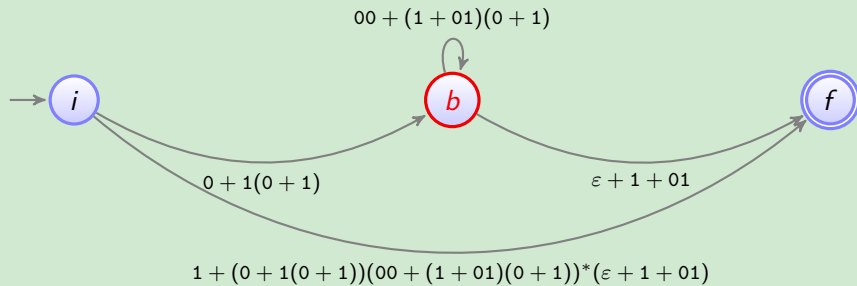
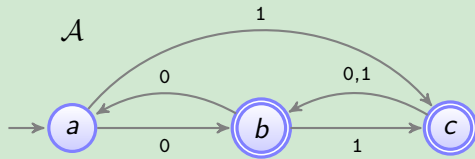


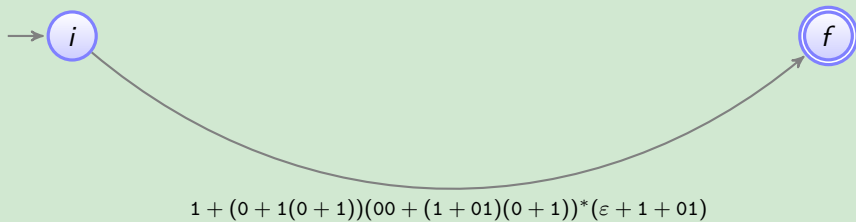
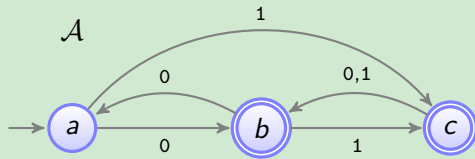
- Falls die Schleife mit dem Ausdruck y nicht existiert, entfällt y^* .
- Falls der Übergang $p-p'$ bereits existiert, wird xy^*z addiert.











$$r = 1 + (0 + 10 + 11)(00 + 10 + 11 + 010 + 011)^*(\epsilon + 1 + 01)$$

Es gilt: $\mathcal{L}(r) = \mathcal{L}(\mathcal{A})$.

Formale Sprachen – Übersicht

- 1 Operationen auf formalen Sprachen
- 2 Reguläre Sprachen
- 3 Reguläre Ausdrücke
- 4 Eigenschaften regulärer Sprachen
- 5 Vom regulären Ausdruck zum Automaten
- 6 Vom Automaten zum regulären Ausdruck
- 7 Kontextfreie Sprachen**
- 8 Beispiele kontextfreier Grammatiken
- 9 Andere Notationen
- 10 Style Guide für Grammatiken
- 11 Grammatiken in freier Wildbahn
- 12 Jenseits der Kontextfreiheit

Grenzen regulärer Sprachen

Was reguläre Ausdrücke und endliche Automaten beschreiben können:

- lexikale Elemente in Programmiersprachen: Identifier, Numerale, ...
- Morphologie von Worten in natürlichen Sprachen: korrekte Fall-/Zahl-/Zeitendungen, ...
- Systeme, die nur ein endliches Gedächtnis besitzen

Was sie nicht beschreiben können:

- geklammerte Ausdrücke in Programmiersprachen
- geschachtelte Programmstrukturen wie
`if ... while ... if ... endif ... endwhile ... endif`
- Schachtelung von Haupt- und Nebensätzen in natürlichen Sprachen
- Systeme mit einem (potentiell) unendlichen Speicher

Sprache der wohlgeklammerten Ausdrücke ist nicht regulär

$\{(), (()), ()(), ((())), (()()), ()()(), \dots\}$

$\{a^n b^n \mid n \geq 0\}$ ist nicht regulär

$\{\epsilon, ab, aabb, aaabbb, aaaabbbb, aaaaabbbbb, \dots\}$

Kontextfreie Grammatiken

- Menge von Ersetzungsregeln
- 2 Arten von Symbolen: Terminale, Nonterminale
- Ausgehend vom Start-Nonterminal werden Nonterminale solange ersetzt, bis nur noch Terminale bleiben.

$$\begin{aligned}\text{Satz} &\rightarrow \text{HwP ZwP} \\ \text{HwP} &\rightarrow \text{Art Hw} \\ \text{ZwP} &\rightarrow \text{Hzw HwP Zw}\end{aligned}$$
$$\begin{aligned}\text{Art} &\rightarrow \text{der} \mid \text{den} \\ \text{Hw} &\rightarrow \text{Hund} \mid \text{Mond} \\ \text{Hzw} &\rightarrow \text{wird} \\ \text{Zw} &\rightarrow \text{anbellen}\end{aligned}$$
$$\begin{aligned}\text{Satz} &\Rightarrow_p \text{HwP ZwP} \\ &\Rightarrow_p \text{Art Hw Hzw HwP Zw} \\ &\Rightarrow_p \text{der Hund wird Art Hw anbellen} \\ &\Rightarrow_p \text{der Hund wird den Mond anbellen}\end{aligned}$$

Generative Grammatiken gehen zurück auf *Noam Chomsky* (Linguist, Philosoph, Aktivist; *1928).

Kontextfreie Grammatik

... wird beschrieben durch ein 4-Tupel $G = \langle V, T, P, S \rangle$, wobei

- V ... Menge der Nonterminalsymbole (Variablen)
- T ... Menge der Terminalsymbole ($V \cap T = \{\}$)
- $P \subseteq V \times (V \cup T)^*$... Produktionen
- $S \in V$... Startsymbol

Schreibweise: $A \rightarrow w$ statt $(A, w) \in P$

$A \rightarrow w_1 \mid \dots \mid w_n$ statt $A \rightarrow w_1, \dots, A \rightarrow w_n$

Ableitbarkeit

... *in einem Schritt*:

$xAy \Rightarrow xwy$, falls $A \rightarrow w \in P$ und $x, y \in (V \cup T)^*$.

