### Albert Ludwigs Universität Freiburg

#### TECHNISCHE FAKULTÄT

### PicoC-Compiler

## Übersetzung einer Untermenge von C in den Befehlssatz der RETI-CPU

BACHELORARBEIT

 $Abgabedatum: 28^{th}$  April 2022

Author: Jürgen Mattheis

Gutachter: Prof. Dr. Scholl

Betreung: M.Sc. Seufert

Eine Bachelorarbeit am Lehrstuhl für Betriebssysteme

#### **ERKLÄRUNG**

Hiermit erkläre ich, dass ich diese Abschlussarbeit selbständig verfasst habe, keine anderen als die angegebenen Quellen/Hilfsmittel verwendet habe und alle Stellen, die wörtlich oder sinngemäß aus veröffentlichten Schriften entnommen wurden, als solche kenntlich gemacht habe. Darüber hinaus erkläre ich, dass diese Abschlussarbeit nicht, auch nicht auszugsweise, bereits für eine andere Prüfung angefertigt wurde.

## Danksagungen

asdf

## Inhaltsverzeichnis

$\mathbf{A}$	bbild	ıngsverzeichnis	I
C	odeve	rzeichnis	II
Ta	abelle	nverzeichnis	III
D	efinit	onsverzeichnis	IV
$\mathbf{G}$	ramn	atikverzeichnis	V
1	Mot	vation	1
	1.1	RETI-Architektur	2
	1.2	Die Sprache PicoC	4
	1.3	Eigenheiten der Sprachen C und PicoC	5
	1.4	Gesetzte Schwerpunkte	11
	1.5	Über diese Arbeit	12
		1.5.1 Still der Schrifftlichen Ausarbeitung	13
		1.5.2 Aufbau der Schrifftlichen Arbeit	13
2	Imp	ementierung	15
	2.1	Lexikalische Analyse	17
		2.1.1 Konkrette Grammatik für die Lexikalische Analyse	17
		2.1.2 Codebeispiel	19
	2.2	Syntaktische Analyse	20
		2.2.1 Umsetzung von Präzidenz und Assoziativität	20
		2.2.2 Konkrette Grammatik für die Syntaktische Analyse	25
		2.2.3 Ableitungsbaum Generierung	27
		2.2.3.1 Codebeispiel	29
		2.2.3.2 Ausgabe des Ableitunsgbaumes	30
		2.2.4 Ableitungsbaum Vereinfachung	30
		2.2.4.1 Codebeispiel	$\frac{32}{33}$
		2.2.5.1 PicoC-Knoten	აა 35
		2.2.5.2 RETI-Knoten	40
		2.2.5.3 Kompositionen von PicoC-Knoten und RETI-Knoten mit besonderer Bedeutung	
		2.2.5.4 Abstrakte Grammatik	43
		2.2.5.5 Codebeispiel	45
		2.2.5.6 Ausgabe des Abstrakten Syntaxbaumes	46
	2.3	Code Generierung	46
		2.3.1 Passes	48
		2.3.1.1 PicoC-Shrink Pass	48
		2.3.1.1.1 Aufgabe	48
		2.3.1.1.2 Abstrakte Grammatik	49
		2.3.1.1.3 Codebeispiel	50
		2.3.1.2 PicoC-Blocks Pass	52
		2.3.1.2.1 Aufgabe	52
		2 3 1 2 2 Abstrakte Grammatik	52

	2.3.1.2.3	Codebeispiel
	2.3.1.3 PicoC-A	ANF Pass
	2.3.1.3.1	Aufgabe
	2.3.1.3.2	Abstrakte Grammatik
	2.3.1.3.3	Codebeispiel
	2.3.1.4 RETI-H	Blocks Pass
	2.3.1.4.1	Aufgabe
	2.3.1.4.2	Abstrakte Grammatik
	2.3.1.4.3	Codebeispiel
	2.3.1.5 RETI-H	Patch Pass
	2.3.1.5.1	Aufgabe
	2.3.1.5.2	Abstrakte Grammatik
	2.3.1.5.3	Codebeispiel
	2.3.1.6 RETI I	Pass
	2.3.1.6.1	Aufgabe
	2.3.1.6.2	Konkrette und Abstrakte Grammatik 6
	2.3.1.6.3	Codebeispiel
Literatur		

## Abbildungsverzeichnis

1.1	Schritte zum Ausführen eines Programmes mit dem GCC
	Stark vereinfachte Schritte zum Ausführen eines Programmes
1.3	Speicherorganisation
1.4	README.md im Github Repository der Bachelorarbeit
2.1	Ableitungsbäume zu den beiden Ableitungen
	Ableitungsbaum nach Parsen eines Ausdrucks
2.3	Ableitungsbaum nach Vereinfachung
2.4	Generierung eines Abstrakten Syntaxbaumes ohne Umdrehen
2.5	Generierung eines Abstrakten Syntaxbaumes mit Umdrehen
2.6	Cross-Compiler Kompiliervorgang ausgeschrieben
2.7	Cross-Compiler Kompiliervorgang Kurzform
	Architektur mit allen Passes ausgeschrieben

## Codeverzeichnis

1.1	Beispiel für Spiralregel
1.2	Ausgabe von Beispiel für Spiralregel $\dots$
1.3	Beispiel für unterschiedliche Ausführung
1.4	Ausgabe des Beispiels für unterschiedliche Ausführung
1.5	Beispiel mit Dereferenzierungsoperator
1.6	Ausgabe des Beispiels mit Dereferenzierungsoperator
1.7	Beispiel dafür, dass Struct kopiert wird 8
1.8	Ausgabe von Beispiel, dass Struct kopiert wird
1.9	Beispiel dafür, dass Zeiger auf Feld übergeben wird
1.10	Ausgabe von Beispiel dafür, dass Zeiger auf Feld übergeben wird
1.11	Beispiel für Deklaration und Definition
	Ausgabe von Beispiel für Deklaration und Definition
1.13	Beispiel für Sichtbarkeitsbereichs
	Ausgabe von Beispiel für Sichtbarkeitsbereichs
2.1	PicoC-Code des Codebeispiels
2.2	Tokens für das Codebeispiel
2.3	Ableitungsbaum nach Ableitungsbaum Generierung
2.4	Ableitungsbaum nach Ableitungsbaum Vereinfachung
2.5	Aus einem vereinfachtem Ableitungsbaum generierter Abstrakter Syntaxbaum
2.6	PicoC Code für Codebespiel
2.7	Abstrakter Syntaxbaum für Codebespiel
2.8	PicoC-Blocks Pass für Codebespiel
2.9	PicoC-ANF Pass für Codebespiel
2.10	RETI-Blocks Pass für Codebespiel
2.11	RETI-Patch Pass für Codebespiel
2.12	RETI Pass für Codebespiel

## **Tabellenverzeichnis**

1.1	Präzidenzregeln von PicoC
2.1	Präzidenzregeln von PicoC
2.2	Zuordnung der Bezeichnungen von Produktionsregeln zu Operatoren
2.3	PicoC-Knoten Teil 1
2.4	PicoC-Knoten Teil 2
2.5	PicoC-Knoten Teil 3
2.6	PicoC-Knoten Teil 4
2.7	RETI-Knoten
2.8	Kompositionen von PicoC-Knoten und RETI-Knoten mit besonderer Bedeutung

## Definitionsverzeichnis

1.1	Imperative Programmierung
1.2	Strukturierte Programmierung
1.3	Prozedurale Programmierung
1.4	Call by Value
1.5	Call by Reference
1.6	Funktionsprototyp
1.7	Deklaration
1.8	Definition
1.9	Sichtbarkeitsbereich (bzw. engl. Scope)
2.1	Metasyntax
2.2	Metasprache
2.3	Erweiterte Backus-Naur-Form (EBNF)
2.4	Dialekt der EBNF aus Lark
2.5	Abstrakte Syntax Form (ASF)
2.6	Earley Parser
2.7	Earley Recognizer
2.8	Label
2.9	Symboltabelle

## Grammatikverzeichnis

$2.1.1$ Konkrette Grammatik der Sprache $L_{PicoC}$ für die Lexikalische Analyse in EBNF	19
$2.2.1$ Undurchdachte Konkrette Grammatik der Sprache $L_{PicoC}$ für die Syntaktische Analyse in	
EBNF, die Operatorpräzidenz nicht beachtet	21
$2.2.2$ Erster Schritt zu einer durchdachten Konkretten Grammatik der Sprache $L_{PicoC}$ für die	
Syntaktische Analyse in EBNF, die Operatorpräzidenz beachtet	22
2.2.3 Beispiel für eine unäre rechtsassoziative Produktion	23
2.2.4 Beispiel für eine unäre linksassoziative Produktion	23
2.2.5 Beispiel für eine binäre linksassoziative Produktion	24
2.2.6 Beispiel für eine binäre linksassoziative Produktion ohne Linksrekursion	24
$2.2.7$ Durchdachte Konkrette Grammatik der Sprache $L_{PicoC}$ in EBNF, die Operatorpräzidenz beachtet	25
$2.2.8$ Konkrette Grammatik der Sprache $L_{PicoC}$ für die Syntaktische Analyse in EBNF, Teil $1$	26
$2.2.9$ Konkrette Grammatik der Sprache $L_{PicoC}$ für die Syntaktische Analyse in EBNF, Teil $2$	27
2.2.10Abstrakte Grammatik der Sprache $L_{PiocC}$	44
2.3.1 Abstrakte Grammatik der Sprache $L_{PiocC\_Shrink}$	50
	53
2.3.3 Abstrakte Grammatik der Sprache $L_{PiocC\_ANF}$	57
	60
	64
2.3.6 Konkrette Grammatik der Sprache $L_{RETI}$ für die Lexikalische Analyse in EBNF	68
1 IUDII V	68
2.3.8 Abstrakte Grammatik der Sprache $L_{RETI}$	69

## 1 Motivation

Als Programmierer kommt man nicht drumherum einen Compiler zu nutzen, er ist geradezu essentiel für den Beruf oder das Hobby des Programmierens. Selbst in der Programmiersprachen  $L_{Python}$ , welche als interpretierte Sprache bekannt ist, wird das in der Programmiersprache  $L_{Python}$  geschriebene Programm vorher zu Bytecode kompiliert, bevor dieser von der Python Virtual Machine (PVM) interpretiert wird.

Compiler, wie der  $GCC^1$  oder  $Clang^2$  werden üblicherweise über eine Commandline-Schnittstelle verwendet, welche es für den Benutzer unkompliziert macht ein Programm, dass in der Programmiersprache geschrieben ist, die der Compiler kompiliert $^3$  zu Maschinencode zu kompilieren.

Meist funktioniert das über schlichtes und einfaches Angeben der Datei, die das Programm enthält, welches kompiliert werden soll, z.B. im Fall des GCC über > gcc program.c -o machine\_code 4. Als Ergebnis erhält man im Fall des GCC die mit der Option -o selbst benannte Datei machine\_code, welche dann zumindest unter Unix über > ./machine\_code ausgeführt werden kann, wenn das Ausführungsrecht gesetzt ist. Das gesamte gerade erläuterte Vorgehen ist in Abbildung 1.1 veranschaulicht.

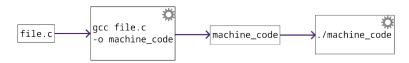


Abbildung 1.1: Schritte zum Ausführen eines Programmes mit dem GCC

Der ganze Kompiliervorgang kann, wie er in Abbildung 1.2 dargestellt ist zu einer Box abstrahiert werden. Der Benutzer gibt ein **Programm** in der Sprache des Compilers rein und erhält **Maschinencode**, den er dann im besten Fall in eine andere Box hineingeben kann, welche die passende **Maschine** oder den passenden **Interpreter** in Form einer **Virtuellen Maschine** repräsentiert, der bzw. die den **Maschinencode** ausführen kann.

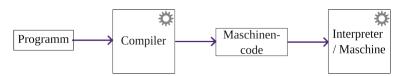


Abbildung 1.2: Stark vereinfachte Schritte zum Ausführen eines Programmes

<sup>&</sup>lt;sup>1</sup>GCC, the GNU Compiler Collection - GNU Project.

 $<sup>^2</sup>$  clang: C++ Compiler.

<sup>&</sup>lt;sup>3</sup>Im Fall des GCC und Clang ist es die Programmiersprache  $L_C$ .

<sup>&</sup>lt;sup>4</sup>Bei mehreren Dateien ist das ganze allerdings etwas komplizierter, weil der GCC beim Angeben aller .c-Dateien nacheinander gcc program\_1.c ... program\_n.c nicht darauf achtet doppelten Code zu entfernen. Beim GCC muss am besten mittels einer Makefile dafür gesorgt werden, dass jede Datei einzeln zu Objectcode (Definition ??) kompiliert wird. Das Kompilieren zu Objectcode geht mittels des Befehls gcc -c program\_1.c ... program\_n.c und alle Objectdateien können am Ende mittels des Linkers mit dem Befehl gcc -o machine\_code program\_1.o ... program\_n.o zusammen gelinkt werden.

Kapitel 1. Motivation 1.1. RETI-Architektur

Der Programmierer muss für das Vorgehen in Abbildung 1.2 nichts über die Theoretischen Grundlagen des Compilerbau wissen, noch wie der Compiler intern umgesetzt ist. In dieser Bachelorarbeit soll diese Compilerbox allerdings geöffnet werden und anhand eines eigenen im Vergleich zum GCC im Funktionsumfang reduzierten Compilers gezeigt werden, wie so ein Compiler unter der Haube stark vereinfacht funktionieren könnte.

Die konkrette Aufgabe besteht darin einen sogenannten PicoC-Compiler zu implementieren, der die Programmiersprache  $L_{PicoC}$ , welche eine Untermenge der Sprache  $L_C$  ist<sup>5</sup> in eine zu Lernzwecken prädestinierte, unkompliziert gehaltene Maschinensprache  $L_{RETI}$  kompilieren kann. Im Unterkapitel 1.1 wird näher auf die RETI-Architektur eingegangen, die der Sprache  $L_{RETI}$  zu Grunde liegt und im Unterkapitel 1.2 wird näher auf die die Sprache  $L_{PicoC}$  eingegangen, welche der PicoC-Compiler zur eben erwähnten Sprache  $L_{RETI}$  kompilieren soll.

#### 1.1 RETI-Architektur

Die RETI-Architektur ist eine zu Lernzwecken für die Vorlesungen P. D. C. Scholl, "Betriebssysteme" und P. D. C. Scholl, "Technische Informatik" entwickelte 32-Bit Architektur, die sich vor allem durch ihre einfache Zugänglichkeit kennzeichnet und deren Maschinensprache  $L_{RETI}$  als Zielsprache des PicoC-Compilers hergenommen wurde. In der Vorlesung P. D. C. Scholl, "Technische Informatik" wird die grundlegende RETI-Architektur erklärt und in der Vorlesung P. D. C. Scholl, "Betriebssysteme" wird diese Architektur erweitert, sodass diese mehr darauf angepasst ist, dass auch komplexere Kontrukte, wie ein Betriebssystem, Interrupts, Prozesse, Funktionen usw. auf nicht zu komplizierte Weise implementiert werden können.

Um den den PicoC-Compiler zu testen war es notwendig einen RETI-Interpreter zu implementieren, der genau die Variante der RETI-Achitektur aus der Vorlesung P. D. C. Scholl, "Betriebssysteme" simuliert.

#### Anmerkung 9

In dieser Bachelorarbeit wird im Folgenden bei der Maschinensprache  $L_{RETI}$  immer von der Variante, welche durch die RETI-Architektur aus der Vorlesung P. D. C. Scholl, "Betriebssysteme" umgesetzt ist ausgegangen.

Die Register der RETI-Architektur werden in Tabelle 2.1 aufgezählt und erläutert. Die Maschinenbefehle und Datenpfade der RETI-Architektur sind im Appendix ?? dokumentiert, da diese nicht explizit zum Verständnis der späteren Kapitel notwendig sind, aber zum vollständigen Verständnis notwendig sind, um die später auftauchenden RETI-Befehle usw. zu verstehen. Der Aufbau der Maschinensprache  $L_{RETI}$  ist durch Grammatik 2.3.6 und Grammatik 2.3.7 zusammengenommen beschrieben. Für genauere Implementierungsdetails ist allerdings auf die Vorlesungen P. D. C. Scholl, "Technische Informatik" und P. D. C. Scholl, "Betriebssysteme" zu verweisen.

2

<sup>&</sup>lt;sup>5</sup>Die der GCC kompilieren kann.

Kapitel 1. Motivation 1.1. RETI-Architektur

Register Kürzel	Register Ausgeschrieben	Aufgabe
PC	Program Counter	Zeigt auf den Maschinenbefehl, der als nächstes ausgeführt werden soll.
ACC	Accumulator	Für Operanden von Operationen oder für temporäre Werte.
IN1	Indexregister 1	Hat dieselbe Aufgabe wie das ACC-Register.
IN2	Indexregister 2	Hat dieselbe Aufgabe wie das ACC-Register.
SP	Stackpointer	Zeigt immer auf die erste freie Speicherzelle am Ende des Stacks, wo als nächstes Speicher allokiert werden kann.
BAF	Begin Aktive Funktion	Zeigt auf den Beginn des Stackframes der aktuell aktiven Funktion.
CS	Codesegment	Zeigt auf den Beginn des Codesegments. Die letzten 10 Bits werden verwendet, um 22 Bit Immediates aufzufüllen. Kann dadurch dazu verwendet werden, festzulegen welcher der 3 Peripheriegeräte <sup>a</sup> in der Memory Map <sup>b</sup> angesprochen werden soll.
DS	Datensegment	Zeigt auf den Beginn des Datensegments.

<sup>&</sup>lt;sup>a</sup> EPROM, UART und SRAM.

Tabelle 1.1: Präzidenzregeln von PicoC

Die RETI-Architektur ermöglicht bei der Ausführung von RETI-Programmen Prozesse zu nutzen. In Abbildung 1.3 ist der Aufbau eines Prozesses im Hauptspeicher der RETI-Architektur zu sehen. Das RETI-Programm nutzt dabei den Stack für temporäre Zwischenergebnisse von Berechnungen und zum Anlegen der Stackframes von Funktionen, welche die Lokalen Variablen und Parameter einer Funktion speichern. Das SP- und BAF-Register erfüllen dabei ihre zugeteilten Aufgaben für den Stack.

Der Abschnitt für die Globalen Statischen Daten ist allgemein dazu da Daten zu beherbergen, die für den Rest der Programmausführung global zugänglich sein sollen, aber auch für die Lokalen Variablen der main-Funktion. Das DS-Register markiert den Anfang des Datensegments und damit auch den Anfang der Globalen Statischen Daten und kann als relativer Orientierungspunkt beim Zugriff und Abspeichern Globaler Statischer Daten dienen. Das CS-Register wird als relativer Orientierungspunkt genutzt, an dem die Ausführung von RETI-Programmen startet und zur Bestimmung der relativen Startadresse, an welcher der RETI-Code einer bestimmten Funktion anfängt.

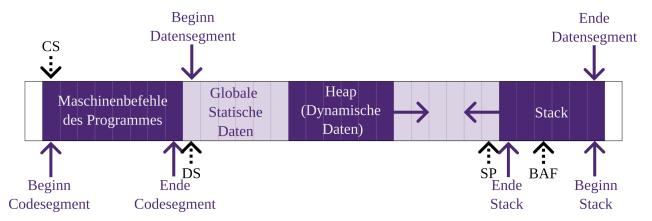


Abbildung 1.3: Speicherorganisation

b Da die Memory Map zum Verständnis der Bachelorarbeit nicht wichtig ist, wird diese nicht mehr als nötig im weiteren Verlauf erläutert.

Kapitel 1. Motivation 1.2. Die Sprache PicoC

Die RETI-Architektur nutzt 3 verschiedene Peripheriegeräte, EPROM, UART und SRAM, die über eine Memory Map<sup>6</sup> den über die Datenpfade der RETI-Architektur ?? ansprechbaren Adressraum von 2<sup>32</sup> Adressen<sup>7</sup> unter sich aufteilen.

Die Ausführung eines Programmes startet auf die einfachste Weise, indem es von einem Startprogramm im EPROM $^8$  aufgerufen wird. Der EPROM wird beim Start einer RETI-CPU als erstes aufgerufen, da nach der Memory Map der erste Adressraum von 0 bis  $2^{30}-1$  dem EPROM zugeordnet ist und das PC-Register initial den Wert 0 hat, also als erstes das Programm ausgeführt wird, welches an Adresse 0 im EPROM anfängt.

Die UART<sup>9</sup> ist eine elektronische Schaltung mit je nach Umsetzung mehr oder weniger Pins, wobei es allerdings immer einen RX- und einen TX-Pin gibt, für jeweils Empfangen<sup>10</sup> und Versenden<sup>11</sup> von Daten gibt. Jeder der Pins wird dabei mit einer anderen Adresse von 2<sup>3</sup> verschiedenen Adressen angsprochen und jeweils 8-Bit können nach den Datenpfaden ?? auf einmal über einen Pin in ein Register der UART geschrieben werden, um versandt zu werden oder von einem Pin empfangen werden. Die UART kann z.B. genutzt werden, um Daten an einen sehr einfach gehaltenen Monitor zu senden, der diese dann anzeigt.

An letzter Stelle muss der SRAM<sup>12</sup> erwähnt werden, bei dem es sich um den Hauptspeicher der RETI-CPU handelt. Der Zugriff auf den Hauptspeicher ist deutlich schneller als z.B. auf ein externes Speichermedium, aber langsamer als der Zugriff auf Register.

#### 1.2 Die Sprache PicoC

Die Sprache  $L_{PicoC}$  ist eine Untermenge der Sprache  $L_C$ , welche

- Einzeilige Kommentare // and Mehrzeilige Kommentare /\* and \*/
- die Primitiven Datentypen int, char und void
- die Abgeleiteten Datentypen Felder (z.B. int ar[3]), Verbunde (z.B. struct st {int attr1; attr2;}) und Zeiger (z.B. int \*pntr)
- if(cond){ }- / else{ }-Statements<sup>13</sup>
- while(cond){ }- und do while(cond){ };-Statements
- Arihmetische Ausdrücke, welche mithilfe der binären Operatoren +, -, \*, /, %, &, |, ^ und unären Operatoren -, ~ umgesetzt sind
- Logische Ausdrücke, welche mithilfe der Relationen ==, !=, <, >, <=, >= und Logischer Verknüpfungen !, &&, || umgesetzt sind
- Zuweisungen, die mit dem Zuweisungsoperator = umgesetzt sind
- Funktionsdeklaration (z.B. int fum(int arg1[3], struct st arg2);), Funktionsdefinition (z.B. int fum(int arg1[3], struct st arg2){}) und Funktionsaufrufe (z.B. fum(ar, st\_var))

<sup>&</sup>lt;sup>6</sup>Da die Memory Map zum Verständnis der Bachelorarbeit nicht wichtig ist, sondern nur bei der Umsetzung des RETI-Interpreters, wird diese nicht näher erläutert als notwendig.

<sup>&</sup>lt;sup>7</sup>Von 0 bis  $2^{32} - 1$ .

<sup>&</sup>lt;sup>8</sup>Kurz für Erasable Programmable Read-Only Memory.

 $<sup>^9 \</sup>mathrm{Kurz}$  für Universal Asynchronous Receiver Transmitter.

<sup>&</sup>lt;sup>10</sup>Engl. Receiving, daher das R.

<sup>&</sup>lt;sup>11</sup>Engl. Transmission, daher das T.

 $<sup>^{12}\</sup>mathrm{Kurz}$  für Static random-access memory.

<sup>&</sup>lt;sup>13</sup>Was die Kombination von if und else, nämlich else if (cond) { } miteinschließt.

beinhaltet. Die ausgegrauten • wurden bereits für das Bachelorprojekt umgesetzt und mussten für die Bachelorarbeit nur an die neue Architektur angepasst werden.

Der Aufbau der Programmiersprache  $L_C$  ist durch Grammatik 2.1.1 und Grammatik 2.2.8 zusammengenommen beschrieben.

#### 1.3 Eigenheiten der Sprachen C und PicoC

Einige Eigenheiten der Programmiersprache  $L_C$ , die genauso ein Teil der Programmiersprache  $L_{PicoC}$  sind, da  $L_{PicoC}$  eine Untermenge von  $L_C$  ist und welche in der Implementierung des PicoC-Compilers in Kapitel Implementierung noch eine wichtige Rolle spielen werden im Folgenden genauer erläutert. Im Folgenden wird immer von der Programmiersprache  $L_{PicoC}$  gesprochen, da es in dieser Bachelorarbeit um diese geht und die folgenden Beispiele für die Ausgaben alle mithilfe des PicoC-Compilers und RETI-Interpreters kompiliert bzw ausgeführt wurden, aber selbiges gilt genauso für  $L_C$  aus bereits erläutertem Grund.

Bei der Programmiersprache  $L_{PicoC}$  handelt es sich im eine imperative (Definition 1.1), strukturierte (Definition 1.2) und prozedurale Programmiersprache (Definition 1.3). Aufgrund dessen, dass es sich bei beiden um Imperative Programmiersprachen handelt ist es wichtig bei der Implementierung die Reihenfolge zu beachten und aufgrund dessen, dass es sich bei beiden um Strukturierte und Prozedurale Programmiersprachen handelt, ist es eine gute Methode bei der Implementierung auf Blöcke<sup>14</sup> zu setzen zwischen denen hin und her gesprungen werden kann und welche in den einzelnen Implementierungsschritten die notwendige Datenstruktur darstellen um Auswahl zwischen Codestücken, Wiederholung von Codestücken und Sprünge zu Blöcken mit entsprechend zu bestimmten Bezeichnern passenden Labeln umzusetzen.

#### Definition 1.1: Imperative Programmierung

Wenn ein Programm aus einer Folge von Befehlen besteht, deren Reihenfolge auch bestimmt in welcher Reihenfolge diese Befehle auf einer Maschine ausgeführt werden.<sup>a</sup>

<sup>a</sup>Thiemann, "Einführung in die Programmierung".

#### Definition 1.2: Strukturierte Programmierung

Wenn ein Programm anstelle von z.B. goto label-Statements Kontrollstruturen, wie z.B. if (cond) {
} else { }, while(cond) { } usw. verwendet, welche dem Programmcode mehr Struktur geben, weil
die Auswahl zwischen Statements und die Wiederholung von Statements eine klare und eindeutige
Struktur hat, welche bei Umsetzung mit einem goto label-Statement nicht so eindeutig erkennbar
wäre und auch nicht umbedingt immer gleich aufgebaut wäre.

<sup>a</sup>Thiemann, "Einführung in die Programmierung".

#### Definition 1.3: Prozedurale Programmierung

Programme werden z.B. mittels Funktionen in überschaubare Unterprogramme bzw. Prozeduren aufgeteilt, die aufrufbar sind. Dies vermeidet einerseits redundanten Code, indem Code wiederverwendbar gemacht wird und andererseits erlaubt es z.B. Codestücke nach ihren Aufgaben zu abstrahieren, den Codestücken wird eine Aufgabe zugeteilt, sie werden zu Unterprogrammen gemacht und fortan über einen Bezeichner aufgerufen, was den Code deutlich überschaubarer macht. da man die Aufgabe eines Codestücks nun nur noch mit seinem Bezeichner assozieren muss.<sup>a</sup>

<sup>&</sup>lt;sup>14</sup>Werden später im Kapitel 2 genauer erklärt.

```
<sup>a</sup>Thiemann, "Einführung in die Programmierung".
```

In  $L_C$  ist die Bestimmung des Datentyps einer Variable etwas komplizierter als in manch anderen Programmiersprachen. Der Grund liegt darin, dass die eckigen [ $\langle i \rangle$ ]-Klammern zur Festlegung der Mächtigkeit eines Feldes hinter der Variable stehen:  $\langle remaining-datatype \rangle \langle var \rangle$ [ $\langle i \rangle$ ], während andere Programmiersprachen die eckigen [ $\langle i \rangle$ ]-Klammern vor die Variable schreiben  $\langle remaining-datatype \rangle$ [ $\langle i \rangle$ ]  $\langle var \rangle$ .

