

Sven Eric Panitz

Lehrbriefe
Programmierung in Java



Ein kleiner Compiler

Ein kleiner Compiler

Sven Eric Panitz

26. Juni 2023

Inhaltsverzeichnis

1	Compilerbau	2
1.1	Phasen eines Compilers	2
1.2	Importierte Klassen	3
2	Abstrakter Syntaxbaum	4
2.1	Operatoren der Sprache	5
2.2	Baumknoten	5
2.2.1	Zahlenliterale	5
2.2.2	Variablen	6
2.2.3	Zuweisung	6
2.2.4	Operator-Ausdrücke	6
2.2.5	Bedingungsausdrücke	7
2.2.6	Schleifen	7
2.2.7	Code-Blöcke	7
2.2.8	Funktionsaufrufe	8
2.2.9	Funktionsdefinitionen	8
2.3	Beispielprogramme	8
2.3.1	Iterative Fakultät	9
2.3.2	Rekursive Fakultät	9
2.3.3	Zwei Parameter	10
2.3.4	Fibonaccizahlen	10
3	Pattern Matching für Baumknoten	11
4	Pretty Printing	12
4.1	Exception freier Rahmen für Writer	12
4.1.1	Stringdarstellung	12
4.1.2	Schreiben des Programms mit Einrückungen	14
5	Interpreter	15

6	X86 Code Generierung	16
6.1	Register	16
6.1.1	Register für Stackzeiger	16
7	Eine main-Methode für Compiler oder Interpreter	27
7.1	Hilfestellungen	27
7.2	Anwender-Frontend des Systems	27

1 Compilerbau

In diesem Projektblatt des Moduls »*Programmierungsmethoden und Techniken*« wird ein Compiler für eine kleine Programmiersprache entwickelt.

Ein Compiler ist ein Computerprogramm, das andere Programme als Eingabe erhält und wiederum Computerprogramme als Ausgabe erzeugt. Wir sprechen also bei einem Compiler über drei beteiligte Programmiersprachen:

- der Quellsprache des Programms, das der Compiler als Eingabe erhält.
- die Zielsprache, in die das Eingabeprogramm übersetzt wird. Dieses ist oft eine Maschinsprache in Form von Assembler.
- die Implementierungssprache, in der der Compiler entwickelt wurde.

Für dieses Projektblatt nutzen wir als Implementierungssprache Java und als Quellsprache eine einfache Sprache mit Bedingungen, Schleifen und Funktionsaufrufen.¹ Zielsprache ist Gnu Assembler (GAS) für x86-64-Architekturen.

1.1 Phasen eines Compilers

Ein Compiler besteht aus mehreren Phasen:

- Die lexikographische Analyse: Hier werden aus einem String die einzelnen Wörter der Sprache erzeugt. Diese werden auch als Token bezeichnet. Wörter sind z.B. Schlüsselwörter, Literale und Symbole der Sprache. Dieser Teil des Compilers wird als Lexer oder auch als Tokenizer bezeichnet. Das Ergebnis der lexikographischen Analyse ist eine Folge von Token.
- Die syntaktische Analyse: Dieser Teil des Compilers wird als Parser bezeichnet. Das Ergebnis ist ein abstrakter Syntaxbaum (kurz AST), der die Struktur des gelesenen Programms darstellt.

¹Wir haben keinen Namen für die Quellsprache. Das Javapaket, in der sie definiert ist heißt `longStack`, da nur mit dem Datentyp `long int` auf dem Stack gearbeitet wird. Das Anwenderprogramm trägt in nostalgischer Anwendung an das Programm *hugs* den Namen TUGS.

- Statische Analysen: In dieser Phase prüft der Compiler den Quelltext auf Konsistenz. Das ist bei statisch getypten Sprache zunächst der Typcheck. Bei einer Sprache wie Java fallen hier eine Vielzahl von Prüfungen an, wie:
 - sind alle Exceptions gefangen oder im throws deklariert?
 - werden Sichtbarkeiten beim Verwenden von Eigenschaften beachtet?
 - gibt es garantiert eine return-Anweisung bei Methoden mit Rückgabewert?
 - werden alle abstrakten Methoden implementiert?
 - gibt es keine zyklische Vererbung?
 - etc etc.
- Code-Generierung: In dieser Phase wird nun der eigentlich Code der Zielsprache generiert.

Viel Compiler haben noch weitere Phasen, wie die Übersetzung des Quelltextes in eine kleinere Kernsprache, den sogenannten Entzuckern von Konstrukten, die den Quelltext bequemer machen, einen abstrakten Zwischencode und unterschiedlichen Optimierungsphasen.

In diesem Projektblatt ist der abstrakte Syntaxbaum vorgegeben. Ein Lexer und ein Parser liegen als Klassen bereits vor.

Auf den Syntaxbaum werden drei Algorithmen realisiert:

- Ein Pretty-Printer, der das Programm, das durch den AST dargestellt wird, wieder als Text formatiert darstellt.
- Ein Interpreter, der Ausdrücke direkt auswertet.
- Eine Code-Generierung, die Assembler generiert, der mit dem gcc zusammen mit anderen C-Programmen zu einem ausführbaren Programm gelinkt werden kann.

Als Hilfsalgorithmus werden in einer Funktion alle lokalen Variablen einer Funktion ermittelt.

1.2 Importierte Klassen

Wir benötigen ein rudimentäres IO, eine Reihe von Container-Klassen und zwei funktionale Schnittstellen.

Wir werden folgende Standardklassen verwenden.

Zunächst einige Klassen des Pakets `java.io`:

AST.java

```
package name.panitz.longStack;

import java.io.IOException;
import java.io.StringWriter;
import java.io.Writer;
import java.io.FileWriter;
import java.io.FileReader;
```

```
import java.io.InputStreamReader;
import java.io.BufferedReader;
```

Dann ein paar Sammlungsklassen des Pakets `java.util`:

AST.java

```
import java.util.HashMap;
import java.util.Iterator;
import java.util.List;
import java.util.Map;
import java.util.Set;
import java.util.TreeSet;
```

Und ein paar funktionale Schnittstellen des Pakets `java.util.function`:

AST.java

```
import java.util.function.BinaryOperator;
import java.util.function.Function;
import java.util.function.Consumer;
```

Um alle in der Aufgabe implementierten Methoden manuell auszuprobieren, ist eine Methode `main` am Ende dieses Papiers definiert.

