

ALBERT LUDWIGS UNIVERSITÄT FREIBURG

TECHNISCHE FAKULTÄT

# PicoC-Compiler

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

BACHELORARBEIT

---

*Abgabedatum:* 28<sup>th</sup> April 2022

*Author:*  
Jürgen Mattheis

*Gutachter:*  
Prof. Dr. Scholl

*Betreuung:*  
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.

# Inhaltsverzeichnis

Abbildungsverzeichnis	I
Codeverzeichnis	II
Tabellenverzeichnis	III
Definitionsverzeichnis	V
Grammatikverzeichnis	VI
<b>1 Einführung</b>	<b>1</b>
1.1 Compiler und Interpreter	1
1.1.1 T-Diagramme	3
1.2 Formale Sprachen	5
1.2.1 Ableitungen	8
1.2.2 Präzidenz und Assoziativität	11
1.3 Lexikalische Analyse	12
1.4 Syntaktische Analyse	15
1.5 Code Generierung	22
1.5.1 Monadische Normalform	23
1.5.2 A-Normalform	24
1.5.3 Ausgabe des Maschinencodes	26
1.6 Fehlermeldungen	27
<b>2 Implementierung</b>	<b>29</b>
2.1 Lexikalische Analyse	31
2.1.1 Konkrete Syntax für die Lexikalische Analyse	31
2.1.2 Codebeispiel	33
2.2 Syntaktische Analyse	33
2.2.1 Umsetzung von Präzidenz und Assoziativität	33
2.2.2 Konkrete Syntax für die Syntaktische Analyse	38
2.2.3 Ableitungsbaum Generierung	40
2.2.3.1 Codebeispiel	41
2.2.3.2 Ausgabe des Ableitungsbaums	42
2.2.4 Ableitungsbaum Vereinfachung	43
2.2.4.1 Codebeispiel	44
2.2.5 Abstrakt Syntax Tree Generierung	45
2.2.5.1 PicoC-Knoten	47
2.2.5.2 RETI-Knoten	52
2.2.5.3 Kompositionen von PicoC-Knoten und RETI-Knoten mit besonderer Bedeutung	53
2.2.5.4 Abstrakte Syntax	55
2.2.5.5 Codebeispiel	56
2.2.5.6 Ausgabe des Abstrakter Syntaxbaum	57
2.3 Code Generierung	57
2.3.1 Passes	59
2.3.1.1 PicoC-Shrink Pass	60
2.3.1.1.1 Aufgabe	60

2.3.1.1.2	Abstrakte Syntax . . . . .	60
2.3.1.1.3	Codebeispiel . . . . .	61
2.3.1.2	PicoC-Blocks Pass . . . . .	63
2.3.1.2.1	Aufgabe . . . . .	63
2.3.1.2.2	Abstrakte Syntax . . . . .	63
2.3.1.2.3	Codebeispiel . . . . .	65
2.3.1.3	PicoC-ANF Pass . . . . .	66
2.3.1.3.1	Aufgabe . . . . .	66
2.3.1.3.2	Abstrakte Syntax . . . . .	67
2.3.1.3.3	Codebeispiel . . . . .	69
2.3.1.4	RETI-Blocks Pass . . . . .	70
2.3.1.4.1	Aufgabe . . . . .	70
2.3.1.4.2	Abstrakte Syntax . . . . .	70
2.3.1.4.3	Codebeispiel . . . . .	71
2.3.1.5	RETI-Patch Pass . . . . .	74
2.3.1.5.1	Aufgabe . . . . .	74
2.3.1.5.2	Abstrakte Syntax . . . . .	74
2.3.1.5.3	Codebeispiel . . . . .	75
2.3.1.6	RETI Pass . . . . .	78
2.3.1.6.1	Aufgabe . . . . .	78
2.3.1.6.2	Konkrete und Abstrakte Syntax . . . . .	78
2.3.1.6.3	Codebeispiel . . . . .	80

# Abbildungsverzeichnis

1.1	Horizontale Übersetzungszwischenschritte zusammenfassen . . . . .	5
1.2	Vertikale Interpretierungszwischenschritte zusammenfassen . . . . .	5
1.3	Veranschaulichung von Linksassoziativität und Rechtsassoziativität . . . . .	12
1.4	Veranschaulichung von Präzidenz . . . . .	12
1.5	Veranschaulichung der Lexikalischen Analyse . . . . .	15
1.6	Veranschaulichung des Unterschieds zwischen Ableitungsbaum und Abstraktem Syntaxbaum. . . . .	19
1.7	Veranschaulichung der Syntaktischen Analyse . . . . .	21
1.8	Codebeispiel für das Trennen von Ausdrücken mit und ohne Nebeneffekten . . . . .	24
1.9	Codebeispiel für das Entfernen Komplexer Ausdrücke aus Operationen . . . . .	26
2.1	Ableitungsbäume zu den beiden Ableitungen . . . . .	35
2.2	Ableitungsbaum nach Parsen eines Ausdrucks . . . . .	43
2.3	Ableitungsbaum nach Vereinfachung . . . . .	44
2.4	Abstrakter Syntaxbaum Generierung ohne Umdrehen . . . . .	46
2.5	Abstrakter Syntaxbaum Generierung mit Umdrehen . . . . .	46
2.6	Cross-Compiler Kompiliervorgang ausgeschrieben . . . . .	58
2.7	Cross-Compiler Kompiliervorgang Kurzform . . . . .	58
2.8	Architektur mit allen Passes ausgeschrieben . . . . .	59

# Codeverzeichnis

2.1	PicoC-Code des Codebeispiels . . . . .	33
2.2	Tokens für das Codebeispiel . . . . .	33
2.3	Ableitungsbaum nach Ableitungsbaum Generierung . . . . .	42
2.4	Ableitungsbaum nach Ableitungsbaum Vereinfachung . . . . .	45
2.5	Aus vereinfachtem Ableitungsbaum generierter Abstrakter Syntaxbaum . . . . .	56
2.6	PicoC Code für Codebeispiel . . . . .	62
2.7	Abstrakter Syntaxbaum für Codebeispiel . . . . .	63
2.8	PicoC-Blocks Pass für Codebeispiel . . . . .	66
2.9	PicoC-ANF Pass für Codebeispiel . . . . .	70
2.10	RETI-Blocks Pass für Codebeispiel . . . . .	74
2.11	RETI-Patch Pass für Codebeispiel . . . . .	78
2.12	RETI Pass für Codebeispiel . . . . .	82

# Tabellenverzeichnis

2.1	Präzidenzregeln von PicoC . . . . .	34
2.2	Zuordnung der Bezeichnungen von Produktionsregeln zu Operatoren . . . . .	36
2.3	PicoC-Knoten Teil 1 . . . . .	48
2.4	PicoC-Knoten Teil 2 . . . . .	49
2.5	PicoC-Knoten Teil 3 . . . . .	50
2.6	PicoC-Knoten Teil 4 . . . . .	51
2.7	RETI-Knoten . . . . .	53
2.8	Kompositionen von PicoC-Knoten und RETI-Knoten mit besonderer Bedeutung . . . . .	54

# Definitionsverzeichnis

1.1	Interpreter . . . . .	1
1.2	Compiler . . . . .	1
1.3	Maschiensprache . . . . .	2
1.4	Immediate . . . . .	2
1.5	Cross-Compiler . . . . .	3
1.6	T-Diagram Programm . . . . .	3
1.7	T-Diagram Übersetzer (bzw. eng. Translator) . . . . .	4
1.8	T-Diagram Interpreter . . . . .	4
1.9	T-Diagram Maschine . . . . .	4
1.10	Symbol . . . . .	5
1.11	Alphabet . . . . .	6
1.12	Wort . . . . .	6
1.13	Formale Sprache . . . . .	6
1.14	Syntax . . . . .	6
1.15	Semantik . . . . .	6
1.16	Formale Grammatik . . . . .	6
1.17	Chromsky Hierarchie . . . . .	7
1.18	Reguläre Grammatik . . . . .	8
1.19	Kontextfreie Grammatik . . . . .	8
1.20	Wortproblem . . . . .	8
1.21	1-Schritt-Ableitungsrelation . . . . .	9
1.22	Ableitungsrelation . . . . .	9
1.23	Links- und Rechtsableitungableitung . . . . .	9
1.24	Linksrekursive Grammatiken . . . . .	9
1.25	Formaler Ableitungsbaum . . . . .	10
1.26	Mehrdeutige Grammatik . . . . .	11
1.27	Assoziativität . . . . .	11
1.28	Präzidenz . . . . .	12
1.29	Pipe-Filter Architekturpattern . . . . .	12
1.30	Pattern . . . . .	13
1.31	Lexeme . . . . .	13
1.32	Lexer (bzw. Scanner oder auch Tokenizer) . . . . .	13
1.33	Bezeichner (bzw. Identifier) . . . . .	14
1.34	Literal . . . . .	15
1.35	Konkrete Syntax . . . . .	16
1.36	Ableitungsbaum (bzw. Konkretter Syntaxbaum, engl. Derivation Tree) . . . . .	16
1.37	Parser . . . . .	16
1.38	Recognizer (bzw. Erkennen) . . . . .	17
1.39	Transformer . . . . .	18
1.40	Visitor . . . . .	18
1.41	Abstrakte Syntax . . . . .	19
1.42	Abstrakter Syntaxbaum (bzw. engl. Abstract Syntax Tree, kurz AST) . . . . .	19
1.43	Pass . . . . .	22
1.44	Reiner Ausdruck (bzw. engl. pure expression) . . . . .	23
1.45	Unreiner Ausdruck . . . . .	23
1.46	Monadische Normalform (bzw. engl. monadic normal form) . . . . .	23
1.47	Location . . . . .	24



1.48	Atomarer Ausdruck . . . . .	25
1.49	Komplexer Ausdruck . . . . .	25
1.50	A-Normalform (ANF) . . . . .	25
1.51	Fehlermeldung . . . . .	27
2.1	Metasyntax . . . . .	29
2.2	Metasprache . . . . .	29
2.3	Erweiterte Backus-Naur-Form (EBNF) . . . . .	29
2.4	Dialekt der EBNF aus Lark . . . . .	30
2.5	Abstrakte Syntax Form (ASF) . . . . .	30
2.6	Earley Parser . . . . .	40
2.7	Earley Recognizer . . . . .	40
2.8	Label . . . . .	52
2.9	Token-Knoten . . . . .	52
2.10	Container-Knoten . . . . .	52
2.11	Symboltabelle . . . . .	67

# Grammatikverzeichnis

1.1	Produktionen des Ableitungsbaums . . . . .	10
2.1.1	Grammatik der Konkreten Syntax der Sprache $L_{PicoC}$ für die Lexikalische Analyse in EBNF	32
2.2.1	Undurchdachte Konkrete Syntax der Sprache $L_{PicoC}$ für die Syntaktische Analyse in EBNF	34
2.2.2	Durchdachte Konkrete Syntax der Sprache $L_{PicoC}$ für die Syntaktische Analyse in EBNF . .	35
2.2.3	Beispiel für eine unäre rechtsassoziative Produktion . . . . .	36
2.2.4	Beispiel für eine unäre linksassoziative Produktion . . . . .	36
2.2.5	Beispiel für eine linksassoziative Produktion . . . . .	37
2.2.6	Beispiel für eine linksassoziative Produktion . . . . .	37
2.2.7	Durchdachte Konkrete Syntax für Operatorpräzidenz in EBNF . . . . .	38
2.2.8	Grammatik der Konkreten Syntax der Sprache $L_{PicoC}$ für die Syntaktische Analyse in EBNF, Teil 1 . . . . .	39
2.2.9	Grammatik der Konkreten Syntax der Sprache $L_{PicoC}$ für die Syntaktische Analyse in EBNF, Teil 2 . . . . .	40
2.2.10	Abstrakte Syntax der Sprache $L_{PlocC}$ . . . . .	55
2.3.1	Abstrakte Syntax der Sprache $L_{PlocC\_Shrink}$ . . . . .	61
2.3.2	Abstrakte Syntax der Sprache $L_{PlocC\_Blocks}$ . . . . .	64
2.3.3	Abstrakte Syntax der Sprache $L_{PlocC\_ANF}$ . . . . .	68
2.3.4	Abstrakte Syntax der Sprache $L_{RETI\_Blocks}$ . . . . .	71
2.3.5	Abstrakte Syntax der Sprache $L_{RETI\_Patch}$ . . . . .	75
2.3.6	Konkrete Syntax der Sprache $L_{RETI}$ für die Lexikalische Analyse in EBNF . . . . .	79
2.3.7	Konkrete Syntax der Sprache $L_{RETI}$ für die Syntaktische Analyse in EBNF . . . . .	79
2.3.8	Abstrakte Syntax der Sprache $L_{RETI}$ . . . . .	79

# 1 Einführung

## 1.1 Compiler und Interpreter

Der wohl wichtigsten zu klärenden Begriffe, sind die eines **Compilers** (Definition 1.2) und eines **Interpreters** (Definition 1.1), da das Schreiben eines Compilers von der **PicoC-Sprache**  $L_{PicoC}$  in die **RETI-Sprache**  $L_{RETI}$  das Thema dieser Bachelorarbeit ist und die Definition eines **Interpreters** genutzt wird, um zu definieren was ein **Compiler** ist. Des Weiteren wurde zur **Qualitätsicherung** ein **RETI-Interpreter** implementiert, um mithilfe des **GCC**<sup>1</sup> und von **Tests** die **Beziehungen** in 1.2.1 zu belegen (siehe Subkapitel ??).

### Definition 1.1: Interpreter

*Interpretiert die Befehle bzw. **Statements** eines Programmes  $P$  direkt.*

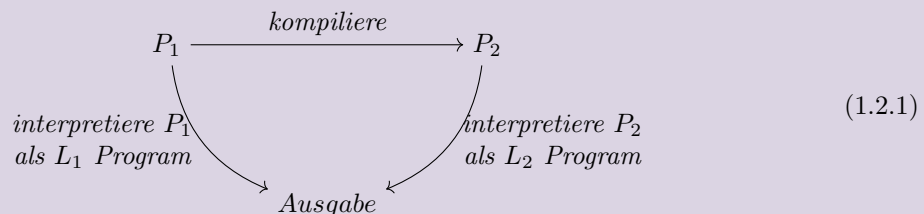
*Auf die Implementierung bezogen arbeitet ein Interpreter auf den compilerinternen **Sub-Bäumen** des **Abstrakter Syntaxbaum** (Definition 1.42) und führt je nach Komposition der **Nodes** des Abstrakter Syntaxbaum, auf die er während des Darüber-Iterierens stösst unterschiedliche Anweisungen aus.<sup>a</sup>*

<sup>a</sup>G. Siek, *Course Webpage for Compilers (P423, P523, E313, and E513)*.

### Definition 1.2: Compiler

***Kompiliert** ein **beliebiges** Program  $P_1$ , welches in einer Sprache  $L_1$  geschrieben ist, in ein Program  $P_2$ , welches in einer Sprache  $L_2$  geschrieben ist.*

*Wobei **Kompilieren** meint, dass ein beliebiges Program  $P_1$  in der Sprache  $L_1$  so in die Sprache  $L_2$  zu einem Programm  $P_2$  übersetzt wird, dass bei beiden Programmen, wenn sie von **Interpretern** ihrer jeweiligen Sprachen  $L_1$  und  $L_2$  **interpretiert** werden, die gleiche **Ausgabe** rauskommt, wie es in Diagramm 1.2.1 dargestellt ist. Also beide Programme  $P_1$  und  $P_2$  die gleiche **Semantik** (Definition 1.15) haben und sich nur **syntaktisch** (Definition 1.14) durch die Sprachen  $L_1$  und  $L_2$ , in denen sie geschrieben stehen unterscheiden.<sup>a</sup>*



<sup>a</sup>G. Siek, *Course Webpage for Compilers (P423, P523, E313, and E513)*.

<sup>1</sup>Sammlung von Compilern für Linux bzw. GNU-Linux, steht für **GNU Compiler Collection**

Im Folgenden wird ein voll ausgeschriebener **Compiler** als  $C_{i.w.k.min}^{o-j}$  geschrieben, wobei  $C_w$  die **Sprache** bezeichnet, die der Compiler als **Input** nimmt und zu einer nicht näher spezifizierten Maschinensprache  $L_{B_i}$  einer Maschine  $M_i$  kompiliert. Fall die Notwendigkeit besteht die **Maschine**  $M_i$  anzugeben, zu dessen **Maschinensprache**  $L_{B_i}$  der Compiler kompiliert, wird das als  $C_i$  geschrieben. Falls die Notwendigkeit besteht die **Sprache**  $L_o$  anzugeben, in der der Compiler selbst geschrieben ist, wird das als  $C^o$  geschrieben. Falls die Notwendigkeit besteht die Version der Sprache, in die der Compiler kompiliert ( $L_{w.k}$ ) oder in der er selbst geschrieben ist ( $L_{o.j}$ ) anzugeben, wird das als  $C_{w.k}^{o-j}$  geschrieben. Falls es sich um einen **minimalen Compiler** handelt (Definition ??) kann man das als  $C_{min}$  schreiben.

Üblicherweise kompiliert ein **Compiler** ein **Program**, dass in einer **Programmiersprache** geschrieben ist zu **Maschinenenncode**, der in **Maschinensprache** (Definition 1.3) geschrieben ist, aber es gibt z.B. auch **Transpiler** (Definition ??) oder **Cross-Compiler** (Definition 1.5). Des Weiteren sind **Maschinensprache** und **Assemblersprache** (Definition ??) voneinander zu unterscheiden.

### Definition 1.3: Maschinensprache

*Programmiersprache, deren mögliche Programme die **hardwarenaheste Repräsentation** eines möglicherweise zuvor hierzu kompilierten bzw. assemblierten Programmes darstellen. Jeder Maschinenbefehl entspricht einer bestimmten **Aufgabe**, die die CPU im **vereinfachten Fall** in einem **Zyklus** der **Fetch- und Execute-Phase**, genauer gesagt in der **Execute-Phase** übernehmen kann oder allgemein in einer **geringen konstanten** Anzahl von Fetch- und Execute Phasen im **Komplexeren Fall**. Die Maschinenbefehle sind meist so designed, dass sie sich innerhalb bestimmter **Wortbreiten**, die 2er Potenzen sind codieren lassen. Im einfachsten Fall innerhalb einer **Speicherzelle** des **Hauptspeichers**.<sup>a,b</sup>*

<sup>a</sup>Viele Prozessorarchitekturen erlauben es allerdings auch z.B. **zwei** Maschinenbefehle in **eine** Speicherzelle des Hauptspeichers zu komprimieren, wenn diese zwei Maschinenbefehle keine Operanden mit zu großen **Immediates** (Definition 1.4) haben.

<sup>b</sup>Scholl, „Betriebssysteme“.

Der **Maschinenenncode**, denn ein üblicher Compiler einer Programmiersprache generiert, enthält seine Folge von Maschinenbefehlen üblicherweise in **binärer Repräsentation**, da diese in erster Linie für die Maschine, die binär arbeitet verständlich sein sollen und nicht für den Programmierer.

Der **PicoC-Compiler**, der den Zweck erfüllt für Studenten ein **Anschauungs- und Lernwerkzeug** zu sein, generiert allerdings Maschinenenncode, der die Maschinenbefehle bzw. RETI-Befehle in **menschenlesbarer Form** mit ausgeschriebenen RETI-Operationen, RETI-Registern und Immediates (Definition 1.4) enthält. Für den **RETI-Interpreter** ist es ebenfalls nicht notwendig, dass der Maschinenenncode, denn der PicoC-Compiler generiert in binärer Darstellung ist, denn es ist für den RETI-Interpreter ebenfalls leichter diese einfach direkt in menschenlesbarer Form zu interpretieren, da der RETI-Interpreter nur die sichtbare Funktionsweise einer RETI-CPU **simulieren** soll und nicht deren mögliche interne Umsetzung<sup>2</sup>.

### Definition 1.4: Immediate

***Konstanter Wert**, der als **Teil eines Maschinenbefehls** gespeichert ist und dessen **Wertebereich** dementsprechend auch durch die die Anzahl an Bits, die ihm innerhalb dieses **Maschinenbefehls** zur Verfügung gestellt sind, **beschränkter** ist als bei sonstigen Werten innerhalb des Hauptspeichers, denen eine ganze Speicherzelle des Hauptspeichers zur Verfügung steht.<sup>a</sup>*

<sup>2</sup>Eine **RETI-CPU** zu bauen, die menschenlesbaren Maschinenenncode in z.B. **UTF-8 Codierung** ausführen kann, wäre dagegen unnötig kompliziert und aufwändig, da Hardware **binär** arbeitet und man dieser daher lieber direkt die binär codierten Maschinenbefehle übergibt, anstatt z.B. eine unnötig **platzverbrauchenden** UTF-8 Codierung zu verwenden, die nur in sehr vielen Schritt einen Befehl verarbeiten kann, da die Register und Speicherzellen des Hauptspeichers üblicherweise nur **32- bzw. 64-Bit Breite** haben.

<sup>a</sup>Ljohhuh, *What is an immediate value?*

### Definition 1.5: Cross-Compiler

Kompiliert auf einer **Maschine**  $M_1$  ein Programm, dass in einer **Sprache**  $L_w$  geschrieben ist für eine **andere Maschine**  $M_2$ , wobei beide Maschinen  $M_1$  und  $M_2$  unterschiedliche **Maschinensprachen**  $B_1$  und  $B_2$  haben.<sup>a,b</sup>

<sup>a</sup>Beim **PicoC-Compiler** handelt es sich um einen **Cross-Compiler**  $C_{PicoC}^{Python}$ .

<sup>b</sup>Earley und Sturgis, „A formalism for translator interactions“.

Ein **Cross-Compiler** ist entweder notwendig, wenn eine Zielmaschine  $M_2$  nicht ausreichend **Rechenleistung** hat, um ein Programm in der Wunschsprache  $L_w$  selbst **zeitnah** zu kompilieren oder wenn noch kein Compiler  $C_w$  für die **Wunschsprache**  $L_w$  und andere Programmiersprachen  $L_o$ , in denen man Programmieren wollen würde existiert, der unter der **Maschinensprache**  $B_2$  einer Zielmaschine  $M_2$  läuft.<sup>3</sup>

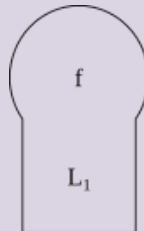
#### 1.1.1 T-Diagramme

Um die Architektur von **Compilern** und **Interpretern** übersichtlich darzustellen eignen sich **T-Diagramme**, deren Spezifikation aus dem Paper Earley und Sturgis, „A formalism for translator interactions“ entnommen ist besonders gut, da diese optimal darauf **zugeschnitten** sind die Eigenheiten von Compilern in ihrer Art der Darstellung unterzubringen.

Die **Notation** setzt sich dabei aus den **Blöcken** für ein Programm (Definition 1.6), einen Übersetzer (Definition 1.7), einen Interpreter (Definition 1.8) und eine Maschine (Definition 1.9) zusammen.

### Definition 1.6: T-Diagramm Programm

Repräsentiert ein **Programm**, dass in der **Sprache**  $L_1$  geschrieben ist und die **Funktion**  $f$  berechnet.<sup>a</sup>



<sup>a</sup>Earley und Sturgis, „A formalism for translator interactions“.

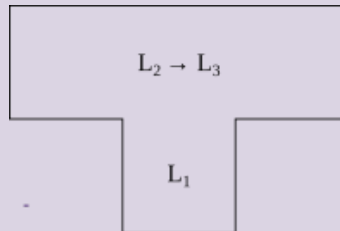
Es ist bei **T-Diagrammen** nicht notwendig beim entsprechenden **Platzhalter**, in den man die genutzte **Sprache** schreibt, den **Namen der Sprache** an ein  $L$  dranzuhängen, weil hier immer eine **Sprache** steht. Es würde in Definition 1.6 also reichen einfach eine 1 hinzuschreiben.

<sup>3</sup>Die an vielen Universitäten und Schulen eingesetzten programmierbaren Roboter von **Lego Mindstorms** nutzen z.B. einen **Cross-Compiler**, um für den programmierbaren Microcontroller eine **C-ähnliche Sprache** in die Maschinensprache des Microcontrollers zu kompilieren, da der Microcontroller selbst nicht genug Rechenleistung besitzt, um ein Programm selbst zeitnah zu kompilieren.

**Definition 1.7: T-Diagramm Übersetzer (bzw. eng. Translator)**

Repräsentiert einen **Übersetzer**, der in der **Sprache**  $L_1$  geschrieben ist und **Programme** von der **Sprache**  $L_2$  in die **Sprache**  $L_3$  kompiliert.

Für den **Übersetzer** gelten genauso, wie für einen **Compiler**<sup>a</sup> die **Beziehungen** in 1.2.1.<sup>b</sup>

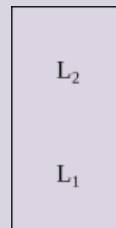


<sup>a</sup>Zwischen den Begriffen **Übersetzung** und **Kompilierung** gibt es einen kleinen Unterschied, **Übersetzung** ist **kleinschrittiger** als **Kompilierung** und ist auch zwischen **Passes** möglich, **Kompilierung** beinhaltet dagegen bereits alle **Passes** in einem Schritt. **Kompilieren** ist also auch **Übersetzen**, aber **Übersetzen** ist nicht immer auch **Kompilieren**.

<sup>b</sup>Earley und Sturgis, „A formalism for translator interactions“.

**Definition 1.8: T-Diagramm Interpreter**

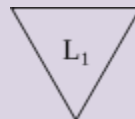
Repräsentiert einen **Interpreter**, der in der **Sprache**  $L_1$  geschrieben ist und **Programme** in der **Sprache**  $L_2$  interpretiert.<sup>a</sup>



<sup>a</sup>Earley und Sturgis, „A formalism for translator interactions“.

**Definition 1.9: T-Diagramm Maschine**

Repräsentiert eine **Maschine**, welche ein **Programm** in **Maschinensprache**  $L_1$  ausführt.<sup>a,b</sup>



<sup>a</sup>Wenn die Maschine **Programme** in einer höheren Sprache als **Maschinensprache** ausführt, ist es auch erlaubt diese Notation zu verwenden, dann handelt es sich um eine **Abstrakte Maschine**, wie z.B. die **Python Virtual Machine** (PVM) oder **Java Virtual Machine** (JVM).

<sup>b</sup>Earley und Sturgis, „A formalism for translator interactions“.

Aus den verschiedenen **Blöcken** lassen sich **Kompositionen** bilden, indem man sie **adjazent** zueinander platziert. Allgemein lässt sich grob sagen, dass **vertikale Adjazents** für **Interpretation** und **horizontale Adjazents** für **Übersetzung** steht.

Sowohl **horizontale** als auch **vertikale Adjazents** lassen sich, wie man in den Abbildungen 1.1 und 1.2 erkennen kann zusammenfassen.