... *in mehreren Schritten*:

$u \xRightarrow{*} v$, falls

- $u = v$, oder
- $u \Rightarrow u'$ und $u' \xRightarrow{*} v$ für ein Wort $u' \in (V \cup T)^*$.

Von Grammatik G generierte Sprache

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

Grammatiken G_1 und G_2 heißen äquivalent, wenn $\mathcal{L}(G_1) = \mathcal{L}(G_2)$ gilt.

$$G = \langle \{S\}, \{a, b\}, \{S \rightarrow \varepsilon \mid a S b\}, S \rangle$$

$$S \xRightarrow{*} aabb, \text{ weil } S \Rightarrow a S b \Rightarrow aa S bb \Rightarrow aa \varepsilon bb = aabb$$

$$\mathcal{L}(G) = \{ a^n b^n \mid n \geq 0 \} = \{ \varepsilon, ab, aabb, aaabbb, \dots \}$$

Kontextfreie Sprachen

Eine Sprache heißt kontextfrei, wenn es eine kontextfreie Grammatik gibt, die sie generiert.

Verschiedene Ableitbarkeitsbegriffe

Linksableitbarkeit: $x A y \Rightarrow_L x w y$ falls $A \rightarrow w \in P$ und $x \in T^*$

(In jedem Schritt wird die linkeste Variable ersetzt.)

Rechtsableitbarkeit: $x A y \Rightarrow_R x w y$ falls $A \rightarrow w \in P$ und $y \in T^*$

(In jedem Schritt wird die rechteste Variable ersetzt.)

Parallelableitbarkeit: $x_0 A_1 x_1 \cdots A_n x_n \Rightarrow_P x_0 w_1 x_1 \cdots w_n x_n$

falls $A_1 \rightarrow w_1, \dots, A_n \rightarrow w_n \in P$ und $x_0, \dots, x_n \in T^*$

(In jedem Schritt werden alle Variablen ersetzt.)

- \Rightarrow_L^* , \Rightarrow_R^* und \Rightarrow_P^* sind eingeschränkte Ableitungsrelationen:
Manche Wörter $w \in (V \cup T)^*$ sind mit \Rightarrow^* herleitbar ($S \Rightarrow^* w$),
aber nicht mit diesen Relationen.
- Sie können aber jedes Wort der Sprache ableiten. Für alle $w \in T^*$ gilt:
 $S \Rightarrow^* w$ gdw. $S \Rightarrow_L^* w$ gdw. $S \Rightarrow_R^* w$ gdw. $S \Rightarrow_P^* w$

$$\begin{aligned}\mathcal{L}(G) &= \{ w \in T^* \mid S \Rightarrow^* w \} = \{ w \in T^* \mid S \Rightarrow_L^* w \} \\ &= \{ w \in T^* \mid S \Rightarrow_R^* w \} = \{ w \in T^* \mid S \Rightarrow_P^* w \}\end{aligned}$$

Formale Sprachen – Übersicht

- 1 Operationen auf formalen Sprachen
- 2 Reguläre Sprachen
- 3 Reguläre Ausdrücke
- 4 Eigenschaften regulärer Sprachen
- 5 Vom regulären Ausdruck zum Automaten
- 6 Vom Automaten zum regulären Ausdruck
- 7 Kontextfreie Sprachen
- 8 Beispiele kontextfreier Grammatiken**
- 9 Andere Notationen
- 10 Style Guide für Grammatiken
- 11 Grammatiken in freier Wildbahn
- 12 Jenseits der Kontextfreiheit

Kontextfreie Grammatiken – Beispiele

Syntax aussagenlogischer Formeln

$G = \langle \{Fm, Op, Var\}, T, P, Fm \rangle$, wobei

$T = \{ \top, \perp, \neg, (,), \wedge, \uparrow, \vee, \downarrow, \Leftrightarrow, \nLeftrightarrow, \Rightarrow, \Leftarrow \} \cup \mathcal{V}$

$P = \{ \begin{array}{l} Fm \rightarrow Var \mid \top \mid \perp \mid \neg Fm \mid (Fm Op Fm) , \\ Op \rightarrow \wedge \mid \uparrow \mid \vee \mid \downarrow \mid \Leftrightarrow \mid \nLeftrightarrow \mid \Rightarrow \mid \Leftarrow , \\ Var \rightarrow A \mid B \mid C \mid \dots \end{array} \}$

$((A \uparrow B) \uparrow (A \uparrow B))$ ist eine aussagenlogische Formel, weil:

$Fm \Rightarrow_L (Fm Op Fm)$	$\Rightarrow_L ((Fm Op Fm) Op Fm)$
$\Rightarrow_L ((Var Op Fm) Op Fm)$	$\Rightarrow_L ((A Op Fm) Op Fm)$
$\Rightarrow_L ((A \uparrow Fm) Op Fm)$	$\Rightarrow_L ((A \uparrow Var) Op Fm)$
$\Rightarrow_L ((A \uparrow B) Op Fm)$	$\Rightarrow_L ((A \uparrow B) \uparrow Fm)$
$\Rightarrow_L ((A \uparrow B) \uparrow (Fm Op Fm))$	$\Rightarrow_L ((A \uparrow B) \uparrow (Var Op Fm))$
$\Rightarrow_L ((A \uparrow B) \uparrow (A Op Fm))$	$\Rightarrow_L ((A \uparrow B) \uparrow (A \uparrow Fm))$
$\Rightarrow_L ((A \uparrow B) \uparrow (A \uparrow Var))$	$\Rightarrow_L ((A \uparrow B) \uparrow (A \uparrow B))$

Kontextfreie Grammatiken – Beispiele

Syntax aussagenlogischer Formeln

$G = \langle \{Fm, Op, Var\}, T, P, Fm \rangle$, wobei

$T = \{\top, \perp, \neg, (,), \wedge, \uparrow, \vee, \downarrow, \Leftrightarrow, \nLeftrightarrow, \Rightarrow, \Leftarrow\} \cup \mathcal{V}$