2 Abstrakter Syntaxbaum

Die Umsetzung des AST und der drei Algorithmen des AST bündeln wir in einer einzigen Java Schnittstelle. Alle Klassen, die wir benötigen, sind innere statische Klassen dieser Schnittstelle. Funktionalitäten sind als `default`-Methoden der Schnittstelle realisiert.

AST.java

```
@SuppressWarnings("unchecked")
public sealed interface AST
    permits AST.IntLiteral, AST.Var, AST.Assign, AST.OpExpr, AST.IfExpr, AST.WhileExpr
        , AST.Sequence, AST.FunCall {
```

Wir haben die Schnittstelle verriegelt, sodass er nur genau die acht angegebenen Klassen geben darf, die die Schnittstelle implementieren. Dieses hilft dabei, dass der javacompiler in switch-Ausdrücken checken kann, das auch alle acht Fälle abgedeckt wurden, ohne dass wir einen default Fall vorsehen müssen.

2.1 Operatoren der Sprache

Bevor wir die eigentlichen Baumknoten des AST definieren, seien die Operatoren der Quellsprache in Form einer Aufzählungsklasse definiert.

AST.java

```
public static enum BinOP{
    add ((x,y)->x+y,"+"),
    sub ((x,y)->x-y,"-"),
    mult((x,y)->x*y,"*"),
    eq  ((x,y)->x==y?"1L:0L","=");

    BinOP(BinaryOperator<Long> op,String name){
        this.op = op;
        this.name = name;
    }
    BinaryOperator<Long> op;
    String name;
}
```

Die vier Werte dieser Aufzählungsklasse beinhalten den textuellen Namen des Operators und die Funktion, die der Operator ausdrückt.²

2.2 Baumknoten

Nun definieren wir Klassen, die die Schnittstelle AST implementieren. Diese Klassen beschreiben als Baumknoten jeweils bestimmte Konstrukte unserer Quellsprache. Anweisungen, die aus mehreren anderen Ausdrücken oder Anweisungen zusammengesetzt sind, haben diese als Kindknoten.

Zusammen mit den Baumknoten geben wir jeweils ein Beispiel für das Sprachkonstrukt im Quelltext.

2.2.1 Zahlenliterale

Der einfachste Ausdruck unserer Quellsprache sind Zahlenliterale. Wir verarbeiten als einzigen Datentyp ganze Zahlen des Typs long, also 64 Bit lange Zahlen.

Der Baumknoten für Zahlenliterale hat keine Kinder und enthält den Wert des Zahlenliterals.

AST.java

```
public static record IntLiteral(long n) implements AST{}
```

²Die Division und Modulo-Operation haben wir bewusst ausgelassen, weil für diese das Erzeugen von Assembler-Code etwas schwieriger ist.

2.2.2 Variablen

Unsere Quellsprache kennt Variablen. Solche Variablen können lokale Variablen einer Funktion sein, oder Parameter der Funktion. Der Baumknoten für Variablen hat keine Kinder und enthält den Variablennamen.

AST.java

```
public static record Var(String name) implements AST{}
```

2.2.3 Zuweisung

In der Quellsprache können Variablen Werte zugewiesen werden. Als Zuweisungsoperator dient der Operator `:=`. Auf der rechten Seite kann ein beliebiger Ausdruck unserer Quellsprache stehen.

`x := 42`

Die Baumknoten der Zuweisung enthalten zwei Kinder. Einen Variablenknoten und einen beliebigen AST für die rechte Seite der Zuweisung.

AST.java

```
public static record Assign(Var v, AST right) implements AST{}
```

2.2.4 Operator-Ausdrücke

Die binären Operatoren haben wir bereits in einer Aufzählungsklasse definiert. Es gibt im Quelltext die Operatoren `+`, `-`, `*` und `=`. Auf weitere Operatoren haben wir der Einfachheit halber verzichtet. Insbesondere die Divisionsoperationen sind im Assembler ein klein wenig komplexer.

Im Quelltext gilt die übliche Operatorpräzedenz und es können Ausdrücke geklammert sein, so dass der folgende Ausdruck ein gültiger Ausdruck ist.

`(17+4)*2=26*2-10`

Ein Baumknoten für einen Operatorausdruck hat als Kinder den linken und den rechten Operanden und ein Objekt für den Operator.

AST.java

```
public static record OpExpr(AST left, BinOP op, AST right) implements AST{}
```


2.2.5 Bedingungsausdrücke

Die Quellsprache kennt if-Bedingungen. Diese haben immer auch den else-Fall zu spezifizieren. Syntaktisch werden die drei Schlüsselwörter `if`, `then` und `else` verwendet, um die Bedingung und die beiden Alternativen zu trennen.

```
if 17+4*2=x-1 then 1 else 42
```

Die Bedingung wird als wahr bezeichnet, wenn sie nicht den Wert 0 hat.

Der Baumknoten hat entsprechend drei Kinder: die Bedingung und die beiden Alternativen.

AST.java

```
public static record IfExpr(AST cond, AST alt1, AST alt2) implements AST{}
```

2.2.6 Schleifen

Die Quellsprache kennt eine typische `while`-Schleife aus einer Bedingung und einem Schleifenrumpf. Zur Trennung werden syntaktisch die Schlüsselwörter `while` und `do` verwendet.

```
while x+y=8 do x:=x-1
```

Entsprechend hat der Baumknoten für Schleifen zwei Kinder. In einem Kind steht der AST für die Bedingung im anderen Kind der AST für den Schleifenrumpf.

AST.java

```
public static record WhileExpr(AST cond, AST body) implements AST{}
```

2.2.7 Code-Blöcke

Ein Code-Block ist die Sequenz von mehreren Ausdrücken.

Syntaktisch sind Code-Blöcke in geschweiften Klammern eingeschlossen und die Ausdrücke des Code-Blocks durch Semikolons getrennt. Das ist anders als in Java. Die Semikolons stehen nur zwischen zwei Ausdrücken, nicht auch noch nach dem letzten Ausdruck. Sie markieren also nicht das Ende eines Ausdrucks.

