



0101seda010100
software engineering dependability

Software Entwicklung 2

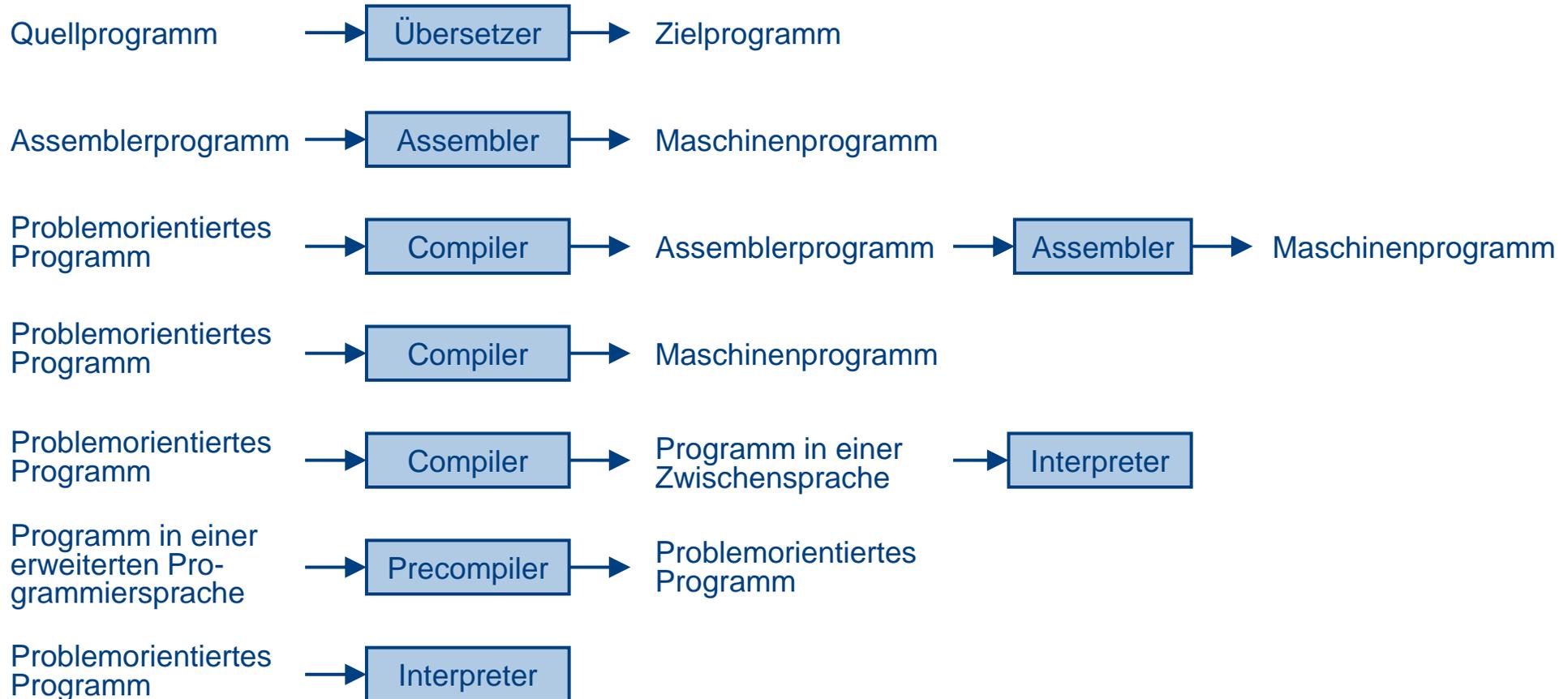
Übersetzerbau 1

- Aufgaben von Übersetzern
- Weitere Anwendungen von Übersetzern
- Arten von Übersetzern
- Grundlagen: Sprachdefinition
- Grammatik: BNF, EBNF, Syntaxdiagramme
- Produktionen, Ketten, Kettenmengen
- Satzform, Satz, Ableitung
- Chomsky-Hierarchie
- Aufbau von Compilern
- Aufgaben der Compiler-Komponenten: Analyseteil
- Aufgaben der Compiler-Komponenten: Syntheseteil
- Einpaß- vs. Mehrpaß-Compiler

- Wissen, welche Arten von Übersetzern es gibt
- BNF, EBNF und Syntaxdiagramme anwenden können und insbesondere ineinander überführen können
- Die Begriffe Terminalsymbol, Nichtterminalsymbol, Kette, Ableitung und Sprache definieren können
- Die Chomsky-Hierarchie anwenden können
- Aufgaben der unterschiedlichen Teile eines Übersetzers kennen

- Die Aufgabe der Übersetzung von Texten einer Sprache A in eine andere Sprache B ist für natürlich gesprochene und geschriebene Sprachen allgemein bekannt.
- Analog ist es in der Informatik die zentrale Aufgabe eines Übersetzers, alle Sätze einer Quellsprache A, die Programme, in gleichbedeutende Sätze einer Zielsprache, die Zielprogramme, zu transformieren.
- Der Begriff Übersetzer ist dabei die übergeordnete Bezeichnung von Werkzeugen für diese Aufgabe. In Abhängigkeit von den Eigenschaften der Quellsprache und den Eigenschaften der Zielsprache werden spezialisierte Begriffe verwendet.
- Mit den gleichen Techniken können Softwarekomponenten realisiert werden, die z.B. Kommandosyntax an Kommandoschnittstellen analysieren.