$P = \{ Fm \rightarrow Var \mid \top \mid \perp \mid \neg Fm \mid (Fm Op Fm) ,$
 $Op \rightarrow \wedge \mid \uparrow \mid \vee \mid \downarrow \mid \Leftrightarrow \mid \nLeftrightarrow \mid \Rightarrow \mid \Leftarrow ,$
 $Var \rightarrow A \mid B \mid C \mid \dots \}$

$((A \uparrow B) \uparrow (A \uparrow B))$ ist eine aussagenlogische Formel, weil:

$Fm \Rightarrow_P (Fm Op Fm)$
 $\Rightarrow_P ((Fm Op Fm) \uparrow (Fm Op Fm))$
 $\Rightarrow_P ((Var \uparrow Var) \uparrow (Var \uparrow Var))$
 $\Rightarrow_P ((A \uparrow B) \uparrow (A \uparrow B))$

Kontextfreie Grammatiken - Beispiele

Reelle Numerale

$G = \langle V, T, P, Real \rangle$, wobei

$V = \{ Real, Scale, Sign, Digits, Digit \}$

$T = \{ 0, \dots, 9, ., E, +, - \}$

und P folgende Produktionen sind:

$Real \rightarrow Digit Digits . Digits Scale$

$Scale \rightarrow \varepsilon \mid E Sign Digit Digits$

$Sign \rightarrow \varepsilon \mid + \mid -$

$Digits \rightarrow \varepsilon \mid Digit Digits$

$Digit \rightarrow 0 \mid \dots \mid 9$

Optionalität:

$A \rightarrow \varepsilon \mid B$ (A steht für optionales B)

Wiederholung:

$A \rightarrow \varepsilon \mid B A$ (A steht für wiederholtes B , auch 0 Mal)

$A \rightarrow B \mid B A$ (A steht für wiederholtes B , mind. 1 Mal)

Kontextfreie Grammatiken - Beispiele

Wohlgeformte Klammerausdrücke

WKA ist die kleinste Menge, sodass

- $\varepsilon \in WKA$
- $(w) \in WKA$, wenn $w \in WKA$
- $w_1 w_2 \in WKA$, wenn $w_1, w_2 \in WKA$

$$G = \langle \{W\}, \{ (,) \}, \{ W \rightarrow \varepsilon \mid (W) \mid W W \}, W \rangle$$

Beispiel einer Ableitung: $W \Rightarrow_P W W$
 $\Rightarrow_P (W)(W)$
 $\Rightarrow_P ()(W W)$
 $\Rightarrow_P ()((W)(W))$
 $\Rightarrow_P ()((())())$

$$\begin{aligned} \mathcal{L}(G) &= \{ \varepsilon, (), (()), ()(), ((())), (()()), ()(), ()()(), \dots \} \\ &= WKA \end{aligned}$$

Induktive Definition vs. kontextfreie Grammatik

Induktive Definition für \mathcal{M}

Mengen $\mathcal{M}, \mathcal{M}_0, \mathcal{A}_1, \mathcal{B}_1, \dots$

\mathcal{M} ist die kleinste Menge, sodass:

$$\mathcal{M}_0 \subseteq \mathcal{M}$$

$$f(x_1, \dots, x_m) \in \mathcal{M}, \text{ falls } x_1 \in \mathcal{A}_1, \dots, x_m \in \mathcal{A}_m$$

$$g(x_1, \dots, x_n) \in \mathcal{M}, \text{ falls } x_1 \in \mathcal{B}_1, \dots, x_n \in \mathcal{B}_n$$

$$\dots \in \mathcal{M}, \text{ falls } \dots$$

$$h(x_1, x_2) \in \mathcal{M}, \text{ falls } x_1, x_2 \in \mathcal{M}$$

$$h(\mathbf{x}, \mathbf{x}) \in \mathcal{M}, \text{ falls } \mathbf{x} \in \mathcal{M}$$

kontextfreie Grammatik für \mathcal{M}

Nonterminale M, M_0, A_1, B_1, \dots

$$\mathcal{M} = \mathcal{L}(\langle \dots, P, M \rangle), \text{ wobei } P:$$

$$M \rightarrow M_0$$

$$M \rightarrow f(A_1, \dots, A_m)$$

$$M \rightarrow g(B_1, \dots, B_n)$$

$$M \rightarrow \dots$$

$$M \rightarrow h(M, M)$$

keine Entsprechung, nicht kontextfrei

Formale Sprachen – Übersicht

- 1 Operationen auf formalen Sprachen
- 2 Reguläre Sprachen
- 3 Reguläre Ausdrücke
- 4 Eigenschaften regulärer Sprachen
- 5 Vom regulären Ausdruck zum Automaten
- 6 Vom Automaten zum regulären Ausdruck
- 7 Kontextfreie Sprachen
- 8 Beispiele kontextfreier Grammatiken
- 9 Andere Notationen**
- 10 Style Guide für Grammatiken
- 11 Grammatiken in freier Wildbahn
- 12 Jenseits der Kontextfreiheit

Backus-Naur-Form (BNF)

Synonym für kontextfreie Produktionen

Historische Notation:

- Nonterminale in spitzen Klammern
- Terminale unter Anführungszeichen
- $::=$ als Trennsymbol

Durch diese Konvention entfällt die Notwendigkeit, die Mengen der Nonterminale (V) und der Terminale (T) anzugeben.

```
 $\langle Real \rangle ::= \langle Digit \rangle \langle Digits \rangle " ." \langle Digits \rangle \langle Scale \rangle$   
 $\langle Digit \rangle ::= "0" \mid \dots \mid "9"$   
 $\langle Digits \rangle ::= \varepsilon \mid \langle Digit \rangle \langle Digits \rangle$   
 $\langle Scale \rangle ::= \varepsilon \mid "E" \langle Sign \rangle \langle Digit \rangle \langle Digits \rangle$   
 $\langle Sign \rangle ::= \varepsilon \mid "+" \mid "-"$ 
```