Der Wert eines Code-Blocks sei der Wert des letzten Ausdrucks in der Folge von Ausdrücken.

```
{x:=5  
;r:=1  
;while (x)  
  {r := r*x  
  ;x := x -1  
  }
```

```
;r  
}
```

Der Baumknoten für Code-Blöcke kann eine beliebige Anzahl von Kindern haben. Daher enthält er eine Liste von Kindern.

AST.java

```
public static record Sequence(List<AST> sts) implements AST{}
```

2.2.8 Funktionsaufrufe

Unsere Quellsprache kennt auch Funktionen, die in herkömmlicher Weise aufgerufen werden können.

```
if x=0 then 1 else x*f(x-1)
```

Der Baumknoten enthält den Namen der aufgerufenen Funktion und eine Liste von Argumenten.

AST.java

```
public static record FunCall(String name, List<AST> args) implements AST{}
```

2.2.9 Funktionsdefinitionen

In der Quellsprache ist es auch möglich, Funktionen zu definieren. Hierzu dient das Schlüsselwort `fun`. Eine Funktion hat dann einen Namen und eine Parameterliste. Ein Gleichheitszeichen läutet den Funktionsrumpf ein.

```
fun fac(x) = if x=0 then 1 else x*f(x-1)
```

Der Name der Funktion und der Name der Argumente sind im Knoten als Strings gespeichert. Der Rumpf ist ein Kindknoten im Baum. Die Klasse, die Funktionsdefinitionen beschreibt, ist kein Teil des AST. Es gibt also vorerst auch keine lokalen Funktionsdefinitionen.

AST.java

```
static record Fun(String name, List<String> args, AST body){}
```

2.3 Beispielprogramme

Um ein besseres Gefühl für die Sprache und den AST zu bekommen, betrachten wir ein paar Beispielprogramme.

2.3.1 Iterative Fakultät

Zunächst eine Umsetzung der Fakultätsfunktion mit Hilfe einer Schleife. Es gibt eine Ergebnisvariable `r` und es wird auf dem Parameter `x` gerechnet.

```
fun fac(x) =  
  {r := 1  
   ;while x  
   do  
     {r := r*x  
      ;x := x-1  
     }  
   ;r  
  }
```

Folgendes ist die AST dieses Programms.

AST.java

```
static Fun factorial= new Fun("fac",List.of("x"),  
  new Sequence(List.of  
    (new Assign(new Var("r"), new IntLiteral(1))  
    ,new WhileExpr  
      (new Var("x")  
      ,new Sequence(List.of  
        (new Assign  
          (new Var("r")  
            ,new OpExpr(new Var("r"), BinOP.mult,new Var("x")))  
          ,new Assign  
            (new Var("x")  
              ,new OpExpr(new Var("x"), BinOP.sub, new IntLiteral(1)))  
        ))  
      )  
    ,new Var("r")  
    ))  
  );
```

2.3.2 Rekursive Fakultät

Die rekursive Version der Fakultätsfunktion lässt sich elegant in einer Zeile schreiben:

```
fun f(x) = if x = 0 then 1 else x*f(x-1)
```

Entsprechend ist auch der AST für dieses Programm etwas kürzer.

AST.java

```
static Fun factorialRek = new Fun("f",List.of("x"),
    new IfExpr(new OpExpr(new Var("x"),BinOP.eq,new IntLiteral(0))
        , new IntLiteral(1)
        , new OpExpr(new Var("x"),BinOP.mult
            ,new FunCall("f"
                , List.of(new OpExpr(new Var("x")
                    , BinOP.sub, new IntLiteral(1)))))));
```

2.3.3 Zwei Parameter

Die folgende Funktion soll dazu dienen, zu testen, ob die Parameter nicht vertauscht werden. Dieses passiert gerne einmal bei der Code-Generierung.

```
fun minus(x,y) = x-y
```

Hier noch der AST für diese Funktion.

AST.java

```
static Fun minus =
    new Fun("minus"
        ,List.of("x","y")
        ,new OpExpr(new Var("x"), BinOP.sub, new Var("y")));
```

2.3.4 Fibonaccizahlen

Als letztes darf die Funktion der Fibonaccizahlen nicht fehlen. Auch diese kommt ganz ohne einen längeren Code-Block aus.

```
fun fic(x) = if x=0 then 0 else if x = 1 then 1 else fib(x-2)+fib(x-1)
```

Und auch hierzu sei der AST in Java definiert.

AST.java

```
static Fun fib
    = new Fun("fib",List.of("x")
        ,new IfExpr(new OpExpr(new IntLiteral(0),BinOP.eq,new Var("x"))
            ,new IntLiteral(0)
            ,new IfExpr(new OpExpr(new IntLiteral(1),BinOP.eq,new Var("x"))
                ,new IntLiteral(1)
```

```

        ,new OpExpr
        (new FunCall("fib"
        ,List.of
        (new OpExpr(new Var("x"),BinOP.sub, new IntLiteral(2))))
        ,BinOP.add
        ,new FunCall("fib"
        ,List.of
        (new OpExpr(new Var("x"),BinOP.sub, new IntLiteral(1))))
        )));

```

3 Pattern Matching für Baumknoten

Mit AST haben wir jetzt eine Baumstruktur mit acht verschiedenen Arten von Baumknoten. Algorithmen auf den Datentyp AST werden in der Regel eine Fallunterscheidung über diese acht Baumknotenarten machen. Java bietet hierfür als *previre feature* ein sogenanntes Pattern Matching in der Switch-Anweisung an. Hier wird in den case-Ausdrücken angegeben, für welche Klasse der implementierenden Klasse der Fall zutrifft. Dabei wird eine neue Variable für das Objekt, auf das das *switch* angewendet wird, eingeführt.

Wir geben ein einfaches Beispiel für eine default Methode in der Schnittstelle AST, die mit einem switch-case und pattern Matching auf die Baumknotenklassen unterscheidet und einen String erzeugt:

AST.java

```

default String whatAreYou() {
    return switch(this){
        case Var v      -> "Variable mit Namen "+v.name;
        case IntLiteral il -> "Literal mit Wert "+il.n;
        case FunCall fc  -> "Aufruf der Funktion: "+fc.name;
        default          -> "Irgend ein anderer Baumknoten";
    };
}

```

Man kann sogar ein sogenanntes *deconstruction pattern* verwenden. dabei werden entsprechend des kanonischen Konstruktors der Recordklasse Variablen für die Argumente des Konstruktors angegeben. Mit diesen kann man dann auf die Attribute des Recors-Objekts zugreifen.

