

Projektarbeit Compilerbau

Jonas Wenner, Jens Marvin Reuter, Noah Hoffmann

27. September 2018

Inhaltsverzeichnis

| | | |
|----------|--|-----------|
| 1 | Anforderungen | 5 |
| 1.1 | Infrastruktur | 5 |
| 1.2 | Programmiersprache | 5 |
| 1.3 | Dokumentation | 6 |
| 1.4 | Präsentation | 6 |
| 2 | Theoretischer Hintergrund | 7 |
| 2.1 | Verallgemeinerung | 7 |
| 2.2 | Lexikalische Analyse | 7 |
| 2.3 | Syntaktische Analyse | 7 |
| 2.4 | Syntaxgesteuerte Übersetzung | 8 |
| 3 | Konzeption der Entwicklung | 9 |
| 3.1 | Allgemeiner Grundsatz | 9 |
| 3.2 | Wahl der Entwicklungswerkzeuge | 9 |
| 3.3 | Herkunft des Namens C= | 9 |
| 3.4 | Sprachdesign | 9 |
| 3.5 | Implementierung | 10 |
| 4 | Verwendete Technologien | 11 |
| 4.1 | Eclipse | 11 |
| 4.2 | TestNG | 11 |
| 4.3 | ANTLR | 11 |
| 4.4 | Jasmin | 12 |
| 4.5 | Java | 12 |
| 4.6 | Java Virtual Machine | 12 |
| 5 | Vorgehen während der Entwicklung | 13 |
| 5.1 | Erste Schritte: Auswertung von arithmetischen Ausdrücken . . | 13 |
| 5.2 | Ausweiten der simplen Grammatik | 21 |
| 5.2.1 | Ergebnisausgabe | 21 |
| 5.3 | Implementierung von Variablen | 22 |
| 5.4 | Implementierung von Funktionen | 22 |
| 5.5 | Implementierung von Konstanten | 22 |
| 5.6 | Implementierung von Bedingten Verzweigungen | 23 |
| 5.7 | Implementierung einer Schleife | 23 |

| | | |
|----------|---|-----------|
| 6 | Aufgetrene Probleme und deren Lösung | 24 |
| 6.1 | Operatorenpriorität bei arithmetischen Ausdrücken | 24 |
| 6.2 | Variablen | 24 |
| 6.2.1 | Redefinition von bereits definierte Variable | 24 |
| 6.2.2 | Zugriff auf undefinierte Variable | 25 |
| 6.3 | Funktionen | 25 |
| 6.3.1 | Zugriff auf undefinierte Funktion | 25 |
| 6.4 | Redefinition von bereits definierter Funktion | 25 |
| 6.4.1 | Gültigkeitsbereiche | 25 |
| 6.4.2 | Funktionen mit gleichem Namen und unterschiedlichen Signaturen | 25 |
| 6.5 | Vorwärtsdeklarationen | 25 |
| 6.6 | Bedingte Verzweigungen | 25 |
| 6.6.1 | Umsetzung in Jasmin mit Hilfe von Labels und Sprung- befehlen | 25 |
| 7 | Verwendung des Compilers | 26 |
| 7.1 | Verwendung der UnsereCompilerbauSprache | 26 |
| 7.2 | Aufrufen des UnserCompilerbauCompilers | 26 |
| 8 | Anhang | 27 |
| 8.1 | Quellenangaben | 27 |
| 8.2 | Quellcode des Compiler(?) | 27 |

Zusammenfassung

Als Abschlussprojekt des Moduls 'Compilerbau' unter Leitung von Thorsten Jakobs, M.Sc., htw saar wurde dieser Compiler zur von uns eigens entwickelten Sprache *C=* (Arbeitstitel) entwickelt, mit den in *1. Anforderungen* beschriebenen Voraussetzungen und Zielen. Die Sprache besitzt eine C-ähnliche Syntax, wurde jedoch auf wesentliche Elemente wie Grundrechenarten, Funktionen, if-Statements und eine Form der Schleife reduziert. Weiterhin vorhanden sind Bezeichner, Konstanten und Variablen. Die Klammerung wird ebenfalls berücksichtigt.

Diese Dokumentation behandelt die Anforderungen an den Compiler, den theoretischen Hintergrund, die Konzeption der Entwicklung, unsere verwendeten Technologien, sowie unsere Vorgehensweise während der Entwicklung und die Probleme, auf die wir gestoßen sind und unsere Lösungsansätze. Ausgangs wird auch die Handhabung des Compilers näher erläutert.

1 Anforderungen

Um zu gewährleisten, dass alle Modulteilnehmer über die gleichen Voraussetzungen verfügen und nach festgelegten Kriterien bewertet werden können, wurden durch den Modulverantwortlichen einige Rahmenbedingungen vorgegeben. Die Prüfleistung besteht dabei aus der Entwicklung einer einfachen Programmiersprache, dem Anfertigen einer Dokumentation sowie der Präsentation der Ergebnisse. Die Anforderungen an das Projekt werden in den folgenden Abschnitten aufgeführt und erläutert.