Werden die eckigen [<i>]-Klammern hinter die Variable geschrieben, ist es schwieriger den Datentyp abzulesen, als auch ein Programm zu implementieren was diesen erkennt. Damit ein Programmierer den Datentyp ablesen kann, kann dieser die Spiralregel verwenden, die unter Clockwise/Spiral Rule nachgelesen werden kann. Werden die eckigen [<i>]-Klammern hinter die Variable geschrieben, wirken diese zum verwechseln ähnlich zum <var>[<ii]-Operator für den Zugriff auf den Index eines Feldes. Wenn ein Ausdruck geschrieben wird, wie int ar[1] = {42} wird, ist dieser vom Ausdruck var[0] = 42 nur durch den Kontext um var[1] bzw. var[0] rum zu unterscheiden.

In Code 1.1 ist ein Beispiel zu sehen, indem die Variable complex\_var den Datentyp "Feld der Mächtigkeit 1 von Feldern der Mächtigkeit 2 von Zeigern auf Felder der Mächtigkeit 2 von Verbunden vom Typ st" hat. Ein Vorteil die eckigen [<i>]-Klammern hinter die Variable zu schreiben ist in der markierten Zeile in Code 1.1 zu sehen. Will man auf ein Element dieses Datentyps zugreifen (\*complex\_var[0][1])[1].attr, so ist der Ausdruck fast genau gleich aufgebaut, wie der Ausdruck für den Datentyp struct st (\*complex\_var[1][2])[2]. Die Ausgabe des Beispiels in Code 1.1 ist in Code 1.2 zu sehen.

```
1 struct st {int attr;};
2
3 void main() {
4   struct st st_var[2] = {{.attr=314}, {.attr=42}};
5   struct st (*complex_var[1][2])[2] = {{&st_var}};
6   print((*complex_var[0][1])[1].attr);
7 }
```

Code 1.1: Beispiel für Spiralregel

```
1 42
```

Code 1.2: Ausgabe von Beispiel für Spiralregel

In  $L_C$  ist die Ausführbarkeit einer Operation oder wie diese Operation ausgeführt wird davon abhängig, was für einen Datentyp die Variable in diesem Kontext der Operation hat. In dem Beispiel in Code 1.3 wird in Zeile 2 ein "Feld der Mächtigkeit 1 von Feldern der Mächtigkeit 2" und Zeile 3 ein "Zeiger auf Felder der Mächtigkeit 2" erstellt. In den markierten Zeilen wird zweimal in Folge die gleiche Operation  ${\bf var}[0][1]$  ausgeführt, allerdings hat die Operation aufgrund der unterschiedlichen Datentypen der Variablen einen unterschiedlichen Effekt.

In Zeile 4 wird ein normaler Zugriff auf den zweiten Eintrag im ersten Eintrag des Felds int ar[1][2] = {{314, 42}} durchgeführt und in Zeile 5 wird allerdings erst dem Zeiger int (\*pntr)[2] = &ar[0]; gefolgt und dann ein Zugriff auf den zweiten Eintrag im ersten Eintrag des Felds int ar[1][2] = {{314, 42}} durchgeführt. Beide Operationen haben, wie in Code 1.4 zu sehen ist die gleiche Ausgabe.

```
1 void main() {
2   int ar[1][2] = {{314, 42}};
3   int (*pntr)[2] = &ar[0];
4   print(ar[0][1]);
5   print(pntr[0][1]);
6 }
```

Code 1.3: Beispiel für unterschiedliche Ausführung

```
1 42 42
```

Code 1.4: Ausgabe des Beispiels für unterschiedliche Ausführung

Eine weitere interessante Eigenheit, die tätsächlich nur in der Untermenge von  $L_C$ , die  $L_{PicoC}$  darstellt gültig ist<sup>15</sup>, ist dass die Operationen  $\langle var \rangle$ [0][1] und  $*(*(\langle var \rangle + 0) + 1)$  aus Code 1.3 und Code 1.5 komplett austauschbar sind. Die Ausgabe in Code 1.4 ist folglich identisch zur Ausgabe in Code 1.6.

```
1 void main() {
2   int ar[1][2] = {{314, 42}};
3   int (*pntr)[2] = &ar[0];
4   print(*(*(ar+0)+1));
5   print(*(*(pntr+0)+1));
6 }
```

Code 1.5: Beispiel mit Dereferenzierungsoperator

```
1 42 42
```

Code 1.6: Ausgabe des Beispiels mit Dereferenzierungsoperator

In der Programmiersprache  $L_{PicoC}$  werden alle Argumente bei einem Funktionsaufruf nach der Call By Value-Strategie (Definition 1.4) übergeben. Ein Beispiel hierfür ist in Code 1.7 zu sehen. Hierbei wird ein Verbund struct st copyable\_ar = {.ar={314, 314}}; <sup>16</sup> an die Funktion fun übergeben. Hierzu wird der Verbund in den Stackframe der aufgerufenen Funktion fun kopiert und an den Parameter fun gebunden.

Wie an der Ausgabe in Code 1.7 zu sehen ist hat die Zuweisung copyable\_ar.ar[1] = 42 an den Parameter struct st copyable\_ar in der aufgerufenen Funktion fun keinen Einfluss auf die übergebene lokale Variable copyable\_ar der aufrufenden Funktion. Der Eintrag an Index 1 im Feld bleibt bei 314.

<sup>&</sup>lt;sup>15</sup>In der Sprache  $L_C$  gibt es einen Unterschied bei der Initialisierung bei z.B. int \*var = "string" und z.B. int var[1] = "string", der allerdings nichts mit den beiden Operatoren zu tuen hat, sondern mit der Initialisierung, bei der die Sprache  $L_C$  verwirrenderweise die eckigen Klammern [] genauso, wie beim Operator für den Zugriff auf einen Arrayindex, vor den Bezeichner schreibt (z.B. var[1]), obwohl es ein Derived Datatype ist.

 $<sup>^{16}</sup>$ Später wird darauf eingegangen, warum der Verbund den Bezeichner  $copyable\_ar$  erhalten hat.

#### Definition 1.4: Call by Value

Z

Es wird eine Kopie des Ergebnisses eines Ausdrucks, welcher ein Argument eines Funktionsaufrufes darstellt an den entsprechenden Parameter der aufgerufenen Funktion gebunden.

Das hat zur Folge, dass bei Übergabe einer Variable als Argument an eine Funktion, diese Variable bei Änderungen am entsprechenden Parameter der aufgerufenen Funktion in der aufrufenden Funktion unverändert bleibt.<sup>a</sup>

<sup>a</sup>Bast, "Programmieren in C".

```
1 struct st {int ar[2];};
2
3 int fun(struct st copyable_ar) {
4   copyable_ar.ar[1] = 42;
5 }
6
7 void main() {
8   struct st copyable_ar = {.ar={314, 314}};
9   print(copyable_ar.ar[1]);
10   fun(copyable_ar);
11   print(copyable_ar.ar[1]);
12 }
```

Code 1.7: Beispiel dafür, dass Struct kopiert wird

```
1 314 314
```

Code 1.8: Ausgabe von Beispiel, dass Struct kopiert wird

In der Programmiersprache  $L_{PicoC}$  gibt es kein Call by Reference (Definition 1.5), allerdings kann der Effekt von Call by Reference mittels Zeigern simuliert werden, wie es in Code 1.11 bei der Funktion fun\_declared\_before und dem Parameter int \*param zu sehen ist. Genau dieser Trick wird bei Feldern verwendet, um nicht das gesamte Feld bei einem Funktionsaufruf in den Stackframe der aufgerufenen Funktion fun kopieren zu müssen.

Wie im Beispiel in Code 1.9 zu sehen ist, wird in der markierten Zeile ein Feld int ar[2] = {314, 314} an die Funktion fun übergeben. Wie in der Ausgabe in Code 1.10 zu sehen ist, hat sich der Eintrag an Index 1 im Feld nach dem Funktionsuaufruf zu 42 geändert. Wird ein Feld direkt als Ausdruck ar ohne z.B. die eckigen []-Klammern für einen Indexzugriff hingeschrieben wird die Adresse des Felds verwendet und nicht z.B. der erste Eintrag des Felds.

Eine Möglichkeit ein Feld als Kopie und nicht als Referenz zu übergeben ist es, wie in Code 1.7 das Feld als Attribut eines Verbundes zu übergeben, wie bei der Variable copyable\_ar.

#### Definition 1.5: Call by Reference

Z

Es wird eine implizite Referenz einer Variable, welche ein Argument eines Funktionsaufrufes darstellt an den entsprechenden Parameter der aufgerufenen Funktion gebunden.

Implizit meint hier, dass der Benutzer einer Programmiersprache mit Call by Reference nicht mitbekommt, dass er das Argument selbst verändert und keine lokale Kopie des Arguments.<sup>a</sup>

<sup>a</sup>Bast, "Programmieren in C".

```
int fun(int ar[2]) {
    ar[1] = 42;
    }

void main() {
    int ar[2] = {314, 314};
    print(ar[1]);
    fun(ar);
    print(ar[1]);
}
```

Code 1.9: Beispiel dafür, dass Zeiger auf Feld übergeben wird

```
1 314 42
```

Code 1.10: Ausgabe von Beispiel dafür, dass Zeiger auf Feld übergeben wird

Ein Programm in der Programmiersprache  $L_{PicoC}$  wird von oben-nach-unten ausgewertet. Ein Problem tritt auf, wenn z.B. eine Funktion verwendet werden soll, die aber erst unter dem entsprechenden Funktionsaufruf definiert (Definition 1.8) wird. Es ist wichtig, dass der Prototyp (Definition 1.6) einer Funktion vorher durch die Funktionsdefinition bekannt ist, damit überprüft werden kann, ob die beim Funktionsaufruf übergebenen Argumente den gleichen Datentyp haben, wie die Parameter des Prototyps und ob die Anzahl Argumente mit der Anzahl Parameter des Prototyps übereinstimmt.

Allerdings lassen sich die Funktionen nicht immer so anordnen, dass jede in einem Funktionsaufruf referenzierte Funktion vorher definiert sein kann. Aus diesem Grund ist es möglich den Prototyp einer Funktion vorher zu deklarieren (Definition 1.7), wie es in den markierten Zeile im Beispiel in Code 1.11 zu sehen ist. Die Ausgabe des Beispiels ist in Code 1.12 zu sehen.

#### Definition 1.6: Funktionsprototyp

Deklaration einer Funktion, welche den Funktionsbezeichner, die Datentypen der einzelnen Funktionsparameter, die Parametereihenfolge und den Rückgabewert einer Funktion spezifiziert. Es ist nicht möglich zwei Funktiondeklarationen mit dem gleichen Funktionsbezeichner zu haben. ab

<sup>&</sup>lt;sup>a</sup>Der Funktionsprototyp ist von der Funktionsignatur zu unterschieden, die in Programmiersprache wie C++ und Java für die Auflösung von Überladung bei z.B. Methoden verwendet wird und sich in manchen Sprachen für den Rückgabewert interessiert und in manchen nicht, je nach Umsetzung. In solchen Sprachen ist es möglich mehrere Methoden oder Funktionen mit dem gleichen Bezeichner zu haben, solange sie sich durch die Datentpyen von Parametern, die Parameterreihenfolge, manchmal auch Scopes und Klassentpyen usw. unterschieden.

<sup>&</sup>lt;sup>b</sup>What is the difference between function prototype and function signature?

#### Definition 1.7: Deklaration

1

Der Datentyp bzw. Prototyp einer Variablen bzw. Funktion, sowie der Bezeichner dieser Variable bzw. Funktion wird dem Compiler mitgeteilt. ab c

 $^a$ Über das Schlüsselwort extern lassen sich in der Programiersprache  $L_C$  Veriablen deklarieren, ohne sie zu definieren.

- <sup>b</sup> Variablen in C und C++, Deklaration und Definition Coder-Welten.de.
- <sup>c</sup>P. D. P. Scholl, "Einführung in Embedded Systems".

#### Definition 1.8: Definition

Dem Compiler wird mitgeteilt, dass zu einem bestimmten Zeitpunkt in der Programmausführung Speicher für eine Variable angelegt werden soll und wo<sup>a</sup> dieser angelegt werden soll. Eine Funktion ist definiert ihr eine relative Anfangsadresse im Hauptspeicher zugewiesen werden kann, aber welcher die Maschinenbefehle für diese Funktion abgespeichert werden können. bc

<sup>a</sup>Im Fall des PicoC-Compilers in den Globalen Statischen Daten oder auf dem Stack.

- $^b$  Variablen in C und C++, Deklaration und Definition Coder-Welten.de.
- <sup>c</sup>P. D. P. Scholl, "Einführung in Embedded Systems".

```
void fun_declared_before(int *param);

int fun_defined(int param) {
   return param + 10;
}

void main() {
   int res = fun_defined(22);
   fun_declared_before(&res);
   print(res);
}

void fun_declared_before(int *param) {
   *param = *param + 10;
}
```

Code 1.11: Beispiel für Deklaration und Definition

```
1 42
```

Code 1.12: Ausgabe von Beispiel für Deklaration und Definition

In  $L_{PicoC}$  lässt sich eine definierte Variable nur innerhalb ihres Sichtbarkeitsbereichs (Definition 1.9) verwenden. Lokale Variablen und Parameter lassen sich nur innerhalb der Funktion in welcher sie deklariert bzw. definiert wurden verwenden. Der Sichtbarkeitsbereich von Lokalen Variablen und Parametern erstreckt sich herbei von der öffnenden {-Klammer bis zur schließenden }-Klammer der Funktionsdefinition, in welcher sie definiert wurden.

Verschiedene Sichtbarkeitsbereiche können dabei identische Bezeichner besitzen. Im Beispiel in Code 1.13 kommt der markierte Bezeichner local\_var in 2 verschiedenen Sichtbarkeitsbereichen vor, doch bezeichnet er 2 unterschiedliche Variablen. Der Parameter param und die Lokale Variable local\_var dürfen nicht

den gleichen Bezeichner haben, da sie sich im gleichen Sichtbarkeitsbereich der Funktion fun\_scope befinden. Die Ausgabe des Beispiels in Code 1.13 ist in Code 1.14 zu sehen.

```
Definition 1.9: Sichtbarkeitsbereich (bzw. engl. Scope)

Bereich in einem Programm, in dem eine Variable sichtbar ist und verwendet werden kann.

Thiemann, "Einführung in die Programmierung".

int fun_scope(int param) {

int local_var = 2;

print(param);

print(local_var);

}

void main() {

int local_var = 4;

fun_scope(local_var);

}
```

Code 1.13: Beispiel für Sichtbarkeitsbereichs

```
1 4 2
```

Code 1.14: Ausgabe von Beispiel für Sichtbarkeitsbereichs

#### 1.4 Gesetzte Schwerpunkte

Ein Schwerpunkt dieser Bachelorarbeit ist es in erster Linie bei der Kompilierung der Programmiersprache  $L_{PicoC}$  in die Maschinensprache  $L_{RETI}$  die Syntax und Semantik der Sprache  $L_C$  identisch nachzuahmen. Der PicoC-Compiler soll die Sprache  $L_{PicoC}$  im Vergleich zu z.B. dem  $GCC^{17}$  ohne merklichen Unterschied<sup>18</sup> komplieren können.

In zweiter Linie soll dabei möglichst immer so Vorgegangen werden, wie es die RETI-Codeschnipsel aus der Vorlesung P. D. C. Scholl, "Betriebssysteme" vorgeben. Allerdings sollten diese bei Inkonsistenzen bezüglich der durch sie selbst vorgegebenen Paradigmen und anderen Umstimmigkeiten angepasst werden, da der erstere Schwerpunkt überwiegt.

Des Weiteren ist die Laufzeit bei Compilern zwar vor allem in der Industrie nicht unwichtig, aber bei Compilern, verglichen mit Interpretern weniger zu gewichten, da ein Compiler bei einem fertig implementierten Programm nur einmal Maschinencode generieren muss und dieser Maschinencode danach fortan ausgeführt wird. Beim einem Compiler ist daher eher zu priorisieren, dass der kompilierte Maschinencode möglichst effizient ist.

 $<sup>^{17}</sup>$ Da die Sprache  $L_{PicoC}$  eine Untermenge von  $L_C$  ist, kann der GCC  $L_{PicoC}$  ebenfalls kompilieren, allerdings nicht in die gewünschte Maschinensprache  $L_{RETI}$ .

<sup>&</sup>lt;sup>18</sup>Natürlich mit Ausnahme der sich unterscheidenden Maschinensprachen zu welchen kompiliert wird und der unterschiedlichen Commandline-Optionen und Fehlermeldungen.

Kapitel 1. Motivation 1.5. Über diese Arbeit

#### 1.5 Über diese Arbeit

Der Quellcode des PicoC-Compilers ist öffentlich unter Link<sup>19</sup> zu finden. In der Datei README.md (siehe Abbildung 1.4) ist unter "Getting Started" ein kleines Einführungstutorial verlinkt. Unter "Usage" ist eine Dokumentation über die verschiedenen Command-line Optionen und verschiedene Funktionalitäten der Shell verlinkt. Deneben finden sich noch weitere Links zu möglicherweise interessanten Dokumenten. Der letzte Commit vor der Abgabe der Bachelorarbeit ist unter Link<sup>20</sup> zu finden.



Abbildung 1.4: README.md im Github Repository der Bachelorarbeit

Die Schrifftliche Ausarbeitung der Bachelorarbeit wurde ebenfalls veröffentlicht, falls Studenten, die den PicoC-Compiler in Zukunft nutzen sich in der Tiefe dafür interessieren, wie dieser unter der Haube funktioniert. Die Schrifftliche Ausarbeitung dieser Bachelorarbeit ist als PDF unter Link<sup>21</sup> zu finden. Die PDF der Schrifftliche Ausarbeitung der Bachleorararbeit wird aus dem Latexquellcode, welcher unter Link<sup>22</sup> veröffentlicht ist automatisch mithife der Github Action Nemec, copy\_file\_to\_another\_repo\_action und der Makefile Ueda, Makefile for LaTeX generiert.

Alle verwendeten Latex Bibliotheken sind unter Link<sup>23</sup> zu finden<sup>24</sup>. Die Grafiken, die nicht mittels der Tikz Bibliothek in Latex erstellt wurden, wurden mithilfe des Vectorgraphikeditors Inkscape<sup>25</sup> erstellt. Falls Interesse besteht Grafiken aus der Schrifftlichen Ausarbeitung der Bachelorarbeit zu verwenden, so sind diese zusammen mit den .svg-Dateien von Inkscape im Ordner /figures zu finden.

Alle weitere verwendete Software, wie verwendete Python Bibliotheken, Vim/Neovim Plugins, Tmux Plugins usw. sind in der README.md unter "References" bzw. direkt unter Link<sup>26</sup> zu finden.

 $<sup>^{19} \</sup>verb|https://github.com/matthejue/PicoC-Compiler.$ 

 $<sup>^{20} \</sup>texttt{https://github.com/matthejue/PicoC-Compiler/tree/bcafedffa9ff3075372554b14f1a1d369af68971.}$ 

<sup>&</sup>lt;sup>21</sup>https://github.com/matthejue/Bachelorarbeit\_out/blob/main/Main.pdf.

<sup>22</sup>https://github.com/matthejue/Bachelorarbeit.

 $<sup>^{23}</sup>$ https://github.com/matthejue/Bachelorarbeit/blob/master/content/Packete\_und\_Deklarationen.tex.

 $<sup>^{24}</sup>$ Jede einzelne verwendete Latex Bibliothek einzeln anzugeben wäre allerdings etwas zu aufwendig

 $<sup>^{25} \</sup>mathrm{Developers}, \ \mathit{Draw Freely} - \mathit{Inkscape}.$ 

 $<sup>^{26}</sup>$ https://github.com/matthejue/PicoC-Compiler/blob/new\_architecture/doc/references.md.

Kapitel 1. Motivation 1.5. Über diese Arbeit

#### 1.5.1 Still der Schrifftlichen Ausarbeitung

In dieser Schrifftliche Ausarbeitung der Bachelorarbeit sind die manche Wörter für einen besseren Lesefluss hervorgehoben. Es ist so gedacht, dass die Hervorgehobenen Wörter beim Lesen sichtbare Ankerpunkte darstellen an denen sich orientiert werden kann, aber auch damit der Inhalt eines vorher gelesener Paragraphs nochmal durch Überfliegen der Hervorgehobenen Wörter in Erinnerung gerufen werden kann.

Bei den Erklärungen wurden darauf geachtet bei jeder der verwendeten Methodiken und jeder Designentscheidung die Frage zu klären, "warum etwas geanu so gemacht wurde und nicht anders", denn wie es im Buch LeFever, *The Art of Explanation* auf eine deutlich ausführlichere Weise dargelegt wird, ist einer der zentralen Fragen, die ein Leser in erster Linie zum wirklichen Verständnis eines Themas beantwortet braucht<sup>27</sup> die Frage des "warum".

Zum Verweis auf Quellen an denen sich z.B. bei der Formulierung von Definitionen orientiert wurde, wurden um den Lesefluss nicht zu stören Fußnoten<sup>28</sup> verwendet. Die meisten Definitionen wurden in eigenen Worten formuliert, damit die Definitionen selbst zueinander konsistent sind, wie auch das in Ihnen verwendete Vokabular. Wurde eine Definition wörtlich aus einer Quelle übernomnen, so wurde die Definition oder der entsprechende Teil in "Anführungszeichen" gesetzt. Beim Verweis auf Quellen außerhalb einer Definitionsbox, wurde allerdings meistens, sofern die Quelle wirklich relevant war auf das Zitieren über Fußnoten verzichtet.

#### 1.5.2 Aufbau der Schrifftlichen Arbeit

Die Schrifftliche Ausarbeitung der Bachelorarbeit ist in 4 Kapitel unterteilt: Motivation, ??, Implementierung und ??.

Im momentanen Kapitel Motivation, wurde ein kurzer Einstieg in das Thema Compilerbau gegeben und die zentrale Aufgabenstellung der Bachelorarbeit erläutert, sowie auf Schwerpunkte und kleinere Teilprobleme, die eines besonderen Fokus bedürfen eingegangen.

Im Kapitel ?? werden die notwendigen Theoretischen Grundlagen eingeführt, die zum Verständnis des Kapitels Implementierung notwendig sind. Das Kapitel soll darüberhinaus aber auch einen Überblick über das Thema Compilerbau geben, sodass nicht nur ein Grundverständnis für das eine spezifische Vorgehen, welches zur Implementierung des PicoC-Compiler verwendet wurde vermittelt wird, sondern auch ein Vergleich zu anderen Vorgehensweisen möglich ist. Die Theoretischen Grundlagen umfassen die wichtigsten Definitionen in Bezug zu Compilern und den verschiedenen Phasen der Kompilierung, welche durch die Unterkapitel Lexikalische Analyse, Syntaktische Analyse und Code Generierung repräsentiert sind.

Des Weiteren wurden für T-Diagramme und Formale Sprachen eigene Unterkapitel erstellt. Für T-Diagramme wurde ein eigenes Unterkapitel erstellt, da sie häufig in der Schrifftlichen Ausarbeitung verwendet werden und die T-Diagramm Notation nicht allgemein bekannt ist. Für Formale Sprachen wurde ein eigenes Unterkapitel erstellt, da für den Gutachter Prof. Dr. Scholl das Thema Formale Sprachen eher fachfremd ist, aber dieses Thema einige zentrale und wichtige Fachbegriffe besitzt, bei denen es wichtig ist die genaue Definition zu haben. Generell wurde im Kapitel Einführung versucht an Erklärungen nicht zu sparren, damit aufgrund dessen, dass das Thema eher fachfremd für Prof. Dr. Scholl ist für das Kapitel Implementierung keine wichtigen Aspekte unverständlich bleiben.

Im Kapitel Implementierung werden die einzelnen Aspekte der Implementierung des PicoC-Compilers, unterteilt in die verschiedenen Phasen der Kompilierung nach dennen das Kapitel Einführung ebenfalls unterteilt ist erklärt. Dadurch, dass Kapitel Implementierung und Kapitel Einführung eine ähliche

<sup>&</sup>lt;sup>27</sup>Vor allem Anfang, wo der Leser wenig über das Thema weiß.

<sup>&</sup>lt;sup>28</sup>Das ist ein Beispiel für eine Fußnote.

Kapitel 1. Motivation 1.5. Über diese Arbeit

Kapiteleinteilung haben, ist es besonders einfach zwischen beiden hin und her zu wechseln. Wenn z.B. eine Definition im Kapitel Einführung gesucht wird, die zum Verständis eines Aspekts in Kapitel Implemenentierung notwendig ist, so kann aufgrund der ähnlichen Kapiteleinteilung die entsprechende Definition analog im Kapitel Einleitung gefunden werden.

Im Kapitel ?? wird ein Überblick über die wichtigsten Funktionalitäten des PicoC-Compilers gegeben, indem anhand kleiner Anleitungen gezeigt wird wie man diese verwendet. Des Weiteren wird darauf eingegangen, wie die Qualitätsicherung für den PicoC-Compiler umgesetzt wurde, also wie gewährleistet wird, dass der PicoC-Compiler funktioniert. Zum Schluss wird noch auf weitere Erweiterungsideen eingegangen, die auch interessant zu implementieren wären.

Im Kapitel ?? werden einige Definitionen und Themen angesprochen, die bei Interesse zur weiteren Vertiefung da sind und unabhänging von den anderen Kapiteln sind. Diese Themen und Definitionen sind dazu da den Bogen von der spezifischen Implementierung des PicoC-Compilers wieder zum allgemeinen Vorgehen bei der Implementierung eines Compilers zu schlagen. Diese Themen und Definitionen passen nicht ins Kapitel Implementierung, da diese selbst nichts mit der Implementierung des PicoC-Compilers zu tuen haben und auch nichts ins Kapitel ??, da dieses nur Theoretische Grundlagen erklärt, die für das Kapitel Implementierung wichtig sind.

Generell wurde in der Schrifftlichen Ausarbeitung immer versucht Parallelen zu Implementierung echter Compiler zu ziehen. Der Zweck des PicoC-Compilers ist es primär ein Lerntool zu sein, weshalb Methoden, wie Liveness Analyse (Definition ??) usw., die in echten Compilern zur Anwendung kommen nicht umgesetzt wurden, da sich an die vorgegebenen Paradigmen aus der Vorlesung P. D. C. Scholl, "Betriebssysteme" gehalten werden sollte.

# 2 Implementierung

In diesem Kapitel wird, nachdem im Kapitel ?? die nötigen theoretischen Grundlagen des Compilerbau vermittelt wurden, nun auf die Implementierung des PicoC-Compilers eingegangen. Aufgeteilt in die selben Kategorien Lexikalische Analyse 2.1, Syntaktische Analyse 2.2 und Code Generierung 2.3, wie in Kapitel ??, werden in den folgenden Unterkapiteln die einzelnen Zwischenschritte vom einem Programm in der Konkretten Syntax der Sprache  $L_{PicoC}$  hin zum einem Programm mit derselben Semantik in der Konkretten Syntax der Sprache  $L_{RETI}$  erklärt.

