

1b**Grammatiken und die Chomsky-Hierarchie**

PD Dr. Jan Johannsen

Institut für Informatik

Stand: 25. März 2026

Basierend auf Folien von PD Dr. David Sabel und Prof. Dr. Jasmin Blanchette



Formale Sprachen darstellen

- ▶ Sei Σ ein **Alphabet**.
- ▶ Eine **Sprache** über Σ ist eine Teilmenge von Σ^* .
- ▶ Für $\Sigma = \{ (,), +, -, *, /, a \}$ sei L_{ArEx} die Sprache aller korrekt geklammerten Ausdrücke.
Z.B. $((a + a) - a) * a \in L_{ArEx}$ aber $(a -) + a \notin L_{ArEx}$.
- ▶ Unsere bisherigen Operationen auf Sprachen (Mengen) können das nicht darstellen.

Benötigt: Formalismus, um L_{ArEx} zu beschreiben

Formale Sprachen darstellen

Anforderungen:

- ▶ Beschreibung muss endlich sein.
- ▶ Sprache selbst muss aber auch unendlich viele Objekte erlauben.

Zwei wesentliche solchen Formalismen sind

- ▶ Grammatiken
- ▶ Automaten.

Grammatik für einen sehr kleinen Teil der deutschen Sprache:

<Satz> → <Subjekt><Prädikat><Objekt>

<Subjekt> → <Artikel><Attribut><Nomen>

<Objekt> → <Artikel><Attribut><Nomen>

<Artikel> → ϵ

<Artikel> → der

<Artikel> → das

<Attribut> → <Adjektiv>

<Attribut> → <Adjektiv><Attribut>

<Adjektiv> → kleine

<Adjektiv> → große

<Adjektiv> → nette

<Adjektiv> → blaue

<Nomen> → Frau

<Nomen> → Mann

<Nomen> → Auto

<Prädikat> → fährt

<Prädikat> → liebt

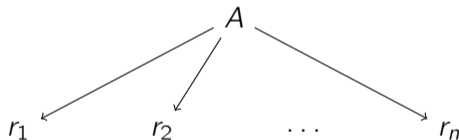
- ▶ **Grammatik** = endliche Menge von Regeln „linke Seite \rightarrow rechte Seite“
- ▶ Symbole in spitzen Klammern wie \langle Artikel \rangle sind **Variablen**, d.h. sie sind **Platzhalter**, die weiter **ersetzt** werden müssen.
- ▶ Z.B. kann

der kleine nette Mann fährt das große blaue Auto

durch die vorige Grammatik abgeleitet werden.

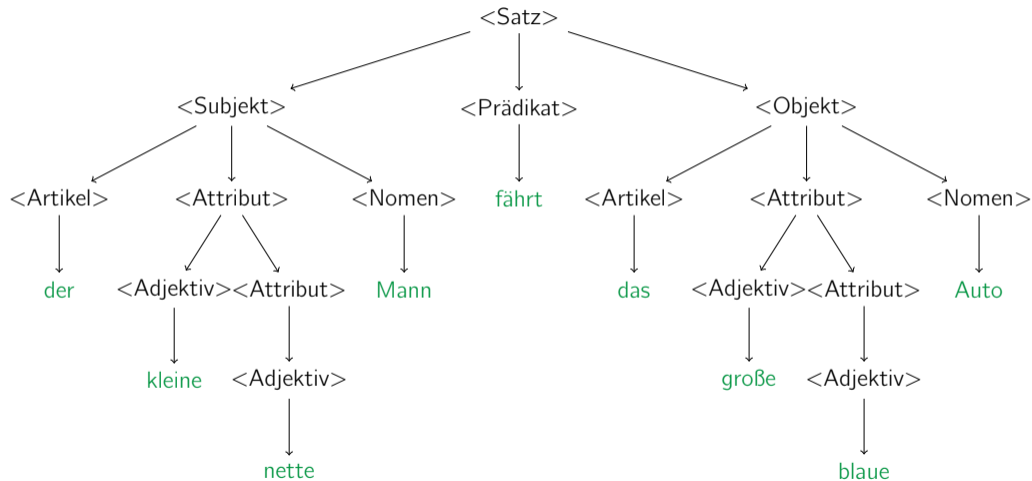
Syntaxbäume

- ▶ Ein **Syntaxbaum** stellt dar, wie ein Satz entsteht.
- ▶ Die Anwendung von $A \rightarrow r_1 r_2 \dots r_n$ wird durch den Teilbaum



dargestellt.

