

# **Kapitel 3.3 Syntaxdefinitionen**

# Syntax

**Syntax:** Festlegung des Satzbaus.

Beispiele syntaktisch falscher deutscher Sätze:

*Kai liest eine Buch. Buch lesen Kai. Kai pr1 &. Kai liest ein Buch, weil ihr ist langweilig.*

Beispiele syntaktisch falscher OCAML-Phrasen (Frasen?):

```
let let x = 3 in 2;;
```

```
let if = 1;;
```

```
2 + * 3;;
```

```
2 : 3;;
```

# Semantik

**Semantik:** Festlegung der Bedeutung eines Satzes.

Semantische Fragen im Deutschen: Worauf bezieht sich ein Relativpronomen? Welchen Einfluss haben Fragepartikel wie “eigentlich”, “denn”? Wann verwendet man welche Zeitform?

Semantische Fragen bei OCAML: Was ist der Wert von Ausdrücken wie

```
let x = 1 in let y = x in let x = 2 in y;;  
let x = 1./0. in 2;;  
let rec f x = f x in if true then 1 else f 0;;
```

*Grundfrage der Semantik* von Programmiersprachen: Welche Wirkung hat ein syntaktisch korrektes Programm?

Aus *historischen Gründen* werden Fragen der Typüberprüfung auch der Semantik zugerechnet.

# Formale Syntax

- Ein *Alphabet* ist eine endliche Menge, deren Elemente *Symbole* genannt werden.
- Eine *Zeichenkette* (auch *Wort*) über einem Alphabet  $\Sigma$  ist eine endliche Folge von Elementen  $\sigma_1, \dots, \sigma_n$  von  $\Sigma$ , wobei  $n \geq 0$ . Man schreibt ein Wort als  $\sigma_1\sigma_2 \dots \sigma_n$ . Der Fall  $n = 0$  bezeichnet das *leere Wort* geschrieben  $\varepsilon$ .
- Die Menge aller Wörter über  $\Sigma$  wird mit  $\Sigma^*$  bezeichnet. Zu zwei Wörtern  $w = \sigma_1 \dots \sigma_n$  und  $w' = \sigma'_1 \dots \sigma'_m$  bildet man die *Verkettung* (Konkatenation)  $ww' = \sigma_1 \dots \sigma_n \sigma'_1 \dots \sigma'_m$ . Es ist  $\varepsilon w = w \varepsilon = w$  und  $(ww')w'' = w(w'w'')$ .
- Eine *formale Sprache* ist eine Teilmenge von  $\Sigma^*$ .

# Beispiele

$\Sigma = \{a, b\}$ .

Wörter über  $\Sigma$ : *aaba, baab, bababa, baba,  $\varepsilon$ , bbbb*.

Sprachen über  $\Sigma$ :  $\emptyset, \{a, b\}, \{a^n b^n \mid n \in \mathbb{N}\}$ .

$\Sigma = \{0, 1, \dots, 9, e, -, +, ., E\}$ .

Wörter über  $\Sigma$ : *-1E98, --2e--, 32.e*

Sprachen über  $\Sigma$ : Syntaktisch korrekte `float` Konstanten,  
 $\{e^n \mid n \text{ gerade}\}$ .

$\Sigma = \{0, \dots, 9, \text{if, then, let, } \dots \}$

Sprache über  $\Sigma$ : alle syntaktisch korrekten OCAML Phrasen.

# Beispiel: OCAML-Bezeichner

- Bezeichner sind Zeichenketten über dem Alphabet  $\Sigma = \{A, \dots, Z, a, \dots, z, 0, \dots, 9, ', \_ \}$ . Bezeichner müssen mit einem Kleinbuchstaben oder dem Zeichen `_` beginnen.
- Ein *Buchstabe* ist ein *Kleinbuchstabe* oder ein *Großbuchstabe*
- Ein *Kleinbuchstabe* ist ein Zeichen `a, \dots, z`.
- Ein *Großbuchstabe* ist ein Zeichen `A, \dots, Z`.
- Eine *Ziffer* ist ein Zeichen `0, \dots, 9`.