1.1 Infrastruktur

Das Ziel ist die Entwicklung eines Compiler-Frontends, welches Quelltext in der von uns entwickelten Programmiersprache zu Jasmin-Code übersetzt. Jasmin ist ein Assembler für die Java Virtual Machine, der Quelltext in einer Assembler-ähnlichen Syntax zu Bytecode übersetzt. Dieser Bytecode kann schließlich von der Java Virtual Machine ausgeführt werden.

1.2 Programmiersprache

Die von uns entwickelte Programmiersprache besitzt folgende Komponenten:

- Bezeichner
- Konstanten
- Variablen
- Arithmetik mit Grundrechenarten
- Klammerung
- Funktionen
- if-else-Anweisung
- while-Schleife

Die Entwicklung wird mit 50% in der Gesamtbewertung gewichtet. In einem späteren Kapitel dieser Dokumentation wird erläutert, wie die Programmiersprache verwendet wird.

1.3 Dokumentation

Es ist eine Dokumentation anzufertigen, die Aufschluss über die Konzeption der Entwicklung, die verwendeten Technologien, das Vorgehen während der Entwicklung, aufgetretene Probleme und Lösungsansätze zu diesen sowie die Verwendung des Compilers gibt.

Die Dokumentation wird mit 30% in der Gesamtbewertung gewichtet.

1.4 Präsentation

In einer circa 20-minütigen Präsentation werden die verwendeten Technologien, wichtige Entwicklungsschritte, sowie aufgetretene Probleme erläutert. Außerdem wird der Compiler vorgeführt.

Die Präsentation wird mit 20% in der Gesamtbewertung gewichtet.

2 Theoretischer Hintergrund

Um die folgenden Abschnitte verständlich zu machen ist eine kurze Erläuterung des theoretischen Hintergrunds sinnvoll.

2.1 Verallgemeinerung

Der zu entwickelnde Compiler nimmt als Eingabe eine Text-Datei mit Instruktionen in der Ausgangssprache. Der Compiler verarbeitet diese Eingabe und gibt als Ausgabe eine Datei mit Anweisungen in der Zielsprache zurück. Diese Verarbeitung kann grob in drei Abschnitte eingeteilt werden.

2.2 Lexikalische Analyse

Der erste Schritt ist dabei die sogenannte lexikalische Analyse. Die Eingabedatei ist für den Computer zunächst nur eine Kette von Zeichen. Diese Zeichen sind in sogenannte Tokens aufzuteilen. Beispielhaft könnte die Eingabe

a = b + c ;

zu den folgenden Tokens aufgelöst werden.

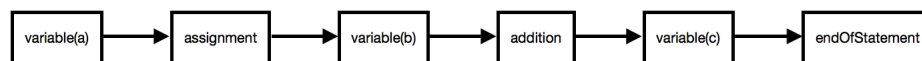


Abbildung 1: Kurzes Beispiel für lexikalische Analyse

Der Teil des Compilers, der diese Aufgabe übernimmt, wird *Lexer* genannt. Der Lexer überprüft für jedes Zeichen der Eingabedatei, welchem Token dieses zugeordnet werden kann. Dies geschieht auf Grundlage von definierten Regeln, welche mit sogenannten *regulären Ausdrücken* formuliert werden. Ein regulärer Ausdruck besteht aus einer Sequenz von Zeichen, die ein bestimmtes Muster definieren. Reguläre Ausdrücke können auch rekursiv verwendet werden, um auf Grundlage einfacher regulärer Ausdrücke komplexere zu bilden. Beispiele für reguläre Ausdrücke werden im Kapitel 5.1. *Erste Schritte: Auswertung von arithmetischen Ausdrücken* aufgeführt.

2.3 Syntaktische Analyse

Im zweiten Schritt der Verarbeitung wird die *syntaktische Analyse* durchgeführt. Dabei wird überprüft, ob die Eingabedatei einer definierten Grammatik entspricht, während bei der lexikalischen Analyse lediglich die übergebene

Zeichenfolge in Tokens aufgeteilt wird, nicht jedoch überprüft wird, ob diese der Grammatik entsprechend zusammengeführt werden können.

Bei einer erfolgreichen Auswertung wird die Folge an Tokens, die bei der lexikalischen Analyse ermittelt wurde, in einen Syntaxbaum übergeführt. Dies ist notwendig, um die Tokens in der richtigen Reihenfolge auszuwerten.

2.4 Syntaxgesteuerte Übersetzung

Die eigentlich Erzeugung von Code in der Ausgangssprache findet zuletzt statt. Hierbei wird jeder Knoten des Syntaxbaum betrachtet und je nach Kontext eine Zeichenkette mit Anweisungen in der Zielsprache generiert. In diesem Schritt findet ebenfalls die semantische Analyse statt, beispielsweise die Überprüfung, ob eine aufgerufene Variable vorher deklariert wurde.

