

Grundlagen der Informatik

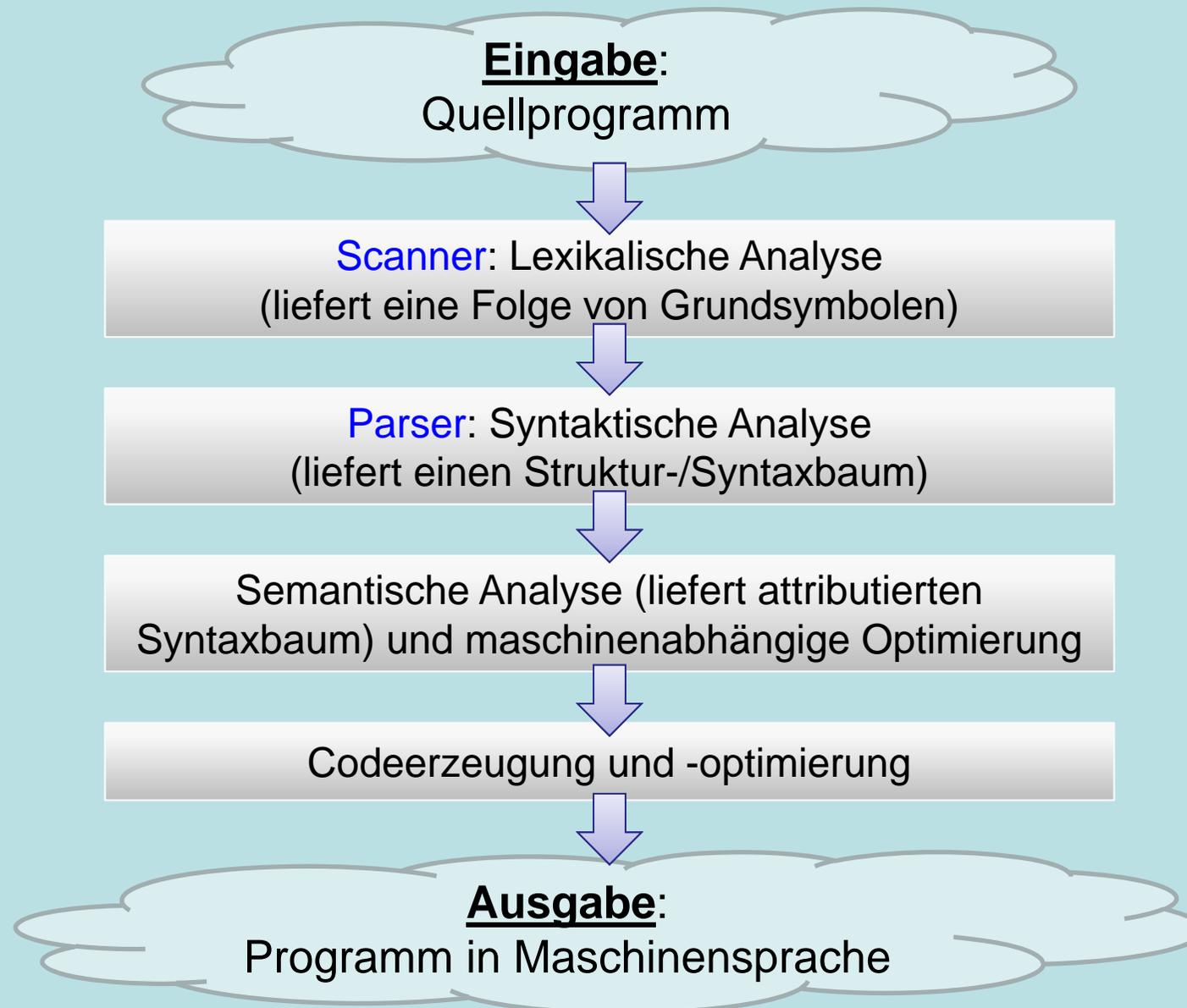
- Lexikalische und Syntaktische Analyse Forts. -

