# Albert Ludwigs Universität Freiburg

#### TECHNISCHE FAKULTÄT

### PicoC-Compiler

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

BACHELORARBEIT

Abgabedatum: 13. September 2022

Autor: Jürgen Mattheis

Gutachter: Prof. Dr. Scholl

Betreung: M.Sc. Seufert

Eine Bachelorarbeit am Lehrstuhl für

Betriebssysteme

#### **ERKLÄRUNG**

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

## Danksagungen

Bevor der Inhalt dieser Schrifftlichen Ausarbeitung der Bachelorarbeit anfängt, will ich einigen Personen noch meinen Dank aussprechen.

Ich schreibe die folgenden Danksagungen nicht auf eine bestimmte Weise, wie es sich vielleicht etabliert haben sollte Danksagungen zu schreiben und verwende auch keine künstlichen Floskeln, wie "mein aufrichtigster Dank" oder "aus tiefstem Herzen", sondern drücke im Folgenden die Dinge nur so aus, wie ich sie auch wirklich meine.

Estmal, ich hatte selten im Studium das Gefühl irgendwo Kunde zu sein, aber bei dieser Bachelorarbeit und dem vorangegangenen Bachelorprojekt hatte ich genau diese Gefühl, obwohl die Verhältnisse eigentlich genau umgekehrt sein sollten. Die Umgang mit mir wahr echt unglaublich nett und unbürokratisch, was ich als keine Selbverständlichkeit ansehe und sehr wertgeschätzt habe.

An erster Stelle will ich zu meinem Betreuer M.Sc. Tobias Seufert kommen, der netterweise auch bereits die Betreuung meines Bachelorprojektes übernommen hatte. Wie auch während des Bachelorprojektes, haben wir uns auch bei den Meetings während der Bachelorarbeit hervorragend verstanden. Dabei ging die Freundlichkeit und das Engagement seitens Tobias weit über das heraus, was man bereits als eine gute Betreuung bezeichnen würde.

Es gibt verschiedene Typen von Menschen, es gibt Leute, die nur genauso viel tun, wie es die Anforderungen verlangen und nichts darüberhinaus tun, wenn es nicht einen eigenen Vorteil für sie hat und es gibt Personen, die sich für nichts zu Schade sind und dies aus einer Philanthropie oder Leidenschafft heraus tun, auch wenn es für sie keine Vorteile hat. Tobias¹ konnte ich während der langen Zeit, die er mein Bachelorprojekt und dann meine Bachelorarbeit betreut hat eindeutig als letzteren Typ Mensch einordnen.

Er war sich nie zu Schade für meine vielen Fragen während der Meetings, auch wenn ich meine Zeit ziemlich oft überzogen habe<sup>2</sup>, er hat sich bei der Korrektur dieser Schrifftlichen Ausarbeitung sogar die Mühe gemacht bei den einzelnen Problemstellen längere, wirklich hilfreiche Textkommentare zu verfassen und obendrauf auch noch Tippfehler usw. angemerkt und war sich nicht zu Schade die Rolle des Nachrichtenübermittlers zwischen mir und Prof. Dr. Scholl zu übernehmen. All dies war absolut keine Selbverständlichkeit, vor allem wenn ich die Betreuung anderer Studenten, die ich kenne mit der vergleiche, die mir zu Teil wurde.

An den Kommentar zu meinem Betreuer Tobias will ich einen Kommentar zu meinem Gutachter Prof. Dr. Scholl anschließen. Wofür ich meinem Gutachter Prof. Dr. Scholl sehr dankbar bin, ist, dass er meine damals sehr ambitionierten Ideen für mögliche Funktionalitäten, die ich in den PicoC-Compiler für die Bachelorarbeit implementierten wollte runtergeschraubt hat. Man erlebt es äußerst selten im Studium, dass Studenten freiwillig weniger Arbeit gegeben wird.

Bei den für die Bachelorarbeit zu implementierenden Funktionalitäten gab es bei der Implementierung viele unerwartete kleine Details, die ich vorher garnicht bedacht hatte, die in ihrer Masse unerwartet viel Zeit zum Implementieren gebraucht haben. Mit den von Prof. Dr. Scholl festgelegten Funktionalitäten für die Bachelorarbeit ist der Zeitplan jedoch ziemlich perfekt aufgegangen. Mit meinen ambitionierten Plänen wäre es bei der Bachelorarbeit dagegeben wohl mit der Zeit äußerst kritisch geworden. Das Prof. Dr. Scholl mir zu

<sup>&</sup>lt;sup>1</sup>Wie auch Prof. Dr. Scholl. Hier geht es aber erstmal um Tobias.

<sup>&</sup>lt;sup>2</sup>Wofür ich mich auch nochmal Entschuldigen will.

seinem eigenen Nachteil $^3$  weniger Arbeit aufgebrummt hat empfand ich als ich eine äußerst nette Geste, die ich sehr geschätzt habe.

Wie mein Betreuer M.Sc. Tobias Seufert und wahrscheinlich auch mein Gutachter Prof. Dr. Scholl im Verlauf dieser Bachelorarbeit und des vorangegangenen Bachelorprojektes gemerkt haben, kann ich schon manchmal ziemlich eigensinnigen sein, bei der Weise, wie ich bestimmte Dinge umsetzen will. Ich habe es sehr geschätzt, dass mir das durchgehen gelassen wurde. Es ist, wie ich die Universitätswelt als Student erlebe bei Arbeitsvorgaben keine Selbverständlichkeit, dass dem Studenten überhaupt die Freiheit und das Vertrauen gegeben wird diese auf seine eigenen Weise umzusetzen.

Vor allem, da mein eigenes Vorgehen größtenteils Vorteile für mich hatte, da ich auf diese Weise am meisten über Compilerbau gelernt hab und eher Nachteile für Prof. Dr. Scholl, da mein eigenes Vorgehen entsprechend mehr Zeit brauchte und ich daher als Bachelorarbeit keinen dazu passenden RETI-Emulator mit Graphischer Anzeige implementieren konnte, da die restlichen Funktionalitäten des PicoC-Compilers noch implementiert werden mussten.

Glücklicherweise gibt es aber doch noch einen passenden RETI-Emulator, der den PicoC-Compiler über seine Kommandozeilenargumente aufruft, um ein PicoC-Programm visuell auf einer RETI-CPU auszuführen. Für dessen Implementierung hat sich Michel Giehl netterweise zur Verfügung gestellt. Daher Danke auch an Michel Giehl, dass er sich mit meinem PicoC-Compiler ausgeinandergesetzt hat und diesen in seinen RETI-Emulator integriert hat, sodass am Ende durch unsere beiden Arbeiten ein anschauliches Lerntool für die kommenden Studentengenerationen entstehen konnte. Vor allem da er auch mir darin vertrauen musste, dass ich mit meinem PicoC-Compiler nicht irgendeinen Misst baue. Der RETI-Emulator von Michel Giehl ist unter Link<sup>5</sup> zu finden.

Mir hat die Implementierung des PicoC-Compilers tatsächlich ziemlich viel Spaß gemacht, da Compilerbau auch in mein perönliches Interessengebiet fällt<sup>6</sup>. Das Aufschreiben dieser Schrifftlichen Ausarbeitung hat mir dagegen eher weniger Spaß gemacht<sup>7</sup>. Wobei ich allerdings sagen muss, dass ich eine große Erleichterung verspüre das ganze Wissen über Compilerbau mal aufgeschrieben zu haben, damit ich mir keine Sorgen machen muss dieses ziemlich nützliche Wissen irgendwann wieder zu vergessen. Es hilft einem auch als Programmierer ungemein weiter zu wissen, wie ein Compiler unter der Haube funktioniert, da man sich so viel besser merken, wie eine bestimmte Funktionalität einer Programmiersprache zu verwenden ist. Manch eine Funktionalität einer Programmiesprache kann in der Verwendung ziemlich wilkürlich erscheinen, wenn man die technische Umsetzung dahinter im Compiler nicht kennt.

Ich wollte mich daher auch noch dafür Bedanken, dass mir ein so ergiebiges und interessantes Thema als Bachelorarbeit vorgeschlagen wurde und vor allem, dass auch das Vertrauen in mich gesteckt wurde, dass ich am Ende auch einen funktionsfähigen, sauber programmierten und gut durchdachten Compiler implementiere.

Zum Schluss nochmal ein abschließendes Danke an meinen Betreuer M.Sc Seufert und meinen Gutachter Prof. Dr. Scholl für die Betreuung und Bereitstellung dieser interessanten Bachelorarbeit und des vorangegangenen Bachelorprojektes und Michel Giehl für das Integrieren des PicoC-Compilers in seinen RETI-Emulator.

<sup>&</sup>lt;sup>3</sup>Der PicoC-Compiler hätte schließlich mehr Funktionalitäten haben können.

<sup>&</sup>lt;sup>4</sup>Vielleicht finde ich ja noch im nächsten Semester während des Betriebssysteme Tutorats noch etwas Zeit einige weitere Features einzubauen oder möglicherweise im Rahmen eines Masterprojektes <sup>3</sup>.

 $<sup>^5</sup>$ https://github.com/michel-giehl/Reti-Emulator.

<sup>&</sup>lt;sup>6</sup>Womit nicht alle Studenten so viel Glück haben.

<sup>&</sup>lt;sup>7</sup>Dieses ständige überlegen, wo man möglicherweise eine Erklärlücke hat, ob man nicht was wichtiges ausgelassen hat usw.