# Formale Definition als BNF-Grammatik

BNF = *Backus-Naur Form*

$\langle \text{Bezeichner} \rangle ::= \{ \langle \text{KlBuchst} \rangle \mid - \} \{ \langle \text{Buchst} \rangle \mid \langle \text{Ziffer} \rangle \mid ' \mid - \}^*$

$\langle \text{Buchst} \rangle ::= \langle \text{KlBuchst} \rangle \mid \langle \text{GrBuchst} \rangle$

$\langle \text{KlBuchst} \rangle ::= a \mid b \mid c \mid d \mid e \mid f \mid g \mid h \mid i \mid j \mid k \mid l \mid m \mid n \mid o$   
 $\mid p \mid q \mid r \mid s \mid t \mid u \mid v \mid w \mid x \mid y \mid z$

$\langle \text{GrBuchst} \rangle ::= A \mid B \mid C \mid D \mid E \mid F \mid G \mid H \mid I \mid J \mid K \mid L \mid M \mid N \mid O$   
 $\mid P \mid Q \mid R \mid S \mid T \mid U \mid V \mid W \mid X \mid Y \mid Z$

$\langle \text{Ziffer} \rangle ::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

# float-Literale

*Literal* = Konstante in einer Programmiersprache.

$$\langle \text{float-Literal} \rangle ::= [\langle \text{Vorzeichen} \rangle] \langle \text{Mantisse} \rangle [\langle \text{Exponent} \rangle]$$
$$\langle \text{Vorzeichen} \rangle ::= - \mid +$$
$$\langle \text{Mantisse} \rangle ::= \{ \langle \text{Ziffer} \rangle \}^+ [ \cdot \{ \langle \text{Ziffer} \rangle \}^* ]$$
$$\langle \text{Exponent} \rangle ::= \{ e \mid E \} [\langle \text{Vorzeichen} \rangle] \{ \langle \text{Ziffer} \rangle \}^+$$

Zusätzliche *Kontextbedingung*: Entweder ein Dezimalpunkt, oder ein e oder ein E muss vorhanden sein.

**Übung:** man verbessere die BNF Darstellung so, dass diese Kontextbedingung wegfallen kann.

# Syntax der BNF

Eine BNF-Grammatik ist ein Quadrupel  $G = (\Sigma, V, S, P)$ .

- $\Sigma$  ist die Menge der Terminalsymbole, meist in Courier gesetzt.
- $V$  ist die Menge der Nichtterminalsymbole, meist in spitze Klammern gesetzt. Im Beispiel:  
$$V = \{\langle \text{float-Literal} \rangle, \langle \text{Mantisse} \rangle, \langle \text{Vorzeichen} \rangle, \langle \text{Exponent} \rangle\}.$$
- $S \in V$  ist ein ausgezeichnetes Nichtterminalsymbol, das *Startsymbol*.  
Im Beispiel:  $S = \langle \text{float-Literal} \rangle$ .
- $P$  ist eine endliche Menge von *Produktionen* der Form  $X ::= \delta$ , wobei  $\delta$  eine *BNF-Satzform*, s.u., ist.

# BNF-Satzformen

- Jedes Symbol in  $V \cup \Sigma$  ist eine BNF Satzform.
- Sind  $\gamma_1, \dots, \gamma_n$  BNF-Satzformen, so auch  $\gamma_1 \mid \dots \mid \gamma_n$  (*Auswahl*).
- Sind  $\gamma_1, \dots, \gamma_n$  BNF-Satzformen, so auch  $\gamma_1 \dots \gamma_n$  (*Verkettung*).
- Ist  $\gamma$  eine BNF Satzform, so auch  $\{\gamma\}$  (*Klammerung*).
- Ist  $\gamma$  eine BNF Satzform, so auch  $\{\gamma\}^*$  (*Iteration*).
- Ist  $\gamma$  eine BNF Satzform, so auch  $\{\gamma\}^+$  (*nichtleere Iteration*).
- Ist  $\gamma$  eine BNF Satzform, so auch  $[\gamma]$  (*Option*).

Eine Satzform der Gestalt  $\gamma_1 \mid \dots \mid \gamma_n$  muss immer geklammert werden, es sei denn, sie tritt unmittelbar als rechte Seite einer Produktion auf.

# Semantik der BNF

Gegeben sei eine EBNF Grammatik  $G = (\Sigma, V, S, P)$ . Jedem Nichtterminalsymbol  $X$  in  $V$  und jeder BNF Satzform  $\delta$  wird eine Sprache  $\mathcal{L}(X)$ , bzw.  $\mathcal{L}(\delta)$  wie folgt zugeordnet.

