

# Theoretische Informatik 1

ZUSAMMENFASSUNG ZUM MODUL FORMALE SPRACHEN UND AUTOMATENTHEORIE

---

---

Simon KÖNIG

# INHALTSVERZEICHNIS

<b>I</b>	<b>Formale Sprachen und Automatentheorie</b>	<b>4</b>
<b>1</b>	<b>Formale Sprachen – Einstieg</b>	<b>5</b>
<b>2</b>	<b>Alle Sprachen</b>	<b>6</b>
<b>3</b>	<b>Rekursiv aufzählbare Sprachen (r.e.): Typ-0</b>	<b>7</b>
<b>4</b>	<b>Entscheidbare Sprachen (REC):</b>	<b>8</b>
<b>5</b>	<b>Kontextsensitive Sprachen (CSL): Typ-1</b>	<b>9</b>
5.1	Algorithmus zur Entscheidbarkeit des Wortproblems: . . . . .	9
5.2	Automatenmodell LBA: . . . . .	9
5.3	Beispiele: . . . . .	9
<b>6</b>	<b>Kontextfreie Sprachen (CFL): Typ-2</b>	<b>10</b>
6.1	Automatenmodell PDA: . . . . .	10
6.2	Sätze zu den kontextfreien Sprachen: . . . . .	10
6.2.1	Pumping-Lemma für Typ-2: . . . . .	10
6.3	Chomsky-Normalform . . . . .	10
6.3.1	Umformungsalgorithmus: . . . . .	10
6.4	Greibach-Normalform . . . . .	11
6.4.1	Umformungsalgorithmus: . . . . .	11
6.5	Beispiele: . . . . .	12
<b>7</b>	<b>Deterministisch Kontextfreie Sprachen (DCFL):</b>	<b>13</b>
7.1	Sätze zu den deterministisch kontextfreien Sprachen: . . . . .	13
7.2	Beispiele: . . . . .	13
<b>8</b>	<b>Reguläre Sprachen (REG): Typ-3</b>	<b>14</b>
8.1	Automatenmodell DEA: . . . . .	14
8.2	Automatenmodell NEA: . . . . .	14
8.3	Reguläre Ausdrücke . . . . .	15
8.4	Sätze zu den regulären Sprachen: . . . . .	15
8.4.1	Pumping-Lemma für Typ-3 . . . . .	15
8.4.2	Myhill-Nerode-Äquivalenz . . . . .	16
8.4.3	Erkennung durch Monoide – Syntaktisches Monoid . . . . .	16
8.5	Beispiele: . . . . .	17
<b>9</b>	<b>Endliche Sprachen (FIN):</b>	<b>18</b>
9.1	Beispiele: . . . . .	18
<b>10</b>	<b>Algorithmen</b>	<b>19</b>
10.1	Konstruktionsalgorithmus für Minimal-DEAs: . . . . .	19

<b>11</b>	<b>Anwendungen der Sätze</b>	<b>21</b>
11.1	Pumping-Lemma für Typ-3 . . . . .	21
11.2	Myhill-Nerode-Äquivalenz . . . . .	21

# 1

## Formale Sprachen und Automatentheorie

# 1: FORMALE SPRACHEN – EINSTIEG

## Definition 1.1: Alphabet

Als Alphabet bezeichnen wir eine endliche, nichtleere Menge, deren Elemente Buchstaben genannt werden.

Dieses wird üblich mit  $\Sigma$  bezeichnet.

## Definition 1.2: Freies Monoid über $\Sigma$

Ein Monoid ist eine Menge mit einer assoziativen Verknüpfung und einem neutralen Element. Die Menge aller endlichen Zeichenketten, die sich aus Elementen von  $\Sigma$  bilden lassen bilden mit der *Konkatenation* ein Monoid  $\Sigma^*$ .

Das leere Wort ( $\epsilon$ ) bildet das neutrale Element.

## Definition 1.3: Grammatik

Eine Grammatik ist ein Quadrupel:

$$G = (V, \Sigma, P, S)$$

$V$  Eine Menge an Zeichen, den Variablen

$\Sigma$  Das Alphabet,  $V \cap \Sigma = \emptyset$

$P$  Die Produktionsmenge,  $P \subseteq (V \cup \Sigma)^+ \times (V \cup \Sigma)^*$

$S$  Das Startsymbol oder die Startvariable,  $S \in V$

## 2: ALLE SPRACHEN

### 3: REKURSIV AUFZÄHLBARE SPRACHEN (r.e.): TYP-0

Nicht entscheidbar in allen Kategorien, Abschluss unter Stern, Vereinigung, Schnitt und Konkatenation.