# Inhaltsverzeichnis

Abbildungsverzeichnis	Ι
Codeverzeichnis	II
Tabellenverzeichnis	III
Definitionsverzeichnis	$\mathbf{V}$
Grammatikverzeichnis	VI
1 Theoretische Grundlagen  1.1 Compiler und Interpreter 1.1.1 T-Diagramme  1.2 Formale Sprachen 1.2.1 Ableitungen 1.2.2 Präzedenz und Assoziativität  1.3 Lexikalische Analyse 1.4 Syntaktische Analyse 1.5 Code Generierung 1.5.1 Monadische Normalform 1.5.2 A-Normalform 1.5.3 Ausgabe des Maschinencodes  1.6 Fehlermeldungen	1 4 6 9 12 13 16 25 27 31 31
Literatur	A

# Abbildungsverzeichnis

1.1	Horinzontale Übersetzungszwischenschritte zusammenfassen.
1.2	Veranschaulichung von Linksassoziativität und Rechtsassoziativität
1.3	Veranschaulichung von Präzedenz
1.4	Veranschaulichung der Lexikalischen Analyse
1.5	Veranschaulichung des Unterschieds zwischen Ableitungsbaum und Abstraktem Syntaxbaum. 23
1.6	Veranschaulichung der Syntaktischen Analyse
1.7	Codebeispiel dafür Code in die Monadische Normalform zu bringen
1.8	Übersicht über Komplexe, Atomare, Unreine und Reine Ausdrücke
1.9	Codebeispiel für das Entfernen Komplexer Ausdrücke aus Operationen

# Codeverzeichnis

# **Tabellenverzeichnis**

1.1	Beispiele für	Lexeme und ihre ents	prechenden Tokens.		1
-----	---------------	----------------------	--------------------	--	---

# Definitionsverzeichnis

1.1	Pipe-Filter Architekturpattern		 			
1.2	Interpreter		 			-
1.3	Compiler		 			6
1.4	Maschinensprache		 			
1.5	Immediate		 			;
1.6	Cross-Compiler		 			
1.7	T-Diagramm Programm					4
1.8	T-Diagramm Übersetzer (bzw. eng. Translator)					4
1.9	T-Diagramm Interpreter					ļ
1.10	Symbol		 			(
1.11	·					(
1.12	Wort					(
	Formale Sprache					,
1.14	Syntax		 			,
	Semantik					,
	Formale Grammatik					,
	Chromsky Hierarchie					8
	Reguläre Grammatik					9
1.19	Kontextfreie Grammatik		 			9
1.20	Wortproblem		 			9
	1-Schritt-Ableitungsrelation					9
	Ableitungsrelation					10
	Links- und Rechtsableitungableitung					10
	Formaler Ableitungsbaum					10
1.25	Mehrdeutige Grammatik		 			12
	Assoziativität					1:
	Präzedenz					1;
1.28	Lexeme		 			13
1.29	Token		 			1
1.30	Lexer (bzw. Scanner oder auch Tokenizer)		 			1
	Literal					1!
1.32	Konkrete Syntax		 			1'
1.33	Konkrete Grammatik		 			18
	Ableitungsbaum (bzw. Konkreter Syntaxbaum, engl. Derivation Tree) .					18
1.35	Parser		 			18
1.36	Erkenner (bzw. engl. Recognizer)		 			19
1.37	Transformer		 			20
1.38	Visitor		 			2
1.39	Abstrakte Syntax		 			2
1.40	Abstrakte Grammatik		 			2
1.41	Abstrakter Syntaxbaum (bzw. engl. Abstract Syntax Tree, kurz AST) .		 			2
	Pass					2!
1.43	Ausdruck (bzw. engl. Expression)		 			2!
1.44	Anweisung (bzw. engl. Statement)		 			20
	Reiner Ausdruck / Reine Anweisung (bzw. engl. pure expression)					20
1.46	Unreiner Ausdruck / Unreine Anweisung		 			20
1.47	Monadische Normalform (bzw. engl. monadic normal form)					2'

1.48	ocation	3
1.49	tomarer Ausdruck	3
1.50	Complexer Ausdruck	)
1.51	-Normalform (ANF)	9
1.52	ehlermeldung $\ldots \ldots \ldots \ldots \ldots 32$	2

# Grammatikverzeichnis

1.1	Grammatik für einen Ableitungsbaum in EBNF	11
1.2	Produktionen für Ableitungsbaum in EBNF	23
1.3	Produktionen für Abstrakten Syntaxbaum in ASF	23

# Theoretische Grundlagen

In diesem Kapitel wird auf die Theoretischen Grundlagen eingegangen, die zum Verständnis der Implementierung in Kapitel ?? notwendig sind. Zuerst wird in Unterkapitel 1.1 genauer darauf eingegangen was ein Compiler und Interpreter eigentlich sind und damit in Verbindung stehende Begriffe und T-Diagramme erklärt. Danach wird in Unterkapitel 1.2 eine kleine Einführung zu einem der Grundpfeiler des Compilerbau, den Formalen Sprachen gegeben. Danach werden die einzelnen Filter des üblicherweise bei der Implementierung von Compilern genutzten Pipe-Filter-Architekturpatterns (Definition 1.1) nacheinander erklärt. Die Filter beinhalten die Lexikalische Analyse 1.3, Syntaktische Analyse 1.4 und Code Generierung 1.5. Zum Schluss wird in Unterkapitel 1.6 darauf eingegangen in welchen Situationen Fehlermeldungen auszugeben sind.

#### Definition 1.1: Pipe-Filter Architekturpattern

1

Ist ein Archikteturpattern, 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. Bei der Weiterverarbeitung können Teile der Eingabe entfernt, hinzugefügt oder vollständig ersetzt werden.

Eine Pipe stellt ein Bindeglied zwischen zwei Filtern dar. ab



<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 Anlehung 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".

#### 1.1 Compiler und Interpreter

Die wohl wichtigsten zu klärenden Begriffe, sind die eines Compilers (Definition 1.3) und eines Interpreters (Definition 1.2), 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.

#### Anmerkung Q

Des Weiteren wurde zur Qualitätsicherung ein RETI-Interpreter implementiert, um mithilfe des GCC<sup>a</sup> und von Tests die Beziehungen in 1.3.1 zu belegen (siehe Unterkapitel ??), weshalb es auch nochmal wichtig ist die Definition eines Interpreters eingeführt zu haben.

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

#### Definition 1.2: Interpreter

Z

Programm, dass die Anweisungen eines Programmes mehr oder weniger direkt ausführt.

In einer konkreten Implementierung arbeitet ein Interpreter auf einem compilerinternen Abstrakten Syntaxbaum (wird später eingeführt unter Definition 1.41) und führt je nach Komposition der Knoten des Abstrakten Syntaxbaumes, auf die er während des Darüber-Iterierens stösst unterschiedliche Aufgaben aus.<sup>a</sup>

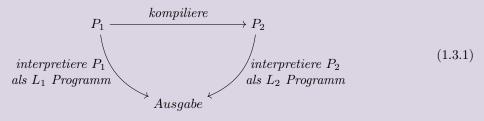
<sup>a</sup>G. Siek, Course Webpage for Compilers (P423, P523, E313, and E513).

#### Definition 1.3: Compiler

Z

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

Dabei muss gelten, dass die beiden Programme  $P_1$  und  $P_2$ , wenn sie von Interpretern ihrer jeweiligen Sprachen  $L_1$  und  $L_2$  interpretiert werden die gleiche Ausgabe haben. Dies ist in Diagramm 1.3.1 dargestellt. Beide Programme  $P_1$  und  $P_2$  sollen die gleiche Semantik (Definition 1.15) haben und unterscheiden sich nur syntaktisch (Definition 1.14).



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

Üblicherweise kompiliert ein Compiler ein Programm, das in einer Programmiersprache geschrieben ist zu einem Maschinenprogramm, das in Maschinensprache (Definition 1.4) geschrieben ist, aber es gibt z.B. auch Transpiler (Definition ?? im ??) oder Cross-Compiler (Definition 1.6)<sup>1</sup>.

#### Definition 1.4: Maschinensprache

Programmiersprache, deren mögliche Programme die hardwarenaheste Repräsentation von zuvor hierzu kompilierten bzw. assemblierten Programmen darstellen.

Jeder Maschinenbefehl entspricht einer bestimmten Aufgabe, welche eine CPU im einfachen Fall in einem Zyklus der Fetch- und Execute-Phase, genauergesagt 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 entworfen, dass sie sich innerhalb bestimmter Wortbreiten, die Zweierpotenzen sind kodieren lassen. Im einfachsten Fall innerhalb einer Speicherzelle des Hauptspeichers.<sup>a</sup>

<sup>&</sup>lt;sup>1</sup>Des Weiteren sind Maschinensprache und Assemblersprache (Definition ?? im ??) voneinander zu unterscheiden.

Die Programme<sup>b</sup> einer Maschinensprache können dabei in verschiedenen Repräsentationen dargestellt werden, wie z.B. in binärer Rerpräsentation, hexadezimaler Repräsentation, aber auch in menschenlesbarer<sup>c</sup> Repräsentation.<sup>d</sup>

Die Folge von Maschinenbefehlen, die ein üblicher Compiler generiert, ist üblicherweise in binärer Repräsentation, da die Maschinenbefehle in erster Linie für die Maschine, die binär arbeitet verständlich sein sollen und nicht für den Programmierer.

#### Anmerkung Q

Der PicoC-Compiler, der den Zweck erfüllt für Studenten ein Anschauungs- und Lernwerkzeug zu sein, generiert allerdings RETI-Code, der die RETI-Befehle in menschenlesbarer Repräsentation mit menschenlesbar ausgeschriebenen RETI-Operationen, RETI-Registern und Immediates (Definition 1.5) enthält. Für den RETI-Interpreter ist es ebenfalls nicht notwendig, dass der RETI-Code, den der PicoC-Compiler generiert, in binärer Repräsentation ist, denn es ist für den RETI-Interpreter ebenfalls leichter diesen einfach direkt in menschenlesbarer Repräsentation zu interpretieren. Der RETI-Interpreter soll nur die sichtbare Funktionsweise einer RETI-CPU simulieren und nicht deren mögliche interne Umsetzung<sup>a</sup>.

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

#### Definition 1.5: Immediate

Z

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

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

#### Definition 1.6: Cross-Compiler

Z

Kompiliert auf einer Maschine  $M_1$  ein Programm, dass in einer Wunschsprache  $L_w$  geschrieben ist für eine andere Maschine  $M_2$ , wobei beide Maschinen  $M_1$  und  $M_2$  unterschiedliche Maschinensprachen  $L_{M_1}$  und  $L_{M_2}$  haben.  $^{ab}$ 

Ein Cross-Compiler ist entweder notwendig, wenn eine Zielmaschine  $M_2$  nicht ausreichend Rechenleistung besitzt, um ein Programm in der Wunschsprache  $L_w$  selbst zeitnah zu kompilieren oder wenn noch keine Compiler  $C_w$  oder  $C_o$  für die Wunschsprache  $L_w$  oder eine andere Programmiersprache  $L_o$ , in welcher der

<sup>&</sup>lt;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 1.5) haben.

 $<sup>^</sup>b$ Bzw. Wörter.

 $<sup>^</sup>c\mathrm{So}$  wie die Programme des PicoC-Compilers dargestellt werden.

<sup>&</sup>lt;sup>d</sup>Scholl, "Betriebssysteme".

