

Einführung in die Computerlinguistik

– Formale Grammatiken

rechtslineare und kontextfreie Grammatiken

Kellerautomaten

Dozentin: Wiebke Petersen

13. Foliensatz

Formale Grammatik

Definition

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

- einer Alphabet von Terminalsymbolen T (häufig auch Σ)
- einer Alphabet von Nichtterminalsymbolen N mit $N \cap T = \emptyset$
- einem Startsymbol $S \in N$
- einer Menge von Regeln/Produktionen
 $P \subseteq \{ \langle \alpha, \beta \rangle \mid \alpha, \beta \in (N \cup T)^* \text{ und } \alpha \notin T^* \}$.

Für eine Regel $\langle \alpha, \beta \rangle$ schreiben wir auch $\alpha \rightarrow \beta$.

S	→	NP VP	VP	→	V	NP	→	D N
D	→	the	N	→	cat	V	→	sleeps

Generiert: the cat sleeps

Konvention: Wir verwenden Großbuchstaben für Nichtterminalsymbole und Kleinbuchstaben für Terminalsymbole

Terminologie

$$G = \langle \{S, NP, VP, V, D, N, EN\}, \{\text{the, cat, peter, chases}\}, S, P \rangle$$

$$P = \left\{ \begin{array}{llll} S & \rightarrow & NP VP & VP & \rightarrow & V NP & NP & \rightarrow & D N \\ NP & \rightarrow & EN & D & \rightarrow & \text{the} & N & \rightarrow & \text{cat} \\ EN & \rightarrow & \text{peter} & V & \rightarrow & \text{chases} & & & \end{array} \right\}$$

“NP VP” ist **in einem Schritt** aus S **ableitbar**

“the cat chases peter” ist **ableitbar** aus S :

Ableitung:

$$\begin{array}{lll} S & \rightarrow & NP VP & \rightarrow & NP V NP & \rightarrow & NP V EN \\ & \rightarrow & NP V \text{ peter} & \rightarrow & NP \text{ chases peter} & \rightarrow & D N \text{ chases peter} \\ & \rightarrow & D \text{ cat chases peter} & \rightarrow & \text{the cat chases peter} & & \end{array}$$

Die Menge aller aus dem Startsymbol S ableitbarer Wörter ist die von der Grammatik G **erzeugte Sprache** $L(G)$.

$$L(G) = \left\{ \begin{array}{ll} \text{the cat chases peter,} & \text{peter chases the cat,} \\ \text{peter chases peter,} & \text{the cat chases the cat} \end{array} \right\}$$

rechtslineare Grammatiken

Definition

Eine Grammatik (N, T, S, P) heißt **rechtslinear**, wenn alle Regeln/Produktionen die folgende Form haben:

$A \rightarrow a$ oder $A \rightarrow aB$ wobei $a \in T \cup \{\varepsilon\}$ und $A, B \in N$.

Eine durch eine rechtslineare Grammatik erzeugte Sprache heißt **rechtslinear**.

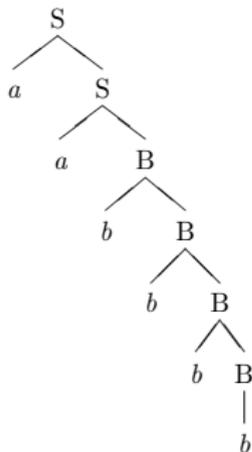
(Vorsicht, häufig werden schärfere [aber äquivalente] Bedingungen gefordert!)