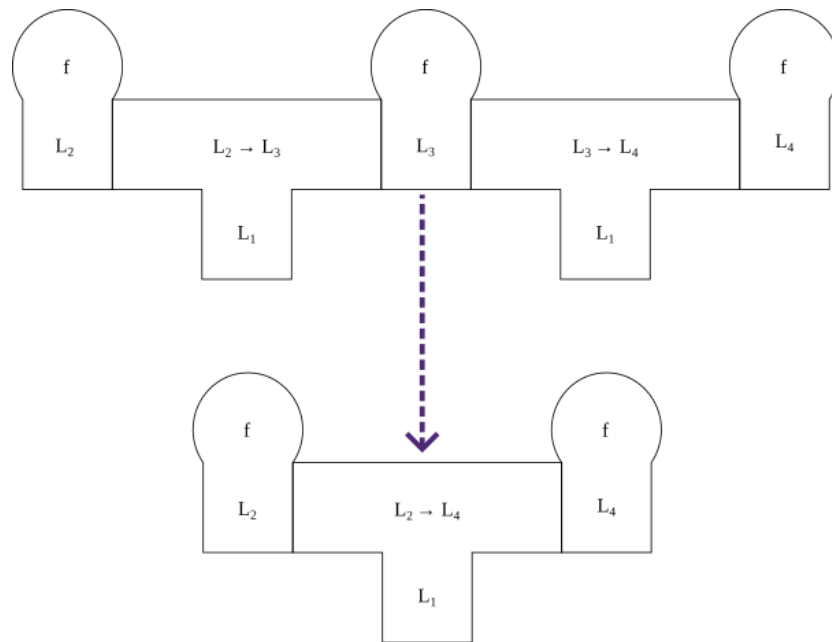


Abbildung 1.1: Horizontale Übersetzungszwischenschritte zusammenfassen

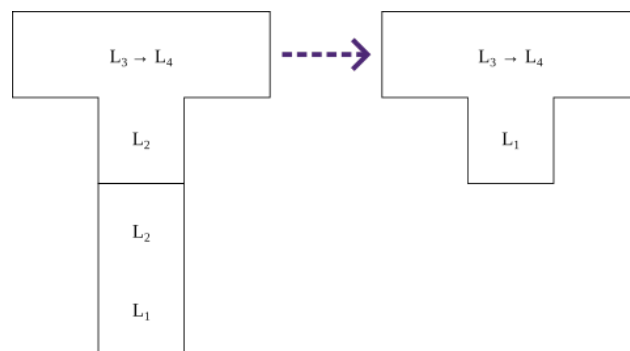


Abbildung 1.2: Vertikale Interpretierungszwischenschritte zusammenfassen

## 1.2 Formale Sprachen

Das **Kompilieren** eines Programmes hat viel mit dem Thema **Formaler Sprachen** (Definition 1.13) zu tun, da bereits das Kompilieren an sich das Übersetzen eines Programmes aus der Sprache  $L_1$  in eine Sprache  $L_2$  ist. Aus diesem Grund ist es wichtig die **Grundlagen Formaler Sprachen**, was die Begriffe **Symbol** (Definition 1.10), **Alphabet** (Definition 1.11), **Wort** (Definition 1.12) usw. beinhaltet vorher eingeführt zu haben.

### Definition 1.10: Symbol

„Ein Symbol ist ein **Element** eines **Alphabets**  $\Sigma$ .“<sup>a</sup>

<sup>a</sup>Nebel, „Theoretische Informatik“.

**Definition 1.11: Alphabet**

„Ein Alphabet ist eine **endliche, nicht-leere** Menge aus **Symbolen** (Definition 1.10).“<sup>a</sup>

<sup>a</sup>Nebel, „Theoretische Informatik“.

**Definition 1.12: Wort**

„Ein Wort  $w = a_1 \dots a_n \in \Sigma^*$  ist eine **endliche Folge** von **Symbolen** aus einem **Alphabet**  $\Sigma$ .“<sup>a</sup>

<sup>a</sup>Nebel, „Theoretische Informatik“.

**Definition 1.13: Formale Sprache**

„Eine Formale Sprache ist eine Menge von **Wörtern** (Definition 1.12) über dem **Alphabet**  $\Sigma$  (Definition 1.11).“<sup>a</sup>

Das Adjektiv „**formal**“ kann dabei weggelassen werden, wenn der **Kontext** indem die **Sprache** verwendet wird **eindeutig** ist, da man das Adjektiv „**formal**“ nur verwendet um den Unterschied zum im normalen Sprachgebrauch verwendeten Begriff einer **Sprache** herauszustellen.

<sup>a</sup>Nebel, „Theoretische Informatik“.

Bei der Übersetzung eines Programmes von einer Sprache  $L_1$  zur Sprache  $L_2$  muss die **Semantik** (Definition 1.15) **gleich** bleiben. Beide Sprachen  $L_1$  und  $L_2$  haben eine **Grammatik** (Definition 1.16), welche diese beschreibt und können verschiedene **Syntaxen** (Definition 1.14) haben.

**Definition 1.14: Syntax**

Die Syntax bezeichnet alles was mit dem **Aufbau** von **Formalen Sprachen** zu tun hat. Die **Grammatik** einer Sprache, aber auch die in **Natürlicher Sprache** ausgedrückten Regeln, welche den Aufbau von Wörtern einer Formalen Sprache beschreiben werden als Syntax bezeichnet. Es kann auch mehrere **verschiedene Syntaxen** für die **gleiche Sprache** geben<sup>a, b</sup>

<sup>a</sup>Z.B. die **Konkrete** und **Abstrakte Syntax**, die später eingeführt werden.

<sup>b</sup>Thiemann, „Einführung in die Programmierung“.

**Definition 1.15: Semantik**

Die Semantik bezeichnet alles was mit der **Bedeutung** von **Formalen Sprachen** zu tun hat.<sup>a</sup>

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

**Definition 1.16: Formale Grammatik**

„Eine Formale Grammatik beschreibt wie **Wörter** einer **Sprache** abgeleitet werden können.“<sup>a</sup>

Das Adjektiv „**formal**“ kann dabei weggelassen werden, wenn der **Kontext** indem die **Grammatik** verwendet wird **eindeutig** ist, da man das Adjektiv „**formal**“ nur verwendet um den Unterschied zum im normalen Sprachgebrauch verwendeten Begriff einer **Grammatik** herauszustellen.

Eine Grammatik wird durch das Tupel  $G = \langle N, \Sigma, P, S \rangle$  dargestellt, wobei:

- $N \hat{=}$  **Nicht-Terminalsymbole**.



- $\Sigma \hat{=} \text{Terminalsymbole}$ , wobei  $N \cap \Sigma = \emptyset^{b,c}$ .
- $P \hat{=} \text{Menge von Produktionsregeln}$   $w \rightarrow v$ , wobei  $w, v \in (N \cup \Sigma)^* \wedge w \notin \Sigma^*$ .<sup>de</sup>
- $S \hat{=} \text{Startsymbol}$ , wobei  $S \in N$ .

Zusätzlich ist es praktisch **Nicht-Terminalsymbole**  $N$ , **Terminalsymbole**  $\Sigma$  und das **leere Wort**  $\varepsilon$  allgemein als Menge der **Grammatiksymbole**  $C = N \cup \Sigma \cup \varepsilon$  zu definieren.

<sup>a</sup>Nebel, „Theoretische Informatik“.

<sup>b</sup>Weil mit ihnen **terminiert** wird.

<sup>c</sup>Kann auch als **Alphabet** bezeichnet werden.

<sup>d</sup> $w$  muss **mindestens** ein **Nicht-Terminalsymbol** enthalten.

<sup>e</sup>Bzw.  $w, v \in V^* \wedge w \notin \Sigma^*$ .

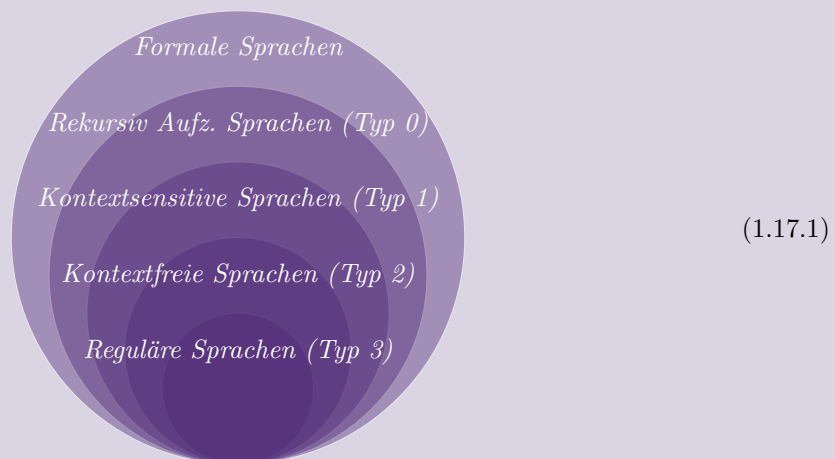
Die gerade definierten **Formale Sprachen** lassen sich des Weiteren in Klassen der **Chromsky Hierarchie** (Definition 1.17) einteilen.

### Definition 1.17: Chromsky Hierarchie

Die Chromsky Hierarchie ist eine Hierarchie in der **Formale Sprachen** nach der **Komplexität** ihrer **Formalen Grammatiken** in verschiedene **Klassen** unterteilt werden. Jede dieser Klassen hat verschiedene **Eigenschaften**, wie **Entscheidungsprobleme**, die in dieser Klasse **entscheidbar** bzw. **unentscheidbar** sind usw.

Eine Sprache  $L_i$  ist in der **Chromsky Hierarchie** vom Typ  $i \in \{0, \dots, 3\}$ , falls sie von einer Grammatik dieses Typs  $i$  erzeugt wird.

Zwischen den Sprachmengen **benachbarter Klassen** in Abbildung 1.17.1 besteht eine **echte Teilmengebeziehung**:  $L_3 \subset L_2 \subset L_1 \subset L_0$ . Jede **Reguläre Sprache** ist auch eine **Kontextfreie Sprache**, aber nicht jede **Kontextfreie Sprache** ist auch eine **Reguläre Sprache**.<sup>a</sup>



<sup>a</sup>Nebel, „Theoretische Informatik“.

Für diese Bachelorarbeit sind allerdings nur die **Spracheklassen** der **Chromsky-Hierarchie** relevant, die von **Regulären** (Definition 1.18) und **Kontextfreien Grammatiken** (Definition 1.19) beschrieben werden.

**Definition 1.18: Reguläre Grammatik**

„Ist eine Grammatik für die gilt, dass alle Produktionen eine der Formen:

$$A \rightarrow cB, \quad A \rightarrow c, \quad A \rightarrow \varepsilon \quad (1.18.1)$$

haben, wobei  $A, B$  **Nicht-Terminalsymbole** sind und  $c$  ein **Terminalsymbol** ist<sup>a,b</sup>.“<sup>c</sup>

<sup>a</sup>Diese Definition einer **Regulären Grammatik** ist **rechtsregulär**, es ist auch möglich diese Definition **linksregulär** zu formulieren, aber diese Details sind für die Bachelorarbeit nicht relevant.

<sup>b</sup>Dadurch, dass die **linke** Seite immer nur ein **Nicht-Terminalsymbol** sein darf ist jede **Reguläre Grammatik** auch eine **Kontextfrei Grammatik**.

<sup>c</sup>Nebel, „Theoretische Informatik“.

**Definition 1.19: Kontextfreie Grammatik**

„Ist eine Grammatik für die gilt, dass alle Produktionen die Form:

$$A \rightarrow v \quad (1.19.1)$$

haben, wobei  $A$  ein **Nicht-Terminalsymbol** ist und  $v$  ein beliebige Folge von **Grammatiksymbolen**<sup>a</sup> ist.“<sup>b</sup>

<sup>a</sup>Also eine beliebige Folge von **Nicht-Terminalsymbolen** und **Terminalsymbolen**.

<sup>b</sup>Nebel, „Theoretische Informatik“.

Ob sich ein Programm überhaupt kompilieren lässt entscheidet sich anhand des **Wortproblems** (Definition 1.20). In einem **Compiler** oder **Interpreter** ist das Wortproblem üblicherweise immer **entscheidbar**. Wenn das Programm ein **Wort** der **Sprache** ist, die der Compiler kompiliert, so klappt das Kompilieren, ist es **kein Wort** der **Sprache**, die der Compiler kompiliert, wird eine **Fehlermeldung** ausgegeben.

**Definition 1.20: Wortproblem**

Ein Entscheidungsproblem, bei dem man zu einem **Wort**  $w \in \Sigma^*$  und einer **Sprache**  $L$  als **Eingabe** 1 oder 0<sup>a</sup> **ausgibt**, je nachdem, ob dieses Wort  $w$  Teil der Sprache  $L$  ist  $w \in L$  oder nicht  $w \notin L$ .<sup>b</sup>

Das Wortproblem kann durch die folgende **Indikatorfunktion**<sup>c</sup> zusammengefasst werden:

$$\mathbb{1}_L : \Sigma^* \rightarrow \{0, 1\} : w \mapsto \begin{cases} 1 & \text{falls } w \in L \\ 0 & \text{sonst} \end{cases} \quad (1.20.1)$$

<sup>a</sup>Bzw. „ja“ oder „nein“ usw., es muss nicht umgedingt 1 oder 0 sein.

<sup>b</sup>Nebel, „Theoretische Informatik“.

<sup>c</sup>Auch **Charakteristische Funktion** genannt.

**1.2.1 Ableitungen**

Um sicher zu wissen, ob ein Compiler ein **Programm**<sup>4</sup> kompilieren kann, ist es möglich das Programm mithilfe der **Grammatik** der **Sprache** des Compilers **abzuleiten**. Hierbei wird zwischen der **1-Schritt-Ableitungsrelation** (Definition 1.21) und der normalen **Ableitungsrelation** (Definition 1.22) unterscheiden.

<sup>4</sup>Bzw. **Wort**.

**Definition 1.21: 1-Schritt-Ableitungsrelation**

„Eine **binäre Relation**  $\Rightarrow$  zwischen Wörtern aus  $(N \cup \Sigma)^*$ , die alle möglichen Wörter  $(N \cup \Sigma)^*$  in Relation zueinander setzt, die sich nur durch das **einmalige** Anwenden einer Produktionsregel voneinander unterscheiden.

Es gilt  $u \Rightarrow v$  **genau dann wenn**  $u = w_1 x w_2$ ,  $v = w_1 y w_2$  **und** es eine Regel  $x \rightarrow y \in P$  gibt, wobei  $w_1, w_2, x, y \in (N \cup \Sigma)^*$ “<sup>a</sup>

<sup>a</sup>Nebel, „Theoretische Informatik“.

**Definition 1.22: Ableitungsrelation**

„Eine **binäre Relation**  $\Rightarrow^*$ , welche der **reflexive, transitive Abschluss** der **1-Schritt-Ableitungsrelation**  $\Rightarrow$  ist. Auf der **rechten** Seite der Ableitungsrelation  $\Rightarrow^*$  steht also ein Wort aus  $(N \cup \Sigma)^*$ , welches durch **beliebig häufiges** Anwenden von Produktionsregeln entsteht.

Es gilt  $u \Rightarrow^* v$  **genau dann wenn**  $u = w_1 \Rightarrow \dots \Rightarrow w_n = v$ , wobei  $n \geq 1$  und  $w_1, \dots, w_n \in (N \cup \Sigma)^*$ .“<sup>a</sup>

<sup>a</sup>Nebel, „Theoretische Informatik“.

Beim Ableiten kann auf verschiedene Weisen vorgegangen werden, dasselbe **Programm**<sup>5</sup> kann z.B. über eine **Linksableitung** als auch eine **Rechtsableitung** (Definition 1.23) abgeleitet werden. Das ist später bei den verschiedenen **Ansätzen** für das **Parsen** eines Programmes in Unterkapitel 1.4 relevant.

**Definition 1.23: Links- und Rechtsableitung**

„In jedem **Ableitungsschritt** wird bei **Typ-3- und Typ-2-Grammatiken** auf das am **weitesten links** (**Linksableitung**) bzw. **rechts** (**Rechtsableitung**) stehende **Nicht-Terminalsymbol** eine Produktionsregel angewandt, bei **Typ-1- und Typ-0-Grammatiken** ist es statt einem **Nicht-Terminalsymbol** die **linke** Seite einer Produktion.

Mit diesem Vorgehen kann man jedes ableitbare Wort generieren, denn dieses Vorgehen entspricht **Tiefensuche von links-nach-rechts**.“<sup>a</sup>

<sup>a</sup>Nebel, „Theoretische Informatik“.

Manche der **Ansätze** für das **Parsen** eines Programmes haben ein Problem, wenn die Grammatik, die zur Entscheidung des **Wortproblems** für das Programm verwendet wird eine **Linksrekursive Grammatik** (Definition 1.24) ist<sup>6</sup>.

**Definition 1.24: Linksrekursive Grammatiken**

Eine **Grammatik** ist **linksrekursiv**, wenn sie ein **Nicht-Terminalsymbol** enthält, dass **linksrekursiv** ist.

Ein **Nicht-Terminalsymbol** ist **linksrekursiv**, wenn das **linkeste Symbol** in einer seiner **Produktionen** es selbst ist oder zu sich selbst gemacht werden kann durch eine Folge von Ableitungen:

$$A \Rightarrow^* Aa,$$

wobei  $a$  eine beliebige Folge von **Terminalsymbolen** und **Nicht-Terminalsymbolen** ist.<sup>a</sup>

<sup>a</sup>Parsing Expressions · Crafting Interpreters.

<sup>5</sup>Bzw. **Wort**.

<sup>6</sup>Für den im **PicoC-Compiler** verwendeten **Earley Parsers** stellt dies allerdings **kein** Problem dar.

Um herauszufinden, ob eine Grammatik **mehrdeutig** (siehe Unterkapitel ??) ist, werden **Ableitungen** als **Formale Ableitungsbäume** (Definition 1.25) dargestellt. **Formale Ableitungsbäume** werden im Unterkapitel 1.4 nochmal relevant, da in der **Syntaktischen Analyse** Ableitungsbäume (Definition 1.36) als eine **compilerinterne Datenstruktur** umgesetzt werden.

#### Definition 1.25: Formaler Ableitungsbaum

Ist ein Baum, in dem die **Konkrete Syntax** eines **Wortes**<sup>a</sup> nach den **Produktionen** der zugehörigen Grammatik, die angewendet werden mussten um das Wort **abzuleiten** zergliedert **hierarchisch** dargestellt wird.

Das Adjektiv „**formal**“ kann dabei weggelassen werden, wenn der **Kontext** in dem der **Ableitungsbaum** verwendet wird **eindeutig** ist, da man das Adjektiv „**formal**“ nur verwendet um den Unterschied zum **compilerinternen Ableitungsbaum** herauszustellen, der den Formalen Ableitungsbaum als **Datenstruktur** zur einfachen Weiterverarbeitung umsetzt.

Den Knoten dieses Baumes sind **Grammatiksymbole**  $C = N \cup \Sigma \cup \varepsilon$  (Definition 1.16) zugeordnet. Die **Inneren Knoten** des Baumes sind **Nicht-Terminalsymbole**  $N$  und die **Blätter** sind entweder **Terminalsymbole**  $\Sigma$  oder das **leere Wort**  $\varepsilon$ .<sup>b</sup>

<sup>a</sup>Z.B. **Programmcode**.

<sup>b</sup>Nebel, „Theoretische Informatik“.

In Abbildung 1.25.2 ist ein Beispiel für einen **Formalen Ableitungsbaum** zu sehen, der sich aus der **Ableitung** 1.25.1 nach den im **Dialekt der Erweiterter Backus-Naur-Form des Lark Parsing Toolkit** (Definition 2.3) angegebenen **Produktionen** 1.1 einer ansonsten nicht näher spezifizierten **Grammatik**  $G = \langle N, \Sigma, P, add \rangle$  ergibt.

$DIG\_NO\_0$	$::=$	"1"		"2"		"3"		"4"		"5"		"6"	$L\_Lex$
		"7"		"8"		"9"							
$DIG\_WITH\_0$	$::=$	"0"		$DIG\_NO\_0$									
$NUM$	$::=$	"0"		$DIG\_NO\_0$	$DIG\_WITH\_0^*$								
$add$	$::=$	$add$	"+"	$mul$		$mul$							$L\_Parse$
$mul$	$::=$	$mul$	"*"	$NUM$		$NUM$							

**Grammatik 1.1:** *Produktionen des Ableitungsbaums*

$$add \Rightarrow mul \Rightarrow mul \text{ "*" } NUM \Rightarrow NUM \text{ "*" } NUM \Rightarrow 4 \text{ "*" } NUM \Rightarrow 4 \text{ "*" } 2 \quad (1.25.1)$$

Bei Ableitungsbäumen gibt es **keine** einheitliche **Regelung**, wie damit umgegangen wird, wenn die **Alternativen** einer Produktion unterschiedliche viele **Nicht-Terminalsymbole** enthalten. Es gibt einmal die Möglichkeit, wie im Ableitungsbaum 1.25.2 von der **Maximalzahl** auszugehen und beim Nicht-Erreichen der Maximalzahl entsprechend der **Differenz zur Maximalzahl** viele **Blätter** mit dem **leeren Wort**  $\varepsilon$  hinzuzufügen.



Eine andere Möglichkeit ist, wie im Ableitungsbaum 1.25.3 nur die vorhandenen **Nicht-Terminalsymbole** als Kinder hinzuzufügen<sup>7</sup>.



Für einen Compiler ist es notwendig, dass die **Grammatik**, welche die **Konkrete Syntax** beschreibt keine **Mehrdeutige Grammatik** (Definition 1.26) ist, denn sonst können unter anderem die **Präferenzregeln** der verschiedenen **Operatoren** nicht gewährleistet werden, wie später in Unterkapitel 2.2.1 an einem Beispiel demonstriert wird.

#### Definition 1.26: Mehrdeutige Grammatik

„Eine Grammatik ist **mehrdeutig**, wenn es ein Wort  $w \in L(G)$  gibt, das mehrere **Ableitungsbäume** zulässt“.<sup>a,b</sup>

<sup>a</sup>Alternativ, wenn es für  $w$  **mehrere** unterschiedliche **Linksableitungen** gibt.

<sup>b</sup>Nebel, „Theoretische Informatik“.

## 1.2.2 Präferenz und Assoziativität

Will man die **Operatoren** aus einer **Programmiersprache** in einer **Grammatik** für eine **Konkrete Syntax** ausdrücken, die **nicht mehrdeutig** ist, so lässt sich das nach einem klaren Schema machen, wenn die **Assoziativität** (Definition 1.27) und **Präferenz** (Definition 1.28) dieser **Operatoren** festgelegt ist. Dieses Schema wird in Unterkapitel 2.2.1 genauer erklärt.

#### Definition 1.27: Assoziativität

„Bestimmt, welcher Operator aus einer Reihe **gleicher** Operatoren **zuerst** ausgewertet wird.“

Es wird grundsätzlich zwischen **linksassoziativen** Operatoren, bei denen der **linke Operator** vor dem **rechten Operator** ausgewertet wird und **rechtsassoziativen** Operatoren, bei denen es genau anders rum ist unterschieden.<sup>a</sup>

<sup>a</sup>Parsing Expressions · Crafting Interpreters.

<sup>7</sup>Diese Option wurde beim **PicoC-Compiler** gewählt.

Bei **Assoziativität** ist z.B. der **Multiplikationsoperator** `*` ein Beispiel für einen **linksassoziativen** Operator und ein **Zuweisungsoperator** `=` ein Beispiel für einen **rechtsassoziativen** Operator. Dies ist in Abbildung 1.3 mithilfe von Klammern `()` veranschaulicht.

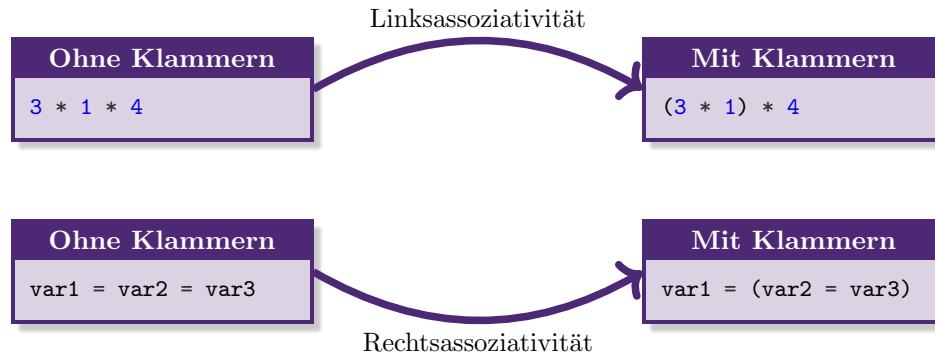


Abbildung 1.3: Veranschaulichung von Linksassoziativität und Rechtsassoziativität

#### Definition 1.28: Präzidenz

„Bestimmt, welcher Operator **zuerst** in einem Ausdruck, der eine Mischung **verschiedener** Operatoren enthält, ausgewertet wird. Operatoren mit einer **höheren Präzidenz**, werden **vor** Operatoren mit **niedrigerer Präzidenz** ausgewertet.“<sup>a</sup>

<sup>a</sup>Parsing Expressions · Crafting Interpreters.

Bei **Präzidenz** ist die Mischung der Operatoren für **Subtraktion** `-` und für **Multiplikation** `*` ein Beispiel für den Einfluss von Präzidenz. Dies ist in Abbildung 1.4 mithilfe der Klammern `()` veranschaulicht. Im Beispiel in Abbildung 1.4 ist bei den beiden **Subtraktionsoperatoren** `-` nacheinander und dem darauffolgenden **Multiplikationsoperator** `*` sowohl **Assoziativität** als auch **Präzidenz** im Spiel.



Abbildung 1.4: Veranschaulichung von Präzidenz

## 1.3 Lexikalische Analyse

Die **Lexikalische Analyse** bildet üblicherweise den ersten Filter innerhalb des **Pipe-Filter Architektur-patterns** (Definition 1.29) bei der Implementierung von Compilern. Die Aufgabe der lexikalischen Analyse ist vereinfacht gesagt, in einem Inputstring, z.B. dem Inhalt einer Datei, welche in **UTF-8** codiert ist, Folgen endlicher Symbole (auch **Wörter** genannt) zu finden, die bestimmte **Pattern** (Definition 1.30) matchen, die durch eine **reguläre Grammatik** spezifiziert sind. Diese Folgen endlicher Symbole werden auch **Lexeme** (Definition 1.31) genannt.

#### Definition 1.29: Pipe-Filter Architekturpattern

Ist ein **Architekturpattern**, welches aus **Pipes** und **Filtern** besteht, wobei der **Ausgang** eines **Filters** der **Eingang** des durch eine **Pipe** verbundenen adjazenten nächsten **Filters** ist, falls es einen gibt.

Ein **Filter** stellt einen Schritt dar, indem eine Eingabe **weiterverarbeitet** wird und **weitergereicht** wird. Bei der **Weiterverarbeitung** können Teile der Eingabe **entfernt**, **hinzugefügt** oder **vollständig ersetzt** werden.

Eine **Pipe** stellt ein **Bindeglied** zwischen zwei **Filtern** dar.<sup>a,b</sup>



<sup>a</sup>Das ein **Bindeglied** eine eigene Bezeichnung erhält, bedeutet allerdings nicht, dass es eine eigene wichtige **Aufgabe** erfüllt. Wie bei vielen **Pattern**, soll mit dem Namen des **Pattern**, in diesem Fall durch das **Pipe** die Anlehnung an z.B. die **Pipes aus Unix**, z.B. `cat /proc/bus/input/devices | less` zum Ausdruck gebracht werden. Und so banal es klingt, sollen manche Bezeichnungen von Pattern auch einfach nur gut klingen.

<sup>b</sup>Westphal, „Softwaretechnik“.

### Definition 1.30: Pattern

**Beschreibung** aller möglichen **Lexeme**, die eine Menge  $\mathbb{P}_T$  bilden und einem bestimmten **Token**  $T$  zugeordnet werden. Die Menge  $\mathbb{P}_T$  ist eine möglicherweise unendliche Menge von **Wörtern**, die sich mit den Produktionen einer **regulären Grammatik**  $G_{Lex}$  einer **regulären Sprache**  $L_{Lex}$  beschreiben lassen<sup>a</sup>, die für die Beschreibung eines **Tokens**  $T$  zuständig sind.<sup>b</sup>

<sup>a</sup>Als Beschreibungswerkzeug können aber auch z.B. reguläre Ausdrücke hergenommen werden.

<sup>b</sup>Thiemann, „Compilerbau“.

### Definition 1.31: Lexeme

Ein **Lexeme** ist ein **Wort** aus dem Inputstring, welches das **Pattern** für eines der **Token**  $T$  einer **Sprache**  $L_{Lex}$  **matched**.<sup>a</sup>

<sup>a</sup>Thiemann, „Compilerbau“.

Diese **Lexeme** werden vom **Lexer** (Definition 1.32) im **Inputstring** identifiziert und **Tokens**  $T$  zugeordnet. Das jeweils nächste **Lexeme** fängt dabei genau nach dem letzten Symbol des **Lexemes** an, das zuletzt vom **Lexer** erkannt wurde. Die **Tokens** (Definition 1.32) sind es, die letztendlich an die **Syntaktische Analyse** weitergegeben werden.

