
Werkzeuge für den
Übersetzerbau

J. Grosch
H. Emmelmann

DR. JOSEF GROSCH

COCOLAB - DATENVERARBEITUNG

GERMANY

Cocktail

Toolbox for Compiler Construction

Werkzeuge für den Übersetzerbau

Josef Grosch, Helmut Emmelmann

Feb. 7, 1995

Document No. 21

Copyright © 1995 Dr. Josef Grosch

Dr. Josef Grosch
CoCoLab - Datenverarbeitung
Breslauer Str. 64c
76139 Karlsruhe
Germany

Phone: +49-721-91537544
Fax: +49-721-91537543
Email: grosch@cocolab.com

Werkzeuge für den Übersetzerbau

Übersicht

Mit Übersetzerbau-Werkzeugen lassen sich Übersetzer für Programmiersprachen weitgehend automatisch generieren. Wir stellen einen Werkzeugkasten vor, welcher die Konstruktion nahezu aller Phasen eines Übersetzers unterstützt. Die Entwurfsziele für diesen Werkzeugkasten waren praktische Brauchbarkeit, deutlich reduzierter Erstellungsaufwand für Übersetzer und hohe Qualität der generierten Übersetzer. Besonders hinsichtlich Effizienz sind die Werkzeuge konkurrenzfähig zur Programmierung von Hand. Zur Zeit können mit den Werkzeugen Übersetzermodule in den Zielsprachen C, C++ und Modula-2 sowie teilweise Ada und Eiffel erzeugt werden. Viele realistische Anwendungen demonstrieren die ausgezeichnete Leistungsfähigkeit der Werkzeuge und zeigen, daß die Werkzeuge die Konstruktion von Übersetzern mit Produktionsqualität erlauben.

1. Aufbau eines Übersetzers

Ein wichtiges Hilfsmittel zur Programmierung eines Computers ist ein Übersetzer (compiler). Ein Übersetzer ist ein Programm, welches ein in einer Programmiersprache geschriebenes Programm in eine Maschinensprache übersetzt. Die Hardware versteht genaugenommen nur aus Nullen und Einsen zusammengesetzte Maschinensprach-Programme. Um einem Computer auch eine für den menschlichen Programmierer besser geeignete höhere Programmiersprache verständlich zu machen, ist eine Übersetzung nötig.

Die Konstruktion eines Übersetzers ist eine anspruchsvolle und aufwendige Aufgabe. Der Bedarf an Übersetzern ist relativ groß, da für jede Programmiersprache und jeden Computer ein eigener Übersetzer notwendig ist. Es lohnt sich daher, nach Methoden zu suchen, die die Erstellung von Compilern zu vereinfachen. Doch bevor wir zu unserem eigentlichen Thema kommen, nämlich der automatischen Generierung von Übersetzern mit Übersetzerbau-Werkzeugen, möchten wir kurz den Aufbau und die prinzipielle Funktionsweise eines Übersetzers erläutern. Die rechte Spalte in Abb. 1 zeigt die Phasen bzw. Module eines Übersetzers.

Die lexikalische Analyse liest das Quellprogramm zeichenweise. Sie faßt die Zeichenfolgen für Bezeichner, Zahlen und Schlüsselwörter zu Grundsymbolen (tokens) zusammen und überliest Zwischenräume und Kommentare.

Die syntaktische Analyse hat als Eingabe eine Folge von Grundsymbolen. Sie überprüft das Quellprogramm auf syntaktische Fehler und rekonstruiert die Struktur des Programms, d. h. sie erkennt den Aufbau der Ausdrücke und Anweisungen sowie deren Zusammenhang. Diese Struktur wird oft in Form eines Syntaxbaums gespeichert.

Die semantische Analyse überprüft die Kontextbedingungen bzw. die Regeln der statischen Semantik und berechnet für die Codegenerierung nötige Eigenschaften. Ein Beispiel für eine Kontextbedingung ist die Vorschrift, daß alle Variablen deklariert sein müssen. Zur statischen Semantik zählen die Analyse der Gültigkeitsbereiche, die Namensanalyse, d. h. die Feststellung der zu einem Bezeichner gehörenden Deklaration, und die Typüberprüfung.

Zur Vereinfachung der gesamten Übersetzungsaufgabe wird diese häufig in zwei Schritte unterteilt. Der Syntaxbaum wird zunächst von einer Transformationsphase in eine Zwischensprache umgewandelt. Diese Zwischensprache ist meist maschinenorientiert jedoch noch maschinenunabhängig. Das niedrige Niveau der Zwischensprache erleichtert dem Codegenerator die Erzeugung der Maschinensprache.

Zu den Aufgaben des Codegenerators zählen die Befehlsauswahl, d. h. die Abbildung der Zwischensprachanweisungen auf Maschinenbefehle, sowie die Speicher- und Registerzuteilung. Die Ausgabe ist schließlich ein binär-codiertes Maschinenprogramm.

Der folgende Abschnitt beschreibt die Vorteile, die Entwurfsziele und den Inhalt des Werkzeugkastens. Abschnitt 3 stellt die gemeinsamen Eigenschaften der Werkzeuge dar. Im Abschnitt 4 wird das von uns bevorzugte Übersetzermodell beschrieben. Der Abschnitt 5 enthält eine kurze Darstellung der einzelnen Werkzeuge. Abschnitt 6 berichtet von den Erfahrungen des Einsatzes der Werkzeuge in einer realistischen Anwendungen. Abschnitt 7 enthält eine Zusammenfassung und beschreibt weiterführende Arbeiten.