- Falls  $x \in \Sigma$ , so ist  $\mathcal{L}(x) = \{x\}$ .
- Falls  $X \in V$ , so ist  $\mathcal{L}(X) = \bigcup_{\delta \in P, X ::= \delta} \mathcal{L}(\delta)$ .
- $\mathcal{L}(\gamma_1 \mid \dots \mid \gamma_n) = \mathcal{L}(\gamma_1) \cup \dots \cup \mathcal{L}(\gamma_n)$ .
- $\mathcal{L}(\gamma_1 \dots \gamma_n) = \{w_1 w_2 \dots w_n \mid w_1 \in \mathcal{L}(\gamma_1), \dots, w_n \in \mathcal{L}(\gamma_n)\}$ .
- $\mathcal{L}(\{\gamma\}^*) = \{w_1 w_2 \dots w_n \mid n \in \mathbb{N}, w_i \in \mathcal{L}(\gamma)\}$
- $\mathcal{L}(\{\gamma\}^+) = \{w_1 w_2 \dots w_n \mid n \in \mathbb{N}, n > 0, w_i \in \mathcal{L}(\gamma)\}$
- $\mathcal{L}([\gamma]) = \mathcal{L}(\gamma) \cup \{\varepsilon\}$

# Semantik durch Ersetzung

**Satz:** Ein Wort  $w$  ist in  $\mathcal{L}(S)$ , wenn man es aus  $S$  durch die folgenden Ersetzungsoperationen erhalten kann:

- Falls  $X ::= \gamma_1 \mid \cdots \mid \gamma_n$  eine Produktion ist, so darf man ein Vorkommen von  $X$  durch eines der  $\gamma_i$  ersetzen.
- Ein Vorkommen von  $\{\gamma_1 \mid \cdots \mid \gamma_n\}$  darf man durch eines der  $\gamma_i$  ersetzen.
- Ein Vorkommen von  $\{\gamma\}^*$  darf man durch  $\overbrace{\{\gamma\}\{\gamma\}\dots\{\gamma\}}^{n\text{-mal}}$  mit  $n \geq 0$  ersetzen.
- Ein Vorkommen von  $\{\gamma\}^+$  darf man durch  $\overbrace{\{\gamma\}\{\gamma\}\dots\{\gamma\}}^{n\text{-mal}}$  mit  $n > 0$  ersetzen.
- Ein Vorkommen von  $\{\gamma\}$  darf man durch  $\gamma$  ersetzen, wenn  $\gamma$  nicht von der Form  $\gamma_1 \mid \cdots \mid \gamma_n$  ist.
- Ein Vorkommen von  $[\gamma]$  darf man durch  $\{\gamma\}$  ersetzen, oder ersatzlos streichen.

# Beispiel

$\langle \text{float-Literal} \rangle \rightarrow [\langle \text{Vorzeichen} \rangle] \langle \text{Mantisse} \rangle [\langle \text{Exponent} \rangle] \rightarrow$   
 $\langle \text{Mantisse} \rangle [\langle \text{Exponent} \rangle] \rightarrow \{ \langle \text{Ziffer} \rangle \}^+ [ \cdot \{ \langle \text{Ziffer} \rangle \}^* ] [\langle \text{Exponent} \rangle] \rightarrow$   
 $\{ \langle \text{Ziffer} \rangle \}^+ \cdot \{ \langle \text{Ziffer} \rangle \}^* [\langle \text{Exponent} \rangle] \rightarrow$   
 $\{ \langle \text{Ziffer} \rangle \}^+ \cdot \{ \langle \text{Ziffer} \rangle \} \{ \langle \text{Ziffer} \rangle \} [\langle \text{Exponent} \rangle] \rightarrow$   
 $\{ \langle \text{Ziffer} \rangle \} \cdot \{ \langle \text{Ziffer} \rangle \} \{ \langle \text{Ziffer} \rangle \} [\langle \text{Exponent} \rangle] \rightarrow 2.71 \langle \text{Exponent} \rangle \rightarrow$   
 $2.71 \{ e \mid E \} [\langle \text{Vorzeichen} \rangle] \{ \langle \text{Ziffer} \rangle \}^+ \rightarrow$   
 $2.71E \langle \text{Vorzeichen} \rangle \{ \langle \text{Ziffer} \rangle \} \{ \langle \text{Ziffer} \rangle \} \{ \langle \text{Ziffer} \rangle \} \rightarrow$   
 $2.71E \{ - \mid + \} 001 \rightarrow 2.71E-001$

Also  $2.71E-001 \in \mathcal{L}(\langle \text{float-Literal} \rangle)$ .

