

.

# Reguläre Sprachen, Ausdrucksstärke

---

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?

- Syntax überprüfen
- "Wörter" überprüfen
- Teile (Leerzeichen) erkennen

# Themen für heute

- Endliche Automaten
- Reguläre Sprachen
- Lexer

# Endliche Automaten

---

# Alphabete

*Σ Sigma*

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

*ε epsilon*

**Def.:** Ein **Wort**  $w$  über einem Alphabet  $\Sigma$  ist eine endliche Folge von Symbolen aus  $\Sigma$ .  $\epsilon$  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  $\Sigma$  ist eine Menge von Wörtern über diesem Alphabet. Sprachen können endlich oder unendlich viele Wörter enthalten.

## Beispiel

Elemente kann man hintereinander hängen

$\Rightarrow$  Sei  $M$  eine Menge

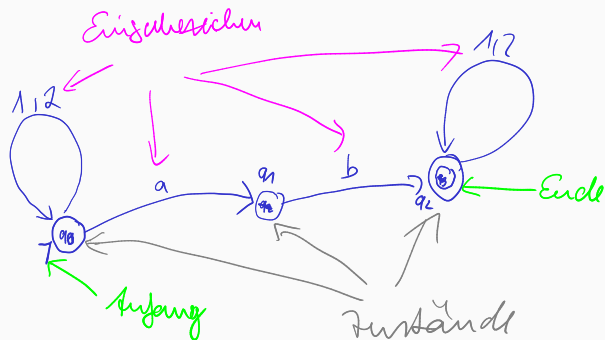
$M^*$  die Menge aller hintereinander gehängten Elemente,  
(endlich viele)  
 $\downarrow$  einschließlich  $\varepsilon$

Menge'sche Hülle

$M^+ : M^* \text{ ohne } \{\varepsilon\}$

# Deterministische endliche Automaten

DFA



state machine

Wenn die Eingabe komplett verarbeitet wurde  
und der Automat in einem Endzustand ist  
wurde das Wort akzeptiert



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

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

# Die Übergangsfunktion

  $\stackrel{1}{=}$   $\delta$  : ein Übergang

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

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

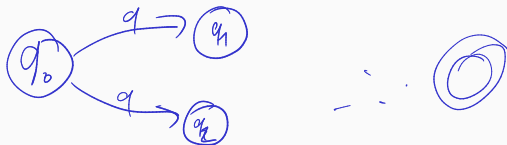
**Def.:** Ein DFA akzeptiert ein Wort  $w \in \Sigma^*$  genau dann, wenn  $\delta^*(q_0, w) \in F$ .

# Beispiel



	a	b
→ q <sub>0</sub>	q <sub>1</sub>	q <sub>1</sub>
q <sub>1</sub>	q <sub>0</sub>	—

# Nichtdeterministische endliche Automaten



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

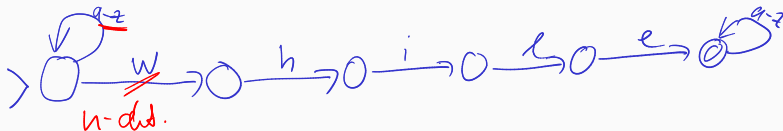
- $Q$  : endliche Menge von **Zuständen**
- $\Sigma$  : Alphabet von **Eingabesymbolen**
- $\delta$  : die (eventuell partielle) Übergangsfunktion  $(Q \times \Sigma) \rightarrow \emptyset \mathcal{P}(Q)$
- $q_0 \in Q$  : der **Startzustand**
- $F \subseteq Q$  : die Menge der **Endzustände**

$\rightarrow$  Potenzmenge  
 $\downarrow$   
Menge aller  
Teilmenge

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

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

# Wozu NFAs im Compilerbau?



Pattern Matching (Erkennung von Schlüsselwörtern, Bezeichnern, ...) 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 (Anzahl d. Zustände)

# Reguläre Sprachen

---

# Reguläre Ausdrücke definieren Sprachen

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

*beschriebenen Sprachen*

- Basis:

- $\epsilon$  und  $\emptyset$  sind reguläre Ausdrücke mit  $L(\epsilon) = \{\epsilon\}$ ,  $L(\emptyset) = \emptyset$

- Sei  $a$  ein Symbol  $\Rightarrow a$  ist ein regex mit  $L(a) = \{a\}$   
 *$\epsilon \Sigma$*

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

- $E + F$  ist ein regex und bezeichnet die Vereinigung  $L(E + F) = L(E) \cup L(F)$

*auch !  
Hintereinanderschließen*

- $EF$  ist ein regex und bezeichnet die Konkatenation  $L(EF) = L(E)L(F)$

- $E^*$  ist ein regex und bezeichnet die Kleene-Hülle  $L(E^*) = (L(E))^*$

- $(E)$  ist ein regex mit  $L((E)) = L(E)$

*[ ]*

*[a,b]*

*a oder b*

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



## Beispiel

regulär	Sprache
01	$\{01\}$
$(a+b)$	$\{a, b\}$
$(a+b)^*$	$\{\epsilon, a, b, aa, ab, ba, bb, aaaa, aab, \dots\}$ unendlich viele jedes endlich lang

