

7. Formale Sprachen und Grammatiken

Computer verwenden zur Verarbeitung von Daten und Informationen künstliche, *formale Sprachen* (Maschinenspr., Assemblersprachen, Programmierspr., Datenbankspr., Wissensrepräsentationsspr., ...)

Vorteile gegenüber natürlichen Sprachen:

- exakte Definition der zulässigen Ausdrücke und ihrer Bedeutung
- keine Kontextabhängigkeit der Bedeutung
- dadurch erheblich leichtere Verarbeitung

Eine Sprache L besteht aus Wörtern über einem zugrundeliegenden Alphabet Σ , es gilt also: L Teilmenge von Σ^* .

Die Syntax einer Sprache legt fest, welche Ausdrücke zur Sprache gehören.

Die Semantik legt fest, was die Ausdrücke bedeuten.

Endliche Sprachen können durch Aufzählen ihrer Elemente definiert werden, für unendliche Sprachen braucht man endliche Beschreibung.

Formale Sprachen

Formale Sprache

Alphabet: geordneter Zeichenvorrat

Formale Sprache: Menge von Zeichenketten, die aus den Symbolen eines beliebigen Alphabets aufgebaut sind.

Zeichenkette: endliche Folge von Symbolen, die ohne Zwischenraum hintereinander geschrieben werden (Verkettung)

Wohlgeformtheit

Definition der *zulässigen Zeichenketten* einer Sprache durch induktive Definition

induktive Definition:

1. definitorische Anfangsbestimmung
2. generierende Bestimmung
3. Abschlußbestimmung

Beispiel Aussagenlogik

Aussagenlogik - Syntax

Aussagenlogische Variable und die logischen Konstanten *wahr* bzw. *falsch* sind *atomare Formeln*. Es gilt:

- (1) Alle atomaren Formeln sind Formeln.
- (2a) Für alle Formeln F und G sind $(F \wedge G)$ und $(F \vee G)$ Formeln.
- (2b) Für jede Formel F ist $\neg F$ eine Formel.
- (3) Nichts ist eine Formel, was nicht als solche durch (1) und (2) bestimmt ist.

Grammatiken

Formale Sprachen können durch Grammatiken beschrieben werden.

Eine Grammatik $G = (T, N, P, S)$ besteht aus folgenden Komponenten:

- einer endlichen Menge N von Variablen (Nichtterminalsymbolen),
- dem endlichen Terminalalphabet T,
- einer endlichen Menge P von Regeln $\alpha \rightarrow \beta$, wobei α und β aus N und T gebildet sind,
- einer Startvariablen S aus N.

Die durch G erzeugte Sprache $L(G)$ wird wie folgt definiert:

Die Zeichenkette v heißt direkt ableitbar aus der Zeichenkette u

(Notation: $u \Rightarrow v$) genau dann wenn $u = xyz$, $v = xy'z$ und $y \rightarrow y' \in P$.

\Rightarrow^* bezeichnet die reflexive und transitive Hülle von \Rightarrow .

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

Ableitungsbaum

1. Die Wurzel des Ableitungsbaums ist mit dem Startsymbol markiert.
2. Ist ein Knoten mit einem Terminalsymbol markiert, dann ist er ein Blatt und hat keinen Sohn.
3. Ist ein Knoten mit einem Nichtterminalsymbol X markiert, und hat er n Söhne, die von links nach rechts mit X_1, X_2, \dots, X_n markiert sind, so enthält die Grammatik die Regel

$$X \rightarrow X_1, X_2, \dots, X_n$$

Die Folge der Blätter von links nach rechts gelesen ergibt ein Wort der Sprache.

Chomsky Grammatik

Induktive Definition. Aufgeführt sind die in der jeweils höheren Sprache zusätzlich gültigen Regeln.

Reguläre Sprachen (Typ 3)

Produktionsregeln der Form

$$\langle X \rangle \rightarrow a|a\langle Y \rangle \quad (\text{alternativ } \langle X \rangle \rightarrow \langle Y \rangle a)$$

Bemerkung: Das Nichtterminalsymbol steht in allen Regeln entweder immer rechts oder immer links

Kontextfreie Sprachen (Typ 2)

Produktionsregeln der Form

$$\langle X \rangle \rightarrow u\langle Y \rangle v$$

Kontextsensitive Sprachen (Typ 1)

Produktionsregeln der Form

$u\langle X \rangle v \rightarrow u\langle Y \rangle v$ wobei die Länge der Zeichenkette auf der linken Seite nicht größer als die der rechten Seite sein darf.