3 Konzeption der Entwicklung

3.1 Allgemeiner Grundsatz

Während der Entwicklung versuchten wir uns möglichst kleine Zwischenziele zu setzen, um die Aufgabe in überschaubare, einfach zu lösende Teilprobleme zu unterteilen. Rückblickend war dies eine sinnvolle Strategie, die zum gewünschten Ziel führte.

3.2 Wahl der Entwicklungswerkzeuge

Bezüglich der Entwicklungswerkzeuge wurden uns verschiedene Möglichkeiten vorgestellt. Die erste Möglichkeit bestand darin das Tool "lex" beziehungsweise dessen Open-Source-Implementierung "flex" als Lexer-Generator und das Tool "yacc" beziehungsweise dessen Open-Source-Implementierung "bison" als Parser-Generator zu nutzen. Die zweite Möglichkeit bestand in der Nutzung des Werkzeuges ANTLR, das die Funktionalitäten der zuvor genannten Werkzeuge in einem Programm zusammenfasst.

Da ANTLR automatisiert auf Grundlage einer Datei, die die von uns formulierte Grammatik der Sprache enthält, Lexer sowie Parser zeitgleich erstellen kann, entschieden wir uns für diese Variante. Der Vorteil besteht dabei darin, dass zwei essentielle Teile des Compilers von einem Tool übernommen werden.

3.3 Herkunft des Namens C=

Für die Programmiersprache wurde der Name C= gewählt (gesprochen: C equal). Ursprünglich sollte die Programmiersprache C - - genannt werden, in Anlehnung an den verkleinerten Sprachumfang der Sprache C sowie der Namensgebung von C++. Da C - - jedoch bereits als Zwischensprache existiert, wurde die im Rahmen des Moduls Compilerbau entwickelte Sprache in Anlehnung an C# (# ist auch als vier aneinandergereihte Additionssymbole zu interpretieren) C= genannt, wobei das = vier aneinandergereihte Subtraktionszeichen darstellt.

3.4 Sprachdesign

Es wurde sich darauf geeinigt, dass die Syntax der von uns entworfenen Ausgangssprache der Syntax der Programmiersprache C ähnlich sein soll, da dies gleich mehrere Vorteile hervorbringt:

- Alle Projektteilnehmer sind mit der Sprache C vertraut

- C ist eine der am weitesten verbreiteten compilierten Sprachen
- Die Sprache C besitzt aufgrund seiner Klammerregeln, sowie anderer sprachlicher Regeln und Konventionen eine gute Lesbarkeit

3.5 Implementierung

Unser Compiler verwendet im Zusammenspiel mit dem von ANTLR generierten Lexer und Parser einen Visitor für die syntaxgesteuerte Übersetzung. Die Alternative dazu besteht in einem Listener, der sich von einem Visitor dahingehend unterscheidet, dass die verschiedenen Methoden die auf die Knoten des Syntaxbaumes angewandt werden, keinen Rückgabewert besitzen und die Ergebnisse dieser somit separat gespeichert werden müssen. Der Vorteil eines Listeners hingegen besteht darin, dass untergeordnete Knoten im Syntaxbaum automatisch bearbeitet werden und nicht wie beim Visitor explizit besucht werden müssen. Nach der Abwägung dieser Vor- und Nachteile entschieden wir uns für einen Visitor, da dieser uns als die effizientere Methode erschien.

4 Verwendete Technologien

4.1 Eclipse

Zur Entwicklung wurde die integrierte Entwicklungsumgebung Eclipse Photon in der aktuellen Version (4.8.0) genutzt. Eclipse bietet mehrere Vorteile, die ein einfacheres und effizienteres Entwickeln erlauben. Eclipse enthält einen Datei-Explorer, mit dessen Hilfe die Navigation durch die verschiedenen Verzeichnisse zusammen mit Editor und Konsole innerhalb eines Fensters geschehen kann. Des Weiteren bietet Eclipse automatische Vorschläge zur Code- und Textvervollständigung an, was die Schreibgeschwindigkeit und damit die Arbeitseffizienz erhöht. Außerdem kann Eclipse durch Plug-Ins erweitert werden. Als Plug-In nutzen wir beispielsweise das Test-Framework TestNG, um verschiedene Testfälle automatisiert zu prüfen.

4.2 TestNG

TestNG ist ein Framework, das dem zu testenden Programm Eingabewerte übergibt und die tatsächlichen Ergebnisse des Programms mit den erwarteten Ergebnissen abgleicht. Dabei können beliebig viele Testszenarien geprüft werden. Es werden sogenannte positive Tests, also solche, bei denen die Eingabe ein erwartetes Ergebnis hervorruft, als auch negative Tests, bei denen geprüft wird, ob nicht vorhergesehene Eingaben mit entsprechenden Fehlermeldungen korrekt behandelt werden, durchgeführt. TestNG wurde von uns in Version 6.10 (PRÜFEN!) genutzt.

Als Beispiel: Der Quelltext

```
println(1+4);
```

sollte als Ergebnis "5" auf der Konsole ausgeben. TestNG übergibt den Quelltext an unseren Compiler, ruft Jasmin auf, führt das Programm aus und gleicht dann die Ergebnisse ab. Wenn das Ergebnis "5" ist, gilt der Test als bestanden, andernfalls als durchgefallen.