Für das Parsen<sup>1</sup> des Programmes in der Konkretten Syntax der Sprache  $L_{PicoC}$  wird das Lark Parsing Toolkit<sup>2 3</sup> verwendet. Das Lark Parsing Toolkit ist eine Bibliothek, die es ermöglicht mittels einer in einem eigenen Dialekt der Erweiterten Back-Naur-Form (Definition 2.3 bzw. für den Dialekt von Lark Definition 2.4) spezifizierten Konkretten Grammatik ein Programm in Konkretter Syntax zu parsen und daraus einen Ableitungsbaum für die kompilerintere Weiterverarbeitung zu generieren.

#### Definition 2.1: Metasyntax

Z

Steht für den Aufbau einer Metasprache (Definition 2.2), der durch eine Grammatik oder Natürliche Sprache beschrieben werden kann.

#### Definition 2.2: Metasprache

Z

Eine Sprache, die dazu genutzt wird andere Sprachen zu beschreiben<sup>a</sup>.

<sup>a</sup>Das "Meta" drückt allgemein aus, dass sich etwas auf einer höheren Ebene befindet. Um über die Ebene sprechen zu können, in der man sich selbst befindet, muss man von einer höheren, außenstehenden Ebene darüber reden.

#### Definition 2.3: Erweiterte Backus-Naur-Form (EBNF)



Die Erweiterte Backus-Naur-Form<sup>a</sup> ist eine Metasyntax (Definition 2.1), die dazu verwendet wird Kontextfreie Grammatiken darzustellen.<sup>bc</sup>

Die Erweiterte Backus-Naur-Form ist zwar standartisiert und die Spezifikation des Standards kann unter  $Link^d$  aufgefunden werden, allerdings werden in der Praxis, wie z.B. in Lark oft eigene Notationen verwendet.

 $<sup>^</sup>a$ Der Name kommt daher, dass es eine Erweiterung der Backus-Naur-Form ist, die hier allerdings nicht weiter erläutert wird.

 $<sup>{}^</sup>b$ Nebel, "Theoretische Informatik".

<sup>&</sup>lt;sup>c</sup>Grammar Reference — Lark documentation.

dhttps://standards.iso.org/ittf/PubliclyAvailableStandards/.

<sup>&</sup>lt;sup>1</sup>Wobei beim **Parsen** auch das **Lexen** inbegriffen ist.

 $<sup>^2</sup> Lark$  - a parsing toolkit for Python.

<sup>&</sup>lt;sup>3</sup>Shinan, lark.

#### Definition 2.4: Dialekt der EBNF aus Lark

7

Das Lark Parsing Toolkit verwendet eine eigene Notation für die Erweiterte Backus-Naur-Form, die sich teilweise in einzelnen Aspekten von der Syntax aus dem Standard unterscheidet und unter Link<sup>a</sup> dokumentiert ist.

Ein wichtiger Unterschied ist z.B., dass dieser Dialekt anstelle von geschweiften Klammern {} für die Darstellung von Wiederholung, den aus regulären Ausdrücken bekannten \*-Quantor optional zusammen mit runden Klammern () verwendet: ()\*.

Um bei einer Produktion auszudrücken, wozu die linke Seite abgeleitet werden kann, also "kann abgeleitet werden zu", wird das ::=-Symbol verwendet.

Das Lark Parsing Toolkit wurde vor allem deswegen gewählt, weil es sehr einfach in der Verwendung ist. Andere derartige Tools, wie z.B. ANTLR<sup>4</sup> sind Parser Generatoren, die zur Konkretten Grammatik einer Sprache einen Parser in einer vorher bestimmten Programmiersprache generieren, anstatt wie das Lark Parsing Toolkit bei Angabe einer Konkretten Grammatik direkt ein Programm in dieser Konkretten Grammatik parsen und einen Ableitungsbaum dafür generieren zu können.

Eine möglichst geringe Laufzeit durch Verwenden der effizientesten Algorithmen zu erreichen war keine der Hauptzielsetzungen für den PicoC-Compiler, da der PicoC-Compiler vor allem als Lerntool konzipiert ist, mit dem Studenten lernen können, wie der Kompiliervorgang von der Programmiersprache  $L_{PicoC}$  zur Maschinensprache  $L_{RETI}$  funktioniert. Eine ausführliche Diskussion zur Priorisierung Laufziet wurde in Unterkapitel ?? geführt. Lark besitzt des Weiteren eine sehr gute Dokumentation Welcome to Lark's documentation! — Lark documentation, sodass anderen Studenten, die den PicoC-Compiler vielleicht in ihr Projekt einbinden wollen, unkompliziert Erweiterungen für den PicoC-Compiler schreiben können.

Neben den Konkretten Grammatiken, die aufgrund der Verwendung des Lark Parsing Toolkit in einem eigenen Dialekt der Erweiterten Back-Naur-Form spezifiziert sind, werden in den folgenden Unterkapiteln die Abstrakten Grammatiken, welche spezifizieren, welche Kompositionen für die Abstrakten Syntaxbäume der verschiedenden Passes erlaubt sind in einer bewusst anderen Notation aufgeschrieben, die allerdings Ähnlichkeit mit dem Dialekt der Erweiterten Backus-Naur-Form aus dem Lark Parsing Toolkit hat.

Die Notation für die Abstrakte Syntax unterscheidet sich bewusst von der Erweiterten Backus-Naur-Form, da in der Abstrakten Syntax Kompositionen von Knoten beschrieben werden, die klar auszumachen sind, wodurch es die Abstrakten Grammatiken nur unnötig verkomplizieren würde, wenn man die Erweiterte Backus-Naur-Form verwenden würde. Es gibt leider keine Standardnotation für Abstrakte Grammatiken, die sich deutlich durchgesetzt hat, daher wird für Abstrakte Grammatiken eine eigene Abstrakte Syntax Form Notation (Definition 2.5) verwendet. Des Weiteren trägt das Verwenden einer unterschiedlichen Notation für Konkrette und Abstrakte Syntax auch dazu bei, dass man beide direkter voneinander unterscheiden kann.

#### Definition 2.5: Abstrakte Syntax Form (ASF)

Die Abstrakte Syntax Form ist eine eigene Metasyntax für Abstrakte Grammatiken, die für diese Bachelorarbeit definiert wurde und sich von dem Dialekt der Backus-Naur-Form des Lark Parsing Toolkit nur dadurch unterschiedet, dass Terminalsymbole nicht von "" engeschlossen sein müssen, da die Knoten in der Abstrakten Syntax, sowieso schon klar auszumachen sind und von anderen Symbolen der Metasprache leicht zu unterschiden sind.

<sup>&</sup>lt;sup>a</sup>https://lark-parser.readthedocs.io/en/latest/grammar.html.

<sup>&</sup>lt;sup>b</sup>Bzw. kann der \*-Quantor auch keinmal wiederholen bedeuteten.

 $<sup>^4</sup>ANTLR.$ 

Letztendlich geht es allerdings nur darum, dass aufgrund der Verwendung des Lark Parsing Toolkit die Konkrette Grammatik in einem eigenen Dialekt der Erweiterter Backus-Naur-Form angegeben sein muss und für das Implementieren der Passes die Abstrakte Grammatik für den Programmierer möglichst einfach verständlich sein sollte, weshalb sich die Abstrake Syntax Form gut dafür eignet.

#### 2.1 Lexikalische Analyse

Für die Lexikalische Analyse ist es nur notwendig eine Konkrette Grammatik zu definieren, die den Teil der Konkretten Syntax beschreibt, der die verschiedenen Pattern für die verschiedenen Token der Sprache  $L_{PicoC}$  beschreibt, also den Teil der für die Lexikalische Analyse wichtig ist. Diese Konkrette Grammatik wird dann vom Lark Parsing Toolkit dazu verwendet ein Programm in Konkretter Syntax zu lexen und daraus Tokens für die Syntaktische Analyse zu erstellen, wie es im Unterkapitel ?? erläutert ist.

#### 2.1.1 Konkrette Grammatik für die Lexikalische Analyse

In der Konkretten Grammatik 2.1.1 für die Lexikalische Analyse stehen großgeschriebene Nicht-Terminalsymbole entweder für einen Tokennamen oder einen Teil der Beschreibung eines Tokennamen. Zum Beispiel handelt es sich bei dem großgeschriebenen Nicht-Terminalsymbol NUM um einen Tokennamen, der durch die Produktion NUM ::= "0" | DIG\_NO\_0 DIG\_WITH\_0\* beschrieben wird und beschreibt, wie ein möglicher Tokenwert, in diesem Fall eine Zahl aufgebaut sein kann. Das ist daran festzumachen, dass das Nicht-Terminalsymbol NUM in keiner anderen Produktion vorkommt, die auf der linken Seite des ::=-Symbols ebenfalls ein großgeschriebenen Nicht-Terminalsymbol hat. Dagegen dient das großgeschriebene Nicht-Terminalsymbol DIG\_NO\_0 aus der Produktion NUM ::= "0" | DIG\_NO\_0 DIG\_WITH\_0\* nur zu Beschreibung von NUM.

Die in der Konkretten Grammatik 2.1.1 für die Lexikalische Analyse definierten Nicht-Terminalsymbole können in der Konkretten Grammatik 2.2.8 für die Syntaktischen Analayse verwendet werden, um z.B. zu beschreiben, in welchem Kontext z.B. eine Zahl NUM stehen darf.

Die in der Konkrette Grammatik vereinzelt kleingeschriebenen Nicht-Terminalsymbole, wie name haben nur den Zweck mehrere Tokennamen, wie NAME | INT\_NAME | CHAR\_NAME unter einem Überbegriff zu sammeln.

In Lark steht eine Zahl .Zahl, die an ein Nicht-Terminalsymbol angehängt ist, dass auf der linken Seite des ::=-Symbols einer Produktion steht für die Priorität der Produktion dieses Nicht-Terminalsymbols. Es gibt den Fall, dass ein Wort von mehreren Produktionen erkannt wird, z.B. wird das Wort int sowohl von der Produktion NAME, als auch von der Produktion INT\_DT erkannt. Daher ist es notwendig für INT\_DT eine Priorität INT\_DT.2 zu setzen<sup>5</sup>, damit das Wort int den Tokennamen INT\_DT zugewiesen bekommt und nicht NAME.

Allerdings muss für den Fall, dass int der Präfix eines Wortes ist, z.B. int\_var noch die Produktion INT\_NAME.3 definiert werden, da der im Lark Parsing Toolkit verwendete Basic Lexer sobald ein Wort von einer Produktion erkannt wird, diesem direkt einen Tokennamen zuordnet, auch wenn das Wort eigentlich von einer anderen Produktion erkannt werden sollte. In diesem Fall würden aus int\_var die Token Token('INT\_DT', 'int'), Token('NAME', '\_var') generiert, anstatt Token(NAME, 'int\_var'). Daher muss die Produktion INT\_NAME.3 eingeführt werden, die immer zuerst geprüft wird. Wenn es sich nur um das Wort int handelt, wird zuerst die Produktion INT\_NAME.3 geprüft, es stellt sich heraus, dass int von der Produktion INT\_NAME.3 nicht erkannt wird, daher wird als nächstes INT\_DT.2 geprüft, welches int erkennt.

Die Implementierung des Basic Lexer aus dem Lark Parsing Toolkit ist unter Link<sup>6</sup> zu finden ist. Diese

<sup>&</sup>lt;sup>5</sup>Es wird immer die höchste Priorität zuerst genommen.

<sup>6</sup>https://github.com/lark-parser/lark/blob/d03f32be7f418dc21cfa45acc458e67fe0580f60/lark/lexer.py

Implementierung ist allerdings zu spezifisch auf Lark zugeschnitten und ist aufgrund dessen, dass sie in der Lage ist nach einer spezifizierten Konkretten Grammatik zu lexen, zu komplex, um sie an dieser Stelle allgemein erklären zu können.

Der Basic Lexer verhält sich allerdings grundlegend so, wie es im Unterkapitel ?? erklärt wurde, allerdinds berücksichtigt der Basic Lexer ebenfalls Priortiäten, sodass für den aktuellen Index im Eingabeprogramm zuerst alle Produktionen der höchsten Priorität geprüft werden. Sobald eine dieser Produktionen ein Wort an dem aktuellen Index im Eingabeprogramm erkennt, bekommt es direkt den Tokenwert dieser Produktion zugewiesen. Weitere Produktionen werden nicht mehr geprüft. Ansonsten werden alle Produktionen der nächstniedrigeren Priorität geprüft usw.

```
/[\wedge \backslash n]*/
COMMENT
                                                  /(. | \n)*? / "*/"
                                                                           L_{-}Comment
                       ::=
                            "//""_{-}"?"#"/[\wedge \setminus n]*/
RETI\_COMMENT.2
                       ::=
                                           "3"
                                    "2"
DIG\_NO\_0
                       ::=
                            "1"
                                                   "4"
                                                           "5"
                                                                  "6"
                                                                            L_Arith
                            "7"
                                    "8"
                                            "9"
DIG\_WITH\_0
                            "0"
                                    DIG\_NO\_0
                       ::=
                            "0"
NUM
                                    DIG_NO_0 DIG_WITH_0*
                       ::=
                            " ".."~"
ASCII\_CHAR
                       ::=
                            "'"ASCII\_CHAR"'"
CHAR
                       ::=
FILENAME
                            ASCII\_CHAR + ".picoc"
                       ::=
                            "a"..."z" | "A"..."Z"
LETTER
                       ::=
                            (LETTER | "_")
NAME
                       ::=
                                 (LETTER | DIG_WITH_0 | "_")*
                            NAME | INT_NAME | CHAR_NAME
name
                       ::=
                            VOID\_NAME
                            " | "
LOGIC\_NOT
                       ::=
                            " \sim "
NOT
                       ::=
                            "&"
REF\_AND
                       ::=
un\_op
                       ::=
                            SUB\_MINUS \mid LOGIC\_NOT \mid NOT
                            MUL\_DEREF\_PNTR \mid REF\_AND
                            "*"
MUL\_DEREF\_PNTR
                       ::=
                            "/"
DIV
                       ::=
                            "%"
MOD
                       ::=
prec1\_op
                            MUL\_DEREF\_PNTR \mid DIV \mid
                                                               MOD
                       ::=
ADD
                            "+"
                       ::=
                            "_"
SUB\_MINUS
                       ::=
                                      SUB\_MINUS
prec2\_op
                       ::=
                            ADD
                            "<"
LT
                       ::=
                                                                            L\_Logic
                            "<="
LTE
                       ::=
                            ">"
GT
                       ::=
GTE
                            ">="
                       ::=
                            LT
                                   LTE \mid GT \mid GTE
rel\_op
                       ::=
EQ
                       ::=
                            "=="
NEQ
                            "!="
                       ::=
                            EQ
                                    NEQ
eq\_op
                       ::=
                            "int"
INT\_DT.2
                                                                            L\_Assign\_Alloc
                       ::=
INT\_NAME.3
                            "int" (LETTER \mid DIG\_WITH\_0 \mid "\_")+
                       ::=
                            "char"
CHAR\_DT.2
                       ::=
CHAR\_NAME.3
                            "char"
                                                  DIG\_WITH\_0 \mid "\_")+
                       ::=
                                   (LETTER
                            "void"
VOID\_DT.2
                       ::=
VOID\_NAME.3
                            "void" (LETTER
                                                 DIG\_WITH\_0
                       ::=
prim_{-}dt
                       ::=
                            INT\_DT
                                         CHAR\_DT
                                                        VOID\_DT
```

Grammatik 2.1.1: Konkrette Grammatik der Sprache  $L_{PicoC}$  für die Lexikalische Analyse in EBNF

#### 2.1.2 Codebeispiel

In den folgenden Unterkapiteln wird das Beispiel in Code 2.1 dazu verwendet die Konstruktion eines Abstrakten Syntaxbaumes in seinen einzelnen Zwischenschritten zu erläutern.

```
1 struct st {int *(*attr)[4][5];};
2
3 void main() {
4   struct st *(*var[3][2]);
5 }
```

Code 2.1: PicoC-Code des Codebeispiels

Die vom Basic Lexer des Lark Parsing Toolkit erkannten Token sind Code 2.2 zu sehen.

Code 2.2: Tokens für das Codebeispiel

#### 2.2 Syntaktische Analyse

In der Syntaktischen Analyse ist es die Aufgabe des Parsers aus einem Programm in Konkretter Syntax unter Verwendung der Tokens aus der Lexikalischen Analyse einen Ableitungsbaum zu generieren. Es ist danach die Aufgabe möglicher Visitors und die Aufgabe des Transformers aus diesem Ableitungsbaum einen Abstrakten Syntaxbaum in Abstrakter Syntax zu generieren.

#### 2.2.1 Umsetzung von Präzidenz und Assoziativität

In diesem Unterkapitel wird eine ähnliche Erklärweise, wie in Parsing Expressions · Crafting Interpreters verwendet. Die Programmiersprache  $L_{PicoC}$  hat dieselben Präzidenzregeln implementiert, wie die Programmiersprache  $L_{C}^{7}$ . Die Präzidenzregeln der verschiedenen Operatoren der Programmiersprache  $L_{PicoC}$  sind in Tabelle 2.1 aufgelistet.

<sup>&</sup>lt;sup>7</sup>C Operator Precedence - cppreference.com.

Präzidenzst	ufe Operatoren	Beschreibung	Assoziativität
1	a()	Funktionsaufruf	
	a[]	Indexzugriff	Links, dann rechts $\rightarrow$
	a.b	Attributzugriff	
2	-a	Unäres Minus	
	!a ~a	Logisches NOT und Bitweise NOT	Rechts, dann links $\leftarrow$
	*a &a	Dereferenz und Referenz, auch	Rechts, dami miks ←
		Adresse-von	
3	a*b a/b a%b	Multiplikation, Division und Modulo	
4	a+b a-b	Addition und Subtraktion	
5	a <b a="" a<="b">b a&gt;=b</b>	Kleiner, Kleiner Gleich, Größer,	
		Größer gleich	
6	a==b a!=b	Gleichheit und Ungleichheit	Links dann rechts
7	a&b	Bitweise UND	Links, dann rechts $\rightarrow$
8	a^b	Bitweise XOR (exclusive or)	
9	a b	Bitweise ODER (inclusive or)	
10	a&&b	Logiches UND	
11	a  b	Logisches ODER	
12	a=b	Zuweisung	Rechts, dann links $\leftarrow$

Tabelle 2.1: Präzidenzregeln von PicoC

Würde man diese Operatoren ohne Beachtung von Präzidenzreglen (Definition ??) und Assoziativität (Definition ??) in eine Konkrette Grammatik verarbeiten wollen, so könnte eine Konkrette Grammatik  $G = \langle N, \Sigma, P, exp \rangle$  mit Produktionen 2.2.1 P dabei rauskommen.

```
CHAR
                                                      "("exp")"
                            NUM
                                                                                            L_Arith
prim_{-}exp
                  exp"["exp"]"
                                 | exp"."name | name"("fun\_args")"
                                                                                            +L_{-}Logic
                  [exp(","exp)*]
"-" | " ~ "
fun\_args
            ::=
                                                                                            + L_-Pntr
                                    | "!" | "*" | "&:"
                                                                                            + L_Array
un\_op
                                                                                            + L_Struct
un_{-}exp
                  un\_op \ exp
                  "*" | "/" | "%" | "+" | "-" | "&" | "<" | "<" | ">=" | "!=" | "=="
                                                                                            + L_{-}Fun
bin\_op
                  "&&" | "||" | " = "
bin_exp
                  exp\ bin\_op\ exp
                 prim_{exp} \mid un_{exp} \mid bin_{exp}
exp
```

Grammatik 2.2.1: Undurchdachte Konkrette Grammatik der Sprache  $L_{PicoC}$  für die Syntaktische Analyse in EBNF, die Operatorpräzidenz nicht beachtet

Die Konkrette Grammatik 2.2.1 ist allerdings mehrdeutig, d.h. verschiedene Linksableitungen in der Konkretten Grammatik können zum selben Wort abgeleitet werden. Z.B. kann das Wort 3 \* 1 && 4 sowohl über die Linksableitung 2.2.1 als auch über die Linksableitung 2.2.2 abgeleitet werden.

exp 
$$\Rightarrow$$
 bin\_exp  $\Rightarrow$  exp bin\_op exp  $\Rightarrow$  bin\_exp bin\_op exp  $\Rightarrow$  exp bin\_op exp bin\_op exp  $\Rightarrow$  3 \* 1 && 4

```
\begin{array}{l} \exp \Rightarrow \operatorname{bin\_exp} \Rightarrow \exp \ \operatorname{bin\_op} \ \exp \ \Rightarrow \operatorname{prim\_exp} \ \operatorname{bin\_op} \ \exp \Rightarrow \operatorname{NUM} \ \operatorname{bin\_op} \ \exp \\ \Rightarrow 3 \ \operatorname{bin\_op} \ \exp \Rightarrow 3 \ * \ \operatorname{exp} \Rightarrow 3 \ * \ \operatorname{exp} \ \Rightarrow 3 \ * \ 1 \ \&\& \ 4 \end{array}
```

Beide Wörter sind gleich, allerdings sind die Ableitungsbäume unterschiedlich, wie in Abbildung 2.1 zu sehen ist.



Abbildung 2.1: Ableitungsbäume zu den beiden Ableitungen

Der linke Baum entspricht Ableitung 2.2.1 und der rechte Baum entspricht Ableitung 2.2.2. Würde man in den Ausdrücken, die von diesen Bäumen darsgestellt sind in Klammern setzen, um die Präzidenz sichtbar zu machen, so würde Ableitung 2.2.1 die Klammerung (3 \* 1) & 4 haben und die Ableitung 2.2.2 die Klammerung 3 \* (1 & 4) haben.

Aus diesem Grund ist es wichtig die Präzidenzregeln und die Assoziativität der Operatoren beim Erstellen der Konkretten Grammatik miteinzubeziehen. Hierzu wird nun Tabelle 2.1 betrachtet. Für jede Präzidenzstufe in der Tabelle 2.1 wird eine eigene Regel erstellt werden, wie es in Grammatik 2.2.2 dargestellt ist. Zudem braucht es eine Produktion prim\_exp für die höchste Präzidenzstufe, welche Literale, wie 'c', 5 oder var und geklammerte Ausdrücke wie (3 & 14) abdeckt.

$prim\_exp$	::=	 $L\_Arith + L\_Array$
$post\_exp$	::=	 $+$ $L_Pntr + L_Struct$
$un\_exp$	::=	 $+ L_{-}Fun$
$arith\_prec1$	::=	
$arith\_prec2$	::=	
$arith\_and$	::=	
$arith\_oplus$	::=	
$arith\_or$	::=	
$rel\_exp$	::=	 $L\_Logic$
$eq\_exp$	::=	
$logic\_and$	::=	
$logic\_or$	::=	
$assign\_stmt$	::=	 $L\_Assign$

Grammatik 2.2.2: Erster Schritt zu einer durchdachten Konkretten Grammatik der Sprache  $L_{PicoC}$  für die Syntaktische Analyse in EBNF, die Operatorpräzidenz beachtet

Einigen Bezeichnungen der Produktionen sind in Tabelle 2.2 ihren jeweiligen Operatoren zugeordnet für welche sie zuständig sind.

Bezeichnung der Produktionsregel	Operatoren
post_exp	a() a[] a.b
un_exp	-a !a ~a *a &a
arith_prec1	a*b a/b a%b
arith_prec2	a+b a-b
$\mathtt{arith}_{\mathtt{-}}\mathtt{and}$	a <b a="" a<="b">b a&gt;=b</b>
arith_oplus	a==b a!=b
arith_or	a&b
rel_exp	a^b
eq_exp	a b
logic_and	a&&b
logic_or	a  b
assign	a=b

Tabelle 2.2: Zuordnung der Bezeichnungen von Produktionsregeln zu Operatoren

Als nächstes müssen die einzelnen **Produktionen** entsprechend der **Ausdrücke** für die sie zuständig sind definiert werden. Jede der **Produktionen** soll nur Ausdrücke **erkennen** können, deren **Präzidenzstufe** die ist, für welche die jeweilige Produktion verantwortlich ist oder deren Präzidenzstufe **höher** ist. Z.B. soll un\_op sowohl den Ausdruck -(3 \* 14) als auch einfach nur (3 \* 14)<sup>8</sup> erkennen können, aber nicht 3 \* 14 ohne Klammern, da dieser Ausdruck eine **geringe Präzidenz** hat. Des Weiteren muss bei Produktionen für Ausdrücke mit **Operatoren** unterschieden werden, ob die Operatoren linksassoziativ oder **rechtsassoziativ**, unär, binär usw. sind.

Bei z.B. der Produktion um\_exp in 2.2.3 für die rechtsassoziativen unären Operatoren -a, !a ~a, \*a und &a ist die Alternative um\_op um\_exp dafür zuständig, dass diese unären Operatoren rechtsassoziativ geschachtelt werden können (z.B. !-~42). Die Alternative post\_exp ist dafür zuständig, dass die Produktion auch terminieren kann und es auch möglich ist auschließlich einen Ausdruck höherer Präzidenz (z.B. 42) zu haben.

$$un\_exp ::= un\_op un\_exp \mid post\_exp$$

Grammatik 2.2.3: Beispiel für eine unäre rechtsassoziative Produktion

Bei z.B. der Produktion post\_exp in 2.2.4 für die linksassoziativen unären Operatoren a(), a[] und a.b sind die Alternativen post\_exp"["logic\_or"]" und post\_exp"."name dafür zuständig, dass diese unären Operatoren linksassoziativ geschachtelt werden können (z.B. ar[3][1].car[4]). Die Alternative name"("fun\_args")" ist für einen einzelnen Funktionsaufruf zuständig. Die Alternative prim\_exp ist dafür zuständig, dass die Produktion nicht nur bei name"("fun\_args")" terminieren kann und es auch möglich ist auschließlich einen Ausdruck der höchsten Präzidenz (z.B. 42) zu haben.

Bei z.B. der Produktion prec2\_exp in 2.2.5 für die binären linksassoziativen Operatoren a+b und a-b ist die Alternative arith\_prec2 prec2\_op arith\_prec1 dafür zuständig, dass mehrere Operationen der Präzidenzstufe 4 in Folge erkannt werden können<sup>9</sup> (z.B. 3 + 1 - 4, wobei - und + beide Präzidenzstufe 4

<sup>&</sup>lt;sup>8</sup>Geklammerte Ausdrücke werden nämlich von prim\_exp erkannt, welches eine höhere Präzidenzstufe hat.

<sup>&</sup>lt;sup>9</sup>Bezogen auf Tabelle 2.1.

haben). Das Nicht-Terminalsymbol arith\_prec1 auf der rechten Seite ermöglicht es, dass zwischen den Operationen der Präzidenzstufe 4 auch Operationen der Präzidenzstufe 3 auftauchen können (z.B. 3 + 1 / 4 - 1, wobei - und + beide Präzidenzstufe 4 haben und / Präzidenzstufe 3). Mit der Alternative arith\_prec1 ist es möglich, dass auschließlich ein Ausdruck höherer Präzidenz erkannt wird (z.B. 1 / 4).