Prof. Dr. Klaus Volbert



Hochschule für angewandte Wissenschaften
Fakultät Informatik und Mathematik

Wintersemester 2010/11
Regensburg, 17./18. November 2010



Lexikalische und Syntaktische Analyse

- Mittels **formaler Sprachen** soll entschieden werden, ob ein Element zu einer vorgegebenen Klasse gehört oder nicht
- Wie kann die Menge aller lexikalisch und syntaktisch korrekten C/C++- Programme beschrieben werden?
- Sei L die Menge aller syntaktisch korrekten C/C++- Programme, wie kann dann für ein Programm p entschieden werden, ob $p \in L$ (**Wortproblem**)?
- Sprachen können durch **Grammatiken** beschrieben werden
- Grammatiken geben Regeln (**Produktionen**) an, nach denen syntaktisch korrekte Sätze aufgebaut werden dürfen
- Anmerkung: Ein großer Teil der Arbeit eines Informatikers besteht in der Definition neuer Sprachen (z. B. Protokolle)

- Ein **Alphabet** Σ ist eine endliche Menge von **Zeichen** ($\Sigma = \{0,1\}$)
- Ein **Wort** über einem Alphabet Σ ist eine endliche Folge von Zeichen aus Σ , alle Wörter der Länge n : Σ^n
- Die **Menge aller Wörter** über Σ ist Σ^*
- Das **leere Wort** ist ε , es gilt: $\varepsilon \in \Sigma^*$
- Wörter können zusammengesetzt werden (**Konkatenation**), es gilt: $\varepsilon u = u\varepsilon = u$, $\Sigma^+ = \Sigma \Sigma^* = \Sigma^* \setminus \{\varepsilon\}$ („mindestens einmal“)
- Die **Länge eines Wortes** ist die Anzahl der Zeichen, aus der sich ein Wort zusammensetzt, Notation: $|u|$, es gilt: $|\varepsilon| = 0$
- Eine **Sprache** L über einem Alphabet Σ ist eine Menge von Wörtern über Σ , d.h. $L \subseteq \Sigma^*$
- Anmerkung:
 - Ein Programm ist ein Wort, d.h. hier weicht der Begriff von dem Wortbegriff in natürlichen Sprachen (z.B. deutsch) ab

- Eine Grammatik G ist beschrieben durch ein 4-Tupel $G = (V, \Sigma, P, S)$ mit:
 - V ist ein endliches Alphabet von **Variablen** (Nicht-Terminale)
 - Σ ist ein endliches Alphabet von **Terminalen**
 - $S \in V$ ist das **Startsymbol**
 - $P \subseteq ((V \cup \Sigma)^+ \times (V \cup \Sigma)^*)$ ist eine endliche Menge von **Produktionen** (Produktionen werden auch **Regeln** genannt)
- **Ableitbarkeit**
 - $b \in (V \cup \Sigma)^*$ ist **(direkt) ableitbar** aus $a \in (V \cup \Sigma)^+$, $a \rightarrow b$, wenn es eine Regel $(u, v) \in P$ und $\alpha, \beta \in (V \cup \Sigma)^*$ gibt: $a = \alpha u \beta$, $b = \alpha v \beta$
 - b ist aus a **(indirekt) ableitbar**, $a \xrightarrow{*} b$, wenn b durch endlich viele direkte Ableitungsschritte aus a erzeugbar ist
- Die von G erzeugte Sprache ist

$$L(G) = \{w \in \Sigma^* \mid S \xrightarrow{*} w\}$$

Beispiele

- $G_1 = (V, \Sigma, P, S)$ mit
 - $V = \{S, A, B\}$
 - $\Sigma = \{0,1\}$
 - $P = \{S \rightarrow ASB, A \rightarrow 0, B \rightarrow 1, S \rightarrow \varepsilon\}$
- Welche Sprache definiert G_1 , also was ist $L = L(G_1)$?
- $G_2 = (V, \Sigma, P, S)$ mit
 - $V = \{S, R, L\}$
 - $\Sigma = \{0,1\}$
 - $P = \{S \rightarrow 01L0, 1L \rightarrow L11, 0L \rightarrow 0R, R1 \rightarrow 1R, R0 \rightarrow L0, L0 \rightarrow 0\}$
- Was ist $L = L(G_2)$?