4.3 ANTLR

ANTLR steht für *ANother Tool for Language Recognition* und ist ein Lexer- und Parsergenerator. ANTLR generiert auf Grundlage einer Grammatik-Datei Java-Code, der einen entsprechenden Lexer sowie ein Template für Teile des Parsers implementiert. Dadurch, dass die von ANTLR generierten Programme aus einer Eingabedatei in der Ausgangssprache einen Syntaxbaum erstellen können, besteht der größte Arbeitsaufwand daraus, Gramma-

tiken für ANTLR zu formulieren und Teile des Parsers (in unserem Fall ein Visitor) zu implementieren. Es wurde die ANTLR-Version 4.7.1 verwendet.

4.4 Jasmin

Jasmin ist ein Assembler für die Java Virtual Machine. Dabei werden Textdateien mit Anweisungen in Jasmin-Syntax zu Bytecode übersetzt. Dieser Bytecode kann von der JVM ausgeführt werden.

4.5 Java

Die Programmiersprache Java wurde für das Projekt gewählt, da der von ANTLR generierte Code ebenfalls in Java ist. Es ist teilweise notwendig von Klassen des ANTLR generierten Code abzuleiten, weshalb die Wahl einer anderen Programmiersprache als Java Problemstellungen wie die Code-Kompatibilität zur anderen Programmiersprache darstellt. Des Weiteren muss die Java Runtime Environment ohnehin genutzt werden, um ANTLR, Jasmin sowie schließlich die Kompilate des C-Compilers auszuführen.

4.6 Java Virtual Machine

Die Java VM ist ein Zwischenschritt beim Ausführen von Java-Code. Eine Java-Quelldatei (.java) wird zunächst durch den Java-Compiler zu Bytecode (.class) übersetzt und dann von der Java VM interpretiert. Dieser Zwischenschritt ermöglicht eine Plattformunabhängigkeit, da die Kompilate (.class-Dateien) nicht maschinenspezifisch übersetzt werden. Plattformabhängiger Maschinencode wird erst von der Java VM generiert, weshalb jeder Bytecode ausgeführt werden kann, solange für die entsprechende Maschine die Java VM verfügbar ist.

In diesem Projekt wird die Java VM jedoch nicht nur zur Übersetzung unseres Code genutzt, sondern auch um die Kompilate unseres eigenen Compilers auszuführen. Diese Kompilate werden von Jasmin zu Bytecode übersetzt, der Bytecode wiederum wird von der Java VM ausgeführt.

Zu beachten ist, dass die Java VM stack-basiert operiert und nicht wie beispielsweise der x86-Befehlssatz mit Registern arbeitet, was bei der Implementierung zu beachten ist und eine besondere Denkweise erfordert.

5 Vorgehen während der Entwicklung

5.1 Erste Schritte: Auswertung von arithmetischen Ausdrücken

Der Vorgehensweise unter 3.1. *Allgemeiner Grundsatz* entsprechend, implementierten wir zuerst die unserer Sicht nach grundlegendste Fähigkeit einer Programmiersprache: das Auswerten von arithmetischen Ausdrücken. Es soll möglich sein, Ausdrücke, die die Grundrechenarten sowie eine beliebig tiefe Schachtelung von Klammern enthalten, zu erkennen und der Operatorenpriorität entsprechend die gelesenen Tokens in einem Syntaxbaum zusammenzuführen. Dazu formulierten wir folgende Grammatik:

```
//Datei: Arithmetic.g4
grammar Arithmetic;

////////////////////////////////////////
// beliebige Folge der Ziffern 0 bis 9
INTEGER : [0-9]+ ;
// ueberspringt Leerzeichen , Tabstops sowie Linefeeds
WS : [ \t\r\n]+ -> skip ;
// oeffnende runde Klammer
LPAREN : '(';
// schliessende runde Klammer
RPAREN : ')';
////////////////////////////////////////
//mathematische Operatoren
PLUSOP : '+';
MINOP : '-';
MULTOP : '*';
DIVOP : '/';
////////////////////////////////////////
//Regeln fuer math. Ausdruecke
expression: INTEGER
        | LPAREN expression RPAREN
        | expression DIVOP expression
        | expression MULTOP expression
        | expression MINOP expression
        | expression PLUSOP expression
        ;
```

////////////////////////////////////

Eine Grammatik für ANTLR hat folgenden Aufbau:

Die Definition der Grammatik beginnt mit dem Schlüsselwort "grammar" sowie dem Namen der Grammatik. Dabei gilt es zu beachten, dass die Datei den gleichen Namen wie die Grammatik selbst sowie die Dateierdung ".g4" besitzt. Diese Deklaration wird mit einem Semikolon abgeschlossen.

Zeilen, die mit "//" beginnen, sind Kommentare und haben keinen Einfluss auf die Grammatik. Sie dienen lediglich als Erläuterungen und zur Formatierung, um die Lesbarkeit zu verbessern.