<sup>&</sup>lt;sup>a</sup>Beim PicoC-Compiler handelt es sich um einen Cross-Compiler  $C_{PicoC}^{Python}$ , der in der Sprache  $L_{Python}$  geschrieben ist und die Sprache  $L_{PicoC}$  kompiliert.

<sup>&</sup>lt;sup>b</sup>Earley und Sturgis, "A formalism for translator interactions".

Wunschcompiler  $C_w$  implementiert ist existieren, die unter der Zielmaschine  $M_2$  laufen.<sup>2</sup>

#### 1.1.1 T-Diagramme

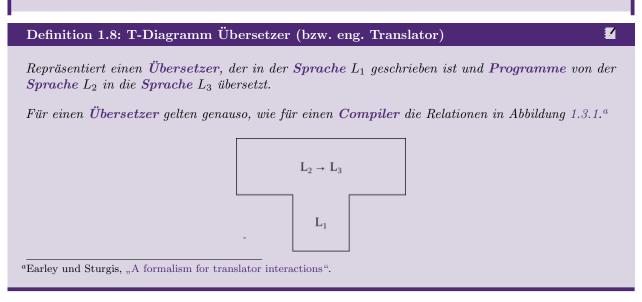
Um die Architektur von Compilern und Interpretern übersichtlich darzustellen eignen sich T-Diagramme, deren Spezifikation aus der Wissenschaftlichen Publikation Earley und Sturgis, "A formalism for translator interactions" entnommen ist besonders gut, da diese optimal darauf zugeschnitten sind die Eigenheiten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

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

# Definition 1.7: T-Diagramm Programm Repräsentiert ein Programm, dass in der Sprache $L_1$ geschrieben ist und die Funktion f berechnet. f $L_1$ \*Earley und Sturgis, "A formalism for translator interactions".

#### Anmerkung Q

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 1.7 also reichen einfach eine 1 hinzuschreiben.



 $<sup>^2</sup>$ Die an vielen Universitäten und Schulen eingesetzten programmierbaren Roboter von Lego Mindstorms nutzen z.B. einen Cross-Compiler, um für den programmierbaren Microcontroller eine zu  $L_C$  ähnliche Sprache in die Maschinensprache des Microcontrollers zu kompilieren, da es schneller geht ein Programm direkt auf der Maschine, auf der man programmiert zu kompilieren.

#### Anmerkung Q

Zwischen den Begriffen Übersetzung und Kompilierung gibt es einen Unterschied.

Übersetzung ist der allgemeinere Begriff und verlangt nur, dass Eingabe und Ausgabe des Übersetzers die gleiche Bedeutung<sup>a</sup> haben müssen, also die Relationen in Abbildung 1.3.1 erfüllt  $\operatorname{sind}^b$ .

Kompilierung beinhaltet dagegen meist auch das Lexen und Parsen oder irgendeine Form von Umwandlung eines Programmes von der Textrepräsentation in eine compilerinterne Datenstruktur und erst dann ein oder mehrere Übersetzungsschritte.

Kompilieren ist also auch Übersetzen, aber Übersetzen ist nicht immer auch Kompilieren.

#### **Definition 1.9: 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>

 $L_2$ 

 $L_1$ 

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

Die horinzontale Adjazenz lässt sich, wie man in Abbildung 1.1 erkennen kann zusammenfassen.

<sup>&</sup>lt;sup>a</sup>Auch Semantik (Definition 1.15) genannt.

<sup>&</sup>lt;sup>b</sup>Und ist auch zwischen Passes (Definition 1.42) möglich.

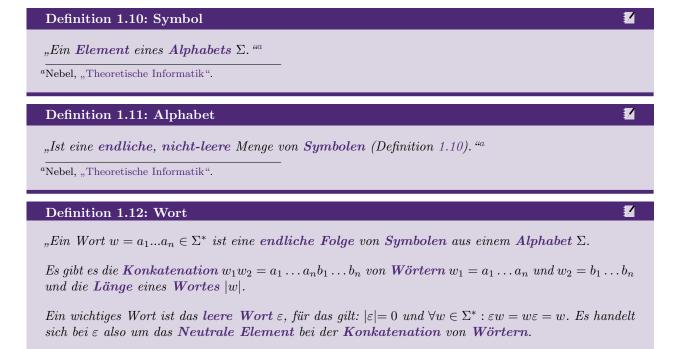
 $<sup>^</sup>a\mathrm{Earley}$  und Sturgis, "A formalism for translator interactions".



Abbildung 1.1: Horinzontale Übersetzungszwischenschritte 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 vorher eingeführt zu haben, was die Begriffe Symbol (Definition 1.10), Alphabet (Definition 1.11) und Wort (Definition 1.12) beinhaltet.



Bei  $\Sigma^*$  handelt es sich um Kleenesche Hülle eines Alphabets  $\Sigma$ , es ist die Sprache aller Wörter, welche durch beliebige Konkatenation von Symbolen aus dem Alphabet  $\Sigma$  gebildet werden können. Die Kleenesche Hülle ist die größte Sprache über  $\Sigma$  und jede Sprache über  $\Sigma$  ist eine Teilmenge davon. Es gilt des Weiteren:  $\varepsilon \in \Sigma^*$ . "a

<sup>a</sup>Nebel, "Theoretische Informatik".

#### Definition 1.13: Formale Sprache

"Menge von Wörtern (Definition 1.12) über dem Alphabet  $\Sigma$  (Definition 1.11). "

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 in verschiedenen Syntaxen (Definition 1.14) dargestellt sein.

#### Definition 1.14: Syntax



Bezeichnet alles was mit dem Aufbau von Wörtern einer Formalen Sprache zu tun hat. Eine Formale Grammatik oder in Natürlicher Sprache ausgedrückte Regeln können die Syntax einer Sprache beschreiben. Es kann auch mehrere verschiedene Syntaxen für die gleiche Sprache geben<sup>a</sup>.<sup>b</sup>

<sup>a</sup>Z.B. die Konkrete (Definition 1.32) und Abstrakte Syntax (Definition 1.39), die später eingeführt werden.

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

#### Definition 1.15: Semantik



Bezeichnet alles was mit der Bedeutung von Wörtern einer Formalen Sprache zu tun hat. a

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

#### Definition 1.16: Formale Grammatik



"Beschreibt, wie Wörter einer Formalen Sprache abgeleitet werden können.

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

- N = Nicht-Terminalsymbole.
- $\Sigma = Terminal symbole$ , wobei  $N \cap \Sigma = \emptyset$ .
- $P = Menge\ von\ Produktionsregeln\ w \to v$ , wobei  $w, v \in (N \cup \Sigma)^*\ und\ w \notin \Sigma^*$ .
- $S \triangleq Startsymbol$ , 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.

Es ist möglich zwei Grammatiken  $G_1$  und  $G_2$  in einer Vereinigungsgrammatik  $G_1 \uplus G_2 = \langle N_1 \cup N_2 \cup \{S\}, \Sigma, P_1 \cup P_2 \cup \{S ::= S_1 \mid S_2\}, S \rangle$  zu vereinigen. "e

"Des Weiteren gibt es die von einer Grammatik erzeugte Sprache:  $L(G) = \{w \in \Sigma^* \mid S \Rightarrow^* w\}$  "f, g

Formale Sprachen lassen sich in Klassen der Chromsky Hierarchie (Definition 1.17) einteilen.

#### Definition 1.17: Chromsky Hierarchie

Eine Hierarchie, in der Formale Sprachen nach der Komplexität ihrer Formalen Grammatiken in verschiedene Klassen unterteilt werden. Jede dieser Klassen hat verschiedene Eigenschaften, wie Entscheidungeprobleme, die in einer Klasse entscheidbar und in einer anderen unentscheidbar sind usw.

Eine Sprache  $L_i$  ist in der Chromsky Hierarchie in der Klasse  $i \in \{0, ..., 3\}$ , falls sie von einer Grammatik dieser Klasse i erzeugt<sup>a</sup> werden kann.

Zwischen den Sprachmengen benachbarter Klassen in Abbildung 1.17.1 besteht eine echte Teilmengenbeziehung:  $L_3 \subset L_2 \subset L_1 \subset L_0$ .

Formale Sprachen

Rekursiv Aufz. Sprachen (Klasse 0)

Kontextsensitive Sprachen (Klasse 1)

Kontextfreie Sprachen (Klasse 2)

Reguläre Sprachen (Klasse 3)

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

<sup>&</sup>lt;sup>a</sup>Weil mit ihnen terminiert wird.

<sup>&</sup>lt;sup>b</sup>Kann auch als **Alphabet** bezeichnet werden.

<sup>&</sup>lt;sup>c</sup>w muss mindestens ein Nicht-Terminalsymbol enthalten.

<sup>&</sup>lt;sup>d</sup>Bzw.  $w, v \in C^*$  und  $w \notin \Sigma^*$ .

<sup>&</sup>lt;sup>e</sup>Die Grammatik des PicoC-Compilers lässt sich in eine Grammatik für die Lexikalische Analyse  $G_{Lex}$  und eine für die Syntaktische Analyse  $G_{Parse}$  unterteilen. Die gesamte Grammatik des PicoC-Compilers steht allerdings vereinigt in einer Datei.

<sup>&</sup>lt;sup>f</sup>Die Nicht-Terminlsymbole in w fallen dabei weg, weil  $w \in \Sigma^*$ 

<sup>&</sup>lt;sup>g</sup>Nebel, "Theoretische Informatik".

<sup>&</sup>lt;sup>a</sup>Erzeugen meint hier, dass genau die Wörter der Sprache sich mit der Grammatik ableiten lassen, keines mehr oder weniger.

<sup>&</sup>lt;sup>b</sup>Z.B. ist jede Reguläre Sprache auch eine Kontextfreie Sprache, aber nicht jede Kontextfreie Sprache ist auch eine Reguläre Sprache.

<sup>&</sup>lt;sup>c</sup>Nebel, "Theoretische Informatik".