### Definition 1.32: Lexer (bzw. Scanner oder auch Tokenizer)

Ein **Lexer** ist eine **partielle Funktion**  $lex : \Sigma^* \rightarrow (N \times W)^*$ , welche ein **Wort** bzw. **Lexeme** aus  $\Sigma^*$  auf ein **Token**  $T$  mit einem **Tokennamen**  $N$  und einem **Tokenwert**  $W$  abbildet, falls dieses **Wort** sich unter der **regulären Grammatik**  $G_{Lex}$ , der **regulären Sprache**  $L_{Lex}$  ableiten lässt bzw. einem der **Pattern** der Sprache  $L_{Lex}$  entspricht.<sup>a</sup>

<sup>a</sup>Thiemann, „Compilerbau“.

Ein **Lexer** ist im Allgemeinen eine **partielle Funktion**, da es Zeichenfolgen geben kann, die kein **Pattern** eines **Tokens** der Sprache  $L_{Lex}$  **matchen**. In Bezug auf eine Implementierung, wird, wenn der Lexer Teil der Implementierung eines Compilers ist, in diesem Fall eine **Fehlermeldung** ausgegeben.

Um Verwirrung verzubäuen ist es wichtig folgende Unterscheidung hervorzuheben:

Wenn von **Symbolen** die Rede ist, so werden in der **Lexikalischen Analyse**, der **Syntaktische**



**Analyse** und der **Code Generierung**, auf diesen verschiedenen Ebenen unterschiedliche Konzepte als Symbole bezeichnet.

In der Lexikalischen Analyse sind einzelne **Zeichen eines Zeichensatzes** die Symbole.

In der Syntaktischen Analyse sind die **Tokennamen** die Symbole.

In der Code Generierung sind die **Bezeichner** (Definition 1.33) von **Variablen, Konstanten und Funktionen** die Symbole<sup>a</sup>.

<sup>a</sup>Das ist der Grund, warum die **Tabelle**, in der Informationen zu **Bezeichnern** gespeichert werden, in Kapitel 2 **Symboltabelle** genannt wird.

### Definition 1.33: Bezeichner (bzw. Identifier)

***Tokenwert**, der eine Konstante, Variable, Funktion usw. innerhalb ihres **Scopes eindeutig** benennt.<sup>a,b</sup>*

<sup>a</sup>Außer wenn z.B. bei Funktionen die Programmiersprache das **Überladen** erlaubt usw. In diesem Fall wird die **Signatur** der Funktion als weiteres Unterscheidungsmerkmal hinzugenommen, damit es eindeutig ist.

<sup>b</sup>Thiemann, „Einführung in die Programmierung“.

Eine weitere Aufgabe der **Lexikalischen Analyse** ist es jegliche für die Weiterverarbeitung unwichtigen Symbole, wie Leerzeichen `␣`, Newline `\n`<sup>8</sup> und Tabs `\t` aus dem Inputstring herauszufiltern. Das geschieht mittels des **Lexers**, der allen für die **Syntaktische Analyse** unwichtigen Zeichen das leere Wort  $\epsilon$  zuordnet. Das ist auch im Sinne der Definition, denn  $\epsilon \in (N \times W)^*$  ist immer der Fall beim **Kleene Stern Operator**  $*$ . Nur das, was für die **Syntaktische Analyse** wichtig ist, soll weiterverarbeitet werden, alles andere wird herausgefiltert.

Der Grund warum nicht einfach nur die **Lexeme** an die **Syntaktische Analyse** weitergegeben werden und der Grund für die Aufteilung des **Tokens** in **Tokenname** und **Tokenwert** ist, weil z.B. die Bezeichner von Variablen, Konstanten und Funktionen beliebige Zeichenfolgen sein können, wie `my_fun`, `my_var` oder `my_const` und es auch viele verschiedenen Zahlen gibt, wie 42, 314 oder 12. Die **Überbegriffe** bzw. **Tokennamen** für beliebige Bezeichner von Variablen, Konstanten und Funktionen und beliebige Zahlen sind aber trotz allem z.B. `NAME` und `NUM`<sup>9</sup>, bzw. wenn man sich nicht Kurzformen sucht `IDENTIFIER` und `NUMBER`. Für **Lexeme**, wie `if` oder `}` sind die **Tokennamen** bzw. Überbegriffe genau die Bezeichnungen, die man diesen Zeichenfolgen geben würde, nämlich `IF` und `RBRACE`.

Ein **Lexeme** ist damit aber nicht immer das gleiche, wie der **Tokenwert**, denn z.B. im Falle von PicoC kann der Wert 99 durch zwei verschiedene **Literale** (Definition 1.34) dargestellt werden, einmal als ASCII-Zeichen `'c'`, dass den entsprechenden Wert in der ASCII-Tabelle hat und des Weiteren auch in Dezimalschreibweise als 99<sup>10</sup>. Der **Tokenwert** ist jedoch der letztendlich verwendete Wert an sich, unabhängig von der Darstellungsform.

Die **Grammatik**  $G_{Lex}$ , die zur Beschreibung der Token  $T$  der Sprache  $L_{Lex}$  verwendet wird ist üblicherweise **regulär**, da ein typischer **Lexer** immer nur **ein Symbol** vorausschaut<sup>11</sup>, sich nichts merken muss und unabhängig davon, was für Symbole davor aufgetaucht sind läuft. Die Grammatik 2.1.1 liefert den Beweis, dass die Sprache  $L_{PicoC\_Lex}$  des **PicoC-Compilers** auf jeden Fall **regulär** ist, da sie fast die Definition 1.18 erfüllt. Einzig die Produktion `CHAR ::= "'ASCII_CHAR'"` sieht problematisch aus, kann allerdings auch

<sup>8</sup>In Unix Systemen wird für Newline das ASCII Symbol **line feed**, in Windows hingegen die ASCII Symbole **carriage return** und **line feed** nacheinander verwendet. Das wird aber meist durch die verwendete Programmiersprache, die man zur Implementierung des Lexers nutzt wegabstrahiert.

<sup>9</sup>Diese **Tokennamen** wurden im **PicoC-Compiler** verwendet, da man beim Programmieren möglichst **kurze** und **leicht verständliche** Bezeichner für seine Nodes haben will, damit unter anderem **mehr Code** in eine Zeile passt.

<sup>10</sup>Die Programmiersprache **Python** erlaubt es z.B. dieser Wert auch mit den Literalen `0b1100011` und `0x63` darzustellen.

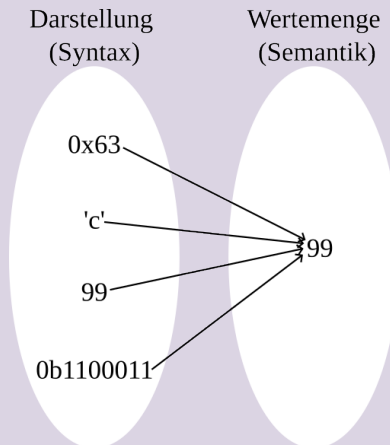
<sup>11</sup>Man nennt das auch einem **Lookahead** von 1



als  $\{\text{CHAR} ::= \text{" "CHAR2}, \text{CHAR2} ::= \text{ASCII\_CHAR" "}\}$  **regulär** ausgedrückt werden<sup>12</sup>. Somit existiert eine **reguläre Grammatik**, welche die **Sprache**  $L_{\text{PicoC\_Lex}}$  beschreibt und damit ist die **Sprache**  $L_{\text{PicoC\_Lex}}$  **regulär**.

#### Definition 1.34: Literal

Eine von möglicherweise vielen weiteren **Darstellungsformen** (als **Zeichenkette**) für ein und denselben **Wert** eines **Datentyps**.<sup>a</sup>



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

Um eine Gesamtübersicht über die **Lexikalische Analyse** zu geben, ist in Abbildung 1.5 die Lexikalische Analyse an einem Beispiel veranschaulicht.

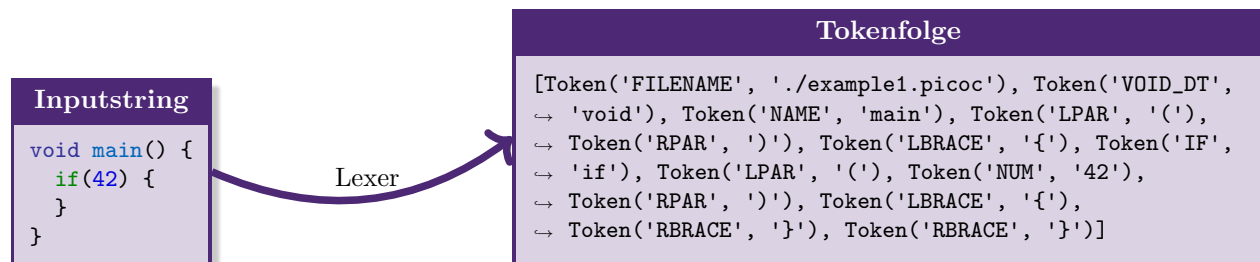


Abbildung 1.5: Veranschaulichung der Lexikalischen Analyse

## 1.4 Syntaktische Analyse

In der **Syntaktischen Analyse** ist für einige Sprachen eine **Kontextfreie Grammatik**  $G_{\text{Parse}}$  notwendig, um diese Sprachen zu beschreiben, da viele Programmiersprachen z.B. für **Funktionsaufrufe**  $\text{fun}(\text{arg})$  und **Codeblöcke**  $\text{if}(1)\{\}$  syntaktische Mittel verwenden, die es notwendig machen sich zu merken, wieviele öffnende runde Klammern '(' bzw. öffnende geschweifte Klammern '{' es momentan gibt, die noch nicht durch eine entsprechende schließende runde Klammer ')' bzw. schließende geschweifte Klammer '}' geschlossen wurden.

Die **Syntax**, in welcher ein **Programm** aufgeschrieben ist, wird auch als **Konkrete Syntax** (Definition 1.35) bezeichnet. In einem Zwischenschritt, dem **Parsen** wird aus diesem Programm mithilfe eines **Parsers** (Defini-

<sup>12</sup>Eine derartige Regel würde nur Probleme bereiten, wenn sich aus `ASCII.CHAR` **beliebig breite** Wörter ableiten lassen.

tion 1.37), ein **Ableitungsbaum** (Definition 1.36) generiert, der als Zwischenstufe hin zum einem **Abstrakter Syntaxbaum** (Definition 1.42) dient. Beim Compilerbau ist es förderlich kleinschrittig vorzugehen, deshalb erst die Generierung des **Ableitungsbaums** und dann erst des **Abstrakten Syntaxbaumes**.

#### Definition 1.35: Konkrete Syntax

*Steht für alles, was mit dem **Aufbau** von **Ableitungsbäumen** zu tun hat, also z.B. was für **Ableitungen** mit den **Grammatiken**  $G_{Lex}$  und  $G_{Parse}$  zusammengekommen möglich sind.*

*Ein **Programm** in seiner **Textrepräsentation**, wie es in einer Textdatei nach den Produktionen der **Grammatiken**  $G_{Lex}$  und  $G_{Parse}$  abgeleitet steht, bevor man es kompiliert, ist in **Konkreter Syntax** aufgeschrieben.<sup>a</sup>*

<sup>a</sup>G. Siek, *Course Webpage for Compilers (P423, P523, E313, and E513)*.

#### Definition 1.36: Ableitungsbaum (bzw. Konkretter Syntaxbaum, engl. Derivation Tree)

*Compilerinterne **Datenstruktur** für den **Formalen Ableitungsbaum** (Definition 1.25) eines in **Konkreter Syntax** geschriebenen Programmes.*

*Die **Konkrete Syntax** nach der sich der **Ableitungsbaum** richtet wird optimalerweise immer so definiert, dass sich möglichst einfach ein **Abstrakter Syntaxbaum** daraus konstruieren lässt.<sup>a</sup>*

<sup>a</sup>JSON parser - Tutorial — Lark documentation.

#### Definition 1.37: Parser

*Ein **Parser** ist ein Programm, dass aus einem Inputstring, der in **Konkreter Syntax** geschrieben ist, eine compilerinterne Darstellung, den **Ableitungsbaum** generiert, was auch als **Parsen** bezeichnet wird.<sup>a, b</sup>*

<sup>a</sup>Es gibt allerdings auch alternative Definitionen, denen nach ein Parser in Bezug auf Compilerbau ein Programm ist, dass einen Inputstring von **Konkreter Syntax** in **Abstrakte Syntax** übersetzt. Im Folgenden wird allerdings die Definition 1.37 verwendet.

<sup>b</sup>JSON parser - Tutorial — Lark documentation.

An dieser Stelle könnte möglicherweise eine Verwirrung entstehen, welche Rolle dann überhaupt ein **Lexer** hier spielt.

In Bezug auf Compilerbau ist ein **Lexer** ein Teil eines **Parsers**. Der **Lexer** ist ausschließlich für die **Lexikalische Analyse** verantwortlich und entspricht z.B., wenn man bei einem Wanderausflug verschiedenen Insekten entdeckt, dem Nachschlagen in einem Insekten**lexikon** und dem Aufschreiben, welchen Insekten man in welcher **Reihenfolge** begegnet ist. Zudem kann man bestimmte **Sehenswürdigkeiten** an denen man während des Ausflugs vorbeikommt ebenfalls festhalten, da es eine Rolle spielen kann in welchem örtlichen **Kontext** man den Insekten begegnet ist<sup>a</sup>.

Der **Parser** vereinigt sowohl die **Lexikalische Analyse**, als auch einen Teil der **Syntaktischen Analyse** in sich und entspricht, um auf das Beispiel zurückzukommen, dem Darstellen von **Beziehungen** zwischen den Insektenbeugnungen in einer für die **Weiterverarbeitung tauglichen Form**<sup>b</sup>.

In der Weiterverarbeitung kann der **Interpreter** das interpretieren und daraus bestimmte Schlüsse ziehen und ein **Compiler** könnte es vielleicht in eine für Menschen leichter entschlüsselbare Sprache kompilieren.

<sup>a</sup>Das würde z.B. der Rolle eines **Semikolon** ; in der Sprache  $L_{PicoC}$  entsprechen.

<sup>b</sup>Z.B. gibt es bestimmte **Wechselbeziehungen** zwischen Insekten, Insekten beeinflussen sich gegenseitig.

Die vom **Lexer** im Inputstring identifizierten **Token** werden in der **Syntaktischen Analyse** vom **Parser** als **Wegweiser** verwendet, da je nachdem, in welcher Reihenfolge die **Token** auftauchen, dies einer anderen Ableitung in der **Grammatik**  $G_{Parse}$  entspricht. Dabei wird in der Grammatik  $L_{Parse}$  nach dem **Tokennamen** unterschieden und nicht nach dem Tokenwert, da es nur von Interesse ist, ob an einer bestimmten Stelle z.B. eine **Zahl** steht und nicht, welchen konkreten Wert diese **Zahl** hat. Der **Tokenwert** ist erst später in der **Code Generierung** in 1.5 wieder relevant.

Ein **Parser** ist genauer gesagt ein erweiterter **Recognizer** (Definition 1.38), denn ein Parser löst das **Wortproblem** (Definition 1.20) für die **Sprache**, die durch die **Konkrete Syntax** beschrieben wird und konstruiert parallel dazu oder im Nachgang aus den Informationen, die während der Ausführung des Recognition Algorithmus gesichert wurden den **Ableitungsbaum**.

#### Definition 1.38: Recognizer (bzw. Erkenner)

*Entspricht dem Maschinenmodell eines **Automaten**. Im Bezug auf Compilerbau entspricht der **Recognizer** einem **Kellerautomaten**, in dem **Wörter** bestimmter **Kontextfreier Sprachen** erkannt werden. Der **Recognizer** ist ein **Algorithmus**, der erkennt, ob ein **Eingabewort** sich mit den **Produktionen der Konkrete Syntax** ableiten lässt, also ob er bzw. es Teil der Sprache ist, die von der **Konkreten Syntax** beschrieben wird oder nicht. Das vom **Recognizer** gelöste Problem ist auch als **Wortproblem** (Definition 1.20) bekannt.<sup>a</sup>*

<sup>a</sup>Thieman, „Compilerbau“.

Für das **Parsen** gibt es grundsätzlich **drei** verschiedene Ansätze:

- **Top-Down Parsing:** Der **Ableitungsbaum** wird von **oben-nach-unten** generiert, also von der **Wurzel** zu den **Blättern**. Dementsprechend fängt die Generierung des **Derivation Tree** mit dem **Startsymbol** der **Grammatik** an und wendet in jedem Schritt eine **Linksableitung** auf die **Nicht-Terminalsymbole** an, bis man **Terminalsymbole** hat, die sich zum gewünschten **Inputstring** abgeleitet haben oder sich herausstellt, dass dieser nicht abgeleitet werden kann.<sup>a</sup>

Der Grund, warum die **Linksableitung** verwendet wird und nicht z.B. die **Rechtsableitung**, ist, weil der **Eingabewert** bzw. der **Inputstring** von **links nach rechts** eingelesen wird, was gut damit zusammenpasst, dass die **Linksableitung** die **Blätter** von **links-nach-rechts** generiert.

Welche der **Produktionen** für ein **Nicht-Terminalsymbol** angewandt wird, wenn es mehrere **Alternativen** gibt, wird entweder durch **Backtracking** oder durch **Vorausschauen** gelöst.

Eine sehr einfach zu implementierende Technik für **Top-Down Parser** ist hierbei der **Rekursive Abstieg** (Definition ??).

Mit dieser Methode ist das Parsen **Linksrekursiver Grammatiken** (Definition 1.24) allerdings nicht möglich, ohne die Grammatik vorher umgeformt zu haben und jegliche **Linksrekursion** aus der **Grammatik** entfernt zu haben, da diese zu **Unendlicher Rekursion** führt.

**Rekursiver Abstieg** kann mit **Backtracking** verbunden werden, um auch Grammatiken parsen zu können, die nicht **LL(k)** (Definition ??) sind. Dabei werden meist nach dem **Depth-First-Search Prinzip** alle **Produktionen** für ein **Nicht-Terminalsymbol** solange durchgegangen bis der gewünschte Inputstring abgeleitet ist oder alle **Alternativen** für einen Schritt abgesucht sind, bis man wieder beim ersten Schritt angekommen ist und da auch alle **Alternativen** abgesucht sind, was dann bedeutet, dass der **Inputstring** sich **nicht** mit der verwendeten Grammatik

ableiten lässt.<sup>b</sup>

Wenn man eine **LL(k)** Grammatik hat, kann man auf **Backtracking verzichten** und es reicht einfach nur immer *k* **Token** im Inputstring **vorauszuschauen**. **Mehrdeutige Grammatiken** sind dadurch ausgeschlossen, weil **LL(k)** keine **Mehrdeutigkeit** zulässt.<sup>c</sup>

- **Bottom-Up Parsing:** Es wird mit dem **Eingabewort** gestartet und versucht **Rechtsableitungen** entsprechend der **Produktionen** der **Konkreten Syntax** rückwärts anzuwenden, bis man beim **Startsymbol** landet.<sup>d</sup>
- **Chart Parsing:** Es wird **Dynamische Programmierung** verwendet und **partielle Zwischenergebnisse** werden in einer **Tabelle** (bzw. einem **Chart**) gespeichert und können **wiederverwendet** werden. Das macht das Parsen **Kontextfreier Grammatiken** effizienter, sodass es nur noch **polynomielle** Zeit braucht, da **Backtracking** nicht mehr notwendig ist<sup>e</sup>. **Chart Parser** können dabei **top-down** oder **bottom-up** Ansätze umsetzen. Da die **Implementierung** von **Chart Parsern** fundamental anders ist als bei **Top-Down** und **Bottom-Up Parsern**, wird diese **Kategorie** von Parsern nochmal **speziell unterschieden** und nicht gesagt, es sei ein **Top-Down Parser** oder **Bottom-Up Parser**, der **Dynamische Programmierung** verwendet.

<sup>a</sup>What is Top-Down Parsing?

<sup>b</sup>Diese Form von Parsing wurde im **PicoC-Compiler** implementiert, als dieser noch auf dem Stand des **Bachelorprojektes** war, bevor er durch den nicht selbst implementierten **Earley Parser** von **Lark** (siehe *Lark - a parsing toolkit for Python*) ersetzt wurde.

<sup>c</sup>Diese Art von Parser ist im **RETI-Interpreter** implementiert, da die **RETI-Sprache** eine besonders simple **LL(1)** Grammatik besitzt. Diese Art von **Parser** wird auch als **Predictive Parser** oder **LL(k) Recursive Descent Parser** bezeichnet, wobei **Recursive Descent** das englische Wort für **Rekursiven Abstieg** ist.

<sup>d</sup>What is Bottom-up Parsing?

<sup>e</sup>Der **Earley Parser**, den **Lark** und damit der **PicoC-Compiler** verwendet fällt unter diese Kategorie.

Der **Abstrakter Syntaxbaum** wird mithilfe von **Transformern** (Definition 1.39) und **Visitors** (Definition 1.40) generiert und ist das Endprodukt der **Syntaktischen Analyse**, welches an die **Code Generierung** weitergegeben wird. Wenn man die gesamte **Syntaktische Analyse** betrachtet, so übersetzt diese ein Programm von der **Konkreten Syntax** in die **Abstrakte Syntax** (Definition 1.41).

#### Definition 1.39: Transformer

Ein Programm, dass von **unten-nach-oben**, nach dem **Breadth First Search** Prinzip alle Knoten des **Ableitungsbaum** besucht und beim Antreffen eines bestimmten Knoten des **Derivation Tree** einen entsprechenden Knoten des **Abstrakter Syntaxbaum** erzeugt und diesen anstelle des Knotens des **Derivation Tree** setzt und so Stück für Stück den **Abstrakter Syntaxbaum** konstruiert.<sup>a</sup>

<sup>a</sup>Transformers & Visitors — Lark documentation.

#### Definition 1.40: Visitor

Ein Programm, dass von **unten-nach-oben**, nach dem **Breadth First Search** Prinzip alle Knoten des **Ableitungsbaum** besucht und in Bezug zu Compilerbau, beim Antreffen eines bestimmten **Knoten** des **Derivation Tree**, diesen **in-place** mit anderen Knoten **tauscht** oder **manipuliert**, um den **Derivation Tree** für die weitere Verarbeitung durch z.B. einen **Transformer** zu vereinfachen.<sup>a,b</sup>

<sup>a</sup>Kann theoretisch auch zur Konstruktion eines **Abstrakter Syntaxbaum** verwendet werden, wenn z.B. eine externe Klasse verwendet wird, welches für die Konstruktion des **Abstrakter Syntaxbaum** verantwortlich ist. Aber dafür ist ein **Transformer** besser geeignet.

<sup>b</sup>Transformers & Visitors — Lark documentation.

**Definition 1.41: Abstrakte Syntax**

Steht für alles, was mit dem **Aufbau** von **Abstrakten Syntaxbäumen** zu tun hat, also z.B. was für Arten von **Kompositionen** mit den **Knoten** eines **Abstrakten Syntaxbaums** möglich sind.

Ein **Abstract Syntax Tree**, der zur Kompilierung eines Wortes<sup>a</sup> generiert wurde, ist nach einer **Abstrakter Syntax** konstruiert.

Jene Produktionen, die in der **Konkreten Syntax** für die Umsetzung von **Präzidenz** notwendig waren, sind in der **Abstrakten Syntax** abgeflacht. Dadurch sind die **Kompositionen**, welche die Knoten im **Abstract Syntax Tree** bilden können **syntaktisch** meist näher zur Syntax von **Maschinenbefehlen**.<sup>b</sup>

<sup>a</sup>Z.B. **Programmcode**.

<sup>b</sup>G. Siek, *Course Webpage for Compilers (P423, P523, E313, and E513)*.

**Definition 1.42: Abstrakter Syntaxbaum (bzw. engl. Abstract Syntax Tree, kurz AST)**

Ist ein **compilerinterne Datenstruktur**, welche eine **Abstraktion** eines dazugehörigen **Ableitungsbaumes** darstellt, in dessen Aufbau auch das Erfordernis eines **leichten Zugriffs** und einer **leichten Weiterverarbeitbarkeit** eingeflossen ist. Bei der Betrachtung eines **Knoten**, der für einen Teil des Programms steht, soll man möglichst schnell die Fragen beantworten können, welche **Funktionalität** der Sprache dieser umsetzt, welche **Bestandteile** er hat und welche **Funktionalität** der Sprache diese Bestandteile umsetzen usw.

Im Gegensatz zum **Formalen Ableitungsbaum**, ergibt es beim **Abstrakten Syntaxbaum** keinen Sinn zusätzlich einen **Formalen Abstrakten Syntaxbaum** zu unterscheiden, da das Konzept eines **Abstrakten Syntaxbaumes** ohne eine Datenstruktur zu sein für sich allein gesehen keine Sinn hat. Wenn von **Abstrakten Syntaxbäumen** die Rede ist, ist immer eine **Datenstruktur** gemeint.

Die **Abstrakte Syntax** nach der sich der **Abstrakte Syntaxbaum** richtet wird optimalerweise immer so definiert, dass der **Abstrakte Syntaxbaum** in den darauffolgenden Verarbeitungsschritten<sup>a</sup> möglichst **einfach weiterverarbeitet** werden kann.

<sup>a</sup>Die verschiedenen **Passes**.

In Abbildung 1.6 wird das Beispiel aus Unterkapitel 1.2.1 fortgeführt, welches den **Arithmetischen Ausdruck**  $4 * 2$  in Bezug auf die Grammatik 1.1, welche die **höhere Präzidenz** der **Multiplikation**  $*$  berücksichtigt in einem **Ableitungsbaum** darstellt. In Abbildung 1.6 wird der Ableitungsbaum zu einem Abstrakten Syntaxbaum **abstrahiert**. Das geschieht bezogen auf das Beispiel aus Unterkapitel 1.2.1, indem jegliche Knoten, die im **Ableitungsbaum** nur existieren, weil die Grammatik so umgesetzt ist, dass es nur **einen** einzigen möglichen **Ableitungsbaum** geben kann **wegabstrahiert** werden.



**Abbildung 1.6:** Veranschaulichung des Unterschieds zwischen Ableitungsbaum und Abstraktem Syntaxbaum.

Die **Baumdatenstruktur** des **Ableitungsbaumes** und **Abstrakten Syntaxbaumes** ermöglicht es die

Operationen, die ein Compiler bzw. Interpreter bei der Weiterverarbeitung des Programmes ausführen muss möglichst **effizient** auszuführen und auf **unkomplizierte** Weise direkt zu erkennen, welche er ausführen muss.

Um eine Gesamtübersicht über die **Syntaktische Analyse** zu geben, ist in Abbildung 1.7 die Syntaktische mit dem Beispiel aus Subkapitel 1.3 fortgeführt.

