

1 Syntaktische Analyse/Lexparse

Übersetzung und Phasen

- **Syntaxanalyse** → Struktur des Programms (AST)
- **Kontextanalyse** → Bedeutungsprüfung (Typen, Gültigkeit)
- **Codegenerierung** → Übersetzung in Zielcode

1.1 Compilerstruktur

Ein-Pass-Compiler:

- Führt alle Phasen gleichzeitig aus
- Keine echte Zwischendarstellung (IR)
- Typisch für kleine Sprachen (z. B. Pascal)

Multi-Pass-Compiler:

- Arbeitet mit mehreren Durchgängen über Quelltext/IR
- Datenweitergabe über IR (AST)
- Bessere Modularität und Optimierung

1.2 Syntaxanalyse – Überblick

- **Scanner (Lexer)**: Wandelt Zeichenfolge → Tokenfolge
- **Parser**: Wandelt Tokenfolge → Abstract Syntax Tree (AST)
- Token = atomares Symbol des Quellprogramms

Kontextfreie Grammatik (CFG)

Eine CFG ist ein 4-Tupel (N, T, P, S) mit:

- N : Nichtterminale
- T : Terminale
- P : Produktionen
- S : Startsymbol

1.3 BNF und EBNF

- **BNF**: Grundform zur Definition von Grammatiken
- **EBNF**: Erweiterung mit regulären Ausdrücken, optionalen und wiederholten Konstrukten
- Beispiel:
Expression ::= primary-Expression (operator primary-Expression)*

1.4 Grammatiktransformationen

- **Gruppierung**: Zusammenfassen gleicher LHS
- **Linksauklammern**: Gemeinsame Präfixe auslagern
- **Linksrekursion beseitigen**: $N ::= X|NY \Rightarrow N ::= X(Y)^*$

- **Ersetzung von Nicht-Terminalen:** falls nur eine Regel existiert

1.5 Parsing-Grundlagen

Parsing

Entscheidung, ob Eingabe zur Grammatik gehört und Aufbau des Syntaxbaumes.

- **Top-Down (z. B. rekursiver Abstieg):** Von Startsymbol zu Terminalen
- **Bottom-Up (z. B. Shift/Reduce):** Von Terminalen zur Wurzel

1.6 Top-Down Parsing

- Aufbau des Syntaxbaums von oben nach unten
- Expandiere jeweils das linke Nichtterminal
- LL(k): Grammatik, bei der mit k Lookahead-Tokens eindeutig entschieden werden kann
- **LL(1)** → wichtigster Fall für rekursiven Abstieg

1.7 Rekursiver Abstieg

Rekursiver Abstieg

Jedes Nichtterminal erhält eine Prozedur `parseN()`, deren Aufrufstruktur dem Parsebaum entspricht.

- **accept(t)** prüft aktuelles Token
- **acceptIt()** akzeptiert aktuelles Token ohne Prüfung
- **currentToken:** vom Scanner geliefert

1.8 Starter- und Folgemengen

- T : Menge der Terminals (Tokens, z.B. `id`, `+`)
- N : Menge der Nichtterminale (Variablen, z.B. `Expression`)
- ε : Das leere Wort (Epsilon)
- $\$$: End-of-File Marker (`eof`)

starters[[X]] (First-Menge)

Menge aller Terminals, mit denen ein aus X abgeleiteter String beginnen kann.

Berechnung:

- Ist $X \in T$ (Terminal):
 $starters[[X]] = \{X\}$
- Ist $X \in N$ (Nichtterminal) mit $X \rightarrow Y_1Y_2\dots$:
 - Füge $starters[[Y_1]] \setminus \{\varepsilon\}$ hinzu.
 - Falls $\varepsilon \in starters[[Y_1]]$, füge auch $starters[[Y_2]]$ hinzu (usw.).
- **Epsilon:** Falls $X \rightarrow \varepsilon$ existiert, ist $\varepsilon \in starters[[X]]$.

follow[[A]] (Folgemenge)

Menge aller Terminate, die in einer Satzform unmittelbar rechts von einem Nichtterminal A stehen können.