Auf die Deklaration dürfen beliebig viele Regeln für die Grammatik folgen. Zur Formulierung der Regeln werden reguläre Ausdrücke genutzt. Beispielsweise besagt die sechste Zeile, dass es eine Regel INTEGER gibt, wobei ein INTEGER sich aus einer beliebigen Folge der Ziffern von 0 bis 9 zusammensetzt. Der reguläre Ausdruck [0-9] gibt an, dass ein beliebiges Zeichen im Bereich von 0 bis 9 vorkommen darf. Das abschließende "+" bedeutet, dass es sich um eine Kette dieser Zeichen, die beliebig lange ist, jedoch mindestens die Länge 1 besitzt, handelt.

Die Regel WS (Whitespace) besagt, dass bestimmte Zeichen, die nur der Formatierung dienen, ignoriert werden, da sie für das Übersetzen einer Quelldatei keine Bedeutung haben.

Die darauffolgenden Regeln sind Aliase für die Zeichen, die in arithmetischen Ausdrücken verwendet werden. Diese Auslagerung steigert unseres Erachtens nach die Lesbarkeit der letzten Regeln dieses Beispiels, sind aber nicht zwingend notwendig.

Die wohl relevanteste Regel trägt den Namen "expression" und legt fest, wie sich ein arithmetischer Ausdruck zusammensetzen kann. Im Vergleich zu den vorherigen Regeln, wurde hier von der Möglichkeit, mehrere alternative Möglichkeiten anzugeben, Gebrauch gemacht. Die Regel besagt, dass ein arithmetischer Ausdruck entweder aus

- einer Zahl,
- einem geklammerten Ausdruck,
- einer Division mit zwei Operanden,
- einer Multiplikation mit zwei Operanden,
- einer Subtraktion mit zwei Operanden
- oder einer Addition mit zwei Operanden

Damit diese Priorität gewährleistet werden kann, sind die Regeln in dieser bestimmten Reihenfolge notiert. ANTLR versucht immer zuerst die "oberste" Regel anzuwenden, daraufhin die darunter stehende usw.. Deshalb wird der folgende Ausdruck

zu diesem Syntaxbaum aufgelöst:

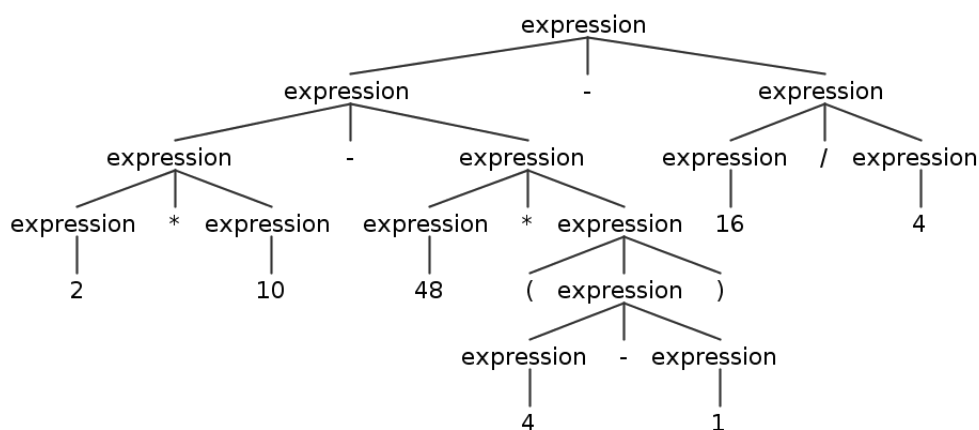


Abbildung 2: Syntaxbaum für beispielhaften arithmetischen Ausdruck

Des Weiteren ist zu beachten, dass arithmetischen Ausdrücke gleicher Operatorenpriorität von links nach rechts ausgewertet werden müssen. Als Beispiel lässt sich hier der Term

nennen. Sowohl die Subtraktion als auch die Addition besitzen die gleiche Bindung. Wenn man hier fälschlicherweise zuerst die Addition, also $5+2$, was 7 ergibt, auswertet und das Ergebnis dieser Operation von 8 subtrahiert, erhält man als Gesamtergebnis 1. In der korrekten Reihenfolge erfolgt zuerst die Subtraktion, also $8-5$, was 3 ergibt, und erst darauf die Addition von 1, was als Gesamtergebnis 4 liefert. Um diese Fehlerquelle auszuschließen, sind in unserer Grammatik die Regel für Subtraktion vor der für Addition und die

Regel für Division vor der für Multiplikation definiert. Da bei reinen Additionen bzw. reinen Multiplikationen die Auswertungsreihenfolge tatsächlich keine Rolle spielt, eine Subtraktion jedoch vor einer Additionen (vgl. Beispiel oben) ausgewertet werden muss, sorgt die Reihenfolge der Regeln in der Grammatik für eine korrekte Auswertung.