Unbeschränkte Sprachen (Typ 0)

Typen von Grammatiken

eine Grammatik heißt

- kontextsensitiv falls für alle Regeln $w_1 \rightarrow w_2$ gilt: $|w_1| \leq |w_2|$
($|w|$: Anzahl der Variablen und Terminalzeichen von w)
 - kontextfrei falls für alle Regeln $w_1 \rightarrow w_2$ gilt: w_1 aus V
 - regulär falls zusätzlich gilt: w_2 Terminalzeichen oder Terminalzeichen gefolgt von Variable
- durch Grammatiktypen erzeugbare Sprachen bilden echte Teilmengen



Chomsky-Hierarchie

Eine Grammatik für die natürliche Sprache

P:

<Satz>	->	<Subjekt> <Prädikat> <Objekt>
<Subjekt>	->	<Artikel> <Attribut> <Substantiv>
<Artikel>	->	ϵ
<Artikel>	->	der
<Artikel>	->	die
<Artikel>	->	das
<Attribut>	->	ϵ
<Attribut>	->	<Adjektiv>
<Attribut>	->	<Adjektiv> <Attribut>
<Adjektiv>	->	kleine
<Adjektiv>	->	bissige
<Adjektiv>	->	große
<Substantiv>	->	Hund
<Substantiv>	->	Katze
<Prädikat>	->	jagt
<Objekt>	->	<Artikel> <Attribut> <Substantiv>

N: {<Satz>, <Subjekt>, <Artikel>, <Attribut>, <Adjektiv>, <Substantiv>, <Prädikat>, <Objekt>}

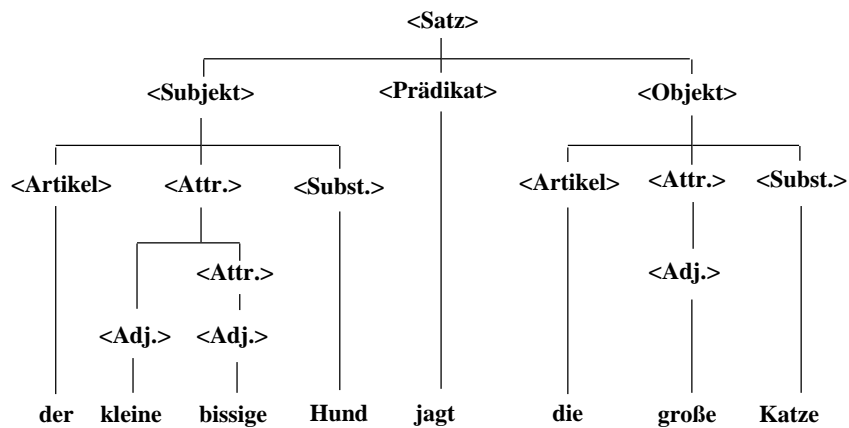
T: {der, die, das, kleine, bissige, große, Hund, Katze, jagt}

S: <Satz>

Eine Ableitung

<Satz> => <Subjekt> <Prädikat> <Objekt>
 => <Artikel> <Attribut> <Substantiv> <Prädikat> <Objekt>
 => der <Artikel> <Substantiv> <Prädikat> <Objekt>
 => der <Adjektiv> <Attribut> <Substantiv> <Prädikat> <Objekt>
 => der kleine <Artikel> <Substantiv> <Prädikat> <Objekt>
 => der kleine <Adjektiv> <Substantiv> <Prädikat> <Objekt>
 => der kleine bissige <Substantiv> <Prädikat> <Objekt>
 => der kleine bissige Hund <Prädikat> <Objekt>
 => der kleine bissige Hund jagt <Objekt>
 => der kleine bissige Hund jagt <Artikel> <Attribut> <Substantiv>
 => der kleine bissige Hund jagt die <Artikel> <Substantiv>
 => der kleine bissige Hund jagt die <Adjektiv> <Substantiv>
 => der kleine bissige Hund jagt die große <Substantiv>
 => der kleine bissige Hund jagt die große Katze

Ein Ableitungsbaum



Backus - Naur - Formen.

Vereinfachte Darstellung kontextfreier Grammatiken durch BNF oder erweiterte BNF (EBNF) sowie durch Syntaxdiagramme (s.u.)

BNF Notation:

- Nichtterminale Symbole durch spitze Klammern $\langle \dots \rangle$ gekennzeichnet.
- **Metasymbole** sind
 - ::= ist definiert als
 - | trennt Alternativen.

EBNF (erweiterte BNF) Notation:

Im Vergleich zu BNF reichere Metasyntax und damit einfachere und besser lesbare Schreibweise.

- **Metasymbole** sind
 - (...) genau eine Alternative kann stehen.
 - [...] Option: Inhalt der Klammer steht einmal oder keinmal.
 - {...} Wiederholung: Inhalt der Klammer steht **mindestens** einmal, d.h.
 $\{A\}=A|AA|AAA|\dots$

Beispiel: $\{(ab)\}$ erzeugt die Menge der Zeichenketten $\{a,b,aa,ab,ba,bb,\dots\}$ denn der Inhalt der Klammer ist der Ausdruck $A=(a|b)$, dessen Auswertung entweder a oder b ergibt.