- Formatierer wie *nroff* und *troff* in UNIX sind in Wirklichkeit Übersetzer, die Text und Formatierkommandos übersetzen
- Precompiler wie *equ* (*Formatierung von Ausdrücken*), *tbl* (*Erzeugung von Tabellen*), *pic* (*Zeichnen von Bildern*) erlauben eine Erweiterung von *nroff* und *troff*.
- Ein *silicon compiler* übersetzt die Layout- und Aufbauspezifikation für VLSI-Schaltkreise in eine VLSI-Schaltkreismaske.
- Abfrageprogramme für Datenbanken. z.B. geschrieben in SQL werden durch Übersetzer in physikalische Zugriffsbefehle auf die Datenbankdateien transformiert.
- Programme geschrieben in einer Roboterprogrammiersprache werden durch einen Compiler in eine Zwischensprache, z.B. IRDATA (*Industrial Robot Data*) übersetzt.



Beispiel:

- Regeln zur Bildung gültiger Sätze:
 - $\text{SATZ} ::= \text{SUBJEKT PRÄDIKAT}$.
 - $\text{SUBJEKT} ::= \text{"Peter"} \mid \text{"Siegfried"}$.
 - $\text{PRÄDIKAT} ::= \text{"arbeitet"} \mid \text{"schläft"}$.
- Gültige Sätze dieser Sprache:
 - Peter arbeitet
 - Siegfried arbeitet
 - Peter schläft
 - Siegfried schläft
- Ungültiger Satz dieser Sprache:
 - Peter arbeitet nicht

Beispiel (erweitert):

- Regeln zur Bildung gültiger Sätze:

- SATZ ::= SUBJEKT PRÄDIKAT ZUSATZ.
- SUBJEKT ::= "Peter" | "Siegfried".
- PRÄDIKAT ::= "arbeitet" | "schläft".
- ZUSATZ ::= "nicht" | ϵ .

- Gültige Sätze dieser Sprache:

- | | |
|----------------------|--------------------------|
| • Peter arbeitet | Peter arbeitet nicht |
| • Siegfried arbeitet | Siegfried arbeitet nicht |
| • Peter schläft | Peter schläft nicht |
| • Siegfried schläft | Siegfried schläft nicht |

- Viertupel aus
 - N: Vokabular der nicht-terminalen Symbole:
z.B. SATZ, SUBJEKT, PRÄDIKAT
 - S: Startsymbol aus N: z.B. SATZ
 - T: Vokabular der terminalen Symbole:
z.B. Peter, Siegfried, arbeitet, schläft
 - P: Menge der Produktionsregeln: siehe Beispiel

- N und T sind disjunkt!

- Vokabular $V = N \cup T$

- 1960 von Backus und Naur zur Beschreibung von Algol 60 eingeführt:
- $\text{syntax} ::= \text{production syntax} \mid \varepsilon$.
- $\text{production} ::= \text{identifrier} ::= \text{expression} \text{ " . "}$.
- $\text{expression} ::= \text{term} \mid \text{expression} \mid \text{term}$.
- $\text{term} ::= \text{factor} \mid \text{term factor}$.
- $\text{factor} ::= \text{identifrier} \mid \text{string}$.

- $\text{identifrier} ::= \text{letter} \mid \text{identifrier letter} \mid \text{identifrier digit}$.
- $\text{string} ::= \text{stringhead} \text{ " " " "}$.
- $\text{stringhead} ::= \text{ " " " " } \mid \text{stringhead letter}$.
- $\text{letter} ::= \text{ "A" } \mid \dots \mid \text{ "Z"}$.
- $\text{digit} ::= \text{ "0" } \mid \dots \mid \text{ "9"}$.