Beispiel:

$G = (\{S, B\}, \{a, b\}, S, P)$ mit

$$P = \left\{ \begin{array}{l} S \rightarrow aS \\ S \rightarrow aB \\ B \rightarrow bB \\ B \rightarrow b \end{array} \right\}$$

G generiert $L(G) = L(a^+b^+)$



Ableitungsbaum $(aabbbb \in L(G))$

rechtslineare Grammatiken und reguläre Sprachen

Theorem

Sei L eine formale Sprache, dann sind die folgenden Aussagen äquivalent:

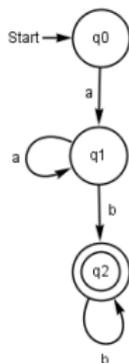
- ① L ist regulär.
- ② Es gibt eine rechtslineare Grammatik G , die L erzeugt.
- ③ Es gibt einen endlichen Automaten A , der L akzeptiert.
- ④ Es gibt einen regulären Ausdruck R , der L beschreibt.

rechtslineare Grammatik:

$$G = (\{S, B\}, \{a, b\}, S, P) \text{ mit}$$

$$P = \left\{ \begin{array}{l} S \rightarrow aS \\ S \rightarrow aB \\ B \rightarrow bB \\ B \rightarrow b \end{array} \right\}$$

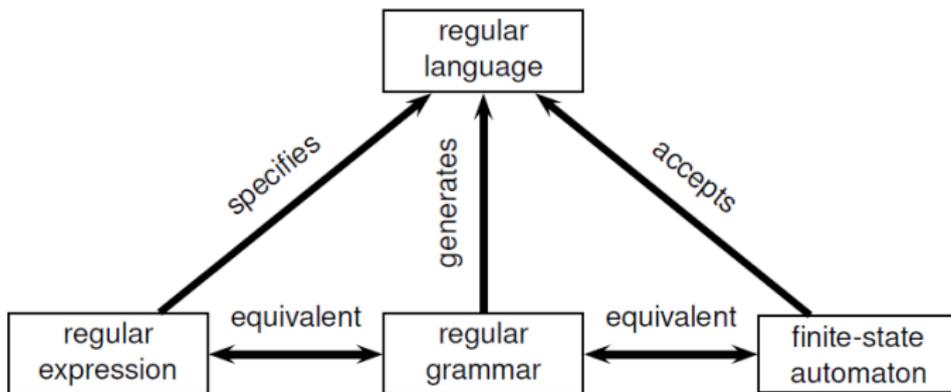
endlicher Automat:



regulärer Ausdruck:

$a^+ b^+$

Zusammenfassung: reguläre Sprachen



kontextfreie Grammatik

Definition

Eine Grammatik (N, T, S, P) heißt **kontextfrei**, wenn alle Regeln/Produktionen die folgende Form haben:

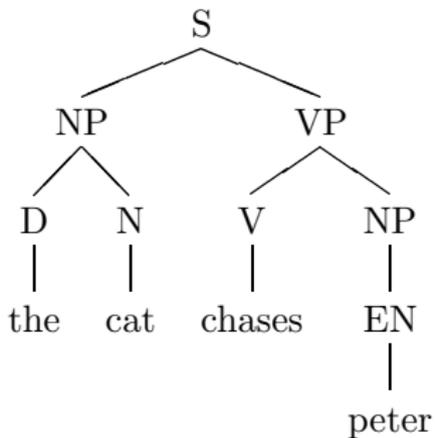
$$A \rightarrow \alpha, \text{ wobei } A \in N \text{ und } \alpha \in (T \cup N)^*.$$

Eine durch eine kontextfreie Grammatik erzeugte Sprache heißt **kontextfrei**.

Die Menge der kontextfreien Sprachen ist eine echte Obermenge der Menge der regulären Sprachen

Beweis: Jede reguläre Sprache ist per Definition auch kontextfrei und es gibt mindestens eine kontextfreie Sprache, nämlich $L(a^n b^n)$ mit $n \geq 0$, die nicht regulär ist. ($S \rightarrow aSb, S \rightarrow \varepsilon$)