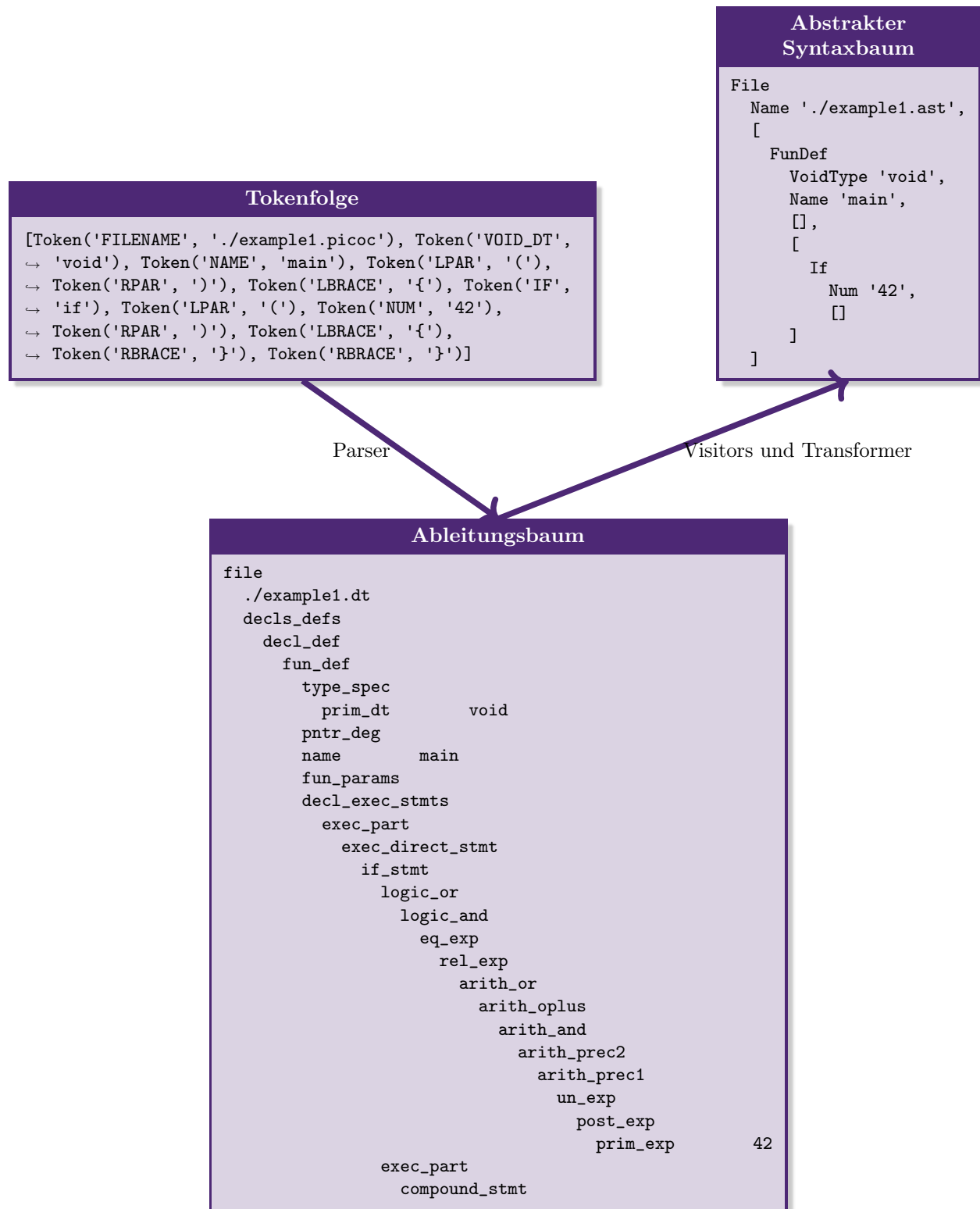


Abbildung 1.7: Veranschaulichung der Syntaktischen Analyse



## 1.5 Code Generierung

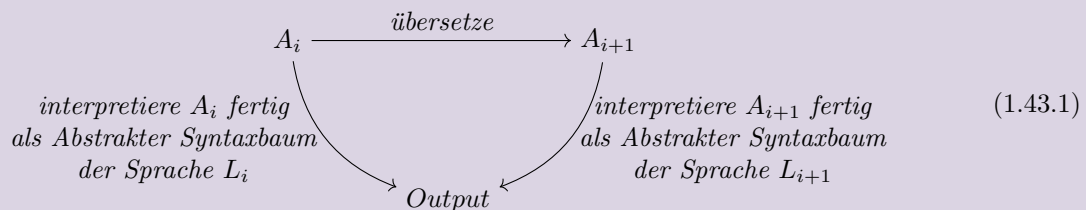
In der **Code Generierung** steht man nun dem Problem gegenüber einen **Abstrakter Syntaxbaum** einer Sprache  $L_1$  in den **Abstrakter Syntaxbaum** einer Sprache  $L_2$  umformen zu müssen. Dieses Problem lässt sich vereinfachen, indem man das Problem in mehrere Schritte unterteilt, die man **Passes** (Definition 1.43) nennt. So wie es auch schon mit dem **Derivation Tree** in der Syntaktischen Analyse gemacht wurde, den man als Zwischenstufe zum **Abstrakter Syntaxbaum** konstruiert hatte. Aus dem Derivation konnte, dann unkompliziert und einfach mit **Transformern** und **Visitors** ein **Abstrakter Syntaxbaum** generiert werden.

Man spricht hier von dem „**Abstrakten Syntaxbaum einer Sprache  $L_1$  (bzw.  $L_2$ )**“ und meint hier mit der Sprache  $L_1$  (bzw.  $L_2$ ) **nicht** die Sprache, welche durch die **Abstrakte Syntax**, nach welcher der **Abstrakte Syntaxbaum** abgeleitet ist beschrieben wird. Es ist damit **immer** die Sprache gemeint, die **kompiliert** werden soll und zu deren Zweck der **Abstrakt Syntax Tree** überhaupt konstruiert wird. Für die tatsächliche Sprache, die durch die **Abstrakt Syntax** beschrieben wird, interessiert man sich nie wirklich explizit. Diese **Redeart** wurde aus der **Quelle** G. Siek, *Course Webpage for Compilers (P423, P523, E313, and E513)* übernommen.

### Definition 1.43: Pass

*Einzelner Übersetzungsschritt in einem Kompiliervorgang von einem beliebigen **Abstrakten Syntaxbaum**  $A_i$  einer Sprache  $L_i$  zu einem **Abstrakten Syntaxbaum**  $A_{i+1}$  einer Sprache  $L_{i+1}$ , der meist **eine** bestimmte **Teilaufgabe** übernimmt, die sich mit keiner **Teilaufgabe** eines anderen **Passes** überschneidet und möglichst wenig **Ähnlichkeit** mit den **Teilaufgaben** anderer **Passes** haben sollte.<sup>ab</sup>*

*Für jeden **Pass** und für einen beliebigen **Abstrakten Syntaxbaum**  $A_i$  gilt ähnlich, wie bei einem **vollständigen Compiler** in 1.43.1, dass:*



*wobei man hier so tut, als gäbe es zwei **Interpreter** für die zwei Sprachen  $L_i$  und  $L_{i+1}$ , welche den jeweiligen **Abstrakten Syntaxbaum**  $A_i$  bzw.  $A_{i+1}$  fertig interpretieren.<sup>cd</sup>*

<sup>a</sup>Ein **Pass** kann mit einem **Transpiler** ?? (Definition ??) verglichen werden, da sich die zwei Sprachen  $L_i$  und  $L_{i+1}$  aufgrund der **Kleinschrittigkeit** meist auf einem ähnlichen **Abstraktionslevel** befinden. Der Unterschied ist allerdings, dass ein **Transpiler** zwei Programme, die in  $L_i$  bzw.  $L_{i+1}$  geschrieben sind kompiliert. Ein **Pass** ist dagegen immer **kleinschrittig** und operiert ausschließlich auf **Abstrakten Syntaxbäumen**, ohne Parsing usw.

<sup>b</sup>Der Begriff kommt aus dem **Englischen** von „passing over“, da der gesamte **Abstrakte Syntaxbaum** in einem **Pass** durchlaufen wird.

<sup>c</sup>**Interpretieren** geht immer von einem Programm in **Konkreter Syntax** aus, wobei der **Abstrakte Syntaxbaum** ein **Zwischenschritt** bei der **Interpretierung** ist.

<sup>d</sup>G. Siek, *Course Webpage for Compilers (P423, P523, E313, and E513)*.

Die von den **Passes** umgeformten **Abstrakter Syntaxbaums** sollten dabei mit jedem **Pass** der **Syntax** von **Maschinenbefehlen** immer ähnlicher werden, bis es schließlich nur noch Maschinenbefehle sind.



### 1.5.1 Monadische Normalform

Hat man es mit einer Sprache zu tun, welche **Unreine Ausdrücke** (Definition 1.45) besitzt, so ist es sinnvoll einen **Pass** einzuführen, der **Reine** (Definition 1.44) und **Unreine Ausdrücke** voneinander **trennt**. Das wird erreicht, indem man aus den Unreinen Ausdrücken **vorangestellte Statements** macht, die man **vor** den jeweiligen reinen Ausdruck, mit dem sie **gemischt** waren stellt. Der Unreine Ausdruck muss als **erstes** ausgeführt werden, für den Fall, dass der **Effekt**, denn ein **Unreiner Ausdruck** hatte den **Reinen Ausdruck**, mit dem er gemischt war in irgendeinerweise beeinflussen könnte.

#### Definition 1.44: Reiner Ausdruck (bzw. engl. pure expression)

*Ein **Reiner Ausdruck** ist ein Ausdruck, der **rein** ist. Das bedeutet, dass dieser Ausdruck **keine Nebeneffekte** erzeugt. Ein **Nebeneffekt** ist eine **Bedeutung**, die ein Ausdruck hat, die sich **nicht** mit **RETI-Code** darstellen lässt.<sup>a,b</sup>*

<sup>a</sup>Sondern z.B. **intern** etwas am **Kompilierprozess** ändert.

<sup>b</sup>G. Siek, *Course Webpage for Compilers (P423, P523, E313, and E513)*.

#### Definition 1.45: Unreiner Ausdruck

*Ein **Unreiner Ausdruck** ist ein Ausdruck, der kein **Reiner Ausdruck** ist.*

Auf diese Weise sind alle **Statements** und **Ausdrücke** in **Monadischer Normalform** (Definition 1.46).

#### Definition 1.46: Monadische Normalform (bzw. engl. monadic normal form)

*Ein **Statement** oder **Ausdruck** ist in **Monadischer Normalform**, wenn er nach einer **Konkreten Syntax** in **Monadischer Normalform** abgeleitet wurde.*

*Eine **Konkrete Syntax** ist in **Monadischer Normalform**, wenn sie **reine Ausdrücke** und **unreine Ausdrücke nicht** miteinander **mischt**, sondern voneinander **trennt**.<sup>a</sup>*

*Eine **Abstrakte Syntax** ist in **Monadischer Normalform**, wenn die **Konkrete Syntax** für welche sie definiert wurde in **Monadischer Normalform** ist.*

<sup>a</sup>G. Siek, *Course Webpage for Compilers (P423, P523, E313, and E513)*.

Ein **Beispiel** für dieses Vorgehen ist in Abbildung 1.8 zu sehen, wo der Einfachheit halber auf die Darstellung in **Abstrakter Syntax** verzichtet wurde und die Codebeispiele in der entsprechenden **Konkreten Syntax**<sup>13</sup> aufgeschrieben wurden.

In der Abbildung 1.8 ist der Ausdruck mit dem **Nebeneffekt** eine Variable zu **allokieren**: `int var`, mit dem Ausdruck für eine **Zuweisung** `exp = 5 % 4` gemischt, daher muss der **Unreine** Ausdruck als eigenständiges Statement **vorangestellt** werden.

<sup>13</sup>Für deren Kompilierung die **Abstrakte Syntax** überhaupt definiert wurde.



Abbildung 1.8: Codebeispiel für das Trennen von Ausdrücken mit und ohne Nebeneffekten

Die Aufgabe eines solchen **Passes** ist es, den **Abstrakter Syntaxbaum** der **Syntax** von **Maschinenbefehlen** anzunähern, indem Subbäume vorangestellt werden, die keine Entsprechung in **RETI-Knoten** haben. Somit wird eine **Seperation** von Subbäumen, die keine Entsprechung in **RETI-Knoten** haben und denen, die eine haben bewerkstelligt wird. Ein **Reiner Ausdruck** ist **Maschinenbefehlen** ähnlicher als ein Ausdruck, indem ein **Reiner** und **Unreiner Ausdruck** gemischt sind. Somit sparrt man sich in der Implementierung **Fallunterscheidungen**, indem die **Reinen Ausdrücke** direkt in **RETI-Code** übersetzt werden können und **nicht** unterschieden werden muss, ob darin **Unreine Ausdrücke** vorkommen.

### 1.5.2 A-Normalform

Im Falle dessen, dass es sich bei der **Sprache**  $L_1$  um eine **höhere Programmiersprache** und bei  $L_2$  um **Maschinensprache** handelt, ist es fast unerlässlich einen **Pass** einzuführen, der **Komplexe Ausdrücke** (Definition 1.49) aus **Statements** und **Ausdrücken** entfernt. Das wird erreicht, indem man aus den Komplexen Ausdrücken **vorangestellte** Statements macht, in denen die **Komplexen Ausdrücke temporären Locations** zugewiesen werden (Definiton 1.47) und dann anstelle des **Komplexen Ausdrucks** auf die jeweilige **temporäre Location** zugegriffen wird.

Sollte in dem **Statement**, indem der **Komplexe Ausdruck** einer **temporären Location** zugewiesen wird, der Komplexe Ausdruck **Teilausdrücke** enthalten, die **komplex** sind, muss die gleiche Prozedur erneut für die **Teilausdrücke** angewandt werden, bis **Komplexe Ausdrücke** nur noch in Statements zur Zuweisung an Locations auftauchen, aber die Komplexen Ausdrücke nur **Atomare Ausdrücke** (Definiton 1.48) enthalten.

Sollte es sich bei dem **Komplexen Ausdruck** um einen **Unreinen Ausdruck** handeln, welcher nur einen **Nebeneffekt** ausführt und sich nicht in **RETI-Befehle** übersetzt, so wird aus diesem ein **vorangestelltes Statement** gemacht, welches einfach nur den **Nebeneffekt** dieses **Unreinen Ausdrucks** ausführt.

#### Definition 1.47: Location

*Kollektiver Begriff für **Variablen**, **Attribute** bzw. **Elemente** von Variablen bestimmter Datentypen, **Speicherbereiche auf dem Stack**, die **temporäre Zwischenergebnisse** speichern und **Register**.*

*Im Grunde genommen alles, was mit einem **Programm zu tun** hat und irgendwo **gespeichert** ist oder als **Speicherort** dient.<sup>a</sup>*

<sup>a</sup>G. Siek, *Course Webpage for Compilers (P423, P523, E313, and E513)*.

Auf diese Weise sind alle **Statements** und **Ausdrücke** in **A-Normalform** (Definition 1.50). Wenn eine **Konkrete Syntax** in **A-Normalform** ist, ist diese auch automatisch in **Monadischer Normalform** (Definition 1.50), genauso, wie ein **Atomarer Ausdruck** auch ein **Reiner Ausdruck** ist (nach Definition 1.48).

**Definition 1.48: Atomarer Ausdruck**

Ein **Atomarer Ausdruck** ist ein Ausdruck, der ein **Reiner Ausdruck** ist und der in eine **Folge von RETI-Befehlen** übersetzt werden kann, die **atomar** ist, also **nicht** mehr weiter in kleinere Folgen von RETI-Befehlen **zerkleinert** werden kann, welche die **Übersetzung** eines anderen Ausdrucks sind.

Also z.B. im Fall der Sprache  $L_{PicoC}$  entweder eine **Variable** `var`, eine **Zahl** `12`, ein **ASCII-Zeichen** `'c'` oder ein **Zugriff auf eine Location**, wie z.B. `stack(1)`.<sup>a</sup>

<sup>a</sup>G. Siek, *Course Webpage for Compilers (P423, P523, E313, and E513)*.

**Definition 1.49: Komplexer Ausdruck**

Ein **Komplexer Ausdruck** ist ein **Ausdruck**, der **nicht atomar** ist, wie z.B. `5 % 4`, `-1`, `fun(12)` oder `int var`.<sup>ab</sup>

<sup>a</sup>`int var` ist eine **Allokation**.

<sup>b</sup>G. Siek, *Course Webpage for Compilers (P423, P523, E313, and E513)*.

**Definition 1.50: A-Normalform (ANF)**

Ein **Statement** oder **Ausdruck** ist in **A-Normalform**, wenn er nach einer **Konkreten Syntax** in **A-Normalform** abgeleitet wurde.

Eine **Konkrete Syntax** ist in **A-Normalform**, wenn sie in **Monadischer Normalform** ist und wenn alle **Komplexen Ausdrücke** nur **Atomare Ausdrücke** enthalten und einer **Location** zugewiesen sind.

Eine **Abstrakte Syntax** ist in **A-Normalform**, wenn die **Konkrete Syntax** für welche sie definiert wurde in **A-Normalform** ist.<sup>abc</sup>

<sup>a</sup>A-Normalization: *Why and How (with code)*.

<sup>b</sup>Bolingbroke und Peyton Jones, „Types are calling conventions“.

<sup>c</sup>G. Siek, *Course Webpage for Compilers (P423, P523, E313, and E513)*.

Ein **Beispiel** für dieses Vorgehen ist in Abbildung 1.9 zu sehen, wo der Einfachheit halber auf die Darstellung in **Abstrakter Syntax** verzichtet wurde und die Codebeispiele in der entsprechenden **Konkreten Syntax**<sup>14</sup> aufgeschrieben wurden.

Der **PicoC-Compiler** nutzt, anders als es geläufig ist keine **Register** und **Graph Coloring** (Definition ??) inklusive **Liveness Analysis** (Definition ??) usw., um Werte von Variablen, temporäre Zwischenergebnisse usw. abzuspeichern, sondern immer nur den **Hauptspeicher**, wobei **temporäre Zwischenergebnisse** auf den **Stack** gespeichert werden.<sup>15</sup>

Aus diesem Grund verwendet das Beispiel in Abbildung 1.9 eine andere Definition für **Komplexe** und **Atomare Ausdrücke**, da dieses Beispiel, um später keine Verwirrung zu erzeugen der Art nachempfunden ist, wie im **PicoC-ANF Pass** der **Abstrakter Syntaxbaum** umgeformt wird. Weil beim PicoC-Compiler **temporäre Zwischenergebnisse** auf den **Stack** gespeichert werden, wird nur noch ein **Zugriffen auf den Stack**, wie z.B. `stack('1')` als **Atomarer Ausdruck** angesehen. Dementsprechend werden **Ausdrücke** für **Zahl** `4`, **Variable** `var` und **ASCII-Zeichen** `'c'` nun ebenfalls zu den **Komplexen Ausdrücken** gezählt.

Im Fall, dass **Register** für z.B. **temporäre Zwischenergebnisse** genutzt werden und der **Maschinen-**

<sup>14</sup>Für deren Kompilierung die **Abstrakte Syntax** überhaupt definiert wurde.

<sup>15</sup>Die in diesem **Paragraph** erwähnten **Begriffe** werden nur grob erläutert, da sie für den **PicoC-Compiler** keine Rolle spielen. Aber sie wurden erwähnt, damit in dieser **Bachelorarbeit** auch das übliche Vorgehen Erwähnung findet und vom Vorgehen beim **PicoC-Compiler** abgegrenzt werden kann.

**befehlssatz** es erlaubt **zwei Register** miteinander zu verrechnen<sup>16</sup>, ist es möglich **Ausdrücke** für **Zahl** 4, **Variable** *var* und **ASCII-Zeichen** 'c' als **atomar** zu definieren, da sie mit einem **Maschinenbefehl** verarbeitet werden können<sup>17</sup>. Werden allerdings keine **Register** für **Zwischenergebnisse** genutzt werden, braucht man **mehrere Maschinenbefehle**, um die Zwischenergebnisse vom **Stack** zu holen, zu **verrechnen** und das Ergebnis wiederum auf den **Stack** zu **speichern** und das SP-Register **anzupassen**. Daher werden die **Ausdrücke** für **Zahl** 4, **Variable** *var* und **ASCII-Zeichen** 'c' als **Komplexe Ausdrücke** gewertet, da sie niemals in einem **Maschinenbefehl** miteinander verrechnet werden können.

Die Statements 4, x, usw. für sich sind in diesem Fall **Statements**, bei denen ein **Komplexer Ausdruck** einer **Location**, in diesem Fall einer **Speicherzelle des Stack** zugewiesen wird, da 4, x usw. in diesem Fall auch als **Komplexe Ausdrücke** zählen. Auf das Ergebnis dieser **Komplexen Ausdrücke** wird mittels **stack(2)** und **stack(1)** zugegriffen, um diese im **Komplexen Ausdruck** **stack(2) % stack(1)** miteinander zu verrechnen und wiederum einer Speicherzelle des Stack zuzuweisen.

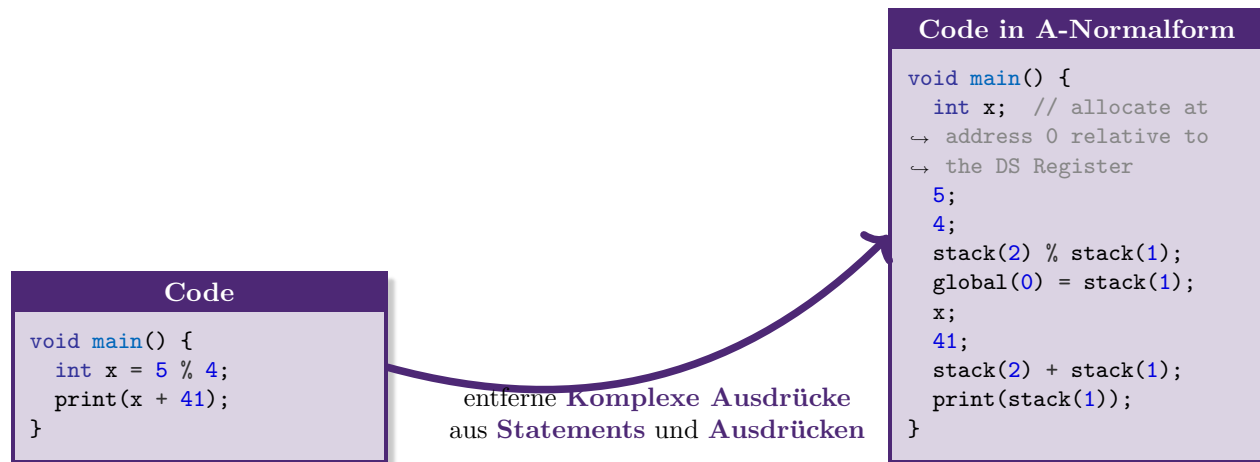


Abbildung 1.9: Codebeispiel für das Entfernen Komplexer Ausdrücke aus Operationen

Ein solcher **Pass** hat vor allem in erster Linie die Aufgabe den **Abstrakt Syntax Tree** der **Syntax** von **Maschinenbefehlen** besonders dadurch anzunähern, dass er auf der Ebene der Konkreten Syntax die Statements **weniger komplex** macht und diese dadurch den ziemlich **einfachen Maschinenbefehlen** syntaktisch ähnlicher sind. Des Weiteren **vereinfacht** dieser Pass die **Implementierung** der nachfolgenden Passes enorm, da Statements z.B. nur noch die Form **global(rel\_addr) = stack(1)** haben, die viel **einfacher verarbeitet** werden kann.

Alle weiteren denkbaren **Passes** sind zu **spezifisch** auf bestimmte **Statements** und **Ausdrücke** ausgelegt, als das sich zu diesen allgemein etwas mit einer **Theorie** dahinter sagen lässt. Alle **Passes**, die zur Implementierung des **PicoC-Compilers** geplant und ausgedacht wurden sind im Unterkapitel 2.3.1 definiert.

### 1.5.3 Ausgabe des Maschinencodes

Nachdem alle **Passes** durchgearbeitet wurden, ist es notwendig aus dem finalen **Abstrakter Syntaxbaum** den eigentlichen **Maschinencode** in **Konkreter Syntax** zu generieren. In üblichen Compilern wird hier für den **Maschinencode** eine **binäre Repräsentation** gewählt<sup>18</sup>. Der Weg von **Abstrakter Syntax** zu **Konkreter Syntax** ist allerdings wesentlich einfacher, als der Weg von der **Konkreten Syntax**

<sup>16</sup>Z.B. **Addieren** oder **Subtraktion** von zwei **Registerinhalten**.

<sup>17</sup>Mit dem **RETI-Befehlssatz** wäre das durchaus möglich, durch z.B. **MULT ACC IN2**.

<sup>18</sup>Da der **PicoC-Compiler** vor allem zu **Lernzwecken** konzipiert ist, wird bei diesem der **Maschinencode** allerdings in einer **menschenlesbaren Repräsentation** ausgegeben.

zur **Abstrakten Syntax**, für die eine gesamte **Syntaktische Analyse**, die eine **Lexikalische Analyse** beinhaltet durchlaufen werden musste.

Jeder **Knoten** des **Abstrakter Syntaxbaums** erhält dazu eine Methode, welche hier `to_string` genannt wird, die eine **Textrepräsentation** seiner selbst und all seiner Knoten mit an den richtigen Stellen passend gesetzten **Semikolons** ; usw. ausgibt. Dabei wird nach dem **Depth-First-Search** Schema der gesamte **Abstract Sybtax Tree** durchlaufen und die Methode `to_string` zur Ausgabe der **Textrepräsentation** der verschiedenen Knoten aufgerufen, die immer wiederum die Methode `to_string` ihrer Kinder aufrufen und die zurückgegebene Textrepräsentation passend **zusammenfügen** und selbst **zurückgeben**.

## 1.6 Fehlermeldungen

Wenn bei einem Compiler ein **unerwünschtes Verhalten** der folgenden **Kategorien**<sup>19</sup> eintritt:

1. der **Parser**<sup>20</sup> entscheidet das **Wortproblem** für ein **Eingabeprogramm**<sup>21</sup> mit 0, also das **Eingabeprogramm** befolgt nicht die **Syntax** der **Sprache** des Compilers<sup>22</sup>.
2. in den **Passes** tritt eine Fall ein, der **nicht** in der **Semantik** der Sprache des Compilers abgedeckt ist, z.B.:
  - eine **Variable** wird **verwendet**, obwohl sie noch **nicht deklariert** ist.
  - bei einem **Funktionsaufruf** werden **mehr** Argumente oder Argumente des **falschen Datentyps** übergeben, als in der **Funktionsdeklaration** oder **Funktionsdefinition** angegeben ist.
3. Während der **Laufzeit** des Compilers tritt ein Ereignis ein, das **nicht** durch die **Semantik** der Sprache des Compilers abgedeckt ist oder das **Betriebssystem** nicht erlaubt, z.B.:
  - eine **nicht erlaubte Operation**, wie **Division durch 0** (z.B.  $42 / 0$ ) soll ausgeführt werden.
  - **Segmentation Fault**: Wenn auf **Speicher zugegriffen** wird, der vom **Betriebssystem geschützt** ist.

oder während des des **Linkens** (Definition ??) etwas nicht zusammenpasst, wie z.B.:

- es gibt **keine** oder **mehr als eine** `main`-Funktion
- eine Funktion, die in einer **Objektdatei** (Definition ??) benötigt wird, wird von **keiner** anderen oder **mehr als einer** Objektdatei bereitgestellt

wird eine **Fehlermeldung** (Definition 1.51) ausgegeben.

### Definition 1.51: Fehlermeldung

*Benachrichtigung beliebiger Form, die einen **Grund** angibt weshalb ein Programm **nicht weiter ausgeführt** werden kann<sup>a</sup>. Das **Ausgeben** einer Fehlermeldung kann dabei auf **verschiedene Weisen** erfolgen, wie z.B.*

- über `stdout` oder `stderr` im einem **Terminal Emulator** oder **richtigen Terminal**<sup>b</sup>.

<sup>19</sup> *Errors in C/C++ - GeeksforGeeks.*

<sup>20</sup> Bzw. der **Recognizer** im Parser.

<sup>21</sup> Bzw. **Wort**.

<sup>22</sup> Bzw. das **Eingabeprogramm** lässt sich **nicht** mit der Grammatik des Compilers **ableiten**.

- über eine *Dialogbox* in einer *Graphischen Benutzerfläche*<sup>c</sup> oder *Zeichenorientierten Benutzerschnittstelle*<sup>d</sup>.
- in ein *Register* oder an eine *spezielle Adresse des Hauptspeichers* wird ein *Wert geschrieben*.
- *Logdatei*<sup>e</sup> auf einem *Speichermedium*.

<sup>a</sup>Dieses Programm kann z.B. ein **Compiler** sein oder ein **Programm**, dass dieser **Compiler** selbst **kompiliert** hat.

<sup>b</sup>Nur unter **Linux**, **Windows** hat sowas nicht.

<sup>c</sup>In engl. **G**raphical **U**ser **I**nterface, kurz **GUI**.

<sup>d</sup>In engl. **T**ext-based **U**ser **I**nterface, kurz **TUI**.