#### Definition 1.18: Reguläre Grammatik

Z

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

$$A \to cB, \qquad A \to c, \qquad A \to \varepsilon$$
 (1.18.1)

haben, wobei A, B Nicht-Terminalsymbole sind und c ein Terminalsymbol ist<sup>ab</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.
- $^b$ 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:

$$A \to v \tag{1.19.1}$$

 $haben,\ wobei\ A\ ein\ Nicht-Terminal symbol\ ist\ und\ v\ ein\ beliebige\ Folge\ von\ Grammatik symbolen^a$   $ist."^b$ 

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

Ob sich ein Programm überhaupt kompilieren lässt, entscheidet sich anhand des Wortproblems (Definition 1.20). In einem Compiler oder Interpreter ist das Wortproblem üblicherweise 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

Z

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

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

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

#### 1.2.1 Ableitungen

Jedes mit einen Compiler kompilierbare Programm kann mithilfe der Grammatik der Sprache des Compilers abgeleitet werden. Hierbei wird zwischen der 1-Schritt-Ableitungsrelation (Definition 1.21) und der normalen Ableitungsrelation (Definition 1.22) unterschieden.

#### Definition 1.21: 1-Schritt-Ableitungsrelation

1

"Eine binäre Relattion  $\Rightarrow$ , welche alle Anordnungen zweier Wörter  $w_1, w_2 \in (N \cup \Sigma)^*$  in Relation zueinander setzt, die sich nur durch das einmalige Anwenden einer Produktionsregel auf  $w_1$  voneinander

<sup>&</sup>lt;sup>b</sup>Nebel, "Theoretische Informatik".

 $<sup>^</sup>a\mathrm{Bzw.}$ "ja" oder "nein" usw., es muss nicht umgedingt 1 oder 0 sein.

<sup>&</sup>lt;sup>b</sup>Auch Charakteristische Funktion genannt.

<sup>&</sup>lt;sup>c</sup>Nebel, "Theoretische Informatik".

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  $w_1, w_2, x, y \in (N \cup \Sigma)^*$  "a

<sup>a</sup>Nebel, "Theoretische Informatik".

#### Definition 1.22: Ableitungsrelation

Z

"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 v aus  $(N \cup \Sigma)^*$ , welches durch beliebig häufiges Anwenden von Produktionsregeln auf ein Wort u entsteht.

Es gilt  $u \Rightarrow^* v$  genau dann wenn  $u = w_1 \Rightarrow \ldots \Rightarrow w_n = v$ , wobei  $n \geq 1$  und  $w_1, \ldots, w_n \in (N \cup \Sigma)^*$ . "a

<sup>a</sup>Nebel, "Theoretische Informatik".

Beim Ableiten kann auf verschiedene Weisen vorgegangen werden. Dasselbe **Programm** 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 und bei den **Linksableitungen** im Unterkapitel ?? relevant.

#### Definition 1.23: Links- und Rechtsableitungableitung



"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. "

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

 $^a\mathrm{Bei}$  Typ-1- und Typ-0-Grammatiken ist es statt einem Nicht-Terminalsymbol die linke Seite einer Produktion.

<sup>b</sup>Nebel, "Theoretische Informatik".

Ob eine Grammatik mehrdeutig (Definition 1.25) ist, kann durch Betrachtung Formaler Ableitungsbäume (Definition 1.24) festgestellt werden. Formale Ableitungsbäume werden im Unterkapitel 1.4 nochmal relevant, da in der Syntaktischen Analyse Formale Ableitungsbäume (Definition 1.34) als eine compilerinterne Datenstruktur umgesetzt werden.

#### Definition 1.24: Formaler Ableitungsbaum



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

Das Adjektiv "formal" kann dabei weggelassen werden, wenn der Kontext indem 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 compilerinterne Datentstruktur umsetzt.

Den Knoten dieses Baumes sind Grammatiksymbole  $C = N \cup \Sigma \cup \varepsilon$  (Definition 1.16) zugeordnet. Den Inneren Knoten des Baumes sind Nicht-Terminalsymbole N zugeordnet und den Blättern sind entweder Terminalsymbole  $\Sigma$  oder das leere Wort  $\varepsilon$  zugeordnet.

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

<sup>b</sup>Nebel, "Theoretische Informatik".

In Abbildung 1.24.2 ist ein Beispiel für einen Formalen Ableitungsbaum zu sehen, der sich aus der Ableitung 1.24.1 nach der im Dialekt der Erweiterten Backus-Naur-Form des Lark Parsing Toolkit (Definition ??) angegebenen Grammatik  $G = \langle N, \Sigma, P, add \rangle$  1.1 ergibt.

$\overline{NUM}$	::=	"4"   "2"
$ADD\_OP$	::=	"+"
$MUL\_OP$	::=	"*"
$\overline{mul}$	::=	mul MUL_OP NUM   NUM L_Parse
add	::=	$add\ ADD\_OP\ mul\  \ mul$

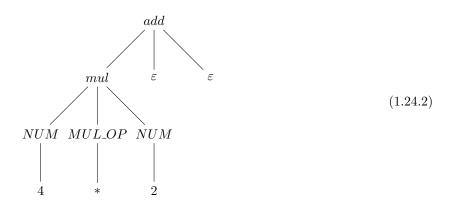
Grammatik 1.1: Grammatik für einen Ableitungsbaum in EBNF

#### Anmerkung 9

Werden die Produktionen einer Grammatik angegeben, wie in Grammatik ??, wird die Angabe dieser Produktionen auch oft als Grammatik bezeichnet, obwohl Grammatiken eigentlich durch ein Tupel  $G = \langle N, \Sigma, P, S \rangle$  dargestellt sind.

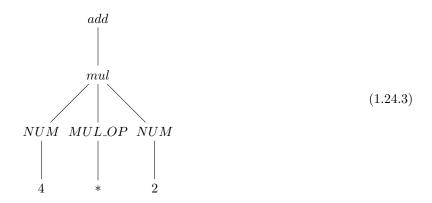
$$add \Rightarrow mul \Rightarrow mul \ MUL\_OP \ NUM \Rightarrow NUM \ MUL\_OP \ NUM \Rightarrow "4" "*" "2"$$
 (1.24.1)

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 1.24.2 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 1.24.3, nur die in den gewählten Produktionen vorhandenen Nicht-Terminalsymbole als Kinder hinzuzufügen<sup>3</sup>.

<sup>&</sup>lt;sup>3</sup>Diese Option wurde beim **PicoC-Compiler** gewählt.



#### 1.2.2 Präzedenz und Assoziativität

Für einen Compiler ist es notwendig, dass die Grammatik des Compilers keine Mehrdeutige Grammatik (Definition 1.25) ist, denn sonst können unter anderem die Assoziativität (Definition 1.26) und die Präzedenz (Definition 1.27) der verschiedenen Operatoren nicht gewährleistet werden, wie später in Unterkapitel ?? an einem Beispiel demonstriert wird. Ein Schema, um die Grammatiken zu definieren, die nicht mehrdeutig sind, wird in Unterkapitel ?? genauer erklärt.

#### Definition 1.25: Mehrdeutige Grammatik



"Eine Grammatik ist mehrdeutig, wenn es ein Wort  $w \in L(G)$  gibt, das mehrere Ableitungsbäume zulässt "ab."

<sup>a</sup>Für die Bedeutung von L(G), siehe Definition 1.16.

 $^b$ Alternativ geht auch die Definition: "Wenn es für w mehrere unterschiedliche Linksableitungen gibt".

<sup>c</sup>Nebel, "Theoretische Informatik".

#### Definition 1.26: Assoziativität



"Bestimmt, welcher Operator aus einer Reihe von Operatoren mit gleicher Präzedenz (Definition 1.27) zuerst ausgewertet wird."

Es wird 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>ab</sup>

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

<sup>b</sup>2.1.7 Vorrangregeln und Assoziativität.

Bei Assoziativität ist z.B. der Multitplikationsoperator \* ein Beispiel für einen linksassoziativen Operator und ein Zuweisungsoperator = ein Beispiel für einen rechtsassoziativen Operator. In Abbildung 1.2 ist ein Beispiel hierfür dargestellt, indem die resultierenden Auswertungsreihenfolgen mithilfe von Klammern () veranschaulicht sind.

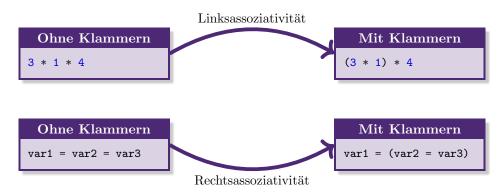


Abbildung 1.2: Veranschaulichung von Linksassoziativität und Rechtsassoziativität.



Die Mischung der Operatoren für Subtraktion '-' und für Multiplikation \*, welche beide eine unterschiedliche Präzedenz haben, ist ein Beispiel für den Einfluss von Präzedenz. In Abbildung 1.3 ist ein Beispiel hierfür dargestellt, indem mithilfe der Klammern () die resultierende Auswertungsreihenfolge veranschaulicht ist. Im Beispiel in Abbildung 1.3 ist durch die beiden Subtraktionsoperatoren '-' nacheinander und den darauffolgenden Multitplikationsoperator \*, sowohl Assoziativität als auch Präzedenz im Spiel.



Abbildung 1.3: Veranschaulichung von Präzedenz.

#### 1.3 Lexikalische Analyse

Die Lexikalische Analyse bildet üblicherweise den ersten Filter innerhalb des Pipe-Filter Architekturpatterns (Definition 1.1) bei der Implementierung von Compilern. Die Aufgabe der Lexikalischen Analyse ist in einem ersten Schritt in einem Eingabewort<sup>4</sup> Lexeme (Definition 1.28) zu finden, die mit einer Grammatik für die Lexikalische Analyse  $G_{Lex}$  abgeleitet werden können.



<sup>&</sup>lt;sup>4</sup>Z.B. dem Inhalt einer Datei, welche in UTF-8 kodiert ist.

Diese Lexeme werden vom Lexer (Definition 1.30) im Eingabewort identifziert und Tokens (Definition 1.29) zugeordnet. Die Tokens sind es, die letztendlich an die Syntaktische Analyse weitergegeben werden.

#### Definition 1.29: Token

Z

Ist ein Tupel (T, V) mit einem Tokentyp T und einem Tokenwert V. Ein Tokentyp T kann hierbei als ein Überbegriff für eine möglicherweise unendliche Menge verschiedener Tokenwerte V verstanden werden<sup>a</sup>.

<sup>a</sup>Z.B. gibt es im PicoC-Compiler viele verschiedene Tokenwerte, wie z.B. 42, 314 oder 12, welche alle unter dem Tokentyp NUM, für Zahl zusammengefasst sind.