Erweiterte Backus-Naur-Form (EBNF)

„Regular Right Part Grammar“ (RRPG)

- erlaubt reguläre Ausdrücke auf der rechten Seite von Produktionen
- Bessere Lesbarkeit durch Vermeidung von Rekursionen und Leerwörtern
- Kompaktere Spezifikationen
- keine echte Erweiterung: EBNF-Notationen lassen sich eliminieren

Elimination der EBNF-Notation:

$A \rightarrow w_1 (w_2) w_3$ entspricht $A \rightarrow w_1 B w_3$
 $B \rightarrow w_2$

$A \rightarrow w_1 \{w_2\} w_3$ entspricht $A \rightarrow w_1 B w_3$
 $B \rightarrow \varepsilon \mid w_2 B$

$A \rightarrow w_1 [w_2] w_3$ entspricht $A \rightarrow w_1 B w_3$
 $B \rightarrow \varepsilon \mid w_2$

Listen von Bezeichnern

EBNF: $IdList \rightarrow Id \{ ", " Id \}$

Ohne EBNF: $IdList \rightarrow Id B$
 $B \rightarrow \varepsilon \mid ", " Id B$

oder $IdList \rightarrow Id \mid Id ", " IdList$

Skalierungsfaktor der reellen Numerale

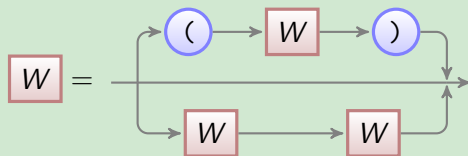
EBNF: $Scale \rightarrow "E" ["+" \mid "-"] Digit \{ Digit \}$

Ohne EBNF: $Scale \rightarrow "E" B_1 Digit B_2$
 $B_1 \rightarrow \varepsilon \mid "+" \mid "-"$
 $B_2 \rightarrow \varepsilon \mid Digit B_2$

Syntaxdiagramme

Mit *rekursiven* Diagrammen können beliebige kontextfreie Grammatiken in EBNF dargestellt werden.

Wohlgeklammerte Ausdrücke


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

Formale Sprachen – Übersicht

- 1 Operationen auf formalen Sprachen
- 2 Reguläre Sprachen
- 3 Reguläre Ausdrücke
- 4 Eigenschaften regulärer Sprachen
- 5 Vom regulären Ausdruck zum Automaten
- 6 Vom Automaten zum regulären Ausdruck
- 7 Kontextfreie Sprachen
- 8 Beispiele kontextfreier Grammatiken
- 9 Andere Notationen
- 10 Style Guide für Grammatiken**
- 11 Grammatiken in freier Wildbahn
- 12 Jenseits der Kontextfreiheit

Style Guide für Grammatiken

Wiederholungsoperator statt Rekursion (wenn möglich).

Schlecht: $A \rightarrow a \mid A A$ oder $A \rightarrow a \mid a A$

Besser: $a \{ a \}$ oder $a +$ oder \dots (je nach gewählter Notation)

Optionsoperator statt Leerwort.

Schlecht: $\varepsilon \mid A$

Besser: $[A]$ oder $A?$ oder \dots (je nach gewählter Notation)

Modularisierung statt Duplizierung.

Schlecht: $Identifier \rightarrow ("a" \mid \dots \mid "z") \{ "a" \mid \dots \mid "z" \mid "0" \mid \dots \mid "9" \}$

Besser: $Identifier \rightarrow Letter \{ Letter \mid Digit \}$
 $Letter \rightarrow "a" \mid \dots \mid "z"$
 $Digit \rightarrow "0" \mid \dots \mid "9"$

Verwendung sprechender Namen.

Schlecht: $X \rightarrow Y \{ Y \mid Z \}$
 $Y \rightarrow "a" \mid \dots \mid "z"$
 $Z \rightarrow "0" \mid \dots \mid "9"$

Besser: Siehe oben.

Klare Unterscheidung zwischen Terminalen, Nonterminalen und Metanotationen.

Schlecht: `\begin{tabular}[[Position]]{Spalte{Spalte}}`

Besser: `"\begin{tabular}" ["[" Position "]"] "{" Spalte { Spalte } "}"`
`"\begin{tabular}" ["[" $\langle Position \rangle$ "]"] "{" $\langle Spalte \rangle$ { $\langle Spalte \rangle$ } "}"`

Inhaltliche Strukturierung statt Minimierung der Nonterminale und Regeln.

Schlecht:

`"\begin{tabular}" ["[" ("t" | "b") "]"] "{" ("l" | "c" | "r") { "l" | "c" | "r" } "}"`

Besser: `... → "\begin{tabular}" [Position] Spalten`

`Position → "[b]" | "[t]"`

`Spalten → "{" Spalte { Spalte } "}"`

`Spalte → "l" | "c" | "r"`

Beispiel: Tabellen in LaTeX

TeX ... Textsatzsystem von Donald E. Knuth

LaTeX ... TeX-Makros für „logisches Markup“ von Leslie Lamport

```
\begin{tabular}[t]{lc}  
Eintrag 11 & Eintrag 12 \\  
Eintrag 21 & \\\br/>  \begin{tabular}{rr}  
    Eintrag 22 & ist selber \\  
    eine       & & Tabelle.  
  \end{tabular} \\  
Eintrag 31 & Eintrag 32  
\end{tabular}
```

Eintrag 11	Eintrag 12
Eintrag 21	Eintrag 22 ist selber eine Tabelle.
Eintrag 31	Eintrag 32

Gesucht: Kontextfreie Grammatik G in EBNF für die Sprache der LaTeX-Tabellen

Beispiel: Tabellen in LaTeX

$G = \langle V, T, P, \textit{Tabelle} \rangle$, wobei:

