

# Formale Sprachen und Automaten

Lischke

SS 2004



# Vorwort

*Dieses Dokument wurde als Skript für die auf der Titelseite genannte Vorlesung erstellt und wird jetzt im Rahmen des Projekts „**Vorlesungsskripte der Fakultät für Mathematik und Informatik**“ weiter betreut. Das Dokument wurde nach bestem Wissen und Gewissen angefertigt. Dennoch garantiert weder der auf der Titelseite genannte Dozent, die Personen, die an dem Dokument mitgewirkt haben, noch die Mitglieder des Projekts für dessen Fehlerfreiheit. Für etwaige Fehler und dessen Folgen wird von keiner der genannten Personen eine Haftung übernommen. Es steht jeder Person frei, dieses Dokument zu lesen, zu verändern oder auf anderen Medien verfügbar zu machen, solange ein Verweis auf die Internetadresse des Projekts <http://uni-skripte.lug-jena.de/> enthalten ist.*

*Diese Ausgabe trägt die Versionsnummer 1701 und ist vom 21. Oktober 2008. Eine (mögliche) aktuellere Ausgabe ist auf der Webseite des Projekts verfügbar.*

*Jeder ist dazu aufgerufen, Verbesserungen, Erweiterungen und Fehlerkorrekturen für das Skript einzureichen bzw. zu melden oder diese selbst einzupflegen – einfach eine E-Mail an die **Mailingliste** [<uni-skripte@lug-jena.de>](mailto:uni-skripte@lug-jena.de) senden. Weitere Informationen sind unter der oben genannten Internetadresse verfügbar.*

*Hiermit möchten wir allen Personen, die an diesem Skript mitgewirkt haben, vielmals danken:*

- *Jörg Sommer* [<joerg@alea.gnuu.de>](mailto:joerg@alea.gnuu.de) (2004)
- *Jens Kubicziel* [<jens@kubicziel.de>](mailto:jens@kubicziel.de) (2006)

# Inhaltsverzeichnis

1 Einführung und Wiederholung

6

# Auflistung der Theoreme

## Sätze

Satz 5	Pumping-Lemma; Bar-Hillit, Perles, Shamir . . . . .	9
--------	---	---

## Definitionen und Festlegungen

# 1 Einführung und Wiederholung

## Definition 1

$L$  heißt (formale) Sprache über dem Alphabet  $X$  (d.h.  $X$  ist endl. Menge)  $:\Leftrightarrow L \subseteq X^*$

$X^*$ ... Vereinigung aller Teilmengen

$\epsilon$ ... leeres Wort

$X^+ = X^* \setminus \{\epsilon\}$

## Definition 2

$G = [X, N, S, R]$  heißt (generative) Grammatik  $:\Leftrightarrow$

1.  $X$  endl. Menge (Menge terminaler Symbole)
2.  $N$  endl. Menge (Menge nichtterminaler oder Hilfsymbole),  $X \cap N = \emptyset$
3.  $S \in N$  Anfangs- oder Startsymbol
4.  $R \subseteq (X \cup N)^* \times (X \cup N)^*$  ( $R$  endl. Regelmengemenge)

Die Paare  $[p, q] \in R$  heißen Regel und werden oft auch als  $p \rightarrow q$  geschrieben.

## Definition 3

$w'$  heißt unmittelbare Ableitung von  $w$  bezgl. einer Grammatik  $G = [X, N, S, R]$   $w \rightarrow w' :\Leftrightarrow$

$$w, w' \in (X \cup N)^* \wedge \exists p \exists q \exists p_1 \exists p_2 (p, q, p_1, p_2 \in (X \cup N)^* \wedge [p, q] \in R \wedge w = p_1 p p_2 \wedge w' = p_1 q p_2)$$

$w'$  heißt Ableitung aus  $w$  bezgl.  $G$ ,  $w \rightarrow^* w' \Leftrightarrow$  oder es gibt endl. viele Wörter  $w_0, w_1, \dots, w_r$  mit  $w_0 = w, w_{i-1} \xrightarrow{G} w_i$  für  $i = 1, \dots, r$  und  $w_r = w'$ . Man nennt dies auch eine Ableitung der Länge  $r$  und schreibt dafür  $w_0 \Rightarrow w_1 \Rightarrow \dots \Rightarrow w_r$  oder  $w_0 \xrightarrow{r} w_r$ .