 $arith\_prec2$  ::=  $arith\_prec2$   $prec2\_op$   $arith\_prec1$  |  $arith\_prec1$ 

Grammatik 2.2.5: Beispiel für eine binäre linksassoziative Produktion

#### Anmerkung Q

Manche Parser<sup>a</sup> haben allerdings ein Problem mit Linksrekursion (Definition ??), wie sie z.B. in der Produktion 2.2.5 vorliegt. Dieses Problem lässt sich allerdings einfach lösen, indem man die Produktion 2.2.5 zur Produktion 2.2.6 umschreibt.

 $arith\_prec2$  ::=  $arith\_prec1$  ( $prec2\_op$   $arith\_prec1$ )\*

Grammatik 2.2.6: Beispiel für eine binäre linksassoziative Produktion ohne Linksrekursion

Die von Produktion 2.2.6 erkannten Ausdrücke sind dieselben, wie für die Produktion 2.2.5, allerdings ist die Produktion 2.2.6 flach gehalten und ruft sich nicht selber auf, sondern nutzt den in der EBNF (Definition 2.3) definierten \*-Operator, um mehrere Operationen der Präzidenzstufe 4 in Folge erkennen zu können (z.B. 3 + 1 - 4, wobei - und + beide Präzidenzstufe 4 haben).

Das Nicht-Terminalsymbol arith\_prec1 erlaubt es, dass zwischen der Folge von Operationen der Präzidenzstufe 4 auch Operationen der Präzidenzstufe 3 auftauchen können (z.B. 3 + 1 / 4 - 1, wobei - und + beide Präzidenzstufe 4 haben und / Präzidenzstufe 3). Da der in der EBNF definierte \*-Operator auch bedeutet, dass das Teilpattern auf das er sich bezieht kein einziges mal vorkommen kann, ist es mit dem linken Nicht-Terminalsymbol arith\_prec1 möglich, dass auschließlich ein Ausdruck höherer Präzidenz erkannt wird (z.B. 1 / 4).

<sup>a</sup>Darunter zählt der Earley Parser, der im PicoC-Compiler verwendet wird nicht.

Alle Operatoren der Sprache  $L_{PicoC}$  sind also entweder binär und linksassoziativ (z.B. a\*b, a-b, a>=b oder a&&b), unär und rechtsassoziativ (z.B. &a oder !a) oder unär und linksassoziativ (z.B. a[] oder a()). Somit ergibt sich die Konkrette Grammatik 2.2.7.

prec1_op prec2_op rel_op eq_op fun_args	::=	" * "   "/"   "%" " + "   " - " " < "   " <= "   " > "   " >= " " == "   "! = " [logic_or("," logic_or)*]	$L\_Misc$
$\begin{array}{c} \hline prim\_exp \\ post\_exp \end{array}$	::=	$name \mid NUM \mid CHAR \mid$ "("logic_or")" $post\_exp$ "["logic_or"]" $\mid post\_exp$ "." $name \mid name$ "("fun_args")" $prim\_exp$	$L_Arith$ + $L_Array$ + $L_Pntr$
un_exp arith_prec1 arith_prec2 arith_and arith_oplus arith_or	::= ::= ::= ::=	un_op_un_exp   post_exp arith_prec1 prec1_op_un_exp   un_exp arith_prec2 prec2_op_arith_prec1   arith_prec1 arith_and "&" arith_prec2   arith_prec2 arith_oplus "\\" arith_and   arith_and arith_or " " arith_oplus   arith_oplus	+ L_Struct + L_Fun
rel_exp eq_exp logic_and logic_or	::= ::= ::=	rel_exp rel_op arith_or   arith_or eq_exp eq_op rel_exp   rel_exp logic_and "&&" eq_exp   eq_exp logic_or "  " logic_and   logic_and	$L\_Logic$
$assign\_stmt$	::=	un_exp "=" logic_or";"	$L_{-}Assign$

Grammatik 2.2.7: Durchdachte Konkrette Grammatik der Sprache  $L_{PicoC}$  in EBNF, die Operatorpräzidenz beachtet

#### 2.2.2 Konkrette Grammatik für die Syntaktische Analyse

Die gesamte Konkrette Grammatik 2.2.8 ergibt sich wenn man die Konkrette Grammatik 2.2.7 um die restliche Syntax der Sprache  $L_{PicoC}$  erweitert, die sich nach einem **ähnlichen Prinzip** wie in Unterkapitel 2.2.7 erläutert ergibt.

Später in der Entwicklung des PicoC-Compilers wurde die Konkrette Grammatik an die aktuellste konstenlos auffindbare Version der echten Konkretten Grammatik ANSI C grammar (Yacc) der Sprache  $L_C$  angepasst<sup>10</sup>, damit es sicherer gewährleistet werden kann, dass der PicoC-Compiler sich genauso verhält, wie geläufige Compiler der Programmiersprache  $L_C$ . Wobei z.B. die Compiler GCC<sup>11</sup> und Clang<sup>12</sup> zu nennen wären.

In der Konkretten Grammatik 2.2.8 für die Syntaktischen Analyse werden einige der Tokennamen aus der Konkretten Grammatik 2.1.1 für die Lexikalischen Analyse verwendet, wie z.B. NUM aber auch name, welches eine Produktion ist, die mehrere Tokennamen unter einem Überbegriff zusammenfasst.

Terminalsymbole, wie ; oder && gehören eigentlich zur Lexikalischen Analyse, jedoch erlaubt das Lark Parsing Toolkit um die Konkrette Grammatik leichter lesbar zu machen einige Terminalsymbole einfach direkt in die Konkrette Grammatik 2.2.8 für die Syntaktische Analyse zu schreiben. Der Tokenname für diese Terminalsymbole wird in diesem Fall vom Lark Parsing Toolkit bestimmt, welches einige sehr häufige verwendete Terminalsymbole, wie; oder && bereits einen eigenen Tokennamen zugewiesen hat.

 $<sup>^{10}</sup>$ An der für die Programmiersprache  $L_{PicoC}$  relevanten Syntax hat sich allerdings über die Jahre nichts verändert, wie die Konkretten Grammatiken für die Syntaktische Analyse  $ANSI\ C\ grammar\ (Lex)$  und Lexikalische Analyse noauthor ansi nodate-2 aus dem Jahre 1985 zeigen.

<sup>&</sup>lt;sup>11</sup>GCC, the GNU Compiler Collection - GNU Project.

 $<sup>^{12}</sup>$  clang: C++ Compiler.

prim_exp post_exp un_exp	::= ::=   ::=	name   NUM   CHAR   "("logic_or")"  array_subscr   struct_attr   fun_call input_exp   print_exp   prim_exp un_op_un_exp   post_exp	$L\_Arith + L\_Array$ + $L\_Pntr + L\_Struct$ + $L\_Fun$
input_exp print_exp arith_prec1 arith_prec2 arith_and arith_oplus arith_or	::= ::= ::= ::= ::=	"input""("")"  "print""("logic_or")"  arith_prec1 prec1_op un_exp   un_exp  arith_prec2 prec2_op arith_prec1   arith_prec1  arith_and "&" arith_prec2   arith_prec2  arith_oplus "\\" arith_and   arith_and  arith_or " " arith_oplus   arith_oplus	
rel_exp eq_exp logic_and logic_or	::= ::= ::=	rel_exp rel_op arith_or   arith_or eq_exp eq_op rel_exp   rel_exp logic_and "&&" eq_exp   eq_exp logic_or "  " logic_and   logic_and	$L\_Logic$
type_spec alloc assign_stmt initializer init_stmt const_init_stmt	::= ::= ::= ::= ::=	<pre>prim_dt   struct_spec type_spec pntr_decl un_exp "=" logic_or";" logic_or   array_init   struct_init alloc "=" initializer";" "const" type_spec name "=" NUM";"</pre>	$L\_Assign\_Alloc$
$pntr\_deg \\ pntr\_decl$	::=	"*"*  pntr_deg array_decl   array_decl	$L_{-}Pntr$
array_dims array_decl array_init array_subscr	::= ::= ::=	("["NUM"]")*  name array_dims   "("pntr_decl")"array_dims  "{"initializer("," initializer) * "}"  post_exp"["logic_or"]"	$L\_Array$
struct_spec struct_params struct_decl struct_init struct_attr	::= ::= ::=	"struct" name (alloc";")+  "struct" name "{"struct_params"}"  "{""."name"="initializer  ("," "."name"="initializer)*"}"  post_exp"."name	$L_{-}Struct$
$if\_stmt$ $if\_else\_stmt$	::=	"if""("logic_or")" exec_part "if""("logic_or")" exec_part "else" exec_part	$L\_If\_Else$
while_stmt do_while_stmt	::=	"while""("logic_or")" exec_part "do" exec_part "while""("logic_or")"";"	$L_{-}Loop$

Grammatik 2.2.8: Konkrette Grammatik der Sprache  $L_{PicoC}$  für die Syntaktische Analyse in EBNF, Teil 1

```
decl\_exp\_stmt
                          alloc";"
                                                                                                   L\_Stmt
                    ::=
decl\_direct\_stmt
                          assign_stmt | init_stmt | const_init_stmt
                    ::=
decl\_part
                          decl\_exp\_stmt \mid decl\_direct\_stmt \mid RETI\_COMMENT
                    ::=
                          "{"exec\_part *"}"
compound\_stmt
                    ::=
                          logic\_or";"
exec\_exp\_stmt
                    ::=
exec\_direct\_stmt
                          if\_stmt \mid if\_else\_stmt \mid while\_stmt \mid do\_while\_stmt
                    ::=
                          assign\_stmt \mid fun\_return\_stmt
                          compound\_stmt \mid exec\_exp\_stmt \mid exec\_direct\_stmt
exec\_part
                    ::=
                          RETI\_COMMENT
                          decl\_part * exec\_part *
decl\_exec\_stmts
                    ::=
                                                                                                   L_{-}Fun
fun\_args
                          [logic\_or("," logic\_or)*]
                    ::=
                          name"("fun\_args")"
fun\_call
                    ::=
fun\_return\_stmt
                          "return" [logic_or]";"
                    ::=
                          [alloc("," alloc)*]
fun\_params
                    ::=
fun\_decl
                          type_spec pntr_deg name" ("fun_params")"
                    ::=
                          type_spec_pntr_deg_name"("fun_params")" "{"decl_exec_stmts"}"
fun_{-}def
                    ::=
                          (struct\_decl
                                           fun\_decl)";"
decl\_def
                                                              fun_{-}def
                                                                                                   L_File
                    ::=
                          decl\_def*
decls\_defs
                    ::=
file
                    ::=
                          FILENAME decls_defs
```

Grammatik 2.2.9: Konkrette Grammatik der Sprache  $L_{PicoC}$  für die Syntaktische Analyse in EBNF, Teil 2

# Anmerkung Q

In der Konkretten Grammatik 2.2.8 sind alle Grammatiksymbole ausgegraut, die das Bachelorprojekt betreffen. Alle nicht ausgegrauten Grammatiksymbole wurden für die Implementierung der neuen Funktionalitäten, welche die Bachelorarbeit betreffen hinzugefügt.

# 2.2.3 Ableitungsbaum Generierung

Die in Unterkapitel 2.2.2 definierte Konkrette Grammatik 2.2.8 lässt sich mithilfe des Earley Parsers (Definition 2.6) von Lark dazu verwenden Code, der in der Sprache  $L_{PicoC}$  geschrieben ist zu parsen um einen Ableitungsbaum zu generieren.

#### Definition 2.6: Earley Parser

Ist ein Algorithmus für das Parsen von Wörtern einer Kontextfreien Sprache, der ein Chart Parser ist, welcher einen mittels Dynamischer Programmierung und dem Top-Down Ansatz arbeitenden Earley Recognizer (Defintion 2.7) nutzt, um einen Ableitungsbaum zu konstruieren.

Zur Konstruktion des Ableitungsbaumes muss dafür gesorgt werden, dass der Earley Recognizer bei der Vervollständigungsoperation Zeiger auf den vorherigen Zustand hinzugefügt, um durch Rückwärtsverfolgen dieser Zeiger die Ableitung wieder nachvollziehen zu können und so einen Ableitungsbaum konstruieren zu können.<sup>a</sup>

<sup>&</sup>lt;sup>a</sup>Earley, "An efficient context-free parsing".

## Definition 2.7: Earley Recognizer

Z

Ist ein Recognizer, der für alle Kontextfreien Sprachen das Wortproblem entscheiden kann und dies mittels Dynamischer Programmierung mit dem Top-Down Ansatz umsetzt.<sup>a</sup>

Eingabe und Ausgabe des Algorithmus sind:

- Eingabe: Eingabewort w und Konkrette Grammatik  $G_{Parse} = \langle N, \Sigma, P, S \rangle$
- Ausgabe: 0 wenn  $w \notin L(G_{Parse})^b$  und 1 wenn  $w \in L(G_{Parse})$

Bevor dieser Algorithmus erklärt wird müssen noch einige Symbole und Notationen erklärt werden:

- $\alpha$ ,  $\beta$ ,  $\gamma$  stellen eine beliebige Folge von Grammatiksymbolen<sup>c</sup> dar
- A und B stellen Nicht-Terminalsymbole dar
- a stellt ein Terminalsymbol dar
- Earley's Punktnotation:  $A := \alpha \bullet \beta$  stellt eine Produktion, in der  $\alpha$  bereits geparst wurde und  $\beta$  noch geparst werden muss
- Die Indexierung ist informell ausgedrückt so umgesetzt, dass die Indices zwischen Tokennamen liegen, also Index 0 vor dem ersten Tokennamen verortet ist, Index 1 nach dem ersten Tokennamen verortet ist und Index n nach dem letzten Tokennamen verortet ist

und davor müssen noch einige Begriffe definiert werden:

- Zustandsmenge: Für jeden der n+1 Indices j wird eine Zustandsmenge Z(j) generiert
- Zustand einer Zustandsmenge: Ist ein Tupel  $(A := \alpha \bullet \beta, i)$ , wobei  $A := \alpha \bullet \beta$  die aktuelle Produktion ist, die bis Punkt  $\bullet$  geparst wurde und i der Index ist, ab welchem der Versuch der Erkennung eines Teilworts des Eingabeworts mithilfe dieser Produktion begann

Der Ablauf des Algorithmus ist wie folgt:

- 1. initialisiere Z(0) mit der Produktion, welches das Startsymbol S auf der linken Seite des ::=-Symbols hat
- 2. es werden in der aktuellen Zustandsmenge Z(j) die folgenden Operationen ausgeführt:
  - Voraussage: Für jeden Zustand in der Zustandsmenge Z(j), der die Form  $(A ::= \alpha \bullet B\gamma, i)$  hat, wird für jede Produktion  $(B ::= \beta)$  in der Konkretten Grammatik, die ein B auf der linken Seite des ::=-Symbols hat ein Zustand  $(B ::= \bullet \beta, j)$  zur Zustandsmenge Z(j) hinzugefügt
  - Überprüfung: Für jeden Zustand in der Zustandsmenge Z(j), der die Form  $(A ::= \alpha \bullet a\gamma, i)$  hat wird der Zustand  $(A ::= \alpha a \bullet \gamma, i)$  zur Zustandsmenge Z(j+1) hinzugefügt
  - Vervollständigung: Für jeden Zustand in der Zustandsmenge Z(j), der die Form  $(B := \beta \bullet, i)$  hat werden alle Zustände in Z(i) gesucht, welche die Form  $(A := \alpha \bullet B\gamma, i)$  haben und es wird der Zustand  $(A := \alpha B \bullet \gamma, i)$  zur Zustandsmenge Z(j) hinzugefügt

bis:

• der Zustand  $(A := \beta \bullet, 0)$  in der Zustandsmenge Z(n) auftaucht, wobei A das Startsym-

```
bol\ S\ ist \Rightarrow w \in L(G_{Parse})
• keine\ Zust\"{a}nde\ mehr\ hinzugef\"{u}gt\ werden\ k\"{o}nnen} \Rightarrow w \not\in L(G_{Parse})

<sup>a</sup>Earley, "An efficient context-free parsing".

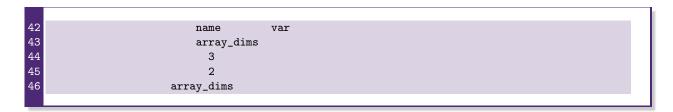
<sup>b</sup>L(G_{Parse}) ist die Sprache, welche durch die Konkrette Grammatik G_{Parse} beschrieben wird.

<sup>c</sup>Also eine Folge von Terminalsymbolen und Nicht-Terminalsymbolen.
```

## 2.2.3.1 Codebeispiel

Der Ableitungsbaum, der mithilfe des Earley Parsers und der Token der Lexikalischen Analyse aus dem Beispiel in Code 2.1 generiert wurde, ist in Code 2.3 zu sehen. Im Code 2.3 wurden einige Zeilen markiert, die später in Unterkapitel 2.2.4.1 zum Vergleich wichtig sind.

```
./verbose_dt_simple_ast_gen_array_decl_and_alloc.dt
     decls_defs
       decl_def
         struct_decl
           name
                        st
 7
8
9
           struct_params
              alloc
                type_spec
10
                  prim_dt
                                  int
                pntr_decl
12
                  pntr_deg
13
                  array_decl
14
                    pntr_decl
                      pntr_deg
16
                      array_decl
17
                        name
                                     attr
18
                        array_dims
19
                    array_dims
20
21
22
       decl_def
23
         fun_def
24
           type_spec
25
                              void
             prim_dt
26
           pntr_deg
27
           name
                        main
28
           fun_params
29
           decl_exec_stmts
30
              decl_part
31
                decl_exp_stmt
32
                  alloc
33
                    type_spec
34
                      struct_spec
35
                        name
                                     st
36
                    pntr_decl
37
                      pntr_deg
38
                      array_decl
39
                        pntr_decl
40
                          pntr_deg
                           array_decl
```



Code 2.3: Ableitungsbaum nach Ableitungsbaum Generierung

### 2.2.3.2 Ausgabe des Ableitunsgbaumes

Die Ausgabe des Ableitungsbaumes wird komplett vom Lark Parsing Toolkit übernommen. Für die Inneren Knoten werden die Nicht-Terminalsymbole, welche in der Konkretten Grammatik den linken Seiten des ::=-Symbols<sup>13</sup> entsprechen hergenommen und die Blätter sind Terminalsymbole, genauso, wie es in der Definition ?? eines Ableitungsbaumes auch schon definiert ist. Die EBNF-Grammatik 2.2.8 des PicoC-Compilers erlaubt es allerdings auch, dass in einem Blatt garnichts  $\varepsilon$  steht, weil es z.B. Produktionen, wie array\_dims ::= ("["NUM"]")\* gibt, in denen auch das leere Wort  $\varepsilon$  abgeleitet werden kann.

Die Ausgabe des Abstrakten Syntaxbaumes ist bewusst so gewählt, dass sie sich optisch vom Ableitungsbaum unterscheidet, indem die Bezeichner der Knoten in UpperCamelCase<sup>14</sup> geschrieben sind, im Gegensatz zum Ableitungsbaum, dessen Innere Knoten im snake\_case geschrieben sind, wie auch die Nicht-Terminalsymbole auf den linken Seiten des ::=-Symbols.

# 2.2.4 Ableitungsbaum Vereinfachung

Der Ableitungsbaum in Code 2.3, dessen Generierung in Unterkapitel 2.2.3.1 besprochen wurde ist noch untauglich, damit aus ihm mittels eines Tramsformers ein Abstrakter Syntaxbaum generiert werden kann. Das Problem ist, dass um den Datentyp einer Variable in der Programmiersprache  $L_C$  und somit auch die Programmiersprache  $L_{PicoC}$  korrekt bestimmen zu können, wie z.B. ein "Feld der Mächtigkeit 3 von Zeigern auf Felder der Mächtigkeit 2 von Integern" int (\*ar[3])[2] die Spiralregel<sup>15</sup> in der Implementeirung des PicoC-Compilers umgesetzt werden muss und das ist nicht alleinig möglich, indem man die entsprechenden Produktionen in der Konkretten Grammatik 2.2.8 der Konkretten Syntax auf eine spezielle Weise passend spezifiziert.

Was man erhalten will, ist ein **entarteter Baum** von **PicoC-Knoten**, an dem man den **Datentyp** direkt ablesen kann, indem man sich einfach über den **entarteten Baum** bewegt, wie z.B. PntrDecl(Num('1'), A rrayDecl([Num('3'),Num('2')],PntrDecl(Num('1'),StructSpec(Name('st'))))) für den Ausdruck struct st \*(\*var[3][2]).

Es sind hierbei mehrere Probleme zu lösen. Hat man den Ausdruck struct st \*(\*var[3][2]) wird dieser zu einem Ableitungsbaum, wie er in Abbildung 2.2 zu sehen ist.

<sup>&</sup>lt;sup>13</sup> Grammar: The language of languages (BNF, EBNF, ABNF and more).

<sup>&</sup>lt;sup>14</sup>Naming convention (programming).

<sup>&</sup>lt;sup>15</sup> Clockwise/Spiral Rule.



Abbildung 2.2: Ableitungsbaum nach Parsen eines Ausdrucks

Dieser Ableitungsbaum für den Ausdruck struct st \*(\*var[3][2]) hat allerdings einen Aufbau welcher durch die Syntax der Zeigerdeklaratoren pntr\_decl(num, datatype) und Felddeklaratoren array\_decl(datatype, nums) bestimmt ist, die spiralähnlich ist. Man würde allerdings gerne einen entarteten Baum erhalten, bei dem der Datentyp immer im zweiten Attribut weitergeht, anstatt abwechselnd im zweiten und ersten, wie beim Zeigerdeklarator pntr\_decl(num, datatype) und Felddeklarator array\_decl(datatype, nums). Daher muss beim FeldDeclarator array\_decl(datatype, nums) immer das erste Attribut datatype mit dem zweiten Attribut nums getauscht werden.

Des Weiteren befindet sich in der Mitte dieser Spirale, die der Ableitungsbaum bildet der Name der Variable name(var) und nicht der innerste Datentyp struct st, da der Ableitungsbaum einfach nur die kompilerinterne Darstellung, die durch das Parsen eines Programms in Konkretter Syntax (z.B. struct st \*(\*var[3][2])) generiert wird darstellt. Der Name der Variable name(var) sollte daher mit dem innersten Datentyp struct st ausgetauscht werden.

In Abbildung 2.3 ist daher zu sehen, wie der **Ableitungsbaum** aus Abbildung 2.2 mithilfe eines **Visitors** (Definition **??**) **vereinfacht** wird, sodass er die gerade erläuterten Ansprüche erfüllt.

Die Implementierung des Visitors aus dem Lark Parsing Toolkit ist unter Link<sup>16</sup> zu finden ist. Diese Implementierung ist allerdings zu spezifisch auf Lark zugeschnitten, um sie an dieser Stelle allgemein erklären zu können. Der Visitor verhält sich allerdings grundlegend so, wie es in Definition ?? erklärt wurde.

 $<sup>^{16}</sup>$ https://github.com/lark-parser/lark/blob/d03f32be7f418dc21cfa45acc458e67fe0580f60/lark/visitors.py.

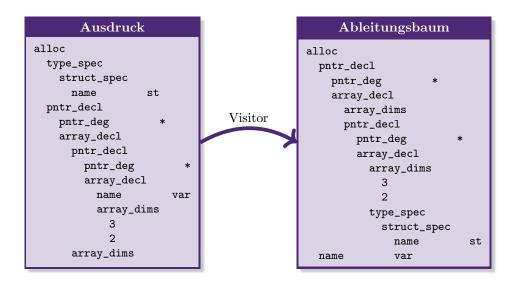


Abbildung 2.3: Ableitungsbaum nach Vereinfachung

## 2.2.4.1 Codebeispiel

In Code 2.4 ist der Ableitungsbaum aus Code 2.3 nach der Vereinfachung mithilfe eines Visitors zu sehen.

```
file
     ./verbose_dt_simple_ast_gen_array_decl_and_alloc.dt_simple
     decls_defs
 4
5
       decl_def
         struct_decl
           name
                        st
 7
8
9
           struct_params
             alloc
               pntr_decl
10
                  pntr_deg
                  array_decl
                    array_dims
                      4
14
                      5
                    pntr_decl
16
                      pntr_deg
17
                      array_decl
18
                        array_dims
19
                        type_spec
20
                          prim_dt
                                           int
               name
                             attr
22
       decl_def
23
         fun_def
24
           type_spec
25
             prim_dt
                              void
26
           pntr_deg
27
           name
                        main
28
           fun_params
           decl_exec_stmts
```

```
decl_part
31
                decl_exp_stmt
32
                   alloc
33
                     pntr_decl
34
                       pntr_deg
35
                       array_decl
36
                          array_dims
37
                         pntr_decl
38
                            pntr_deg
39
                            array_decl
40
                              array_dims
41
                                3
42
                                2
43
                              type_spec
44
                                 struct_spec
45
                                                 st
                                   name
46
                     name
                                   var
```

Code 2.4: Ableitungsbaum nach Ableitungsbaum Vereinfachung

# 2.2.5 Generierung des Abstrakten Syntaxbaumes

Nachdem der Ableitungsbaum in Unterkapitel 2.2.4 vereinfacht wurde, ist der vereinfachte Ableitungsbaum in Code 2.4 nun dazu geeignet, um mit einem Transformer (Definition ??) einen Abstrakten Syntaxbaum aus ihm zu generieren. Würde man den vereinfachten Ableitungsbaum des Ausdrucks struct st \*(\*var[3][2]) auf passende Weise in einen Abstrakten Syntaxbaum umwandeln, so würde dabei ein Abstrakter Syntaxbaum wie in Abbildung 2.4 rauskommen.

Die Implementierung des **Transformers** aus dem **Lark Parsing Toolkit** ist unter Link<sup>17</sup> zu finden ist. Diese Implementierung ist allerdings **zu spezifisch** auf Lark zugeschnitten, um sie an dieser Stelle allgemein erklären zu können. Der **Transformer** verhält sich allerdings grundlegend so, wie es in Definition **??** erklärt wurde.

Den Teilbaum, der den Datentyp darstellt würde man von von oben-nach-unten<sup>18</sup> als "Zeiger auf einen Zeiger auf ein Feld der Mächtigkeit 2 von Feldern der Mächtigkeit 3 von Verbunden des Typs st" lesen, also genau anders herum, als man den Ausdruck struct st \*(\*var[3][2]) mit der Spiralregel lesen würde. Bei der Spiralregel fängt man beim Ausdruck struct st \*(\*var[3][2]) bei der Variable var an und arbeitet sich dann auf "Spiralbahnen", von innen-nach-außen durch den Ausdruck, um herauszufinden, dass dieser Datentyp ein "Feld der Mächtigkeit 3 von Feldern der Mächtigkeit 2 von Zeigern auf einen Zeiger auf einen Verbund vom Typ st" ist.

<sup>&</sup>lt;sup>17</sup>https://github.com/lark-parser/lark/blob/d03f32be7f418dc21cfa45acc458e67fe0580f60/lark/visitors.py.

<sup>&</sup>lt;sup>18</sup>In der Informatik wachsen Bäume von oben-nach-unten, von der Wurzel zur den Blättern, bzw. in diesem Beispiel von links-nach-rechts.