$P = \{ \textit{Tabelle} \rightarrow "\backslash\textit{begin}\{\textit{tabular}\}" [\textit{Position}] \textit{Spalten}$
 \textit{Zeilen}
 $"\backslash\textit{end}\{\textit{tabular}\}" ,$

$\textit{Position} \rightarrow "[\textit{b}] " | "[\textit{t}] " ,$

$\textit{Spalten} \rightarrow "\{ " \textit{Spalte} \{ \textit{Spalte} \} " " ,$

$\textit{Spalte} \rightarrow "\textit{l} " | "\textit{c} " | "\textit{r} " ,$

$\textit{Zeilen} \rightarrow \textit{Zeile} \{ "\backslash\backslash " \textit{Zeile} \} ,$

$\textit{Zeile} \rightarrow \textit{Eintrag} \{ "&" \textit{Eintrag} \} ,$

$\textit{Eintrag} \rightarrow \{ \textit{Tabelle} | \textit{Zeichen} \} ,$

$\textit{Zeichen} \rightarrow "0" | \dots | "9" | "a" | \dots | "Z" | "." | "\square" \}$

$V = \{ \textit{Tabelle}, \textit{Position}, \textit{Spalten}, \textit{Spalte}, \textit{Zeilen}, \textit{Zeile}, \textit{Eintrag}, \textit{Zeichen} \}$

$T = \{ "0", \dots, "9", "a", \dots, "z", "A", \dots, "Z",$
 $".", "\square", "\{", "\}", "[", "]", "&", "\backslash" \}$

Nur eine Rekursion für Schachtelung der Tabellen notwendig!

Formale Sprachen – Übersicht

- 1 Operationen auf formalen Sprachen
- 2 Reguläre Sprachen
- 3 Reguläre Ausdrücke
- 4 Eigenschaften regulärer Sprachen
- 5 Vom regulären Ausdruck zum Automaten
- 6 Vom Automaten zum regulären Ausdruck
- 7 Kontextfreie Sprachen
- 8 Beispiele kontextfreier Grammatiken
- 9 Andere Notationen
- 10 Style Guide für Grammatiken
- 11 Grammatiken in freier Wildbahn**
- 12 Jenseits der Kontextfreiheit

N. Wirth: Programming in Modula-2

Appendix 1

The Syntax of Modula-2

```
1  ident = letter (letter | digit).
2  number = integer | real.
3  integer = digit (digit) | octalDigit (octalDigit) ("B"|"C") |
4    digit (hexDigit) "H".
5  real = digit (digit) "." (digit) [ScaleFactor].
6  ScaleFactor = "E" | "e" | "+" | "-" digit (digit).
7  hexDigit = digit | "A"|"B"|"C"|"D"|"E"|"F".
8  digit = octalDigit | "0"|"1"|"2"|"3"|"4"|"5"|"6"|"7".
9  octalDigit = "0"|"1"|"2"|"3"|"4"|"5"|"6"|"7".
10 string = "" (character) "" | "" (character) "".
11 qualident = ident ("." ident).
12 ConstantDeclaration = ident "=" ConstExpression.
13 ConstExpression = SimpleConstExpr [relation SimpleConstExpr].
14 relation = "<" | "<=" | ">" | ">=" | "=" | "<=" | ">=" | IN.
15 SimpleConstExpr = ["+"|"-"] ConstTerm (AddOperator ConstTerm).
16 AddOperator = "+" | "-" | OR.
17 ConstTerm = ConstFactor (MulOperator ConstFactor).
18 MulOperator = "*" | "/" | DIV | MOD | AND | "&".
19 ConstFactor = qualident | number | string | set |
20   "(" ConstExpression ")" | NOT ConstFactor.
21 set = [qualident] "(" (element ("," element) ")".
22 element = ConstExpression ["." ConstExpression].
23 TypeDeclaration = ident "=" type.
24 type = SimpleType | ArrayType | RecordType | SetType |
25   PointerType | ProcedureType.
26 SimpleType = qualident | enumeration | SubrangeType.
27 enumeration = "(" identList ")".
28 identList = ident ("," ident).
29 SubrangeType = "[" ConstExpression "," ConstExpression "]" .
30 ArrayType = ARRAY SimpleType (":" SimpleType) OF type.
31 RecordType = RECORD FieldListSequence END.
32 FieldListSequence = FieldList (";" FieldList).
33 FieldList = [identList ":" type]
34   CASE [ident ":"] qualident OF variant ("|" variant)
35     [ELSE FieldListSequence] END.
36 variant = CaseLabelList ":" FieldListSequence.
37 CaseLabelList = CaseLabels ("," CaseLabels).
38 CaseLabels = ConstExpression ["." ConstExpression].
39 SetType = SET OF SimpleType.
40 PointerType = POINTER TO type.
41 ProcedureType = PROCEDURE [FormalTypeList].
42 FormalTypeList = "(" ("[" [VAR] FormalType
43   ("," [VAR] FormalType)] ")" (";" [":"qualident].
44 VariableDeclaration = identList ":" type.
```