# 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  $L(E) = L(A)$ .

} Algorithmen zur  
Umwandlung

# Formale Grammatiken

Satz  $\rightarrow$  Subjekt Prädikat Objekt

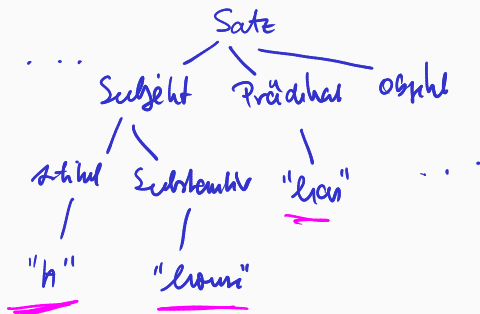
Subjekt  $\rightarrow$  Artikel Substantiv | Pronomen

Prädikat  $\rightarrow$  "is" | "has" | "know" | ...

Objekt  $\rightarrow$  Artikel Substantiv | ...

Artikel  $\rightarrow$  "A" | "An" | "The"

Substantiv  $\rightarrow$  "house" | "roof" | ...



# Formale Definition formaler Grammatiken

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

- $N$ : endliche Menge von **Nichtterminalen** *Subjekt, Prädikat,*
- $T$ : endliche Menge von **Terminalen**,  $N \cap T = \emptyset$  *"has", "is" usw.*
- $S \in N$ : **Startsymbol** *Satz*
- $P$ : endliche Menge von **Produktionen** der Form

$X \rightarrow Y$  mit  $X \in (N \cup T)^* N (N \cup T)^*$ ,  $Y \in (N \cup T)^*$



# Ableitungen

"A" Substantiv Präzision abgeleitet

**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  $G$ .

gamma

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

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

- Basis:  $\forall \alpha \in (N \cup T)^* \alpha \xRightarrow{*} \alpha$  (Jede Zeichenkette leitet sich selbst ab.)
- Induktion: Wenn  $\alpha \xRightarrow{*} \beta$  und  $\beta \xRightarrow{*} \gamma$  dann  $\alpha \xRightarrow{*} \gamma$   
→ endlich viele Schritte  
↳ ein Schritt

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

S.O.

**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 *exor*
  - $X \rightarrow Ya$  mit  $X \in N, a \in T, Y \in N$  (*linksreguläre Grammatik*)
- $X \rightarrow \epsilon$  ist erlaubt

Satz  $\rightarrow$  "A" A | "The" A | "An" A

A  $\rightarrow$  "house" B | "royal" B | ...

B  $\rightarrow$  "is" C | "has" C | ...

C  $\rightarrow$  ...



**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**.

## Reguläre Sprachen

- einfache Struktur
- Matchen von Symbolen (z. B. Klammern) nicht möglich, da die fixe Anzahl von Zuständen eines DFAs die Erkennung solcher Sprachen verhindert.



# Wozu reguläre Sprachen im Compilerbau?

- Reguläre Ausdrücke

*if while ; : + ...*

- definieren Schlüsselwörter und alle weiteren Symbole einer Programmiersprache, z. B. den Aufbau von Gleitkommazahlen
- werden (oft von einem Generator) in DFAs umgewandelt
- sind die Basis des *Scanners* oder *Lexers* *Zerfäher*

# Lexer

---

# Ein Lexer ist mehr als ein DFA

- Ein **Lexer**

- wandelt mittels DFAs aus regulären Ausdrücken die Folge von Zeichen der Quelldatei in eine Folge von sog. Token um
- bekommt als Input eine Liste von Paaren aus regulären Ausdrücken und Tokennamen, z. B. ("while", WHILE)
- Kommentare und Strings müssen richtig erkannt werden. (Schachtelungen)
- liefert Paare von Token und deren Werte, sofern benötigt, z. B. (WHILE, \_), oder (IDENTIFIER, "radius") oder (INTEGERZAHL, "334")

*Besuchen*

# Wie geht es weiter?

## ■ Ein Parser

- führt mit Hilfe des Tokenstreams vom Lexer die Syntaxanalyse durch
- basiert auf einer sog. kontextfreien Grammatik, deren Terminale die Token sind
- liefert die syntaktische Struktur in Form eines Ableitungsbaums (**syntax tree**, **parse tree**), bzw. einen **AST** (abstract syntax tree) ohne redundante Informationen im Ableitungsbaum (z. B. Semikolons)
- liefert evtl. Fehlermeldungen

*/ des Lexers*

*(Der Lexer liefert auch Fehlermeldungen)*

## 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.