Mit einem *deconstruction pattern* kann die obige Funktion auch wie folge geschrieben werden:

AST.java

```

default String whatAreYou2() {
    return switch(this){
        case Var(var name) -> "Variable mit Namen "+name;

```

```

        case IntLiteral(var n) -> "Literal mit Wert "+n;
        case FunCall(var name, var args) -> "Aufruf der Funktion: "+name;
        default                    -> "Irgend ein anderer Baumknoten";
    };
}

```

4 Pretty Printing

In diesem Abschnitt wird in Form einer Aufgabe nun der erste Algorithmus für den AST geschrieben. Die erste Aufgabe wird sein, aus dem AST wieder den Quelltext unserer Quellsprache zu schreiben.

4.1 Exception freier Rahmen für Writer

Hierzu schreiben wir uns eine Hilfsklasse, die es vereinfachen soll, Text in einen Writer zu schreiben. Diese Klasse bietet einen kleinen Rahmen um einen Writer herum. Beim Schreiben werden eventuelle Exceptions gefangen und wieder als `RuntimeException` geworfen. Dieses ist notwendig, weil die Funktionen, die im Pattern-Matching verwendet werden, keine allgemeinen Ausnahmen werfen dürfen.

Zusätzlich bietet die Hilfsklasse die Möglichkeit, eine neue Zeile mit einer Einrückung zu schreiben. Die Einrückung kann erhöht und wieder verringert werden.

Als dritte Zusatzfunktion generiert die Klasse die Folge der natürlichen Zahlen. Diese werden wir bei der Code-Generierung verwenden, um dort unterschiedliche Labelnamen zu generieren. Die komplette Hilfsklasse befindet sich in Abbildung 1.

4.1.1 Stringdarstellung

Die Hilfsklasse `ExWriter` kommt bei unseren ersten Algorithmus für den AST in Einsatz. Für einfache Test wird ein `StringWriter` in die Klasse `ExWriter` gesteckt.

AST.java

```

default String show() {
    var r = new ExWriter(new StringWriter());
    show(r);
    return r.w.toString();
}

```

Die folgende Methode schreibt ein Funktionsdefinition:

AST.java

```

static class ExWriter{
    Writer w;
    String indent = "";
    void addIndent(){indent=indent+"  ";}
    void subIndent(){indent=indent.substring(2);}
    void nl(){write("\n"+indent);}

    int lbl=0;
    int next(){return lbl++;}

    public ExWriter(Writer w) {
        this.w = w;
    }
    void lnwrite(Object o) {
        nl();
        write(o);
    }
    void write(Object o) {
        try {
            w.write(o+"");
        } catch (IOException e) {
            throw new RuntimeException(e);
        }
    }
}

```

Abbildung 1: Hilfsklasse zum formatierten Schreiben mit abgefangenen Ausnahmen.

AST.java

```

static String show(Fun fd) {
    var r = new ExWriter(new StringWriter());
    r.write("fun ");
    r.write(fd.name);
    r.write("(");
    var first=true;
    for (var arg:fd.args){
        if (first) first=false;else r.write(", ");
        r.write(arg);
    }
    r.write(") = ");
    r.addIndent();
    r.nl();
    fd.body.show(r);
    return r.w.toString();
}

```

4.1.2 Schreiben des Programms mit Einrückungen

Jetzt können wir alle Knoten des AST mit Einrückungen schön schreiben.

Aufgabe 1 Vervollständigen Sie die Methode `show`, so dass das durch den AST dargestellte Programm vollständig als Quelltext geschrieben wird. Nutzen Sie dabei die Syntax, wie sie informell mit den Baumklassen vorgestellt wurde.

Berücksichtigen Sie, dass bei Operatorausdrücken wegen der Operatorpräzedenzen eventuell Klammern gesetzt werden müssen.

AST.java

```
default void show(ExWriter r) {
    switch (this){
        case Var(var name)      -> {r.write(name);}
        case IntLiteral(var n)  -> {r.write(n+"");}
        default -> {r.write("not yet implemented. Show: "+this);}
    };
}
```

Zum Ausprobieren der Lösung kann die `main`-Methode verwendet werden. Hier eine Beispielsitzung:

```
$ java -cp classes/ name.panitz.longStack.AST -i t1.ls
TUGS (Tiny Usable Great System) :? for help
> :s fac
fun fac(x) =
  if x=0
  then 1
  else x*fac(x-1)
> :s facIt
fun facIt(x) =
  {r := 1
  ;while not(x=0)
  do
    {r := r*x
    ;x := x-1
    }
  ;r
  }
> :q
```


5 Interpreter

Als nächstes soll ein Programm unserer Quellsprache direkt ausgewertet werden. Der einzige Datentyp, der verarbeitet wird, sind dabei ganze 64-Bit Zahlen. Wir überladen die default-Methode `ev` (für *evaluate*).

Die erste Version wertet einen Ausdruck aus, ohne dabei auf Funktionsdefinitionen zuzugreifen. Das Ergebnis ist dann die Zahl, die der Ausdruck berechnet.

AST.java

```
default long ev() {return ev(List.of());}
```

Diese ruft die überladene Version auf, die als zusätzlichen Parameter eine Abbildung von Strings auf Funktionsdefinitionen enthält.

AST.java

```
default long ev(List<Fun> funs) {  
    HashMap<String, Fun> fs = new HashMap<>();  
    funs.forEach(fun->fs.put(fun.name, fun));  
    return ev(fs, new HashMap<>());  
}
```

Und diese schließlich ruft die Version auf, die zusätzlich eine Umgebung (*environment*) enthält, in der für Variablennamen aktuell Werte gespeichert sind. Aus dieser Umgebung können die aktuellen Werte für Variablen während der Auswertung gelesen werden und bei einer Zuweisung können dort neue Werte abgespeichert werden.

Aufgabe 2 Ergänzen Sie die fehlenden Fälle, zur Auswertung eines Ausdrucks.