## 4: ENTSCHEIDBARE SPRACHEN (REC):

$$(REC) \subset (r.e.)$$



## 5: KONTEXTSENSITIVE SPRACHEN (CSL): TYP-1

$$(\text{CSL}) \subset (\text{REC}) \subset (\text{r.e.})$$

Wortproblem entscheidbar. Abschluss unter allen Operationen.

### 5.1 Algorithmus zur Entscheidbarkeit des Wortproblems:

Die Funktion  $\text{Abl}_n(X)$  wird iteriert angewendet, bis sich entweder  $X$  nicht mehr ändert ( $w \notin L(G)$ ) oder das gesuchte Wort  $w$  in  $X$  enthalten ist ( $w \in L(G)$ ).

Dabei ist  $n$  die Länge des gesuchten Worts  $w$ , also  $|w|$ .

Die Funktion  $\text{Abl}_n(X)$  ist für eine Grammatik  $G$  wie folgt definiert:

$$\text{Abl}_n(X) := X \cup \{w \in (V \cup \Sigma)^* \mid |w| \leq n \wedge \exists y \in X : y \Rightarrow_G w\}$$

### 5.2 Automatenmodell LBA:

Die Kontextsensitiven Sprachen werden von linear beschränkten Automaten akzeptiert.

### 5.3 Beispiele:

- $L_1 = \{a^n b^n c^n \mid n \geq 1\}$  (drei und mehr gleiche Exponenten)
-

## 6: KONTEXTFREIE SPRACHEN (CFL): TYP-2

$$(CFL) \subset (CSL) \subset (REC) \subset (r.e.)$$

Wort- und Leerheitsproblem entscheidbar.

Beschreib- und erkennbar durch einen nichtdeterministischen Kellerautomaten (PDA).

### 6.1 Automatenmodell PDA:

### 6.2 Sätze zu den kontextfreien Sprachen:

- Jede kontextfreie Sprache über einem unären Alphabet ist regulär!
- Die Klasse der kontextfreien Sprachen ist abgeschlossen unter Sternoperation, Vereinigung und Konkatenation.
- Das Wortproblem ( $\mathcal{O}(n^3)$ ) sowie das Leerheitsproblem sind entscheidbar.

#### 6.2.1 Pumping-Lemma für Typ-2:

Sei  $L \subseteq \Sigma^*$  eine kontextfreie Sprache, dann gibt es eine Zahl  $n$  so, dass für alle  $z \in L$  mit  $|z| \geq n$  eine Zerlegung mit  $z = uvwxy$  in  $u, v, w, x, y \in \Sigma^*$  existiert für die die drei Bedingungen erfüllt sind:

- $|vx| \geq 1$
- $|vwx| \leq n$
- $\forall i \in \mathbb{N} : uv^iwx^iy \in L$

### 6.3 Chomsky-Normalform

Eine Typ-2 Grammatik  $(V, \Sigma, P, S)$  ist in Chomsky-Normalform (CNF), wenn gilt:

$$\forall (u, v) \in P : v \in V^2 \cup \Sigma$$

Zu jeder Typ-2 Grammatik existiert eine Grammatik  $G'$  in CNF für die gilt  $L(G) = L(G')$ !

#### 6.3.1 Umformungsalgorithmus:

1. Zunächst wollen wir erreichen, dass folgendes gilt:  $(u, v) \in P \Rightarrow (|v| > 1 \vee v \in \Sigma)$

(a) **Ringableitungen entfernen:**

Eine Ringableitung liegt vor, wenn es Variablen  $A_1, \dots, A_r$  gibt, die sich im Kreis in einander ableiten lassen, d.h. es gibt Regeln  $A_i \rightarrow A_{i+1}$  und  $A_r \rightarrow A_1$ .

Um dies loszuwerden, werden alle Variablen  $A_i$  durch eine neue Variable  $A$  ersetzt. Überflüssige Regeln wie  $A \rightarrow A$  werden gelöscht.