<sup>e</sup>In engl. **log file**.

# 2 Implementierung

In diesem Kapitel wird, nachdem im Kapitel 1 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 1, werden in den folgenden Unterkapiteln die einzelnen **Zwischenschritte** vom einem **Programm** in der **Konkreten Syntax** der Sprache  $L_{PicoC}$  hin zum einem Programm mit derselben **Semantik** in der **Konkreten Syntax** der Sprache  $L_{RETI}$  erklärt.

Für das Parsen<sup>1</sup> des Programmes in der **Konkreten 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 eines in einem **eigenen Dialekt** der **Erweiterten Back-Naur-Form** (Definition 2.3 bzw. für den Dialekt von Lark Definition 2.4) spezifizierten Grammatik der **Konkreten Syntax** ein Programm in ebendieser **Konkreten Syntax** zu parsen und daraus einen **Ableitungsbaum** für die kompilierintere Weiterverarbeitung zu generieren.

## Definition 2.1: Metasyntax

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

## Definition 2.2: Metasprache

*Eine Metasprache ist 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 **Kontextfreier Grammatiken** darzustellen.<sup>b c</sup>*

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

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

<sup>b</sup>Nebel, „Theoretische Informatik“.

<sup>c</sup>Grammar Reference — Lark documentation.

<sup>d</sup><https://standards.iso.org/ittf/PubliclyAvailableStandards/>.

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

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

<sup>3</sup>Shinan, lark.



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

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 für die Grammatiken des PicoC-Compilers 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: `()*`.<sup>b</sup>

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

<sup>b</sup>Bzw. kann der \*-Quantor auch **keinmal** wiederholen bedeuten.

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 **Konkreten Syntax** einer Sprache einen **Parser** in einer vorher bestimmten Programmiersprache generieren, anstatt wie das **Lark Parsing Toolkit** bei Angabe einer **Konkreten Syntax** direkt ein Programm in dieser Konkreten Syntax **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 Laufzeit 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 **Konkreten Syntaxen<sup>5</sup>**, die aufgrund der Verwendung des **Lark Parsing Toolkit** in einem **eigenen Dialekt** der **Erweiterter Back-Naur-Form** spezifiziert sind, werden in den folgenden Unterkapiteln die **Abstrakte Syntaxen**, welche spezifizieren, welche Kompositionen für die **Abstrakter Syntaxbaums** der verschiedenen **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 Grammatik nur **unnötig verkomplizieren** würde, wenn man die **Erweiterte Backus-Naur-Form** verwenden würde. Es gibt leider **keine** Standardnotation für die **Abstrakte Syntax**, die sich deutlich durchgesetzt hat, daher wird für die Abstrakte Syntaxen eine eigene **Abstract Syntax Form Notation** (Definition 2.5) verwendet. Des Weiteren trägt das Verwenden einer **unterschiedlichen Notation** für **Konkrete** 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 die **Grammatiken von Abstrakten Syntaxen**, die für diese Bachelorarbeit definiert wurde und sich von dem **Dialekt der Backus-Naur-Form** des **Lark Parsing Toolkit** nur dadurch unterscheidet, dass **Terminalsymbole** nicht von `"` eingeschlossen sein müssen, da die **Knoten** in der **Abstrakten Syntax**, sowieso schon **klar auszumachen** sind und von anderen Symbolen der **Metasprache** leicht zu unterscheiden sind.

Letztendlich geht es allerdings nur darum, dass aufgrund der Verwendung des **Lark Parsing Toolkit** die **Konkrete Syntax** in einem eigenen Dialekt der **Erweiterter Backus-Naur-Form** angegeben sein muss

<sup>4</sup>**ANTLR**.

<sup>5</sup>Der **Plural** von Syntax ist **Syntaxen**, wie es in Quelle *Syntax* verifiziert werden kann.



und für das Implementieren der Passes die **Abstrakte Syntax** für den **Programmierer** möglichst **einfach verständlich** sein sollte, weshalb sich die **Abstrakte Syntax Form** gut dafür eignet.

## 2.1 Lexikalische Analyse

Für die **Lexikalische Analyse** ist es nur notwendig eine Grammatik zu definieren, die den Teil der **Konkreten 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 Grammatik wird dann vom **Lark Parsing Toolkit** dazu verwendet ein Programm in **Konkreter Syntax** zu lexen und daraus Tokens für die **Syntaktische Analyse** zu erstellen, wie es im Unterkapitel 1.3 erläutert ist.

### 2.1.1 Konkrete Syntax für die Lexikalische Analyse

In der 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 „**kann abgeleitet werden zu**“-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 Grammatik 2.1.1 definierten **Nicht-Terminalsymbole** können in der Grammatik 2.2.8 der **Konkreten Syntax** für die **Syntaktischen Analyse** verwendet werden, um z.B. zu beschreiben, in welchem Kontext z.B. eine Zahl NUM stehen darf.

Die in der Konkrete Syntax 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 „**kann abgeleitet werden zu**“-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>6</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.

Der **Basic Lexer** des **Lark Parsing Toolkit** funktioniert grundlegend so wie es im Unterkapitel 1.3 erklärt wurde, allerdings berücksichtigt der **Basic Lexer** ebenfalls **Prioritäten**, 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

<sup>6</sup>Es wird immer die **höchste** Priorität **zuerst** genommen.

**Tokenwert** dieser Produktion zugewiesen, weitere Produktionen werden **nicht** mehr geprüft. Ansonsten werden alle Produktionen der **nächstniedrigeren** Priorität geprüft usw.

<i>COMMENT</i>	::=	"//"/[ $\backslash$ n]*/   "/*"/( $\cdot$   $\backslash$ n)*?/"*/"	<i>L_Comment</i>
<i>RETI_COMMENT.2</i>	::=	"//"" $\backslash$ "?"#[ $\backslash$ n]*/	
<i>DIG_NO_0</i>	::=	"1"   "2"   "3"   "4"   "5"   "6"   "7"   "8"   "9"	<i>L_Arith</i>
<i>DIG_WITH_0</i>	::=	"0"   <i>DIG_NO_0</i>	
<i>NUM</i>	::=	"0"   <i>DIG_NO_0 DIG_WITH_0*</i>	
<i>ASCII_CHAR</i>	::=	" $\backslash$ "." $\cdot$ "." $\sim$ "	
<i>CHAR</i>	::=	"'" <i>ASCII_CHAR</i> "'"	
<i>FILENAME</i>	::=	<i>ASCII_CHAR</i> + ". <i>picoc</i> "	
<i>LETTER</i>	::=	"a" $\cdot$ " $\cdot$ "z"   "A" $\cdot$ " $\cdot$ "Z"	
<i>NAME</i>	::=	( <i>LETTER</i>   " $\cdot$ ") ( <i>LETTER</i>   <i>DIG_WITH_0</i>   " $\cdot$ ")*	
<i>name</i>	::=	<i>NAME</i>   <i>INT_NAME</i>   <i>CHAR_NAME</i>   <i>VOID_NAME</i>	
<i>LOGIC_NOT</i>	::=	"!"	
<i>NOT</i>	::=	" $\sim$ "	
<i>REF_AND</i>	::=	"&"	
<i>un_op</i>	::=	<i>SUB_MINUS</i>   <i>LOGIC_NOT</i>   <i>NOT</i>   <i>MUL_DEREF_PNTR</i>   <i>REF_AND</i>	
<i>MUL_DEREF_PNTR</i>	::=	"*"	
<i>DIV</i>	::=	"/"	
<i>MOD</i>	::=	"%"	
<i>prec1_op</i>	::=	<i>MUL_DEREF_PNTR</i>   <i>DIV</i>   <i>MOD</i>	
<i>ADD</i>	::=	"+"	
<i>SUB_MINUS</i>	::=	"-"	
<i>prec2_op</i>	::=	<i>ADD</i>   <i>SUB_MINUS</i>	
<i>LT</i>	::=	"<"	<i>L_Logic</i>
<i>LTE</i>	::=	"<="	
<i>GT</i>	::=	">"	
<i>GTE</i>	::=	">="	
<i>rel_op</i>	::=	<i>LT</i>   <i>LTE</i>   <i>GT</i>   <i>GTE</i>	
<i>EQ</i>	::=	"=="	
<i>NEQ</i>	::=	"!="	
<i>eq_op</i>	::=	<i>EQ</i>   <i>NEQ</i>	
<i>INT_DT.2</i>	::=	"int"	<i>L_Assign_Alloc</i>
<i>INT_NAME.3</i>	::=	"int" ( <i>LETTER</i>   <i>DIG_WITH_0</i>   " $\cdot$ ")+	
<i>CHAR_DT.2</i>	::=	"char"	
<i>CHAR_NAME.3</i>	::=	"char" ( <i>LETTER</i>   <i>DIG_WITH_0</i>   " $\cdot$ ")+	
<i>VOID_DT.2</i>	::=	"void"	
<i>VOID_NAME.3</i>	::=	"void" ( <i>LETTER</i>   <i>DIG_WITH_0</i>   " $\cdot$ ")+	
<i>prim_dt</i>	::=	<i>INT_DT</i>   <i>CHAR_DT</i>   <i>VOID_DT</i>	

**Grammatik 2.1.1:** Grammatik der Konkreten Syntax 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 **Abstrakter Syntaxbaums** 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.

```
1 [Token('FILENAME', './verbose_dt_simple_ast_gen_array_decl_and_alloc.picoc'), Token('STRUCT',
  ↳ 'struct'), Token('NAME', 'st'), Token('LBRACE', '{'), Token('INT_DT', 'int'),
  ↳ Token('MUL_DEREF_PNTR', '*'), Token('LPAR', '('), Token('MUL_DEREF_PNTR', '*'),
  ↳ Token('NAME', 'attr'), Token('RPAR', ')'), Token('LSQB', '['), Token('NUM', '4'),
  ↳ Token('RSQB', ']'), Token('LSQB', '['), Token('NUM', '5'), Token('RSQB', ']'),
  ↳ Token('SEMICOLON', ';'), Token('RBRACE', '}'), Token('SEMICOLON', ';'), Token('VOID_DT',
  ↳ 'void'), Token('NAME', 'main'), Token('LPAR', '('), Token('RPAR', ')'), Token('LBRACE',
  ↳ '{'), Token('STRUCT', 'struct'), Token('NAME', 'st'), Token('MUL_DEREF_PNTR', '*'),
  ↳ Token('LPAR', '('), Token('MUL_DEREF_PNTR', '*'), Token('NAME', 'var'), Token('LSQB',
  ↳ '['), Token('NUM', '3'), Token('RSQB', ']'), Token('LSQB', '['), Token('NUM', '2'),
  ↳ Token('RSQB', ']'), Token('RPAR', ')'), Token('SEMICOLON', ';'), Token('RBRACE', '}')]
```

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 **Konkreter 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 **Abstrakter Syntaxbaum** in **Abstrakter Syntax** zu generieren.

### 2.2.1 Umsetzung von Präzedenz und Assoziativität

Die Programmiersprache  $L_{PicoC}$  hat dieselben **Präzidenzregeln** implementiert, wie die Programmiersprache  $L_C$ <sup>7</sup>. Die **Präzidenzregeln** der Programmiersprache  $L_{PicoC}$  sind in Tabelle 2.1 aufgelistet.

<sup>7</sup>C Operator Precedence - [cppreference.com](http://cppreference.com).

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

Tabelle 2.1: Präzidenzregeln von *PicoC*

Würde man diese **Operatoren** ohne Beachtung von **Präzidenzregeln** (Definiton 1.28) und **Assoziativität** (Definition 1.27) in eine Grammatik verarbeiten wollen, so könnte eine Grammatik, wie Grammatik 2.2.1 dabei rauskommen.

<i>prim_exp</i>	::=	<i>name</i>   <i>NUM</i>   <i>CHAR</i>   "(" <i>exp</i> ")"	<i>L_Arith</i> +
<i>un_op</i>	::=	"-"   "~"   "!"   "*"   "&"	
<i>un_exp</i>	::=	<i>un_op exp</i>	
<i>bin_op</i>	::=	"*"   "/"   "%"   "+"   "-"   "&"   "^"   " "   "<"   "<="   ">"   ">="   "!="   "=="   "&&"   "  "	
<i>bin_exp</i>	::=	<i>exp bin_op exp</i>	
<i>exp</i>	::=	<i>prim_exp</i>   <i>un_exp</i>   <i>bin_exp</i>	

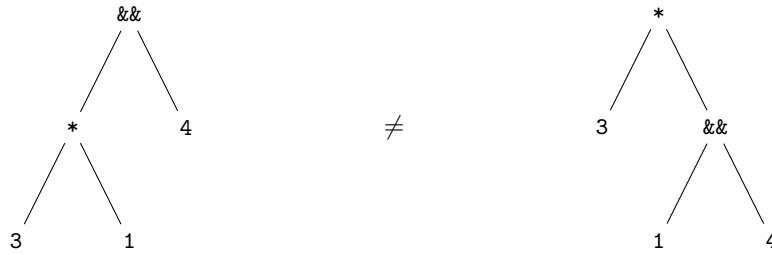
**Grammatik 2.2.1:** Undurchdachte Konkrete Syntax der Sprache  $L_{PicoC}$  für die Syntaktische Analyse in EBNF

Die Grammatik 2.2.1 ist allerdings **mehrdeutig**, d.h. verschiedene **Linksableitungen** in der Grammatik können zum selben **Wort** abgeleitet werden. Z.B. kann das Wort `3 * 1 && 4` sowohl über die **Linksableitung** 2.5.1 als auch über die **Linksableitung** 2.5.2 abgeleitet werden.

$$\begin{aligned}
 \text{exp} &\Rightarrow \text{bin\_exp} \Rightarrow \text{exp bin\_op exp} \Rightarrow \text{bin\_exp bin\_op exp} \\
 &\Rightarrow \text{exp bin\_op exp bin\_op exp} \Rightarrow^* 3 * 1 \&\& 4
 \end{aligned}
 \tag{2.5.1}$$

$$\begin{aligned}
 \text{exp} &\Rightarrow \text{bin\_exp} \Rightarrow \text{exp bin\_op exp} \Rightarrow \text{prim\_exp bin\_op exp} \Rightarrow \text{NUM bin\_op exp} \\
 &\Rightarrow 3 \text{ bin\_op exp} \Rightarrow 3 * \text{exp} \Rightarrow 3 * \text{bin\_exp} \Rightarrow 3 * \text{exp bin\_exp} \Rightarrow^* 3 * 1 \&\& 4
 \end{aligned}
 \tag{2.5.2}$$

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 2.5.1 und der **rechte Baum** entspricht Ableitung 2.5.2. Würde man in den Ausdrücken, die von diesen Bäumen dargestellt sind in **Klammern** setzen, um die **Präzedenz** sichtbar zu machen, so würde Ableitung 2.5.1 die Klammerung  $(3 * 1) \&\& 4$  haben und die Ableitung 2.5.2 die Klammerung  $3 * (1 \&\& 4)$  haben.

Aus diesem Grund ist es wichtig die **Präzidenzregeln** und die **Assoziativität** der Operatoren beim Erstellen der 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.

<code>prim_exp</code>	::=	...	<code>L_Arith + L_Array</code>
<code>post_exp</code>	::=	...	<code>+ L_Pntr + L_Struct</code>
<code>un_exp</code>	::=	...	<code>+ L_Fun</code>
<code>arith_prec1</code>	::=	...	
<code>arith_prec2</code>	::=	...	
<code>arith_and</code>	::=	...	
<code>arith_oplus</code>	::=	...	
<code>arith_or</code>	::=	...	
<code>rel_exp</code>	::=	...	<code>L_Logic</code>
<code>eq_exp</code>	::=	...	
<code>logic_and</code>	::=	...	
<code>logic_or</code>	::=	...	
<code>assign_stmt</code>	::=	...	<code>L_Assign</code>

**Grammatik 2.2.2:** Durchdachte Konkrete Syntax der Sprache  $L_{PicoC}$  für die Syntaktische Analyse in EBNF

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
arith_and	a<b a<=b a>b a>=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 **un\_exp** in 2.2.3 für die **rechtsassoziativen unären Operatoren**  $-a$ ,  $!a$ ,  $\sim a$ ,  $*a$  und  $\&a$  ist die **Alternative** **un\_op un\_exp** dafür zuständig, dass diese unären Operatoren **rechtsassoziativ** geschachtelt werden können (z.B.  $! \sim 42$ ). Die Alternative **post\_exp** ist dafür zuständig, dass die Produktion auch **terminieren** kann und es auch möglich ist ausschließlich 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 ausschließlich einen Ausdruck der **höchsten Präzidenz** (z.B. 42) zu haben.

$$post\_exp ::= post\_exp\ ["logic\_or"] \mid post\_exp\ ."name" \mid name\ ("fun\_args") \mid prim\_exp$$

**Grammatik 2.2.4:** Beispiel für eine unäre linksassoziative Produktion

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

<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 ausschließlich ein Ausdruck **höherer Präzidenz** erkannt wird (z.B.  $1 / 4$ ).

$$\text{arith\_prec2} ::= \text{arith\_prec2 } \text{prec2\_op } \text{arith\_prec1} \mid \text{arith\_prec1}$$

**Grammatik 2.2.5:** *Beispiel für eine linksassoziative Produktion*

Manche **Parser**<sup>a</sup> haben allerdings ein Problem mit **Linksrekursion** (Definition 1.24), 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.

$$\text{arith\_prec2} ::= \text{arith\_prec1 } (\text{prec2\_op } \text{arith\_prec1})^*$$

**Grammatik 2.2.6:** *Beispiel für eine linksassoziative Produktion*

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 ausschließlich 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 Grammatik 2.2.7.

<i>prec1_op</i>	::=	"*"   "/"   "%"	<i>L_Misc</i>
<i>prec2_op</i>	::=	"+"   "-"	
<i>rel_op</i>	::=	"<"   "<="   ">"   ">="	
<i>eq_op</i>	::=	"=="   "!="	
<i>fun_args</i>	::=	[ <i>logic_or</i> ("," <i>logic_or</i> )*]	
<i>prim_exp</i>	::=	<i>name</i>   <i>NUM</i>   <i>CHAR</i>   "(" <i>logic_or</i> ")"	<i>L_Arith</i>
<i>post_exp</i>	::=	<i>post_exp</i> [" <i>logic_or</i> "]   <i>post_exp</i> ." <i>name</i> "   <i>name</i> (" <i>fun_args</i> ")"	+ <i>L_Array</i>
		<i>prim_exp</i>	+ <i>L_Pntr</i>
<i>un_exp</i>	::=	<i>un_op</i> <i>un_exp</i>   <i>post_exp</i>	+ <i>L_Struct</i>
<i>arith_prec1</i>	::=	<i>arith_prec1</i> <i>prec1_op</i> <i>un_exp</i>   <i>un_exp</i>	+ <i>L_Fun</i>
<i>arith_prec2</i>	::=	<i>arith_prec2</i> <i>prec2_op</i> <i>arith_prec1</i>   <i>arith_prec1</i>	
<i>arith_and</i>	::=	<i>arith_and</i> "&" <i>arith_prec2</i>   <i>arith_prec2</i>	
<i>arith_oplus</i>	::=	<i>arith_oplus</i> "^" <i>arith_and</i>   <i>arith_and</i>	
<i>arith_or</i>	::=	<i>arith_or</i> " " <i>arith_oplus</i>   <i>arith_oplus</i>	
<i>rel_exp</i>	::=	<i>rel_exp</i> <i>rel_op</i> <i>arith_or</i>   <i>arith_or</i>	<i>L_Logic</i>
<i>eq_exp</i>	::=	<i>eq_exp</i> <i>eq_op</i> <i>rel_exp</i>   <i>rel_exp</i>	
<i>logic_and</i>	::=	<i>logic_and</i> "&&" <i>eq_exp</i>   <i>eq_exp</i>	
<i>logic_or</i>	::=	<i>logic_or</i> "  " <i>logic_and</i>   <i>logic_and</i>	
<i>assign_stmt</i>	::=	<i>un_exp</i> "=" <i>logic_or</i> ";"	<i>L_Assign</i>

**Grammatik 2.2.7:** Durchdachte Konkrete Syntax für Operatorpräzedenz in EBNF

## 2.2.2 Konkrete Syntax für die Syntaktische Analyse

Die gesamte Grammatik 2.2.8, welche die **Konkrete Syntax** der Sprache  $L_{PicoC}$  für die **Syntaktische Analyse** beschreibt ergibt sich wenn man die 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 **Konkrete Syntax** an die **aktuellste konsistentlos auffindbare Version** der echten 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 Grammatik 2.2.8, welche die **Konkrete Syntax** der Sprache  $L_{PicoC}$  für die **Syntaktische Analyse** beschreibt, werden einige der **Tokennamen** aus der Grammatik 2.1.1 der **Konkreten Syntax** 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 Grammatik leichter lesbar zu machen einige **Terminalsymbole** einfach direkt in die Grammatik 2.2.8 der **Konkreten Syntax** 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 **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 Grammatiken für die **Syntaktische Analyse** *ANSI C grammar (Lex)* und **Lexikalische Analyse** *noauthor'ansi'nodate-2* aus dem Jahre 1985 zeigen.

<sup>11</sup>*GCC, the GNU Compiler Collection - GNU Project.*

<sup>12</sup>*clang: C++ Compiler.*



<i>prim_exp</i>	::=	<i>name</i>   <i>NUM</i>   <i>CHAR</i>   "(" <i>logic_or</i> ")"	<i>L_Arith</i> + <i>L_Array</i>
<i>post_exp</i>	::=	<i>array_subscr</i>   <i>struct_attr</i>   <i>fun_call</i>	+ <i>L_Pntr</i> + <i>L_Struct</i>
		<i>input_exp</i>   <i>print_exp</i>   <i>prim_exp</i>	+ <i>L_Fun</i>
<i>un_exp</i>	::=	<i>un_op un_exp</i>   <i>post_exp</i>	
<i>input_exp</i>	::=	"input" "(" ")"	
<i>print_exp</i>	::=	"print" "(" <i>logic_or</i> ")"	
<i>arith_prec1</i>	::=	<i>arith_prec1 prec1_op un_exp</i>   <i>un_exp</i>	
<i>arith_prec2</i>	::=	<i>arith_prec2 prec2_op arith_prec1</i>   <i>arith_prec1</i>	
<i>arith_and</i>	::=	<i>arith_and</i> "&" <i>arith_prec2</i>   <i>arith_prec2</i>	
<i>arith_oplus</i>	::=	<i>arith_oplus</i> "^" <i>arith_and</i>   <i>arith_and</i>	
<i>arith_or</i>	::=	<i>arith_or</i> " " <i>arith_oplus</i>   <i>arith_oplus</i>	
<i>rel_exp</i>	::=	<i>rel_exp rel_op arith_or</i>   <i>arith_or</i>	<i>L_Logic</i>
<i>eq_exp</i>	::=	<i>eq_exp eq_op rel_exp</i>   <i>rel_exp</i>	
<i>logic_and</i>	::=	<i>logic_and</i> "&&" <i>eq_exp</i>   <i>eq_exp</i>	
<i>logic_or</i>	::=	<i>logic_or</i> "  " <i>logic_and</i>   <i>logic_and</i>	
<i>type_spec</i>	::=	<i>prim_dt</i>   <i>struct_spec</i>	<i>L_Assign_Alloc</i>
<i>alloc</i>	::=	<i>type_spec pntr_decl</i>	
<i>assign_stmt</i>	::=	<i>un_exp</i> "=" <i>logic_or</i> ";"	
<i>initializer</i>	::=	<i>logic_or</i>   <i>array_init</i>   <i>struct_init</i>	
<i>init_stmt</i>	::=	<i>alloc</i> "=" <i>initializer</i> ";"	
<i>const_init_stmt</i>	::=	"const" <i>type_spec name</i> "=" <i>NUM</i> ";"	
<i>pntr_deg</i>	::=	"*"	<i>L_Pntr</i>
<i>pntr_decl</i>	::=	<i>pntr_deg array_decl</i>   <i>array_decl</i>	
<i>array_dims</i>	::=	("[" <i>NUM</i> "]" ) *	<i>L_Array</i>
<i>array_decl</i>	::=	<i>name array_dims</i>   "(" <i>pntr_decl</i> ")" <i>array_dims</i>	
<i>array_init</i>	::=	"{" <i>initializer</i> ("," <i>initializer</i> ) * "}"	
<i>array_subscr</i>	::=	<i>post_exp</i> "[" <i>logic_or</i> "]"	
<i>struct_spec</i>	::=	"struct" <i>name</i>	<i>L_Struct</i>
<i>struct_params</i>	::=	( <i>alloc</i> ";" ) +	
<i>struct_decl</i>	::=	"struct" <i>name</i> "{" <i>struct_params</i> "}"	
<i>struct_init</i>	::=	"{" " ." <i>name</i> "=" <i>initializer</i> ("," " ." <i>name</i> "=" <i>initializer</i> ) * "}"	
<i>struct_attr</i>	::=	<i>post_exp</i> "." <i>name</i>	
<i>if_stmt</i>	::=	"if" "(" <i>logic_or</i> ")" <i>exec_part</i>	<i>L_If_Else</i>
<i>if_else_stmt</i>	::=	"if" "(" <i>logic_or</i> ")" <i>exec_part</i> "else" <i>exec_part</i>	
<i>while_stmt</i>	::=	"while" "(" <i>logic_or</i> ")" <i>exec_part</i>	<i>L_Loop</i>
<i>do_while_stmt</i>	::=	"do" <i>exec_part</i> "while" "(" <i>logic_or</i> ")" ";"	

**Grammatik 2.2.8:** Grammatik der Konkreten Syntax der Sprache  $L_{\text{PicoC}}$  für die Syntaktische Analyse in EBNF, Teil 1

<code>decl_exp_stmt</code>	<code>::= alloc";"</code>	<i>L_Stmt</i>
<code>decl_direct_stmt</code>	<code>::= assign_stmt   init_stmt   const_init_stmt</code>	
<code>decl_part</code>	<code>::= decl_exp_stmt   decl_direct_stmt   RETI_COMMENT</code>	
<code>compound_stmt</code>	<code>::= "{" exec_part * "}"</code>	
<code>exec_exp_stmt</code>	<code>::= logic_or";"</code>	
<code>exec_direct_stmt</code>	<code>::= if_stmt   if_else_stmt   while_stmt   do_while_stmt</code>	
	<code>  assign_stmt   fun_return_stmt</code>	
<code>exec_part</code>	<code>::= compound_stmt   exec_exp_stmt   exec_direct_stmt</code>	
	<code>  RETI_COMMENT</code>	
<code>decl_exec_stmts</code>	<code>::= decl_part * exec_part*</code>	
<code>fun_args</code>	<code>::= [logic_or(" ", logic_or)*]</code>	<i>L_Fun</i>
<code>fun_call</code>	<code>::= name("fun_args")</code>	
<code>fun_return_stmt</code>	<code>::= "return" [logic_or];"</code>	
<code>fun_params</code>	<code>::= [alloc(" ", alloc)*]</code>	
<code>fun_decl</code>	<code>::= type_spec ptr_deg name("fun_params")</code>	
<code>fun_def</code>	<code>::= type_spec ptr_deg name("fun_params") "{" decl_exec_stmts "}"</code>	
<code>decl_def</code>	<code>::= (struct_decl   fun_decl);"   fun_def</code>	<i>L_File</i>
<code>decls_defs</code>	<code>::= decl_def*</code>	
<code>file</code>	<code>::= FILENAME decls_defs</code>	

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

In der 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 **Konkrete Syntax**, die von der Grammatik 2.2.8 beschrieben wird 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 **top-down** arbeitenden **Earley Recognizer** (Defintion 2.7) nutzt, um einen **Ableitungsbaum** zu konstruieren.

#### Definition 2.7: Earley Recognizer

Ist ein **Recognizer** der für **alle Kontextfreien Sprachen** das **Wortproblem** entscheiden kann.

