

Reguläre Sprachen, Ausdrucksstärke (Teil 1)

BC George (HSBI)

Unless otherwise noted, this work is licensed under CC BY-SA 4.0.

Motivation

Was muss ein Compiler wohl als erstes tun?

- Einlesen
- Bsp. von Quellcode
 - Wörter
 - Zeichen

Themen für heute

- Lexer Scanner Zerhöler
- Endliche Automaten
- Reguläre Sprachen

Endliche Automaten

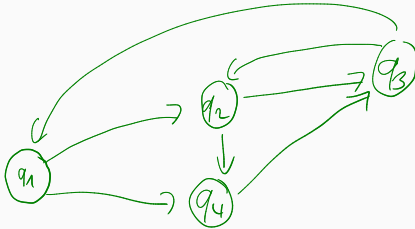
Def.: Ein **Alphabet** Σ ist eine endliche, nicht-leere Menge von Symbolen. Die Symbole eines Alphabets heißen *Buchstaben*.

Def.: Ein **Wort** w über einem Alphabet Σ ist eine endliche Folge von Symbolen aus Σ . ϵ ist das leere Wort. Die *Länge* $|w|$ eines Wortes w ist die Anzahl von Buchstaben, die es enthält (Kardinalität).

Def.: Eine **Sprache** L über einem Alphabet Σ ist eine Menge von Wörtern über diesem Alphabet. Sprachen können endlich oder unendlich viele Wörter enthalten.

epsilon

State machine

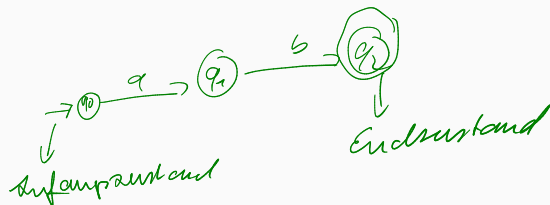


Deterministische endliche Automaten

Bestimmte State machines:

- Eingaben bestimmen Zustandsübergänge
- Zustandsübergänge sind eindeutig
- Es gibt Anfang(szustand) und End(zuständ)e

genau einen



Automat läuft einmal durch

$$Q = \{q_0, q_1, q_2\}$$

Wie definieren wir das formal?

- Menge der Zustände
- Eingabealphabet
- Anfangszustand
- Endzustände
- Übergangsfkt

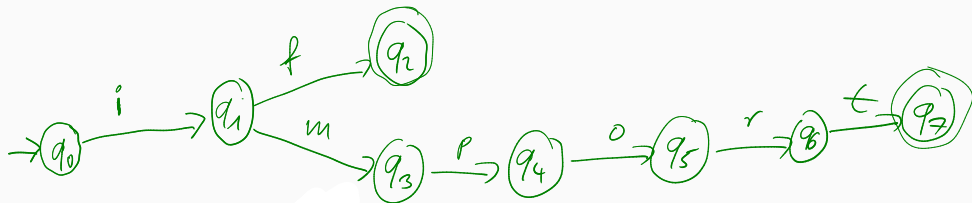
Def.: Deterministischer endlicher Automat

Def.: Ein **deterministischer endlicher Automat** (DFA) ist ein 5-Tupel $A = (Q, \Sigma, \delta, q_0, F)$ mit

- Q : endliche Menge von **Zuständen**
- Σ : Alphabet von **Eingabesymbolen**
- δ : die (eventuell partielle) **Übergangsfunktion** $(Q \times \Sigma) \rightarrow Q$, δ kann partiell sein
- $q_0 \in Q$: der **Startzustand**
- $F \subseteq Q$: die Menge der **Endzustände**

großes
Sigma
/
delta

$$\begin{aligned} \mathbb{N} &\rightarrow \mathbb{R} \\ x &\mapsto x+1 \end{aligned}$$



ifm \rightarrow nicht akzeptiert

impo \rightarrow "

$$L^*(q_0, if) = L(\underbrace{L^*(q_0, i)}_{L(q_0, q)}, f)$$

q_0
 $L(q_0, i)$

Eingabewörter statt Buchstaben

0 mal oder endlich
oft angewendet

Def.: Wir definieren $\delta^* : (Q \times \Sigma^*) \rightarrow Q$: induktiv wie folgt:

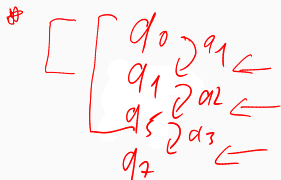
- Basis: $\delta^*(q, \epsilon) = q \quad \forall q \in Q$

$q_i \in \Sigma \quad w = a_1 \dots a_n$

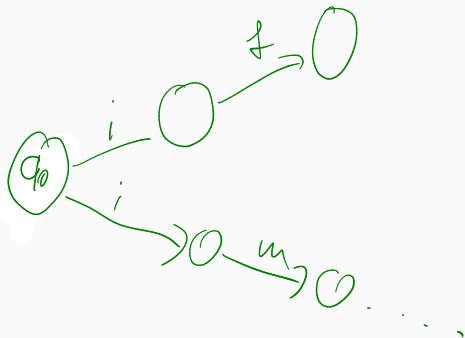
- Induktion: $\delta^*(q, a_1, \dots, a_n) = \delta(\delta^*(q, a_1, \dots, a_{n-1}), a_n)$