AST.java

```
default long ev(Map<String, Fun> fs, Map<String, Long> env) {  
    return switch(this){  
        case IntLiteral il -> il.n;  
        case Var v -> env.get(v.name);  
        case OpExpr ae -> 42L; /*TODO*/  
        case IfExpr ie -> 42L; /*TODO*/  
        case Assign as -> 42L; /*TODO*/  
        case WhileExpr we -> 42L; /*TODO*/  
        case FunCall fc -> {yield 42L; /*TODO*/ }  
        case Sequence sq -> {  
            var r = 42L;  
            /*TO DO*/;  
            yield r;  
        }  
    }  
}
```

```
};  
}
```

Mit Hilfe der main-Methode sollten Sie jetzt in der Lage sein, beliebige Ausdrücke auswerten zu lassen. Hier ein Beispielsitzung:

```
$ java -cp classes/ name/panitz/longStack/AST -i src/t1.ls  
TUGS (Tiny Usable Great System) :? for help  
> 17+4*2  
25  
> (17+4)*2  
42  
> fac(5)  
120  
> fib(10)  
55  
> :q  
$
```

6 X86 Code Generierung

Wir werden in diesem Kapitel nun Assembler Code generieren. Der generierte Assembler ist in der Syntax für den Gnu Assembler. Er ist für eine 64-Bit Maschine mit X86 Befehlssatz. Der generierte Assembler kann von C-Programmen aus aufgerufen werden. Er hält sich also insbesondere an die Aufrufkonvention (engl. calling convention) des gcc.

Der Assembler arbeitet mit Registern und mit einem Stack. Gerechnet wird in der Regel mit Werten, die in Registern gespeichert sind. Zwischenergebnisse und ebenso lokale Variablen und Funktionsargumente werden auf dem Stack gespeichert.

6.1 Register

Wir betrachten zunächst die relevanten Register des Assemblers.

6.1.1 Register für Stackzeiger

Die beiden Register `%rsp` (Stackpointer) und `%rbp` (Basepointer) werden ausschließlich dafür verwendet, um auf Adressen des Stacks zu verweisen.

Im `%rsp` steht die aktuelle Adresse des ersten freien Platzes auf dem Stack notiert. Dabei ist übrigens zu beachten, dass der Stack nach unten wächst.

Im `%rbp` steht eine Adresse im Stack zu der relativ die Argumente und lokalen Variablen der Funktion zu finden sind, die aktuell ausgeführt wird. Es ist quasi der Zeiger auf die Basis des gerade ausgeführten Funktionsaufrufs auf dem Stack.

Zum Lösen der Aufgabe werden sie diese beiden Register nicht manuell verändern. Sie werden aber über den Basepointer auf die Argumente und lokalen Variablen einer Funktion auf dem Stack zugreifen.

Rechenregister Wir werden ausschließlich Rechenoperationen auf den beiden Registern `%rax` und `%rbx` ausführen. Dem Register `%rax` kommt dabei die besondere Rolle zu, dass es immer das zuletzt berechnete Ergebnis enthalten wird. Insbesondere nach Ende der Ausführung einer Funktion steht das Ergebnis der Funktion in diesem Register.

Argumentregister Sechs weitere Register werden verwendet, um Argumente an eine Funktion zu übergeben. Hat eine Funktion mehr als 6 Argumente, so werden weitere Argumente auf den Stack gespeichert. Funktionen, die ihrerseits im Rumpf andere Funktionen aufrufen, benötigen die sechs Register, um deren Argumente dort zu übergeben. Daher werden wir zu Beginn einer Funktion alle die Argumente, die in einem Register übergeben wurden, auch direkt auf den Stack ablegen.

Die sechs Register für die Argumentübergabe sind namentlich in folgendem Array genannt.

AST.java

```
String[] rs = {"%rdi", "%rsi", "%rdx", "%rcx", "%r8", "%r9"};
```

Instruktionen Nun betrachten wir die relevanten Assembler-Anweisungen für diese Aufgabe und deren Argumente.

Argumente Die Anweisungen des Assemblers, die wir verwenden, haben 0, 1, oder 2 Argumente. Diese werden in vier verschiedene Formen auftreten.

Zahlenkonstanten Zahlenkonstanten können direkt als Literale geschrieben werden. Ihnen wird das Dollarzeichen vorangestellt. So bezeichnet also `$42` die Zahl 42.

Registerwerte Register sind Speicherzellen, in denen Zahlen stehen. Soll in einer Anweisung direkt mit einer solchen Speicherstelle gearbeitet werden, um dort eine Zahl hineinzuschreiben, oder herauszulesen, so wird direkt der Name des Registers notiert, zum Beispiel `%rax`.

Man nennt dieses auch direkte Adressierung.

Registeradressen Wenn die Zahl, die in einem Register gespeichert ist, erst die Adresse einer weiteren Speicherzelle ist, mit der gearbeitet werden soll, so schreibt man das Register in Klammern also z.B. (`%rax`). Dieses entspricht der Pointer-Dereferenzierung in C. Man nennt dieses die indirekte Adressierung.

Man kann sogar noch einen Schritt weiter gehen. Dem geklammerten Register kann noch eine Zahl vorangestellt werden, also zum Beispiel `8(%rax)`. Dieses bedeutet: nimm die Zahl aus Register `%rax` als Adresse im Hauptspeicher. Addiere auf diese Adresse noch die acht drauf, gehe also 8 Byte weiter im Hauptspeicher und arbeite mit dieser Speicherzelle.

Label Es gibt Sprungbefehle, die zu einem im Assemblercode gesetzten Label springen. Dann ist das Argument der Anweisung ein Label. Label beginnen bei uns mit einem Punkt. Sie können einen beliebigen Bezeichner haben. Wir werden für Label den Buchstaben L benutzen und dann durchnummerieren, also `.L1`, `.L2`, `.L3`...

Stack Instruktionen

pushq Mit der Anweisung `pushq` können neue Daten auf den Stack gelegt werden. Es bewirkt, dass dabei auch das Register `%rsp` geändert wird. Die Anweisung hat ein Argument, mit dem spezifiziert wird, welche Daten auf den Stack gelegt werden sollen.

Wir geben ein paar Beispiele:

- Um zum Beispiel die Konstante 42 auf den Stack zu legen, verwendet man:
`pushq $42`
- Um den Wert, der im Register `%rax` steht, auf den Stack zu legen, verwendet man: `pushq %rax`
- Um den Wert, der im Stack zwei Plätze oberhalb der Adresse in `%rbp` gespeichert ist, auf den Stack zu legen, verwendet man:
`pusq 16(%rbp)`

Die 16, weil wir ausschließlich mit 8-Byte Daten operieren.

popq Um Daten, die direkt als oberstes auf dem Stack liegen, wieder von Stack zu nehmen und woanders zu speichern, dient die Anweisung `popq`. Das Argument ist dann das Ziel, in dem die Daten gespeichert werden. So wird mit `popq %rax` das oberste Stackelement in das Register `%rax` gespeichert.