- 1977 von Wirth zur Verbesserung der Lesbarkeit eingeführt:
- $\text{syntax} ::= \{\text{production}\}.$
- $\text{production} ::= \text{identifizier} "::=" \text{expression} ".".$
- $\text{expression} ::= \text{term} \{ "|" \text{term}\}.$
- $\text{term} ::= \text{faktor} \{\text{faktor}\}.$
- $\text{faktor} ::= \text{identifizier} | \text{string} | "(" \text{expression} ")" | "[" \text{expression} "]" | "{" \text{expression} "}".$

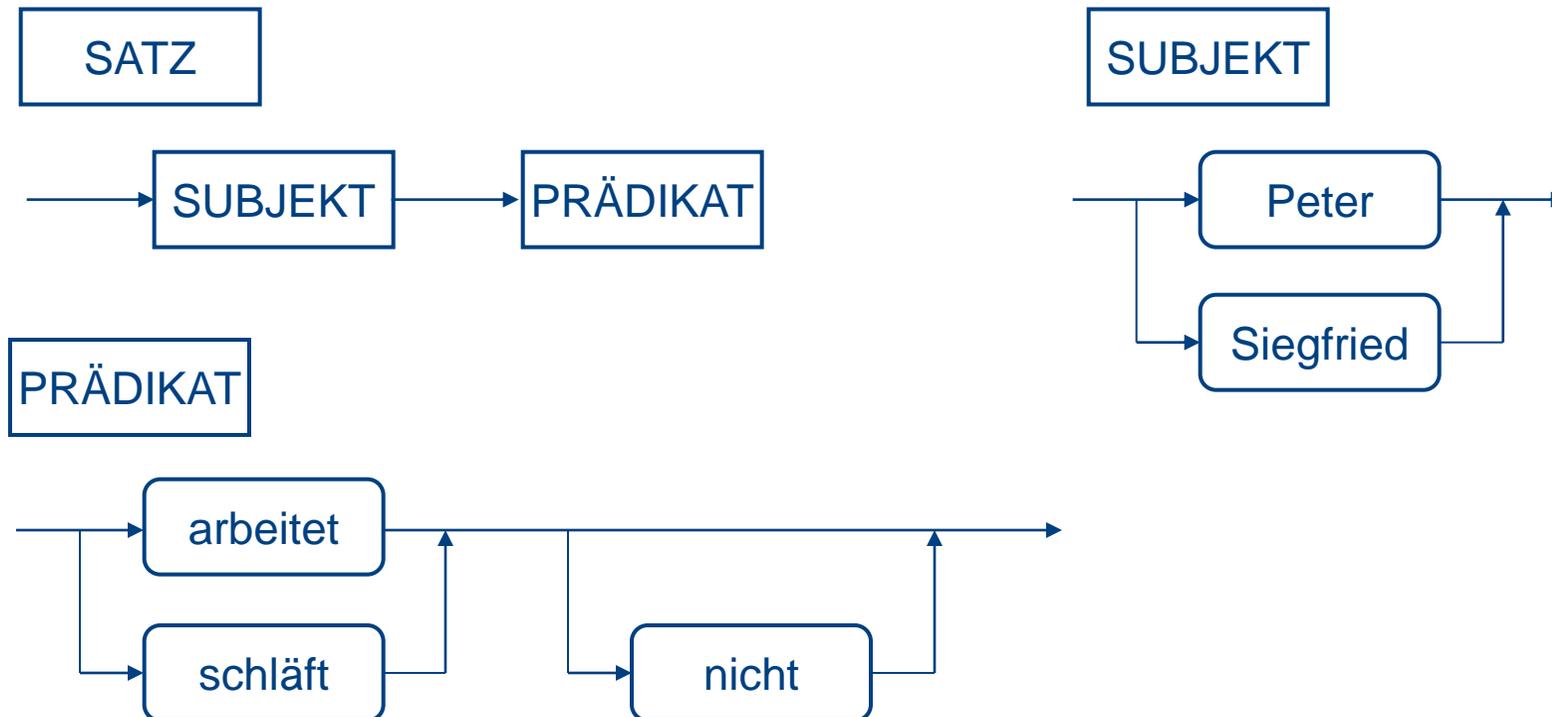
- $\text{identifizier} ::= \text{letter} \{\text{letter} | \text{digit}\}.$
- $\text{string} ::= "" \{\text{character}\} "".$
- $\text{letter} ::= "A" | \dots | "Z".$
- $\text{digit} ::= "0" | \dots | "9".$

- ::= : „Zuweisungsoperator“ zwischen linker und rechter Seite einer Produktion
- | : Alternative
- { } : Wiederholung (kein mal, einmal und beliebig oft)
- [] : Option (kein mal oder einmal)
- () : Klammerung

- Beispiel: Sprache, die die Sätze ε , a, aa, aaa, ... enthält.
 - BNF: $S ::= A.$
 $A ::= aA \mid \varepsilon.$
 - EBNF: $S ::= A.$
 $A ::= \{a\}.$

- Das erweiterte Beispiel in EBNF:
 - SATZ ::= SUBJEKT PRÄDIKAT.
 - SUBJEKT ::= "Peter" | "Siegfried".
 - PRÄDIKAT ::= ("arbeitet" | "schläft") ["nicht"].