Syntaxbaum zum Beispiel



Definition einer Grammatik

Definition

Eine **Grammatik** ist ein 4-Tupel $G = (V, \Sigma, P, S)$ wobei:

- ▶ V ist eine endliche Menge von **Variablen** (alternativ **Nichtterminalen**)
- ▶ Σ (mit $V \cap \Sigma = \emptyset$) ist ein **Alphabet** von **Zeichen** (alternativ **Terminalen**)
- ▶ P ist eine endliche Menge von **Produktionen** (alternativ **Regeln**)
von der Form $\ell \rightarrow r$ wobei $\ell \in (V \cup \Sigma)^+$ und $r \in (V \cup \Sigma)^*$
- ▶ $S \in V$ ist das **Startsymbol** (alternativ **Startvariable**).

Produktionen mit ε auf der rechten Seite heißen **ε -Produktionen**.

Manchmal genügt es, die Produktionen P alleine zu notieren
(wenn klar ist, was V , Σ und S sind).

Beispiel für eine Grammatik

$G = (V, \Sigma, P, E)$ mit

$V = \{E, M, Z\},$

$\Sigma = \{+, *, 1, 2, (,)\}$ und

$P = \{E \rightarrow M,$

$E \rightarrow E + M,$

$M \rightarrow Z,$

$M \rightarrow M * Z,$

$Z \rightarrow 1,$

$Z \rightarrow 2,$

$Z \rightarrow (E)\}$

Definition

Sei $G = (V, \Sigma, P, S)$ eine Grammatik.

Eine **Satzform** ist ein Wort aus $(V \cup \Sigma)^*$.

Satzform u **geht unter** Grammatik G **unmittelbar in** Satzform v **über**, $u \Rightarrow_G v$, wenn

$$u = w_1 \ell w_2 \text{ und } v = w_1 r w_2 \text{ mit } \ell \rightarrow r \in P$$

- ▶ Wenn G klar ist, schreiben wir $u \Rightarrow v$ statt $u \Rightarrow_G v$.
- ▶ \Rightarrow^* ist die reflexiv-transitive Hülle von \Rightarrow . Sie ist definiert durch folgende Regeln (und nur diese):
 - ▶ falls $u \Rightarrow v$, dann ist $u \Rightarrow^* v$
 - ▶ $u \Rightarrow^* u$
 - ▶ falls $u \Rightarrow^* v$ und $v \Rightarrow^* w$, dann ist $u \Rightarrow^* w$.

Definition

Sei $G = (V, \Sigma, P, S)$ eine Grammatik.

Eine Folge (w_0, w_1, \dots, w_n) mit $w_0 = S$, $w_n \in \Sigma^*$ und $w_{i-1} \Rightarrow w_i$ für $i = 1, \dots, n$ heißt **Ableitung** von w_n .

Statt (w_0, \dots, w_n) schreiben wir auch $w_0 \Rightarrow \dots \Rightarrow w_n$.

Beispiel für eine Ableitung

$G = (V, \Sigma, P, E)$ mit $V = \{E, M, Z\}$ und $\Sigma = \{+, *, 1, 2, (,)\}$ und
 $P = \{E \rightarrow M, E \rightarrow E + M, M \rightarrow Z, M \rightarrow M * Z, Z \rightarrow 1, Z \rightarrow 2, Z \rightarrow (E)\}$

Eine Ableitung von $(2 + 1) * (2 + 2)$:

$E \Rightarrow M \Rightarrow M * Z \Rightarrow Z * Z \Rightarrow Z * (E) \Rightarrow Z * (E + M)$
 $\Rightarrow (E) * (E + M) \Rightarrow (E) * (E + Z) \Rightarrow (E + M) * (E + Z)$
 $\Rightarrow (M + M) * (E + Z) \Rightarrow (M + M) * (M + Z)$
 $\Rightarrow (M + M) * (Z + Z) \Rightarrow (M + M) * (Z + 2)$
 $\Rightarrow (M + Z) * (Z + 2) \Rightarrow (M + Z) * (2 + 2)$
 $\Rightarrow (Z + Z) * (2 + 2) \Rightarrow (2 + Z) * (2 + 2) \Rightarrow (2 + 1) * (2 + 2)$