158 A1. The Syntax of Modula-2

```
45 designator = qualident ("[" ident | "[" ExpList "]" | ")" .
46 ExpList = expression ("," expression).
47 expression = SimpleExpression [relation SimpleExpression].
48 SimpleExpression = ["+"|"-"] term (AddOperator term).
49 term = factor (MulOperator factor).
50 factor = number | string | set | designator [ActualParameters] |
51   "(" expression ")" | NOT factor.
52 ActualParameters = "(" [ExpList] ")" .
53 statement = [assignment | ProcedureCall |
54   IFStatement | CaseStatement | WhileStatement |
55   RepeatStatement | LoopStatement | ForStatement |
56   WithStatement | EXIT | RETURN [expression] ].
57 assignment = designator "=" expression.
58 ProcedureCall = designator [ActualParameters].
59 StatementSequence = statement (";" statement).
60 IFStatement = IF expression THEN StatementSequence
61   [ELSEIF expression THEN StatementSequence]
62   [ELSE StatementSequence] END.
63 CaseStatement = CASE expression OF case ("|" case)
64   [ELSE StatementSequence] END.
65 case = CaseLabelList ":" StatementSequence.
66 WhileStatement = WHILE expression DO StatementSequence END.
67 RepeatStatement = REPEAT StatementSequence UNTIL expression.
68 ForStatement = FOR ident "=" expression TO expression
69   [BY ConstExpression] DO StatementSequence END.
70 LoopStatement = LOOP StatementSequence END.
71 WithStatement = WITH designator DO StatementSequence END.
72 ProcedureDeclaration = ProcedureHeading ":" block ident.
73 ProcedureHeading = PROCEDURE ident [FormalParameters].
74 block = (declaration) (BEGIN StatementSequence) END.
75 declaration = CONST (ConstantDeclaration ";" ) |
76   TYPE (TypeDeclaration ";" ) |
77   VAR (VariableDeclaration ";" ) |
78   ProcedureDeclaration ";" | ModuleDeclaration ";".
79 FormalParameters =
80   "(" [FPSection (";" FPSection)] ")" [";" qualident].
81 FPSection = [VAR] identList ":" FormalType.
82 FormalType = [ARRAY OF] qualident.
83 ModuleDeclaration =
84   MODULE ident [priority] ":" {import} [export] block ident.
85 priority = "(" ConstExpression ")".
86 export = EXPORT [QUALIFIED] identList ":".
87 import = [FROM ident] IMPORT identList ":".
88 DefinitionModule = DEFINITION MODULE ident ":" {import}
89   [export] (definition) END ident ":".
90 definition = CONST (ConstantDeclaration ";" ) |
91   TYPE (ident ["=" type] ";" ) |
92   VAR (VariableDeclaration ";" ) |
93   ProcedureHeading ":".
94 ProgramModule =
95   MODULE ident [priority] ":" {import} block ident ":".
```

N. Wirth: Programming in Modula-2

Appendix 1

The Syntax of Modula-2

```
1 ident = letter {letter | digit}.
2 number = integer | real.
3 integer = digit {digit} | octalDigit {octalDigit}
4         digit {hexDigit} "H".
5 real = digit {digit} "." {digit} [ScaleFactor].
6 ScaleFactor = "E" ["+"|"−"] digit {digit}.
7 hexDigit = digit | "A"|"B"|"C"|"D"|"E"|"F".
8 digit = octalDigit | "8"|"9".
9 octalDigit = "0"|"1"|"2"|"3"|"4"|"5"|"6"|"7".
10 string = "" {character} "" | "" {character} "
11 qualified = ident ("." ident).
12 ConstantDeclaration = ident "=" ConstExpr.
13 ConstExpr = SimpleConstExpr [relation
14 relation = "<" | "<=" | ">" | ">=" | "=" | "<=" | ">=" | "="].
15 SimpleConstExpr = ["+"|"−"] ConstTerm (AddOperator
16 AddOperator = "+" | "−" | OR.
17 ConstTerm = ConstFactor (MulOperator Const
18 MulOperator = "*" | "/" | DIV | MOD | AND.
19 ConstFactor = qualified | number | string |
20 ("[" ConstExpr "]" | NOT ConstExpr).
21 set = [qualified] "(" {element "(" element
22 element = ConstExpr ["." ConstExpr].
23 TypeDeclaration = ident "=" type.
24 type = SimpleType | ArrayType | RecordType | SetType |
25 PointerType | ProcedureType.
26 SimpleType = qualified | enumeration | SubrangeType.
27 enumeration = "(" identList ")".
28 identList = ident ("." ident).
29 SubrangeType = "[" ConstExpr ":" ConstExpr "]"
30 ArrayType = ARRAY SimpleType "(" SimpleType ")" OF type.
31 RecordType = RECORD FieldListSequence END.
32 FieldListSequence = FieldList (";" FieldList).
33 FieldList = [identList ":" type].
34 CASE [ident ":"] qualified OF variant ("|" variant)
35 [ELSE FieldListSequence] END.
36 variant = CaseLabelList ":" FieldListSequence.
37 CaseLabelList = CaseLabels (";" CaseLabels).
38 CaseLabels = ConstExpr ["." ConstExpr].
39 SetType = SET OF SimpleType.
40 PointerType = POINTER TO type.
41 ProcedureType = PROCEDURE [FormalTypeList].
42 FormalTypeList = "(" ["[VAR] FormalType"
43 ("[" VAR] FormalType)"] ")" (";" FormalTypeList).
44 VariableDeclaration = identList ":" type.
```

```
1 ident = letter {letter | digit}.
2 number = integer | real.
3 integer = digit {digit} | octalDigit {octalDigit} ("B"|"C") |
4         digit {hexDigit} "H".
5 real = digit {digit} "." {digit} [ScaleFactor].
6 ScaleFactor = "E" ["+"|"−"] digit {digit}.
7 hexDigit = digit | "A"|"B"|"C"|"D"|"E"|"F".
8 digit = octalDigit | "8"|"9".
9 octalDigit = "0"|"1"|"2"|"3"|"4"|"5"|"6"|"7".
```

```
72 ProcedureDeclaration = ProcedureHeading ":" block ident.
73 ProcedureHeading = PROCEDURE ident [FormalParameters].
74 block = (declaration) (BEGIN StatementSequence) END.
75 declaration = CONST (ConstantDeclaration ":" ) |
76 TYPE (TypeDeclaration ":" ) |
77 VAR (VariableDeclaration ":" ) |
78 ProcedureDeclaration ":" | ModuleDeclaration ":".
79 FormalParameters =
80 "(" [FPSection {";" FPSection}] ")" [";" qualified].
81 FPSection = [VAR] identList ":" FormalType.
82 FormalType = [ARRAY OF] qualified.
83 ModuleDeclaration =
84 MODULE ident [priority] ":" {import} [export] block ident.
85 priority = "(" ConstExpr ")" .
86 export = EXPORT [QUALIFIED] identList ":".
87 import = [FROM ident] IMPORT identList ":".
88 DefinitionModule = DEFINITION MODULE ident ":" {import}
89 [export] (definition) END ident ":".
90 definition = CONST (ConstantDeclaration ":" ) |
91 TYPE (ident ["=" type] ":" ) |
92 VAR (VariableDeclaration ":" ) |
93 ProcedureHeading ":".
94 ProgramModule =
95 MODULE ident [priority] ":" {import} block ident ":".
```

CHAPTER 18

Syntax

THIS chapter presents a grammar for the Java programming language. The grammar presented piecemeal in the preceding chapters (§2.3) is much better for exposition, but it is not well suited as a basis for a parser. The grammar presented in this chapter is the basis for the reference implementation. Note that it is not an LL(1) grammar, though in many cases it minimizes the necessary look ahead.