(b) **Variablen anordnen:**

Man legt eine Ordnung der Variablen fest:  $V = \{B_1, B_2, \dots, B_n\}$ , hierfür muss gelten:

$$A_i \rightarrow A_j \in P \Leftrightarrow i < j$$

Falls dies nicht gilt, müssen Abkürzungen verwendet werden, also alle Produktionen von  $A_j$  werden eingesetzt:

$$P = (P \setminus \{A_i \rightarrow A_j\}) \cup \{(A_i, w) \mid (A_j, w) \in P\}$$

2. Jetzt gilt für jede Regel  $(u, v) \in P$  entweder  $v \in \Sigma$  oder  $|v| \geq 2$ .

Für letztere Regeln werden nun Pseudoterminal eingeführt. Es werden neue Variablen und Produktionen für jedes Terminalsymbol hinzugefügt, z.B.  $V_a \rightarrow a$ .

3. **Letzer Schritt:** Alle rechten Seiten mit  $|v| > 2$  müssen nun noch auf Länge 2 gekürzt werden. Hierfür werden wiederum neue Variablen eingefügt:

$$A \rightarrow C_1 C_2 C_3$$

Wird gekürzt zu

$$\begin{aligned} A &\rightarrow C_1 D_1 \\ D_1 &\rightarrow C_2 C_3 \end{aligned}$$

## 6.4 Greibach-Normalform

Eine Typ-2 Grammatik  $(V, \Sigma, P, S)$  ist in Greibach-Normalform (GNF), wenn gilt:

$$\forall (u, v) \in P : v \in \Sigma V^*$$

Zu jeder Typ-2 Grammatik existiert eine Grammatik  $G'$  in GNF für die gilt  $L(G) = L(G')$ !

### 6.4.1 Umformungsalgorithmus:

1. **Mh?**

2. **Beseitigung von Linksrekursion:**

Alle Produktionsregeln sind von der Form:

$$A \rightarrow A\alpha_1 \mid \dots \mid A\alpha_k \mid \beta_1 \mid \dots \mid \beta_l$$

Diese können durch diese  $2k + 2l$  Regeln ersetzt werden:

$$\begin{aligned} A &\rightarrow \beta_1 \mid \dots \mid \beta_l \\ A &\rightarrow \beta_1 B \mid \dots \mid \beta_l B \\ B &\rightarrow \alpha_1 \mid \dots \mid \alpha_k \\ B &\rightarrow \alpha_1 B \mid \dots \mid \alpha_k B \end{aligned}$$

Nun sind keinerlei Linksrekursionen mehr vorhanden!

## 6.5 Beispiele:

- $L_1 = \{a^n b^n \mid n \geq 1\}$  (zwei gleiche Exponenten)
- $L_2 = \{ww^R \mid w \in \Sigma^*\}$  (unmarkierte Palindrome)
- Korrekt geklammerte arithmetische Ausdrücke (Dyck-Wörter)

## 7: DETERMINISTISCH KONTEXTFREIE SPRACHEN (DCFL):

$$(\text{DCFL}) \subset (\text{CFL}) \subset (\text{CSL}) \subset (\text{REC}) \subset (\text{r.e.})$$

Wort-, Leerheits- und Äquivalenzproblem entscheidbar. Abschluss nur unter Komplement.  
Beschreib- und erkennbar durch einen endlichen deterministischen Kellerautomaten (DPDA).

### 7.1 Sätze zu den deterministisch kontextfreien Sprachen:

- Eine deterministisch kontextfreie Sprache geschnitten mit einer regulären Sprache ist wieder eine deterministisch kontextfreie Sprache.

$$L_{(\text{DCFL})} \cap L_{(\text{REG})} \in (\text{DCFL})$$

### 7.2 Beispiele:

- $L_1 = \{w\$w^R \mid w \in \Sigma^*\}$  (markierte Palindrome)

## 8: REGULÄRE SPRACHEN (REG): TYP-3

$$(\text{REG}) \subset (\text{DCFL}) \subset (\text{CFL}) \subset (\text{CSL}) \subset (\text{REC}) \subset (\text{r.e.})$$