#### Definition 1.30: Lexer (bzw. Scanner oder auch Tokenizer)

**Z** 

Ein Lexer ist eine Totale Funktion<sup>a</sup>: lex:  $\Sigma^* \to (T \times V)^*, w \mapsto (t_1, v_1) \dots (t_n, v_n)$ , welche ein Eingabewort<sup>b</sup>  $w \in \Sigma^*$  auf eine endliche Folge von Tokens  $(t_1, v_n) \dots (t_n, v_n) \in (T \times V)^*$  abbildet.

Die Definitionsmenge der Totalen Funktion lex erlaubt nur Wörter, die sich mit der Grammatik für die Lexikalische Analyse  $G_{Lex}$  ableiten lassen.

Ist das Abbilden eines Eingabeworts w auf eine Folge von Tokens  $(t_1, v_n) \dots (t_n, v_n)$  nicht möglich, da das Eingabewort Teilwörter enthält, die sich nicht mit der Grammatik für Lexikalische Analyse  $G_{Lex}$  ableiten lassen, so wird in diesem Fall eine Fehlermeldung (Definition 1.52) ausgegeben.

#### Anmerkung Q

Um Verwirrung vorzubeugen ist es wichtig die kontextabhängigen unterschiedliche Bedeutungen des Wortes Symbol hervorzuheben:

Wenn von Symbolen die Rede ist, so werden in der Lexikalischen Analyse, der Syntaktischen Analyse und der Code Generierung unterschiedliche Konzepte als Symbole bezeichnet.

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

In der Syntaktischen Analyse sind die Tokentypen die Symbole.

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

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

Eine weitere Aufgabe der Lexikalischen Analyse ist es jegliche für die Syntaktische Analyse unwichtigen Symbole, wie Leerzeichen  $_{-}$ , Newline  $\n^5$  und Tabs  $\t$  aus dem Eingabewort w herauszufiltern. Das geschieht im Lexer, indem dieser für alle unwichtigen Symbole bzw. Folgen von Symbolen<sup>6</sup> kein Token in der Folge von Tokens  $(t_1, v_n) \dots (t_n, v_n)$  vorsieht.

<sup>&</sup>lt;sup>a</sup>Alternativ könnte man die Funktion lex auch als Partielle Funktion definieren, aber das würde zum Ausdruck bringen, dass der Lexer bei einem Eingabewort, das sich nicht mit der Grammatik für die Lexikalische Analyse  $G_{Lex}$  ableiten lässt trotzdem durchchläuft. Die Funktion lex als Totale Funktion zu definieren drückt eher aus, dass ein Eingabewort, das nicht in der Definitionsmenge liegt zu einer Fehlermeldung führt.

<sup>&</sup>lt;sup>b</sup>Z.B. Quellcode eines Eingabeprogramms.

<sup>&</sup>lt;sup>c</sup>Thiemann, "Compilerbau".

<sup>&</sup>lt;sup>5</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 Porgrammiersprache, die man zur Inplementierung des Lexers nutzt wegabstrahiert.

<sup>&</sup>lt;sup>6</sup>Bzw. Teilwörter des Eingabeworts.

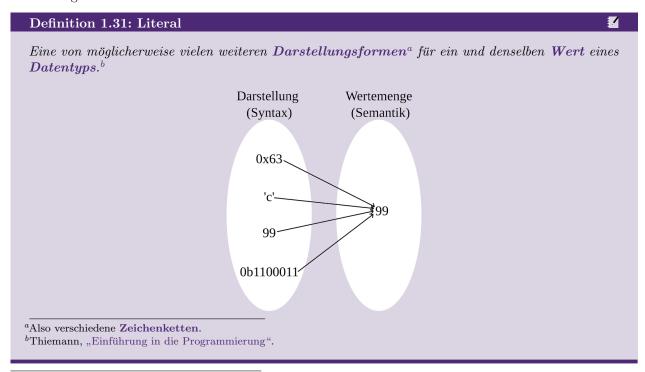
Der Grund, warum nicht einfach nur Lexeme an die Syntaktische Analyse weitergegeben werden und der Grund für die Aufteilung von Tokens in Tokentyp T und Tokenwert V, ist, weil z.B. die Bezeichner von Variablen, Konstanten und Funktionen und auch Zahlen beliebige Folgen von Symbolen sein können. Später in der Syntaktischen Analyse in Unterkapitel 1.4 ist es nur relevant, ob an einer bestimmten Stelle ein bestimmter Tokentyp T, z.B. eine Zahl NUM steht und der Tokenwert V ist erst wieder in der Code Generierung in Unterkapitel 1.5 relevant.

Wie es in Tabelle 1.1 zu sehen ist, gibt es für verschiedene Bezeichner, wie z.B. my\_fun, my\_var oder my\_const und verschiedene Zahlen, wie z.B. 42, 314 oder 12 passende Tokentypen NAME und NUM<sup>7 8</sup>, die einen Überbegriff darstellen. Für Lexeme, wie if oder } sind die Tokentypen dagegen genau die Bezeichnungen, die man diesen beiden Folgen von Symbolen geben würde, nämlich IF und RBRACE.

Lexeme	Token
42, 314	Token('NUM', '42'), Token('NUM', '314')
<pre>my_fun, my_var, my_const</pre>	Token('NAME', 'my_fun'), Token('NAME', 'my_var'), Token('NAME', 'my_const')
<b>if</b> , }	Token('IF', 'if'), Token('RBRACE', '}')
99, 'c'	Token('NUM', '99'), Token('CHAR', '99')

Tabelle 1.1: Beispiele für Lexeme und ihre entsprechenden Tokens.

Ein Lexeme ist nicht immer das gleiche, wie der Tokenwert V, denn wie in Tabelle 1.1 zu sehen ist, kann z.B. im Fall von  $L_{PicoC}$  der Tokenwert 99 durch zwei verschiedene Literale (Definition 1.31) dargestellt werden. Einmal kann der Tokenwert 99 als ASCII-Zeichen 'c' dargestellt werden, das dann als Tokenwert den entsprechenden Index in der ASCII-Tabelle erhält, nämlich 99 und des Weiteren auch in Dezimalschreibweise als 99<sup>9</sup>. Der Tokenwert ist der letztendlich verwendete Wert an sich, unabhängig von der Darstellungsform.



<sup>&</sup>lt;sup>7</sup>Bzw. wenn man sich nicht Kurzformen sucht, wären IDENTIFIER und NUMBER die passenden Bezeichnungen.

<sup>&</sup>lt;sup>8</sup>Die Bezeichnungen der Tokentypen wurden im PicoC-Compiler so gewählt, da man beim Programmieren möglichst kurze und leicht verständliche Bezeichner haben will, damit unter anderem auch mehr Code in eine Zeile passt.

 $<sup>^9</sup>$ Die Programmiersprache  $L_{Python}$  erlaubt es z.B. den Wert 99 auch mit den Literalen 0b1100011 und 0x63 darzustellen.

Die Grammatik für die Lexikalische Analyse  $G_{Lex}$  ist üblicherweise regulär, da ein typischer Lexer immer nur ein Symbol vorausschaut<sup>10</sup> und sich nichts merkt, also unabhängig davon, was für Symbole und wie oft bestimmte Symbole davor aufgetaucht sind funktioniert.

Der Grund, weshalb ein Lexer relativ unkompliziert zu implementieren ist, ist weil dieser nur Lexeme erkennen muss, welche durch eine Reguläre Grammatik beschrieben sind. Parser, welche im nächsten Unterkapitel 1.4 erklärt werden, sind dagegen deutlich komplexer zu implementieren, da diese bei den meisten Programmiersprachen eine Kontextfreie Grammatik umsetzen müssen<sup>11</sup>.

#### Anmerkung Q

Auch für den PicoC-Compiler lässt sich aus der im Dialekt der Backus-Naur-Form des Lark Parsing Toolkit (Definition ??) spezifizierten Grammatik für die Lexikalische Analyse ??  $G_{PicoC\_Lex}$  schlussfolgern, dass diese Grammatik eine Reguläre Grammatik ist, da alle ihre Produktionen die Definition 1.18 einer Regulären Grammatik erfüllen.

Produktionen mit Alternative, wie z.B.  $DIG\_WITH\_0 ::= "0" \mid DIG\_NO\_0$  in Grammatik ?? sind unproblematisch, denn sie können immer auch als  $\{DIG\_WITH\_0 ::= "0", DIG\_WITH\_0 ::= DIG\_NO\_0\}$  umgeschrieben werden und z.B.  $DIG\_WITH\_0*$ , (LETTER |  $DIG\_WITH\_0 \mid "\_")+$  und " $\_"."\sim"$  in Grammatik ?? können alle zu Alternativen umgeschrieben werden, wie es in Definition ?? gezeigt wird und diese Alternativen können wie gerade gezeigt umgeformt werden, um ebenfalls regulär zu sein. Somit existiert mit der Grammatik ?? eine Reguläre Grammatik, welche die Sprache  $L_{PicoC\_Lex}$  beschreibt und damit ist die Sprache  $L_{PicoC\_Lex}$  nach der Chromsky Hierarchie (Definition 1.17) regulär.

Um eine Gesamtübersicht über die Lexikalische Analyse zu geben, ist in Abbildung 1.4 die Lexikalische Analyse an einem Beispiel veranschaulicht.

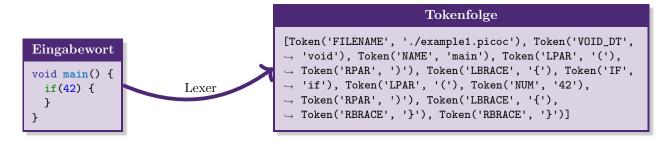


Abbildung 1.4: Veranschaulichung der Lexikalischen Analyse.

#### Anmerkung Q

Das Symbol  $\hookrightarrow$  zeigt in Codebeispielen einen Zeilenumbruch an, wenn eine Zeile zu lang ist.

#### 1.4 Syntaktische Analyse

In der Syntaktischen Analyse ist für einige Sprachen eine Kontextfreie Grammatik  $G_{Parse}$  notwendig, um diese Sprachen zu beschreiben, da viele Programmiersprachen z.B. für Funktionsaufrufe fun(arg) und Codeblöcke if (1) {} syntaktische Mittel verwenden, die es notwendig machen sich zu merken, wieviele

<sup>&</sup>lt;sup>10</sup>Man nennt das auch einem Lookahead von 1.

<sup>&</sup>lt;sup>11</sup>Hierzu kann z.B. der Earley Erkenner in Definition ?? im ?? als Beispiel genannt werden, der ein Bestandteil eines Earley Parsers (Definition ??) ist und bereits deutlich komplexer ist als ein typischer Lexer.