- Anstelle von Produktionen in BNF oder EBNF können auch Syntaxdiagramme angegeben werden:



- Jedes Syntaxdiagramm beschreibt eine Produktionsregel der Sprache
 - Die Produktionsregeln beschreiben, welche Symbolfolgen der Sprache legal bzw. erlaubt sind
- Notationen
 - Syntaxdiagramme
 - Erweiterte Backus Naur Form (EBNF)
 - Backus Naur Form (BNF).

- (Symbol-)Kette: Endliche Folge σ von Symbolen über einem Vokabular V
 - $\sigma_1 = \text{BeginSym ReadSym LParen <idlist> RParen SemiColonEndSym EofSym}$
 - $\sigma_2 = \text{BeginSym <StatementList> EndSym EofSym}$
- Eine Kette kann aus nicht-terminalen und terminalen Symbolen bestehen.

- Abkürzungen
 - a^n : Kette aus n gleichen Symbolen a
 - $a^3 = aaa$
 - ε leere Kette, d.h. die Kette aus null Symbolen
 - a^+ : Menge $\{a^n : n \geq 1\}$
 - $a^+ = \{a, aa, aaa, aaaa, \dots\}$
 - a^* : Menge $\{a^n : n \geq 0\}$
 - $a^* = \{\varepsilon, a, aa, aaa, aaaa, \dots\}$.
 - Offensichtlich ist $a^* = a^+ \cup \{\varepsilon\}$.

- V^+ : Menge aller nichtleeren Ketten, die sich aus den Symbolen von V bilden lassen
 - $V^+ = V \cup V^2 \cup \dots$
- V^* : Menge aller Ketten, einschließlich der leeren Kette, die sich aus den Symbolen von V bilden lassen
 - $V^* = \{\varepsilon\} \cup V^+$
- Die Menge V ist immer endlich
- Die Mengen a^+ , a^* , V^+ , und V^* sind immer unendlich
- Für $V = \{a, b, c\}$ sind
 - $V^+ = \{a, b, c, aa, ab, ac, ba, bb, bc, ca, cb, cc, aaa, \dots\}$
 - $V^* = \{\varepsilon, a, b, c, aa, ab, ac, ba, bb, bc, ca, cb, cc, aaa, \dots\}$

- Sprache = Menge von Symbolketten, die von der Syntax als **wohlgeformt definiert** sind.
- Produktionsregeln beschreiben Ersetzungen von Symbolketten durch Symbolketten
 - Ausgangspunkt: Startsymbol S
 - Wiederholtes Ersetzen gemäß den Regeln der Grammatik durch Symbolketten, bis eine Symbolkette entsteht, die nur noch terminale Symbole enthält.

- Beginn mit dem Syntaxdiagramm **Start** bzw. der Produktion, auf deren linker Seite das Nichtterminal-Symbol **Start** vorkommt
- Alle Syntaxdiagramme bzw. Produktionsregeln solange anwenden, bis nur noch terminale Symbole vorhanden sind
- Jede Sequenz terminaler Symbole, die so erzeugt werden kann, ist gültig
- Umgekehrt gilt, dass jede Sequenz terminaler Symbole, die nicht durch eine Sequenz von Ersetzungen nicht-terminaler Symbole erzeugt werden kann, illegal ist.

- Rekursiv anwendbare Analyseroutinen, die durch den Syntaxbaum absteigen, den sie bei der Programmbearbeitung erkennen
 - Einfachste Syntaxanalyse-Technik
 - wird in vielen Compilern verwendet
- Grundidee
 - Zu jedem nicht-terminalen Symbol existiert eine (gleichnamige) Analyseprozedur, die jede terminale Symbolfolge erkennt, die durch das nicht-terminale Symbol erzeugt werden kann.
 - In einer Analyse-Prozedur können sowohl nicht-terminale als auch terminale Symbole auftauchen.

- Satzformen
 - Das Startsymbol S selbst und alle Symbolketten, die aus S durch die Anwendung der Produktionen hervorgehen
- Sätze
 - Satzformen, die nur aus terminalen Symbolen bestehen
- Sind α und β zwei Satzformen und geht β aus α durch Anwendung einer Produktion hervor, dann $\alpha \Rightarrow \beta$. (Ableitung)

• Direkte Ableitung

- Eine Kette α produziert direkt eine Kette β
 - geschrieben $\alpha \Rightarrow \beta$,
- wenn es irgendwelche Ketten ω_1 und ω_2 gibt, so daß
 - $\alpha = \omega_1 \sigma \omega_2$
 - $\beta = \omega_1 \tau \omega_2$ und
 - die Produktion $\sigma ::= \tau$ zur Grammatik gehört
- β heißt dann eine direkte Ableitung von α