2. Werkzeugkasten

Die Erstellung eines Übersetzers von Hand ist eine sehr anspruchsvolle und aufwendige Aufgabe. Durch den Einsatz von Übersetzerbau-Werkzeugen läßt sich dieser Aufwand reduzieren. Im folgenden stellen wir einen Werkzeugkasten zur Übersetzer-Generierung vor, welcher für nahezu jede Übersetzerphase Werkzeuge enthält. Diese sind für den Einsatz in realistischen Übersetzerprojekten konzipiert.

Im allgemeinen akzeptieren die Werkzeuge als Eingabe eine Spezifikation, die in einer werkzeug-spezifischen Sprache geschrieben ist. Sie produzieren als Ausgabe ein Programm-Modul in einer Zielsprache (C, C++ oder Modula-2 und teilweise Ada oder Eiffel). Deshalb kann ein Werkzeug als generische Lösung eines Teilproblems in einem Übersetzer gesehen werden, wobei mit Hilfe einer Spezifikation eine konkrete Lösung gewonnen wird.

Die Benutzung von Werkzeugen hat gegenüber der Programmierung von Hand mehrere Vorteile: Der zur Konstruktion eines Übersetzers notwendige Aufwand wird wesentlich verringert. An Stelle eines Programms wird eine viel kürzere Spezifikation entwickelt. Die Werkzeuge können eine Spezifikation in vielfacher Weise auf Konsistenz überprüfen. Das Schreiben automatisch prüfbarer Spezifikationen verringert die Anzahl möglicher Fehler und erhöht so die Zuverlässigkeit des resultierenden Übersetzers.

Die wichtigsten Entwurfsziele für den Werkzeugkasten waren:

- praktische Brauchbarkeit für realistische Programmiersprachen
- automatische Generierung von Übersetzern mit Produktionsqualität
- wesentliche Steigerung der Übersetzerbau-Produktivität
- mit Handprogrammierung vergleichbare Qualität der erzeugten Übersetzer

Mit dieser Zielsetzung sollte die praktische Einsatzfähigkeit des Werkzeugkastens in realistischen Übersetzerbauprojekten erreicht werden. Daher wurde auch die Konkurrenzfähigkeit zur Handprogrammierung betont. Wir meinen, daß die hohe Produktivität und Zuverlässigkeit nicht durch eine geringere Codequalität oder Effizienz des resultierenden Compilers erkaufte werden muß.

Der Werkzeugkasten enthält folgende Werkzeuge:

| | |
|-------|---|
| Rex | Generator für lexikalische Analysatoren |
| Lark | LR(1) Parser-Generator mit Backtracking |
| Ell | LL(1) Parser-Generator |
| Ast | Generator für abstrakte Syntaxbäume |
| Ag | Generator für Attributauswerter |
| Puma | Transformation attributierter Syntaxbäume |
| Reuse | Bibliothek wiederverwendbarer Module |

Alle Werkzeuge wurden ursprünglich in Modula-2 programmiert und laufen unter dem Betriebssystem UNIX. Unter Verwendung des Modula-2 nach C Übersetzers *Mtc* [Mar90] (siehe Abschnitt 6), konnte von den Programmen automatisch eine C-Version erstellt werden. Zur Zeit erzeugen die Werkzeuge Module in den Zielsprachen C, C++ oder Modula-2 und teilweise auch Ada oder Eiffel.

3. Gemeinsame Eigenschaften

Unsere Entwurfsziele führten zu einigen für alle Werkzeuge gemeinsamen Entwurfsentscheidungen. Nahezu jedes Werkzeug benötigt eine Programmiersprache, mit der der Benutzer gewisse Aktionen, Bedingungen oder Berechnungen spezifizieren kann. Das trifft offensichtlich für Attributgrammatiken zu, aber auch der Transformations-Generator muß Attribute und Bedingungen auswerten. Sogar die Parser-Generatoren brauchen eine solche Sprache zur Spezifikation semantischer Aktionen.

Wir entschieden uns dafür, direkt die Zielsprache (nämlich C, C++ oder Modula-2 usw.) zu verwenden. Deshalb können Spezifikationen Abschnitte mit Zielsprachanweisungen enthalten. Abgesehen von geringfügigen Ersetzungen wird dieser Text unverändert in die erzeugten Module kopiert. Der Nachteil dieser Methode ist, daß die in der Zielsprache geschriebenen Teile nicht vollständig von den Werkzeugen überprüft werden können. Zum Beispiel kann das Attributgrammatik-Werkzeug nicht überprüfen, ob Attributberechnungen keine Seiteneffekte haben. Andererseits wird damit eine sehr große Flexibilität erreicht, da die volle Ausdruckskraft der Zielsprache zur Verfügung steht. Ebenso wird die praktische Brauchbarkeit drastisch erhöht, da die Einbeziehung anderer, eventuell handgeschriebener Komponenten leicht möglich ist. Schließlich führt es zu einfachen Werkzeugen und einfachen Spezifikationssprachen.

Unsere Erfahrung mit früheren Werkzeugen zeigte, daß während der Konstruktion realistischer Übersetzer eine Reihe kleiner Sonderprobleme auftritt, die nicht mit den Werkzeugen gelöst werden können. Deswegen sind Schlupflöcher nötig, also Möglichkeiten, die es dem Werkzeugbenutzer erlauben, leicht handgeschriebene Teile einzufügen. Diese Schlupflöcher tragen auch dazu bei, die Werkzeuge zu vereinfachen, da man nicht gezwungen ist, für jedes Problem sofort eine Lösung bereitzustellen. Das Schlupfloch kann benutzt werden solange bis eine wirklich gute Lösung gefunden wird, welche man in ein Werkzeug einbauen kann.