## Definition 4

$L$  heißt die von der Grammatik  $G[\dots]$  erzeugte Sprache,  $L = L(G) :\Leftrightarrow L = \{w : S \xrightarrow{*}_G w\} \cap X^*$

## Definition 5

Zwei Grammatiken  $G_1$  und  $G_2$  heißen äquivalent  $G_1 \sim G_2 :\Leftrightarrow L(G_1) = L(G_2)$

**Definition 6**

Eine Grammatik  $G = [X, N, S, R]$  heißt *Normalformgrammatik*  $:\Leftrightarrow R \subseteq N^+ \times (X \cup N)^*$

**Satz 1**

Zu jeder Grammatik gibt es eine äquivalente Normalformgrammatik.

BEWEIS:

Es sei  $G = [X, N, S, R]$  beliebige Grammatik und zunächst  $[\alpha, \beta] \notin R$  für jedes Wort  $p$

$X' := x' : x' \in X$  wobei  $X' \cap X = X' \cap N = \emptyset$

$R'$  ... Menge aller Regeln, die aus der Regel der Menge  $R$  entstehen, in dem darin alle Buchstaben  $x \in X$  durch die entsprechenden Buchstaben  $x' \in X'$  ersetzt werden.

$R'' := \{[x', x] : x \in X\}$  und  $G' := [X, N \cup X', S', R' \cup R''] \rightarrow G' \sim G$

Ist  $[e, p] \in R$  für gewisse  $p$ , so sei  $R'$  wie oben mit Ausnahme aller Regeln  $[e, p']$ .

Sei  $E \notin X' \cup N$  und  $R_E := \{[B, BE] : B \in X \cup N\} \cup \{[B, EB] : B \in N\} \cup \{[E, p'] : [e, p] \in R \wedge p' \text{ entsteht aus } p \text{ durch Ersetzen aller } x \text{ in } X \text{ durch } x \in X'\}$ ,

$G_E := [N, N \cup X' \cup \{E\}, S, R' \cup R'' \cup R_E] \rightarrow G_E \sim G$ . ■

Avram Noam Chomsky

**Definition 7**

Eine Grammatik  $G = [X, N, S, R]$  (ohne Einschränkung an die Regelmenge  $R$ ) heißt Grammatik vom *Typ 0*

$G$  heißt Grammatik vom *Typ 1* oder nicht verkürzende Grammatik  $:\Leftrightarrow \forall p \forall q ([p, q] \in R \rightarrow l(p) \leq l(q))$

$G$  heißt Grammatik vom *Typ 1'* oder Kontextgrammatik  $\Leftrightarrow \forall p \forall q ([p, q] \in R \rightarrow \exists p_1 \exists p_2 \exists u \exists B (p_1, p_2, u \leq (X \cup N)^* \wedge B \in N \wedge u \neq e \wedge p = p_1 B p_2 \wedge q = p_1 \cup p_2)$

$G$  heißt Grammatik vom *Typ 2* oder kontextfreie Grammatik  $:\Leftrightarrow \forall p \forall q ([p, q] \in R) \rightarrow p \in N \wedge q \in (X \cup N)^+$

$G$  heißt Grammatik vom *Typ 3*, Automaten- oder rechtslineare Grammatik  $:\Leftrightarrow \forall p, q ([p, q] \in R \rightarrow p \in N \wedge q \in X \cdot N \cup X)$

Eine Sprache  $L$  heißt vom *Typ  $i$* , wenn sie von einer Grammatik vom *Typ  $i$*  erzeugt wird.  $\text{rund}L_i$  ist die Klasse aller Sprachen vom *Typ  $i$*  ( $i \in \{0, 1, 1', 2, 3\}$ ).