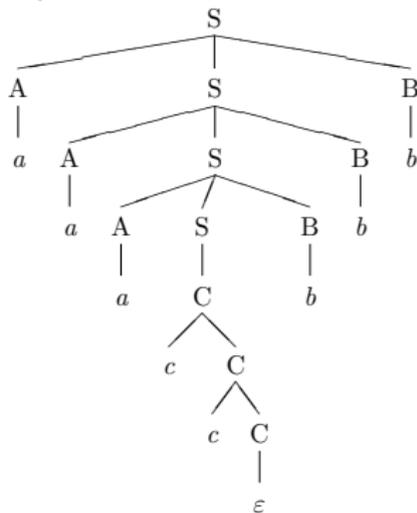
Ableitungsbaum



Beispiel einer kontextfreien Sprache

$$G = \langle \{S, A, B, C\}, \{a, b, c\}, S, P \rangle$$

$$P = \left\{ \begin{array}{ll} S \rightarrow ASB & S \rightarrow C \\ A \rightarrow a & B \rightarrow b \\ C \rightarrow cC & C \rightarrow \varepsilon \end{array} \right\}$$

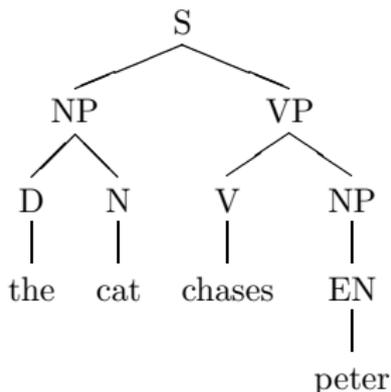


Übung: Welche Sprache generiert diese Grammatik? Können Sie eine äquivalente Grammatik angeben, die mit weniger Nichtterminalsymbolen auskommt?

Linksableitung

Gegeben eine kontextfreie Grammatik G . Eine Ableitung bei der stets das am weitesten links stehende nichtterminale Symbol ersetzt wird, heißt **Linksableitung**

$S \rightarrow NP VP \quad \rightarrow D N VP \quad \rightarrow \text{the } N VP$
 $\rightarrow \text{the cat } VP \quad \rightarrow \text{the cat } V NP \quad \rightarrow \text{the cat chases } NP$
 $\rightarrow \text{the cat chases } EN \quad \rightarrow \text{the cat chases peter}$



Zu jeder Linksableitung gibt es genau einen Ableitungsbaum und zu jedem Ableitungsbaum gibt es genau eine Linksableitung.

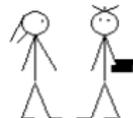
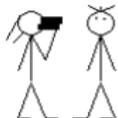
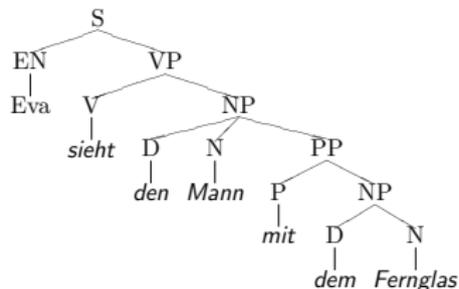
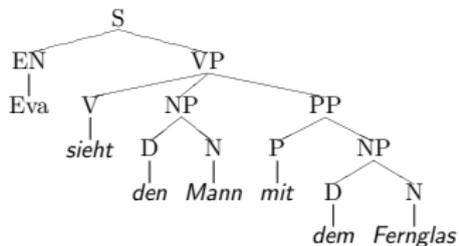
ambige Grammatik

Eine Grammatik G heißt **ambig**, wenn es für ein Wort $w \in L(G)$ mehr als eine Linksableitung gibt.

$G = (N, T, NP, P)$ mit $N = \{S, EN, NP, VP, PP, D, N, P\}$,

$T = \{\text{Eva, sieht, den, Mann, mit, dem, Fernglas}\}$,

$$P = \left\{ \begin{array}{lll} S \rightarrow EN VP & VP \rightarrow V NP & VP \rightarrow V NP PP \\ NP \rightarrow D N & NP \rightarrow D N PP & PP \rightarrow P NP \\ EN \rightarrow \text{Eva} & P \rightarrow \text{mit} & V \rightarrow \text{sieht} \\ D \rightarrow \text{den} & D \rightarrow \text{dem} & N \rightarrow \text{Mann} \\ N \rightarrow \text{Fernglas} & & \end{array} \right\}$$



Kellerautomaten

Ziel: Automatenmodell mit dem genau die kontextfreien Sprachen akzeptiert werden können (analog zu endlichen Automaten und regulären Sprachen).