• Mehrfache Ableitung

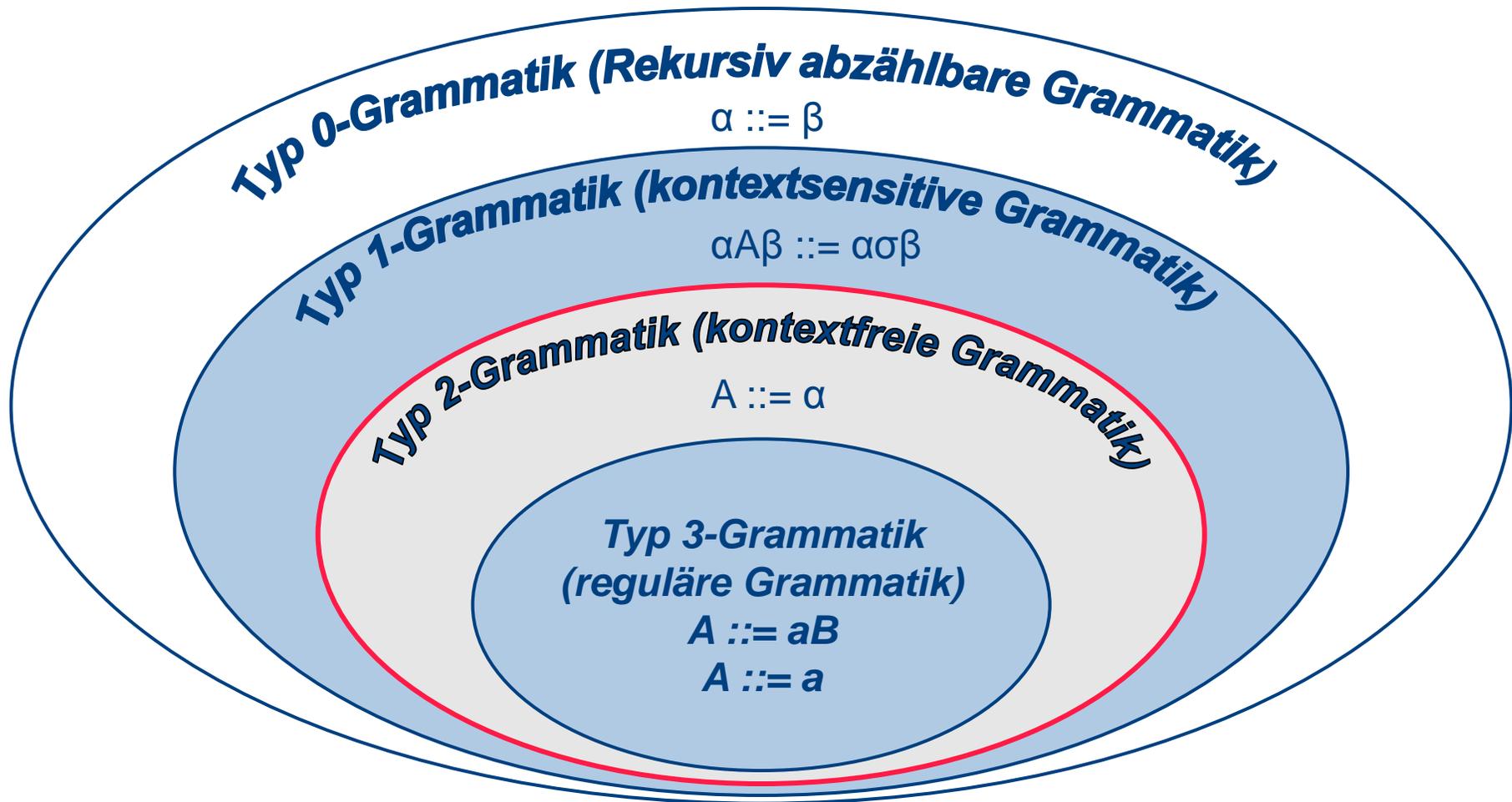
- Eine Kette α produziert eine Kette β
 - geschrieben $\alpha \Rightarrow^+ \beta$,
 - wenn es eine Folge direkter Ableitungen $\alpha = \omega_0 \Rightarrow \omega_1 \Rightarrow \omega_2 \Rightarrow \dots \Rightarrow \omega_n = \beta$ mit $n \geq 1$ gibt.
- Diese Folge heißt mehrfache Ableitung der Länge n

