

## Grammatiken und die Chomsky-Hierarchie

Prof. Dr. David Sabel

LFE Theoretische Informatik



Letzte Änderung der Folien: 3. Mai 2022

## Formale Sprachen darstellen (2)

Anforderungen:

- **Endliche** Beschreibung
- Sprache selbst muss aber auch unendlich viele Objekte erlauben

Zwei wesentliche solchen Formalismen sind

- Grammatiken
- Automaten

## Formale Sprachen darstellen

- Sei  $\Sigma$  ein Alphabet.
- Eine **Sprache über  $\Sigma$**  ist eine Teilmenge von  $\Sigma^*$ .
- Z.B. 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

## Grammatiken

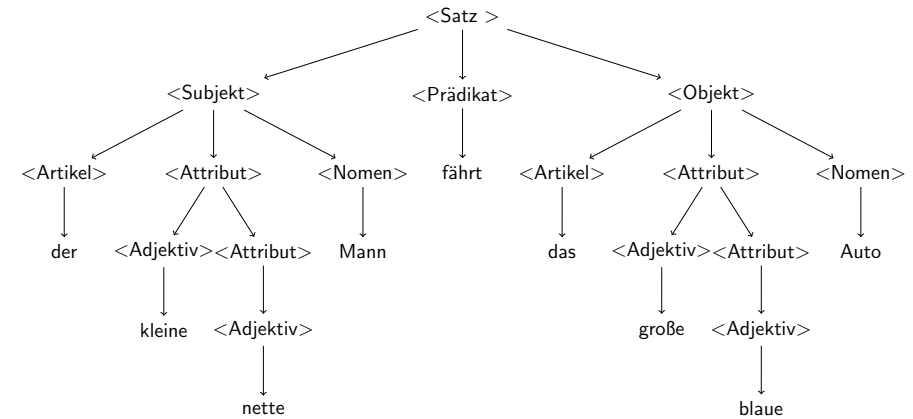
Grammatik für einen sehr kleinen Teil der deutschen Sprache:

<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> → Mann  
 <Nomen> → Auto  
 <Prädikat> → fährt  
 <Prädikat> → liebt

## Grammatiken

- 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 obige Grammatik abgeleitet werden

## Syntaxbaum zum Beispiel



## Definition einer Grammatik

### Definition (Grammatik)

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

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

Manchmal genügt es,  $P$  alleine zu notieren  
(wenn klar ist, was Variablen, Zeichen und Startsymbol sind)

## Beispiel für eine Grammatik

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

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

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

$$P = \{E \rightarrow M,$$

$$E \rightarrow E + M,$$

$$M \rightarrow Z,$$

$$M \rightarrow M * Z,$$

$$Z \rightarrow 1,$$

$$Z \rightarrow 2,$$

$$Z \rightarrow (E)\}$$

## Ableitung

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

### Ableitungsschritt $\Rightarrow_G$

Für Satzformen  $u, v$  (d.h. Worte aus  $(V \cup \Sigma)^*$ ) sagen wir:

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

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

- Wenn  $G$  klar ist, schreiben wir  $u \Rightarrow v$  statt  $u \Rightarrow_G v$
- $\Rightarrow_G^*$  sei die reflexiv-transitive Hülle von  $\Rightarrow_G$

### Ableitung

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

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

$$P = \left\{ \begin{array}{l} E \rightarrow M, \quad E \rightarrow E + M, \quad M \rightarrow Z, \quad M \rightarrow M * Z, \\ Z \rightarrow 1, \quad Z \rightarrow 2, \quad Z \rightarrow (E) \end{array} \right\}$$

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

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

## Beispiel: Ableitungen sind nicht eindeutig

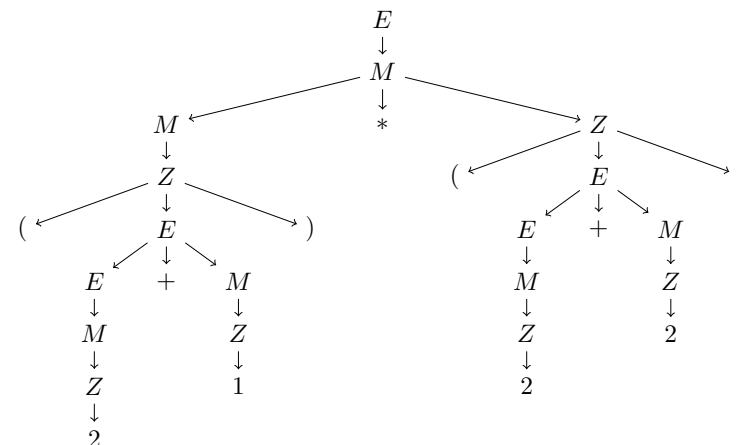
Ableitung von letzter Folie (keine Linksableitung):