Zustand nach der Verarbeitung von $a_1 \dots a_{n-1}$

Def.: Ein DFA akzeptiert ein Wort $w \in \Sigma^*$ genau dann, wenn $\delta^*(q_0, w) \in F$.



s. Folie 9



Def.: Nichtdeterministischer Automat

Def.: Ein **nichtdeterministischer endlicher Automat** (NFA) ist ein 5-Tupel $A = (Q, \Sigma, \delta, q_0, F)$ mit

- Q : endliche Menge von **Zuständen**
- Σ : Alphabet von **Eingabesymbolen**
- δ : die (eventuell partielle) **Übergangsfunktion** $(Q \times \Sigma) \rightarrow \cancel{Q}$
- $q_0 \in Q$: der **Startzustand**
- $F \subseteq Q$: die Menge der **Endzustände**

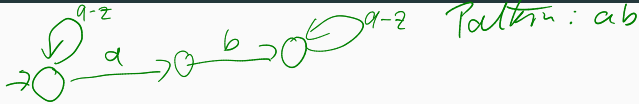
*→ Menge aller
Teilungen
von Q*
 $\mathcal{P}(Q)$

Def.: Sei A ein DFA oder ein NFA. Dann ist $L(A)$ die von A akzeptierte Sprache, d. h.

$$L(A) = \{\text{Wörter } w \mid \delta^*(q_0, w) \in F\}$$

$$w = a_1 \dots a_n \quad a_i \in \Sigma$$

Wozu NFAs im Compilerbau?



Pattern Matching (Erkennung von Schlüsselwörtern, Bezeichnen, ...) geht mit NFAs.

NFAs sind so nicht zu programmieren, aber:

Satz: Eine Sprache L wird von einem NFA akzeptiert $\Leftrightarrow L$ wird von einem DFA akzeptiert.

D. h. es existieren Algorithmen zur

- Umwandlung von NFAs in DFAs
- Minimierung von DFAs

Reguläre Sprachen

Reguläre Ausdrücke definieren Sprachen

Def.: Induktive Definition von **regulären Ausdrücken** (regex) und der von ihnen repräsentierten Sprache L :

- Basis:

- ϵ und \emptyset sind reguläre Ausdrücke mit $L(\epsilon) = \{\epsilon\}$, $L(\emptyset) = \emptyset$ — keine Sprache
- Sei a ein Symbol $\Rightarrow a$ ist ein regex mit $L(a) = \{a\}$

- Induktion: Seien E, F reguläre Ausdrücke. Dann gilt:

- $E + F$ ist ein regex und bezeichnet die Vereinigung $L(E + F) = L(E) \cup L(F)$ = | oder ,
 - EF ist ein regex und bezeichnet die Konkatenation $L(EF) = L(E)L(F)$
 - E^* ist ein regex und bezeichnet die Kleene-Hülle $L(E^*) = (L(E))^*$ oder endlich oft
 - (E) ist ein regex mit $L((E)) = L(E)$
- $(ab)^* = \{\epsilon, ab, abab, ababab, \dots\}$

Vorrangregeln der Operatoren für reguläre Ausdrücke: *, Konkatenation, +

(E^+ : endlich oft)

Beispiel

$1 (0-9)^* 2$ beschreibt alle Wörter, die mit 1 beginnen
und mit 2 enden, dazwischen stehen 0-9,
0 oder endlich viele

0^*

0 oder endlich viele Nullen

Wichtige Identitäten

Satz: Sei A ein DFA $\Rightarrow \exists$ regex R mit $L(A) = L(R)$.

Satz: Sei E ein regex $\Rightarrow \exists$ DFA A mit $L(E) = L(A)$.