Ein analoges Beispiel zu Divisionen und Multiplikationen ist

| |
|---------|
| $8/2*4$ |
|---------|

Auch hier ergibt sich ein ähnliches Problem wie beim vorherigen Beispiel. Wird zuerst die Multiplikation ($2*4$) ausgeführt und erst danach die Division (also $8/8$ in diesem Fall), ist das Endergebnis 1 und nicht wie in der richtigen Reihenfolge 16.

Nach diesen Schritten sind wir also in der Lage einen Syntaxbaum zu erstellen. Der nötige Programmcode dafür wird von ANTLR automatisiert erstellt. Dabei ist die Eingabe dieses Programmes der auszuwertende Ausdruck. Die Ausgabe ist der Syntaxbaum, der zu Testzwecken auch als Grafik (vgl. Abbildung oben) ausgegeben werden kann.

Der nächste wichtige Schritt der Übersetzung besteht nun in der Auswertung dieses Baumes. Dazu wird der Baum als Datenstruktur betrachtet. ANTLR liefert dabei mehrere Methoden, die auf Instanzen der Klasse "ParseTree" angewandt werden können. Jedes Token, das in der Grammatik definiert wurde und durch einen Knoten im Baum repräsentiert wird, besitzt eine sogenannte Visit-Methode. Diese Methode gibt eine Kette mit Zeichen zurück, wobei diese Zeichenkette die Anweisungen in der Zielsprache (Jasmin) enthält. Das Abarbeiten dieser Visit-Methoden in der richtigen Reihenfolge bildet die Grundlage für die Übersetzung, da hiernach alle Instruktionen in der Zielsprache zusammengesetzt sind. Was genau in einer Visit-Methode passiert, wird vom Entwickler festgelegt. ANTLR stellt so gesehen nur eine Vorlage zur Verfügung.

Um diese korrekte Reihenfolge zu gewährleisten, muss eine Anfangsregel (Startaxiom) festgelegt sein. In diesem Beispiel wurde festgelegt, dass der Programmstart - sprich der Wurzelknoten des Baumes - eine "expression" sein muss. Das bedeutet, dass zunächst die Visit-Methode des Wurzelknoten aufgerufen wird. Damit nun auch die inneren Knoten des Syntaxbaums berücksichtigt werden, ist der Aufruf einer weiteren von ANTLR generierten Methode notwendig. Jeder Knoten des Baums besitzt eine visitChildren()-Methode, welche die entsprechenden Visit-Methoden der untergeordneten Knoten aufruft.

Um diese rekursive Vorgehensweise verständlicher zu machen, folgt ein Beispiel.

```
/**@brief
 * verarbeitet Additionen
 */
public String visitAddition(AdditionContext ctx) {
    return visitChildren(ctx) + "\n"
        + "iadd\n";
}
```

Die Tatsache, dass es eine Methode `visitAddition` mit diesem Eingabeparameter und diesem Rückgabewert gibt, geht auf ANTLR zurück. Der Funktionsrumpf wurde jedoch von uns verfasst.

Zunächst werden die Kind-Knoten der Addition, sprich die Operanden, besucht. Wenn die Operanden aus weiteren mathematischen Operationen bestehen, werden zunächst diese aufgerufen, um die beliebige Länge von Ausdrücken zu ermöglichen. Handelt es sich bei einem Operanden um eine Zahl, wird die Methode `visitNumber()` aufgerufen.

```
/**@brief
 * verarbeitet ganze Zahlen
 */
public String visitNumber(NumberContext ctx) {
    return "ldc " + ctx.getChild(0) + "\n";
}
```

Der Befehl *ldc* steht für *load constant* und ist die Jasmin-Instruktion, um einen Wert auf den Stack zu legen. `ctx.getChild(0)` sorgt dafür, dass der Wert aus dem entsprechenden Knoten aus dem Baum entnommen wird. Der Befehl *iadd* in der `visitAddition()`-Methode steht für *integer addition* und ist die Jasmin-Anweisung, zwei ganzzahlige Werte vom Stack zu nehmen, diese zu addieren und schließlich das Ergebnis wieder auf den Stack zu legen.

Der Compiler-interne Ablauf für die Übersetzung des Ausdruck

2 + 4

wäre also wie folgt.

Nach der lexikalischen und syntaktischen Analyse liegt folgender Syntaxbaum vor:

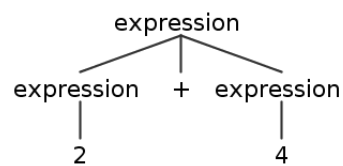


Abbildung 3: Syntaxbaum für beispielhaften arithmetischen Ausdruck

Zunächst wird die visitAddition-Methode aufgerufen, da der Operator der expression das "+" Zeichen ist. Die visitChildren-Methode gibt nun die Zeichenkette

ldc 2
ldc 4

zurück, da die untergeordneten Expressions nur Zahlen enthalten. Dem fügt die visitAddition-Methode wiederum den Befehl

iadd

hinzu, um die Addition durchzuführen.

In der JavaVM werden diese Instruktionen wie folgt verarbeitet: Im Initialzustand ist der Stack leer.