Die Werkzeuge sind größtenteils von einander unabhängig. Dies wird dadurch erzielt, daß keiner der erzeugten Module eine festgelegte Ausgabe besitzt. Stattdessen wird diese Ausgabe mittels Anweisungen der Zielsprache spezifiziert und kann somit beliebig gewählt werden. Die Unabhängigkeit der Werkzeuge sorgt für große Freiheiten beim Übersetzerentwurf. Eine Ausnahme bilden die Werkzeuge *Ag* und *Puma*, denn sie basieren auf den mit *Ast* spezifizierten Syntaxbäumen. Deshalb hängen diese Werkzeuge von *Ast* ab, und alle drei Werkzeuge sind für Übersetzer zugeschnitten, die einen attributierten abstrakten Syntaxbaum benutzen.

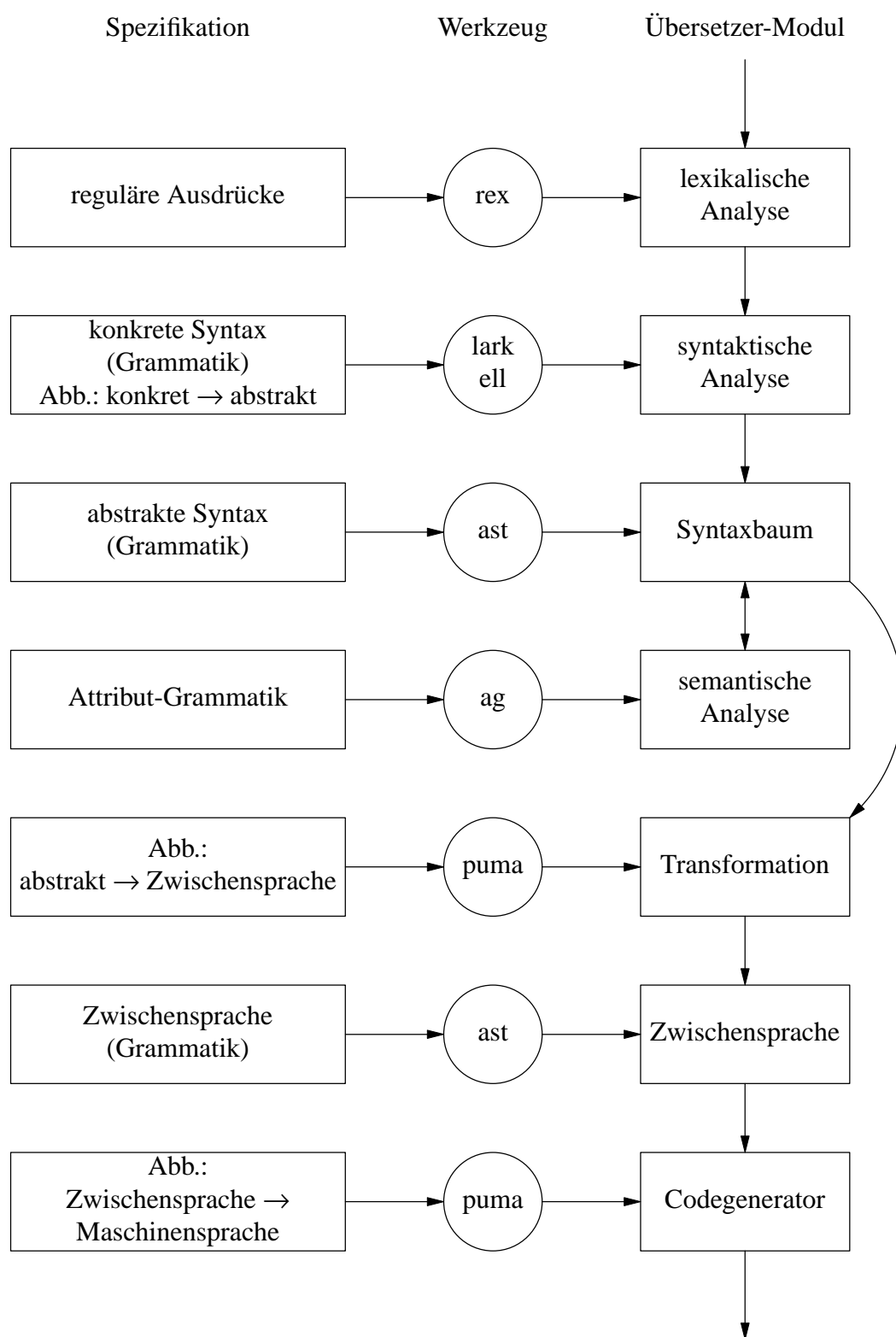


Abb. 1: Übersetzer-Modell

4. Übersetzer-Modell

Obwohl die Werkzeuge kein bestimmtes Übersetzer-Modell erzwingen, möchten wir das von uns bevorzugte Modell vorstellen. Wir meinen, daß dieses am besten von den Werkzeugen unterstützt wird. Wir betrachten die semantische Analyse nach wie vor als den schwierigsten Teil eines Übersetzers. Deshalb gehen wir für die semantische Analyse und die Erzeugung einer Zwischensprache von der abstrakten Syntax aus. Wir bauen den abstrakten Syntaxbaum explizit auf, welcher während der semantischen Analyse eventuell mit Attributen ergänzt wird. Neben der abstrakten Syntax, welche als erste, hohe Zwischensprache betrachtet werden kann, bevorzugen wir die Verwendung einer zweiten, niederen Zwischensprache als Schnittstelle zum Codegenerator. Dies bringt Vorteile in der Optimierung und der mustergesteuerten Codeauswahl mit sich.