Ableitungen sind nicht eindeutig

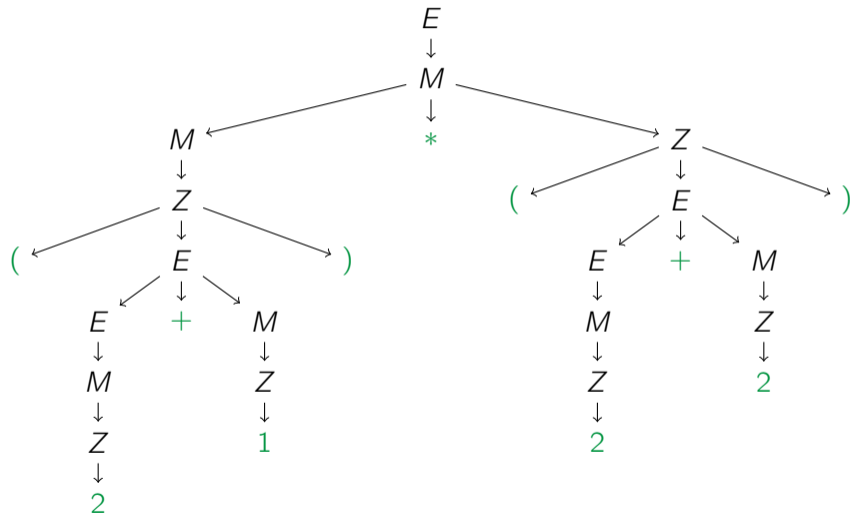
Ableitung von letzter Folie:

$$\begin{aligned} E &\Rightarrow M \Rightarrow M * Z \Rightarrow Z * Z \Rightarrow Z * (E) \Rightarrow Z * (E + M) \\ &\Rightarrow (E) * (E + M) \Rightarrow (E) * (E + Z) \Rightarrow (E + M) * (E + Z) \\ &\Rightarrow (M + M) * (E + Z) \Rightarrow (M + M) * (M + Z) \\ &\Rightarrow (M + M) * (Z + Z) \Rightarrow (M + M) * (Z + 2) \\ &\Rightarrow (M + Z) * (Z + 2) \Rightarrow (M + Z) * (2 + 2) \\ &\Rightarrow (Z + Z) * (2 + 2) \Rightarrow (2 + Z) * (2 + 2) \Rightarrow (2 + 1) * (2 + 2) \end{aligned}$$

Linksableitung (ersetzt immer die linkeste Variable):

$$\begin{aligned} E &\Rightarrow M \Rightarrow M * Z \Rightarrow Z * Z \Rightarrow (E) * Z \\ &\Rightarrow (E + M) * Z \Rightarrow (M + M) * Z \Rightarrow (Z + M) * Z \\ &\Rightarrow (2 + M) * Z \Rightarrow (2 + Z) * Z \Rightarrow (2 + 1) * Z \Rightarrow (2 + 1) * (E) \\ &\Rightarrow (2 + 1) * (E + M) \Rightarrow (2 + 1) * (M + M) \Rightarrow (2 + 1) * (Z + M) \\ &\Rightarrow (2 + 1) * (2 + M) \Rightarrow (2 + 1) * (2 + Z) \Rightarrow (2 + 1) * (2 + 2) \end{aligned}$$

Syntaxbaum zu beiden Ableitungen



Nichtdeterminismus beim Ableiten

Für eine Satzform u kann es verschiedene Satzformen v geben mit $u \Rightarrow_G v$.

Quellen des Nichtdeterminismus:

- ▶ Wähle **welche Produktion** $\ell \rightarrow r$ aus P angewendet wird.
- ▶ Wähle die **Position des Teilworts** ℓ in u , das durch r ersetzt wird.

Aber: Es gibt **nur endlich viele Satzformen** v für jeden Schritt.

Definition

Die von einer Grammatik $G = (V, \Sigma, P, S)$ erzeugte Sprache $L(G)$ ist

$$L(G) := \{w \in \Sigma^* \mid S \Rightarrow_G^* w\}$$

Beispiele für erzeugte Sprache

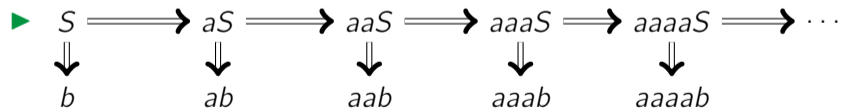
$$G_1 = (\{S\}, \{a\}, \{S \rightarrow aS\}, S)$$

$$L(G_1) = ? \emptyset$$

- ▶ $S \Rightarrow aS \Rightarrow aaS \Rightarrow \dots$ endet nie.
- ▶ Andere Ableitungen gibt es nicht.
- ▶ Daher sind keine Wörter aus $\{a\}^*$ ableitbar.

$$G_2 = (\{S\}, \{a, b\}, \{S \rightarrow aS, S \rightarrow b\}, S)$$

$$L(G_2) = ? \{a^n b \mid n \in \mathbb{N}\}$$



- ▶ Für alle $i \in \mathbb{N}$ gilt $S \Rightarrow^i a^i S \Rightarrow a^i b$.