Chomsky-Hierarchie

- Eine Grammatik $G = (V, \Sigma, P, S)$ heißt vom Typ Chomsky-0.
- G heißt **kontextsensitiv** (vom Typ Chomsky-1), wenn für jede Regel $a \rightarrow b$ gilt: $|a| \leq |b|$ (Ausnahme: $A \rightarrow \varepsilon$ wird zugelassen)
- G heißt **kontextfrei** (vom Typ Chomsky-2), wenn alle Regeln die Form $A \rightarrow b$ mit $A \in V$ haben
- G heißt **regulär** (vom Typ Chomsky-3), wenn alle Regeln die Form $A \rightarrow \varepsilon$ oder $A \rightarrow aB$ mit $a \in \Sigma$ und $A, B \in V$ haben
- Eine Sprache heißt kontextsensitiv, kontextfrei oder regulär, wenn sie durch eine entspr. Grammatik erzeugt werden kann
- Sei L_i die Menge der Typ Chomsky- i Sprachen, dann gilt:

$$L_3 \subseteq L_2 \subseteq L_1 \subseteq L_0$$

(Anmerkung: Chomsky-Hierarchie ist sogar echt)

- Welche besondere Struktur haben **reguläre Sprachen**?

Beispiel einer regulären Sprache

- Reguläre Grammatik über $\Sigma = \{0,1\}$ für

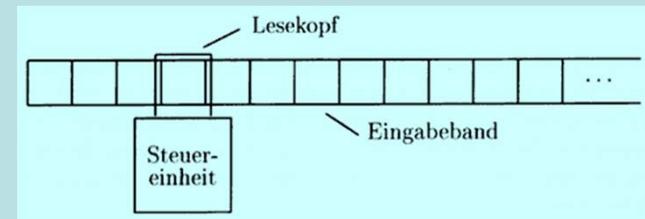
$$L_3 = \{ w \mid w \text{ enthält eine gerade Anzahl Einsen} \}$$

- $G_3 = (V, \Sigma, P, S)$ mit
 - $V = \{S, A\}$
 - $\Sigma = \{0,1\}$
 - $P = \{S \rightarrow 0S, S \rightarrow 1A, A \rightarrow 0A, A \rightarrow 1S, S \rightarrow \varepsilon\}$

- Reguläre Sprachen können auch durch reguläre Ausdrücke beschrieben werden
- Reguläre Ausdrücke über Σ sind wie folgt definiert:
 - $\emptyset = \{\}$ und $\{\varepsilon\}$ sind regulär
 - $\{a\}$ ist regulär für jedes $a \in \Sigma$
 - Sind L_1, L_2 regulär, dann auch $L_1 + L_2, L_1 \cdot L_2, L_1^*$ mit
 - $L_1 + L_2$ Vereinigung von L_1 und L_2
 - $L_1 \cdot L_2$ Verkettung von L_1 und L_2 ($= \{w_1w_2 \mid w_1 \in L_1, w_2 \in L_2\}$)
 - L_1^* **Kleene'sche Abschluss:** $\bigcup_{i \geq 0} L_1^i$ (...was ist L_1^+ ?)
 - Jeder reguläre Ausdruck kann durch Anwendung der oben angegebenen Regeln erzeugt werden
- Prioritäten: 1) * , 2) \cdot , 3) $+$
- Beispiel:
 - $L_{\text{DezKonst}} = (1+2+3+4+5+6+7+8+9) \cdot (0+1+2+3+4+5+6+7+8+9)^* + 0$
- Anwendung in Werkzeugen (Editoren) und Sprachen (perl, php, ...)

Deterministischer endlicher Automat

- Ein deterministischer endlicher Automat (DEA/DFA) ist beschrieben durch ein 5-Tupel $M = (Q, \Sigma, \delta, q_0, F)$ mit:
 - Q ist eine endliche, nichtleere Menge von **Zuständen**
 - Σ ist ein endliches, nicht leeres **Eingabealphabet**
 - $\delta: Q \times \Sigma \rightarrow Q$ ist die **Übergangsfunktion**
 - $q_0 \in Q$ ist der **Startzustand**
 - $F \subseteq Q$ ist die Menge der **akzeptierenden Endzustände**



- **Akzeptanz**

- Ein DFA **akzeptiert** eine Eingabe x_1, \dots, x_n , wenn gilt:

$$q_0 x_1, \dots, x_n \xrightarrow{*} x_1, \dots, x_n q \text{ mit } q \in F$$

- Alternativ durch induktive Fortsetzung von δ auf $Q \times \Sigma^* \rightarrow Q$:

$$\delta(q, \varepsilon) = q \text{ für } n > 0, \text{ sonst } \delta(q, x_1, \dots, x_n) = \delta(\delta(q, x_1, \dots, x_{n-1}), x_n)$$