### 8.1 Automatenmodell DEA:

Ein deterministischer endlicher Automat ist ein 5-Tupel

$$M = (Z, \Sigma, \delta, z_0, E)$$

$Z$  endliche Zustandsmenge

$\Sigma$  Eingabealphabet

$\delta$  Überföhrungsfunktion  $\delta : Z \times \Sigma \rightarrow Z$

$z_0$  Startzustand,  $z_0 \in Z$

$E$  Endzustandsmenge,  $E \subseteq Z$

Es lässt sich außerdem eine erweiterte Funktion  $\hat{\delta}$  definieren:

$$\hat{\delta} : Z \times \Sigma^* \rightarrow Z$$

Mit den folgenden Eigenschaften:

$$\hat{\delta}(z, \epsilon) = z$$

$$\hat{\delta}(z, ax) = \hat{\delta}(\delta(z, a), x)$$

Die von einem deterministischen Automaten  $M$  akzeptierte Sprache ist

$$T(M) = \left\{ w \in \Sigma^* \mid \hat{\delta}(z_0, w) \in E \right\}$$

### 8.2 Automatenmodell NEA:

Ein nichtdeterministischer endlicher Automat ist ein 5-Tupel

$$M = (Z, \Sigma, \delta, S, E)$$

$Z$  endliche Zustandsmenge

$\Sigma$  Eingabealphabet

$\delta$  Überföhrungsfunktion  $\delta : Z \times \Sigma \rightarrow \text{Pot}(Z)$

$S$  Startzustandsmenge,  $S \subseteq Z$

$E$  Endzustandsmenge,  $E \subseteq Z$

Der NEA ist formal stärker als der DEA, sie akzeptieren jedoch beide die selbe Sprachklasse. Die akzeptierte Sprache eines nichtdeterministischen endlichen Automaten ist:

$$T(M) = \left\{ w \in \Sigma^* \mid \hat{\delta}(S, w) \cap E \neq \emptyset \right\}$$

## 8.3 Reguläre Ausdrücke

Die regulären Sprachen lassen sich zusätzlich zu den zwei Automatenmodellen auch durch sog. reguläre Ausdrücke beschreiben. Eine Definition für die Syntax der regulären Ausdrücke ist:

- $\emptyset$  und  $\epsilon$  sind reguläre Ausdrücke.
- $a$  ist ein regulärer Ausdruck für alle  $a \in \Sigma$
- Wenn  $\alpha$  und  $\beta$  reguläre Ausdrücke sind, dann sind auch  $\alpha\beta$ ,  $(\alpha|\beta)$  und  $(\alpha)^*$  reguläre Ausdrücke.

Die Semantik der regulären Ausdrücke ist ebenso induktiv bestimmt:

- $L(\emptyset) = \emptyset$  und  $L(\epsilon) = \{\epsilon\}$
- $L(a) = \{a\}$  für jedes  $a \in \Sigma$
- $L(\alpha\beta) = L(\alpha)L(\beta)$ ,  $L(\alpha|\beta) = L(\alpha) \cup L(\beta)$ ,  $L((\alpha)^*) = L(\alpha)^*$

## 8.4 Sätze zu den regulären Sprachen:

- Typ-3 Sprachen können *nicht* inhärent mehrdeutig sein, da sich zu jeder Sprache ein Minimalautomat bilden lässt.
- Die Klasse der Typ-3 Sprachen ist unter allen boole'schen Operationen, Sternoperation und der Konkatenation abgeschlossen.
- Für reguläre Sprachen ist das Wortproblem (in Linearzeit), das Leerheitsproblem, das Äquivalenzproblem sowie das Schnittproblem entscheidbar.

### 8.4.1 Pumping-Lemma für Typ-3

Für jede reguläre Sprache  $L$  gibt es ein  $n \in \mathbb{N}$ , so dass für jedes  $x \in L$  mit  $|x| \geq n$  eine Zerlegung in drei Teile existiert:  $x = uvw$ , so dass die drei Bedingungen erfüllt sind:

- $|v| \geq 1$
- $|uv| \leq n$
- $\forall i \in \mathbb{N} : uv^i w \in L$  („Pump-Bedingung“)

Gilt die Negation dieser Aussage, also

$$\forall n \in \mathbb{N} : \exists x \in L, |x| \geq n : \forall u, v, w \in \Sigma^*, x = uvw, |v| \geq 1, |uv| \leq n : \exists i \in \mathbb{N} : uv^i w \notin L$$

so ist  $L$  nicht regulär!