Die Chomsky-Hierarchie

Noam Chomsky teilte die Grammatiken in Typen 0 bis 3 nach Art der erlaubten Regeln.

Definition

Sei $G = (V, \Sigma, P, S)$ eine Grammatik.

- ▶ G ist immer vom Typ 0.
- ▶ G ist vom Typ 1 (alternativ kontextsensitiv), wenn:
für alle $\ell \rightarrow r \in P$ gilt $|\ell| \leq |r|$.
- ▶ G ist vom Typ 2 (alternativ kontextfrei), wenn:
 G ist vom Typ 1 und für alle $\ell \rightarrow r \in P$ gilt $\ell \in V$.
- ▶ G ist vom Typ 3 (alternativ regulär), wenn:
 G ist vom Typ 2 und für alle $A \rightarrow r \in P$ gilt $r = a$ oder $r = aA'$ für
 $a \in \Sigma, A' \in V$ (d.h. die rechten Seiten sind Satzformen aus $\Sigma \cup \Sigma V$).

Definition

Für $i \in \{0, 1, 2, 3\}$ nennt man eine formale Sprache $L \subseteq \Sigma^*$ vom Typ i , falls es eine Typ i -Grammatik G gibt, sodass $L(G) = L$ gilt.

Spricht man von dem Typ einer formalen Sprache, so ist meistens der größtmögliche Typ gemeint.

Beispiele für die Chomsky-Hierarchie

$G_1 = (\{S\}, \{a, b\}, \{S \rightarrow aS, S \rightarrow b\}, S)$ ist vom Typ 3 (regulär).

$G_2 = (\{E, M, Z\}, \{+, *, 1, 2, (,)\}, P, E)$ mit
 $P = \{E \rightarrow M, E \rightarrow E + M, M \rightarrow Z, M \rightarrow M * Z,$
 $Z \rightarrow 1, Z \rightarrow 2, Z \rightarrow (E)\}$ ist vom Typ 2 (kontextfrei).

$G_3 = (\{S, B, C\}, \{a, b, c\}, P, S)$ mit
 $P = \{S \rightarrow aSBC, S \rightarrow aBC, CB \rightarrow BC, aB \rightarrow ab,$
 $bB \rightarrow bb, bC \rightarrow bc, cC \rightarrow cc\}$ ist vom Typ 1 (kontextsensitiv).

$G_4 = (\{S, T, A, B, \$\}, \{a, b\}, P, S)$ mit
 $P = \{S \rightarrow \$T\$, T \rightarrow aAT, T \rightarrow bBT, T \rightarrow \epsilon, \$a \rightarrow a\$,$
 $\$b \rightarrow b\$, Aa \rightarrow aA, Ab \rightarrow bA, Ba \rightarrow aB, Bb \rightarrow bB,$
 $A\$ \rightarrow \$a, B\$ \rightarrow \$b, \$\$ \rightarrow \epsilon\}$ ist vom Typ 0.

Sonderregeln für ε -Produktionen

- ▶ Das leere Wort ε kann bisher nicht für Typ 1, 2, 3-Grammatiken erzeugt werden. Die Produktion $S \rightarrow \varepsilon$ erfüllt die Typ 1-Bedingung $|S| \leq |\varepsilon|$ nicht.

Daher:

1. Sonderregel: ε -Produktion in Typ 1, 2, 3-Grammatiken

Eine Grammatik $G = (V, \Sigma, P, S)$ vom Typ 1, 2 oder 3 darf eine Produktion $S \rightarrow \varepsilon \in P$ enthalten, vorausgesetzt, dass S auf keiner rechten Seite einer Produktion in P vorkommt.

Zudem:

2. Sonderregel: ε -Produktionen in Typ 2, 3-Grammatiken

Eine Grammatik $G = (V, \Sigma, P, S)$ vom Typ 2 oder 3 darf Produktionen von der Form $A \rightarrow \varepsilon \in P$ enthalten, wo $A \in V \setminus \{S\}$.

Begründung in der nächsten Vorlesung (nur FSK).