The grammar below uses the following BNF-style conventions:

- $[x]$ denotes zero or one occurrences of x .
- $\{x\}$ denotes zero or more occurrences of x .
- x / y means one of either x or y .

Identifier:
IDENTIFIER

QualifiedIdentifier:
Identifier $\{ \cdot \text{ } \textit{Identifier} \}$

QualifiedIdentifierList:
QualifiedIdentifier $\{ \cdot \text{ } \textit{QualifiedIdentifier} \}$

EnumBody:
 $\{ \textit{EnumConstants} \} [\cdot \text{ } \textit{EnumBodyDeclarations}] \}$

EnumConstants:
EnumConstant
EnumConstants $\cdot \text{ } \textit{EnumConstant}$

EnumConstant:
 $[\textit{Annotations}] \textit{Identifier} [\textit{Arguments}] [\textit{ClassBody}]$

EnumBodyDeclarations:
 $\cdot \text{ } [\textit{ClassBodyDeclaration}]$

AnnotationTypeBody:
 $[\textit{AnnotationTypeElementDeclarations}] \}$

AnnotationTypeElementDeclarations:
AnnotationTypeElementDeclaration
AnnotationTypeElementDeclarations *AnnotationTypeElementDeclaration*

AnnotationTypeElementDeclaration:
 $\{ \textit{Modifier} \} \textit{AnnotationTypeElementRest}$

AnnotationTypeElementRest:
Type Identifier AnnotationMethodOrConstantRest \cdot
ClassDeclaration
InterfaceDeclaration
EnumDeclaration
AnnotationTypeDeclaration

AnnotationMethodOrConstantRest:
AnnotationMethodRest
ConstantDeclaratorsRest

AnnotationMethodRest:
 $() [\{ \cdot \}] [\textit{default ElementValue}]$

CHAPTER 18

EnumBody:
{ [*EnumConstants*] [.] [*EnumBodyDeclarations*] }

EnumConstants:
EnumConstant

The grammar below uses the following BNF-style conventions:

- *[x]* denotes zero or one occurrences of *x*.
- *{x}* denotes zero or more occurrences of *x*.
- *x | y* means one of either *x* or *y*.

AnnotationTypeElementDeclaration:
{ [*Modifier*] [*AnnotationTypeElementRest*] }

AnnotationTypeElementDeclaration:
{ [*Modifier*] *AnnotationTypeElementRest* }

AnnotationTypeElementRest:
Type Identifier AnnotationMethodOrConstantRest ;
ClassDeclaration
InterfaceDeclaration
EnumDeclaration
AnnotationTypeDeclaration

AnnotationMethodOrConstantRest:
AnnotationMethodRest
ConstantDeclaratorsRest

AnnotationMethodRest:
() [(*l*)] [*default* *ElementValue*]

THIS chapter presents a grammar for the Java programming language. The grammar presented piecemeal in the preceding chapters (§§ 13.1–13.10) for exposition, but it is not well suited as a basis for a parser. The grammar in this chapter is the basis for the reference implementation. Note that the LL(1) grammar, though in many cases it minimizes the necessary

The grammar below uses the following BNF-style conventions:

- *[x]* denotes zero or one occurrences of *x*.
- *{x}* denotes zero or more occurrences of *x*.
- *x | y* means one of either *x* or *y*.

Identifier:
IDENTIFIER

QualifiedIdentifier:
Identifier [. *Identifier*]

QualifiedIdentifierList:
QualifiedIdentifier { , *QualifiedIdentifier* }

CHAPTER 18

Syntax

THIS chapter presents a grammar for the Java programming language. The grammar presented piecemeal in the preceding chapters (§2.3) is much better for exposition, but it is not well suited as a basis for a parser. The grammar presented in this chapter is the basis for the reference implementation. Note that it is not an LL(1) grammar, though in many cases it minimizes the necessary look ahead.

The grammar below uses the following BNF-style conventions

- $[x]$ denotes zero or one occurrences of x .
- $\{x\}$ denotes zero or more occurrences of x .
- x/y means one of either x or y .

Identifier:
IDENTIFIER

QualifiedIdentifier:
Identifier $\{ \cdot \textit{Identifier} \}$

QualifiedIdentifierList:
QualifiedIdentifier $\{ , \textit{QualifiedIdentifier} \}$

EnumBody:
 $\{ [\textit{EnumConstants}] [\cdot] [\textit{EnumBodyDeclarations}] \}$

EnumConstants:
EnumConstant
EnumConstants $, \textit{EnumConstant}$

EnumConstant:
 $[\textit{Annotations}] \textit{Identifier} [\textit{Arguments}] [\textit{ClassBody}]$

EnumBodyDeclarations:
 $;\{ \textit{ClassBodyDeclaration} \}$

AnnotationTypeBody:
 $\{ [\textit{AnnotationTypeDeclaration}] [\cdot] [\textit{AnnotationTypeDeclarations}] \}$

QualifiedIdentifier:
Identifier $\{ \cdot \textit{Identifier} \}$

QualifiedIdentifierList:
QualifiedIdentifier $\{ , \textit{QualifiedIdentifier} \}$

$() [\{ \{ \} \} [\textit{default ElementValue}]]$

CHAPTER 18

Syntax

THIS chapter presents a grammar for the Java programming language. The grammar presented piecemeal in the preceding chapters (§2.3) is much better for exposition, but it is not well suited as a basis for a parser. The grammar presented in this chapter is the basis for the reference implementation. Note that it is not an LL(1) grammar, though in many cases it minimizes the necessary look ahead.