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

Linksableitung: ersetzt immer das linkeste Nichtterminal

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

## Syntaxbaum (zu beiden Ableitungen)



## Nichtdeterminismus beim Ableiten

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

Quellen des Nichtdeterminismus:

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

Aber: Es gibt **nur endliche viele**  $v_i$  für jeden Schritt!

## Beispiele

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

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

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

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

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

$$\begin{array}{ccccccccc} S' & \Longrightarrow & aS' & \Longrightarrow & aaS' & \Longrightarrow & aaaS' & \Longrightarrow & aaaaS' & \Longrightarrow & \dots \\ \Downarrow & & \Downarrow & & \Downarrow & & \Downarrow & & \Downarrow & & \\ \bullet & & b & & ab & & aab & & aaab & & aaaab \end{array}$$

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

## Erzeugte Sprache

### Erzeugte Sprache einer Grammatik

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

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

## Die Chomsky-Hierarchie

Noam Chomsky teilte die Grammatiken in Typen 0 bis 3:

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

### $G$ ist vom Typ 0

$G$  ist automatisch vom Typ 0.

### $G$ ist vom Typ 1 (kontextsensitive Grammatik), wenn ...

für alle  $(\ell \rightarrow r) \in P$ :  $|\ell| \leq |r|$ .

### $G$ ist vom Typ 2 (kontextfreie Grammatik), wenn ...

$G$  ist vom Typ 1 und für alle  $(\ell \rightarrow r) \in P$  gilt:  $\ell = A \in V$

### $G$ ist vom Typ 3 (reguläre Grammatik), wenn ...

$G$  ist vom Typ 2 und für alle  $(A \rightarrow r) \in P$  gilt:  $r = a$  oder  $r = aA'$  für  $a \in \Sigma, A' \in V$  (die rechten Seiten sind Worte aus  $(\Sigma \cup (\Sigma V))^*$ )

## Typ $i$ -Sprachen

### Definition

Für  $i = 0, 1, 2, 3$  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 stets der größtmögliche Typ gemeint.

Bemerkung: Die Definition erlaubt Aussagen der Form:

Typ  $i + k$ -Sprachen sind eine Teilmenge der Typ  $i$ -Sprachen, da jede Typ  $i + k$ -Grammatik auch eine Typ  $i$ -Grammatik ist.

## Beispiele

- $G_1 = (\{S\}, \{a, b\}, \{S \rightarrow aS, S \rightarrow b\}, S)$  ist regulär (Typ 3)
- $G_2 = (\{E, M, Z\}, \{+, *, 1, 2, (, )\}, P, E)$  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 kontextfrei (Typ 2)
- $G_3 = (\{S, B, C\}, \{a, b, c\}, P, S)$  mit  
 $P = \{S \rightarrow aSBC, S \rightarrow aBC, CB \rightarrow BC, aB \rightarrow ab, bB \rightarrow bb, bC \rightarrow bc, cC \rightarrow cc\}$  ist kontextsensitiv (Typ 1)  
Beachte  $L(G_3) = \{a^n b^n c^n \mid n \in \mathbb{N}_{>0}\}$
- $G_4 = (\{S, T, A, B, \$\}, \{a, b\}, P, S)$  mit  
 $P = \{S \rightarrow \$T\$, T \rightarrow aAT, T \rightarrow bBT, T \rightarrow \varepsilon, \$a \rightarrow a\$, \$b \rightarrow b\$, Aa \rightarrow aA, Ab \rightarrow bA, Ba \rightarrow aB, Bb \rightarrow bB, A\$ \rightarrow \$a, B\$ \rightarrow \$b, \$\$ \rightarrow \varepsilon\}$  ist vom Typ 0