Stack der Java Virtual Machine



Abbildung 4: Zustand des Kellerspeicher der JavaVM

Der Befehl *ldc 2* legt den Wert *2* auf den Stack.

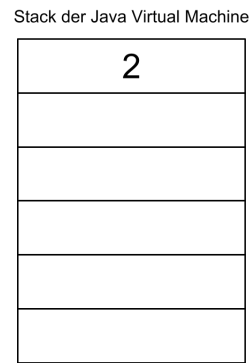


Abbildung 5: Zustand des Kellerspeicher der JavaVM

Der Befehl *ldc 4* legt den Wert *4* auf den Stack.

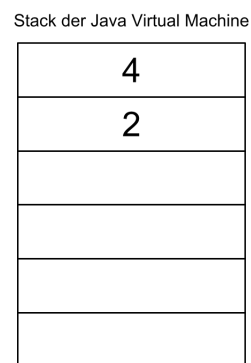


Abbildung 6: Zustand des Kellerspeicher der JavaVM

Der Befehl *iadd* nimmt zwei Werte vom Stack und legt das Ergebnis der Addition dieser Werte wieder auf den Stack

Stack der Java Virtual Machine



Abbildung 7: Zustand des Kellerspeicher der JavaVM

5.2 Ausweiten der simplen Grammatik

Nachdem die Sprache grundlegende Funktionalität erhielt, versuchten wir die weiteren Bestandteile, die durch die Anforderungen festgelegt wurden, zu implementieren.

5.2.1 Ergebnisausgabe

Ein Programm ohne Ausgabe ergibt keinen Sinn, da nie das berechnete Ergebnis betrachtet werden kann. Eine Ausgabefunktion zu implementieren stellte sich nicht als all zu großes Problem heraus, da Jasmin die Fähigkeit besitzt Objekte und Methoden in der Java Library zu nutzen. Wir implementierten dies Ausgabefunktion wie folgt:

```
/**@brief
 * vearbeitet den Aufruf der println-Funktion
 */
public String visitPrintln(PrintlnContext ctx) {
    return    ";calling println\n" +
              //legt ein System.out Objekt auf den Stack
              "getstatic java/lang/System/out Ljava/io/PrintStream;\n" +
              //Argument der print-Funktion
              visit(ctx.argument) +
              //ruft die Methode println des System.out-Objekts auf
              "invokevirtual java/io/PrintStream/println(I)V\n\n";
}
```

Es wird zuerst ein System.out-Objekt auf den Stack gelegt, danach der auszugebende Wert. Schließlich folgt der Aufruf der Methode println(), die das Argument sowie das System.out-Objekt vom Stack wieder herunternimmt. Der Aufruf in unserer Programmiersprache lautet:

```
println( argument );
```

5.3 Implementierung von Variablen

Bestandteile - Deklaration(int [name]): prüfen, ob Variable mit diesem Namen bereits existiert falls ja, Exception (variable already defined) falls nicht, neuer Eintrag mit Namen und Index in Variablentabelle - Wertzuweisung ([name] = [expression]) auflösen des Namen nach Index in Tabelle falls nicht gefunden, Exception (undefined variable) falls gefunden, generieren der Instruktionen visit(ctx.expr), legt zuzuweisenden Wert auf Stack istore [index]", speichert in JavaVM-Variablentabelle - Aufruf ([name] als Teil einer Expression) auflösen des Namen nach Index in Tabelle falls nicht gefunden, Exception (undefined variable) falls gefunden, generieren der Instruktionen iload [index]", kopiert Wert aus der JavaVM-Variablentabelle und legt diesen auf Stack

5.4 Implementierung von Funktionen

Bevor der eigentliche Visitor arbeitet, wird der gesamte Baum auf Funktionsdeklarationen untersucht. Dies sorgt dafür, dass Funktionen an beliebigen Stellen im Quellcode definiert werden können und keine Vorwärtsdeklarationen notwendig sind. Die gefundenen Funktionen werden in einer Liste gespeichert.

Bestandteile - Deklaration prüfen, ob Funktion mit diesem Namen und dieser Signatur bereits existiert falls ja, Exception (function already defined) falls nicht, neuer Eintrag mit Namen und Signatur in Funktionsliste

- Aufruf prüfen, ob Funktion in Funktionsliste enthalten ist falls nein, Exception (undefined function) falls ja, Kopieren von alter Variablentabelle und Konstruieren von neuer (ermöglicht scopes) Zusammensetzen von Argumenten und instructions eigentlicher Aufruf invokestatic [functionname] ([ggf. parameter])"Wiederherstellen von alter Variablentabelle

5.5 Implementierung von Konstanten

Intern Erweiterung von Variablen. Zusätzliche Schritte: - Deklaration - hasBeenAssigned-Flag auf false setzen

- Wertzuweisung - Unterscheidung, ob Konstante oder Variable falls Konstante: Prüfen, ob Konstante bereits Wert zugewiesen wurde falls ja, Exception (constant reassign) falls nicht, Zuweisung durchführen und entsprechendes Flag true setzen