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

- $\alpha, \beta, \gamma$  stellen eine **beliebige Folge** von **Grammatiksymbolen**<sup>a</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 geparkt** wurde und  $\beta$  **noch geparkt** werden muss

und davor müssen noch einige **Begriffe definiert** werden:

- **Zustandsmenge:** Für jeden **Index** des **Eingabeworts** wird eine **Zustandsmenge** generiert. Die **Zustandsmenge** für die **Index**  $j$  im **Eingabewort** wird mit  $Z(j)$  bezeichnet.
- **Zustand der Zustandsmenge:** Ist ein **Tupel**  $(A ::= \alpha \bullet \beta, i)$ , wobei  $A ::= \alpha \bullet \beta$  die **momentane Produktion** ist, die erkannt werden soll, die bis **Punkt**  $\bullet$  geparkt wurde und  $i$  der **Index** im **Eingabewort** ist, ab welchem der Versuch der Erkennung dieser **Produktion** begann.

Der **Ablauf** des Algorithmus ist wie folgt:

- **Input:** Wort  $w$  und **Grammatik**  $G = \langle N, \Sigma, P, A \rangle$
  - **Output:** 0 wenn  $w \notin L(G)$ <sup>b</sup> und 1 wenn  $w \in L(G)$
1. **initialisiere**  $Z(0)$  mit der **Produktion**, welches das **Startsymbol**  $A$  auf der **linken Seite** des „kann abgeleitet werden zu“-Symbols  $::=$  hat
  2. es werden **wiederholt** 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 **Grammatik**, die ein  $B$  auf der **linken Seite** des „kann abgeleitet werden zu“-Symbols  $::=$  hat ein **Zustand**  $(B ::= \bullet\beta, j)$  zur **Zustandsmenge**  $S(j)$  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  $A ::= \beta$  die **Produktion** ist, welche das **Startsymbol**  $A$  auf der **linken Seite** des „kann abgeleitet werden zu“-Symbols  $::=$  hat  $\Rightarrow w \in L(G)$
- **keine Zustände** mehr einer **Zustandsmenge** hinzugefügt werden können  $\Rightarrow w \notin L(G)$

<sup>a</sup>Also eine Folge von **Terminalsymbolen** und **Nicht-Terminalsymbolen**.

<sup>b</sup> $L(G)$  ist die **Sprache**, welche durch die **Grammatik**  $G$  beschrieben wird.

### 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.

```

1 file
2 ./verbose_dt_simple_ast_gen_array_decl_and_alloc.dt
3 decls_defs
4   decl_def
5     struct_decl
6       name      st
7       struct_params
8       alloc
9       type_spec
10      prim_dt      int
11      pntr_decl
12      pntr_deg      *
13      array_decl
14      pntr_decl
15      pntr_deg      *
16      array_decl
17      name      attr
18      array_dims
19      array_dims
20      4
21      5
22   decl_def
23   fun_def
24     type_spec
25     prim_dt      void
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      *
41     array_decl
42     name      var
43     array_dims
44     3
45     2
46     array_dims

```

Code 2.3: Ableitungsbaum nach Ableitungsbaum Generierung

### 2.2.3.2 Ausgabe des Ableitungsbaums

Die Ausgabe des **Ableitungsbaums** wird komplett vom **Lark Parsing Toolkit** übernommen. Für die **Inneren Knoten** werden die **Nicht-Terminalsymbole**, welche in der Grammatik den **linken Seiten** des „kann abgeleitet werden zu“-**Symbols**  $::=$ <sup>13</sup> entsprechen hergenommen und die **Blätter** sind **Terminalsymbole**, genauso, wie es in der Definition 1.36 eines **Ableitungsbaums** auch schon definiert

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

ist. Die **EBNF-Grammatik 2.2.8** des **PicoC-Compilers** erlaubt es allerdings auch, dass in einem **Blatt** gar nichts  $\varepsilon$  steht, weil es z.B. **Produktionen**, wie `array_dims ::= ("NUM")*` gibt, in denen auch das **leere Wort**  $\varepsilon$  abgeleitet werden kann.

Die Ausgabe des **Abstrakter Syntaxbaum** ist bewusst so gewählt, dass sie sich optisch vom **Ableitungsbaum** unterscheidet, indem die Bezeichner der **Knoten** in **UpperCamelCase** geschrieben sind, im Gegensatz zum **Ableitungsbaum**, dessen **Innere Knoten** im **snake\_case** geschrieben sind, wie auch die **Nicht-Terminalsymbole** auf den **linken Seiten** des „kann abgeleitet werden zu“-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 **Transformers** 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 „**Array der Mächtigkeit 3 von Pointern auf Arrays der Mächtigkeit 2 von Integeren**“ `int (*ar[3])[2]` die **Spiralregel**<sup>14</sup> in der Implementierung des **PicoC-Compilers** umgesetzt werden muss und das ist nicht alleinig möglich, indem man die entsprechenden **Produktionen** in der Grammatik 2.2.8 der **Konkreten 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'), ArrayDecl([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.

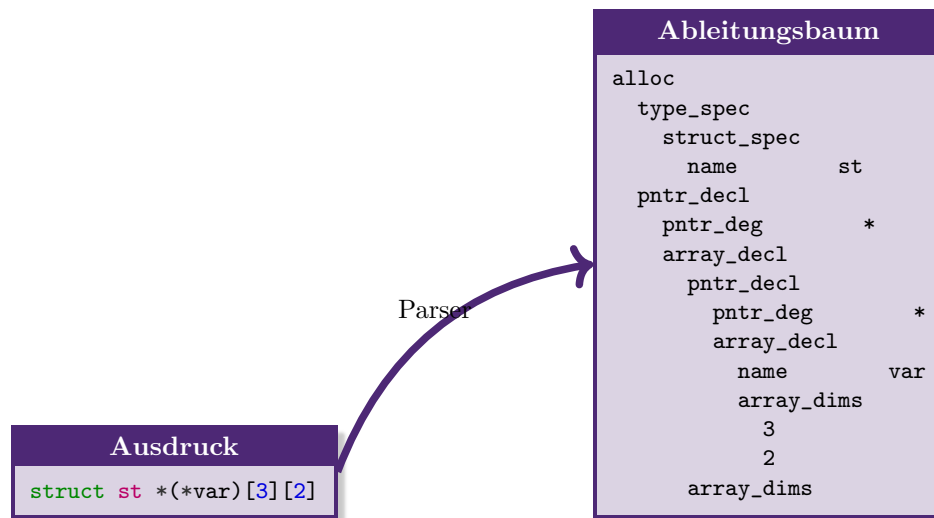


Abbildung 2.2: Ableitungsbaum nach Parsen eines Ausdrucks

Dieser **Ableitungsbaum** für den Ausdruck `struct st *(*var[3][2])` hat allerdings einen Aufbau welcher durch die **Syntax** der **Pointerdeklaratoren** `pntr_decl(num, datatype)` und **Arraydeklaratoren** `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

<sup>14</sup> Clockwise/Spiral Rule.

im **zweiten** und **ersten**, wie beim **Pointerdeklarator** `pntr_decl(num, datatype)` und **Arraydeklarator** `array_decl(datatype, nums)`. Daher muss beim **ArrayDeclarator** `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 **Konkreter 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 1.40) **vereinfacht** wird, sodass er die gerade erläuterten Ansprüche erfüllt.

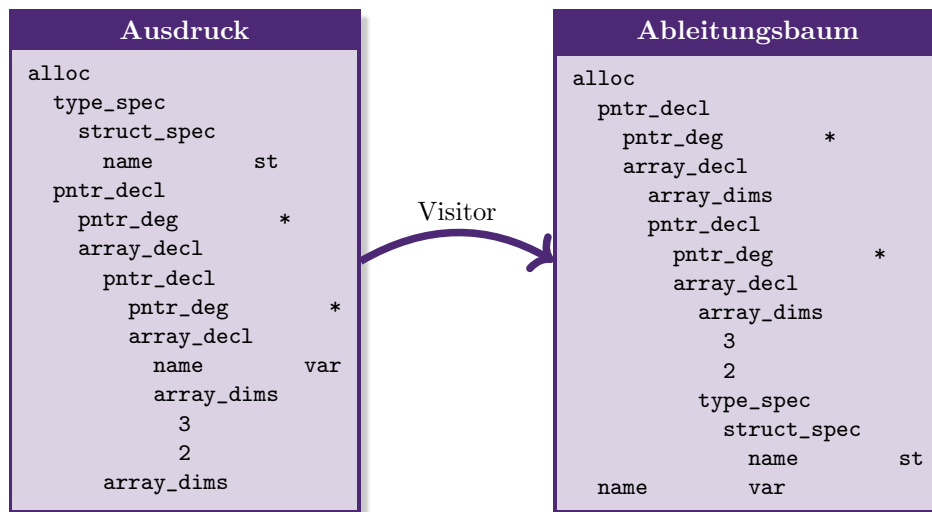


Abbildung 2.3: Ableitungsbaum nach Vereinfachung

#### 2.2.4.1 Codebeispiel

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

```

1 file
2 ./verbose_dt_simple_ast_gen_array_decl_and_alloc.dt_simple
3 decls_defs
4 decl_def
5 struct_decl
6 name      st
7 struct_params
8 alloc
9 pntr_decl
10 pntr_deg      *
11 array_decl
12 array_dims
13 4
14 5
15 pntr_decl

```

```

16      ptr_deg      *
17      array_decl
18      array_dims
19      type_spec
20      prim_dt      int
21      name      attr
22  decl_def
23  fun_def
24      type_spec
25      prim_dt      void
26      ptr_deg
27      name      main
28      fun_params
29      decl_exec_stmts
30      decl_part
31      decl_exp_stmt
32  alloc
33      ptr_decl
34      ptr_deg      *
35      array_decl
36      array_dims
37      ptr_decl
38      ptr_deg      *
39      array_decl
40      array_dims
41      3
42      2
43      type_spec
44      struct_spec
45      name      st
46      name      var

```

Code 2.4: Ableitungsbaum nach Ableitungsbaum Vereinfachung

### 2.2.5 Abstrakt Syntax Tree Generierung

Nachdem der **Derivation Tree** in Unterkapitel 2.2.4 vereinfacht wurde, ist der **vereinfachte Ableitungsbaum** in Code 2.4 nun dazu geeignet, um mit einem **Transformer** (Definition 1.39) einen **Abstrakter Syntaxbaum** aus ihm zu generieren. Würde man den **vereinfachten Ableitungsbaum** des Ausdrucks `struct st *(*var[3][2])` auf passende Weise in einen **Abstrakter Syntaxbaum** umwandeln, so würde dabei ein **Abstrakter Syntaxbaum** wie in Abbildung 2.4 rauskommen.

Den Teilbaum, der den Datentyp darstellt würde man von **oben-nach-unten**<sup>15</sup> als „**Pointer auf einen Pointer auf ein Array der Mächtigkeit 2, 3 von Structs 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 „**Array der Mächtigkeit 3, 2 von Pointern auf einen Pointer auf einen Struct vom Typ st**“ ist.

<sup>15</sup>In der Informatik wachsen Bäume von **oben-nach-unten**, von der **Wurzel** zur den **Blättern**, bzw. in diesem Beispiel von **links-nach-rechts**.

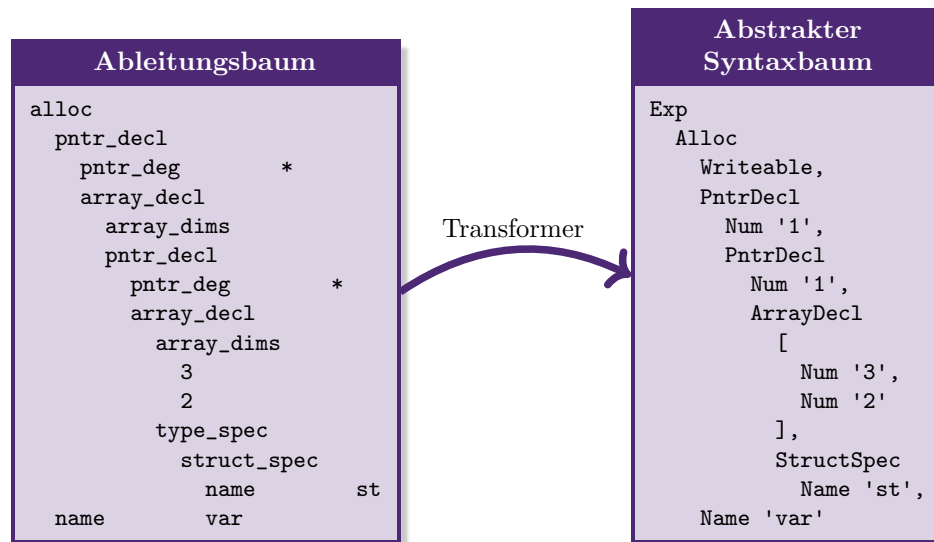


Abbildung 2.4: Abstrakter Syntaxbaum Generierung ohne Umdrehen

Dieser **Abstrakter Syntaxbaum** ist für die Weiterverarbeitung ungeeignet, denn für die **Adressberechnung** für eine Aneinanderreihung von Zugriffen auf **Pointerelemente**, **Arrayelemente** oder **Structattribute**, welche in Unterkapitel ?? genauer erläutert wird, will man den Datentyp in **umgekehrter Reihenfolge**. Aus diesem Grund muss der **Transformer** bei der Konstruktion des **Abstrakter Syntaxbaum** 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.

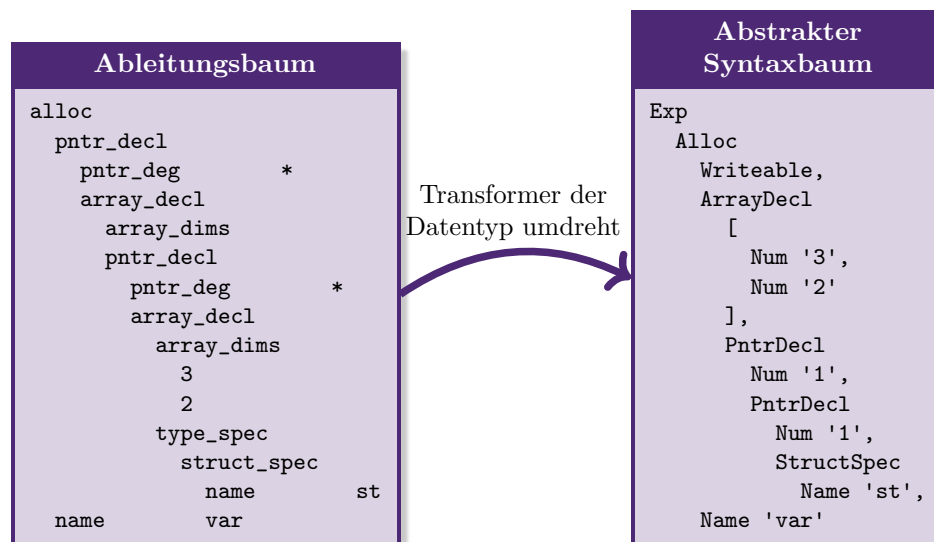


Abbildung 2.5: Abstrakter Syntaxbaum Generierung mit Umdrehen

Die Weiterverarbeitung des **Abstrakter Syntaxbaums** geschieht mithilfe von **Passes**, welche im Unterkapitel 1.5 genauer beschrieben werden. Da die Knoten des **Abstrakter Syntaxbaum** anders als beim **Ableitungsbaum** nicht die gleichen Bezeichnungen haben wie **Produktionen** der Grammatik der **Kon-**



**kretten Syntax** 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 **Abstrakter Syntaxbaums** vorkommen.

Des Weiteren gibt die **Abstrakte Syntax** die durch die Grammatik 2.2.1 in Unterkapitel 2.2.5.4 beschrieben wird 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>16</sup>. Alle **PicoC-Knoten**, die in den von den verschiedenen Passes generierten **Abstrakter Syntaxbaums** vorkommen sind in Tabelle 2.3 mit einem **Bschreibungstext** dokumentiert.

<sup>16</sup>Z.B. steht der **PicoC-Knoten** `Name(str)` für einen **Bezeichner**. Anstatt diesen Knoten in englisch `Identifizier(str)` zu nennen, wurde dieser als `Name(str)` gewählt, da `Name(str)` **kürzer** ist und **intuitiver verständlich**.

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

Tabelle 2.3: PicoC-Knoten Teil 1

Piocc-Knoten	Beschreibung
Exp(exp, datatype, error_data)	Container für einen <b>beliebigen Ausdruck</b> , dessen Ergebnis auf den <b>Stack</b> soll. Zudem besitzt er 2 versteckte Attribute, wobei <b>datatype</b> im <b>RETI Blocks Pass</b> wichtig ist und <b>error_data</b> für <b>Fehlermeldungen</b> wichtig ist.
Stack(num)	Container, der für das <b>temporäre</b> Ergebnis einer Berechnung, das <b>num</b> Speicherzellen relativ zum <b>Stackpointer Register SP</b> steht.
Stackframe(num)	Container, der für eine Variable steht, die <b>num</b> Speicherzellen relativ zum <b>Begin-Aktive-Funktion Register BAF</b> steht.
Global(num)	Container, der für eine Variable steht, die <b>num</b> Speicherzellen relativ zum <b>Datensegment Register DS</b> steht.
StackMalloc(num)	Container, der für das <b>Allokieren</b> von <b>num</b> Speicherzellen auf dem <b>Stack</b> steht.
PntrDecl(num, datatype)	Container, der für den <b>Pointerdatatype</b> steht: <code>&lt;prim_dt&gt; *&lt;var&gt;</code> , wobei das <b>Attribut</b> <b>num</b> die <b>Anzahl zusammengefasster Pointer</b> angibt und <b>datatype</b> der Datentyp ist, auf den der oder die <b>Pointer</b> zeigen.
Ref(exp, datatype, error_data)	Container, der für die Anwendung des <b>Referenz-Operators</b> <code>&amp;&lt;var&gt;</code> steht und die <b>Adresse</b> einer <b>Location</b> (Definition 1.47) auf den Stack schreiben soll, die über <b>exp</b> eingegrenzt wird. Zudem besitzt er 2 versteckte Attribute, wobei <b>datatype</b> im <b>RETI Blocks Pass</b> wichtig ist und <b>error_data</b> für <b>Fehlermeldungen</b> wichtig ist.
Deref(lhs, exp)	Container für den <b>Indezzugriff</b> auf einen <b>Array- oder Pointerdatatype</b> : <code>&lt;var&gt;[&lt;i&gt;]</code> , wobei <b>exp1</b> eine angehängte weitere <code>Subscr(exp1, exp2)</code> , <code>Deref(exp1, exp2)</code> , <code>Attr(exp, name)</code> oder ein <code>Name('var')</code> sein kann und <b>exp2</b> der Index ist auf den zugegriffen werden soll.
ArrayDecl(nums, datatype)	Container, der für den <b>Arraydatatype</b> steht: <code>&lt;prim_dt&gt; &lt;var&gt;[&lt;i&gt;]</code> , wobei das <b>Attribut</b> <b>nums</b> eine Liste von <code>Num('x')</code> ist, die die <b>Dimensionen</b> des Arrays angibt und <b>datatype</b> der Datentyp ist, der über das Anwenden von <code>Subscript()</code> auf das Array zugreifbar ist.
Array(exps, datatype)	Container für den <b>Initializer</b> eines <b>Arrays</b> , dessen Einträge <b>exps</b> weitere Initializer für eine <b>Array-Dimension</b> oder ein Initializer für ein <b>Struct</b> oder ein <b>Logischer Ausdruck</b> sein können, z.B. <code>{1, 2}</code> , <code>{3, 4}</code> . Des Weiteren besitzt er ein verstecktes Attribut <b>datatype</b> , welches für den <b>PicoC-ANF Pass</b> Informationen transportiert, die für <b>Fehlermeldungen</b> wichtig sind.
Subscr(exp1, exp2)	Container für den <b>Indezzugriff</b> auf einen <b>Array- oder Pointerdatatype</b> : <code>&lt;var&gt;[&lt;i&gt;]</code> , wobei <b>exp1</b> eine angehängte weitere <code>Subscr(exp1, exp2)</code> , <code>Deref(exp1, exp2)</code> oder <code>Attr(exp, name)</code> Operation sein kann oder ein <code>Name('var')</code> sein kann und <b>exp2</b> der Index ist auf den zugegriffen werden soll.
StructSpec(name)	Container für einen selbst definierten <b>Structdatatype</b> : <code>struct &lt;name&gt;</code> , wobei das <b>Attribut</b> <b>name</b> festlegt, welchen <b>selbst definierte</b> Structdatatype dieser Container-Knoten repräsentiert.
Attr(exp, name)	Container für den <b>Attributzugriff</b> auf einen <b>Structdatatype</b> : <code>&lt;var&gt;.&lt;attr&gt;</code> , wobei <b>exp1</b> eine angehängte weitere <code>Subscr(exp1, exp2)</code> , <code>Deref(exp1, exp2)</code> oder <code>Attr(exp, name)</code> Operation sein kann oder ein <code>Name('var')</code> sein kann und <b>name</b> das Attribut ist, auf das zugegriffen werden soll.

Tabelle 2.4: PicoC-Knoten Teil 2

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

Tabelle 2.5: PicoC-Knoten Teil 3

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

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

Jeder **Knoten** hat darüberhinaus auch noch 2 **Attribute** **value** und **position**, wobei **value** bei einem **Token-Knoten** (Definition 2.9) dem **Tokenwert** des Tokens, welches es ersetzt entspricht und bei **Container-Knoten** (Definition 2.10) unbesetzt ist. Das **Attribut position** wird später für Fehlermeldungen gebraucht.

**Definition 2.9: Token-Knoten**

Ersetzt ein *Token* bei der Generierung des *Abstrakter Syntaxbaum*, damit der Zugriff auf Knoten des Abstrakter Syntaxbaum möglichst *simpel* ist und keine vermeidbaren Fallunterscheidungen gemacht werden müssen.

*Token-Knoten* entsprechen im Abstrakter Syntaxbaum *Blättern*.<sup>a</sup>

<sup>a</sup>Thiemann, „Compilerbau“.

**Definition 2.10: Container-Knoten**

Dient als *Container* für andere *Container-Knoten* und *Token-Knoten*. Die *Container-Knoten* werden optimalerweise immer so gewählt, dass sie *mehrere Produktionen der Konkreten Syntax* abdecken, die einen *gleichen Aufbau* haben und sich auch unter einem *Überbegriff* zusammenfassen lassen.<sup>a</sup>

*Container-Knoten* entsprechen im Abstrakter Syntaxbaum *Inneren Knoten*.<sup>b</sup>

<sup>a</sup>Wie z.B. die verschiedenen **Arithmetischen Ausdrücke**, wie z.B. **1 % 3** und **Logischen Ausdrücke**, wie z.B. **1 && 2 < 3**, die einen gleichen Aufbau haben mit immer einer **Operation** in der Mitte haben und 2 **Operanden** auf beiden Seiten und sich unter dem Überbegriff **Binäre Operationen** zusammenfassen lassen.

<sup>b</sup>Thiemann, „Compilerbau“.

**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 erläuterten 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 **Abstrakter Syntaxbaums** vorkommen sind in Tabelle 2.2.5.1 mit einem **Beschreibungstext** dokumentiert.

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

<sup>a</sup> 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 und nicht bereits in der **Abstrakten Syntax 2.2.8** enthalten sind.



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



Um die obige Tabelle 2.8 nicht mit unnötig viel repetitiven Inhalt zu füllen, wurden die zahlreichen Kompositionen **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 Syntax

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

<i>stmt</i>	::=	<i>SingleLineComment(str, str)</i>   <i>RETIComment()</i>	<i>L_Comment</i>
<i>un_op</i>	::=	<i>Minus()</i>   <i>Not()</i>	<i>L_Arith</i>
<i>bin_op</i>	::=	<i>Add()</i>   <i>Sub()</i>   <i>Mul()</i>   <i>Div()</i>   <i>Mod()</i>   <i>Oplus()</i>   <i>And()</i>   <i>Or()</i>	
<i>exp</i>	::=	<i>Name(str)</i>   <i>Num(str)</i>   <i>Char(str)</i>   <i>BinOp(&lt;exp&gt;, &lt;bin_op&gt;, &lt;exp&gt;)</i>   <i>UnOp(&lt;un_op&gt;, &lt;exp&gt;)</i>   <i>Call(Name('input'), Empty())</i>   <i>Call(Name('print'), &lt;exp&gt;)</i>	
<i>stmt</i>	::=	<i>Exp(&lt;exp&gt;)</i>	
<i>un_op</i>	::=	<i>LogicNot()</i>	<i>L_Logic</i>
<i>rel</i>	::=	<i>Eq()</i>   <i>NEq()</i>   <i>Lt()</i>   <i>LtE()</i>   <i>Gt()</i>   <i>GtE()</i>	
<i>bin_op</i>	::=	<i>LogicAnd()</i>   <i>LogicOr()</i>	
<i>exp</i>	::=	<i>Atom(&lt;exp&gt;, &lt;rel&gt;, &lt;exp&gt;)</i>   <i>ToBool(&lt;exp&gt;)</i>	
<i>type_qual</i>	::=	<i>Const()</i>   <i>Writeable()</i>	<i>L_Assign_Alloc</i>
<i>datatype</i>	::=	<i>IntType()</i>   <i>CharType()</i>   <i>VoidType()</i>	
<i>exp</i>	::=	<i>Alloc(&lt;type_qual&gt;, &lt;datatype&gt;, Name(str))</i>	
<i>stmt</i>	::=	<i>Assign(&lt;exp&gt;, &lt;exp&gt;)</i>	
<i>datatype</i>	::=	<i>PntrDecl(Num(str), &lt;datatype&gt;)</i>	<i>L_Pntr</i>
<i>exp</i>	::=	<i>Deref(&lt;exp&gt;, &lt;exp&gt;)</i>   <i>Ref(&lt;exp&gt;)</i>	
<i>datatype</i>	::=	<i>ArrayDecl(Num(str)+, &lt;datatype&gt;)</i>	<i>L_Array</i>
<i>exp</i>	::=	<i>Subscr(&lt;exp&gt;, &lt;exp&gt;)</i>   <i>Array(&lt;exp&gt;+)</i>	
<i>datatype</i>	::=	<i>StructSpec(Name(str))</i>	<i>L_Struct</i>
<i>exp</i>	::=	<i>Attr(&lt;exp&gt;, Name(str))</i>   <i>StructAssign(Name(str), &lt;exp&gt;+)</i>	
<i>decl_def</i>	::=	<i>StructDecl(Name(str),</i> <i>Alloc(Writeable(), &lt;datatype&gt;, Name(str))+)</i>	
<i>stmt</i>	::=	<i>If(&lt;exp&gt;, &lt;stmt&gt;*)</i>   <i>IfElse(&lt;exp&gt;, &lt;stmt&gt;*, &lt;stmt&gt;*)</i>	<i>L_If_Else</i>
<i>stmt</i>	::=	<i>While(&lt;exp&gt;, &lt;stmt&gt;*)</i>   <i>DoWhile(&lt;exp&gt;, &lt;stmt&gt;*)</i>	<i>L_Loop</i>
<i>exp</i>	::=	<i>Call(Name(str), &lt;exp&gt;*)</i>	<i>L_Fun</i>
<i>stmt</i>	::=	<i>Return(&lt;exp&gt;)</i>	
<i>decl_def</i>	::=	<i>FunDecl(&lt;datatype&gt;, Name(str),</i> <i>Alloc(Writeable(), &lt;datatype&gt;, Name(str))*)</i>   <i>FunDef(&lt;datatype&gt;, Name(str),</i> <i>Alloc(Writeable(), &lt;datatype&gt;, Name(str))* &lt;stmt&gt;*)</i>	
<i>file</i>	::=	<i>File(Name(str), &lt;decl_def&gt;*)</i>	<i>L_File</i>

**Grammatik 2.2.10:** Abstrakte Syntax der Sprache  $L_{PicoC}$

Man spricht hier von der „**Abstrakten Syntax der Sprache  $L_{PicoC}$** “ und meint hier mit der Sprache  $L_{PicoC}$  **nicht** die Sprache, welche durch die **Abstrakte Syntax** beschrieben wird. Es ist damit **immer** die Sprache gemeint, die **kompiliert** werden soll und zu deren Zweck die **Abstrakt Syntax** überhaupt definiert wird. Für die tatsächliche Sprache, die durch die **Abstrakt Syntax** beschrieben wird, interessiert man sich nie wirklich explizit. Diese **Redeart** wurde aus der **Quelle** G. Siek, *Course Webpage for Compilers (P423, P523, E313, and E513)* übernommen.