Abbildung 2.4: Generierung eines Abstrakten Syntaxbaumes ohne Umdrehen

Dieser Abstrakte Syntaxbaum ist für die Weiterverarbeitung ungeeignet, denn für die Adressberechnung für eine Aneinandereihung von Zugriffen auf Zeigerelemente, Feldelemente oder Verbundattribute, welche in Unterkapitel ?? genauer erläutert wird, will man den Datentyp in umgekehrter Reihenfolge. Aus diesem Grund muss der Transformer bei der Konstruktion des Abstrakten Syntaxbaumes zusätzlich dafür sorgen, dass jeder Teilbaum, der für einen Datentyp steht umgedreht wird. Auf diese Weise kommt ein Abstrakter Syntaxbaum mit richtig rum gedrehtem Datentyp, wie in Abbildung 2.5 zustande, der für die Weiterverarbeitung geeignet ist.

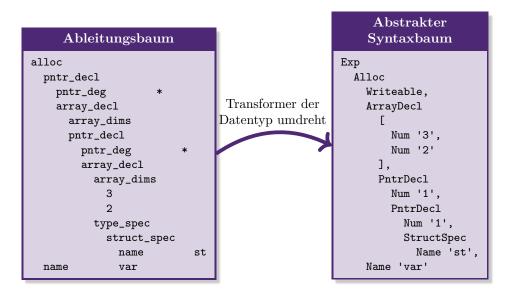


Abbildung 2.5: Generierung eines Abstrakten Syntaxbaumes mit Umdrehen

Die Weiterverarbeitung des Abstrakten Syntaxbaumes geschieht mithilfe von Passes, welche im Unterkapitel?? genauer beschrieben werden. Da die Knoten des Abstrakten Syntaxbaumes anders als beim Ableitungsbaum nicht die gleichen Bezeichnungen haben wie Produktionen der Konkretten Grammatik

ist es in den folgenden Unterkapiteln 2.2.5.1, 2.2.5.2 und 2.2.5.3 notwendig die Bedeutung der einzelnen PicoC-Knoten, RETI-Knoten und bestimmter Kompositionen dieser Knoten zu dokumentieren, die alle in den unterschiedlichen von den Passes umgeformten Abstrakten Syntaxbäumen vorkommen.

Des Weiteren gibt die Abstrakte Grammatik 2.2.10 in Unterkapitel 2.2.5.4 Aufschluss darüber welche Kompositionen von PicoC-Knoten, neben den bereits in Tabelle 2.2.10 definierten Kompositionen mit Bedeutung insgesamt überhaupt möglich sind.

#### 2.2.5.1 PicoC-Knoten

Bei den PicoC-Knoten handelt es sich um Knoten, die irgendeinen Ausdruck aus der Sprache  $L_{PicoC}$  darstellen. Für die PicoC-Knoten wurden möglichst kurze und leicht verständliche Bezeichner gewählt, da auf diese Weise bei der Implementierung der einzelnen Passes möglichst viel Code in eine Zeile passt und dieser Code auch durch leicht verständliche Bezeichner von Knoten intuitiv verständlich sein sollte<sup>19</sup>. Alle PicoC-Knoten, die in den von den verschiedenen Passes generierten Abstrakten Syntaxbäumen vorkommen sind in Tabelle 2.3 mit einem Beschreibungstext dokumentiert.

<sup>&</sup>lt;sup>19</sup>Z.B. steht der PicoC-Knoten Name(str) für einen Bezeichner. Anstatt diesen Knoten in englisch Identifier(str) zu nennen, wurde dieser als Name(str) gewählt, da Name(str) kürzer ist und inuitiver verständlich.

PiocC-Knoten	Beschreibung
Name(val)	Ein Bezeichner, z.B. my_fun, my_var usw., aber da es keine gute Kurzform für Identifier() (englisches Wort für Bezeichner) gibt, wurde dieser Knoten Name() genannt.
Num(val)	Eine <b>Z</b> ahl, z.B. 42, -3 usw.
Char(val)	Ein Zeichen der ASCII-Zeichenkodierung, z.B. 'c', '*' usw.
<pre>Minus(), Not(), DerefOp(), RefOp(), LogicNot()</pre>	Die unären Operatoren un_op: -a, ~a, *a, &a !a.
<pre>Add(), Sub(), Mul(), Div(), Mod(), Oplus(), And(), Or(), LogicAnd(), LogicOr()</pre>	Die binären Operatoren bin_op: a + b, a - b, a * b, a / b, a % b, a % b, a % b, a   b.
Eq(), NEq(), Lt(), LtE(), Gt(), GtE()	Die Relationen rel: a == b, a != b, a < b, a <= b, a > b, a >= b.
<pre>Const(), Writeable()</pre>	Die Type Qualifier type_qual: const, was für ein nicht beschreibbare Konstante steht und das nicht Angeben von const, was für einen beschreibbare Variable steht.
<pre>IntType(), CharType(), VoidType()</pre>	Die Type Specifier für Primitiven Datentypen, die in der Abstrakten Syntax, um eine intuitive Bezeichnung zu haben einfach nur als Datentypen datatype eingeordnet werden: int, char, void.
Placeholder()	Platzhalter für einen Knoten, der diesen später ersetzt.
<pre>BinOp(exp, bin_op, exp)</pre>	Container für eine binäre Operation mit 2 Expressions: <exp1> <bin_op> <exp2></exp2></bin_op></exp1>
UnOp(un_op, exp)	Container für eine <b>unäre Operation</b> mit einer Expression: <un_op> <exp>.</exp></un_op>
Exit(num)	Container für einen Exit Code, der vor der Beendigung in das ACC Register geschrieben wird und steht für die Beendigung des laufenden Programmes.
Atom(exp, rel, exp)	Container für eine binäre Relation mit 2 Expressions: <exp1> <re1> <exp2></exp2></re1></exp1>
ToBool(exp)	Container für einen Arithmetischen Ausdruck, wie z.B. 1 + 3 oder einfach nur 3, der nicht nur 1 oder 0 als Ergebnis haben kann und daher bei einem Ergebnis $x>1$ auf 1 abgebildet wird.
Alloc(type_qual, datatype, name, local_var_or_param)	Container für eine Allokation <type_qual> <datatype> <name> mit den notwendigen Knoten type_qual, datatype und name, die alle für einen Eintrag in der Symboltabelle notwendigen Informationen enthalten. Zudem besitzt er ein verstecktes Attribut local_var_or_param, dass die Information trägt, ob es sich bei der Variable um eine Lokale Variable oder einen Parameter handelt.</name></datatype></type_qual>
Assign(lhs, exp)	Container für eine <b>Zuweisung</b> , wobei 1hs ein Subscr(exp1, exp2), Deref(exp1, exp2), Attr(exp, name) oder Name('var') sein kann und exp ein beliebiger <b>Logischer Ausdruck</b> sein kann: 1hs = exp.

Tabelle 2.3: PicoC-Knoten Teil 1

PiocC-Knoten	Beschreibung
<pre>Exp(exp, datatype, error_data)</pre>	Container für einen beliebigen Ausdruck, dessen Ergebnis auf den Stack soll. Zudem besitzt er 2 versteckte Attribute, wobei datatype im RETI Blocks Pass wichtig ist und error_data für Fehlermeldungen wichtig ist.
Stack(num)	Container, der für das temporäre Ergebnis einer Berechnung, das num Speicherzellen relativ zum Stackpointer Register SP steht.
Stackframe(num)	Container, der für eine Variable steht, die num Speicherzellen relativ zum Begin-Aktive-Funktion Register BAF steht.
Global(num)	Container, der für eine Variable steht, die num Speicherzellen relativ zum Datensegment Register DS steht.
StackMalloc(num)	Container, der für das Allokieren von num Speicherzellen auf dem Stack steht.
PntrDecl(num, datatype)	Container, der für den Zeigerdatentyp steht: <pre><pre><pre><pre>*<var></var></pre>, wobei das Attribut num die Anzahl zusammen- gefasster Zeiger angibt und datatype der Datentyp ist, auf den der oder die Zeiger zeigen.</pre></pre></pre>
Ref(exp, datatype, error_data)	Container, der für die Anwendung des Referenz-Operators & <var> steht und die Adresse einer Location (Definition ??) auf den Stack schreiben soll, die über exp eingegrenzt wird. Zudem besitzt er 2 versteckte Attribute, wobei datatype im RETI Blocks Pass wichtig ist und error_data für Fehlermeldungen wichtig ist.</var>
Deref(lhs, exp)	Container für den Indexzugriff auf einen Feld- oder Zei- gerdatentyp: <var>[<i>], wobei exp1 eine angehängte weite- re Subscr(exp1, exp2), Deref(exp1, exp2), Attr(exp, name) oder ein Name('var') sein kann und exp2 der Index ist auf den zugegriffen werden soll.</i></var>
ArrayDecl(nums, datatype)	Container, der für den Felddatentyp steht: <prim_dt> <var>[<i>], wobei das Attribut nums eine Liste von Num('x') ist, die die Dimensionen des Feld angibt und datatype der Datentyp ist, der über das Anwenden von Subscript() auf das Feld zugreifbar ist.</i></var></prim_dt>
Array(exps, datatype)	Container für den Initializer eines Feldes, dessen Einträge exps weitere Initializer für eine Feld-Dimension oder ein Initializer für ein Struct oder ein Logischer Ausdruck sein können, z.B. {{1, 2}, {3, 4}}. Des Weiteren besitzt er ein verstecktes Attribut datatype, welches für den PicoC-ANF Pass Informationen transportiert, die für Fehlermeldungen wichtig sind.
Subscr(exp1, exp2)	Container für den Indexzugriff auf einen Feld- oder Zeigerdatentyp: <var>[<i>], wobei exp1 eine angehängte weitere Subscr(exp1, exp2), Deref(exp1, exp2) oder Attr(exp, name) Operation sein kann oder ein Name('var') sein kann und exp2 der Index ist auf den zugegriffen werden soll.</i></var>
StructSpec(name)	Container für einen selbst definierten Structtyp: struct <name>, wobei das Attribut name festlegt, welchen selbst definierten Structtyp dieser Knoten repräsentiert.</name>
Attr(exp, name)	Container für den Attributzugriff auf einen Structdatentyp: <var>.<attr>, wobei exp1 eine angehängte weitere Subscr(exp1, exp2), Deref(exp1, exp2) oder Attr(exp, name) Operation sein kann oder ein Name('var') sein kann und name das Attribut ist, auf das zugegriffen werden soll.</attr></var>

PiocC-Knoten	Beschreibung
Struct(assigns, datatype)	Container für den Initializer eines Structs, z.B {. <attr1>={1, 2}, .<attr2>={3, 4}}, dessen Eintrag assigns eine Liste von Assign(1hs, exp) ist mit einer Zuordnung eines Attributezeichners, zu einem weiteren Initializer für eine Feld-Dimension oder zu einem Initializer für ein Struct oder zu einem Logischen Ausdruck. Des Weiteren besitzt er ein verstecktes Attribut datatype, welches für den PicoC-ANF Pass Informationen transportiert, die für Fehlermeldungen wichtig sind.</attr2></attr1>
StructDecl(name, allocs)	Container für die Deklaration eines selbstdefinierten Structtyps, z.B. struct <var> {<datatype> <attr1>; <datatype> <attr2>;};, wobei name der Bezeichner des Structtyps ist und allocs eine Liste von Bezeichnern der Attribute des Structtyps mit dazugehörigem Datentyp, wofür sich der Knoten Alloc(type_qual, datatype, name) sehr gut als Container eignet.</attr2></datatype></attr1></datatype></var>
If(exp, stmts)	Container für ein If Statement if( <exp>) { <stmts> } in- klusive Condition exp und einem Branch stmts, indem eine Liste von Statements stehen kann oder ein einzelnes GoTo(Name('block.xyz')).</stmts></exp>
<pre>IfElse(exp, stmts1, stmts2)</pre>	Container für ein If-Else Statement if( <exp>) { <stmts2> } else { <stmts2> } inklusive Codition exp und 2 Branches stmts1 und stmts2, die zwei Alternativen Darstellen in denen jeweils Listen von Statements oder GoTo(Name('block.xyz'))'s stehen können.</stmts2></stmts2></exp>
While(exp, stmts)	Container für ein While-Statement while( <exp>) { <stmts> } inklusive Condition exp und einem Branch stmts, indem eine Liste von Statements stehen kann oder ein einzelnes GoTo(Name('block.xyz')).</stmts></exp>
DoWhile(exp, stmts)	Container für ein Do-While-Statement do { <stmts> } while(<exp>); inklusive Condition exp und einem Branch stmts, indem eine Liste von Statements stehen kann oder ein einzelnes GoTo(Name('block.xyz')).</exp></stmts>
Call(name, exps)	Container für einen Funktionsaufruf: fun_name(exps), wobei name der Bezeichner der Funktion ist, die aufgerufen werden soll und exps eine Liste von Argumenten ist, die an die Funktion übergeben werden soll.
Return(exp)	Container für ein Return-Statement: return <exp>, wobei das Attribut exp einen Logischen Ausdruck darstellt, dessen Ergebnis vom Return-Statement zurückgegeben wird.</exp>
FunDecl(datatype, name, allocs)	Container für eine Funktionsdeklaration, z.B. <datatype> <fun_name>(<datatype> <param1>, <datatype> <param2>), wobei datatype der Rückgabewert der Funktion ist, name der Bezeichner der Funktion ist und allocs die Parameter der Funktion sind, wobei der Knoten Alloc(type_spec, datatype, name) als Cotainer für die Parameter dient.</param2></datatype></param1></datatype></fun_name></datatype>

Tabelle 2.5: PicoC-Knoten Teil 3

PiocC-Knoten	Beschreibung
FunDef(datatype, name, allocs,	Container für eine Funktionsdefinition, z.B. <datatype></datatype>
stmts_blocks)	<pre><fun_name>(<datatype> <param/>) {<stmts>}, wobei datatype der Rückgabewert der Funktion ist, name der Bezeichner der Funktion ist, allocs die Parameter der Funktion sind, wobei der Knoten Alloc(type_spec, datatype, name) als Con- tainer für die Parameter dient und stmts_blocks eine Liste von Statemetns bzw. Blöcken ist, welche diese Funktion beinhaltet.</stmts></datatype></fun_name></pre>
NewStackframe(fun_name, goto_after_call)	Container für die Erstellung eines neuen Stackframes und Speicherung des Werts des BAF-Registers der aufrufenden Funktion und der Rücksprungadresse nacheinander an den Anfang des neuen Stackframes. Das Attribut fun name stehte dabei für den Bezeichner der Funktion, für die ein neuer Stackframe erstellt werden soll. Das Attribut fun name dient später dazu den Block dieser Funktion zu finden, weil dieser für den weiteren Kompiliervorang wichtige Information in seinen versteckte Attributen gespeichert hat. Des Weiteren ist das Attribut goto_after_call ein GoTo(Name('addr@next_instr')), welches später durch die Adresse des Befehls, der direkt auf den Sprungbefehl folgt, ersetzt wird.
RemoveStackframe()	Container für das Entfernen des aktuellen Stackframes durch das Wiederherstellen des im noch aktuellen Stackframe gespeicherten Werts des BAF-Registes der aufrufenden Funktion und das Setzen des SP-Registers auf den Wert des BAF-Registesr vor der Wiederherstellung.
File(name, decls_defs_blocks)	Container für alle Funkionen oder Blöcke, welche eine Datei als Ursprung haben, wobei name der Dateiname der Datei ist, die erstellt wird und decls_defs_blocks eine Liste von Funktionen bzw. Blöcken ist.
Block(name, stmts_instrs, instrs_before, num_instrs, param_size, local_vars_size)	Container für Statements, der auch als Block bezeichnet wird, wobei das Attribut name der Bezeichner des Labels (Definition 2.8) des Blocks ist und stmts_instrs eine Liste von Statements oder Befehlen. Zudem besitzt er noch 3 versteckte Attribute, wobei instrs_before die Zahl der Befehle vor diesem Block zählt, num_instrs die Zahl der Befehle ohne Kommentare in diesem Block zählt, param_size die voraussichtliche Anzahl an Speicherzellen aufaddiert, die für die Parameter der Funktion belegt werden müssen und local_vars_size die voraussichtliche Anzahl an Speicherzellen aufaddiert, die für die lokalen Variablen der Funktion belegt werden müssen.
GoTo(name)	Container für ein Goto zu einem anderen Block, wobei das Attribut name der Bezeichner des Labels des Blocks ist zu dem Gesprungen werden soll.
SingleLineComment(prefix, content)	Container für einen Kommentar, den der Compiler selber während des Kompiliervorangs erstellt, der im RETI-Interpreter selbst später nicht sichtbar sein wird, aber in den Immediate-Dateien, welche die Abstrakten Syntaxbäume nach den verschiedenen Passes enthalten.
RETIComment(value)	Container für einen Kommentar im Code der Form: // # comment, der im RETI-Intepreter später sichtbar sein wird und zur Orientierung genutzt werden kann, allerdings in einer tatsächlichen Implementierung einer RETI-CPU nicht umsetzbar ist und auch nicht sinnvoll wäre umzusetzen. Der Kommentar ist im Attribut value, welches jeder Knoten besitzt gespeichert.

#### Definition 2.8: Label

/

Durch einen Bezeichner eindeutig zuordenbares Sprungziel im Programmcode.<sup>a</sup>

<sup>a</sup>Thiemann, "Compilerbau".

# Anmerkung Q

Die ausgegrauten Attribute der PicoC-Knoten sind versteckte Attribute, die nicht direkt bei der Erstellung der PicoC-Knoten mit einem Wert initialisiert werden, sondern im Verlauf der Kompilierung beim Durchlaufen der verschiedenen Passes etwas zugewiesen bekommen, dass im weiteren Kompiliervorgang Informationen transportiert, die später im Kompiliervorgang nicht mehr so leicht zugänglich wären.

Jeder Knoten hat darüberhinaus auch noch 2 Attribute value und position, wobei value bei einem Knoten, der einen Blatt darstellt dem Tokenwert des Tokens, welches es ersetzt entspricht und Inneren Knoten unbesetzt ist. Das Attribut position wird für Fehlermeldungen gebraucht.

#### 2.2.5.2 RETI-Knoten

Bei den RETI-Knoten handelt es sich um Knoten, die irgendeinen Ausdruck aus der Sprache  $L_{RETI}$  darstellen. Für die RETI-Knoten wurden aus bereits in Unterkapitel 2.2.5.1 erläutertem Grund, genauso wie für die RETI-Knoten möglichst kurze und leicht verständliche Bezeichner gewählt. Alle RETI-Knoten, die in den von den verschiedenen Passes generierten Abstrakten Syntaxbäumen vorkommen sind in Tabelle 2.2.5.1 mit einem Beschreibungstext dokumentiert.

RETI-Knoten	Beschreibung
Program(name, instrs)	Container für alle Befehle: <name> <instrs>, wobei name</instrs></name>
	der Dateiname der Datei ist, die erstellt wird und instrs
	eine Liste von Befehlen ist.
<pre>Instr(op, args)</pre>	Container für einen Befehl: <op> <args>, wobei op eine</args></op>
· 17 · 0 ·	Operation ist und args eine Liste von Argumenten
	für dieser Operation.
<pre>Jump(rel, im_goto)</pre>	Container für einen Sprungbefehl: JUMP <rel> <im>, wo-</im></rel>
•	bei rel eine Relation ist und im_goto ein Immediate
	Value Im(val) für die Anzahl an Speicherzellen, um
	die relativ zum Sprungbefehl gesprungen werden soll
	oder ein GoTo(Name('block.xyz')), das später im RETI-
	Patch Pass durch einen passenden Immediate Value
	ersetzt wird.
Int(num)	Container für einen Interruptaufruf: INT <im>, wobei num</im>
	die Interrruptvektornummer (IVN) für die passende
	Speicherzelle in der Interruptvektortabelle ist, in der
	die Adresse der Interrupt-Service-Routine (ISR) steht.
Call(name, reg)	Container für einen <b>Prozeduraufruf</b> : CALL <name> <reg>,</reg></name>
	wobei name der Bezeichner der Prozedur, die aufgerufen
	werden soll ist und reg ein Register ist, das als Argu-
	ment an die Prozedur dient. Diese Operation ist in der
	Betriebssysteme Vorlesung $^a$ nicht deklariert, sondern wur-
	de dazuerfunden, um unkompliziert ein CALL PRINT ACC
	oder CALL INPUT ACC im RETI-Interpreter simulieren zu
	können.
Name(val)	Bezeichner für eine Prozedur, z.B. PRINT oder INPUT oder
	den Programnamen, z.B. PROGRAMNAME. Dieses Argu-
	ment ist in der Betriebssysteme Vorlesung <sup>a</sup> nicht dekla-
	riert, sondern wurde dazuerfunden, um Bezeichner, wie
	PRINT, INPUT oder PROGRAMNAME schreiben zu können.
Reg(reg)	Container für ein Register.
Im(val)	Ein Immediate Value, z.B. 42, -3 usw.
Add(), Sub(), Mult(), Div(), Mod(),	Compute-Memory oder Compute-Register Operatio-
Oplus(), Or(), And()	nen: ADD, SUB, MULT, DIV, OPLUS, OR, AND.
Addi(), Subi(), Multi(), Divi(), Modi(), Oplusi(), Ori(), Andi()	Compute-Immediate Operationen: ADDI, SUBI, MULTI, DIVI, MODI, OPLUSI, ORI, ANDI.
Load(), Loadin(), Loadi()	Load Operationen: LOAD, LOADIN, LOADI.
Store(), Storein(), Move()	Store Operationen: STORE, STOREIN, MOVE.
Lt(), LtE(), Gt(), GtE(), Eq(), NEq(),	Relationen: <, <=, >, >=, ==, !=, _NOP.
Always(), NOp()	
Rti()	Return-From-Interrupt Operation: RTI.
Pc(), In1(), In2(), Acc(), Sp(), Baf(),	Register: PC, IN1, IN2, ACC, SP, BAF, CS, DS.
Cs(), Ds()	

<sup>&</sup>lt;sup>a</sup> P. D. C. Scholl, "Betriebssysteme"

Tabelle 2.7: RETI-Knoten

# 2.2.5.3 Kompositionen von PicoC-Knoten und RETI-Knoten mit besonderer Bedeutung

In Tabelle 2.8 sind jegliche Kompositionen von PicoC-Knoten und RETI-Knoten aufgelistet, die eine besondere Bedeutung haben.

Komposition	Beschreibung
Ref(Global(Num('addr')))	Speichert Adresse der Speicherzelle, die Num('addr') Speicherzellen relativ zum Datensegment Register DS steht auf den Stack.
Ref(Stackframe(Num('addr')))	Speichert Adresse der Speicherzelle, die Num('addr') Speicherzellen relativ zum Begin-Aktive-Funktion Register BAF steht auf den Stack.
<pre>Ref(Subscr(Stack(Num('addr1')), Stack(Num('addr2'))))</pre>	Berechnet die nächste Adresse aus der Adresse, die an Speicherzelle Stack(Num('addr1')) steht und dem Subscript Index, der an Speicherzelle Stack(Num('addr2')) steht und speichert diese auf den Stack. Die Berechnung ist abhängig davon ob der Datentyp ArrayDecl(datatype) oder PntrDecl(datatype) ist. Der Datentyp ist ein verstecktes Attribut von Ref(exp).
<pre>Ref(Attr(Stack(Num('addr1')), Name('attr')))</pre>	Berechnet die nächste Adresse aus der Adresse, die an Speicherzelle Stack(Num('addr1')) steht und dem Attributnamen Name('attr') und speichert diese auf den Stack. Zur Berechnung ist der Name des Struct in StructSpec(Name('st')) notwendig, dessen Attribut Name('attr') ist. StructSpec(Name('st')) ist ein verstecktes Attribut von Ref(exp).
<pre>Assign(Stack(Num('size'))), Global(Num('addr')))</pre>	Schreibt Num('size') viele Speicherzellen, die ab Global(Num('addr')) relativ zum Datensegment Register DS stehen, versetzt genauso auf den Stack.
Assign(Stack(Num('size')), Stackframe(Num('addr')))	Schreibt Num('size') viele Speicherzellen, die ab Stackframe(Num('addr')) relativ zum Begin-Aktive-Funktion Register BAF stehen, versetzt genauso auf den Stack.
<pre>Exp(Global(Num('addr'))</pre>	Speichert Inhalt der Speicherzelle, die Num('addr') Speicherzellen relativ zum Datensegment Register DS steht auf den Stack.
<pre>Exp(Stackframe(Num('addr'))</pre>	Speichert Inhalt der Speicherzelle, die Num('addr') Speicherzellen relativ zum Begin-Aktive-Funktion Register BAF steht auf den Stack.
<pre>Exp(Stack(Num('addr')))</pre>	Speichert Inhalt der Speicherzelle, die Num('addr') Speicherzellen relativ zum Stackpointer Register SP steht auf den Stack.
Assign(Stack(Num('addr1')), Stack(Num('addr2')))	Speichert Inhalt der Speicherzelle Stack(Num('addr2')), die Num('addr2') Speicherzellen relativ zum Stackpoin- ter Register SP steht an der Adresse in der Speicherzelle, die Num('addr1') Speicherzellen relativ zum Stackpoin- ter Register SP steht.
<pre>Assign(Global(Num('addr')), Stack(Num('size')))</pre>	Schreibt Num('size') viele Speicherzellen, die auf dem Stack stehen, versetzt genauso auf die Speicherzellen ab Num('addr') relativ zum Datensegment Register DS.
<pre>Assign(Stackframe(Num('addr')), Stack(Num('size')))</pre>	Schreibt Num('size') viele Speicherzellen, die auf dem Stack stehen, versetzt genauso auf die Speicherzellen ab Num('addr') relativ zum Begin-Aktive-Funktion Register BAF.
<pre>Exp(Reg(reg))</pre>	Schreibt den aktuellen Wert des Registers reg auf den Stack.
<pre>Instr(Loadi(), [Reg(Acc()), GoTo(Name('addr@next_instr'))])</pre>	Lädt in das Register ACC die Adresse des Befehls, der in diesem Kontext direkt nach dem Sprung zum Block einer anderen Funktion steht.

Tabelle 2.8: Kompositionen von PicoC-Knoten und RETI-Knoten mit besonderer Bedeutung

# Anmerkung Q

Um die obige Tabelle 2.8 nicht mit unnötig viel repetetiven Inhalt zu füllen, wurden die zahlreichen Kompostionen ausgelassen, bei denen einfach nur exp durch  $Stack(Num('x')), x \in \mathbb{N}$  ersetzt wurde.

Zudem sind auch jegliche Kombinationen ausgelassen, bei denen einfach nur eine **Expression** an ein Exp(exp) bzw. Ref(exp) drangehängt wurde.