**ABER:** Das Pumping-Lemma gibt keine Charakterisierung der Typ-3 Sprachen an! Es gibt auch Sprachen, die nicht Typ-3 sind, die Aussage des Lemmas aber trotzdem erfüllen! Das Pumping-Lemma gibt also nur eine Möglichkeit, zu Beweisen, dass eine Sprache *nicht* regulär ist! (Siehe Abschnitt 11.1)

### 8.4.2 Myhill-Nerode-Äquivalenz

Mit der Myhill-Nerode-Äquivalenz ist es möglich nachzuweisen, ob eine Sprache regulär ist.

$$xR_L y \iff [\forall w \in \Sigma^* : xw \in L \iff yw \in L]$$

Bzw. anhand eines DEA (dies führt zu einer Verfeinerung von  $R_L$ )

$$xR_M y \iff [\delta(z_0, x) = \delta(z_0, y)]$$

Es gilt:

$$xR_M y \Rightarrow \forall w \in \Sigma^* : \delta(z_0, xw) = \delta(z_0, yw) \Rightarrow xR_L y$$

Die Sprache  $L \subseteq \Sigma^*$  ist genau dann regulär, wenn der Index der Myhill-Nerode-Äquivalenz  $R_L$  endlich ist.

### 8.4.3 Erkennung durch Monoide – Syntaktisches Monoid

Sei  $L \subseteq \Sigma^*$  eine formale Sprache und  $M$  ein Monoid.

$M$  erkennt  $L$ , wenn eine Teilmenge  $A \subseteq M$  und ein Homomorphismus  $\varphi : \Sigma^* \rightarrow M$  existiert, so dass gilt:

$$\begin{aligned} L &= \varphi^{-1}(A) & \text{d.h. } w \in L &\iff \varphi(w) \in A \\ L &= \varphi^{-1}(\varphi(L)) & \text{d.h. } w \in L &\iff \varphi(w) \in \varphi(L) \end{aligned}$$

Weiter kann man für eine konkrete Sprache  $L$  die *syntaktische Kongruenz* definieren:

$$w_1 \equiv_L w_2 \iff [\forall x, y \in \Sigma^* : xw_1y \in L \iff xw_2y \in L]$$

Basierend auf dieser Kongruenz definieren wir das Quotientenmonoid der Kongruenz, dessen Elemente die Äquivalenzklassen sind.

Das Quotientenmonoid oder auch syntaktisches Monoid bezüglich der syntaktischen Kongruenz wird mit

$$\text{Synt}(L) := (\Sigma^* / \equiv_L)$$

bezeichnet.

Für jede Sprache  $L$  gibt es ein syntaktisches Monoid das  $L$  mit dem Homomorphismus

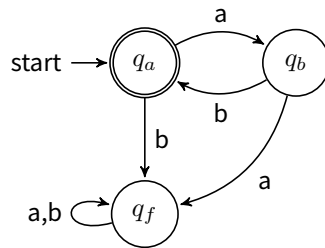
$$\varphi : L \rightarrow \text{Synt}(L), w \mapsto [w]$$

erkennt. Ist  $|\text{Synt}(L)|$  endlich, so ist  $L$  regulär bzw. erkennbar.



## 8.5 Beispiele:

- $L_1 = \Sigma^*$
- $L_2 = L(a(a|b)^*bb(ab)^*)$
- Automat  $M$  mit  $L_3 = T(M) = \{(ab)^n \mid n \in \mathbb{N}_0\}$ :



## 9: ENDLICHE SPRACHEN (FIN):

$$(\text{FIN}) \subset (\text{REG}) \subset (\text{DCFL}) \subset (\text{CFL}) \subset (\text{CSL}) \subset (\text{REC}) \subset (\text{r.e.})$$

Alle endlichen Sprachen sind regulär.