Lösung: Hinzunahme eines unbeschränkten Speichers in Form eines Stapels, von dessen Spitze etwas genommen und auf dessen Spitze etwas abgelegt werden kann.

Kellerautomaten entstehen aus endlichen Automaten durch

- Hinzunahme eines Kelleralphabets
- Erweiterung der Transitionen (es muss das Lesen und Ersetzen der Kellerspitze realisiert werden)

Informelles Beispiel

Wir betrachten die Sprache $\{a^i b^i \mid i > 0\}$.

Das Akzeptieren eines Eingabewortes geschieht wie folgt:

1. **Aufbau des Kellers:** für jedes gelesene a lege ein Symbol auf dem Keller ab (wir nehmen das Symbol Z)
2. **Abbau des Kellers:** für jedes gelesene b nehme ein Symbol Z vom Keller herunter
3. **durch zwei Kontrollzustände** Sorge dafür, dass Aufbau und Abbau nur in dieser Reihenfolge möglich sind (Aufbau mit q_0 , Abbau mit q_1 , keine Rückkehr nach q_0)
4. **Akzeptiere**, wenn am Ende des Eingabewortes der Kellerboden erreicht ist.

Der Keller realisiert in diesem Fall eine Zählervariable.

Transition

oberstes Symbol auf dem Stack
(wird entfernt)
- pop up -

neuer Zustand

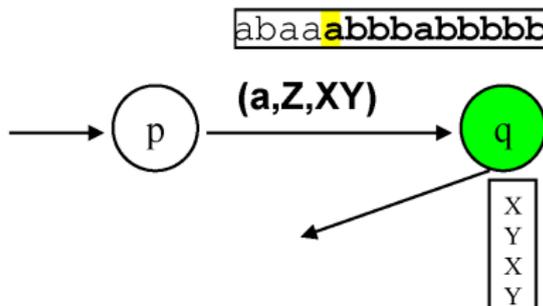
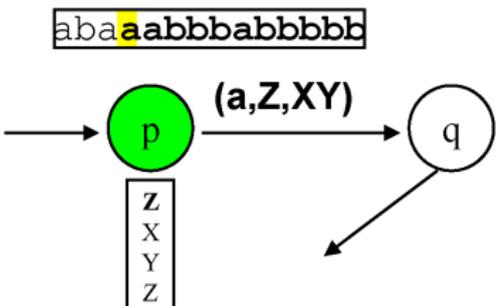
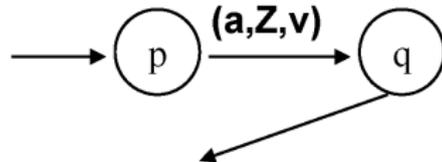
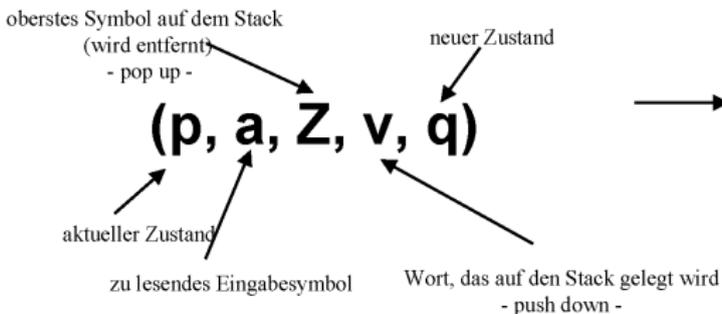
(p, a, Z̄, v, q)