## 2.2.5.4 Abstrakte Grammatik

Die Abstrakte Syntax der Sprache  $L_{PicoC}$  wird durch die Abstrakte Grammatik 2.2.10 beschrieben.

stmt	::=	$SingleLineComment(str, str) \mid RETIComment()$	$L_{-}Comment$
un_op bin_op exp	::= - - - - - - - ::=	$\begin{array}{c cccc} Minus() &   & Not() \\ Add() &   & Sub() &   & Mul() &   & Div() &   & Mod() \\ Oplus() &   & And() &   & Or() \\ Name(str) &   & Num(str) &   & Char(str) \\ BinOp(\langle exp\rangle, \langle bin\_op\rangle, \langle exp\rangle) &   & Call(Name('input'), Empty()) \\ UnOp(\langle un\_op\rangle, \langle exp\rangle) &   & Call(Name('print'), \langle exp\rangle) \\ Exp(\langle exp\rangle) & \end{array}$	$L\_Arith$
$un\_op$ $rel$ $bin\_op$ $exp$	::= ::= ::=	$\begin{array}{c cccc} LogicNot() \\ Eq() &   & NEq() &   & Lt() &   & LtE() &   & Gt() &   & GtE() \\ LogicAnd() &   & LogicOr() \\ Atom(\langle exp \rangle, \langle rel \rangle, \langle exp \rangle) &   & ToBool(\langle exp \rangle) \end{array}$	$L\_Logic$
type_qual datatype exp stmt	::= ::= ::=	$Const() \mid Writeable() \\ IntType() \mid CharType() \mid VoidType() \\ Alloc(\langle type\_qual \rangle, \langle datatype \rangle, Name(str)) \\ Assign(\langle exp \rangle, \langle exp \rangle)$	$L\_Assign\_Alloc$
$\begin{array}{c} \hline datatype \\ exp \end{array}$	::=	$PntrDecl(Num(str), \langle datatype \rangle)$ $Deref(\langle exp \rangle, \langle exp \rangle) \mid Ref(\langle exp \rangle)$	$L\_Pntr$
$\begin{array}{c} \hline datatype \\ exp \end{array}$	::=	$\begin{array}{c c} ArrayDecl(Num(str)+,\langle datatype\rangle) \\ Subscr(\langle exp\rangle,\langle exp\rangle) &   & Array(\langle exp\rangle+) \end{array}$	L_Array
datatype exp decl_def	::= ::=   ::=	$StructSpec(Name(str)) \\ Attr(\langle exp \rangle, Name(str)) \\ Struct(Assign(Name(str), \langle exp \rangle) +) \\ StructDecl(Name(str), \\ Alloc(Writeable(), \langle datatype \rangle, Name(str)) +) \\$	$L\_Struct$
stmt	::=	$If(\langle exp \rangle, \langle stmt \rangle *) $ $IfElse(\langle exp \rangle, \langle stmt \rangle *, \langle stmt \rangle *)$	$L\_If\_Else$
stmt	::=	$While(\langle exp \rangle, \langle stmt \rangle *) \\ DoWhile(\langle exp \rangle, \langle stmt \rangle *)$	$L\_Loop$
$exp$ $stmt$ $decl\_def$	::= ::= ::=	$Call(Name(str), \langle exp \rangle *) \\ Return(\langle exp \rangle) \\ FunDecl(\langle datatype \rangle, Name(str), \\ Alloc(Writeable(), \langle datatype \rangle, Name(str)) *) \\ FunDef(\langle datatype \rangle, Name(str), \\ Alloc(Writeable(), \langle datatype \rangle, Name(str)) *, \langle stmt \rangle *) \\$	L_Fun
$\overline{file}$	::=	$File(Name(str), \langle decl\_def \rangle *)$	$L\_File$

Grammatik 2.2.10: Abstrakte Grammatik der Sprache L<sub>PiocC</sub>

# Anmerkung Q

In dieser Schrifftlichen Ausarbeitung der Bachelorarbeit wird oft von der "Abstrakten Grammatik der Sprache  $L_{PicoC}$ " agesprochen. Gemeint ist mit der Sprache  $L_{PicoC}$  nicht die Sprache, welche durch die Abstrakte Grammatik beschrieben wird. Es ist damit immer die Sprache gemeint, die kompiliert werden soll und zu deren Zweck die Abstrakt Grammatik überhaupt definiert wird. Für die tatsächliche Sprache, die durch die Abstrakt Grammatik beschrieben wird, interessiert man sich nicht explizit. Diese Konvention wurde aus der Quelle G. Siek,  $Course\ Webpage\ for\ Compilers$ 

```
(P423,\ P523,\ E313,\ and\ E513) übernommen. {}^a\overline{\rm Manchmal} auch von der Abstrakten Syntax der Sprache L_{PicoC}.
```

## 2.2.5.5 Codebeispiel

In Code 2.5 ist der Abstrakte Syntaxbaum zu sehen, der aus dem vereinfachten Ableitungsbaum aus Code 2.4 mithilfe eines Transformers generiert wurde.

```
File
     Name './verbose_dt_simple_ast_gen_array_decl_and_alloc.ast',
 4
5
       StructDecl
         Name 'st',
 6
7
8
9
            Alloc
              Writeable,
              PntrDecl
                Num '1',
11
                ArrayDecl
12
13
                     Num '4',
14
                     Num '5'
                  ],
16
                  PntrDecl
17
                     Num '1',
18
                     IntType 'int',
19
              Name 'attr'
20
         ],
       FunDef
22
         VoidType 'void',
23
         Name 'main',
24
          [],
25
          [
            Exp
26
27
              Alloc
28
                Writeable,
29
                ArrayDecl
30
31
                     Num '3',
                     Num '2'
33
                  ],
34
                  {\tt PntrDecl}
35
                     Num '1',
36
                     PntrDecl
                       Num '1',
37
38
                       StructSpec
39
                         Name 'st',
40
                Name 'var'
41
         ]
42
     ]
```

Code 2.5: Aus einem vereinfachtem Ableitungsbaum generierter Abstrakter Syntaxbaum

### 2.2.5.6 Ausgabe des Abstrakten Syntaxbaumes

Ein Teilbaum eines Abstrakten Syntaxbaumes kann entweder in der Konkretten Syntax der Sprache, für dessen Kompilierung er generiert wurde oder in der Abstrakten Syntax, die beschreibt, wie der Abstrakte Syntaxbaum selbst aufgebaut sein darf ausgegeben werden.

Das Ausgeben eines Abstrakten Syntaxbaumes wird im PicoC-Compiler über die Magische Methode  $\_repr\_()^{20}$  der Programmiersprache  $L_{Python}$  umgesetzt. Sobald ein PicoC-Knoten oder RETI-Knoten ausgegeben werden soll, gibt seine Magische Methode  $\_repr\_()$  eine nach der Abstrakten oder Konkretten Syntax aufgebaute Textrepräsentation seiner selbst und all seiner Knoten mit an den richtigen Stellen passend gesetzten runden öffnenden ( und schließenden ) Klammern, sowie Kommas ', ', Semikolons ; usw. zur Darstellung der Hierarchie und zur Abtrennung zurück. Dabei wird nach dem Prinzip der Tiefensuche der gesamte Abstrakte Syntaxbaum durchlaufen und die Magische  $\_repr\_()$ -Methode der verschiedenen Knoten aufgerufen, die immer jeweils die  $\_repr\_()$ -Methode ihrer Kinder aufrufen und die zurückgegebene Textrepräsentation passend zusammenfügen und selbst zurückgeben.

Beim PicoC-Compiler wurden Abstrakte und Konkrette Syntax miteinander gemischt. Für PicoC-Knoten wurde die Abstrakte Syntax verwendet, da Passes schließlich auf Abstrakten Syntaxbäumen operieren. Bei RETI-Knoten wurde die Konkrette Syntax verwendet, da Maschienenbefehle in Konkretter Syntax schließlich das Endprodukt des Kompiliervorgangs sein sollen. Da die Abstrakte Syntax von RETI-Knoten so simpel ist, macht es kaum einen Unterschied in der Erkennbarkeit, bis auf fehlende gescheifte Klammern () usw., ob man die RETI-Knoten in Abstrakter oder Konkretter Syntax schreibt. Daher kann man auch einfach gleich die RETI-Knoten in Konkretter Syntax ausgeben und muss nicht beim letzten Pass daran denken, am Ende die Konkrette, statt der Abstrakten Syntax für die RETI-Knoten auszugeben.

# 2.3 Code Generierung

Nach der Generierung eines Abstrakten Syntaxbaumes als Ergebnis der Lexikalischen und Syntaktischen Analyse in Unterkapitel ??, wird in diesem Kapitel auf Basis der verschiedenen Kompositionen von PicoC-Knoten und RETI-Knoten im Abstrakten Syntaxbaum das gewünschte Endprodukt des PicoC-Compilers, der RETI-Code generiert.

Man steht nun dem Problem gegenüber einen Abstrakten Syntaxbaum der Sprache  $L_{PicoC}$ , der durch die Abstrakte Grammatik 2.2.10 spezifiziert ist in einen entsprechenden Abstrakten Syntaxbaum der Sprache  $L_{RETI}$  umzuformen. Das ganze lässt sich, wie in Unterkapitel ?? bereits beschrieben vereinfachen, indem man dieses Problem in mehrere Passes (Definition ??) herunterbricht.

Beim PicoC-Compiler handelt es sich um einen Cross-Compiler (Definiton ??). Damit RETI-Code erzeugt werden kann, der auf der RETI-Architektur läuft, muss erst, wie im T-Diagram (siehe Unterkapitel ??) in Abbildung 2.6 zu sehen ist, der Python-Code des PicoC-Compilers mittels eines Compilers, der z.B. auf einer X<sub>86\_64</sub>-Architektur laufen könnte zu Bytecode kompiliert werden. Dieser Bytecode wird dann von der Python-Virtual-Machine (PVM) interpretiert, welche wiederum auf einer X<sub>86\_64</sub>-Architektur laufen könnte. Und selbst dieses T-Diagram könnte noch ausführlicher ausgedrückt werden, indem nachgeforscht wird, in welcher Sprache eigentlich die Python-Virtual-Machine geschrieben war, bevor sie zu X<sub>86\_64</sub> kompiliert wurde usw.

<sup>&</sup>lt;sup>20</sup>Spezielle Methode, die immer aufgerufen wird, wenn das Object, dass in Besitz dieser Methode ist als String mittels print() oder zur Repräsentation ausgegeben werden soll.

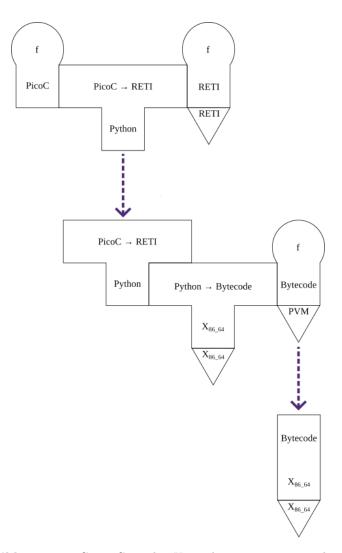


Abbildung 2.6: Cross-Compiler Kompiliervorgang ausgeschrieben

Dieses längliche T-Diagram in Abbildung 2.6 lässt sich zusammenfassen, sodass man das T-Diagram in Abbildung 2.7 erhält, in welcher direkt angegeben ist, dass der PicoC-Compiler in  $X_{86\_64}$ -Maschienensprache geschrieben ist.

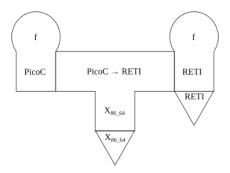


Abbildung 2.7: Cross-Compiler Kompiliervorgang Kurzform

Nachdem der Kompilierprozess des PicoC-Compiler im vertikalen nun genauer angesehen wurde, wird der

Kompilierprozess im Folgenden im horinzontalen, auf der Ebene der verschiedenen Passes genauer betrachtet. Die Abbildung 2.8 gibt einen guten Überblick über alle Passes und wie diese in der Pipe-Architektur (Definition ??) des PicoC-Compilers aufeinanderfolgen. In der Pipe-Architektur nutzt der jeweils nächste Pass den generierten Abstrakten Syntaxbaum des vorherigen Passes oder der Syntaktischen Analyse, um einen eigenen Abstrakten Syntaxbaum in seiner eigenen Sprache zu generieren.

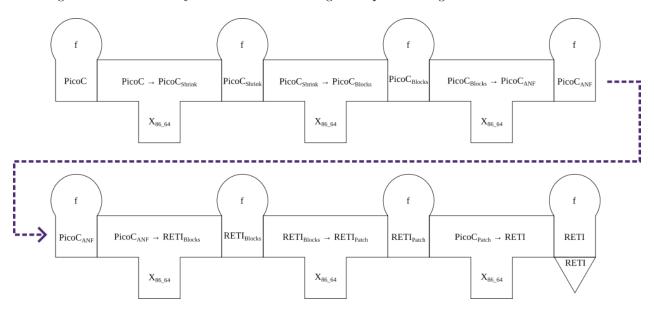


Abbildung 2.8: Architektur mit allen Passes ausgeschrieben

Im Unterkapitel 2.3.1 werden die unterschiedlichen Passes des PicoC-Compilers erklärt. In den darauffolgenden Unterkapiteln ??, ??, ?? und ?? zu Zeigern, Feldern, Verbunden und Funktionen werden einzelne Aspekte, die Thema dieser Bachelorarbeit sind genauer betrachtet und erklärt, die im Unterkapitel 2.3.1 nicht ausreichend vertieft wurden. Viele der verwendenten Ansätze zur Lösung dieser Probleme basieren auf der Vorlesung P. D. C. Scholl, "Betriebssysteme" und wurden in dieser Bachelorarbeit weiter ausgearbeitet, wo es nötig war, sodass diese mit dem PicoC-Compiler auch in der Praxis implementiert werden konnten.

## 2.3.1 Passes

Im Folgenden werden die verschiedenen Passes des PicoC-Compilers für die Generierung von RETI-Code besprochen. Viele dieser Passes haben Aufgaben, die eher unter die Themenbereiche des Bachelorprojekts fallen. Allerdings ist das Verständnis der Passes auch für das Verständnis der veschiedenen Aspekte<sup>21</sup> der Bachelorarbeit wichtig.

Auf jedes Detail der einzelnen Passes wird in diesem Unterkapitel allerdings nicht eingegangen, da diese einerseits in den Unterkapiteln ??, ??, ?? und ?? zu Zeigern, Feldern, Verbunden und Funktionen im Detail erklärt sind und andererseits viele Aufgaben dieser Passes eher dem Bachelorprojekt zuzurechnen sind.

#### 2.3.1.1 PicoC-Shrink Pass

#### 2.3.1.1.1 Aufgabe

Der Aufgabe des PicoC-Shrink Pass ist in Unterkapitel ?? ausführlich an einem Beispiel erklärt. Kurzgefasst hat der PicoC-Shrink Pass die Aufgabe, die Eigenheit auszunutzen, dass der Dereferenzierungoperator \*pntr und die damit einhergehende Zeigerarithmetik \*(pntr + i) sich in der Untermenge der Sprache  $L_C$ ,

<sup>&</sup>lt;sup>21</sup>In kurz: Zeiger, Felder, Verbunde und Funktionen.

welche die Sprache  $L_{PicoC}$  darstellt genau gleich verhält, wie der Operator für den Zugriff auf den Index eines Feldes ar[i].

Daher wandelt der PicoC-Shrink Pass alle Verwendungen des Knoten Deref(exp, i) im jeweiligen Abstrakten Syntaxbaum in Knoten Subscr(exp, i) um, sodass sich dadurch viele vermeidbare Fallunterscheidungen und doppelter Code bei der Implementierung vermeiden lassen. Man lässt die Derefenzierung \*(var + i) einfach von den Routinen für einen Zugriff auf einen Feldindex var[i] übernehmen.

#### 2.3.1.1.2 Abstrakte Grammatik

Die Abstrakte Grammatik 2.3.1 der Sprache  $L_{PicoC\_Shrink}$  ist fast identisch mit der Abstrakten Grammatik 2.2.10 der Sprache  $L_{PicoC}$ , nach welcher der erste Abstrakte Syntaxbaum in der Syntaktischen Analyse generiert wurde. Der einzige Unterschied liegt darin, dass es den Knoten Deref(exp, exp) in Abstrakten Grammatik 2.3.1 nicht mehr gibt. Das liegt daran, dass dieser Pass alle Vorkommnisse des Knoten Deref(exp, exp) durch den Knoten Subscr(exp, exp) auswechselt, der ebenfalls nach der Abstrakten Grammatik der Sprache  $L_{PicoC}$  definiert ist.

stmt	::=	$SingleLineComment(str, str) \mid RETIComment()$	$L_{-}Comment$
un_op bin_op exp	::=	$\begin{array}{c cccc} Minus() &   & Not() \\ Add() &   & Sub() &   & Mul() &   & Div() &   & Mod() \\ Oplus() &   & And() &   & Or() \\ Name(str) &   & Num(str) &   & Char(str) \\ BinOp(\langle exp\rangle, \langle bin\_op\rangle, \langle exp\rangle) & & UnOp(\langle un\_op\rangle, \langle exp\rangle) &   & Call(Name('print'), \langle exp\rangle) \\ Call(Name('print'), \langle exp\rangle) &   & Call(Name('print'), \langle exp\rangle) \end{array}$	$L\_Arith$
$\frac{stmt}{}$	::=	$Exp(\langle exp \rangle)$	
un_op rel bin_op exp	::= ::= ::=	$LogicNot() \\ Eq() \mid NEq() \mid Lt() \mid LtE() \mid Gt() \mid GtE() \\ LogicAnd() \mid LogicOr() \\ Atom(\langle exp \rangle, \langle rel \rangle, \langle exp \rangle) \mid ToBool(\langle exp \rangle)$	$L\_Logic$
type_qual datatype exp stmt	::= ::= ::=	$Const() \mid Writeable()$ $IntType() \mid CharType() \mid VoidType()$ $Alloc(\langle type\_qual \rangle, \langle datatype \rangle, Name(str))$ $Assign(\langle exp \rangle, \langle exp \rangle)$	$L\_Assign\_Alloc$
$\begin{array}{c} datatype \\ exp \end{array}$	::=	$PntrDecl(Num(str), \langle datatype \rangle)$ $Deref(\langle exp \rangle, \langle exp \rangle) \mid Ref(\langle exp \rangle)$	$L\_Pntr$
$\begin{array}{c} datatype \\ exp \end{array}$	::=	$\begin{array}{ll} ArrayDecl(Num(str)+,\langle datatype\rangle) \\ Subscr(\langle exp\rangle,\langle exp\rangle) &   & Array(\langle exp\rangle+) \end{array}$	$L\_Array$
datatype exp decl_def	::= ::=   ::=	StructSpec(Name(str)) $Attr(\langle exp \rangle, Name(str))$ $Struct(Assign(Name(str), \langle exp \rangle)+)$ $StructDecl(Name(str), \langle datatype \rangle, Name(str))+)$	$L\_Struct$
stmt	::=	$If(\langle exp \rangle, \langle stmt \rangle *)$ $IfElse(\langle exp \rangle, \langle stmt \rangle *, \langle stmt \rangle *)$	$L\_If\_Else$
stmt	::=	$While(\langle exp \rangle, \langle stmt \rangle *) \\ DoWhile(\langle exp \rangle, \langle stmt \rangle *)$	$L\_Loop$
$exp$ $stmt$ $decl\_def$	::= ::=	$Call(Name(str), \langle exp \rangle *)$ $Return(\langle exp \rangle)$ $FunDecl(\langle datatype \rangle, Name(str),$ $Alloc(Writeable(), \langle datatype \rangle, Name(str)) *)$ $FunDef(\langle datatype \rangle, Name(str),$ $Alloc(Writeable(), \langle datatype \rangle, Name(str)) *, \langle stmt \rangle *)$	$L\_Fun$
file	::=	$File(Name(str), \langle decl\_def \rangle *)$	$L$ _ $File$

Grammatik 2.3.1: Abstrakte Grammatik der Sprache  $L_{PiocC\_Shrink}$ 

# Anmerkung 9

Der rot markierte Knoten bedeutet, dass dieser im Vergleich zur voherigen Abstrakten Grammatik nicht mehr da ist.

## 2.3.1.1.3 Codebeispiel

In den nächsten Unterkapiteln wird das Beispiel in Code 2.6 zur Anschauung der verschiedenen Passes

verwendet. Im Code 2.6 ist in der Funktion faculty ein iterativer Algorithmus implementiert, der die Fakultät eines übergebenen Arguments berechnet. Der Algorithmus basiert auf einem Beispielprogramm aus der Vorlesung P. D. C. Scholl, "Betriebssysteme", welcher in der Vorlesung allerdings rekursiv implementiert ist.

Dieser rekursive Algoirthmus ist allerdings kein gutes Anschaungsbeispiel, dass viele der Aufgaben der verschiedenen Passes bei der Kompilierung veranschaulicht hätte. Viele Aufgaben der Passes, wie z.B. bei der Kompilierung von if-, if-else-, while- und do-while-Statements wären im Beispiel aus der Vorlesung nicht enthalten gewesen. Daher wurde das Beispiel aus der Vorlesung zu einem iterativen Algorithmus 2.6 umgeschrieben, um if- und while-Statemtens zu enthalten.

Beide Varianten des Algorithmus wurden zum Testen des PicoC-Compilers verwendet und sind als Tests im Ordner /tests unter Link<sup>22</sup>, unter den Testbezeichnungen example\_faculty\_rec.picoc und example\_faculty\_it.picoc zu finden.

Die Codebeispiele in diesem und den folgenden Unterkapiteln dienen allerdings nur als Anschauung des jeweiligen Passes, der in diesem Unterkapitel beschrieben wird und werden nicht im Detail erläutert, da viele Details der Passes später in den Unterkapiteln ??, ??, ?? und ?? zu Zeigern, Feldern, Verbunden und Funktionen mit eigenen Codebeispielen erklärt werden und alle sonstigen Details dem Bachelorprojekt zuzurechnen sind.

```
based on a example program from Christoph Scholl's Operating Systems lecture
  int faculty(int n){
4
    int res = 1;
    while (1) {
      if (n == 1) {
         return res;
9
      res = n * res;
10
          n-1;
11
12 }
13
  void main() {
15
    print(faculty(4));
16 }
```

Code 2.6: PicoC Code für Codebespiel

In Code 2.7 sieht man den Abstrakten Syntaxbaum, der in der Syntaktischen Analyse generiert wurde.

```
1 File
2  Name './example_faculty_it.ast',
3  [
4   FunDef
5   IntType 'int',
6   Name 'faculty',
7  [
```

 $<sup>^{22}</sup>$ https://github.com/matthejue/PicoC-Compiler/tree/new\_architecture/tests.

```
Alloc(Writeable(), IntType('int'), Name('n'))
 9
         ],
10
         Γ
11
           Assign(Alloc(Writeable(), IntType('int'), Name('res')), Num('1')),
12
           While
13
             Num '1',
14
             Γ
               Ιf
16
                 Atom(Name('n'), Eq('=='), Num('1')),
17
18
                   Return(Name('res'))
19
20
               Assign(Name('res'), BinOp(Name('n'), Mul('*'), Name('res')))
21
               Assign(Name('n'), BinOp(Name('n'), Sub('-'), Num('1')))
22
23
         ],
24
       FunDef
25
         VoidType 'void',
26
         Name 'main',
27
         [],
28
         Γ
29
           Exp(Call(Name('print'), [Call(Name('faculty'), [Num('4')])))
30
31
     ]
```

Code 2.7: Abstrakter Syntaxbaum für Codebespiel

Im PicoC-Shrink-Pass ändert sich nichts im Vergleich zum Abstrakten Syntaxbaum in Code 2.7, da das Codebeispiel keine Dereferenzierung enthält.

#### 2.3.1.2 PicoC-Blocks Pass

#### 2.3.1.2.1 Aufgabe

Die Aufgabe des PicoC-Blocks Passes ist es die Knoten If(exp, stmts), IfElse(exp, stmts1, stmts2), While(exp, stmts) und DoWhile(exp, stmts) mithilfe von Block(name, stmts\_instrs-, GoTo(lable)- und IfElse(exp, stmts1, stmts2)-Knoten umzusetzen. Der IfElse(exp, stmts1, stmts2)-Knoten wird zur Umsetzung der Bedingung verwendet und es wird, je nachdem, ob die Bedingung wahr oder falsch ist mithilfe der GoTo(label)-Knoten in einen von zwei alternativen Branches gesprungen oder ein Branch erneut aufgerufen usw.

#### 2.3.1.2.2 Abstrakte Grammatik

Zur Umsetzung dieses Passes ist es notwendig die Abstrakte Grammatik 2.3.1 der Sprache  $L_{PicoC\_Shrink}$  um die Knoten zu erweitern, die im Unterkapitel 2.3.1.2.1 erwähnt wurden. Die Knoten If(exp, stmts), While(exp, stmts) und DoWhile(exp, stmts) gibt es nicht mehr, da sie durch Block(name, stmts\_instrs-, GoTo(lable)-und IfElse(exp, stmts1, stmts2)-Knoten ersetzt wurden. Die Funktionsdefinition FunDef( $\langle datatype \rangle$ , Name(str), Alloc(Writeable(),  $\langle datatype \rangle$ , Name(str))\*,  $\langle block \rangle$ \*) ist nun ein Container für Blöcke Block(Name(str),  $\langle stmt \rangle$ \*) und keine Statements stmt mehr. Das resultiert in der Abstrakten Grammatik 2.3.2 der Sprache  $L_{PicoC\_Blocks}$ .

stmt	::=	$SingleLineComment(str, str) \mid RETIComment()$	$L\_Comment$
un_op	::=	$Minus() \mid Not()$	$L\_Arith$
$bin\_op$	::=	$Add() \mid Sub() \mid Mul() \mid Div() \mid Mod()$ $Oplus() \mid And() \mid Or()$	
exp	::=	$Name(str) \mid Num(str) \mid Char(str)$	
		$BinOp(\langle exp \rangle, \langle bin\_op \rangle, \langle exp \rangle)$	
		$UnOp(\langle un\_op \rangle, \langle exp \rangle) \mid Call(Name('input'), Empty())$ $Call(Name('print'), \langle exp \rangle)$	
stmt	::=	$Exp(\langle exp \rangle)$	
un_op	::=	LogicNot()	$L\_Logic$
rel	::=	$Eq() \mid NEq() \mid Lt() \mid LtE() \mid Gt() \mid GtE()$	
$bin\_op$	::=	$LogicAnd() \mid LogicOr()$	
exp	::=	$Atom(\langle exp \rangle, \langle rel \rangle, \langle exp \rangle) \mid ToBool(\langle exp \rangle)$	
type_qual	::=	Const()   Writeable()	$L\_Assign\_Alloc$
datatype	::=	$IntType() \mid CharType() \mid VoidType()$	
exp $stmt$	::=	$Alloc(\langle type\_qual \rangle, \langle datatype \rangle, Name(str))$ $Assign(\langle exp \rangle, \langle exp \rangle)$	
			T. D. (
datatype	::=	$PntrDecl(Num(str), \langle datatype \rangle)$ $Ref(\langle exp \rangle)$	$L\_Pntr$
exp	::=	1 3 4 5 7 7	- ·
datatype	::=	$ArrayDecl(Num(str)+,\langle datatype \rangle)$	$L\_Array$
exp	::=	$Subscr(\langle exp \rangle, \langle exp \rangle) \mid Array(\langle exp \rangle +)$	
datatype	::=	StructSpec(Name(str))	$L\_Struct$
exp	::=	$Attr(\langle exp \rangle, Name(str))$ $Struct(Assign(Name(str), \langle exp \rangle)+)$	
$decl\_def$	::=	Struct(Assign(Name(str), (exp))+) StructDecl(Name(str), (exp))+)	
acciacj	••-	$Alloc(Writeable(), \langle datatype \rangle, Name(str))+)$	
stmt	::=	$If(\langle exp \rangle, \langle stmt \rangle *)$	L_If_Else
		$IfElse(\langle exp \rangle, \langle stmt \rangle *, \langle stmt \rangle *)$	
stmt	::=	$While(\langle exp \rangle, \langle stmt \rangle *)$	$L\_Loop$
		$DoWhile(\langle exp \rangle, \langle stmt \rangle *)$	
exp	::=	$Call(Name(str), \langle exp \rangle *)$	L_Fun
stmt	::=	$Return(\langle exp \rangle)$	
$decl\_def$	::=	$FunDecl(\langle datatype \rangle, Name(str),$	
	1	$Alloc(Writeable(), \langle datatype \rangle, Name(str))*)$	
		$FunDef(\langle datatype \rangle, Name(str), \\ Alloc(Writeable(), \langle datatype \rangle, Name(str))*, \langle block \rangle*)$	
11 1			I D1 1
$block \\ stmt$	::=	$Block(Name(str), \langle stmt \rangle *)$ GoTo(Name(str))	$L\_Blocks$
	-::=	(	T 77:1
file	::=	$File(Name(str), \langle decl\_def \rangle *)$	$L\_File$

Grammatik 2.3.2: Abstrakte Grammatik der Sprache  $L_{PiocC\_Blocks}$ 

# Anmerkung Q

Alles rot markierte bedeutet, es wurde entfernt oder abgeändert. Alles ausgegraute bedeutet, es hat sich im Vergleich zur letzten Abstrakten Grammatik nichts geändert. Alle normal in schwarz geschriebenen Knoten wurden neu hinzugefügt.