- Damit: Ein DFA akzeptiert $x = x_1, \dots, x_n \Leftrightarrow \delta(q_0, x) \in F$

Beispiel (Automat verstehen)

- $M = (Q, \Sigma, \delta, q_0, F)$ mit

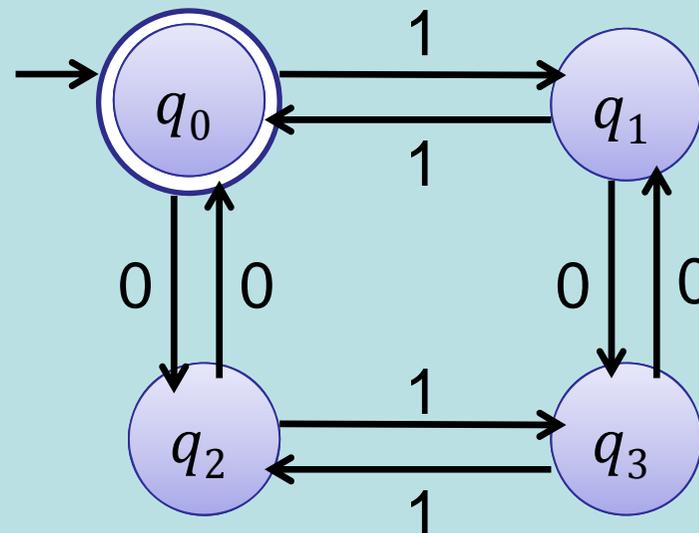
- $Q = \{q_0, \dots, q_3\}$

- $\Sigma = \{0,1\}$

- $F = \{q_0\}$

δ	q_0	q_1	q_2	q_3
0	q_2	q_3	q_0	q_1
1	q_1	q_0	q_3	q_2

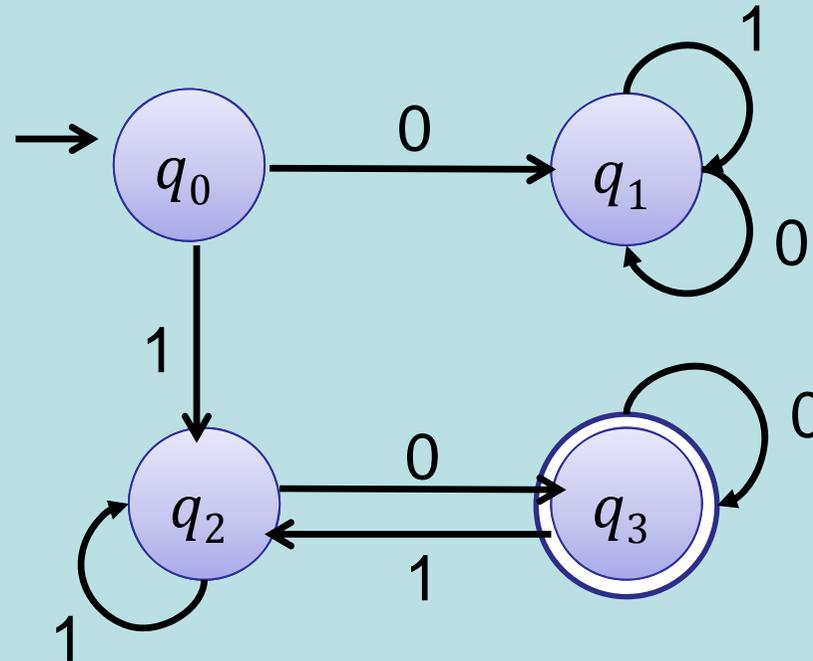
- Grafische Darstellung (Gerichteter Graph mit Knoten und Kanten):



- Welche Sprache akzeptiert der DFA M ?

Beispiel (Automat entwerfen)

$L = \{ w \mid w \text{ beginnt mit } 1 \text{ und endet mit } 0 \}$



δ	q_0	q_1	q_2	q_3
0	q_1	q_1	q_3	q_3
1	q_2	q_1	q_2	q_2