aktueller Zustand

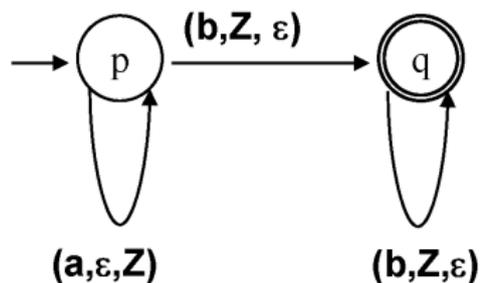
zu lesendes Eingabesymbol

Wort, das auf den Stack gelegt wird
- push down -

Transition

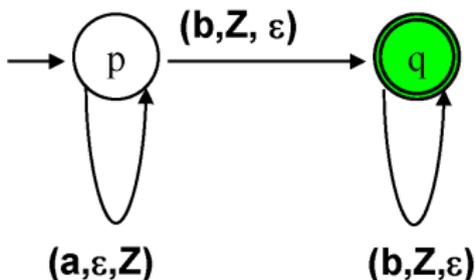


Beispiel eines Kellerautomaten



dieser Kellerautomat akzeptiert
die Sprache $a^n b^n$

Arbeitsweise eines Kellerautomaten



Der Automat befindet
sich in einem Endzustand!

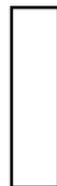
Darum akzeptiert der
Kellerautomat das Wort!

(noch) zu lesendes Wort:



Das Wort ist abgearbeitet!

aktueller Stack:



Der Stack ist leer!

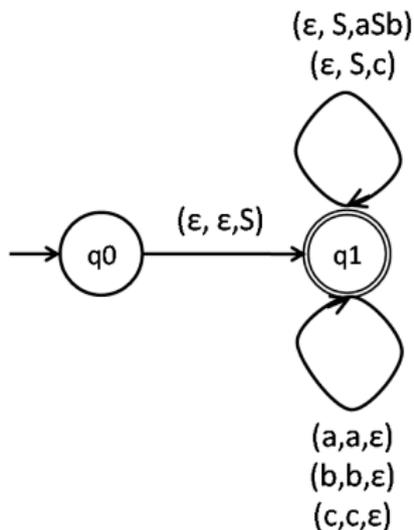
Von kontextfreien Grammatiken zu Kellerautomaten

- Für das Akzeptieren einer kontextfreien Sprache genügt ein Kellerautomat mit nur zwei Zuständen, wobei die einzige Aufgabe des Startzustands darin besteht, das Startsymbol S der Grammatik in den Keller zu legen.
- Während der eigentlichen Rechnung befindet sich der Automat permanent in dem selben Zustand (dem Nichtstartzustand), die Rechnung findet nur in dem Keller statt.
- Der konstruierte Automat vollzieht zwei verschiedene Arbeitsschritte:
 - Nicht-Leseschritt mit Bezug auf Grammatikregel $B \rightarrow \beta$: Ersetze die Kellerspitze B mit β (Expansion).
 - Leseschritte:
Lies ein $a \in T$ der Eingabekette und entferne a von der Kellerspitze (Scan).

Beispiel: Von kontextfreien Grammatiken zu Kellerautomaten

Grammatik: $(\{S\}, \{a, b, c\}, S, \{S \rightarrow aSb, S \rightarrow c\})$

generierte Sprache: $L(a^n cb^n)$



Hausaufgaben

- ① Sei L die Sprache, die aus allen nichtleeren Wörtern über dem Alphabet $\{a, b\}$ besteht, in denen auf jedes a unmittelbar ein b folgt. Beispiele für Wörter dieser Sprache: $bbbab$, $abababab$, bb , $babbbbab$.
- (a) geben Sie eine rechtslineare Grammatik G an, die L erzeugt und zeichnen Sie den Ableitungsbaum für das Wort $bbababb$
- (b) geben Sie einen endlichen Automaten A an, der L akzeptiert.
- (c) geben Sie einen regulären Ausdruck R an, der L beschreibt.
- ② Geben sie jeweils eine kontextfreie Grammatik zu den folgenden Sprachen an:
- (a) $L_1 = \{a^i b^j \mid i > j > 0\}$
- (b) $L_2 = \{w \in \{a, b\}^* \mid w \text{ ist ein Palindrom}\}$

Wählen Sie pro Sprache ein Wort, das mindestens die Länge 5 hat, und zeichnen Sie den Ableitungsbaum in Bezug auf Ihre Grammatik.