### 2.2.5.5 Codebeispiel

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

```

1 File
2   Name './verbose_dt_simple_ast_gen_array_decl_and_alloc.ast',
3   [
4     StructDecl
5       Name 'st',
6       [
7         Alloc
8           Writeable,
9           PtrDecl
10            Num '1',
11            ArrayDecl
12              [
13                Num '4',
14                Num '5'
15              ],
16            PtrDecl
17              Num '1',
18              IntType 'int',
19            Name 'attr'
20          ],
21      FunDef
22        VoidType 'void',
23        Name 'main',
24        [],
25        [
26          Exp
27            Alloc
28              Writeable,
29              ArrayDecl
30                [
31                  Num '3',
32                  Num '2'
33                ],
34              PtrDecl
35                Num '1',
36              PtrDecl
37                Num '1',
38              StructSpec
39                Name 'st',
40              Name 'var'
41            ]
42          ]

```

**Code 2.5:** *Aus vereinfachtem Ableitungsbaum generierter Abstrakter Syntaxbaum***2.2.5.6 Ausgabe des Abstrakter Syntaxbaum**

Ein **Knoten** eines **Abstrakter Syntaxbaum** kann entweder in der **Konkreter Syntax** der Sprache, für dessen Kompilierung er generiert wurde oder in der **Abstrakter Syntax**, die beschreibt, wie der Abstrakter Syntaxbaum selbst aufgebaut sein darf ausgegeben werden.

Das Ausgeben eines **Abstrakter Syntaxbaums** wird im **PicoC-Compiler** über die **Magische Methode** `__repr__()`<sup>17</sup> der Programmiersprache **Python** umgesetzt. Sobald ein **PicoC-Knoten** oder **RETI-Knoten** ausgegeben werden soll, gibt seine Magische Methode `__repr__()` eine nach der **Abstrakten** oder **Konkreten 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 **Depth-First-Search** Schema der gesamte **Abstract Sybtax Tree** 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**.

Im **PicoC-Compiler** wurden **Abstrakte** und **Konkrete Syntax** miteinander gemischt. Für **PicoC-Knoten** wurde die **Abstrakte Syntax** verwendet, da Passes schließlich auf **Abstrakter Syntaxbaums** operieren. Bei **RETI-Knoten** wurde die **Konkrete Syntax** verwendet, da **Maschinenbefehle** in **Konkreter 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 geschiefte Klammern ( ) usw., ob man die **RETI-Knoten** in **Abstrakter** oder **Konkreter Syntax** schreibt. Daher kann man auch einfach gleich die **RETI-Knoten** in **Konkreter Syntax** ausgeben und muss nicht beim letzten **Pass** daran denken, am Ende die **Konkrete**, statt der **Abstrakten Syntax** für die **RETI-Knoten** auszugeben.

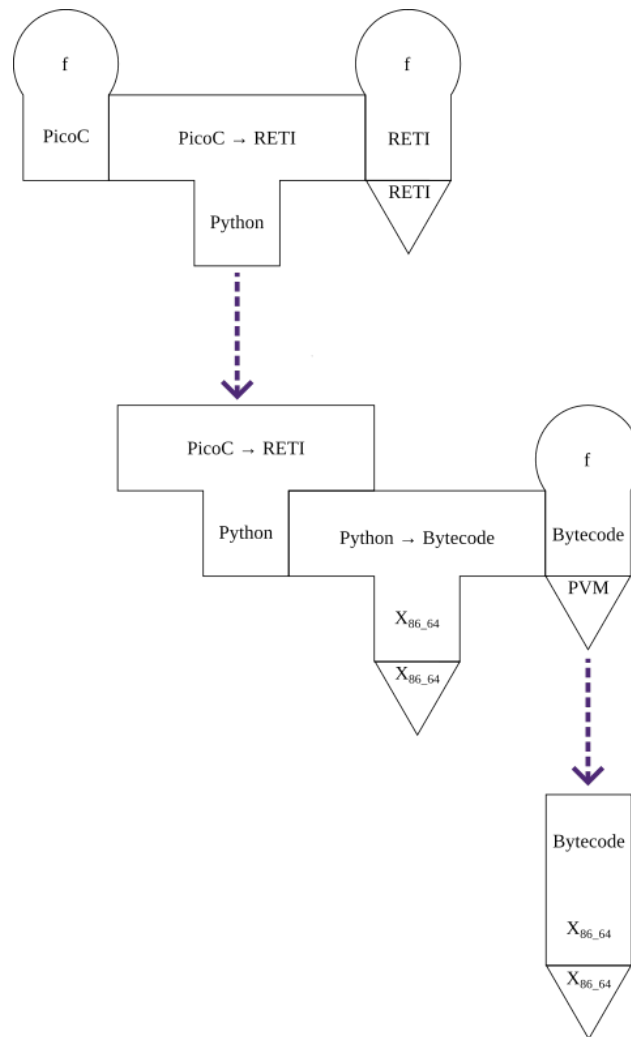
**2.3 Code Generierung**

Nach der Generierung eines **Abstrakter Syntaxbaum** als Ergebnis der **Lexikalischen** und **Syntaktischen Analyse** in Unterkapitel 1.4, wird in diesem Kapitel mit den verschiedenen **Kompositionen** von **PicoC-Knoten** und **RETI-Knoten** im **Abstrakter Syntaxbaum** als Basis das gewünschte Endprodukt des **PicoC-Compilers**, der **RETI-Code** generiert.

Man steht nun dem Problem gegenüber einen **Abstrakter Syntaxbaum** der Sprache  $L_{PicoC}$ , der durch die **Abstrakte Syntax** in Grammatik 2.2.10 spezifiziert ist in einen entsprechenden **Abstrakter Syntaxbaum** der Sprache  $L_{RETI}$  umzuformen. Das ganze lässt sich, wie in Unterkapitel 1.5 bereits beschrieben vereinfachen, indem man dieses Problem in mehrere **Passes** (Definition 1.43) herunterbricht.

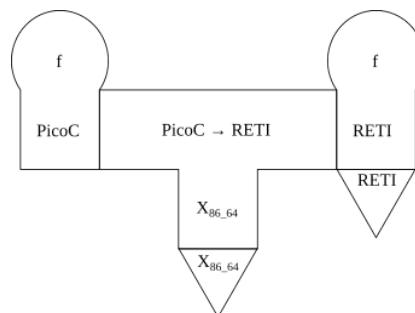
Beim **PicoC-Compiler** handelt es sich um einen **Cross-Compiler** (Definiton 1.5). Damit **RETI-Code** erzeugt werden kann, der auf der **RETI-Architektur** läuft, muss erst, wie im **T-Diagramm** (siehe Unterkapitel 1.1.1) in Abbildung 2.6 zu sehen ist, der **Python-Code** des **PicoC-Compilers** mittels eines Compilers, der z.B. auf einer  $X_{86,64}$ -Architektur laufen könnte zu **Bytecode** kompiliert werden. Dieser **Bytecode** wird dann von der **Python-Virtual-Machine** (PVM) interpretiert, welche wiederum auf einer  $X_{86,64}$ -Architektur laufen könnte. Und selbst dieses **T-Diagramm** könnte noch ausführlicher ausgedrückt werden, indem nachgeforscht wird, in welcher Sprache eigentlich die **Python-Virtual-Machine** geschrieben war, bevor sie zu  $X_{86,64}$  kompiliert wurde usw.

<sup>17</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-Diagramm** in Abbildung 2.6 lässt sich zusammenfassen, sodass man das **T-Diagramm** in Abbildung 2.7 erhält, in welcher direkt angegeben ist, dass der **PicoC-Compiler** in  $X_{86\_64}$ -Maschinensprache 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 **horizontalen**, 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 1.29) des **PicoC-Compilers** aufeinanderfolgen. In der **Pipe-Architektur** nutzt der jeweils nächste **Pass** den generierten **Abstrakter Syntaxbaum** des vorherigen Passes oder der Syntaktischen Analyse, um einen eigenen **Abstrakter Syntaxbaum** in seiner eigenen **Sprache** zu generieren.

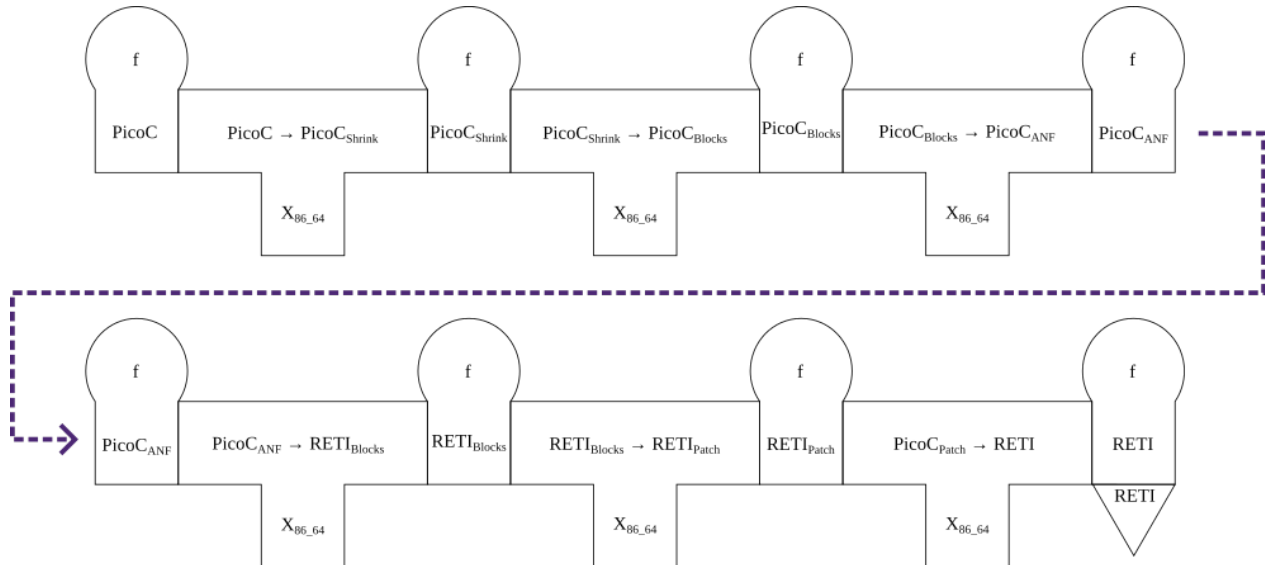


Abbildung 2.8: Architektur mit allen Passes ausgeschrieben

Im Unterkapitel 2.3.1 werden die unterschiedlichen **Passes** des PicoC-Compilers erklärt. In den darauffolgenden Unterkapiteln ??, ??, ?? und ?? zu **Pointern**, **Arrays**, **Structs** 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 verwendeten **Ansätze** zur Lösung dieser Probleme basieren auf der Vorlesung 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.

Um die verschiedenen Aspekte besser erklären zu können, werden **Codebeispiele** verwendet, in welchen ein kleines repräsentatives **PicoC-Programm** für einen spezifischen Aspekt in wichtigen **Zwischenstadien der Kompilierung** gezeigt wird<sup>18</sup>. Die **Codebeispiele** wurden alle mit dem **PicoC-Compiler** kompiliert und danach **nicht** mehr **verändert**, also genauso, wie der **PicoC-Compiler** sie kompiliert aus den Dateien in dieses Dokument eingelesen. Alle hier zur Repräsentation verwendeten **PicoC-Programme** lassen sich unter dem **Link**<sup>19</sup> finden und mithilfe der im Ordner `/code_examples` beiliegenden `Makefile` und dem Befehl `> make compile-all` genauso **kompilieren**, wie sie hier dargestellt sind<sup>20</sup>.

### 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 verschiedenen Aspekte<sup>21</sup> der

<sup>18</sup>Also die verschiedenen in den **Passes** generierten **Abstrakter Syntaxbaums**, sofern der **Pass** für den gezeigten Aspekt relevant ist.

<sup>19</sup>[https://github.com/matthejue/Bachelorarbeit/tree/master/code\\_examples](https://github.com/matthejue/Bachelorarbeit/tree/master/code_examples).

<sup>20</sup>Es wurden zu diesem Zweck spezielle neue **Command-line Optionen** erstellt, die bestimmte Kommentare **herausfiltern** und manche Container-Knoten **einzeilig** machen, damit die generierten **Abstrakter Syntaxbaums** in den verschiedenen Zwischenstufen der Kompilierung **nicht** zu langgestreckt und **überfüllt** mit Kommentaren sind.

<sup>21</sup>In kurz: **Pointer**, **Arrays**, **Structs** und **Funktionen**.

**Bachelorarbeit** wichtig.

Auf jedes Detail der einzelnen **Passes** wird in diesem Unterkapitel allerdings nicht eingegangen, da diese einerseits in den Unterkapiteln ??, ??, ?? und ?? zu **Pointern**, **Arrays**, **Structs** 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 **Pointer Arithmetik** **\*(pntr + i)** sich in der Untermenge der Sprache  $L_C$ , welche die Sprache  $L_{PicoC}$  darstellt genau gleich verhält, wie der **Operator** für den **Zugriff** auf den **Index** eines **Arrays** **ar[i]**.

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

#### 2.3.1.1.2 Abstrakte Syntax

Die **Abstrakte Syntax** der Sprache  $L_{PicoC\_Shrink}$  in Tabelle 2.3.1 ist fast **identisch** mit der **Abstrakten Syntax** der Sprache  $L_{PicoC}$  in Tabelle 2.2.10, nach welcher der **erste** Abstrakter Syntaxbaum in der **Syntaktischen Analyse** generiert wurde. Der einzige **Unterschied** liegt darin, dass es den Knoten **Deref(exp, exp)** in Tabelle 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 bereits in der **Abstrakten Syntax** der Sprache  $L_{PicoC}$  definiert ist.

<i>stmt</i>	::=	<i>SingleLineComment</i> ( <i>str</i> , <i>str</i> )   <i>RETIComment</i> ()	<i>L_Comment</i>
<i>un_op</i>	::=	<i>Minus</i> ()   <i>Not</i> ()	<i>L_Arith</i>
<i>bin_op</i>	::=	<i>Add</i> ()   <i>Sub</i> ()   <i>Mul</i> ()   <i>Div</i> ()   <i>Mod</i> ()   <i>Oplus</i> ()   <i>And</i> ()   <i>Or</i> ()	
<i>exp</i>	::=	<i>Name</i> ( <i>str</i> )   <i>Num</i> ( <i>str</i> )   <i>Char</i> ( <i>str</i> )   <i>BinOp</i> ( <i>&lt;exp&gt;</i> , <i>&lt;bin_op&gt;</i> , <i>&lt;exp&gt;</i> )   <i>UnOp</i> ( <i>&lt;un_op&gt;</i> , <i>&lt;exp&gt;</i> )   <i>Call</i> ( <i>Name</i> ('input'), <i>Empty</i> ())   <i>Call</i> ( <i>Name</i> ('print'), <i>&lt;exp&gt;</i> )	
<i>stmt</i>	::=	<i>Exp</i> ( <i>&lt;exp&gt;</i> )	
<i>un_op</i>	::=	<i>LogicNot</i> ()	<i>L_Logic</i>
<i>rel</i>	::=	<i>Eq</i> ()   <i>NEq</i> ()   <i>Lt</i> ()   <i>LtE</i> ()   <i>Gt</i> ()   <i>GtE</i> ()	
<i>bin_op</i>	::=	<i>LogicAnd</i> ()   <i>LogicOr</i> ()	
<i>exp</i>	::=	<i>Atom</i> ( <i>&lt;exp&gt;</i> , <i>&lt;rel&gt;</i> , <i>&lt;exp&gt;</i> )   <i>ToBool</i> ( <i>&lt;exp&gt;</i> )	
<i>type_qual</i>	::=	<i>Const</i> ()   <i>Writeable</i> ()	<i>L_Assign_Alloc</i>
<i>datatype</i>	::=	<i>IntType</i> ()   <i>CharType</i> ()   <i>VoidType</i> ()	
<i>exp</i>	::=	<i>Alloc</i> ( <i>&lt;type_qual&gt;</i> , <i>&lt;datatype&gt;</i> , <i>Name</i> ( <i>str</i> ))	
<i>stmt</i>	::=	<i>Assign</i> ( <i>&lt;exp&gt;</i> , <i>&lt;exp&gt;</i> )	
<i>datatype</i>	::=	<i>PntrDecl</i> ( <i>Num</i> ( <i>str</i> ), <i>&lt;datatype&gt;</i> )	<i>L_Pntr</i>
<i>exp</i>	::=	<b><i>Deref</i></b> ( <i>&lt;exp&gt;</i> , <i>&lt;exp&gt;</i> )   <i>Ref</i> ( <i>&lt;exp&gt;</i> )	
<i>datatype</i>	::=	<i>ArrayDecl</i> ( <i>Num</i> ( <i>str</i> ) +, <i>&lt;datatype&gt;</i> )	<i>L_Array</i>
<i>exp</i>	::=	<i>Subscr</i> ( <i>&lt;exp&gt;</i> , <i>&lt;exp&gt;</i> )   <i>Array</i> ( <i>&lt;exp&gt;</i> +)	
<i>datatype</i>	::=	<i>StructSpec</i> ( <i>Name</i> ( <i>str</i> ))	<i>L_Struct</i>
<i>exp</i>	::=	<i>Attr</i> ( <i>&lt;exp&gt;</i> , <i>Name</i> ( <i>str</i> ))   <i>StructAssign</i> ( <i>Name</i> ( <i>str</i> ), <i>&lt;exp&gt;</i> ) +)	
<i>decl_def</i>	::=	<i>StructDecl</i> ( <i>Name</i> ( <i>str</i> ), <i>Alloc</i> ( <i>Writeable</i> (), <i>&lt;datatype&gt;</i> , <i>Name</i> ( <i>str</i> )) +)	
<i>stmt</i>	::=	<i>If</i> ( <i>&lt;exp&gt;</i> , <i>&lt;stmt&gt;</i> *)   <i>IfElse</i> ( <i>&lt;exp&gt;</i> , <i>&lt;stmt&gt;</i> *, <i>&lt;stmt&gt;</i> *)	<i>L_If_Else</i>
<i>stmt</i>	::=	<i>While</i> ( <i>&lt;exp&gt;</i> , <i>&lt;stmt&gt;</i> *)   <i>DoWhile</i> ( <i>&lt;exp&gt;</i> , <i>&lt;stmt&gt;</i> *)	<i>L_Loop</i>
<i>exp</i>	::=	<i>Call</i> ( <i>Name</i> ( <i>str</i> ), <i>&lt;exp&gt;</i> *)	<i>L_Fun</i>
<i>stmt</i>	::=	<i>Return</i> ( <i>&lt;exp&gt;</i> )	
<i>decl_def</i>	::=	<i>FunDecl</i> ( <i>&lt;datatype&gt;</i> , <i>Name</i> ( <i>str</i> ), <i>Alloc</i> ( <i>Writeable</i> (), <i>&lt;datatype&gt;</i> , <i>Name</i> ( <i>str</i> ))*)   <i>FunDef</i> ( <i>&lt;datatype&gt;</i> , <i>Name</i> ( <i>str</i> ), <i>Alloc</i> ( <i>Writeable</i> (), <i>&lt;datatype&gt;</i> , <i>Name</i> ( <i>str</i> ))**, <i>&lt;stmt&gt;</i> *)	
<i>file</i>	::=	<i>File</i> ( <i>Name</i> ( <i>str</i> ), <i>&lt;decl_def&gt;</i> *)	<i>L_File</i>

**Grammatik 2.3.1:** Abstrakte Syntax der Sprache *L<sub>PiocC-Shrink</sub>*

Der **rot** markierte Knoten bedeutet, dass dieser im Vergleich zur vorherigen **Abstrakten Syntax** 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 Scholl, „Betriebssysteme“, welcher in der Vorlesung allerdings **rekursiv** implementiert ist.

Dieser **rekursive** Algorithmus 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**-Statemens 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 **Pointern**, **Arrays**, **Structs** und **Funktionen** mit eigenen **Codebeispielen** erklärt werden und alle sonstigen Details dem **Bachelorprojekt** zuzurechnen sind.

```

1 // based on a example program from Christoph Scholl's Operating Systems lecture
2
3 int faculty(int n){
4     int res = 1;
5     while (1) {
6         if (n == 1) {
7             return res;
8         }
9         res = n * res;
10        n = n - 1;
11    }
12 }
13
14 void main() {
15     print(faculty(4));
16 }

```

**Code 2.6:** PicoC Code für Codebespiel

In Code 2.7 sieht man den **Abstrakter 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       [
8         Alloc(Writable(), IntType('int'), Name('n'))
9       ],

```

<sup>22</sup>[https://github.com/matthejue/PicoC-Compiler/tree/new\\_architecture/tests](https://github.com/matthejue/PicoC-Compiler/tree/new_architecture/tests).



```

10  [
11    Assign(Alloc(Writable(), IntType('int'), Name('res')), Num('1')),
12    While
13      Num '1',
14      [
15        If
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    [
29      Exp(Call(Name('print'), [Call(Name('faculty'), [Num('4')])]))
30    ]
31  ]

```

Code 2.7: Abstrakter Syntaxbaum für Codebeispiel

Im **PicoC-Shrink-Pass** ändert sich nichts im Vergleich zum **Abstrakter 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(label)-` 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 Syntax

Zur Umsetzung dieses Passes ist es notwendig die **Abstrakte Syntax** der Sprache  $L_{PicoC\_Shrink}$  in Tabelle 2.3.1 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(label)-` und `IfElse(exp, stmts1, stmts2)-Knoten` ersetzt wurden. Die **Funktionsdefinition** `FunDef(<datatype>, Name(str), Alloc(Writable(), <datatype>), Name(str))*`, `<block>*` ist nun ein Container für **Blöcke** `Block(Name(str), <stmt>*)` und keine Statements `stmt` mehr. Das resultiert in der **Abstrakten Syntax** der Sprache  $L_{PicoC\_Blocks}$  in Tabelle 2.3.2.

<i>stmt</i>	::=	<i>SingleLineComment(str, str)</i>   <i>RETIComment()</i>	<i>L_Comment</i>
<i>un_op</i>	::=	<i>Minus()</i>   <i>Not()</i>	<i>L_Arith</i>
<i>bin_op</i>	::=	<i>Add()</i>   <i>Sub()</i>   <i>Mul()</i>   <i>Div()</i>   <i>Mod()</i>   <i>Oplus()</i>   <i>And()</i>   <i>Or()</i>	
<i>exp</i>	::=	<i>Name(str)</i>   <i>Num(str)</i>   <i>Char(str)</i>   <i>BinOp(&lt;exp&gt;, &lt;bin_op&gt;, &lt;exp&gt;)</i>   <i>UnOp(&lt;un_op&gt;, &lt;exp&gt;)</i>   <i>Call(Name('input'), Empty())</i>   <i>Call(Name('print'), &lt;exp&gt;)</i>	
<i>stmt</i>	::=	<i>Exp(&lt;exp&gt;)</i>	
<i>un_op</i>	::=	<i>LogicNot()</i>	<i>L_Logic</i>
<i>rel</i>	::=	<i>Eq()</i>   <i>NEq()</i>   <i>Lt()</i>   <i>LtE()</i>   <i>Gt()</i>   <i>GtE()</i>	
<i>bin_op</i>	::=	<i>LogicAnd()</i>   <i>LogicOr()</i>	
<i>exp</i>	::=	<i>Atom(&lt;exp&gt;, &lt;rel&gt;, &lt;exp&gt;)</i>   <i>ToBool(&lt;exp&gt;)</i>	
<i>type_qual</i>	::=	<i>Const()</i>   <i>Writeable()</i>	<i>L_Assign_Alloc</i>
<i>datatype</i>	::=	<i>IntType()</i>   <i>CharType()</i>   <i>VoidType()</i>	
<i>exp</i>	::=	<i>Alloc(&lt;type_qual&gt;, &lt;datatype&gt;, Name(str))</i>	
<i>stmt</i>	::=	<i>Assign(&lt;exp&gt;, &lt;exp&gt;)</i>	
<i>datatype</i>	::=	<i>PntrDecl(Num(str), &lt;datatype&gt;)</i>	<i>L_Pntr</i>
<i>exp</i>	::=	<i>Ref(&lt;exp&gt;)</i>	
<i>datatype</i>	::=	<i>ArrayDecl(Num(str)+, &lt;datatype&gt;)</i>	<i>L_Array</i>
<i>exp</i>	::=	<i>Subscr(&lt;exp&gt;, &lt;exp&gt;)</i>   <i>Array(&lt;exp&gt;+)</i>	
<i>datatype</i>	::=	<i>StructSpec(Name(str))</i>	<i>L_Struct</i>
<i>exp</i>	::=	<i>Attr(&lt;exp&gt;, Name(str))</i>   <i>StructAssign(Name(str), &lt;exp&gt;)+)</i>	
<i>decl_def</i>	::=	<i>StructDecl(Name(str),</i> <i>Alloc(Writeable(), &lt;datatype&gt;, Name(str))+)</i>	
<i>stmt</i>	::=	<i>If(&lt;exp&gt;, &lt;stmt&gt;*)</i>   <i>IfElse(&lt;exp&gt;, &lt;stmt&gt;*, &lt;stmt&gt;*)</i>	<i>L_If_Else</i>
<i>stmt</i>	::=	<i>While(&lt;exp&gt;, &lt;stmt&gt;*)</i>   <i>DoWhile(&lt;exp&gt;, &lt;stmt&gt;*)</i>	<i>L_Loop</i>
<i>exp</i>	::=	<i>Call(Name(str), &lt;exp&gt;*)</i>	<i>L_Fun</i>
<i>stmt</i>	::=	<i>Return(&lt;exp&gt;)</i>	
<i>decl_def</i>	::=	<i>FunDecl(&lt;datatype&gt;, Name(str),</i> <i>Alloc(Writeable(), &lt;datatype&gt;, Name(str))*)</i>   <i>FunDef(&lt;datatype&gt;, Name(str),</i> <i>Alloc(Writeable(), &lt;datatype&gt;, Name(str))*</i> , <i>&lt;block&gt;*)</i>	
<i>block</i>	::=	<i>Block(Name(str), &lt;stmt&gt;*)</i>	<i>L_Blocks</i>
<i>stmt</i>	::=	<i>GoTo(Name(str))</i>	
<i>file</i>	::=	<i>File(Name(str), &lt;decl_def&gt;*)</i>	<i>L_File</i>

### Grammatik 2.3.2: Abstrakte Syntax der Sprache $L_{Pioc\_Blocks}$

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

Die **Abstrakte Syntax** soll im Gegensatz zur **Konkreten Syntax** 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>.

Man bezeichnet hier die **Abstrakte Syntax** als „**Abstrakte Syntax der Sprache**  $L_{Picoc\_Blocks}$ “. Diese Sprache  $L_{Picoc\_Blocks}$  wird durch eine **Konkrete Syntax** beschrieben, die allerdings nicht weiter relevant ist, da in den **Passes** nur **Abstrakter Syntaxbaums** umgeformt werden. Es ist hierbei nur wichtig zu wissen, dass die **Abstrakte Syntax** theoretisch zur Kompilierung der Sprache  $L_{Picoc\_Blocks}$  definiert ist, also die Sprache  $L_{PicoC\_Blocks}$  nicht die Sprache ist, die von der **Abstrakten Syntax** beschrieben ist.

<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 **Abstract-Syntax-Tree** 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)-** oder **DoWhile(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.