öffnende runde Klammern ( bzw. öffnende geschweifte Klammern { es zu einem bestimmten Zeitpunkt gibt, die noch nicht durch eine entsprechende schließende runde Klammer ) bzw. schließende geschweifte Klammer } geschlossen wurden. Diese syntaktischen Mittel lassen sich nicht mehr mit einer Regulären Grammatik (Definition 1.18) beschreiben, sondern es braucht eine Kontextfreie Grammatik (Definition 1.19) hierfür, die es erlaubt zwischen zwei Terminalsymbolen ein Nicht-Terminalsymbol abzuleiten.

#### Anmerkung Q

Für den PicoC-Compiler lässt sich aus der im Dialekt der Backus-Naur-Form des Lark Parsing Toolkit (Definition ??) spezifizierten Grammatik für die Syntaktische Analyse  $?? G_{PicoC\_Parse}$  schlussfolgern, dass diese eine Kontextfreie Grammatik, aber keine Reguläre Grammatik ist, da alle ihre Produktionen die Definition 1.19 einer Kontextfreien Grammatik erfüllen, aber nicht die Definition 1.18 einer Regulären Grammatik.

Es lässt sich sehr leicht erkennen, dass die Grammatik ?? eine Kontextfreie Grammatik ist, da alle Produktionen auf der linken Seite des ::=-Symbols immer nur ein einzelnes Nicht-Terminalsymbol haben und sich auf der rechten Seite eine beliebige Folge von Grammatiksymbolen<sup>a</sup> befindet.

Es lässt sich wiederum sehr einfach erkennen, dass die Grammatik ?? keine Reguläre Grammatik ist, denn z.B. bei der Produktion  $if\_stmt := "if""("logic\_or")" exec\_part$  ist das Nicht-Terminalsymbol  $logic\_or$  von den Terminalsymbolen für eine öffnende Klammer { und eine schließende Klammer } eingeschlossen, was mit einer Regulären Grammatik nicht ausgedrückt werden kann.

Somit existiert mit der Grammatik ?? eine Kontextfreie Grammatik, die allerdings keine Reguläre Grammatik ist, welche die Sprache  $L_{PicoC\_Parse}$  beschreibt. Hierdurch ist die Sprache  $L_{PicoC\_Parse}$  nach der Chromsky Hierarchie (Definition 1.17) kontextfrei, aber nicht regulär.

 $^a$ Also eine beliebige Folge von Nicht-Terminalsymbolen und Terminalsymbolen.

Die Syntax, in welcher ein Programm vor dem kompilieren in einer Textdatei aufgeschrieben ist, wird auch als Konkrete Syntax (Definition 1.32) bezeichnet. In einem Zwischenschritt, dem Parsen, wird aus diesem Programm mithilfe eines Parsers (Definition 1.35) ein Ableitungsbaum (Definition 1.34) generiert, der als Zwischenstufe hin zum einem Abstrakten Syntaxbaum (Definition 1.41) dient. Beim Compilerbau ist es förderlich kleinschrittig vorzugehen, deshalb erst die Generierung des Ableitungsbaums und dann erst des Abstrakten Syntaxbaums.

#### Definition 1.32: Konkrete Syntax

**Z** 

Bezeichnet alles, was mit dem Aufbau von Wörtern<sup>a</sup> zu tun hat, die nach einer Konkreten Grammatik  $G_{Lex} \uplus G_{Parse}$  (Definition 1.33) abgeleitet wurden.

Die Konkrete Syntax ist die Teilmenge der gesamten Syntax einer Sprache, welche für die Lexikalische und Syntaktische Analyse relevant ist. In der gesamten Syntax einer Sprache<sup>b</sup> kann es z.B. Wörter geben, welche die gesamte Syntax nicht einhalten, die allerdings korrekt nach einer Konkreten Grammatik abgeleitet sind<sup>c</sup>.

Ein Programm, wie es in einer Textdatei nach der Konkreten Grammatik<sup>d</sup> abgeleitet steht, bevor man es kompiliert, ist in Konkreter Syntax aufgeschrieben.<sup>e</sup>

<sup>&</sup>lt;sup>a</sup>Bzw. Programmen.

 $<sup>^</sup>b$ Vor allem bei Programmiersprachen.

<sup>&</sup>lt;sup>c</sup>Wenn ein Programm z.B. nicht deklarierte Variablen hat und aufgrund dessen nicht kompiliert werden kann, hält dieses die gesamte Syntax nicht ein, kann allerdings so nach einer Konkreten Grammatik abgeleitet werden. Solche Details werden üblicherweise nicht in eine Konkrete Grammatik mitaufgenommen.

<sup>&</sup>lt;sup>d</sup>Vereinigungsgrammatik, wie in Definition 1.16 erklärt.

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

Um einen kurzen Begriff für die Grammatik zu haben, welche die Konkrete Syntax einer Sprache beschreibt, wird diese im Folgenden als Konkrete Grammatik (Definition 1.33) bezeichnet.

#### Definition 1.33: Konkrete Grammatik

Z

Grammatik, welche die Konkrete Syntax einer Sprache beschreibt und die Grammatiken  $G_{Lex}$  und  $G_{Parse}$  miteinander vereinigt:  $G_{Lex} \uplus G_{Parse}^{\ a}$ .

<sup>a</sup>Vereinigungsgrammatik, wie in Definition 1.16 erklärt.

#### Definition 1.34: Ableitungsbaum (bzw. Konkreter Syntaxbaum, engl. Derivation Tree)

Compilerinterne Datenstruktur für den Formalen Ableitungsbaum (Definition 1.24) eines in Konkreter Syntax geschriebenen Programmes.

Die Blätter, die beim Formalen Ableitungsbaum das leere Wort  $\varepsilon$  oder Terminalsymbole  $\Sigma$  einer Konkreten Grammatik  $G = \langle N, \Sigma, P, S \rangle$  sein können, sind in dieser Datenstruktur Tokens (T, W) und die Inneren Knoten entsprechen weiterhin Nicht-Terminalsymbolen N einer Konkreten Grammatik  $G = \langle N, \Sigma, P, S \rangle$ . In dieser Datenstruktur wird allerdings nur der Teil eines Formalen Ableitungsbaumes dargestellt, der den Ableitungen einer Grammatik  $G_{Parse}$  entspricht. Die Tokens sind in der Syntaktischen Analyse ein atomarer Grundbaustein<sup>a</sup>, daher sind die Ableitungen der Grammatik  $G_{Lex}$  uninteressant.<sup>b</sup>

<sup>a</sup>Nicht mehr weiter teilbar.

 $^b JSON\ parser$  - Tutorial —  $Lark\ documentation$ .

Die Konkrete Grammatik nach der ein Ableitungsbaum konstruiert wird, ist optimalerweise immer so definiert, dass sich möglichst einfach aus dem Ableitungsbaum ein Abstrakter Syntaxbaum konstruieren lässt.

#### Definition 1.35: Parser



Ein Parser ist ein Programm, dass aus einem Eingabewort<sup>a</sup>, welches in Konkreter Syntax geschrieben ist eine compilerinterne Datenstruktur, den Ableitungsbaum generiert<sup>b</sup>. Dies wird auch als Parsen bezeichnet.<sup>c</sup>

<sup>a</sup>Z.B. ein **Programm** in einer **Textdatei**.

<sup>b</sup>Es gibt allerdings auch alternative Definitionen, denen nach ein Parser ein Programm ist, dass ein Eingabewort von Konkreter Syntax in Abstrakte Syntax übersetzt. Im Folgenden wird allerdings die hiesige Definition 1.35 verwendet. <sup>c</sup>JSON parser - Tutorial — Lark documentation.

#### Anmerkung Q

In Bezug auf Compilerbau ist ein Lexer ein Teil eines Parsers und ist auschließlich für die Lexikalische Analyse verantwortlich. In einem alltäglicheren Szenario entspricht lexen z.B. bei einem Wanderausflug dem Nachschlagen in einem Insektenlexikon und dem Aufschreiben, welcher Insektenart man in welcher Reihenfolge begegnet ist, mit jeweils einem Bild des konkreten Exemplars, dem man 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 gerade angefangene Beispiel zurückzukommen, dem Weiterverarbeiten der Beziehungen zwischen den Insektenbegegnungen unter Berücksichtigung des örtlichen Kontexts $^b$  in eine taugliche Form $^c$ .

In der Weiterverarbeitung könnte ein Interpreter die in eine taugliche Form gebrachten Daten interpretieren und daraus bestimmte Schlüsse ziehen und ein Compiler könnte die Daten vielleicht in eine für Menschen leichter verständliche Sprache kompilieren<sup>d</sup>.

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

 $^b\mathrm{Das}$ entspricht z.B. Semikolons ;, die vorhin bereits als Beispiel genannt wurden.

<sup>c</sup>Z.B. gibt es bestimmte Wechselbeziehungen zwischen Insekten, Insekten beinflussen sich gegenseitig und ihre Umwelt.

 $^d \mathrm{Normalerweise}$  kompiliert man in eine für die CPU verständliche Sprache.

Die vom Lexer aus dem Eingabewort generierten Tokens werden in der Syntaktischen Analyse vom Parser als Wegweiser verwendet, da je nachdem, in welchem Kontext bestimmte Tokens (T, V) mit einem bestimmten Tokentyp T auftauchen, dies einer anderen Ableitung in der Grammatik  $G_{Parse}$  entspricht, die zum Parsen eines Eingabeworts notwendig ist. Dabei wird in der Konkreten Grammatik  $G_{Parse}$  nach Tokentypen unterschieden und nicht nach Tokenwerten, 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 wieder relevant.

Ein Parser ist genauergesagt ein erweiterter Erkenner (Definition 1.36), denn ein Parser löst das Wortproblem (Definition 1.20) für die Sprache, in welcher das Programm, welches kompiliert werden soll geschrieben ist und konstruiert parallel dazu oder im Nachgang aus den Informationen, die während der Ausführung des Erkennungsalgorithmus<sup>12</sup> gesammelt wurden den Ableitungsbaum.

#### Definition 1.36: Erkenner (bzw. engl. Recognizer)



Entspricht in seiner Funktion einem Kellerautomaten<sup>a</sup>, in dem Wörter bestimmter Kontextfreier Sprachen erkannt werden.