$a \mid b \in 10$ if $a \mid b \in 10$

Trennzeichen: Delimiter

Automaten akzeptieren Sprachen,
regex beschreiben Sprachen

$r_1 = \text{if}$
 $r_2 = \text{url}$
 $r_3 = \text{import}$

$r_1 = \{$

Formale Grammatiken

Syntax einer Sprache

Start=
Symbol

Satz

gültig oder
noch

$S \rightarrow P \ O$
Artikel Substantiv

$S \rightarrow 'The' \mid 'A' \mid 'An'$
Artikel

↑
oder

↑
oder

$S \rightarrow 'street' \mid 'book' \mid 'house' \mid 'door' \dots$
Substantiv

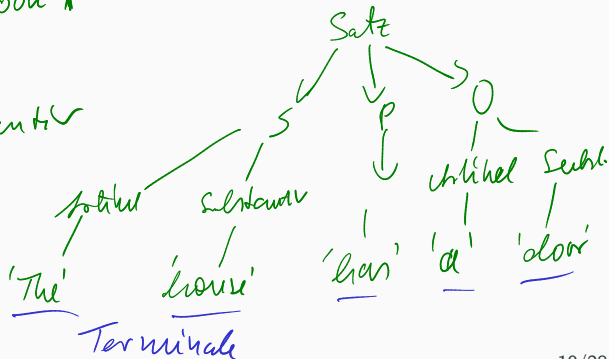
$P \rightarrow 'has' \mid 'is'$
P

$O \rightarrow \text{Artikel Substantiv}$
O

Nichtterminale

Ableitungsbaum

Produktionen



Def.: Eine *formale Grammatik* ist ein 4-Tupel $G = (N, T, P, S)$ aus

- N : endliche Menge von **Nichtterminalen**
- T : endliche Menge von **Terminalen**, $N \cap T = \emptyset$
- $S \in N$: **Startsymbol**
- P : endliche Menge von **Produktionen** der Form

$$\underline{X \rightarrow Y} \text{ mit } X \in (N \cup T)^* \underline{N} (N \cup T)^*, Y \in (N \cup T)^*$$

Def.: Sei $G = (N, T, P, S)$ eine Grammatik, sei $\alpha A \beta$ eine Zeichenkette über $(N \cup T)^*$ und sei $A \rightarrow \gamma$ eine Produktion von G .

Wir schreiben: $\alpha A \beta \Rightarrow \alpha \gamma \beta$ ($\alpha A \beta$ leitet $\alpha \gamma \beta$ ab).

*eine Produktion anwenden
Nichtterminale:
Großbuchstaben*

Def.: Wir definieren die Relation \Rightarrow^* induktiv wie folgt:

- Basis: $\forall \alpha \in (N \cup T)^* \alpha \Rightarrow^* \alpha$ (Jede Zeichenkette leitet sich selbst ab.)
- Induktion: Wenn $\alpha \Rightarrow^* \beta$ und $\beta \Rightarrow \gamma$ dann $\alpha \Rightarrow^* \gamma$

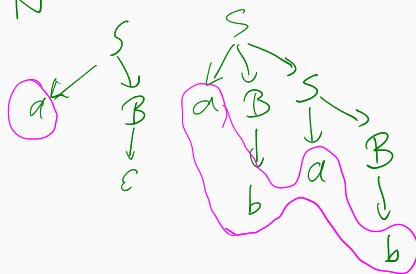
Def.: Sei $G = (N, T, P, S)$ eine formale Grammatik. Dann ist $L(G) = \{\text{Wörter } w \text{ über } \underline{T} \mid S \Rightarrow^* w\}$ die von G erzeugte Sprache.

Grammatiken erzeugen Sprachen

Beispiel

$$\begin{aligned} P \{ & S \rightarrow aBS \mid aB \\ & B \rightarrow aB \mid b \mid \varepsilon \end{aligned}$$

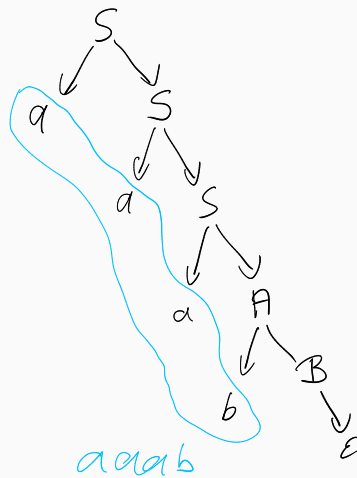
$G = (\underbrace{\{S, B\}}_N, \underbrace{\{a, b\}}_T, P, S)$



Def.: Eine **reguläre (oder type-3-) Grammatik** ist eine formale Grammatik mit den folgenden Einschränkungen:

- Alle Produktionen sind entweder von der Form
 - $X \rightarrow aY$ mit $X \in N, a \in T, Y \in N$ (*rechtsreguläre Grammatik*) oder *oder*
 - $X \rightarrow Ya$ mit $X \in N, a \in T, Y \in N$ (*linksreguläre Grammatik*)
- $X \rightarrow \epsilon$ ist erlaubt

$$\begin{aligned} S &\rightarrow aS \mid aA \\ A &\rightarrow bB \\ B &\rightarrow \epsilon \end{aligned}$$



Satz: Die von endlichen Automaten akzeptierte Sprachklasse, die von regulären Ausdrücken beschriebene Sprachklasse und die von regulären Grammatiken erzeugte Sprachklasse sind identisch und heißen **reguläre Sprachen**.

Wrap-Up

- Definition und Aufgaben von Lexern
- DFAs und NFAs
- Reguläre Ausdrücke
- Reguläre Grammatiken
- Zusammenhänge zwischen diesen Mechanismen und Lexern, bzw. Lexergeneratoren



Unless otherwise noted, this work is licensed under CC BY-SA 4.0.

Last modified: 64341f7 (lecture: split regular expressions into two lessons, 2025-08-29)