Daten verschieben Die beiden Stackbefehle verschieben schon Daten, allerdings immer vom Stack oder auf den Stack. Um zwischen beliebigen Speicherstellen Daten zu verschieben, gibt es die Anweisung `movq`. Sie hat zwei Argumente: die Quelle und das Ziel der Daten. Die beiden Argumente werden durch ein Komma getrennt.

Wir geben ein paar Beispiele:

- `movq %rax, %rbx` schreibt die in Register `%rax` gespeicherte Zahl in das Register `%rbx`.
- `movq (%rax), %rbx` schreibt die Zahl, die in der Speicherzelle gespeichert ist, deren Adresse in Register `%rax` steht, in das Register `%rbx`.

Rechnen

Arithmetische Operationen Für die drei arithmetischen Operatoren unserer Sprache gibt es die direkten drei Anweisungen, die jeweils zwei Parameter haben: `imulq`, `subq` und `addq`. Sie verrechnen die Daten der beiden Argumente und speichern das Ergebnis in das Register des zweiten Arguments. So addiert z.B. `addq %rbx, %rax` die Zahlen in den Registern `%rbx` und `%rax` und speichert das Ergebnis in Register `%rax`.

Vergleichsoperatoren Zwei Werte können mit der Anweisung `cmpq` verglichen werden. Das Ergebnis ist in Form eine Subtraktion. Für zum Beispiel `cmpq %rax, %rbx` gilt: Wenn die Zahl in `%rax` größer ist, ist das Ergebnis positiv, wenn sie kleiner ist negativ und bei zwei gleichen Zahlen ist es 0.

Sprünge

bedingte Sprünge Für bedingte Sprünge verwenden wir ausschließlich die Anweisung `jne`. Sie hat einen Parameter. Dieser ist ein Label. Wenn die vorhergehende `cmpq`-Anweisung nicht die Gleichheit ergeben hat, dann wird als nächste Anweisung zum angegebenen Label gesprungen.

unbedingte Sprünge Wenn man unter jeder Bedingung zu einem bestimmten Label springen will, kann man die Anweisung `jmp` verwenden, die das entsprechende Label als Argument hat.

Berechnung der lokalen Variablen Bevor wir beginnen können, Assembler zu generieren, benötigen wir die Menge aller lokalen Variablen innerhalb einer Funktion. Hierfür ist eine eigene Methode zu schreiben, die diese aus einem AST extrahiert.

AST.java

```
default Set<String> getVars(){
    var r = new TreeSet<String>();
    getVars(r);
    return r;
}
```

Lokale Variablen brauchen mindestens eine Zuweisung innerhalb der Funktionsrumpfs. Somit wird bei jeder Zuweisung eine potentielle weitere lokale Variable auf der linken Seite gefunden. Ansonsten ist der gante AST zu traversieren, um nach weiteren Zuweisungen zu suchen.

Aufgabe 3 Ergänzen Sie die fehlenden Fälle, um die Menge der Variablen eines Programms zu berechnen.

AST.java

```
default void getVars(TreeSet<String> r) {
    switch (this){
        case Assign    as -> {r.add(as.v.name);as.right.getVars(r);}

        /* TODO add missing cases */

        default        -> {}
    };
}
```

X86 Code für Funktionen Wir können nun beginnen, Assembler zu schreiben. Da Funktionen mit der Parameterübergabe für den Anfänger etwas zu komplex sind, ist die Code-Generierung für Funktionsdefinitionen und Funktionsaufrufe bereits umgesetzt. Alle anderen Sprachkonstrukte unserer Quellsprache sind als Aufgabe umzusetzen.

In diesen Abschnitt folgt die Beschreibung der Code-Generierung für Funktionsdefinitionen. Für eine Funktionsdefinition ist Assembler-Code zu generieren und in ein Writer-Objekt zu schreiben.

AST.java

```
static void asm(Fun f,ExWriter r) {
```

Funktionskopf Den Beginn einer Funktion beschreiben im Assembler ein paar globale Labels, in denen der Funktionsname gesetzt wird.

AST.java

```
r.nl();
r.lnwrite(".globl  "+f.name);
r.lnwrite(".type   "+f.name+" , @function");
r.lnwrite(f.name+":");
r.addIndent();
```

Stackpointer verwalten Die erste Aufgabe ist es, die Stackpointer zu verwalten. Hierzu wird der Wert des aktuellen Basepointers auf den Stack geschrieben. Der aktuelle Stackpointer wird dann während der Ausführung der Funktion zum Basepointer.

AST.java

```
r.lnwrite("pushq %rbp");  
r.lnwrite("movq %rsp, %rbp");
```

Am Ende, beim Verlassen der Funktion wird wieder der alte Zustand der beiden Register hergestellt werden.

Parameter von Registern auf den Stack Die Aufrufkonvention lädt die ersten 6 Argumente in die Register, die in der Reihung rs benannt sind. Die weiteren Argumente befinden sich auf dem Stack. Dieses ist eine effizientere Aufrufvariante, als wenn alle Argumente auf den Stack gelegt werden. Allerdings funktioniert diese Variante nur, wenn unsere Funktion selbst keine andere Funktion aufruft. Um hier einheitlich arbeiten zu können legen wir somit als erstes auch die Argumente, die in Registern übergeben werden, auf den Stack.

Dabei bauen wir eine Abbildung von Variablen und Parameternamen zu Zahlen auf. Diese Abbildung sagt für eine Variable, wie viel Schritte vom Basepointer entfernt, sich die Variable auf dem Stack gespeichert befindet.

Als erstes schauen wir, wie viele Argumente über Register übergeben wurden und legen die Abbildung an.

AST.java

```
var registerArgs = Math.min(rs.length, f.args.size());  
var env = new HashMap<String,Integer>();
```

Nun werden nacheinander die Argumente aus den Registern auf den Stack gelegt und in der Abbildung vermerkt, wie viele Speicherplätze vom Basepointer sie entfernt liegen. Das erste Argument liegt direkt über den Basepointer. Deshalb beginnen wir mit -8.