```

1 File
2   Name './example_faculty_it.picoc_blocks',
3   [
4     FunDef
5       IntType 'int',
6       Name 'faculty',
7       [
8         Alloc(Writable(), IntType('int'), Name('n'))
9       ],
10      [
11        Block
12          Name 'faculty.6',
13          [
14            Assign(Alloc(Writable(), 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          [
21            IfElse
22              Num '1',
23              [
24                GoTo(Name('while_branch.4'))
25              ],
26              [
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')),
36             [
37                 GoTo(Name('if.3'))
38             ],
39             [
40                 GoTo(Name('if_else_after.2'))
41             ]
42     ],
43     Block
44         Name 'if.3',
45         [
46             Return(Name('res'))
47         ],
48     Block
49         Name 'if_else_after.2',
50         [
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
56         Name 'while_after.1',
57         []
58 ],
59 FunDef
60     VoidType 'void',
61     Name 'main',
62     [],
63     [
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 Codebeispiel

### 2.3.1.3 PicoC-ANF Pass

#### 2.3.1.3.1 Aufgabe

Die Aufgabe des **PicoC-ANF Passes** ist es den **Abstrakter Syntaxbaum** der Sprache  $L_{PicoC\_Blocks}$  in die **Abstrakte Syntax** der Sprache  $L_{PicoC\_ANF}$  umzuformen, welche in **A-Normalform** (Definition 1.50) und damit auch in **Monadischer Normalform** (Definition 1.46) ist. Um Wiederholung zu vermeiden wird zur Erklärung der **A-Normalform** auf Unterkapitel 1.5.2 verwiesen.

Zudem wird eine **Symboltabelle** (Definition 2.11) 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 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. `identifizier@scope`).<sup>24</sup>

#### Definition 2.11: 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 **Kompilervorgangs** 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 Syntax

Zur Umsetzung dieses Passes ist es notwendig die **Abstrakte Syntax** der Sprache  $L_{PicoC\_Blocks}$  in Tabelle 2.3.2 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 Syntax** der Sprache  $L_{PicoC\_ANF}$  in Grammatik 2.3.3 neue Knoten eingeführt werden.

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

<sup>23</sup>Bzw. der **Bezeichner**

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

<i>stmt</i>	::=	<i>SingleLineComment</i> ( <i>str</i> , <i>str</i> )   <i>RETIComment</i> ()	<i>L_Comment</i>
<i>un_op</i>	::=	<i>Minus</i> ()   <i>Not</i> ()	<i>L_Arith</i>
<i>bin_op</i>	::=	<i>Add</i> ()   <i>Sub</i> ()   <i>Mul</i> ()   <i>Div</i> ()   <i>Mod</i> ()   <i>Oplus</i> ()   <i>And</i> ()   <i>Or</i> ()	
<i>exp</i>	::=	<i>Name</i> ( <i>str</i> )   <i>Num</i> ( <i>str</i> )   <i>Char</i> ( <i>str</i> )   <i>Global</i> ( <i>Num</i> ( <i>str</i> ))   <i>Stackframe</i> ( <i>Num</i> ( <i>str</i> ))   <i>Stack</i> ( <i>Num</i> ( <i>str</i> ))   <i>BinOp</i> ( <i>Stack</i> ( <i>Num</i> ( <i>str</i> )), <i>bin_op</i> , <i>Stack</i> ( <i>Num</i> ( <i>str</i> )))   <i>UnOp</i> ( <i>un_op</i> , <i>Stack</i> ( <i>Num</i> ( <i>str</i> )))   <i>Call</i> ( <i>Name</i> ('input'), <i>Empty</i> ())   <i>Call</i> ( <i>Name</i> ('print'), <i>exp</i> )   <i>Exp</i> ( <i>exp</i> )	
<i>un_op</i>	::=	<i>LogicNot</i> ()	<i>L_Logic</i>
<i>rel</i>	::=	<i>Eq</i> ()   <i>NEq</i> ()   <i>Lt</i> ()   <i>LtE</i> ()   <i>Gt</i> ()   <i>GtE</i> ()	
<i>bin_op</i>	::=	<i>LogicAnd</i> ()   <i>LogicOr</i> ()	
<i>exp</i>	::=	<i>Atom</i> ( <i>Stack</i> ( <i>Num</i> ( <i>str</i> )), <i>rel</i> , <i>Stack</i> ( <i>Num</i> ( <i>str</i> )))   <i>ToBool</i> ( <i>Stack</i> ( <i>Num</i> ( <i>str</i> )))	
<i>type_qual</i>	::=	<i>Const</i> ()   <i>Writeable</i> ()	<i>L_Assign_Alloc</i>
<i>datatype</i>	::=	<i>IntType</i> ()   <i>CharType</i> ()   <i>VoidType</i> ()	
<i>exp</i>	::=	<i>Alloc</i> ( <i>type_qual</i> , <i>datatype</i> , <i>Name</i> ( <i>str</i> ))	
<i>stmt</i>	::=	<i>Assign</i> ( <i>Global</i> ( <i>Num</i> ( <i>str</i> )), <i>Stack</i> ( <i>Num</i> ( <i>str</i> )))   <i>Assign</i> ( <i>Stackframe</i> ( <i>Num</i> ( <i>str</i> )), <i>Stack</i> ( <i>Num</i> ( <i>str</i> )))   <i>Assign</i> ( <i>Stack</i> ( <i>Num</i> ( <i>str</i> )), <i>Global</i> ( <i>Num</i> ( <i>str</i> )))   <i>Assign</i> ( <i>Stack</i> ( <i>Num</i> ( <i>str</i> )), <i>Stackframe</i> ( <i>Num</i> ( <i>str</i> )))	
<i>datatype</i>	::=	<i>PntrDecl</i> ( <i>Num</i> ( <i>str</i> ), <i>datatype</i> )   <i>Ref</i> ( <i>Global</i> ( <i>str</i> ))   <i>Ref</i> ( <i>Stackframe</i> ( <i>str</i> ))   <i>Ref</i> ( <i>Subscr</i> ( <i>exp</i> , <i>exp</i> )   <i>Ref</i> ( <i>Attr</i> ( <i>exp</i> , <i>Name</i> ( <i>str</i> ))))	<i>L_Pntr</i>
<i>datatype</i>	::=	<i>ArrayDecl</i> ( <i>Num</i> ( <i>str</i> )+, <i>datatype</i> )	<i>L_Array</i>
<i>exp</i>	::=	<i>Subscr</i> ( <i>exp</i> , <i>Stack</i> ( <i>Num</i> ( <i>str</i> )))   <i>Array</i> ( <i>exp</i> )+)	
<i>datatype</i>	::=	<i>StructSpec</i> ( <i>Name</i> ( <i>str</i> ))	<i>L_Struct</i>
<i>exp</i>	::=	<i>Attr</i> ( <i>exp</i> , <i>Name</i> ( <i>str</i> ))   <i>Struct</i> ( <i>Assign</i> ( <i>Name</i> ( <i>str</i> ), <i>exp</i> ))+)	
<i>decl_def</i>	::=	<i>StructDecl</i> ( <i>Name</i> ( <i>str</i> ), <i>Alloc</i> ( <i>Writeable</i> (), <i>datatype</i> , <i>Name</i> ( <i>str</i> ))+)	
<i>stmt</i>	::=	<i>IfElse</i> ( <i>Stack</i> ( <i>Num</i> ( <i>str</i> )), <i>stmt</i> *, <i>stmt</i> *)	<i>L_If_Else</i>
<i>exp</i>	::=	<i>Call</i> ( <i>Name</i> ( <i>str</i> ), <i>exp</i> *)	<i>L_Fun</i>
<i>stmt</i>	::=	<i>StackMalloc</i> ( <i>Num</i> ( <i>str</i> ))   <i>NewStackframe</i> ( <i>Name</i> ( <i>str</i> ), <i>GoTo</i> ( <i>str</i> ))   <i>Exp</i> ( <i>GoTo</i> ( <i>Name</i> ( <i>str</i> )))   <i>RemoveStackframe</i> ()   <i>Return</i> ( <i>Empty</i> ())   <i>Return</i> ( <i>exp</i> )	
<i>decl_def</i>	::=	<i>FunDecl</i> ( <i>datatype</i> , <i>Name</i> ( <i>str</i> ) <i>Alloc</i> ( <i>Writeable</i> (), <i>datatype</i> , <i>Name</i> ( <i>str</i> ))*   <i>FunDef</i> ( <i>datatype</i> , <i>Name</i> ( <i>str</i> ), <i>Alloc</i> ( <i>Writeable</i> (), <i>datatype</i> , <i>Name</i> ( <i>str</i> ))* <i>block</i> *)	
<i>block</i>	::=	<i>Block</i> ( <i>Name</i> ( <i>str</i> ), <i>stmt</i> *)	<i>L_Blocks</i>
<i>stmt</i>	::=	<i>GoTo</i> ( <i>Name</i> ( <i>str</i> ))	
<i>file</i>	::=	<i>File</i> ( <i>Name</i> ( <i>str</i> ), <i>block</i> *)	<i>L_File</i>
<i>symbol_table</i>	::=	<i>SymbolTable</i> ( <i>symbol</i> *)	<i>L_Symbol_Table</i>
<i>symbol</i>	::=	<i>Symbol</i> ( <i>type_qual</i> , <i>datatype</i> , <i>name</i> , <i>val</i> , <i>pos</i> , <i>size</i> )	
<i>type_qual</i>	::=	<i>Empty</i> ()	
<i>datatype</i>	::=	<i>BuiltIn</i> ()   <i>SelfDefined</i> ()	
<i>name</i>	::=	<i>Name</i> ( <i>str</i> )	
<i>val</i>	::=	<i>Num</i> ( <i>str</i> )   <i>Empty</i> ()	
<i>pos</i>	::=	<i>Pos</i> ( <i>Num</i> ( <i>str</i> ), <i>Num</i> ( <i>str</i> ))   <i>Empty</i> ()	
<i>size</i>	::=	<i>Num</i> ( <i>str</i> )   <i>Empty</i> ()	

### 2.3.1.3.3 Codebeispiel

In Code 2.9 sieht man den **Abstract-Syntax-Tree** des **PiocC-ANF Passes** für das aus Unterkapitel 2.6 weitergeführte Beispiel, indem alle 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',
3   [
4     Block
5       Name 'faculty.6',
6       [
7         // Assign(Name('res'), Num('1'))
8         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')),
18        IfElse
19          Stack
20            Num '1',
21            [
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              ]
45        ],

```

```

46 Block
47   Name 'if.3',
48   [
49     // Return(Name('res'))
50     Exp(Stackframe(Num('1')))
51     Return(Stack(Num('1')))
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')))
61     // Assign(Name('n'), BinOp(Name('n'), Sub('-'), Num('1')))
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',
70   [
71     Return(Empty())
72   ],
73 Block
74   Name 'main.0',
75   [
76     StackMalloc(Num('2'))
77     Exp(Num('4'))
78     NewStackframe(Name('faculty'), GoTo(Name('addr@next_instr')))
79     Exp(GoTo(Name('faculty.6')))
80     RemoveStackframe()
81     Exp(ACC)
82     Exp(Call(Name('print'), [Stack(Num('1'))]))
83     Return(Empty())
84   ]
85 ]

```

Code 2.9: PicoC-ANF Pass für Codebeispiel

#### 2.3.1.4 RETI-Blocks Pass

##### 2.3.1.4.1 Aufgabe

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

##### 2.3.1.4.2 Abstrakte Syntax

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



sind  $\text{Exp}(\text{GoTo}(\text{str}))$ ,  $\text{Block}(\text{Name}(\text{str}), \langle \text{instr} \rangle^*)$  und  $\text{File}(\text{Name}(\text{str}), \langle \text{block} \rangle^*)$ , da das gesamte Konzept mit den **Blöcken** erst im **RETI-Pass** in Unterkapitel 2.3.8 aufgelöst wird.

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

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

### 2.3.1.4.3 Codebeispiel

In Code 2.10 sieht man den **Abstract-Syntax-Tree** des **RETI-Blocks Passes** für das aus Unterkapitel 2.6 weitergeführte Beispiel, indem die **Statements**, die durch entsprechende **PicoC-Knoten** im **Abstrakt Syntax Tree** der Sprache  $L_{PicoC\_ANF}$  in Grammatik 2.3.3 repräsentiert waren nun durch ihre entsprechenden **RETI-Knoten** ersetzt werden.

```

1 File
2   Name './example_faculty_it.reti_blocks',
3   [
4     Block
5       Name 'faculty.6',
6       [
7         # // Assign(Name('res'), Num('1'))
8         # Exp(Num('1'))
9         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',
22      [
23        # // 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
35     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     [
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',
77     [
78         # // Assign(Name('res'), BinOp(Name('n'), Mul('*'), Name('res')))
79         # Exp(Stackframe(Num('0')))
80         SUBI SP 1;
81         LOADIN BAF ACC -2;

```

```

82     STOREIN SP ACC 1;
83     # Exp(Stackframe(Num('1')))
84     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;
100    LOADIN BAF ACC -2;
101    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;
112    # Assign(Stackframe(Num('0')), Stack(Num('1')))
113    LOADIN SP ACC 1;
114    STOREIN BAF ACC -2;
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    [
122      # Return(Empty())
123      LOADIN BAF PC -1;
124    ],
125  Block
126    Name 'main.0',
127    [
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')))
135      MOVE BAF ACC;
136      ADDI SP 3;
137      MOVE SP BAF;
138      SUBI SP 4;

```

```

139     STOREIN BAF ACC 0;
140     LOADI ACC GoTo(Name('addr@next_instr'));
141     ADD ACC CS;
142     STOREIN BAF ACC -1;
143     # Exp(GoTo(Name('faculty.6'))))
144     Exp(GoTo(Name('faculty.6'))))
145     # RemoveStackframe()
146     MOVE BAF IN1;
147     LOADIN IN1 BAF 0;
148     MOVE IN1 SP;
149     # Exp(ACC)
150     SUBI SP 1;
151     STOREIN SP ACC 1;
152     LOADIN SP ACC 1;
153     ADDI SP 1;
154     CALL PRINT ACC;
155     # Return(Empty())
156     LOADIN BAF PC -1;
157 ]
158 ]

```

Code 2.10: RETI-Blocks Pass für Codebeispiel

Wenn der **Abstrakter Syntaxbaum** ausgegeben wird, ist die Darstellung nicht ausschließlich in **Abstrakter Syntax**, da die **RETI-Knoten** aus bereits im Unterkapitel 2.2.5.6 vermitteltem Grund in **Konkreter Syntax** ausgegeben 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 **Abstrakter Syntaxbaums**, 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>25</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**  $-2^{21}$  oder **größer**  $2^{21} - 1$  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**  $-(2^{31})$  oder **größer**  $2^{31} - 1$  ist, wird eine Fehlermeldung **TooLargeLiteral** ausgegeben.

#### 2.3.1.5.2 Abstrakte Syntax

<sup>25</sup>Das fällt unter die Themenbereiche des **Bachelorprojekts** und wird daher **nicht** genauer erläutert.

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

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

**Grammatik 2.3.5:** Abstrakte Syntax der Sprache  $L_{RETI\_Patch}$

### 2.3.1.5.3 Codebeispiel

In Code 2.11 sieht man den **Abstract-Syntax-Tree** des **Pioco-Patch Passes** für das aus Unterkapitel 2.6 weitergeführte Beispiel. Durch den **RETI-Patch Pass** wurde hier ein `start.<nummer>-Block`<sup>26</sup> eingesetzt, da die `main`-Funktion **nicht** die **erste** Funktion ist. Des Weiteren wurden durch diesen Pass einzelne `GoTo(Name(str))-Statements` entfernt<sup>27</sup>, die nur einen Sprung um **eine** Position entsprochen hätten.

```

1 File
2   Name './example_faculty_it.reti_patch',
3   [
4     Block
5       Name 'start.7',
6       [
7         # // Exp(GoTo(Name('main.0')))
8         Exp(GoTo(Name('main.0')))
9       ],
10    Block
11      Name 'faculty.6',
12      [
13        # // Assign(Name('res'), Num('1'))
14        # Exp(Num('1'))
15        SUBI SP 1;
16        LOADI ACC 1;
17        STOREIN SP ACC 1;
18        # Assign(Stackframe(Num('1')), Stack(Num('1')))
19        LOADIN SP ACC 1;
20        STOREIN BAF ACC -3;

```

<sup>26</sup>Dieser **Block** wurde im Code 2.8 markiert.

<sup>27</sup>Diese **entfernten** `GoTo(Name(str))`'s wurden ebenfalls im Code 2.8 markiert.

```

21     ADDI SP 1;
22     # // While(Num('1'), [])
23     # 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'))
31         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;
49         # 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'))
72         # Exp(Stackframe(Num('1')))
73         SUBI SP 1;
74         LOADIN BAF ACC -3;
75         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;
89         # Exp(Stackframe(Num('1')))
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;
101        STOREIN BAF ACC -3;
102        ADDI SP 1;
103        # // Assign(Name('n'), BinOp(Name('n'), Sub('-'), Num('1')))
104        # Exp(Stackframe(Num('0')))
105        SUBI SP 1;
106        LOADIN BAF ACC -2;
107        STOREIN SP ACC 1;
108        # Exp(Num('1'))
109        SUBI SP 1;
110        LOADI ACC 1;
111        STOREIN SP ACC 1;
112        # 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;
117        ADDI SP 1;
118        # Assign(Stackframe(Num('0')), Stack(Num('1')))
119        LOADIN SP ACC 1;
120        STOREIN BAF ACC -2;
121        ADDI SP 1;
122        # Exp(GoTo(Name('condition_check.5')))
123        Exp(GoTo(Name('condition_check.5')))
124    ],
125 Block
126     Name 'while_after.1',
127     [
128         # Return(Empty())
129         LOADIN BAF PC -1;
130     ],
131 Block
132     Name 'main.0',
133     [
134         # StackMalloc(Num('2'))

```

```

135     SUBI SP 2;
136     # Exp(Num('4'))
137     SUBI SP 1;
138     LOADI ACC 4;
139     STOREIN SP ACC 1;
140     # NewStackframe(Name('faculty'), GoTo(Name('addr@next_instr')))
141     MOVE BAF ACC;
142     ADDI SP 3;
143     MOVE SP BAF;
144     SUBI SP 4;
145     STOREIN BAF ACC 0;
146     LOADI ACC GoTo(Name('addr@next_instr'));
147     ADD ACC CS;
148     STOREIN BAF ACC -1;
149     # Exp(GoTo(Name('faculty.6')))
150     Exp(GoTo(Name('faculty.6')))
151     # RemoveStackframe()
152     MOVE BAF IN1;
153     LOADIN IN1 BAF 0;
154     MOVE IN1 SP;
155     # Exp(ACC)
156     SUBI SP 1;
157     STOREIN SP ACC 1;
158     LOADIN SP ACC 1;
159     ADDI SP 1;
160     CALL PRINT ACC;
161     # Return(Empty())
162     LOADIN BAF PC -1;
163 ]
164 ]

```

Code 2.11: RETI-Patch Pass für Codebeispiel

### 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 Konkrete und Abstrakte Syntax

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

Ausgegeben wird der finale **Maschinencode** allerdings in **Konkreter Syntax**, die sich aus den Grammatiken 2.3.6 und 2.3.7 für jeweils die **Lexikalische** und **Syntaktische Analyse** zusammensetzt. Der Grund, warum die **Konkrete Syntax** 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.



<i>dig_no_0</i>	::=	"1"   "2"   "3"   "4"   "5"   "6"	<i>L_Program</i>
		"7"   "8"   "9"	
<i>dig_with_0</i>	::=	"0"   <i>dig_no_0</i>	
<i>num</i>	::=	"0"   <i>dig_no_0</i> <i>dig_with_0</i> *   "-" <i>dig_no_0</i> *	
<i>letter</i>	::=	"a"..."Z"	
<i>name</i>	::=	<i>letter</i> ( <i>letter</i>   <i>dig_with_0</i>   _)*	
<i>reg</i>	::=	"ACC"   "IN1"   "IN2"   "PC"   "SP"	
		"BAF"   "CS"   "DS"	
<i>arg</i>	::=	<i>reg</i>   <i>num</i>	
<i>rel</i>	::=	"=="   "!="   "<"   "<="   ">"	
		">="   "_NOP"	

**Grammatik 2.3.6:** Konkrete Syntax der Sprache *L<sub>RETI</sub>* für die Lexikalische Analyse in EBNF

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

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

<i>reg</i>	::=	<i>ACC</i> ()   <i>IN1</i> ()   <i>IN2</i> ()   <i>PC</i> ()   <i>SP</i> ()   <i>BAF</i> ()	<i>L_RETI</i>
		<i>CS</i> ()   <i>DS</i> ()	
<i>arg</i>	::=	<i>Reg</i> ( <i>&lt;reg&gt;</i> )   <i>Num</i> ( <i>str</i> )	
<i>rel</i>	::=	<i>Eq</i> ()   <i>NEq</i> ()   <i>Lt</i> ()   <i>LtE</i> ()   <i>Gt</i> ()   <i>GtE</i> ()	
		<i>Always</i> ()   <i>NOp</i> ()	
<i>op</i>	::=	<i>Add</i> ()   <i>Addi</i> ()   <i>Sub</i> ()   <i>Subi</i> ()   <i>Mult</i> ()   <i>Multi</i> ()	
		<i>Div</i> ()   <i>Divi</i> ()   <i>Mod</i> ()   <i>Modi</i> ()   <i>Oplus</i> ()   <i>Oplusi</i> ()	
		<i>Or</i> ()   <i>Ori</i> ()   <i>And</i> ()   <i>Andi</i> ()	
		<i>Load</i> ()   <i>Loadin</i> ()   <i>Loadi</i> ()   <i>Store</i> ()   <i>Storein</i> ()   <i>Move</i> ()	
<i>instr</i>	::=	<i>Instr</i> ( <i>&lt;op&gt;</i> , <i>&lt;arg&gt;</i> +)   <i>Jump</i> ( <i>&lt;rel&gt;</i> , <i>Num</i> ( <i>str</i> ))   <i>Int</i> ( <i>Num</i> ( <i>str</i> ))	
		<i>RTI</i> ()   <i>Call</i> ( <i>Name</i> ('print'), <i>&lt;reg&gt;</i> )   <i>Call</i> ( <i>Name</i> ('input'), <i>&lt;reg&gt;</i> )	
		<i>SingleLineComment</i> ( <i>str</i> , <i>str</i> )	
		<i>Instr</i> ( <i>Loadi</i> (), [ <i>Reg</i> ( <i>Acc</i> ()), <i>GoTo</i> ( <i>Name</i> ( <i>str</i> ))])   <i>Jump</i> ( <i>Eq</i> (), <i>GoTo</i> ( <i>Name</i> ( <i>str</i> )))	
<i>program</i>	::=	<i>Program</i> ( <i>Name</i> ( <i>str</i> ), <i>&lt;instr&gt;</i> *)	
<i>instr</i>	::=	<i>Exp</i> ( <i>GoTo</i> ( <i>str</i> ))   <i>Exit</i> ( <i>Num</i> ( <i>str</i> ))	<i>L_PicoC</i>
<i>block</i>	::=	<i>Block</i> ( <i>Name</i> ( <i>str</i> ), <i>&lt;instr&gt;</i> *)	
<i>file</i>	::=	<i>File</i> ( <i>Name</i> ( <i>str</i> ), <i>&lt;block&gt;</i> *)	

**Grammatik 2.3.8:** Abstrakte Syntax der Sprache *L<sub>RETI</sub>*

### 2.3.1.6.3 Codebeispiel

Nach dem **RETI-Pass** ist das Programm komplett in **RETI-Knoten** übersetzt, die allerdings in ihrer **Konkreten 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** ausschließlich aus **RETI-Ausdrücken** bestehen<sup>28</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 **Abstrakter Syntaxbaum**, wenn er in eine Datei ausgegeben wird, direkt **RETI-Code** in **menschenlesbarer Repräsentation** erzeugt.

```

1 # // 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;
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;

```

<sup>28</sup>Wie z.B. `LOADI ACC GoTo(Name('addr@next_instr')), Exp(GoTo(Name('main.0')))` und `JUMP== GoTo(Name('if_else_after.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;
02 # Exp(Num('4'))
03 SUBI SP 1;
04 LOADI ACC 4;
05 STOREIN SP ACC 1;
06 # NewStackframe(Name('faculty'), GoTo(Name('addr@next_instr'))
07 MOVE BAF ACC;
08 ADDI SP 3;
09 MOVE SP BAF;
10 SUBI SP 4;
11 STOREIN BAF ACC 0;
12 LOADI ACC 80;
13 ADD ACC CS;
14 STOREIN BAF ACC -1;
15 # Exp(GoTo(Name('faculty.6'))
16 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;
25 ADDI SP 1;
26 CALL PRINT ACC;
27 # Return(Empty())
28 LOADIN BAF PC -1;
```

**Code 2.12:** *RETI Pass für Codebespiel*

# Literatur

## Online

- *A-Normalization: Why and How (with code)*. URL: <https://matt.might.net/articles/a-normalization/> (besucht am 23.07.2022).
- *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](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).
- *Errors in C/C++ - GeeksforGeeks*. URL: <https://www.geeksforgeeks.org/errors-in-cc/> (besucht am 10.05.2022).
- *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).
- *JSON parser - Tutorial — Lark documentation*. URL: [https://lark-parser.readthedocs.io/en/latest/json\\_tutorial.html](https://lark-parser.readthedocs.io/en/latest/json_tutorial.html) (besucht am 09.07.2022).
- Ljohhuh. *What is an immediate value?* 4. Apr. 2018. URL: <https://reverseengineering.stackexchange.com/q/17671> (besucht am 13.04.2022).
- *Parsing Expressions · Crafting Interpreters*. URL: <https://www.craftinginterpreters.com/parsing-expressions.html> (besucht am 09.07.2022).
- *Transformers & Visitors — Lark documentation*. URL: <https://lark-parser.readthedocs.io/en/latest/visitors.html> (besucht am 09.07.2022).
- *Welcome to Lark's documentation! — Lark documentation*. URL: <https://lark-parser.readthedocs.io/en/latest/> (besucht am 31.07.2022).

- *What is Bottom-up Parsing?* URL: <https://www.tutorialspoint.com/what-is-bottom-up-parsing> (besucht am 22.06.2022).
- *What is Top-Down Parsing?* URL: <https://www.tutorialspoint.com/what-is-top-down-parsing> (besucht am 22.06.2022).

## 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).

## Artikel

- Earley, J. und Howard E. Sturgis. „A formalism for translator interactions“. In: *CACM* (1970). DOI: [10.1145/355598.362740](https://doi.org/10.1145/355598.362740).

## Vorlesungen

- Nebel, Prof. Dr. Bernhard. „Theoretische Informatik“. Vorlesung. Vorlesung. Universität Freiburg, 2020. URL: [http://gki.informatik.uni-freiburg.de/teaching/ss20/info3/index\\_de.html](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](https://abs.informatik.uni-freiburg.de/src/teach_main.php?id=157) (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).
- Westphal, Dr. Bernd. „Softwaretechnik“. Vorlesung. Vorlesung. Universität Freiburg, 2021. URL: <https://swt.informatik.uni-freiburg.de/teaching/SS2021/swtv1> (besucht am 19.07.2022).

## Sonstige Quellen

- Bolingbroke, Maximilian C. und Simon L. Peyton Jones. „Types are calling conventions“. In: *Proceedings of the 2nd ACM SIGPLAN symposium on Haskell - Haskell '09*. the 2nd ACM SIGPLAN symposium. Edinburgh, Scotland: ACM Press, 2009, S. 1. ISBN: 978-1-60558-508-6. DOI: [10.1145/1596638.1596640](https://doi.org/10.1145/1596638.1596640). URL: <http://portal.acm.org/citation.cfm?doid=1596638.1596640> (besucht am 23.07.2022).
- *Lark - a parsing toolkit for Python*. 26. Apr. 2022. URL: <https://github.com/lark-parser/lark> (besucht am 28.04.2022).
- Shinan, Erez. *lark: a modern parsing library*. Version 1.1.2. URL: <https://github.com/lark-parser/lark> (besucht am 31.07.2022).

- *Syntax*. In: *Wiktionary*. Page Version ID: 9196998. 7. Juni 2022. URL: <https://de.wiktionary.org/w/index.php?title=Syntax&oldid=9196998> (besucht am 31.07.2022).