EnumConstants:
EnumConstant
EnumConstants , EnumConstant

EnumConstant:
[Annotations] Identifier [Arguments] [ClassBody]

EnumBodyDeclarations:
; {ClassBodyDeclaration}

EnumBody:
{ [EnumConstants] [.] [EnumBodyDeclarations] }

EnumConstants:
EnumConstant
EnumConstants , EnumConstant

EnumConstant:
[Annotations] Identifier [Arguments] [ClassBody]

EnumBodyDeclarations:
; {ClassBodyDeclaration}

AnnotationTypeBody:
{ {AnnotationTypeElementDeclaration} }

AnnotationTypeElementDeclarations:
AnnotationTypeElementDeclaration
AnnotationTypeElementDeclarations AnnotationTypeElementDeclaration

AnnotationTypeElementDeclaration:
AnnotationTypeElementRest

AnnotationTypeElementRest:
MethodOrConstantRest ;

MethodOrConstantRest:
MethodRest ;

MethodRest:
Value

N. Wirth: Programming in Modula-2

Appendix 4

Syntax Diagrams

ident



number



integer



real



ScaleFactor

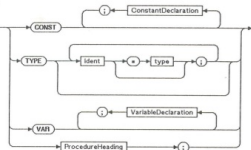


180 A4. Syntax Diagrams

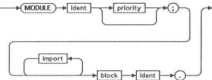
DefinitionModule



definition



ProgramModule



CompilationUnit



N. Wirth: Programming in Modula-2

180 A4. Syntax Diagrams

Appendix 4

Syntax Diagrams

ident



number



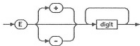
integer



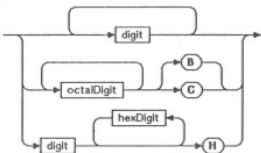
real



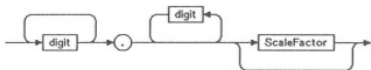
ScaleFactor



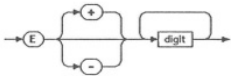
integer



real



ScaleFactor

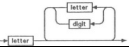


N. Wirth: Programming in Modula-2

Appendix 4

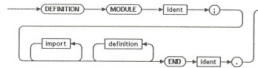
Syntax Diagrams

ident

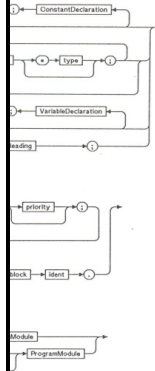
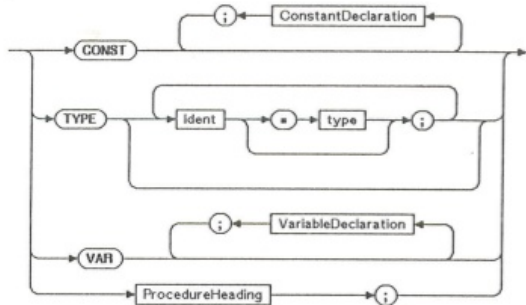


180 A4. Syntax Diagrams

DefinitionModule



definition



Formale Sprachen – Übersicht

- 1 Operationen auf formalen Sprachen
- 2 Reguläre Sprachen
- 3 Reguläre Ausdrücke
- 4 Eigenschaften regulärer Sprachen
- 5 Vom regulären Ausdruck zum Automaten
- 6 Vom Automaten zum regulären Ausdruck
- 7 Kontextfreie Sprachen
- 8 Beispiele kontextfreier Grammatiken
- 9 Andere Notationen
- 10 Style Guide für Grammatiken
- 11 Grammatiken in freier Wildbahn
- 12 Jenseits der Kontextfreiheit**

Jenseits der Kontextfreiheit

Formale Sprachen

$\{ a^n b^n c^n \mid n \geq 0 \}$ ist nicht kontextfrei.

Formale Sprachen

$\{ ww \mid w \in \{a, b, c\}^* \}$ ist nicht kontextfrei.

Programmiersprachen

Typen- und Deklarationsbedingungen nicht oder nur schwer kontextfrei darstellbar.

Natürliche Sprachen (Schweizer Dialekt)

Mer d'chind em Hans es huus
haend wele laa hälle aastriche.

Wir wollten die Kinder dem Hans
helfen lassen das Haus anzustreichen.

- Reguläre Sprache besitzen die gleich Ausdrucksstärke wie DEAs und NEAs.
- Gibt es auch einen passenden Automaten für Kontextfreie Sprachen?
⇨ Ja, nichtdeterministische Kellerautomaten (pushdown automata)
- Kellerautomat = endlicher Automat + Kellerspeicher (Stack)
 - Kellerspeicher ermöglicht unendlich viele „Zustände“
 - Automat sieht nur das oberste Symbol im Keller (und Zustandsübergänge hängen davon ab)
 - Bei einem Übergang können neue Symbole oben auf den Keller geschrieben werden oder das oberste Symbol kann gelöscht werden.

Andere Arten von Grammatiken

- Wir haben nur kontextfreie Grammatiken (Typ-2 Grammatik) betrachtet
- Es gibt ausdrucksstärkere Grammatiken:
 - uneingeschränkte formale Grammatik (Typ-0)
 - kontextsensitive Grammatik (Typ-1)
- ... und einfachere Arten von Grammatiken:
 - reguläre Grammatiken (Typ-3) definieren die regulären Sprachen und besitzen damit dieselbe Ausdrucksstärke wie endliche Automaten und reguläre Ausdrücke.

- Grammatik definieren formale Sprachen
- Kontextfreie Grammatiken
 - ausdrucksstärker als reguläre Ausdrücke und endliche Automaten
 - darstellbar als rekursive Syntaxdiagramme
- Style Guide für Grammatiken
 - EBNF kombiniert reguläre Ausdrücke und kontextfreie Produktionen