Die Abstrakte Grammatik soll im Gegensatz zur Konkretten Grammatik meist nur vom Programmierer verstanden werden, der den Compiler implementiert und sollte daher vor allem einfach verständlich sein und stellt daher eine Obermenge aller tatsächlich möglichen Kompositionen von Knoten dar<sup>a</sup>.

<sup>a</sup>D.h. auch wenn dort **exp** als **Attribut** steht, kann dort nicht jeder Knoten, der sich aus der **Produktion exp** ergibt auch wirklich eingesetzt werden.

#### 2.3.1.2.3 Codebeispiel

In Code 2.8 sieht man den Abstrakten Syntaxbaum des PiocC-Blocks Passes für das aus Unterkapitel 2.6 weitergeführte Beispiel, indem nun eigene Blöcke für die Funktion faculty und die main-Funktion erstellt werden, in denen die ersten Statements der jeweiligen Funktionen bis zum letzten Statement oder bis zum ersten Auftauchen eines If(exp, stmts)-, IfElse(exp, stmts1, stmts2)-, While(exp, stmts)-Knoten stehen. Je nachdem, ob ein If(exp, stmts)-, IfElse(exp, stmts1, stmts2)-, While(exp, stmts)- oder DoWhile(exp, stmts)- Knoten auftaucht, werden für die Bedingung und mögliche Branches eigene Blöcke erstellt.

```
Name './example_faculty_it.picoc_blocks',
 4
       FunDef
 5
         IntType 'int',
 6
         Name 'faculty',
           Alloc(Writeable(), IntType('int'), Name('n'))
 9
         ],
10
         Γ
11
           Block
12
             Name 'faculty.6',
13
14
               Assign(Alloc(Writeable(), IntType('int'), Name('res')), Num('1'))
15
               // While(Num('1'), [])
16
               GoTo(Name('condition_check.5'))
17
             ],
18
           Block
19
             Name 'condition_check.5',
20
             Ε
21
               IfElse
22
                 Num '1',
23
                  Γ
24
                    GoTo(Name('while_branch.4'))
25
                 ],
26
                  Γ
27
                    GoTo(Name('while_after.1'))
28
                 ]
29
             ],
30
           Block
31
             Name 'while_branch.4',
32
33
               // If(Atom(Name('n'), Eq('=='), Num('1')), []),
34
               IfElse
35
                  Atom(Name('n'), Eq('=='), Num('1')),
                  Ε
```

```
GoTo(Name('if.3'))
38
                  ],
39
                  Ε
                    GoTo(Name('if_else_after.2'))
                  ]
42
             ],
43
           Block
44
              Name 'if.3',
45
46
                Return(Name('res'))
47
              ],
48
           Block
49
              Name 'if_else_after.2',
50
              Γ
51
                Assign(Name('res'), BinOp(Name('n'), Mul('*'), Name('res')))
52
                Assign(Name('n'), BinOp(Name('n'), Sub('-'), Num('1')))
53
                GoTo(Name('condition_check.5'))
54
             ],
55
           Block
              Name 'while_after.1',
56
57
58
         ],
59
       FunDef
60
         VoidType 'void',
61
         Name 'main',
62
         [],
63
         Γ
64
           Block
65
              Name 'main.0',
66
67
                Exp(Call(Name('print'), [Call(Name('faculty'), [Num('4')])))
68
69
         ]
70
     ]
```

Code 2.8: PicoC-Blocks Pass für Codebespiel

## 2.3.1.3 PicoC-ANF Pass

#### 2.3.1.3.1 Aufgabe

Die Aufgabe des PicoC-ANF Passes ist es den Abstrakten Syntaxbaum der Sprache  $L_{PicoC\_Blocks}$  in die Abstrakte Grammatik der Sprache  $L_{PicoC\_ANF}$  umzuformen, welche in A-Normalform (Definition ??) und damit auch in Monadischer Normalform (Definition ??) ist. Um Wiederholung zu vermeiden wird zur Erklärung der A-Normalform auf Unterkapitel ?? verwiesen.

Zudem wird eine Symboltabelle (Definition 2.9) eingeführt. In der Symboltabelle wird beim Anlegen eines neuen Eintrags für eine Variable zunächst eine Adresse zugewiesen, die dem Wert einer von zwei Countern rel\_global\_addr und rel\_stack\_addr entspricht. Der Counter rel\_global\_addr ist für Variablen in den Globalen Statischen Daten und der Counter rel\_stack\_addr ist für Variablen auf dem Stackframe. Einer der beiden Counter wird entsprechend der Größe der angelegten Variable hochgezählt.

Kommt im Programmcode an einer späteren Stelle diese Variable Name('symbol') vor, so wird mit dem Symbol<sup>23</sup> als Schlüssel in der Symboltabelle nachgeschlagen und anstelle des Name(str)-Knotens die in

<sup>&</sup>lt;sup>23</sup>Bzw. der Bezeichner

der Symboltabelle nachgeschlagene Adresse in einem Global(Num('addr'))- bzw. Stackframe(Num('addr'))- Knoten eingesetzt eingefügt. Ob der Global(Num('addr'))- oder der Stackframe(Num('addr'))-Knoten zum Einsatz kommt, entscheidet sich anhand des Scopes (z.B. @scope), der in der Symboltabelle an den Bezeichner drangehängt ist (z.B. identifier@scope).<sup>24</sup>

#### Definition 2.9: Symboltabelle

Eine über ein Assoziatives Feld umgesetzte Datenstruktur, die notwendig ist, um das Konzept einer Variablen in einer Sprache umzusetzen. Diese Datenstruktur ordnet jedem Symbol<sup>a</sup> einer Variablen, Konstanten oder Funktion aus einem Programm, Informationen, wie die Adresse, die Position im Programmcode oder den Datentyp zu.

Die Symboltabelle muss nur während des Kompiliervorgangs im Speicher existieren, da die Einträge in der Symboltabelle beeinflussen, was für Maschinencode generiert wird und dadurch im Maschinencode bereits die richtigen Adressen usw. angesprochen werden und es die Symboltabelle selbst nicht mehr braucht.

<sup>a</sup>In einer Symboltabelle werden Bezeichner als Symbole bezeichnet.

#### 2.3.1.3.2 Abstrakte Grammatik

Zur Umsetzung dieses Passes ist es notwendig die Abstrakte Grammatik 2.3.2 der Sprache  $L_{PicoC\_Blocks}$  in die A-Normalform zu bringen. Darunter fällt es unter anderem, dafür zu sorgen, dass Komplexe Knoten, wie z.B. BinOp(exp, bin\_op, exp) nur Atomare Knoten, wie z.B. Stack(Num(str)) enthalten können. Des Weiteren werden auch Funktionen und Funktionsaufrufe aufgelöst, sodass u.a. die Blöcke Block(Name(str), stmt\*) nun direkt im File(Name(str), block\*)-Knoten liegen usw., was in Unterkapitel ?? genauer erklärt wird. Die Symboltabelle ist ebenfalls als Abstrakter Syntaxbaum umgesetzt, wofür in der Abstrakten Grammatik 2.3.3 der Sprache  $L_{PicoC\_ANF}$  der Sprache  $L_{PicoC\_ANF}$  neue Knoten eingeführt werden.

Das ganze resultiert in der Abstrakten Grammatik 2.3.3 der Sprache  $L_{PicoC\_ANF}$ .

<sup>&</sup>lt;sup>24</sup>Die Umsetzung von **Scopes** wird in Unterkapitel ?? genauer beschrieben.

```
RETIComment()
                                                                                                               L_{-}Comment
stmt
                        SingleLineComment(str, str)
                 ::=
                                                                                                               L_Arith
un\_op
                 ::=
                        Minus()
                                        Not()
bin\_op
                 ::=
                        Add()
                                  Sub()
                                                 Mul() \mid Div() \mid
                                                                           Mod()
                                                 |Or()
                        Oplus()
                                   And()
                        Name(str) \mid Num(str) \mid Char(str) \mid Global(Num(str))
exp
                        Stackframe(Num(str)) \mid Stack(Num(str))
                        BinOp(Stack(Num(str)), \langle bin\_op \rangle, Stack(Num(str)))
                        UnOp(\langle un\_op \rangle, Stack(Num(str))) \mid Call(Name('input'), Empty())
                        Call(Name('print'), \langle exp \rangle)
                        Exp(\langle exp \rangle)
                        LogicNot()
                                                                                                               L\_Logic
un\_op
                 ::=
                        Eq() \mid NEq() \mid Lt() \mid LtE() \mid Gt()
rel
                                                                                    GtE()
                 ::=
                        LogicAnd()
                                           LogicOr()
bin\_op
                 ::=
                        Atom(Stack(Num(str)), \langle rel \rangle, Stack(Num(str)))
exp
                 ::=
                        ToBool(Stack(Num(str)))
type\_qual
                        Const()
                                       Writeable()
                                                                                                               L\_Assign\_Alloc
                 ::=
                        IntType() \mid CharType() \mid VoidType()
datatype
                 ::=
exp
                        Alloc(\langle type\_qual \rangle, \langle datatype \rangle, Name(str))
                  ::=
                        Assign(Global(Num(str)), Stack(Num(str)))
stmt
                  ::=
                        Assign(Stackframe(Num(str)), Stack(Num(str)))
                        Assign(Stack(Num(str)), Global(Num(str)))
                        Assign(Stack(Num(str)), Stackframe(Num(str)))
                        PntrDecl(Num(str), \langle datatype \rangle)
                                                                                                               L_{-}Pntr
datatype
                 ::=
                        Ref(Global(str)) \mid Ref(Stackframe(str))
                                                       |Ref(Attr(\langle exp \rangle, Name(str)))|
                        Ref(Subscr(\langle exp \rangle, \langle exp \rangle))
                        ArrayDecl(Num(str)+, \langle datatype \rangle)
                                                                                                               L_-Array
datatupe
                 ::=
                        Subscr(\langle exp \rangle, Stack(Num(str)))
                                                                    Array(\langle exp \rangle +)
exp
                 ::=
datatype
                        StructSpec(Name(str))
                                                                                                               L_-Struct
                 ::=
                        Attr(\langle exp \rangle, Name(str))
exp
                  ::=
                        Struct(Assign(Name(str), \langle exp \rangle) +)
decl\_def
                        StructDecl(Name(str),
                  ::=
                              Alloc(Writeable(), \langle datatype \rangle, Name(str)) +)
                        IfElse(Stack(Num(str)), \langle stmt \rangle *, \langle stmt \rangle *)
                                                                                                               L_If_Else
stmt
                 ::=
                        Call(Name(str), \langle exp \rangle *)
                                                                                                               L_{-}Fun
                 ::=
exp
                        StackMalloc(Num(str)) \mid NewStackframe(Name(str), GoTo(str))
stmt
                 ::=
                        Exp(GoTo(Name(str))) \mid RemoveStackframe()
                        Return(Empty()) \mid Return(\langle exp \rangle)
decl\_def
                        FunDecl(\langle datatype \rangle, Name(str))
                 ::=
                              Alloc(Writeable(), \langle datatype \rangle, Name(str))*)
                        FunDef(\langle datatype \rangle, Name(str),
                              Alloc(Writeable(), \langle datatype \rangle, Name(str))*, \langle block \rangle*)
block
                        Block(Name(str), \langle stmt \rangle *)
                                                                                                               L\_Blocks
                 ::=
stmt
                        GoTo(Name(str))
                  ::=
                                                                                                               L_File
file
                        File(Name(str), \langle block \rangle *)
symbol\_table
                        SymbolTable(\langle symbol \rangle *)
                                                                                                               L\_Symbol\_Table
                 ::=
                        Symbol(\langle type\_qual \rangle, \langle datatype \rangle, \langle name \rangle, \langle val \rangle, \langle pos \rangle, \langle size \rangle)
symbol
                 ::=
                        Empty()
type\_qual
                 ::=
datatype
                 ::=
                        BuiltIn()
                                         SelfDefined()
                        Name(str)
name
                  ::=
val
                        Num(str)
                                          Empty()
                 ::=
                        Pos(Num(str), Num(str))
                                                          \perp Empty()
pos
                  ::=
                        Num(str)
                                         Empty()
size
                                                                                                                                57
```

#### 2.3.1.3.3 Codebeispiel

In Code 2.9 sieht man den Abstrakten Syntaxbaum des PiocC-ANF Passes für das aus Unterkapitel 2.6 weitergeführte Beispiel, indem alls Statements und Ausdrücke in A-Normalform sind. Die IfElse(exp, stmts, stmts)-Knoten sind hier in A-Normalform gebracht worden, indem ihre Komplexe Bedingung vorgezogen wurde und das Ergebnis der Komplexen Bedingung einer Location zugewiesen ist und sie selbst das Ergebnis über den Atomaren Ausdruck Stack(Num(str)) vom Stack lesen: IfElse(Stack(Num(str)), stmts, stmts). Funktionen sind nur noch über die Labels von Blöcken zu erkennen, die den gleichen Bezeichner haben, wie die ursprüngliche Funktion und es lässt sich nur durch das Nachverfolgen der GoTo(Name('label'))-Knoten nachvollziehen, was ursprünglich zur Funktion gehörte.

```
1
  File
 2
     Name './example_faculty_it.picoc_mon',
 4
       Block
         Name 'faculty.6',
 6
 7
8
           // Assign(Name('res'), Num('1'))
           Exp(Num('1'))
 9
           Assign(Stackframe(Num('1')), Stack(Num('1')))
10
           // While(Num('1'), [])
11
           Exp(GoTo(Name('condition_check.5')))
12
         ],
13
       Block
14
         Name 'condition_check.5',
15
16
           // IfElse(Num('1'), [], [])
17
           Exp(Num('1')),
           IfElse
18
19
             Stack
20
                Num '1',
21
             Ε
22
                GoTo(Name('while_branch.4'))
23
             ],
24
             25
                GoTo(Name('while_after.1'))
26
27
         ],
28
       Block
29
         Name 'while_branch.4',
30
31
           // If(Atom(Name('n'), Eq('=='), Num('1')), [])
32
           // IfElse(Atom(Name('n'), Eq('=='), Num('1')), [], [])
33
           Exp(Stackframe(Num('0')))
34
           Exp(Num('1'))
35
           Exp(Atom(Stack(Num('2')), Eq('=='), Stack(Num('1')))),
36
           IfElse
37
             Stack
38
                Num '1',
39
40
                GoTo(Name('if.3'))
41
             ],
42
             [
43
                GoTo(Name('if_else_after.2'))
44
             ]
         ],
```

```
Block
47
         Name 'if.3',
48
           // Return(Name('res'))
           Exp(Stackframe(Num('1')))
           Return(Stack(Num('1')))
51
         ],
52
53
       Block
54
         Name 'if_else_after.2',
55
56
           // Assign(Name('res'), BinOp(Name('n'), Mul('*'), Name('res')))
57
           Exp(Stackframe(Num('0')))
58
           Exp(Stackframe(Num('1')))
59
           Exp(BinOp(Stack(Num('2')), Mul('*'), Stack(Num('1'))))
60
           Assign(Stackframe(Num('1')), Stack(Num('1')))
           // Assign(Name('n'), BinOp(Name('n'), Sub('-'), Num('1')))
61
62
           Exp(Stackframe(Num('0')))
63
           Exp(Num('1'))
64
           Exp(BinOp(Stack(Num('2')), Sub('-'), Stack(Num('1'))))
65
           Assign(Stackframe(Num('0')), Stack(Num('1')))
66
           Exp(GoTo(Name('condition_check.5')))
67
         ],
68
       Block
69
         Name 'while_after.1',
71
           Return(Empty())
72
         ],
73
       Block
         Name 'main.0',
74
75
           StackMalloc(Num('2'))
           Exp(Num('4'))
           NewStackframe(Name('faculty'), GoTo(Name('addr@next_instr')))
           Exp(GoTo(Name('faculty.6')))
80
           RemoveStackframe()
81
           Exp(ACC)
           Exp(Call(Name('print'), [Stack(Num('1'))]))
82
83
           Return(Empty())
84
         ]
85
     ]
```

Code 2.9: Pico C-ANF Pass für Codebespiel

#### 2.3.1.4 RETI-Blocks Pass

#### 2.3.1.4.1 Aufgabe

Die Aufgabe des RETI-Blocks Passes ist es die Statements in den Blöcken, die durch PicoC-Knoten im Abstrakten Syntaxbaum der Sprache  $L_{PicoC\_ANF}$  dargestellt sind durch ihren entsprechenden RETI-Knoten zu ersetzen.

#### 2.3.1.4.2 Abstrakte Grammatik

Die Abstrakte Grammatik 2.3.4 der Sprache  $L_{RETI\_Blocks}$  ist verglichen mit der Abstrakten Grammatik 2.3.3 der Sprache  $L_{PicoC\_ANF}$  stark verändert, denn der Großteil der PicoC-Knoten wird in diesem Pass durch entsprechende RETI-Knoten ersetzt. Die einzigen verbleibenden PicoC-Knoten sind Exp(GoTo(str)),

Block(Name(str), (instr)\*) und File(Name(str), (block)\*), da das gesamte Konzept mit den Blöcken erst im RETI-Pass in Unterkapitel 2.3.8 aufgelöst wird.

```
ACC()
                           IN1()
                                        IN2()
                                                    PC()
                                                                SP()
                                                                           BAF()
                                                                                                              L_{-}RETI
reg
        ::=
              CS() \mid DS()
               Reg(\langle reg \rangle) \mid Num(str)
arg
        ::=
                                   | Lt() | LtE() | Gt() | GtE()
rel
               Eq()
                     |NEq()|
               Always() \mid NOp()
                                       Sub() \mid Subi() \mid Mult() \mid Multi()
               Add()
                           Addi()
op
                          Divi() \mid Mod() \mid Modi() \mid Oplus() \mid Oplusi()
               Div()
               Or() \mid Ori() \mid And() \mid Andi()
               Load() | Loadin() | Loadi() | Store() | Storein() | Move()
instr
               Instr(\langle op \rangle, \langle arg \rangle +) \mid Jump(\langle rel \rangle, Num(str)) \mid Int(Num(str))
               RTI() \quad | \quad Call(Name('print'), \langle reg \rangle) \quad | \quad Call(Name('input'), \langle reg \rangle)
               SingleLineComment(str, str)
               Instr(Loadi(), [Reg(Acc()), GoTo(Name(str))]) \mid Jump(Eq(), GoTo(Name(str)))
               Exp(GoTo(str))
instr
                                                                                                              L_{-}PicoC
block
               Block(Name(str), \langle instr \rangle *)
        ::=
               File(Name(str), \langle block \rangle *)
file
```

Grammatik 2.3.4: Abstrakte Grammatik der Sprache  $L_{RETI\_Blocks}$ 

#### 2.3.1.4.3 Codebeispiel

In Code 2.10 sieht man den Abstrakten Syntaxbaum des RETI-Blocks Passes für das aus Unterkapitel 2.6 weitergeführte Beispiel, indem die Statements, die durch entsprechende PicoC-Knoten im Abstrakten Syntaxbaum der Sprache  $L_{PicoC\_ANF}^{25}$  repräsentiert waren nun durch ihre entsprechennden RETI-Knoten ersetzt werden.

```
File
 2
    Name './example_faculty_it.reti_blocks',
     Γ
      Block
         Name 'faculty.6',
 6
7
8
9
           # // Assign(Name('res'), Num('1'))
           # Exp(Num('1'))
           SUBI SP 1;
10
           LOADI ACC 1;
11
           STOREIN SP ACC 1;
12
           # Assign(Stackframe(Num('1')), Stack(Num('1')))
13
           LOADIN SP ACC 1;
14
           STOREIN BAF ACC -3;
15
           ADDI SP 1;
16
           # // While(Num('1'), [])
17
           # Exp(GoTo(Name('condition_check.5')))
18
           Exp(GoTo(Name('condition_check.5')))
19
         ],
20
       Block
21
         Name 'condition_check.5',
         Γ
```

<sup>&</sup>lt;sup>25</sup>Beschrieben durch die Grammatik 2.3.3.

```
# // IfElse(Num('1'), [], [])
24
           # Exp(Num('1'))
25
           SUBI SP 1;
26
           LOADI ACC 1;
27
           STOREIN SP ACC 1;
28
           # IfElse(Stack(Num('1')), [], [])
29
           LOADIN SP ACC 1;
30
           ADDI SP 1;
31
           JUMP== GoTo(Name('while_after.1'));
32
           Exp(GoTo(Name('while_branch.4')))
33
         ],
34
       Block
         Name 'while_branch.4',
36
37
           # // If(Atom(Name('n'), Eq('=='), Num('1')), [])
38
           # // IfElse(Atom(Name('n'), Eq('=='), Num('1')), [], [])
39
           # Exp(Stackframe(Num('0')))
40
           SUBI SP 1;
41
           LOADIN BAF ACC -2;
42
           STOREIN SP ACC 1;
43
           # Exp(Num('1'))
44
           SUBI SP 1;
45
           LOADI ACC 1;
46
           STOREIN SP ACC 1;
47
           LOADIN SP ACC 2;
48
           LOADIN SP IN2 1;
49
           SUB ACC IN2;
50
           JUMP== 3;
51
           LOADI ACC 0;
52
           JUMP 2;
53
           LOADI ACC 1;
54
           STOREIN SP ACC 2;
55
           ADDI SP 1;
56
           # IfElse(Stack(Num('1')), [], [])
57
           LOADIN SP ACC 1;
58
           ADDI SP 1;
59
           JUMP== GoTo(Name('if_else_after.2'));
60
           Exp(GoTo(Name('if.3')))
61
         ],
62
       Block
63
         Name 'if.3',
64
         Ε
65
           # // Return(Name('res'))
66
           # Exp(Stackframe(Num('1')))
67
           SUBI SP 1;
68
           LOADIN BAF ACC -3;
69
           STOREIN SP ACC 1;
70
           # Return(Stack(Num('1')))
71
           LOADIN SP ACC 1;
72
           ADDI SP 1;
73
           LOADIN BAF PC -1;
74
        ],
75
       Block
76
         Name 'if_else_after.2',
           # // Assign(Name('res'), BinOp(Name('n'), Mul('*'), Name('res')))
           # Exp(Stackframe(Num('0')))
```

```
SUBI SP 1;
81
           LOADIN BAF ACC -2;
82
           STOREIN SP ACC 1;
83
           # Exp(Stackframe(Num('1')))
           SUBI SP 1;
85
           LOADIN BAF ACC -3;
86
           STOREIN SP ACC 1;
87
           # Exp(BinOp(Stack(Num('2')), Mul('*'), Stack(Num('1'))))
88
           LOADIN SP ACC 2;
89
           LOADIN SP IN2 1;
90
           MULT ACC IN2;
91
           STOREIN SP ACC 2;
92
           ADDI SP 1;
93
           # Assign(Stackframe(Num('1')), Stack(Num('1')))
94
           LOADIN SP ACC 1;
95
           STOREIN BAF ACC -3;
96
           ADDI SP 1;
97
           # // Assign(Name('n'), BinOp(Name('n'), Sub('-'), Num('1')))
98
           # Exp(Stackframe(Num('0')))
99
           SUBI SP 1;
00
           LOADIN BAF ACC -2;
L01
           STOREIN SP ACC 1;
102
           # Exp(Num('1'))
103
           SUBI SP 1;
104
           LOADI ACC 1;
105
           STOREIN SP ACC 1;
106
           # Exp(BinOp(Stack(Num('2')), Sub('-'), Stack(Num('1'))))
107
           LOADIN SP ACC 2;
108
           LOADIN SP IN2 1;
109
           SUB ACC IN2;
110
           STOREIN SP ACC 2;
111
           ADDI SP 1;
           # Assign(Stackframe(Num('0')), Stack(Num('1')))
113
           LOADIN SP ACC 1;
           STOREIN BAF ACC -2;
114
115
           ADDI SP 1;
116
           # Exp(GoTo(Name('condition_check.5')))
117
           Exp(GoTo(Name('condition_check.5')))
118
         ],
119
       Block
120
         Name 'while_after.1',
121
         Ε
           # Return(Empty())
123
           LOADIN BAF PC -1;
124
         ],
L25
       Block
126
         Name 'main.0',
L27
128
           # StackMalloc(Num('2'))
129
           SUBI SP 2;
130
           # Exp(Num('4'))
131
           SUBI SP 1;
132
           LOADI ACC 4;
133
           STOREIN SP ACC 1;
134
           # NewStackframe(Name('faculty'), GoTo(Name('addr@next_instr')))
L35
           MOVE BAF ACC;
           ADDI SP 3;
```

```
MOVE SP BAF;
           SUBI SP 4;
.39
           STOREIN BAF ACC 0;
           LOADI ACC GoTo(Name('addr@next_instr'));
           ADD ACC CS;
           STOREIN BAF ACC -1;
43
           # Exp(GoTo(Name('faculty.6')))
L44
           Exp(GoTo(Name('faculty.6')))
L45
           # RemoveStackframe()
L46
           MOVE BAF IN1;
L47
           LOADIN IN1 BAF 0;
48
           MOVE IN1 SP;
49
           # Exp(ACC)
150
           SUBI SP 1;
151
           STOREIN SP ACC 1;
152
           LOADIN SP ACC 1;
153
           ADDI SP 1;
154
           CALL PRINT ACC;
L55
            # Return(Empty())
156
           LOADIN BAF PC -1;
L57
158
```

Code 2.10: RETI-Blocks Pass für Codebespiel

# Anmerkung 9

Wenn der Abstrakte Syntaxbaum ausgegeben wird, ist die Darstellung nicht auschließlich in Abstrakter Syntax, da die RETI-Knoten aus bereits im Unterkapitel 2.2.5.6 vermitteltem Grund in Konkretter Syntax ausgeben werden.