Es ist ein Algorithmus, der erkennt, ob ein Eingabewort sich mit der Konkreten Grammatik einer Sprache ableiten lässt. Der Algorithmus überprüft also, ob das Eingabewort Teil der Sprache ist, die von der Konkreten Grammatik beschrieben wird oder nicht. Ein Erkenner löst folglich das Wortproblem (Definition 1.20).<sup>b</sup>

 $^a$ Automat mit dem Kontextfreie Grammatiken erkannt werden.

#### Anmerkung Q

Für das Parsen gibt es grundsätzlich drei geläufige Ansätze, die unterschieden werden:

• Top-Down Parsing: Der Algorithmus arbeitet von oben-nach-unten, also anschaulich von der Wurzel zu den Blättern eines im Nachhinein oder parallel dazu generierten Ableitungsbaums. Dementsprechend fängt der Algorithmus mit dem Startsymbol einer Konkreten Grammatik an und wendet in jedem Schritt eine Linksableitung auf die Nicht-Terminalsymbole an, bis man die Folge von Terminalsymbolen hat, die sich zum gewünschten Eingabewort abgeleitet hat oder sich herausstellt, dass dieses nicht abgeleitet werden kann.<sup>a</sup>

Der Grund, warum Linksableitungen verwendet werden und nicht z.B. Rechtsableitungen, ist, weil versucht wird das Eingabewort von links nach rechts mithilfe von Terminalsymbolen aus Ableitungen nachzubilden. Das passt gut damit zusammen, dass durch die Linksableitung

<sup>&</sup>lt;sup>b</sup>Thiemann, "Compilerbau".

<sup>&</sup>lt;sup>12</sup>Bzw. engl. recognition algorithm.

die Terminalsymbole von links-nach-rechts erzeugt werden.

Welche der Produktionen für ein Nicht-Terminalsymbol angewandt wird, wenn es mehrere Alternativen gibt, wird entweder durch Backtracking oder durch Vorausschauen entschieden. Eine sehr einfach zu implementierende Technik für Top-Down Parser ist hierbei der Rekursive Abstieg (Definition ??). Für diese Methode muss allerdings, wenn die Konkrete Grammatik eine Linksrekursive Grammatik (Definition ??) ist, diese umgeformt werden, um jegliche Linksrekursion aus dieser zu entfernen, da diese sonst zu unendlicher Rekursion führt.

Rekursiver Abstieg kann mit Backtracking verbunden werden, um auch Konkrete Grammatiken parsen zu können, die nicht LL(k) (Definition ??) sind. Dabei werden meist nach dem Prinzip der Tiefensuche alle Produktionen für ein Nicht-Terminalsymbol solange durchgegangen, bis das gewünschte Eingabewort abgeleitet ist oder alle Alternativen für das momentane Nicht-Terminalsymbol abgesucht sind. Das wird solange durchgeführt, bis man wieder beim Startsymbol angekommen ist und da auch alle Alternativen abgesucht sind. Dies bedeutet dann, dass das Eingabewort sich nicht mit der verwendeten Konkreten Grammatik ableiten lässt.<sup>b</sup>

Wenn man eine LL(k)-Grammatik hat, ist Backtracking nicht notwendig und es reicht einfach nur immer k Tokens im Eingabewort vorauszuschauen<sup>c</sup>. Mehrdeutige Grammatiken sind dadurch ausgeschlossen, weil LL(k) keine Mehrdeutigkeit zulässt.<sup>d</sup>

- Bottom-Up Parsing: Es wird mit dem Eingabewort gestartet und versucht Rechtsableitungen entsprechend der Produktionen einer Konkreten Grammatik rückwärts anzuwenden, bis man beim Startsymbol landet.<sup>e</sup>
- Chart Parsing: Es wird Dynamische Programmierung verwendet, indem partielle Zwischenergebnisse in einer Tabelle<sup>f</sup> gespeichert werden und wiederverwendet werden können. Das macht das Parsen Kontextfreier Grammatiken effizienter, sodass es nur noch polynomielle Zeit braucht, da Backtracking nicht mehr notwendig ist. Chart Parser können dabei top-down oder bottom-up Ansätze umsetzen<sup>g</sup> <sup>h</sup>.

Der Abstrakte Syntaxbaum wird mithilfe von Transformern (Definition 1.37) und Visitors (Definition 1.38) 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 1.39).

#### Definition 1.37: Transformer

Ein Programm, das von unten-nach-oben<sup>a</sup> nach dem Prinzip der Breitensuche alle Knoten des Ableitungsbaum besucht und beim Antreffen eines bestimmten Knoten des Ableitungsbaumes je nach Kontext einen entsprechenden Knoten des Abstrakten Syntaxbaumes erzeugt und diesen anstelle

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

 $<sup>^</sup>b$ Diese Methode des Parsens 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 Webseite Lark - a parsing toolkit for Python) ersetzt wurde.

 $<sup>^{</sup>c}$ Das wird auch als Lookahead von k bezeichnet.

<sup>&</sup>lt;sup>d</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>&</sup>lt;sup>e</sup> What is Bottom-up Parsing?

<sup>&</sup>lt;sup>f</sup>Bzw. einem Chart.

 $<sup>^</sup>g$ 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.

 $<sup>^</sup>h$ Der Earley Parser von Lark, welchen der PicoC-Compiler verwendet, fällt unter diese Kategorie.

des Knotens des Ableitungsbaumes setzt und so Stück für Stück den Abstrakten Syntaxbaum konstruiert.<sup>b</sup>

<sup>a</sup>In der Informatik wachsen Bäume von oben-nach-unten, von der Wurzel zur den Blättern.

#### Definition 1.38: Visitor

Ein Programm, das von unten-nach-oben<sup>a</sup>, nach dem Prinzip der Breitensuche alle Knoten des Ableitungsbaumes besucht und beim Antreffen eines bestimmten Knoten des Ableitungsbaumes, diesen in-place mit anderen Knoten tauscht oder manipuliert, um den Ableitungbaum für die weitere Verarbeitung durch z.B. einen Transformer zu vereinfachen. bc

<sup>a</sup>In der Informatik wachsen Bäume von oben-nach-unten, von der Wurzel zur den Blättern.

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

 $^c$  Transformers  $\, \mathcal{C} \,$  Visitors — Lark documentation.

#### Definition 1.39: Abstrakte Syntax



Steht für alles, was mit dem Aufbau von Abstrakten Syntaxbäumen zu tun hat.

Ein Abstrakter Syntaxbaum, der zur Kompilierung eines Wortes<sup>a</sup> generiert wurde befindet sich in Abstrakter Syntax.<sup>b</sup>

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

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

Um einen kurzen Begriff für die Grammatik, welche die Abstrakte Syntax einer Sprache beschreibt zu haben, wird diese im Folgenden als Abstrakte Grammatik (Definition 1.40) bezeichnet.

#### Definition 1.40: Abstrakte Grammatik



Grammatik, die eine Abstrakte Syntax beschreibt, also beschreibt was für Arten von Kompositionen mit den Knoten eines Abstrakten Syntaxbaumes möglich sind.

Jene Produktionen, die in der Konkreten Grammatik für die Umsetzung von Präzedenz notwendig waren, sind in der Abstrakten Grammatik abgeflacht. Dadurch sind die Kompositionen, welche die Knoten im Abstrakten Syntaxbaum bilden können syntaktisch meist näher an der Syntax von Maschinenbefehlen.

#### Definition 1.41: 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.

Die Knoten des Abstrakten Syntaxbaumes enthalten dabei verschiedene Attribute, welche wichtigen Informationen für den Kompiliervorang und Fehlermeldungen enthalten.<sup>a</sup>

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

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

#### Anmerkung Q

In dieser Bachelorarbeit wird häufig von der "Abstrakten Syntax", der "Abstrakten Grammatik" bzw. dem "Abstrakten Syntaxbaum" einer "Sprache" L gesprochen. Gemeint ist hier mit der Sprache L nicht die Sprache, welche durch die Abstrakte Grammatik, 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 Abstrakte Syntaxbaum überhaupt konstruiert wird. Für die tatsächliche Sprache, die durch die Abstrakte Grammatik beschrieben wird, interessiert man sich nie wirklich explizit. Diese Konvention wurde aus dem Buch G. Siek, Course Webpage for Compilers (P423, P523, E313, and E513) übernommen.

<sup>a</sup>Bzw. es ist die Sprache, welche durch die Konkrete Grammatik beschrieben wird.

Im Abstrakten Syntaxbaum können theoretisch auch die Tokens aus der Lexikalischen Analyse weiterverwendet werden, allerdings ist dies nicht empfehlenswert. Es ist zum empfehlen die Tokens durch eigene entsprechende Knoten umzusetzen, damit der Zugriff auf Knoten des Abstrakten Syntaxbaumes immer einheitlich erfolgen kann und auch, da manche Tokens des Abstrakten Syntaxbaum noch nicht optimal benannt sind. Manche "Symbole" werden in der Lexikalischen Analyse mehrfach verwendet, wie z.B. das Symbol - in  $L_{PicoC}$ , welches für die binäre Subtraktionsoperation als auch die unäre Minusoperation verwendet wurde. Der verwendete Tokentyp dieses Symbols lautet im PicoC-Compiler SUB\_MINUS. Da in der Syntaktischen Analyse beide Operationen nur in bestimmten Kontexten vorkommen, lassen sie sich unterscheiden und dementsprechend können für beide Operationen jeweils zwei seperate Knoten erstellt werden. Im Fall des PicoC-Compilers sind es die Knoten Sub() und Minus().

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

Die Abstrakte Grammatik nach der ein Abstrakter Syntaxbaum konstruiert ist wird optimalerweise immer so definiert, dass der Abstrakte Syntaxbaum in den darauffolgenden Verarbeitungsschritten<sup>13</sup> möglichst einfach weiterverarbeitet werden kann.

Auf der linken Seite in Abbildung ?? wird das Beispiel 1.24.2 aus Unterkapitel 1.2.1 fortgeführt. Dieses Beispiel stellt den Arithmetischen Ausdruck 4 \* 2 in Bezug auf die Konkrete Grammatik  $1.2^{14}$ , welche die höhere Präzedenz der Multipikation \* berücksichtigt in einem Ableitungsbaum dar. Allerdings handelt es sich bei diesem Ableitungsbaum nicht um einen Formalen Ableitungsbaum, sondern um eine compilerinterne Datenstruktur für einen solchen. Dementsprechend sind die Blätter nun Tokens, die mithilfe der Grammatik  $L_{Lex}$  generiert wurden, womit die Darstellung von Ableitungen sich auf die Grammatik  $L_{Parse}$  beschränkt.