- Wenn G eine Grammatik mit dem Startsymbol S ist, dann wird eine Kette α Satzform genannt, wenn $S \Rightarrow^* \alpha$ gilt.
- Ein Satz ist eine Satzform, die nur aus terminalen Symbolen besteht.
- Eine Sprache $L(G)$ ist die Menge aller Sätze, die sich aus dem Startsymbol ableiten lassen:

$$L(G) = \{ \alpha \mid S \Rightarrow^+ \alpha \wedge \alpha \in T^* \}$$

- Die Chomsky-Hierarchie ist eine Aufteilung von Sprachen in Klassen:
- Reguläre Sprachen (Typ 3):
 - Können durch endliche Automaten dargestellt werden
 - Werden durch Scanner analysiert
 - Dürfen nur Produktionen der folgenden Form enthalten:
 $A \rightarrow a$ und $A \rightarrow a B$ bzw. $A \rightarrow a$ und $A \rightarrow B a$.
(Kleinbuchstaben stellen Terminalsymbole dar, Großbuchstaben repräsentieren nichtterminale Symbole)
- Kontextfreie Sprachen (Typ 2)
 - Können durch Kellerautomaten dargestellt werden
 - Werden durch Parser analysiert
 - Auf der linken Seite jeder Produktion muß stets genau ein Nichtterminal-Symbol stehen,
z.B.: $A \rightarrow a B b$
 - Die rechte Seite darf aus einer beliebigen Kombination von terminalen und nichtterminalen Symbolen bestehen

- Kontextsensitive Sprachen (Typ 1)
 - Produktionen haben die Form $\alpha \gamma \beta \rightarrow \alpha \sigma \beta$ mit $\alpha, \beta \in T^*$, $\sigma \in V^+$, $\gamma \in N^+$
 - Die rechte Seite von Produktionen muß mindestens so lang sein wie die linke Seite (längenmonoton)
 - Auf der rechten Produktionsseite darf nur dann das leere Symbol (ϵ) stehen, wenn die linke Produktionsseite das Startsymbol S ist, und dieses auf keiner rechten Produktionsseite in der Grammatik auftritt: $S \rightarrow \epsilon$
- Rekursiv aufzählbare Sprachen (Typ 0)
 - Jede wohlgeformte Grammatik, in der die Terminalsymbole, Nichtterminal-Symbole und das Startsymbol definiert sind, ist vom Chomsky-Typ 0