Regeln:

1. **Startsymbol S :** Enthält immer das Ende-Zeichen (\$).
2. **Rechter Nachbar ($B \rightarrow \alpha A \beta$):**
Alles aus $starters[[\beta]]$ (außer ϵ) kommt zu $follow[[A]]$.
3. **Eltern-Vererbung ($B \rightarrow \alpha A$ oder $\beta \Rightarrow^* \epsilon$):**
Wenn A am Ende steht (oder β wegfallen kann), erbt A alles aus $follow[[B]]$.

Hinweis: Folgemengen enthalten niemals ϵ , können aber \$ enthalten.

LL(1)-Konfliktfreiheit

Eine Grammatik ist LL(1), wenn für jede Produktion mit Alternativen $A \rightarrow \alpha | \beta$ gilt:

- **Disjunkte Starter (Auswahl-Konflikt):**
Die Alternativen dürfen nicht mit demselben Terminal beginnen.
 $starters[[\alpha]] \cap starters[[\beta]] = \emptyset$
- **Disjunkte Folge bei Epsilon (Nullable-Konflikt):**
Falls $\alpha \Rightarrow^* \epsilon$ (d.h. α kann verschwinden), darf die Folgemenge von A keinen gemeinsamen Start mit β haben.
 $starters[[\beta]] \cap follow[[A]] = \emptyset$

1.8.1 Bottom-Up Parsing (Shift & Reduce)

Beim Bottom-Up Parsing (z.B. LR-Parser) wird die Eingabe von links gelesen und versucht, sie auf das Startsymbol zu reduzieren.

Operationen:

- **Shift:** Schiebe das nächste Eingabe-Symbol auf den Stack.
- **Reduce:** Ersetze eine Symbolfolge oben auf dem Stack durch das entsprechende Nichtterminal (gemäß Grammatikregel).

Beispiel-Trace (aus Übung): Grammatik: $S \rightarrow Ac$, $A \rightarrow aA | b$. Eingabe: $aabc$

Stack	Eingabe	Aktion
ϵ	$aabc$	Shift a
a	abc	Shift a
aa	bc	Shift b
ab	c	Reduce ($A \rightarrow b$)
aaA	c	Reduce ($A \rightarrow aA$)
aA	c	Reduce ($A \rightarrow aA$)
A	c	Shift c
Ac	ϵ	Reduce ($S \rightarrow Ac$)
S	ϵ	Akzeptiert

1.9 AST (Abstract Syntax Tree)

- Strukturierte Repräsentation des Programms
- Parser erzeugt AST-Knoten beim rekursiven Abstieg
- Jede Grammatikregel entspricht einer AST-Unterklasse

AST-Aufbau

- **Abstrakte Basisklasse:** AST
- **Subklassen:** Command, Expression, Declaration, TypeDenoter
- Terminalknoten (z. B. Identifier, Operator) speichern tatsächlichen Text

1.10 Scanner (Lexikalische Analyse)

Scanner

Wandelt Zeichen → Tokens anhand regulärer Ausdrücke (REs).

Aufgaben:

- Entfernt Whitespace, Kommentare
- Liefert Token(kind, spelling, position)
- Nutzt endlichen Automaten oder rekursiven Abstieg

Beispiel EBNF Mini-Triangle:

```
Identifier ::= Letter (Letter — Digit)*  
Integer-Literal ::= Digit Digit*  
Operator ::= +| - | * | / | < | > | =
```

1.11 Automatisierung

- Scanner-Generatoren: **JLex, JFlex**
- Parser-Generatoren: **ANTLR (LL*), JavaCC (LL(k))**

1.12 Typische Fehler

- Linksrekursion nicht entfernt
- Linksausklammern vergessen
- Anfangs-/Folgemengen überschneiden sich → nicht LL(1)
- Schlüsselwörter nicht vom Identifier getrennt

1.13 Zusammenfassung – Wichtigste Punkte

- CFG-Grundlagen (BNF/EBNF)
- Transformationen: Gruppierung, Linksausklammern, Rekursionsbeseitigung
- LL(1)-Parsing und Bedingungen
- Rekursiver Abstieg: Struktur und Methoden
- AST-Struktur: Klassenhierarchie, Terminal- und Nichtterminalknoten
- Scanner: REs, endliche Automaten, Schlüsselworterkennung