Abbildung 1 zeigt das von uns bevorzugte Übersetzermodell. Die rechte Spalte enthält die wichtigsten Module eines Übersetzers. Die linke Spalte zeigt die dazu notwendigen Spezifikationen. Die dazwischen liegenden Werkzeuge werden von den Spezifikationen gesteuert und erzeugen die einzelnen Module. Die Pfeile stellen den Datenfluß dar, teils zur Generierungszeit und teils zur Übersetzungszeit.

5. Die Werkzeuge

Die folgenden Abschnitte stellen kurz die einzelnen Werkzeuge des Werkzeugkastens vor. Wir beschreiben nur die Eigenschaften der Werkzeuge. Für weitere Einzelheiten, die Spezifikationstechniken oder für Beispiele sei der Leser auf die existierenden, werkzeug-spezifischen Dokumente verwiesen.

5.1. Rex

Rex (regular expression tool) ist ein Generator für lexikalische Analysatoren [Gro88, Gro89, Groa]. Seine Spezifikationen basieren auf regulären Ausdrücken und beliebigen semantischen Aktionen, die in einer der Zielsprachen C, C++, Modula-2, Ada oder Eiffel geschrieben werden. Immer wenn in der Eingabe des erzeugten lexikalischen Analysators eine einem regulären Ausdruck entsprechende Zeichenkette erkannt wurde, werden die zugehörigen Aktionen ausgeführt. Falls zur eindeutigen Erkennung der Symbole der Kontext betrachtet werden muß, so kann der rechte Kontext durch einen zusätzlichen regulären Ausdruck spezifiziert werden, und der linke Kontext wird mit sogenannten Start-Zuständen behandelt. Falls mehrere reguläre Ausdrücke auf die aktuelle Eingabe zutreffen, so wird der Ausdruck mit der längsten Zeichenkette bevorzugt. Falls es immer noch mehrere Möglichkeiten gibt, so wird der zuerst in der Spezifikation stehende Ausdruck gewählt.

Die erzeugten lexikalischen Analysatoren berechnen automatisch Zeile und Spalte für jedes erkannte Symbol und enthalten einen Mechanismus für Include-Dateien. Bezeichner und Schlüsselwörter können effizient in Groß- oder Kleinbuchstaben normalisiert werden. Es gibt vordefinierte Regeln, um Leerstellen, Tabulatoren und Zeilenwechsel zu überlesen. Die generierten lexikalischen Analysatoren sind tabellengesteuerte, deterministische endliche Automaten. Die Tabellen werden mit der sogenannten Kammvektortechnik komprimiert [ASU86].

Die herausragende Eigenschaft von *Rex* ist seine Geschwindigkeit. Die lexikalischen Analysatoren verarbeiten etwa 2 Millionen Symbole (tokens) pro Minute ohne Hashing von Bezeichnern und 1,5 Millionen Symbole pro Minute mit Hashing (auf einer SPARC station ELC). Dies ist die vierfache Geschwindigkeit gegenüber mit *Lex* [Les75] generierten lexikalischen Analysatoren. In typischen Fällen besitzen mit *Rex* generierte Analysatoren ein Viertel der Größe derer von *Lex*. Normalerweise benötigt *Rex* nur 1/10 der Zeit von *Lex* zum Generieren eines

lexikalischen Analysators.

5.2. Lark

Lark ist ein Parser-Generator der primär LALR(1)- und LR(1)-Grammatiken verarbeitet [Gro88, Grob]. Mit der vorhandenen Backtracking-Einrichtung können auch Parser für wesentlich mächtigere Grammatikklassen generiert werden. Die Grammatikregeln können mit semantischen Aktionen versehen werden, die direkt in einer Zielsprache formuliert sind. Immer wenn der erzeugte Parser eine Grammatikregel erkennt, wird die zugehörige semantische Aktion ausgeführt. Der Generator stellt einen Mechanismus zur S-Attributierung zur Verfügung, d. h. synthetisierte Attribute können während der Zerteilung berechnet werden. Als Testhilfen gibt es eine Protokollierung der Parserschritte und eine graphische Visualisierung auf der Basis von X-Windows.

Im Falle von LR-Konflikten liefert *Lark* nicht wie andere Generatoren nur Information über aus Mengen von Situationen bestehende Zustände, sondern druckt einen Ableitungsbaum, der wesentlich nützlicher zur Analyse des Konflikts ist. Konflikte können durch die Angabe von Priorität und Assoziativität für Operatoren und Produktionen oder durch die Angabe von sogenannten syntaktischen und semantischen Prädikaten gelöst werden. Die generierten Parser beinhalten eine automatische Fehlerbehandlung mit Fehlermeldungen und -reparatur. Zur Fehlerbehandlung wird die vollständige, rücksetzungsfreie Methode von Röhrich [Röh76, Röh80, Röh82] verwendet. Die Parser sind tabellengesteuert und wie im Falle von *Rex* werden die Tabellen mit der Kammvektortechnik komprimiert. Der Generator verwendet den in [DeP82] beschriebenen Algorithmus zur Berechnung der Vorschaumengen. Zur Zeit können Zerteiler in den Zielsprachen C, C++, Modula-2, Ada und Eiffel erzeugt werden.