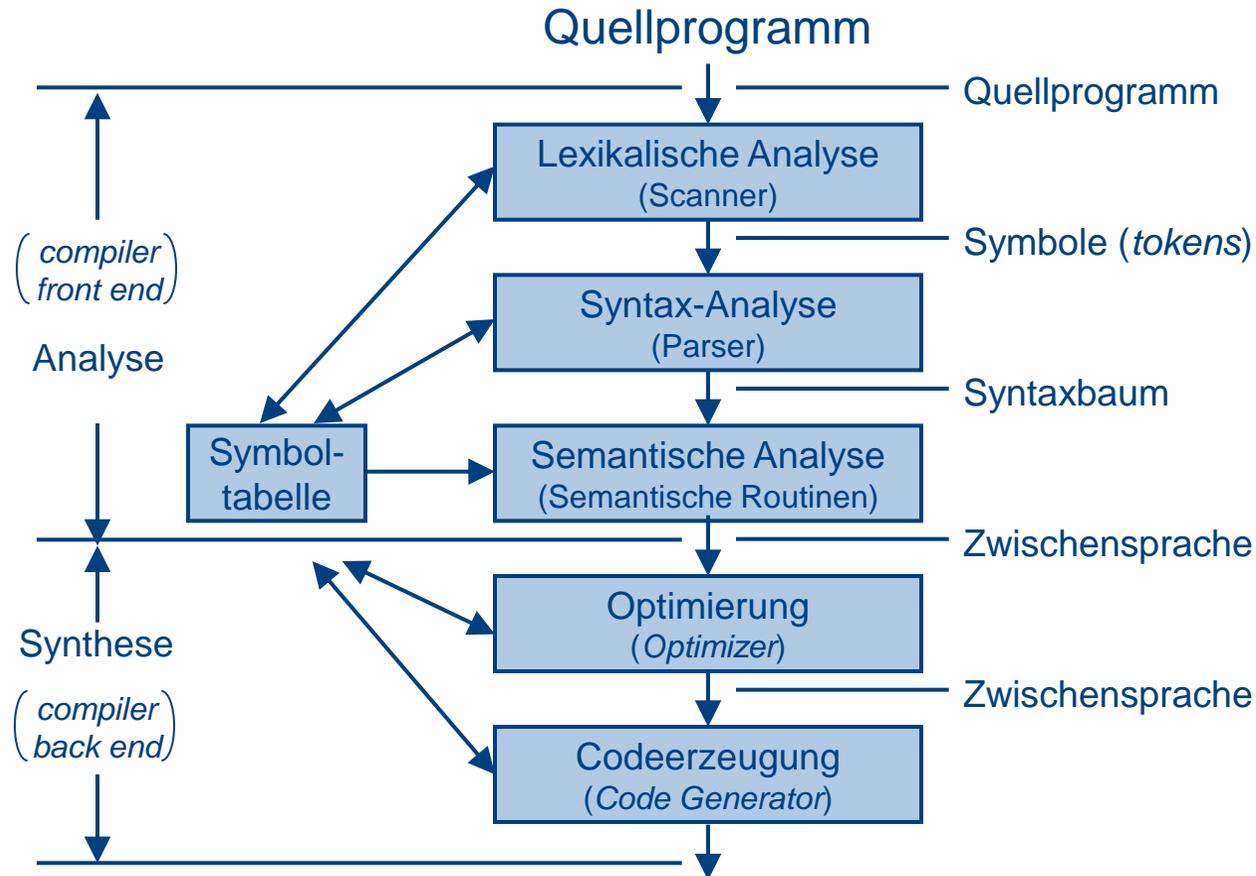


$\alpha =$ beliebige Kette aus N und T

26

- Für den Compilerbau sind nur reguläre und kontextfreie Sprachen interessant
- Reguläre Grammatiken:
 - Definition der Bezeichner und Schlüsselworte, sowie von Konstanten und anderen Symbolen der Sprache
 - Scanner: Zerlegung des Eingabetextes in Symbole
- Kontextfreie Grammatiken:
 - Festlegung der Regeln, wie aus Symbolen korrekte Sätze der Sprache gebildet werden
 - Parser: Syntaxanalyse

- Kontextsensitivität: Übersetzung eines deutschen Textes ins Englische
- *Ich mag Montage nicht.*
- Version 1: *I don't like Mondays.*
- Version 2: *I don't like assembly.*



- **Scanner: Durchführung der lexikalischen Analyse:**

- Das Quellprogramm wird zeichenweise gelesen. Die individuellen Zeichen werden zu Symbolen (*tokens*) zusammengefaßt, die die Eingabe für den Parser darstellen.
- Beispiele für Symbole sind Bezeichner, Zahlen, Operatoren, Begrenzer wie Klammern und Semikolon.

- **Parser: Durchführung der syntaktischen Analyse:**

- Das aus Symbolen bestehende Programm wird in seine grammatischen Bestandteile zerlegt und die Struktur als Syntaxbaum dargestellt.

- Semantische Routinen: Durchführung der semantischen Analyse:
 - Überprüfung der statischen Semantik jedes Sprachkonstrukts, wie Gültigkeitsbereiche von Bezeichnern, Korrespondenz zwischen Deklaration und Anwendung von Bezeichnern, Typverträglichkeit von Operatoren in Ausdrücken. Ist das Konstrukt semantisch korrekt, so wird es in eine Zwischensprache transformiert.
 - Die semantische Analyse kann mit der Syntaxanalyse gemeinsam stattfinden. In diesem Fall verschmelzen diese beiden Analyseteile zu einem Analyseteil.
 - Erfolgt die semantische Analyse separat, dann wird der Syntaxbaum, das Ergebnis der Syntaxanalyse, bei der semantischen Analyse durch semantische Informationen ergänzt.
 - Als Ergebnis liefert der Analyseteil des Compilers,
 - die Feststellung, ob das Quellprogramm korrekt ist
 - eine speicherinterne Darstellung des Quellprogramms, die für den Syntheseteil gut geeignet ist,
 - speicherinterne Tabellen, die für die weitere Verarbeitung der Zwischendarstellung benötigt werden.

Compiler

Beispiel: Lexikalische und syntaktische Analyse

Quellprogramm

Prämie:= 100 + Dienstjahre * Leistungsfaktor;