5.6 Implementierung von Bedingten Verzweigungen

Aufteilen in conditionInstructions onTrueInstructions onFalseInstructions
Zusammensetzen - conditionInstructions, Bedingung wird überprüft falls Bedingung zu true evaluiert wird Sprung zu onTrueInstructions, onFalseInstructions überspringen falls Bedingung zu false evaluiert wird Ausführen von onFalseInstructions, onTrueInstructions überspringen

Labels für Sprungbefehle werden Labels benötigt. Um diese eindeutig benennen zu können wird gezählt, wie oft if-else vorkommt

Evaluiere die Bedingung: wie in C als int (in unserem Fall also expression).
0 == false, jeder andere Wert == true -> expression regel entsprechend angepasst, möglich sind - Vergleiche (<, <=, >, >=, ==) - Logik-Gatter (& &, ||, !)

5.7 Implementierung einer Schleife

Aufteilen in conditionInstructions onTrueInstructions

Zusammensetzen [conditionLabel] conditionInstructions Wenn conditionInstructions zu true evaluiert werden Sprung zu onTrueLabel Wenn conditionInstruction zu false evaluiert werden Sprung zu endOfWhileLabel

[onTrueLabel] onTrueInstructions Sprung zu conditionLabel

[endOfWhileLabel]

Labels für Sprungbefehle werden Labels benötigt. Um diese eindeutig benennen zu können wird gezählt, wie oft while vorkommt

6 Aufgetrene Probleme und deren Lösung

6.1 Operatorenpriorität bei arithmetischen Ausdrücken

vgl. Kapitel oben

6.2 Variablen

Ein wichtiger Bestandteil jeder Programmiersprache ist die Möglichkeit Variablen zu nutzen. Damit eine Variable im späteren Programmverlauf genutzt werden kann, muss sie zunächst deklariert werden. In unserer Sprache erfolgt dies durch die Angabe [Datentyp] [Name];. Beispielsweise legt der Aufruf

```
int x;
```

eine Integer-Variable mit dem Namen x an.

Als nächstes sollte eine Wertzuweisung erfolgen:

```
x = 42;
```

Nachdem diese obligatorischen Schritte vollzogen wurde, kann der Wert der Variable beliebig oft ausgelesen oder verändert werden.

Damit diese Features in die Zielsprache zu übertragen werden, nutzt unser Compiler eine Variablentabelle. Jasmin besitzt die Möglichkeit den Wert der oben auf dem Stack liegt zwischenzuspeichern. Der Befehl

```
astore <var-num>
```

nimmt den Wert vom Stack und speichert in am Index <var-num> in der Variablentabelle.

Mit dem Befehl

```
aload <var-num>
```

wird die Variable an der Position <var-num> wieder auf den Stack gelegt. Jasmin akzeptiert nur ganze Zahlen > 0 , jedoch keine Zeichenketten, für <var-num> Deshalb muss unser Compiler den Variablenname, den der Nutzer der Sprache wählt, zu einem Index auflösen.

6.2.1 Redefinition von bereits definierte Variable

Wenn im Syntaxbaum eine Variablendeklaration gefunden wird, wird zunächst überprüft, ob eine Variable mit diesem Namen bereits vorhanden ist. Für diesen Abgleich wird intern eine HashMap verwendet. Ist der Name noch

nicht vorhanden, wird er der HashMap hinzugefügt, wobei der Schlüssel der Name der Variable und der Wert die aktuelle Größe der Tabelle ist. Dies ist notwendig, um den Jasmin-Befehlen `astore` und `aload` einen Index zu übergeben.

6.2.2 Zugriff auf undefinierte Variable

Ein möglicher Fehler seitens des Benutzer unserer Sprache ist, dass eine Variable aufgerufen wird, obwohl diese zuvor nicht definiert wurde. Beim Aufruf einer Variablen wird deshalb zunächst überprüft, ob diese in der Variablen-Map vorhanden ist. Falls nicht wird eine entsprechende Exception ausgelöst.

6.3 Funktionen

6.3.1 Zugriff auf undefinierte Funktion

6.4 Redefinition von bereits definierter Funktion

6.4.1 Gültigkeitsbereiche

6.4.2 Funktionen mit gleichem Namen und unterschiedlichen Signaturen

6.5 Vorwärtsdeklarationen

6.6 Bedingte Verzweigungen

6.6.1 Umsetzung in Jasmin mit Hilfe von Labels und Sprungbefehlen

7 Verwendung des Compilers

7.1 Verwendung der UnsereCompilerbauSprache

- println-Funktion - Variablen deklarieren, definieren, aufrufen - Konstanten Funktionen deklarieren, aufrufen - Arithmetik - Aussgaenlogik - if-else-statements - while-Schleife

7.2 Aufrufen des UnserCompilerbauCompilers

- java -jar UCC.jar sourceCode.txt (- java -jar jasmin.jar output.j) - java output

8 Anhang

8.1 Quellenangaben

8.2 Quellcode des Compiler(?)