# Kontextbedingungen

Manche syntaktische Bedingungen lassen sich mit BNF nur schwer oder gar nicht formulieren.

Man gibt daher manchmal zusätzliche *Kontextbedingungen* an, denen die syntaktisch korrekten Wörter zusätzlich genügen müssen.

Beispiele:

- Bezeichner dürfen nicht zu den Schlüsselwörtern gehören wie z.B. `let`, `if`, etc.
- float-Literale müssen `.`, `e`, oder `E` enthalten.

Andere Bedingungen, wie korrekte Typisierung oder rechtzeitige Definition von Bezeichnern werden, wie schon gesagt, der Semantik zugerechnet.

# Varianten

Häufig wird statt  $\{\gamma\}^*$  nur  $\{\gamma\}$  geschrieben. Für die Klammerung verwendet man dann runde Klammern.

Die spitzen Klammern zur Kennzeichnung der Nichtterminalsymbole werden oft weggelassen.

Steht kein Courier Zeichensatz zur Verfügung, so schließt man die Terminalsymbole in “Anführungszeichen” ein.

# Ableitungsbäume

Ableitungen in einer BNF lassen sich grafisch durch *Ableitungsbäume* darstellen.

Beispiele von Ableitungsbäumen finden Sie in [Kröger].

Diese Ableitungsbäume sind für die Festlegung der Semantik von Bedeutung.

Ein *Parser* berechnet zu einem vorgegebenen Wort einen Ableitungsbaum, falls das Wort in  $\mathcal{L}(S)$  ist, und erzeugt eine Fehlermeldung, falls nicht.

Diese Aufgabe bezeichnet man als *Syntaxanalyse*.

# Syntaxdiagramme

Man kann BNF Syntaxdefinitionen auch grafisch durch *Syntaxdiagramme* darstellen.

Terminalzeichen werden als Kreise dargestellt.

Für jedes Nichtterminalzeichen  $X$  (als Kasten dargestellt) ein Diagramm.

Eine Zeichenreihe aus  $\mathcal{L}(X)$  erhält man, indem man im Diagramm für  $X$  einen beliebigen Weg vom Anfang zum Ausgang wählt, dabei jedes auftretende Terminalzeichen aufsammelt und jedes auftretende Nichtterminalzeichen  $Y$  durch ein entsprechend gefundenes Wort aus  $\mathcal{L}(Y)$  ersetzt.

Siehe [Kröger] für Beispiele von Syntaxdiagrammen.

# OCAML-Sprachdefinition (Fragment)

$\langle \text{toplevel-command} \rangle ::= \text{let } \langle \text{let-binding} \rangle \{ \text{and } \langle \text{let-binding} \rangle \}^* ; ;$   
|  $\text{let rec } \langle \text{let-binding} \rangle \{ \text{and } \langle \text{let-binding} \rangle \}^* ; ;$   
|  $\langle \text{expression} \rangle ; ;$

$\langle \text{expression} \rangle ::= \langle \text{constant} \rangle$   
|  $( \langle \text{expression} \rangle )$   
|  $\langle \text{expression} \rangle \langle \text{infix-op} \rangle \langle \text{expression} \rangle$   
|  $\text{if } \langle \text{expression} \rangle \text{ then } \langle \text{expression} \rangle \text{ else } \langle \text{expression} \rangle$   
|  $\text{let } \langle \text{let-binding} \rangle \{ \text{and } \langle \text{let-binding} \rangle \}^* \text{ in } \langle \text{expression} \rangle$   
|  $\text{let rec } \langle \text{let-binding} \rangle \{ \text{and } \langle \text{let-binding} \rangle \}^* \text{ in } \langle \text{expression} \rangle$   
|  $\text{function } \langle \text{ident} \rangle \text{ -> } \langle \text{expression} \rangle$   
|  $\{ \langle \text{expression} \rangle \}^+$   
|  $\langle \text{expression} \rangle , \langle \text{expression} \rangle$