Mit *Lark* erzeugte Parser sind zwei bis drei mal schneller als mit *Yacc* [Joh75] erzeugte. Sie erreichen eine Geschwindigkeit von 2 Millionen Zeilen pro Minute ohne Berücksichtigung der lexikalischen Analyse. Die Größe der Parser ist gegenüber *Yacc* leicht erhöht, denn die Geschwindigkeit und die komfortable Fehlerbehandlung sind nicht ganz umsonst.

Die Eingabesprachen von *Rex* und *Lark* sind hinsichtlich der Syntax gegenüber *Lex* und *Yacc* lesbarer gestaltet. Mit Hilfe zweier, hier nicht näher beschriebener Präprozessoren können *Rex* und *Lark* auch Eingaben für *Lex* und *Yacc* verarbeiten. Dadurch sind unsere Werkzeuge in Bezug auf die Benutzerschnittstelle kompatibel mit den UNIX-Werkzeugen.

5.3. Ell

Ell ist ein LL(1) Parser-Generator, der Grammatiken, die in extended BNF geschrieben sind, verarbeitet [Gro88, Gro90a, GrV]. Während der Zerteilung kann eine L-Attributierung ausgewertet werden. Die erzeugten Parser beinhalten eine automatische Fehlerbehandlung mit Fehlermeldungen und -reparatur wie *Lark*. Die Parser sind nach dem Verfahren des rekursiven Abstiegs implementiert und erzielen ebenfalls eine Geschwindigkeit von 2 Millionen Zeilen pro Minute auf einer SPARC station ELC. Die möglichen Zielsprachen sind C, C++ und Modula-2.

5.4. Ast

Ast ist ein Generator für abstrakte Syntaxbäume [Gro91, Groa]. Er generiert Programm-Module oder abstrakte Datentypen zur Bearbeitung attributierter Bäume. Neben Bäumen können auch attributierte Graphen bearbeitet werden. Den Knoten dieser Datenstrukturen können beliebig viele Attribute von beliebigem Typ zugeordnet werden. Die Spezifikationen für dieses Werkzeug basieren auf erweiterten kontextfreien Grammatiken. Sie können als gemeinsame Notation sowohl für konkrete und abstrakte Syntax als auch für attributierte Bäume und Graphen betrachtet werden. Ein

Erweiterungsmechanismus stellt einfache und mehrfache Vererbung zur Verfügung. Intern werden die Bäume durch verzeigte Verbunde gespeichert. Zahlreiche Operationen für Bäume und Graphen können auf Anforderung von *Ast* erzeugt werden: Sogenannte Knotenkonstruktoren kombinieren Aggregatschreibweise mit Speicherverwaltung. Lese- und Schreibprozeduren übertragen Graphen aus/in Dateien in lesbarem ASCII- oder internem Binärformat. Die Reihenfolge von Teilbäumen in einer Liste kann umgekehrt werden. Es werden Prozeduren für häufig benutzte Traversierungsstrategien wie *top down* oder *bottom up* zur Verfügung gestellt. Ein interaktiver *Graph-Browser* erlaubt die Inspektion von Graphen in lesbarer Weise und unterstützt so den Programmtest.

5.5. Ag

Ag ist ein Generator für Attributauswerter [Gro90b, Grob, Groc]. Er verarbeitet geordnete Attributgrammatiken (OAGs) [Kas80], wohl-definierte Attributgrammatiken (WAGs) und sogenannte *higher order* Attributgrammatiken (HAGs) [VSK89, Vog93]. Er basiert auf der abstrakten Syntax oder genauer gesagt auf den von *Ast* erzeugten Baummodulen. Deshalb ist die Baumstruktur völlig bekannt. Den Terminalen und Nichtterminalen können beliebig viele Attribute zugeordnet werden. Diese werden mit den Typen der Zielsprache getypt. Dabei sind auch baumwertige Attribute möglich. *Ag* erlaubt regellokale Attribute und bietet einen Erweiterungsmechanismus an, welcher einfache und mehrfache Vererbung für Attribute und Attributberechnungen zur Verfügung stellt. Dieser gestattet ebenfalls die Elimination von Kettenregeln. Die Attributberechnungen werden in der Zielsprache formuliert und sollten in einem funktionalen Stil gehalten sein. Es ist möglich externe Funktionen von getrennt übersetzten Modulen aufzurufen. Die Verwendung nicht-funktionaler Anweisungen und von Seiteneffekten ist möglich, verlangt allerdings sorgfältige Überlegung. Die Syntax der Spezifikationssprache ist im Hinblick auf die Unterstützung kompakter, modularer und lesbarer Dokumente entworfen worden. Eine Attributgrammatik kann aus mehreren Modulen bestehen, wobei die kontextfreie Grammatik nur einmal spezifiziert wird. Es gibt Kurzschreibweisen für Kopierregeln und gefädelte Attribute womit viele triviale Attribut-Berechnungen weggelassen werden können. Die erzeugten Attributauswerter sind sehr effizient, da sie unter Verwendung von rekursiven Prozeduren direkt codiert sind. Die Speicherung der Attribute wird optimiert indem Attribute als lokale Variable und Prozedurparameter implementiert werden, wenn ihre Lebenszeit innerhalb eines Besuches liegt.