### 9.1 Beispiele:

- $L_1 = \emptyset$  (leere Sprache ist endlich)
- $L_2 = \Sigma$  (nur die Buchstaben)
- $L_3 = \{aaa, baba\}$  (z.B. nur zwei Wörter)

## 10: ALGORITHMEN

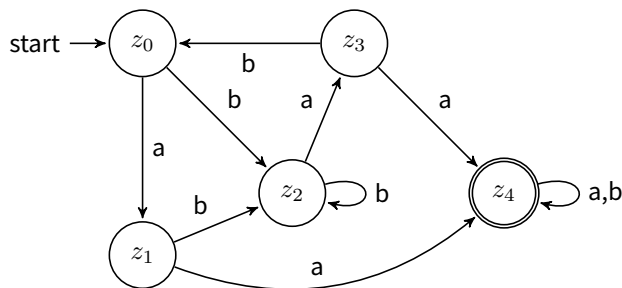
### 10.1 Konstruktionsalgorithmus für Minimal-DEAs:

Mit dem Beweis zur Myhill-Nerode-Äquivalenz wird ein Automat definiert, dieser ist isomorph zum Minimalautomaten. Der Index der Myhill-Nerode-Äquivalenz ist genau die Anzahl der Zustände des Minimalautomaten.

Man kann mit einem einfachen Algorithmus aus einem beliebigen DEA den Minimalautomaten erzeugen:

Wir ermitteln algorithmisch, welche Zustände nicht äquivalent sind und verschmelzen die übrig bleibenden. Nicht äquivalent sind Zustände, bei denen bei Eingabe eines Worts vom einen aus ein Endzustand erreicht wird, vom anderen jedoch nicht.

So sind im ersten Schritt Zustandspaare aus Endzustand und Nichtendzustand nicht äquivalent und werden markiert.



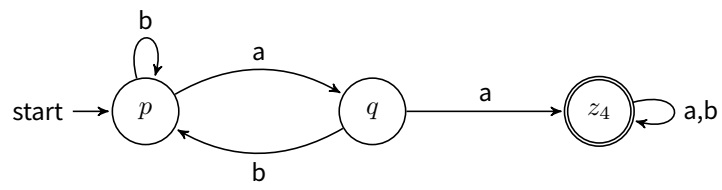
Die Paare  $\{z_i, z_4\}$  mit  $i = 0, 1, 2, 3$  werden markiert:

$z_1$				
$z_2$				
$z_3$				
$z_4$	€	€	€	€
	$z_0$	$z_1$	$z_2$	$z_3$

Durch Testen der Zustandspaare erhält man die untenstehende Tabelle. Hierbei wurden Zeugen für die Antivalenz eingetragen.

$z_1$	a			
$z_2$		a		
$z_3$	a		a	
$z_4$	€	€	€	€
	$z_0$	$z_1$	$z_2$	$z_3$

Damit lassen sich Zustände zusammenfassen:  $p = \{z_0, z_2\}$ ,  $q = \{z_1, z_3\}$ . Der Minimalautomat ist also:



$$T(M) = \{w a a w' \mid w, w' \in \{a, b\}^*\}$$

## 11: ANWEDUNGEN DER SÄTZE

### 11.1 Pumping-Lemma für Typ-3

Für die Sprache  $L = \{a^n b^n \mid n \geq 1\}$ :

Zunächst wählt man ein Wort  $x$ , mit  $|x| \geq n$ :

$$x = a^n b^n \in L, \quad |a^n b^n| = 2n > n$$

Nun zu einer beliebigen Zerlegung  $x = uvw$ , für die die Bedingungen gelten:

$$\begin{aligned} x &= uvw = a^n b^n \\ u &= a^{n-l-k} \quad v = a^l \quad w = a^k b^n \text{ mit } l \geq 1 \\ x &= (a^{n-l-k})(a^l)(a^k b^n) \end{aligned}$$

Pumpt man nun  $v$  mit  $v = 0$ :

$$\begin{aligned} x &= (a^{n-l-k})(a^l)^i (a^k b^n) = (a^{n-l-k})(a^l)^0 (a^k b^n) \\ x &= (a^{n-l-k})(a^k b^n) = a^{n-l} b^n \notin L, \text{ da } l \geq 1 \end{aligned}$$

### 11.2 Myhill-Nerode-Äquivalenz

Für die Sprache  $L = \{w \in \{a, b\}^* \mid w \text{ enthält das Teilwort } abb\}$ :

Finden der Äquivalenzklassen:

$$\begin{aligned} [\epsilon] &= \{\epsilon, b, bb, \dots\} \\ [a] &= \{a, ba, bba, aba, \dots\} = \{wa \mid w \in \{a, b\}^*\} \\ [ab] &= \{ab\} = \{wab \mid w \in \{a, b\}^*\} \\ [abb] &= \{abb, \dots\} = \{wabbw' \mid w, w' \in \{a, b\}^*\} \end{aligned}$$