#### 2.3.1.5 RETI-Patch Pass

#### 2.3.1.5.1 Aufgabe

Die Aufgabe des RETI-Patch Passes ist das Ausbessern (engl. to patch) des Abstrakten Syntaxbaumes, durch:

- das Einfügen eines start.<nummer>-Blockes, welcher ein GoTo(Name('main')) zur main-Funktion enthält, wenn in manchen Fällen die main-Funktion nicht die erste Funktion ist und daher am Anfang zur main-Funktion gesprungen werden muss.
- das Entfernen von GoTo()'s, deren Sprung nur eine Adresse weiterspringen würde.
- das Voranstellen von RETI-Knoten, die vor jeder Division Instr(Div(), args) prüfen, ob, nicht durch 0 geteilt wird.<sup>26</sup>
- das Überprüfen darauf, ob bestimmte Immediates Im(str) in Befehlen, wie z.B. Jump(rel, Im(str)), Instr(Loadin(), [reg, reg, Im(str)]), Instr(Loadi(), [reg, Im(str)]) usw. kleiner -2<sup>21</sup> oder größer 2<sup>21</sup> 1 sind. Im Fall dessen, dass es so ist, muss der gewünschte Zahlenwert durch Bitshiften und Anwenden von bitweise Oder berechnet werden. Im Fall, dessen, dass der Immediate allerdings kleiner -(2<sup>31</sup>) oder größer 2<sup>31</sup> 1 ist, wird eine Fehlermeldung TooLargeLiteral ausgegeben.

<sup>&</sup>lt;sup>26</sup>Das fällt unter die Themenbereiche des Bachelorprojekts und wird daher nicht genauer erläutert.

#### 2.3.1.5.2 Abstrakte Grammatik

Die Abstrakte Grammatik 2.3.5 der Sprache  $L_{RETI\_Patch}$  ist im Vergleich zur Abstrakten Grammatik 2.3.4 der Sprache  $L_{RETI\_Blocks}$  kaum verändert. Es muss nur ein Knoten Exit() hinzugefügt werden, der im Falle einer Division durch 0 die Ausführung des Programs beendet.

```
\mid BAF()
               ACC() \mid IN1()
                                    | IN2() | PC()
                                                                SP()
                                                                                                               L\_RETI
reg
               CS() \mid DS()
               Reg(\langle reg \rangle) \mid Num(str)
arg
               Eq() \mid NEq() \mid Lt() \mid LtE() \mid Gt() \mid GtE()
rel
               Always() \mid NOp()
                                        Sub() \mid Subi() \mid Mult() \mid Multi()
                          Addi()
               Add()
op
                          Divi() \mid Mod() \mid Modi() \mid Oplus() \mid Oplusi()
               Div()
               Or() \mid Ori() \mid And() \mid Andi()
               Load() \mid Loadin() \mid Loadi() \mid Store() \mid Storein() \mid Move()
               Instr(\langle op \rangle, \langle arg \rangle +) \mid Jump(\langle rel \rangle, Num(str)) \mid Int(Num(str))
instr
               RTI() \mid Call(Name('print'), \langle reg \rangle) \mid Call(Name('input'), \langle reg \rangle)
               SingleLineComment(str, str)
               Instr(Loadi(), [Reg(Acc()), GoTo(Name(str))]) \mid Jump(Eq(), GoTo(Name(str)))
               Exp(GoTo(str)) \mid Exit(Num(str))
                                                                                                               L_{-}PicoC
instr
        ::=
block
               Block(Name(str), \langle instr \rangle *)
        ::=
file
               File(Name(str), \langle block \rangle *)
        ::=
```

Grammatik 2.3.5: Abstrakte Grammatik der Sprache L<sub>RETI\_Patch</sub>

#### 2.3.1.5.3 Codebeispiel

In Code 2.11 sieht man den Abstrakten Syntaxbaum des PiocC-Patch Passes für das aus Unterkapitel 2.6 weitergeführte Beispiel. Durch den RETI-Patch Pass wurde hier ein start. <nummer>-Block<sup>27</sup> eingesetzt, da die main-Funktion nicht die erste Funktion ist. Des Weiteren wurden durch diesen Pass einzelne GoTo(Name(str))-Statements entfernt<sup>28</sup>, die nur einen Sprung um eine Position entsprochen hätten.

```
File
    Name './example_faculty_it.reti_patch',
     Γ
       Block
         Name 'start.7',
 7
8
9
           # // Exp(GoTo(Name('main.0')))
           Exp(GoTo(Name('main.0')))
         ],
10
       Block
11
         Name 'faculty.6',
12
13
           # // Assign(Name('res'), Num('1'))
           # Exp(Num('1'))
15
           SUBI SP 1;
16
           LOADI ACC 1:
           STOREIN SP ACC 1;
17
18
           # Assign(Stackframe(Num('1')), Stack(Num('1')))
```

 $<sup>^{27}\</sup>mathrm{Dieser}$ Block wurde im Code 2.8 markiert.

<sup>&</sup>lt;sup>28</sup>Diese entfernten GoTo(Name(str))'s' wurden ebenfalls im Code 2.8 markiert.

```
LOADIN SP ACC 1;
20
           STOREIN BAF ACC -3;
21
           ADDI SP 1;
           # // While(Num('1'), [])
           # Exp(GoTo(Name('condition_check.5')))
24
           # // not included Exp(GoTo(Name('condition_check.5')))
25
         ],
26
       Block
27
         Name 'condition_check.5',
28
29
           # // IfElse(Num('1'), [], [])
30
           # Exp(Num('1'))
           SUBI SP 1;
32
           LOADI ACC 1;
33
           STOREIN SP ACC 1;
34
           # IfElse(Stack(Num('1')), [], [])
35
           LOADIN SP ACC 1;
36
           ADDI SP 1;
37
           JUMP== GoTo(Name('while_after.1'));
38
           # // not included Exp(GoTo(Name('while_branch.4')))
39
         ],
40
       Block
41
         Name 'while_branch.4',
42
43
           # // If(Atom(Name('n'), Eq('=='), Num('1')), [])
44
           # // IfElse(Atom(Name('n'), Eq('=='), Num('1')), [], [])
45
           # Exp(Stackframe(Num('0')))
46
           SUBI SP 1;
47
           LOADIN BAF ACC -2;
48
           STOREIN SP ACC 1;
           # Exp(Num('1'))
50
           SUBI SP 1;
51
           LOADI ACC 1;
52
           STOREIN SP ACC 1;
53
           LOADIN SP ACC 2;
54
           LOADIN SP IN2 1;
55
           SUB ACC IN2;
56
           JUMP== 3;
57
           LOADI ACC 0;
58
           JUMP 2;
59
           LOADI ACC 1;
60
           STOREIN SP ACC 2;
61
           ADDI SP 1;
62
           # IfElse(Stack(Num('1')), [], [])
63
           LOADIN SP ACC 1;
64
           ADDI SP 1;
65
           JUMP== GoTo(Name('if_else_after.2'));
66
           # // not included Exp(GoTo(Name('if.3')))
67
         ],
68
       Block
69
         Name 'if.3',
70
71
           # // Return(Name('res'))
           # Exp(Stackframe(Num('1')))
           SUBI SP 1;
           LOADIN BAF ACC -3;
           STOREIN SP ACC 1;
```

```
76
           # Return(Stack(Num('1')))
77
           LOADIN SP ACC 1;
78
           ADDI SP 1;
79
           LOADIN BAF PC -1;
80
         ],
81
       Block
82
         Name 'if_else_after.2',
83
84
           # // Assign(Name('res'), BinOp(Name('n'), Mul('*'), Name('res')))
85
           # Exp(Stackframe(Num('0')))
86
           SUBI SP 1;
87
           LOADIN BAF ACC -2;
88
           STOREIN SP ACC 1;
           # Exp(Stackframe(Num('1')))
89
90
           SUBI SP 1;
91
           LOADIN BAF ACC -3;
92
           STOREIN SP ACC 1;
93
           # Exp(BinOp(Stack(Num('2')), Mul('*'), Stack(Num('1'))))
94
           LOADIN SP ACC 2:
95
           LOADIN SP IN2 1;
96
           MULT ACC IN2;
97
           STOREIN SP ACC 2;
98
           ADDI SP 1;
99
           # Assign(Stackframe(Num('1')), Stack(Num('1')))
100
           LOADIN SP ACC 1;
01
           STOREIN BAF ACC -3:
           ADDI SP 1:
102
103
           # // Assign(Name('n'), BinOp(Name('n'), Sub('-'), Num('1')))
104
           # Exp(Stackframe(Num('0')))
           SUBI SP 1;
106
           LOADIN BAF ACC -2;
L07
           STOREIN SP ACC 1;
108
           # Exp(Num('1'))
109
           SUBI SP 1;
L10
           LOADI ACC 1;
111
           STOREIN SP ACC 1;
12
           # Exp(BinOp(Stack(Num('2')), Sub('-'), Stack(Num('1'))))
113
           LOADIN SP ACC 2;
114
           LOADIN SP IN2 1;
115
           SUB ACC IN2;
116
           STOREIN SP ACC 2;
17
           ADDI SP 1;
18
           # Assign(Stackframe(Num('0')), Stack(Num('1')))
L19
           LOADIN SP ACC 1;
120
           STOREIN BAF ACC -2;
l21
           ADDI SP 1;
122
           # Exp(GoTo(Name('condition_check.5')))
123
           Exp(GoTo(Name('condition_check.5')))
124
         ],
L25
       Block
L26
         Name 'while_after.1',
L27
128
           # Return(Empty())
L29
           LOADIN BAF PC -1;
130
         ],
l31
       Block
132
         Name 'main.0',
```

```
Γ
            # StackMalloc(Num('2'))
            SUBI SP 2;
            # Exp(Num('4'))
            SUBI SP 1;
137
138
            LOADI ACC 4;
139
            STOREIN SP ACC 1;
L40
            # NewStackframe(Name('faculty'), GoTo(Name('addr@next_instr')))
L41
            MOVE BAF ACC;
142
            ADDI SP 3;
143
           MOVE SP BAF;
44
            SUBI SP 4;
45
            STOREIN BAF ACC 0;
146
           LOADI ACC GoTo(Name('addr@next_instr'));
47
            ADD ACC CS;
148
            STOREIN BAF ACC -1;
149
            # Exp(GoTo(Name('faculty.6')))
150
            Exp(GoTo(Name('faculty.6')))
L51
            # RemoveStackframe()
152
            MOVE BAF IN1;
153
           LOADIN IN1 BAF O;
154
            MOVE IN1 SP;
155
            # Exp(ACC)
156
            SUBI SP 1;
157
            STOREIN SP ACC 1;
158
            LOADIN SP ACC 1;
            ADDI SP 1;
159
160
            CALL PRINT ACC;
161
            # Return(Empty())
62
            LOADIN BAF PC -1;
163
         ]
164
     ]
```

Code 2.11: RETI-Patch Pass für Codebespiel

#### 2.3.1.6 RETI Pass

#### 2.3.1.6.1 Aufgabe

Die Aufgabe des RETI-Patch Passes ist es die GoTo(Name(str))-Knoten in den den Knoten Instr(Loadi(), [reg, GoTo(Name(str))]), Jump(Eq(), GoTo(Name(str))) und Exp(GoTo(Name(str))) durch eine entsprechende Adresse zu ersetzen, die entsprechende Distanz oder einen entsprechenden Sprungbefehl mit passender Distanz Jump(Always(), Im(str(distance))). Die Distanz- und Adressberechnung wird in Unterkapitel ?? genauer mit Formeln erklärt.

#### 2.3.1.6.2 Konkrette und Abstrakte Grammatik

Die Abstrakte Grammatik 2.3.8 der Sprache  $L_{RETI}$  hat im Vergleich zur Abstrakten Grammatik 2.3.5 der Sprache  $L_{RETI\_Patch}$  nur noch auschließlich **RETI-Knoten**. Alle **RETI-Knoten** stehen nun in einem Program(Name(str), instr)-Knoten.

Ausgegeben wird der finale Maschinencode allerdings in Konkretter Syntax, die durch die Konkretten Grammatiken 2.3.6 und 2.3.7 für jeweils die Lexikalische und Syntaktische Analyse beschrieben wird. Der Grund, warum die Konkrette Grammatik der Sprache  $L_{RETI}$  auch nochmal in einen Teil für die Lexikalische und Syntaktische Analyse unterteilt ist, hat den Grund, dass für die Bachelorarbeit zum

Testen des PicoC-Compilers ein RETI-Interpreter implementiert wurde, der den RETI-Code lexen und parsen muss, um ihn später interpretieren zu können.

```
"6"
dig\_no\_0
                                                                       L_Program
                 "7"
                          "8"
                                  "<sub>9"</sub>
                 "0"
dig_with_0
            ::=
                          dig\_no\_0
                 "0"
                         dig\_no\_0 dig\_with\_0* | "-"dig\_no\_0*
num
            ::=
                 "a"..."Z"
letter
            ::=
                 letter(letter \mid dig\_with\_0 \mid \_)*
name
                 "ACC"
                              "IN1" | "IN2"
                                                |"PC""|"SP"
reg
            ::=
                             "CS" | "DS"
                 "BAF"
arg
            ::=
                 reg
                         num
                            "!=" | "<" | "<=" | ">"
rel
            ::=
                  ">="
                            "\_NOP"
```

Grammatik 2.3.6: Konkrette Grammatik der Sprache  $L_{RETI}$  für die Lexikalische Analyse in EBNF

```
"ADDI" reg num |
                                               "SUB" \ reg \ arg
                                                                     L_Program
instr
        ::=
             "ADD" reg arg
             "SUBI" reg num | "MULT" reg arg | "MULTI" reg num
             "DIV" reg arg | "DIVI" reg num | "MOD" reg arg
             "MODI" reg num | "OPLUS" reg arg | "OPLUSI" reg num
             "OR" reg arg | "ORI" reg num
             "AND" reg arg | "ANDI" reg num
             "LOAD" reg num | "LOADIN" arg arg num
             "LOADI" reg num
             "STORE" reg num | "STOREIN" arg argnum
             "MOVE" req req
             "JUMP"rel\ num\ |\ INT\ num\ |\ RTI
             "CALL" "INPUT" reg | "CALL" "PRINT" reg
             name (instr";")*
program
        ::=
```

Grammatik 2.3.7: Konkrette Grammatik der Sprache L<sub>RETI</sub> für die Syntaktische Analyse in EBNF

```
ACC() \mid IN1()
                                                                                                                 L\_RETI
                                           IN2()
                                                        PC()
                                                                   SP()
                                                                              BAF()
reg
                   CS()
                             DS()
                   Reg(\langle reg \rangle) \mid Num(str)
            ::=
arg
rel
                             NEq() \mid Lt()
                                                    LtE()
                                                                Gt() \mid GtE()
                  Always() \mid NOp()
                   Add()
                              Addi()
                                           Sub() \mid Subi() \mid Mult() \mid Multi()
op
                             Divi() \mid Mod() \mid Modi() \mid Oplus() \mid Oplusi()
                   Div()
                   Or() \mid Ori() \mid And() \mid Andi()
                            | Loadin() | Loadi() | Store() | Storein() | Move()
                   Load()
                  Instr(\langle op \rangle, \langle arg \rangle +) \mid Jump(\langle rel \rangle, Num(str)) \mid Int(Num(str))
instr
                   RTI() \mid Call(Name('print'), \langle reg \rangle) \mid Call(Name('input'), \langle reg \rangle)
                  SingleLineComment(str, str)
                   Instr(Loadi(), [Reg(Acc()), GoTo(Name(str))]) \mid Jump(Eq(), GoTo(Name(str)))
program
            ::=
                   Program(Name(str), \langle instr \rangle *)
                   Exp(GoTo(str)) \mid Exit(Num(str))
                                                                                                                 L-PicoC
instr
            ::=
block
                   Block(Name(str), \langle instr \rangle *)
            ::=
                   File(Name(str), \langle block \rangle *)
file
            ::=
```

Grammatik 2.3.8: Abstrakte Grammatik der Sprache  $L_{RETI}$ 

#### 2.3.1.6.3 Codebeispiel

Nach dem RETI-Pass ist das Programm komplett in RETI-Knoten übersetzt, die allerdings in ihrer Konkretten Syntax ausgegeben werden, wie in Code 2.12 zu sehen ist. Es gibt keine Blöcke mehr und die RETI-Befehle in diesen Blöcken wurden zusammengesetzt, wie sie in den Blöcken angeordnet waren. Die letzten Nicht-RETI-Befehle oder RETI-Befehle, die nicht auschließlich aus RETI-Ausdrücken bestehen<sup>29</sup>, die sich in den Blöcken befunden haben, wurden durch RETI-Befehle ersetzt.

Der Program(Name(str), instr)-Knoten, indem alle RETI-Knoten stehen gibt alleinig die RETI-Knoten, die er beinhaltet aus und fügt ansonsten nichts hinzu, wodurch der Abstrakte Syntaxbaum, wenn er in eine Datei ausgegeben wird, direkt RETI-Code in menschenlesbarer Repräsentation erzeugt.

```
# // Exp(GoTo(Name('main.0')))
 2 JUMP 67;
 3 # // Assign(Name('res'), Num('1'))
 4 # Exp(Num('1'))
 5 SUBI SP 1;
 6 LOADI ACC 1;
 7 STOREIN SP ACC 1;
 8 # Assign(Stackframe(Num('1')), Stack(Num('1')))
 9 LOADIN SP ACC 1;
10 STOREIN BAF ACC -3;
11 ADDI SP 1;
12 # // While(Num('1'), [])
13 # Exp(GoTo(Name('condition_check.5')))
14 # // not included Exp(GoTo(Name('condition_check.5')))
15 # // IfElse(Num('1'), [], [])
16 # Exp(Num('1'))
17 SUBI SP 1;
18 LOADI ACC 1;
```

<sup>&</sup>lt;sup>29</sup>Wie z.B. LOADI ACC GoTo(Name('addr@next\_instr')), Exp(GoTo(Name('main.0'))) und JUMP== GoTo(Name('if\_else\_after.2')).

```
19 STOREIN SP ACC 1;
20 # IfElse(Stack(Num('1')), [], [])
21 LOADIN SP ACC 1;
22 ADDI SP 1;
23 JUMP== 54;
24 # // not included Exp(GoTo(Name('while_branch.4')))
25 # // If(Atom(Name('n'), Eq('=='), Num('1')), [])
26 # // IfElse(Atom(Name('n'), Eq('=='), Num('1')), [], [])
27 # Exp(Stackframe(Num('0')))
28 SUBI SP 1;
29 LOADIN BAF ACC -2;
30 STOREIN SP ACC 1;
31 # Exp(Num('1'))
32 SUBI SP 1;
33 LOADI ACC 1;
34 STOREIN SP ACC 1;
35 LOADIN SP ACC 2;
36 LOADIN SP IN2 1;
37 SUB ACC IN2:
38 JUMP== 3;
39 LOADI ACC 0;
40 JUMP 2;
41 LOADI ACC 1;
42 STOREIN SP ACC 2;
43 ADDI SP 1;
44 # IfElse(Stack(Num('1')), [], [])
45 LOADIN SP ACC 1;
46 ADDI SP 1;
47 JUMP== 7;
48 # // not included Exp(GoTo(Name('if.3')))
49 # // Return(Name('res'))
50 # Exp(Stackframe(Num('1')))
51 SUBI SP 1;
52 LOADIN BAF ACC -3;
53 STOREIN SP ACC 1;
54 # Return(Stack(Num('1')))
55 LOADIN SP ACC 1;
56 ADDI SP 1;
57 LOADIN BAF PC -1;
58 # // Assign(Name('res'), BinOp(Name('n'), Mul('*'), Name('res')))
59 # Exp(Stackframe(Num('0')))
60 SUBI SP 1;
61 LOADIN BAF ACC -2;
62 STOREIN SP ACC 1;
63 # Exp(Stackframe(Num('1')))
64 SUBI SP 1;
65 LOADIN BAF ACC -3;
66 STOREIN SP ACC 1;
67 # Exp(BinOp(Stack(Num('2')), Mul('*'), Stack(Num('1'))))
68 LOADIN SP ACC 2;
69 LOADIN SP IN2 1;
70 MULT ACC IN2;
71 STOREIN SP ACC 2;
72 ADDI SP 1;
73 # Assign(Stackframe(Num('1')), Stack(Num('1')))
74 LOADIN SP ACC 1;
75 STOREIN BAF ACC -3;
```

```
76 ADDI SP 1:
77 # // Assign(Name('n'), BinOp(Name('n'), Sub('-'), Num('1')))
78 # Exp(Stackframe(Num('0')))
79 SUBI SP 1;
80 LOADIN BAF ACC -2;
81 STOREIN SP ACC 1;
82 # Exp(Num('1'))
83 SUBI SP 1;
84 LOADI ACC 1;
85 STOREIN SP ACC 1;
86 # Exp(BinOp(Stack(Num('2')), Sub('-'), Stack(Num('1'))))
87 LOADIN SP ACC 2;
88 LOADIN SP IN2 1;
89 SUB ACC IN2;
90 STOREIN SP ACC 2;
91 ADDI SP 1;
92 # Assign(Stackframe(Num('0')), Stack(Num('1')))
93 LOADIN SP ACC 1;
94 STOREIN BAF ACC -2;
95 ADDI SP 1;
96 # Exp(GoTo(Name('condition_check.5')))
97 JUMP -58;
98 # Return(Empty())
99 LOADIN BAF PC -1;
00 # StackMalloc(Num('2'))
01 SUBI SP 2;
102 # Exp(Num('4'))
103 SUBI SP 1;
104 LOADI ACC 4;
05 STOREIN SP ACC 1;
06 # NewStackframe(Name('faculty'), GoTo(Name('addr@next_instr')))
107 MOVE BAF ACC;
08 ADDI SP 3;
109 MOVE SP BAF;
110 SUBI SP 4;
11 STOREIN BAF ACC 0;
12 LOADI ACC 80;
13 ADD ACC CS;
14 STOREIN BAF ACC -1;
115 # Exp(GoTo(Name('faculty.6')))
116 JUMP -78;
17 # RemoveStackframe()
18 MOVE BAF IN1;
19 LOADIN IN1 BAF 0;
20 MOVE IN1 SP;
21 # Exp(ACC)
22 SUBI SP 1;
23 STOREIN SP ACC 1;
24 LOADIN SP ACC 1;
125 ADDI SP 1;
26 CALL PRINT ACC;
27 # Return(Empty())
128 LOADIN BAF PC -1;
```

Code 2.12: RETI Pass für Codebespiel

# Literatur

### Online

- ANSI C grammar (Lex). URL: https://www.lysator.liu.se/c/ANSI-C-grammar-1.html (besucht am 29.07.2022).
- ANSI C grammar (Yacc). URL: http://www.quut.com/c/ANSI-C-grammar-y.html (besucht am 29.07.2022).
- ANTLR. URL: https://www.antlr.org/ (besucht am 31.07.2022).
- C Operator Precedence cppreference.com. URL: https://en.cppreference.com/w/c/language/operator\_precedence (besucht am 27.04.2022).
- clang: C++ Compiler. URL: http://clang.org/ (besucht am 29.07.2022).
- Clockwise/Spiral Rule. URL: https://c-faq.com/decl/spiral.anderson.html (besucht am 29.07.2022).
- Developers, Inkscape Website. *Draw Freely Inkscape*. URL: https://inkscape.org/ (besucht am 03.08.2022).
- GCC, the GNU Compiler Collection GNU Project. URL: https://gcc.gnu.org/ (besucht am 13.07.2022).
- Grammar Reference Lark documentation. URL: https://lark-parser.readthedocs.io/en/latest/grammar.html (besucht am 31.07.2022).
- Grammar: The language of languages (BNF, EBNF, ABNF and more). URL: https://matt.might.net/articles/grammars-bnf-ebnf/ (besucht am 30.07.2022).
- Parsing Expressions · Crafting Interpreters. URL: https://www.craftinginterpreters.com/parsing-expressions.html (besucht am 09.07.2022).
- Variablen in C und C++, Deklaration und Definition Coder-Welten.de. URL: https://www.coder-welten.de/einstieg/variablen-in-c-3.html (besucht am 11.08.2022).
- Welcome to Lark's documentation! Lark documentation. URL: https://lark-parser.readthedocs.io/en/latest/ (besucht am 31.07.2022).
- What is the difference between function prototype and function signature? SoloLearn. URL: https://www.sololearn.com/Discuss/171026/what-is-the-difference-between-function-prototype-and-function-signature/ (besucht am 18.07.2022).

### Bücher

- G. Siek, Jeremy. Course Webpage for Compilers (P423, P523, E313, and E513). 28. Jan. 2022. URL: https://iucompilercourse.github.io/IU-Fall-2021/ (besucht am 28.01.2022).
- LeFever, Lee. The Art of Explanation: Making your Ideas, Products, and Services Easier to Understand. 1. Aufl. Wiley, 20. Nov. 2012.

## Artikel

• Earley, Jay. "An efficient context-free parsing". In: 13 (1968). URL: https://web.archive.org/web/20040708052627/http://www-2.cs.cmu.edu/afs/cs.cmu.edu/project/cmt-55/lti/Courses/711/Class-notes/p94-earley.pdf (besucht am 10.08.2022).

# Vorlesungen

- Bast, Prof. Dr. Hannah. "Programmieren in C". Vorlesung. Vorlesung. Universität Freiburg, 2020. URL: https://ad-wiki.informatik.uni-freiburg.de/teaching/ProgrammierenCplusplusSS2020 (besucht am 09.07.2022).
- Nebel, Prof. Dr. Bernhard. "Theoretische Informatik". Vorlesung. Vorlesung. Universität Freiburg, 2020.
   URL: http://gki.informatik.uni-freiburg.de/teaching/ss20/info3/index\_de.html (besucht am 09.07.2022).
- Scholl, Prof. Dr. Christoph. "Betriebssysteme". Vorlesung. Vorlesung. Universität Freiburg, 2020. URL: https://abs.informatik.uni-freiburg.de/src/teach\_main.php?id=157 (besucht am 09.07.2022).
- "Technische Informatik". Vorlesung. Vorlesung. Universität Freiburg, 3. Aug. 2022. (Besucht am 03.08.2022).
- Scholl, Prof. Dr. Philipp. "Einführung in Embedded Systems". Vorlesung. Vorlesung. Universität Freiburg, 2021. URL: https://earth.informatik.uni-freiburg.de/uploads/es-2122/ (besucht am 09.07.2022).
- Thiemann, Prof. Dr. Peter. "Compilerbau". Vorlesung. Vorlesung. Universität Freiburg, 2021. URL: http://proglang.informatik.uni-freiburg.de/teaching/compilerbau/2021ws/ (besucht am 09.07.2022).
- — "Einführung in die Programmierung". Vorlesung. Vorlesung. Universität Freiburg, 2018. URL: http://proglang.informatik.uni-freiburg.de/teaching/info1/2018/ (besucht am 09.07.2022).

# Sonstige Quellen

- Lark a parsing toolkit for Python. 26. Apr. 2022. URL: https://github.com/lark-parser/lark (besucht am 28.04.2022).
- Naming convention (programming). In: Wikipedia. Page Version ID: 1100066005. 24. Juli 2022. URL: https://en.wikipedia.org/w/index.php?title=Naming\_convention\_(programming)&oldid=1100066005 (besucht am 30.07.2022).

- Nemec, Devin. copy\_file\_to\_another\_repo\_action. original-date: 2020-08-24T19:25:58Z. 27. Juli 2022. URL: https://github.com/dmnemec/copy\_file\_to\_another\_repo\_action (besucht am 03.08.2022).
- Shinan, Erez. lark: a modern parsing library. Version 1.1.2. URL: https://github.com/lark-parser/lark (besucht am 31.07.2022).
- Ueda, Takahiro. *Makefile for LaTeX*. original-date: 2018-07-06T15:01:24Z. 10. Mai 2022. URL: https://github.com/tueda/makefile4latex (besucht am 03.08.2022).