$\langle \text{infix-op} \rangle ::= + \mid - \mid * \mid / \mid +. \mid -. \mid *. \mid /. \mid \&\& \mid || \mid \text{mod}$   
|  $< \mid > \mid <> \mid <= \mid >= \mid <>$

$\langle \text{let-binding} \rangle ::= \{ \langle \text{ident} \rangle \}^+ = \langle \text{expression} \rangle$

# Parsergeneratoren

Erinnerung: *Parser* :  $\Sigma^* \rightarrow \text{Ableitungsbäume} \cup \text{Fehlermeldungen}$ .

Ein *Parsergenerator* erzeugt aus einer BNF-Grammatik automatisch einen Parser.

Der bekannteste Parsergenerator heißt “yacc” (*yet another compiler-compiler*). Er erzeugt aus einer BNF-Grammatik einen in der Programmiersprache C geschriebenen Parser.

Für OCAML gibt es “ocamlyacc”. Es wird ein OCAML-Programm erzeugt.

# Geschichte

Grammatikformalismen, die syntaktisch korrekte Wörter durch einen Erzeugungsprozess definieren (wie die BNF) heißen *generative Grammatiken*.

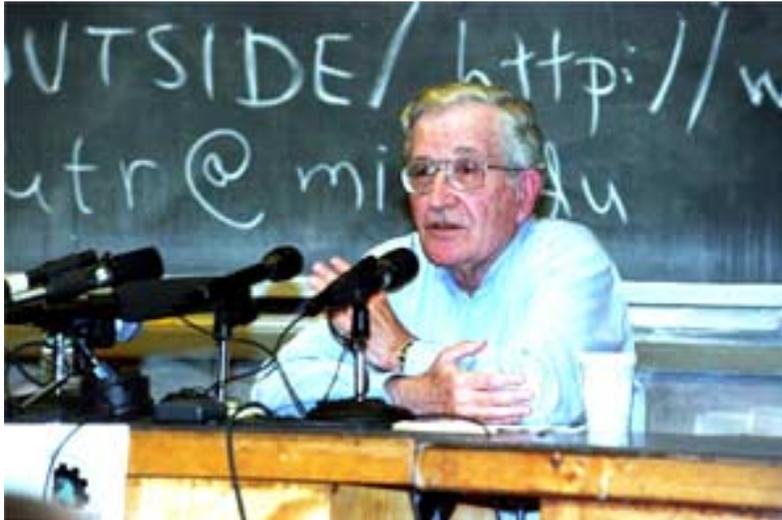
Sie gehen auf den Sprachforscher NOAM CHOMSKY (1928– ) zurück.

Eine *kontextfreie Chomsky-Grammatik* ist eine BNF-Grammatik ohne die Konstrukte  $\{ \}^*$ ,  $[ ]$ ,  $|$ ,  $\{ \}^+$ . Man kann jede BNF-Grammatik durch eine kontextfreie Chomsky-Grammatik simulieren.

Chomsky betrachtet u.a. auch kontextsensitive Grammatiken.

Die Backus-Naur-Form wurde von den Informatikern JOHN BACKUS und PETER NAUR entwickelt, die die Bedeutung von Chomskys generativen Grammatiken für Programmiersprachensyntax erkannten.

# Noam Chomsky



NOAM CHOMSKY

$\langle \text{Bezeichner} \rangle ::= \langle \text{KlBuchst} \rangle \langle \text{Folge} \rangle$

$\langle \text{Bezeichner} \rangle ::= - \langle \text{Folge} \rangle$

$\langle \text{Folge} \rangle ::= \langle \text{Zeichen} \rangle \langle \text{Folge} \rangle$

$\langle \text{Folge} \rangle ::= \varepsilon$

$\langle \text{Zeichen} \rangle ::= ' \quad '$

$\langle \text{Zeichen} \rangle ::= -$

$\langle \text{Zeichen} \rangle ::= \langle \text{Buchst} \rangle$

$\langle \text{Zeichen} \rangle ::= \langle \text{Ziffer} \rangle$

$\langle \text{Buchst} \rangle ::= \langle \text{KlBuchst} \rangle$

$\langle \text{Buchst} \rangle ::= \langle \text{GrBuchst} \rangle$