5.6. Puma

Puma ist ein Werkzeug zur Transformation und Manipulation von attributierten Syntaxbäumen [Gro92, Grod]. Es basiert auf Pattern-Matching und rekursiven Prozeduren. Die erzeugten Transformations-Module haben als Eingabe einen attributierten Baum. Eine reine Transformation bildet diesen auf eine Ausgabe beliebiger Art ab. Die Ausgabe kann ein neuer Baum sein, eine lineare Zwischensprache wie z. B. P-Code, ein Quellprogramm z. B. in Pascal, Assemblercode, Binärcode oder eine Folge von Prozeduraufrufen. Bei einer Modifikation wird der Eingabebaum verändert, wobei die Möglichkeiten von der Berechnung und Speicherung von Attributen bis hin zur Änderung der Baumstruktur reichen. Anwendungsgebiete für Transformationen sind die semantische Analyse, die Erzeugung von Zwischensprachen aus abstrakten Syntaxbäumen, Optimierer für interne Baumstrukturen jeden Niveaus, Quelle-Quelle-Übersetzung und Code-Generierung. *Puma* arbeitet mit dem Werkzeug *Ast* zusammen, wodurch bereits die Definition, Erzeugung und Speicherung von Bäumen unterstützt werden. *Puma* fügt eine kompakte Schreibweise für die Analyse und Synthese von Bäumen hinzu. Das Pattern-Matching kann als die Beschreibung von Entscheidungstabellen angesehen werden.

Die Spezifikation einer Transformation ist regelbasiert. Eine Regel besteht im wesentlichen aus einem oder mehreren Mustern, welche Baumfragmente beschreiben, Bedingungen und einer Folge von Anweisungen. Wenn die Muster mit den als Parameter übergebenen Bäumen zusammenpassen und die Bedingungen erfüllt sind, dann werden die zugehörigen Anweisungen ausgeführt. Es können mehrere Transformationen spezifiziert werden. Die Teilbäume eines Musters können in beliebiger Reihenfolge transformiert werden. Sie können mehrmals mit der selben oder mit verschiedenen Transformationen bearbeitet werden.

Puma erlaubt die implizite Deklaration von Variablen, führt eine Typprüfung in Bezug auf Bäume durch und überprüft die "single assignment"-Beschränkung für Variablen. Die Ausgabe ist der Quellcode eines Programm-Moduls in einer der Zielsprachen C, C++, oder Modula-2. Dieses Modul implementiert die spezifizierten Transformations-Routinen. Es kann leicht in beliebigen Programm-Code integriert werden. Die erzeugten Routinen sind optimiert durch Elimination von gemeinsamen Teilausdrücken und Elimination von "tail"-Rekursion. Das Pattern-Matching ist durch direkten Code implementiert und dadurch effizient realisiert.

5.7. Reuse

Reuse ist eine Bibliothek wiederverwendbarer Module hauptsächlich für den Einsatz im Übersetzerbau [Groe, Grof]. Sie enthält Module oder abstrakte Datentypen, die fast in jedem Übersetzer gebraucht werden:

- eine dynamische Speicherverwaltung
- ein Modul für dynamische und flexible Felder
- ein Modul zur Speicherung variabel langer Zeichenketten
- ein Modul zur Zeichenkettenbearbeitung
- eine Bezeichnertabelle, welche Zeichenketten unter Verwendung eines Hashverfahrens eindeutig auf ganze Zahlen abbildet
- Module für oft verwendete Datenstrukturen wie Mengen von ganzen Zahlen oder binäre Relationen zwischen ganzen Zahlen ohne Beschränkung des Definitionsbereichs.

6. Erfahrungsbericht

Eine erste große und realistische Anwendung des Werkzeugkastens war die Generierung eines Modula-2 nach C Übersetzers [Mar90]. Das *Mtc* genannte Programm übersetzt Modula-2 Programme in lesbaren C Code ohne Einschränkung (sogar geschachtelte Prozeduren und Module). Es ist weitgehend automatisch generiert und folgt dem in Abschnitt 4 vorgeschlagenen Übersetzer-Modell. Anstelle einer Zwischensprache erzeugt das Programm C Code und benötigt deshalb keinen Codegenerator zur Ausgabe von Maschinencode. Es enthält so viel von der semantischen Analyse wie für die Aufgabe gebraucht wird. Die semantische Analyse ist relativ vollständig und enthält die Behandlung der Gültigkeitsbereiche, Namensanalyse und Typbestimmung. Es fehlt die Überprüfung von Kontextbedingungen, da davon ausgegangen wird, daß nur korrekte Programme übersetzt werden. Tabelle 1 enthält die Größen der Spezifikationen und der generierten Quell-Module. Der Entwurf und die Implementierung von *Mtc* wurden im Rahmen einer Diplomarbeit mit einem Aufwand von 6 Mannmonaten durchgeführt. Das Programm ist stabil und es übersetzt regelmäßig mehr als 100.000 Zeilen Modula-2 nach C.