**Definition 8**

Ist  $L$  eine Sprache vom *Typ  $i$* , so heißt auch  $L \cup \{e\}$  Sprache vom *Typ  $i$*  ( $i \in \{0, 1, 1', 2, 3\}$ )

Ist  $G = [X, N, S, R]$  eine Grammatik vom *Typ  $i$*  und  $S' \notin X \cap N$ , so heißt auch  $G' := [X, N \cup S, S, \{[S', e], [S', S]\} \cup R]$  Grammatik vom *Typ  $i$*  ( $i \in \{0, 1, 1', 2, 3\}$ )

## 1 Einführung und Wiederholung

### Satz 2

Jede Grammtik  $G = [X, N, S, R]$  mit  $R \subseteq N \times (X \cup N)^*$  gibt es eine äquivalente Grammatik  $G = [X, N, S, R]$  mit  $R \subseteq N \times (X^*N \cup X^*)$  gibt es eine äquivalente rechtslineare Grammatik.

BEWEIS:

1.  $G = [X, N, S, R]$  ist Grammatik mit  $R \subseteq N \times (X \cup N)^*$

Falls keine Regel  $[A, e] \in R$ , so  $G$  kontextfrei gemäß D7, andernfalls konstruiere  $G'$  nach Algo

Algo:

- setze  $W_0 := A : [A, e] \in R$  und  $W_{i+1} := W_i \cup A : \exists q([A, q] \in R \wedge q \in W_i^*)$  solange bis  $W_{i+1} = W_i$
- Dann  $W := W_i, G' = [X, N, S, R']$  mit  $R' := [p, q'] : q \neq e \wedge \exists q([p, q] \in R \wedge (q' = q \vee q \text{ entsteht aus } q'))$   $G'$  kontextfrei gemäß D7 und  $A \in W \xleftrightarrow{G} A \in N \wedge A \xrightarrow{G} e$ .

Behauptung:  $L(G') = L(G)\{e\} \subseteq B \rightarrow w' \rightarrow [B, w'] \in R' \rightarrow \text{ex. } w = b_1 \dots b_n$  mit  $b_r \in X \cup N, [B, w] \in R \wedge w' = b_{i_1} \dots b_{i_k}$ , wobei  $b_{i_1} \dots b_{i_k}$  Teilfolge von  $b_1 \dots b_n$  und  $b \in W$  für  $b \in \{b_1, \dots, b_n\} \setminus \{b_{i_1}, \dots, b_{i_k}\} \rightarrow B \xrightarrow{G} w'$ .

$p \in L(G') \rightarrow S \xrightarrow{G'} p \rightarrow S \xrightarrow{G} p \rightarrow p \in L(G). L \notin L(G')$ , da nicht verkürzend  $\rightarrow L(G') \subseteq L(G) \setminus \{e\}$

rückw  $\subseteq: G' := [X, N, S, R \cup R'] \rightarrow L(G) \subseteq L(G') \rightarrow L(G) \setminus \{e\} \subseteq L(G') \setminus \{e\} \rightarrow$  reicht z. z.  $L(G'') \setminus \{e\} \subseteq L(G)$

Entsprechend ex.  $w \in L(G) \setminus \{e\}, w \notin L(G')$ .

Es sei  $w_0$

$\Rightarrow_G w_1 \Rightarrow \dots \Rightarrow_G w_t$  eine Ableitung von  $w$  aus  $S$  bzgl.  $G''$  kürzester Länge  $w_0 \Rightarrow S, w_t = w$

Da  $w \in L(G'') \setminus L(G')$  auf wenigstens einmal eine Regel aus  $R \setminus R'$  angewandt worden sein, d. h. eine Regel  $[A, e]$  etwa bei  $w_i \Rightarrow_{G''} w_{i+1} \rightarrow A \in W, w_i = w_{i+1}Aw_{i+2}, w_{i+1} = w_1w_2 \rightarrow ej < imitw_j = w_{j+1}Bw_{j+2}, w_{j+1} = w_{j+1}q_jAq_jw_{j+2}, q_1q_2 \neq q \rightarrow [B, q_1Aq_2] \in R', w_{j+1}q_1 \Rightarrow_{G''}^{i-(i+1)} w_1, q_1w_{j+2} \Rightarrow_{G''}^r w_{j+2}, 0 \leq i \leq i - (j+1) \rightarrow [B, q_1q_2] \in R' \rightarrow w_q \xrightarrow{G}^j w_1Bw_2 \rightarrow_G w_{j+1}q_1q_2w_{j+2} \rightarrow_{G''}^{i-j+1} w_1w_2 \rightarrow_{G''}^{t+(i+1)}$  ist Ableitung von  $w$  aus  $S$  der Länge  $t-1$  Widerspruch dazu, dass  $t$  eine stets kürzeste Länge ist.  $\rightarrow L(G'') \setminus \{e\} \subseteq L(G')$