AST.java

```
var sp=-8;  
for (var i=0;i<registerArgs;i++) {  
    r.lnwrite("movq "+rs[i]+"", "+sp+"(%rbp)");
```

```

    env.put(f.args.get(i), sp);
    sp = sp-8;
}

```

Lokale Variablen in Umgebung vermerken Nun sind die per Register übergebenen Argumente auf dem Stack und in der Umgebung vermerkt. Als nächstes sind Stackpositionen der lokalen Variablen zu vermerken. Dieses sind die Variablen, denen etwas im Funktionsrumpf zugewiesen wird, ohne die Argumente.

Wir haben zwar noch keine initialen Werte, für diese Variablen, aber müssen der Stackpointer entsprechend weiter setzen, damit dieser Platz auf dem Stack nicht durch weitere push-Befehle überschrieben wird.

AST.java

```

var vs = f.body.getVars();
vs.removeAll(f.args);
for (var v:vs) {
    env.put(v, sp);
    sp = sp-8;
}
r.lnwrite("subq $" + (-sp) + ", %rsp");

```

Stackparameter in Umgebung vermerken Gibt es mehr Argumente, als in den Registern übergeben werden kann, liegen die weiteren Argumenten bereits auf dem Stack. Auch dieses ist noch in der Umgebung zu vermerken.

AST.java

```

sp = 16;
for (var i=registerArgs; i<f.args.size(); i++) {
    env.put(f.args.get(i), sp);
    sp+=8;
}

```

Code für Funktionsrumpf Jetzt kann der eigentliche Assembler-Code für den Funktionsrumpf generiert werden.

AST.java

```
f.body.asm(Map.of(),env,r);
```

Stackpointer Zurücksetzen beim Verlassen Beim Verlassen der Funktion, wird der alte Zustand des Stackpointers und des Basepointers wieder hergestellt. Die Funktion schließt mit der Anweisung `ret`, die dafür sorgt, dass zurück zur Aufrufstelle des Funktionsaufrufs gesprungen wird.

AST.java

```
r.lnwrite("movq %rbp, %rsp");  
r.lnwrite("popq %rbp");  
r.lnwrite("ret");  
r.subIndent();  
}
```

Code für Ausdrücke Jetzt kommt der spannende teil. Für alle Ausdrücke der Quellsprache ist Assembler zu generieren.

Für einfache Tests, ohne Funktionsaufrufe und Variablen ist die folgende Methode vorgesehen.

AST.java

```
default String asm() {  
    var r = new ExWriter(new StringWriter());  
    asm(new HashMap<>(),new HashMap<>(),r);  
    return r.w.toString();  
}
```

Eine Statische Methode sei für ein ganzes Programm, das aus einer Liste von Funktionsdefinitionen besteht, gegeben.

AST.java

```
static String asm(List<Fun> fs) {  
    var r = new ExWriter(new StringWriter());  
    fs.forEach(f->asm(f,r));  
    return r.w.toString();  
}
```

Die eigentliche Funktion zur Codegenerierung hat drei Argumente:

- Der Abbildung, die für Funktionsnamen die Funktionsdefinition bereit hält.³
- Die Umgebung: das ist die Abbildung von Variablennamen auf Positionen relativ zum Basepointer.
- Das Objekt, in das der Assembler-Code textuell geschrieben wird.

AST.java

```
default void asm(Map<String, Fun> fs, Map<String, Integer> e, ExWriter r){  
    switch (this){
```

Der Code, der für einen Ausdruck generiert wird, soll immer Code sein, der das Ergebnis für den Ausdruck berechnet, so dass dieses anschließend im Register %rax gespeichert ist.

Zahlenkonstanten Der Assembler-Code für Zahlenkonstanten ist dann denkbar einfach. Die entsprechende Konstante wird in das Register %rax gespeichert.

AST.java

```
case IntLiteral il -> {r.lnwrite("movq $" + il.n + ", %rax");}
```

Variablenzugriff Bei einem Variablenzugriff wird in der Umgebung e die relative Stack-Adresse vom Basepointer ausgehend nachgeschlagen. Aus dieser wird dann vom Stack der Wert der Variablen in das Register %rax geladen

AST.java

```
case Var v -> {r.lnwrite("movq " + e.get(v.name) + "(%rbp), %rax");}
```

Soweit die ersten beiden Fälle für Ausdrücke. Die weiteren werden jetzt als Aufgaben formuliert. Um ein wenig Hinweise für bestimmte Fälle zu bekommen, kann man eine C-Funktion schreiben, die ein äquivalentes Konstrukt enthält, und sich mit dem gcc durch die Option -S den Assembler generieren lassen. Dort lässt sich dann analog betrachten, was für ein Code zu generieren ist.

³Diese wird derzeit nicht verwendet.

Zuweisungen Bei einer Zuweisung ist der Code für die linke Seite der Zuweisung zu generieren. Dieser Code wird dafür sorgen, dass im Register `%rax` das Ergebnis der linken Seite steht. Dieses ist dann auf dem Stack an der Position zu speichern, die für die Variable auf der linken Seite der Zuweisung in der Umgebung vermerkt ist.

Aufgabe 4 Schreiben Sie in der Code-Generierung das Pattern-Case für die Zuweisung.

Operator-Ausdrücke Für einen Operatorausdruck müssen wir für die beiden Operanden Code erzeugen.

Es empfiehlt sich wie folgt vorzugehen:

- Erzeuge Code für den rechten Operanden.
- Push das Ergebnis des rechten Operanden auf den Stack.
- Erzeuge Code für den linken Operanden.
- Hole das Ergebnis des rechten Operanden vom Stack und lade es in Register `%rbx`.
- Generiere die Anweisung für den Operator. Für den Operator `=` ist dabei folgende Anweisungsfolge zu generieren:

```
cmpq %rax, %rbx
sete %al
movzbl %al, %eax
```

Aufgabe 5 Schreiben Sie in der Code-Generierung das Pattern-Case für Operatorausdrücke.

Code-Blöcke Bei einer Sequenz von Ausdrücken in Form eines Code-Blocks ist für diese nacheinander Code zu generieren. Damit steht abschließend das Ergebnis des letzten Ausdrucks der Sequenz in Register `%rax`.

Aufgabe 6 Schreiben Sie in der Code-Generierung das Pattern-Case für Code-Blöcke.