| Phase | Spezifikation | | | Quell-Modul | | | Werkzeug | |
|---------------|---------------|------|-------|-------------|-------|-------|----------|---------------|
| | formal | Code | Summe | Def. | Impl. | Summe | Name | Referenzen |
| Lex. Analyse | 392 | 133 | 525 | 56 | 1320 | 1376 | Rex | [Gro88, Groa] |
| Syntaxanalyse | 951 | 88 | 1039 | 81 | 3007 | 3088 | Ell | [Gro88, GrV] |
| Syntaxbaum | 189 | 51 | 240 | 579 | 2992 | 3571 | Ast | [Groa] |
| Symboltabelle | 115 | 938 | 1053 | 413 | 1475 | 1888 | Ast | [Groa] |
| Sem. Analyse | 1886 | 151 | 2037 | 9 | 3288 | 3297 | Ag | [Grob] |
| Codegenerator | 2793 | 969 | 3762 | 47 | 7309 | 7356 | Puma | [Gro92] |
| Wiederverw. | - | - | - | 819 | 2722 | 3541 | Reuse | [Groe, Grof] |
| Sonstiges | - | - | - | 698 | 3153 | 3851 | | |
| Summe | 6326 | 2330 | 8656 | 2702 | 25266 | 27968 | | |

Tabelle 1: Umfang der Spezifikationen und der Quellmodule von *Mtc*

Die Größe des Binärprogramms beträgt 585 K Bytes. Es läuft mit einer Geschwindigkeit von etwa 2900 Grundsymbolen (*tokens*) pro Sekunde oder 715 Zeilen pro Sekunde auf einem SUN SPARC station ELC. Diese Zahlen berücksichtigen nur die Größe der übersetzten Module. Wenn man zusätzlich die (transitiv) importierten Definitionsmodule berücksichtigt, die ebenfalls lexikalisch, syntaktisch und semantisch analysiert werden, so erreicht man 4700 Grundsymbole pro Sekunde oder 1790 Zeilen pro Sekunde. Zum Vergleich die Zahlen für zwei Übersetzer des Rechner-Herstellers: Der C-Übersetzer läuft mit einer Geschwindigkeit von 390 Zeilen pro Sekunde (ohne *Header*-Dateien) bzw. 650 Zeilen pro Sekunde (mit *Header*-Dateien) und der Modula-2-Übersetzer mit 190 Zeilen pro Sekunde. Die Laufzeit von *Mtc* ist folgendermaßen verteilt:

| | |
|----------------------------------|------|
| lex. + syn. Analyse + Baumaufbau | 42 % |
| semantische Analyse | 33 % |
| C Codegenerierung | 25 % |

Die semantische Analyse verbringt 95% der Zeit mit der Berechnung von Attributen mittels vom Benutzer spezifizierten Anweisungen und nur 5% für die Baumtraversierung bzw. für Besuchaktionen. Für 11 Knotentypen sind fünf Besuche notwendig.

Mtc braucht ungefähr 300 K Bytes dynamischen Speicher pro 1000 Quellzeilen zur Speicherung des abstrakten Syntaxbaums, der Attribute und der Symboltabelle ohne Optimierung der Attributspeicherung. Dies ist bei den heutigen Speicherkapazitäten problemlos möglich. Es zeigt, daß es im Gegensatz zu der in der Literatur vertretenen Meinung möglich ist, alle Attribute im Baum zu speichern. Wir tun dies sogar mit dem sogenannten Umgebungsattribut. Dies wird möglich, indem wir die Symboltabelle als abstrakten Datentyp in der Zielsprache programmieren. Die Implementierung ist zeit- und speichereffizient durch die Ausnutzung von Zeigersemantik und *structure sharing*.

7. Zusammenfassung und Ausblick

Wir haben einen Werkzeugkasten mit Übersetzerbau-Werkzeugen vorgestellt, womit sich Übersetzer für Programmiersprachen weitgehend automatisch generieren lassen. Die Übersetzerbau-Werkzeuge unterstützen die Konstruktion nahezu aller Übersetzerphasen. Die Werkzeuge sind sehr mächtig, flexibel und weitgehend unabhängig von einander. Besonders hervorzuheben sind die

praktische Brauchbarkeit der Werkzeuge, der deutlich reduzierte Erstellungsaufwand für Übersetzer und die hohe Qualität der generierten Übersetzer. Von der Effizienz her sind die Werkzeuge konkurrenzfähig zur Programmierung von Hand. Sie unterstützen einen breiten Bereich von Übersetzerstrukturen und erlauben die Konstruktion von Übersetzern mit Produktionsqualität. Viele realistische Anwendungen zeigen die ausgezeichnete Leistungsfähigkeit der Werkzeuge.

Die Übersetzerbau-Werkzeuge eignen sich für viele Aufgabenstellungen, die über die Konstruktion von reinen Übersetzern hinausgehen. Sie gestatten beispielsweise die Implementierung von Präprozessoren, die Spracherweiterungen und Sprachdialekte auf Standardsprachen abbilden. Wie es das Anwendungsbeispiel zeigt, lassen sich Umsetzer von einer Quellsprache in eine andere erstellen. Weiterhin ist etwa die Generierung von Prüfprogrammen für Programmierkonventionen möglich.

Die Werkzeuge lassen sich immer wieder verbessern und vervollständigen: Zur Zeit wird beispielsweise der Parser-Generator *Lark* um die Verarbeitung von LR(k) Grammatiken erweitert. Der Generator für abstrakte Syntaxbäume *Ast* wird um einen graphischen Browser für Bäume und Graphen ergänzt. Die Optimierungsphase eines Übersetzers sollte selbstverständlich auch unterstützt werden. Dies kann entweder durch einen wiederverwendbaren sprachunabhängigen Optimierer, durch einen parameterisierbaren Optimierer oder durch einen Optimierergenerator geschehen.