Ist  $q_1q_2 = e[B, A] \in R' \rightarrow B \in W \rightarrow \alpha j' < jnitw_j = w_{j+1}Cw_{j+2}w_{j+1} = w_{j+1} \dots \rightarrow$  Schluss analog mit  $BX$  statt  $AB$

$e \notin L(G) \rightarrow G \sim G'$

$e \in L(G) \rightarrow \text{ex. kontextfreie gram. für } L(G'') \cup e \text{ nach D7!}$  ■



**Satz 3**

$REG = L_3 \subset L_2 \subset L_1 = L_{1'} \subset REC \subset RE = L_0$

**Definition 9**

$7'' CS := L_1 = L_{1'}$  heißt Klasse der kintextsensitiven oder kontextsprachen

$CF := L_2$  heißt Klasse der kontextfreien oder Chompsky-Spachen.

$REG := L_3$  heißt Klasse der regulären, Automaten- oder rechstlinearen Sprachen.

**Satz 4**

$REG \subset CF \subset CS \subset REC \subset RE = L_0$

BEWEIS:

$CF \subset noteq$

$L := \{a^n b^n c^n : n \in \mathbb{N} \setminus \{0\}\} G = [a, b, c, S, A, B, S, [S, aSBA], [A, abA], [AB, BA], [bB, bb], [A, c]] \longrightarrow$   
 $L(G) = L$  es ist leicht ersichtlich, dass  $G$   $L$  erzeugt und  $G$  kontextsensitiv ist. ■

**Satz 5 (Pumping-Lemma; Bar-Hillit, Perles, Shamir)**

Zu jeder kontextfreien Sprache  $L$  existiert eine Zahl  $n \in \mathbb{N}$ , so dass jedes Wort  $w \in L$  mit  $l(w) > m$  in der Fom  $w = w_1 w_2 w_3 w_4 w_5$  dargestellt werden kann, wobei  $w_2 w_4 \neq \epsilon$ ,  $l(w_2 w_3 w_4) < m$  mit  $w_1 w_2 w_3 w_4 w_5 \in L$  für jedes  $i \in \mathbb{N}$ .

**Satz 6**

Jede kontextfreie Sprache über einbuchstabigem Alphabet ist regulär.

BEWEIS:

Sei  $L \in CE, L \subset a^*$ .

Nach Satz 4 existiert  $n \in \mathbb{N}$ , so dass für jedes Wort  $a^l \in L$  mit  $l > m$  folgt:  $l = i + j$  mit  $i \leq j \leq m$  mit  $a^{i+kj} \in L$  für jedes  $k \in \mathbb{N}$

$L_0 := \{a^i : 1 \leq i \leq m\} \cap L$  Falls  $L \setminus L_0 \neq \emptyset$  sei  $l_1 := \min\{l : a^l \in L \setminus L_0\}$

Nach Satz 4 ist  $l_1 = i_1 + j_1$  mit  $1 \leq j_1 \leq m, 0 \leq i_1 \leq m$  mit  $a^{i_1+kj_1} : k \in \mathbb{N} \subseteq L$

$L_1 := \{a^{i_1+kj_1} : k \in \mathbb{N}\}$

Seien  $L_1, L_2, \dots, L_{i'}$  konstruiert mit  $L \setminus \bigcup_{i=0}^{i'} L_i, sol_{i'+1} := \min\{l : a^l \in L \setminus \bigcup_{r=0}^{i'} L_r\} \longrightarrow l_{i+j} = i + j$  mit  $1 \leq j \leq m, 0 \leq i$  und  $L_{i+1} := \{a^{i+kj} : k \in \mathbb{N}\} \subseteq L$ .

Nach endl. vielen Schritten muss  $L \setminus \bigcup_{r=0}^i L_r = \emptyset$  eintreten  $\longrightarrow L = \bigcup_{r=0}^i L_r, L_0, \dots, L_i \in REG \longrightarrow L \in REG$  ■

(A4)  $CF_1 \subset CS_1 \subset REC_1 \subset RE_1$