While-Schleifen Bei der Codegenerierung für Schleifen sind zwei Label zu generieren. Eines für den Beginn des Schleifenrumpfes und eines für den Beginn des Bedingungstests.

Es empfiehlt sich folgendes Vorgehen:

- Erzeuge zwei neue Label `l1` und `l2`.
- Generiere den Befehl der zu Label `l1` springt.
- Generiere im Assembler das Label `l2`.
- Generiere Code für den Schleifenrumpf.
- Generiere im Assembler das Label `l1`.

- Generiere Code für die Schleifenbedingung.
- Generiere den Code, der prüft, ob die Schleifenbedingung 0 als Ergebnis hatte: `cmpq $0, %rax`.
- Generiere den bedingten Sprung zum Label 12, wenn der Test zuvor nicht wahr war.

Aufgabe 7 Schreiben Sie in der Code-Generierung das Pattern-Case für Schleifen.

If-Bedingungen Auch für die Realisierung der Bedingungsausdrücke sind zwei Sprunglabel 11 und 12 zu erzeugen. Label 11 steht zu Beginn der then-Alternative, das Label 12 am Ende des Gesamtausdrucks.

Gehen Sie wie folgt vor:

- Generieren Sie Code für die Bedingung.
- Generieren Sie die Anweisung, die das Ergebnis der Bedingung mit 0 vergleicht.
- Generieren Sie einen bedingten Sprung zu Label 11, wenn das letzte Ergebnis nicht gleich ergab.
- Generieren Sie Code für die else-Alternative.
- Generieren Sie einen Sprung zu Label 12.
- Schreiben Sie das Label 12.
- Generieren Sie Code für die then-Alternative.
- Schreiben Sie das Label 12.

Aufgabe 8 Schreiben Sie in der Code-Generierung das Pattern-Case für Bedingungen.

Funktionsaufrufe Nicht als Aufgabe zu realisieren sind Funktionsaufrufe. Hierbei wird Code für die Argumente generiert. Diese werden zunächst in die Argumentregister geladen, wenn es weitere Argumente sind, werden sie auf den Stack gelegt. Als letztes wird die Anweisung `call` mit dem entsprechenden Funktionsnamen generiert

AST.java

```
case FunCall fc ->
{for (int i=0;i<Math.min(5, fc.args.size());i++) {
    fc.args.get(i).asm(fs, e, r);
    r.lnwrite("movq %rax, "+rs[i]);
}
;fc.args.stream().skip(rs.length).forEach(arg -> {
    arg.asm(fs, e, r);
    r.lnwrite("pushq %rax");
});
```

```

        ;r.lnwrite("call "+fc.name)
        ;}
    };
}

```

7 Eine main-Methode für Compiler oder Interpreter

Um alles zu testen, den Compiler, den Pretty-Printer und den Interpreter, dient diese main-Methode,

7.1 Hilfestellungen

Es gibt in der Interpreterschleife eine Hilfestellung.

AST.java

```

static String[] help =
{ "<expr>          evaluate <expr>"
, "'q'           quits the interpreter"
, "'defs'        System.out.println(); shows defined function names"
, "'s <funname>' prints function definition of <funname>"
, "':?'         this help"};

static void printHelp(){
    for (var h:help) System.out.println(h);
}

```

7.2 Anwender-Frontend des Systems

Die Methode kennt zwei Modi. Die Compilierung und die Interpreterschleife. Beide lesen eine Datei mit Funktionsdefinitionen ein.

AST.java

```

public static void main(String[] args) throws Exception {
    if (args.length==0){
        System.out.println
            ("usage: java name.panitz.longStack.AST [-i] filename");
        System.out.println
            ("  where -i starts interpreter "
            +"otherwise assembler code is generated");
    }
}

```

```
    return;  
}
```

Zunächst wird geprüft, ob Assembler generiert werden soll oder eine Interpreter-Session gestartet werden soll.

AST.java

```
var interpreter = args[0].equals("-i");  
var funDefs =  
    LongStackParser.parseFunDefs(args[0].equals("-i"?args[1]:args[0]));
```

Zum Start des Interpreters wird auf die Hilfe verwiesen.

AST.java

```
if (interpreter){  
    var in = new BufferedReader(new InputStreamReader(System.in));  
    System.out.println("TUGS (Tiny Usable Great System) :? for help");
```

Es folgt die eigentliche Lesen-Evaluieren-Drucken-Schleife (*read eval print loop* kurz *REPL*).

Die nächste Zeile wird gelesen:

AST.java

```
while (true){  
    System.out.print("> ");  
    var ln = in.readLine();
```

Soll das Programm beendet werden?

AST.java

```
if (ln.equals(":q")) break;
```

Braucht jemand Hilfe?

AST.java

```
if (ln.equals(":?")){printHelp();continue;}
```

Sollen die Funktionsdefinitionen aufgelistet werden?

AST.java

```
if (ln.equals(":defs")){
    funDefs.parallelStream()
        .forEach(fd->System.out.println(fd.name));
    continue;
}
var showFunction = ln.startsWith(":s ");
try{
```

Soll eine Funktion angezeigt werden?

AST.java

```
if (showFunction){
    var funname = ln.substring(2).trim();
    Fun fundef = funDefs.stream()
        .reduce
            (null
            , (r, fd) -> (r == null && fd.name.equals(funname)) ? fd : r
            , (fd1, fd2) -> fd1 == null ? fd2 : fd1);
    if (fundef != null)
        System.out.println(show(fundef)); ^^I
    else System.out.println("unknown function: "+funname);
}
```

Es soll also schließlich ein Ausdruck ausgewertet werden:

AST.java

```
}else{
    System.out.println
        (LongStackParser.parseExpr(ln).ev(funDefs));
}
}catch (Exception e){
    System.out.println(e);
}
}
```

Wenn keine REPL gestartet wurde, sondern Assembler generiert werden soll:

AST.java

```
^^I
}else{
    var out = new FileWriter
        (args[0].substring(0, args[0].lastIndexOf('.') + ".s");
    var o = new ExWriter(out);
    for (var fun: funDefs) asm(fun, o);
}
```

```
        o.nl();  
        out.close();  
    }  
}
```




Syntaxbaum einer Sprache
Pattern Matching
Pretty Printing
Interpretieren
Assembler Erzeugen