Lexikalische Analyse

Id AssignOp Number PlusOp Id TimesOp Id SemiColon

Id/1 AssignOp Number/100 PlusOp Id/2 TimesOp Id/3 SemiColon

Symbole

Syntaxanalyse

*Symbole und
Symbolklassen*

Id/1 AssignOp Number/100 PlusOp Id/2 TimesOp Id/3 SemiColon

Variable

Term

Term

Ausdruck

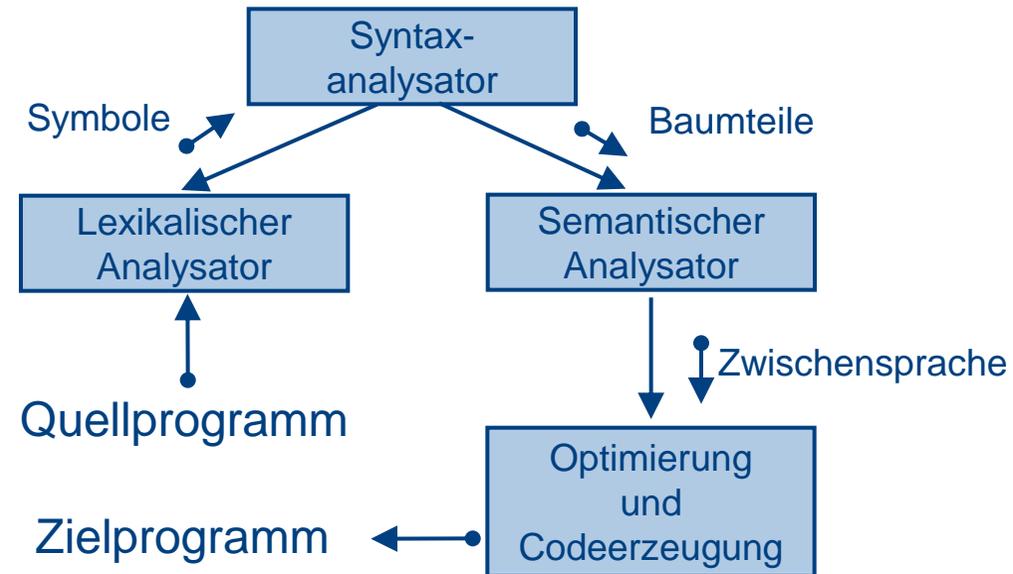
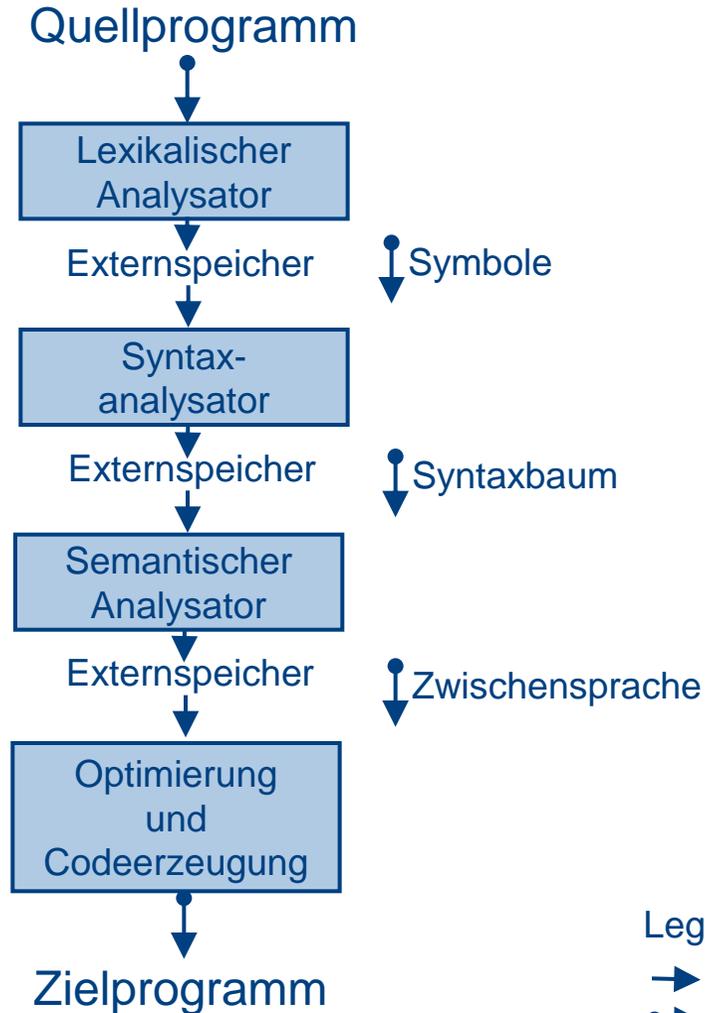
Syntaxbaum

Zuweisung

Symboltabelle
1: Prämie
2: Dienstjahre
3: Leistungsfaktor

- Optimierer:
 - Transformation des Programms in der Zwischensprache in ein funktional äquivalentes Programm in der Zwischensprache, das hinsichtlich bestimmter Kriterien (z.B. Geschwindigkeit) verbessert ist.
 - Der Optimierer kann bei vielen Compilern auch deaktiviert werden.
- Code-Generator:
 - Erzeugung des Zielprogramms aus der Zwischendarstellung.
- Symboltabelle:
 - Sowohl der Analyseteil als auch der Syntheseteil des Compilers greifen auf die sogenannte Symboltabelle zu. In der Symboltabelle werden die Bezeichner mit ihren Attributen gespeichert.
 - Falls keine Optimierung gefordert wird, so kann die Code-Generierung in die semantischen Routinen integriert werden (keine Zwischendarstellung)

- Mehrpaß-Compiler:
 - Nacheinanderausführung der in der statischen Compilerstruktur aufgeführten Komponenten
 - Man spricht dann von einem Mehrpaß-Compiler, da eine Komponente nach der anderen ausgeführt wird. Als Ergebnis jeder Komponente entsteht ein Programm in einer Zwischensprache, das auf einen externen Speicher geschrieben und vom nächsten Paß wieder gelesen wird
- Einpaß-Compiler:
 - Zeitlich verzahnte Ausführung der einzelnen Compilerkomponenten
 - Der Parser ruft den Scanner auf, wenn er das nächste Quellsymbol benötigt. Er ruft die semantischen Routinen auf, wenn er eine syntak- tisch korrekte Konstruktion an sie weitergeben will.



Legende:

- Steuerfluss
- Datenfluss