BNF-Beispiel

\emptyset Bezeichner, also Namen für Variablen, Prozeduren oder Konstanten, können wie folgt beschrieben werden:

```
<Bezeichner> ::= <Buchstabe>|
               <Bezeichner><Buchstabe>|
               <Bezeichner><Ziffer>

<Buchstabe> ::= A|B|C|D|E|F|G|H|I|J|K|L|M|N|O|P|Q|R|S|T|U|V|W|X|Y|Z
              a|b|c|d|e|f|g|h|i|j|k|l|m|n|o|p|q|r|s|t|u|v|w|x|y|z

<Ziffer> ::= 1|2|3|4|5|6|7|8|9|0
```

Die Wertebereiche der Metavariablen sind also:

```
<Ziffer> : {1, ..., 9, 0}
<Buchstabe> : {A, ..., Z, a, ..., z}
<Bezeichner> : alle endlichen, nichtleeren Zeichenketten über der Menge
               {1, ..., 9, 0}  $\cup$  {A, ..., Z, a, ..., z}, die mit einem Buchstaben
               beginnen, z.B.: "A", "Hans", "A12", "ergo", etc.
```

\emptyset Zulässige Variablenvereinbarungen, also zum Beispiel Zeichenketten der Form "VAR x : INTEGER;" , können mit Hilfe der Backus-Naur-Form so spezifiziert werden:

```
<VarVereinb> ::= VAR <Vereinbarungen>
<Vereinbarungen> ::= <Vereinbarung><Vereinbarungen>
<Vereinbarung> ::= <Bezeichnerliste><Typ>;
<Bezeichnerliste> ::= <Bezeichner><Bezeichnerliste>,<Bezeichner>
<Typ> ::= INTEGER | REAL | BOOLEAN |
```

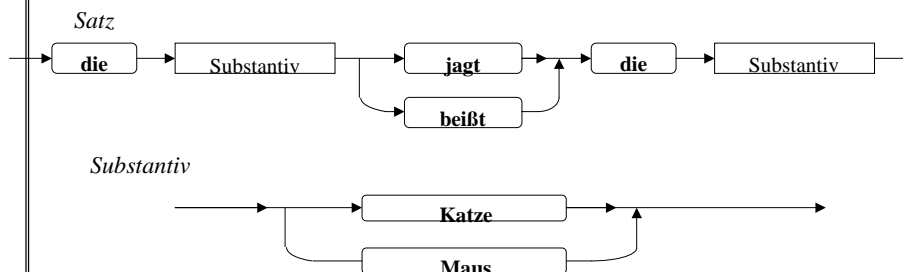
Erweiterte Backus-Naur-Form (EBNF)

Zusätzliche Sprachkonstrukte sind u.a.:

- ∅ Ein Metaausdruck in eckigen Klammern, also ein Ausdruck der Gestalt $[\text{Metaausdruck}]$, bezeichnet Zeichenketten, die optional sind und somit auch weggelassen werden können.
- ∅ Ein Metaausdruck in geschweiften Klammern, also ein Ausdruck der Form $\{ \text{Metaausdruck} \}$, beschreibt Zeichenketten, die beliebig oft wiederholt werden dürfen (0, 1, 2, ...-mal). Dieses Sprachkonstrukt verkörpert also eine Art Iteration.
- ∅ Gebräuchlich ist ferner ein Ausdruck der Form $\{ \text{Metaausdruck} \}_a^b$, mit dem festgelegt wird, daß die durch Metaausdruck definierten Zeichenketten mindestens a-mal, höchstens b-mal wiederholt werden müssen bzw. dürfen. Dabei wird $0 \leq a \leq b$ vorausgesetzt, und es ist auch $b = *$ erlaubt, um auszudrücken, daß die Anzahl der Iterationen nach oben unbeschränkt ist.

Syntaxdiagramme

- ∅ Die Syntax einer (kontextfreien) Programmiersprache läßt sich durch sog. Syntaxdiagramme beschreiben. Sie erlauben die Ableitung aller syntaktisch korrekten Sätze der Sprache
- ∅ Prinzipielle Anwendung
 - ⇒ Rechteck entspricht syntaktischer Variable (Nicht-Terminalsymbol) in BNF
 - ⇒ Rechteck mit <Name> wird bei Durchlauf ersetzt durch Syntaxdiagramm mit <Name>
 - ⇒ Langrunde/Kreis mit Symbol produziert (gibt aus) bei Durchlauf das entsprechende Symbol
- ∅ Beispiel



Umwandlung der Darstellungsform

- θ Kein prinzipieller Unterschied hinsichtlich der Ausdrucksmächtigkeit von Syntaxdiagrammen und EBNF
- θ Automatische Umwandlung möglich

- θ Vorteile der Syntaxdiagramme
 - ⇒ hohe Anschaulichkeit
 - ⇒ geringe Anzahl der Regeln (Diagramme); typischerweise sind mehrere Produktionen der EBNF in einem Syntaxdiagramm zusammengefasst

- θ Vorteile einer BNF-Grammatik
 - ⇒ leichte Erstellung und Änderung
 - ⇒ erzwingt eine klare Struktur
 - ⇒ direkte Verwendung in Übersetzer-erzeugenden Systemen möglich