Danksagung

Wir danken allen die zur Entstehung des Werkzeugkastens und dieses Aufsatzes durch aktive Mitarbeit oder durch ihre Ideen beigetragen haben: Michael Besser, Carsten Gerlhof, Bob Gray, Eduard Klein, Rudolf Landwehr, Matthias Martin, Thomas Müller, F. W. Schröer, Dirk Schwartz-Hertzner, Doris Vielsack, Bertram Vielsack und William M. Waite.

Literatur

- [ASU86] A. V. Aho, R. Sethi and J. D. Ullman, *Compilers: Principles, Techniques, and Tools*, Addison Wesley, Reading, MA, 1986.
- [DeP82] F. DeRemer and T. J. Pennello, Efficient Computation of LALR(1) Look-Ahead Sets, *ACM Trans. Prog. Lang. and Systems* 4, 4 (Oct. 1982), 615-649.
- [Gro88] J. Grosch, Generators for High-Speed Front-Ends, *LNCS 371*, (Oct. 1988), 81-92, Springer Verlag.
- [Gro89] J. Grosch, Efficient Generation of Lexical Analysers, *Software—Practice & Experience* 19, 11 (Nov. 1989), 1089-1103.
- [Gro90a] J. Grosch, Efficient and Comfortable Error Recovery in Recursive Descent Parsers, *Structured Programming 11*, (1990), 129-140.
- [Gro90b] J. Grosch, Object-Oriented Attribute Grammars, in *Proceedings of the Fifth International Symposium on Computer and Information Sciences (ISCIS V)*, A. E. Harmanci and E. Gelenbe (ed.), Cappadocia, Nevsehir, Turkey, Oct. 1990, 807-816.
- [Gro91] J. Grosch, Tool Support for Data Structures, *Structured Programming 12*, (1991), 31-38.
- [Gro92] J. Grosch, Transformation of Attributed Trees Using Pattern Matching, *LNCS 641*, (Oct. 1992), 1-15, Springer Verlag.

- [Groa] J. Grosch, Rex - A Scanner Generator, Cocktail Document No. 5, CoCoLab Germany.
- [Grob] J. Grosch, Lark - An LR(1) Parser Generator With Backtracking, Cocktail Document No. 32, CoCoLab Germany.
- [GrV] J. Grosch and B. Vielsack, The Parser Generators Lalr and Ell, Cocktail Document No. 8, CoCoLab Germany.
- [Groa] J. Grosch, Ast - A Generator for Abstract Syntax Trees, Cocktail Document No. 15, CoCoLab Germany.
- [Grob] J. Grosch, Ag - An Attribute Evaluator Generator, Cocktail Document No. 16, CoCoLab Germany.
- [Groc] J. Grosch, Multiple Inheritance in Object-Oriented Attribute Grammars, Cocktail Document No. 28, CoCoLab Germany.
- [Grod] J. Grosch, Puma - A Generator for the Transformation of Attributed Trees, Cocktail Document No. 26, CoCoLab Germany.
- [Groe] J. Grosch, Reusable Software - A Collection of Modula-Modules, Cocktail Document No. 4, CoCoLab Germany.
- [Grof] J. Grosch, Reusable Software - A Collection of C-Modules, Cocktail Document No. 30, CoCoLab Germany.
- [Joh75] S. C. Johnson, Yacc — Yet Another Compiler-Compiler, Computer Science Technical Report 32, Bell Telephone Laboratories, Murray Hill, NJ, July 1975.
- [Kas80] U. Kastens, Ordered Attribute Grammars, *Acta Inf.* 13, 3 (1980), 229-256.
- [Les75] M. E. Lesk, LEX — A Lexical Analyzer Generator, Computing Science Technical Report 39, Bell Telephone Laboratories, Murray Hill, NJ, 1975.
- [Mar90] M. Martin, Entwurf und Implementierung eines Übersetzers von Modula-2 nach C, Diplomarbeit, GMD Forschungsstelle an der Universität Karlsruhe, Feb. 1990.
- [Röh76] J. Röhrich, Syntax-Error Recovery in LR-Parsers, in *Informatik-Fachberichte*, vol. 1, H.-J. Schneider and M. Nagl (ed.), Springer Verlag, Berlin, 1976, 175-184.
- [Röh80] J. Röhrich, Methods for the Automatic Construction of Error Correcting Parsers, *Acta Inf.* 13, 2 (1980), 115-139.
- [Röh82] J. Röhrich, Behandlung syntaktischer Fehler, *Informatik Spektrum* 5, 3 (1982), 171-184.
- [VSK89] H. H. Vogt, S. D. Swierstra and M. F. Kuiper, Higher Order Attribute Grammars, *SIGPLAN Notices* 24, 7 (July 1989), 131-145.
- [Vog93] H. H. Vogt, *Higher Order Attribute Grammars*, PhD Thesis, University of Utrecht, Feb. 1993.

Inhalt

| | | |
|------|------------------------------------|----|
| | Übersicht | 1 |
| 1. | Aufbau eines Übersetzers | 1 |
| 2. | Werkzeugkasten | 2 |
| 3. | Gemeinsame Eigenschaften | 3 |
| 4. | Übersetzer-Modell | 5 |
| 5. | Die Werkzeuge | 5 |
| 5.1. | Rex | 5 |
| 5.2. | Lark | 6 |
| 5.3. | Ell | 6 |
| 5.4. | Ast | 6 |
| 5.5. | Ag | 7 |
| 5.6. | Puma | 7 |
| 5.7. | Reuse | 8 |
| 6. | Erfahrungsbericht | 8 |
| 7. | Zusammenfassung und Ausblick | 9 |
| | Danksagung | 10 |
| | Literatur | 10 |