Auf der rechten Seite in Abbildung ?? wird der Ableitungsbaum zu einem Abstrakten Syntaxbaum abstrahiert, der nach der Abstrakten Grammatik 1.3 konstruiert ist. Die Abstrakte Grammatik ist hierbei in Abstrakter Syntaxform (Definition ??) angegeben. In der Abstrakten Grammatik 1.3 sind jegliche Produktionen wegabstrahiert, die in der Konkreten Grammatik 1.2 so umgesetzt sind, damit diese Präzidenz beachtet und nicht mehrdeutig ist. Aus diesem Grund gibt es nur noch einen allgemeinen Knoten für binäre Operationen  $BinOp(\langle exp \rangle, \langle bin\_op \rangle, \langle exp \rangle)$ .

<sup>&</sup>lt;sup>13</sup>Den verschiedenen Passes.

<sup>&</sup>lt;sup>14</sup>Die Konkrette Grammatik ist hierbei im Dialekt der Erweiterter Backus-Naur-Form des Lark Parsing Toolkit (Definition ??) angegeben.

```
"4"
                        "2"
NUM
           ::=
                "+"
ADD\_OP
           ::=
           ::=
                "*"
MUL\_OP
                mul\ MUL\_OP\ NUM
                                         NUM
                                                 L\_Parse
mul
                add\ ADD\_OP\ mul
add
                                      mul
           ::=
```

Grammatik 1.2: Produktionen für Ableitungsbaum in EBNF

Grammatik 1.3: Produktionen für Abstrakten Syntaxbaum in ASF

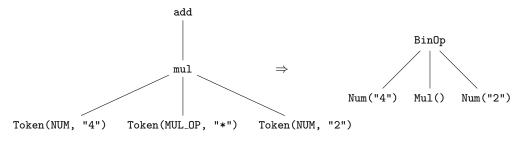


Abbildung 1.5: Veranschaulichung des Unterschieds zwischen Ableitungsbaum und Abstraktem Syntaxbaum.

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, sind in Abbildung 1.6 die einzelnen Zwischenschritte von den Tokens der Lexikalischen Analyse zum Abstrakten Syntaxbaum anhand des fortgeführten Beispiels aus Unterkapitel 1.3 veranschaulicht. In Abbildung 1.6 werden die Darstellungen des Ableitungsbaumes und des Abstrakten Syntaxbaumes verwendet, wie sie vom PicoC-Compiler ausgegeben werden. In der Darstellung des PicoC-Compilers stellen die verschiedenen Einrückungen die verschiedenen Ebenen dieser Bäume dar. Die Bäume wachsen von der Wurzel von links-nach-rechts zu den Blättern.

#### Abstrakter Syntaxbaum File Name './example1.ast', FunDef VoidType 'void', Tokenfolge Name 'main', [], [Token('FILENAME', './example1.picoc'), Token('VOID\_DT', Ε → 'void'), Token('NAME', 'main'), Token('LPAR', '('), Ιf → Token('RPAR', ')'), Token('LBRACE', '{'), Token('IF', Num '42', $_{\hookrightarrow}$ 'if'), Token('LPAR', '('), Token('NUM', '42'), → Token('RPAR', ')'), Token('LBRACE', '{'), ] → Token('RBRACE', '}'), Token('RBRACE', '}')] ] Parser Visitors und Transformer Ableitungsbaum file ./example1.dt decls\_defs decl\_def fun\_def type\_spec prim\_dt void pntr\_deg name main fun\_params decl\_exec\_stmts exec\_part exec\_direct\_stmt if\_stmt logic\_or logic\_and eq\_exp rel\_exp arith\_or arith\_oplus arith\_and arith\_prec2 arith\_prec1 un\_exp post\_exp 42 prim\_exp exec\_part compound\_stmt

Abbildung 1.6: Veranschaulichung der Syntaktischen Analyse.

#### 1.5 Code Generierung

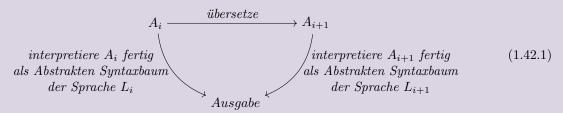
In der Code Generierung steht man nun dem Problem gegenüber einen Abstrakten Syntaxbaum einer Sprache  $L_1$  in den Abstrakten Syntaxbaum einer Sprache  $L_2$  umformen zu müssen. Dieses Problem lässt sich vereinfachen, indem man das Problem in mehrere Schritte unterteilt, die man Passes (Definition 1.42) nennt. So wie es auch schon mit dem Ableitungsbaum in der Syntaktischen Analyse gemacht wurde, den man als Zwischenstufe zum Abstrakten Syntaxbaum kontstruiert hatte. Aus dem Ableitungsbaum konnte dann unkompliziert und einfach mit Transformern und Visitors ein Abstrakter Syntaxbaum generiert werden.

#### Definition 1.42: 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.

Für jeden Pass und für einen beliebigen Abstrakten Syntaxbaum  $A_i$  gilt ähnlich, wie bei einem vollständigen Compiler in 1.42.1, 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.  $^{cd}$ 

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

#### 1.5.1 Monadische Normalform

Zum Verständnis dieses Kapitels sind die Begriffe **Ausdruck** (Definition 1.43) und **Anweisung** (Definition 1.44) wichtig.

#### Definition 1.43: Ausdruck (bzw. engl. Expression)



Code, der eine semantische Bedeutung hat und in einem bestimmten Kontext ausgewertet werden kann, um einen Wert zu liefern oder etwas zu deklarieren.

<sup>&</sup>lt;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 auschließlich auf Abstrakten Syntaxbäumen, ohne Parsing usw.

<sup>&</sup>lt;sup>b</sup>Der Begriff kommt aus dem Englischen von "passing over", da der gesamte Abstrakte Syntaxbaum in einem Pass durchlaufen wird.

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

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

 $Aufgebaut\ sind\ Ausdr\"{u}cke\ meist\ aus\ Kostanten,\ Variablen,\ Funktionsaufrufen,\ Operatoren\ usw.^{ab}$ 

<sup>a</sup>Ein Ausdruck ist z.B 21 \* 2;.

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

#### Definition 1.44: Anweisung (bzw. engl. Statement)

1

Code, der eine Vorschrifft darstellt, die ausgeführt werden soll und als ganzes keinen Wert liefert und nichts deklariert. Eine Anweisung kann sich jedoch aus ein oder mehreren Ausdrücken zusammensetzen, die dies tun.

Anweisungen sind zentrale Elemente Imperativer Programmiersprachen, die sich zu einem großen Teil aus Folgen von Anweisungen zusammensetzen.

In Maschinensprachen werden Anweisungen häufig als Befehle bezeichnet. ab

<sup>a</sup>Eine Anweisung ist z.B int var = 21 \* 2;.

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

Hat man es mit einer Programmiersprache zu tun, deren Programme Unreine Anweisungen (Definition 1.46) besitzen, so ist es sinnvoll einen Pass einzuführen, der Unreine Ausdrücke von den Anweisungen trennt, damit diese zu Reinen Anweisungen (Definition 1.45) werden. Das wird erreicht, indem man aus jedem Unreinen Ausdruck einen vorangestellten Ausdruck macht, den man vor die jeweilige Anweisung setzt, mit welcher der Unreine Ausdruck gemischt war. Der Unreine Ausdruck muss als erstes ausgeführt werden, für den Fall, dass der Effekt, den ein Unreiner Ausdruck hat die Reine Anweisung, mit der er gemischt war in irgendeinerweise beeinflussen könnte.

#### Definition 1.45: Reiner Ausdruck / Reine Anweisung (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 Maschinencode darstellen lässt. Sondern z.B. intern etwas am weiteren Kompiliervorgang ändert<sup>a</sup>.

Eine Reine Anweisung ist eine Anweisung, bei der alle Ausdrücke aus denen sich die Anweisung unter anderem zusammensetzt rein sind.<sup>b</sup>

<sup>a</sup>Z.B. ist die Allokation von Variablen **int var** kein Reiner Ausdruck. Eine Allokation bestimmt den Wert einiger Immediates im finalen Maschinencode, aber entspricht keiner Folge von Maschinenbefehlen.

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

#### Definition 1.46: Unreiner Ausdruck / Unreine Anweisung



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

Eine Unreine Anweisung ist eine Anweisung, bei der mindestens einer der Ausdrücke aus denen sich die Anweisung unter anderem zusammensetzt unrein ist.<sup>a</sup>

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

Auf diese Weise sind alle Anweisungen in Monadischer Normalform (Definiton 1.47).

#### Definition 1.47: Monadische Normalform (bzw. engl. monadic normal form)

7

Code ist in Monadischer Normalform, wenn dieser nach einer Grammatik in Monadischer Normalform abgeleitet wurde.

Eine Konkrete Grammatik ist in Monadischer Normalform, wenn alle ableitbaren Anweisungen rein sind. Oder sehr allgemein ausgedrückt, wenn Reines und Unreines klar voneinander getrennt ist.<sup>a</sup>

Eine Abstrakte Grammatik ist in Monadischer Normalform, wenn die Konkrete Grammatik 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 das Vorgehen, Code in die Monadische Normalform zu bringen, ist in Abbildung 1.7 zu sehen. Der Einfachheit halber wurde auf die Darstellung in Abstrakter Syntax verzichtet, welche allerdings zum großen Teil in dieser Schrifftlichen Ausarbeitung der Bachelorarbeit verwendet wird. Die Codebeispiele in Abbildung 1.7 sind daher in Konkreter Syntax<sup>15</sup> aufgeschrieben.

Links in der Abbildung 1.7 ist der Ausdruck mit dem Nebeneffekt, eine Variable zu definieren: int var, mit dem Ausdruck für eine Zuweisung exp = 5 % 4 gemischt: int var = 5 % 4. Der Unreine Definitionsausdruck int var muss daher vorangestellt werden, wie es rechts in Abbildung 1.7 dargestellt ist<sup>16</sup>.



Abbildung 1.7: Codebeispiel dafür Code in die Monadische Normalform zu bringen.

Die Aufgabe eines solchen Passes ist es, den Abstrakten Syntaxbaum der Syntax von Maschinenbefehlen anzunähren, indem Subbäume vorangestellt werden, die keine Entsprechung in Maschinenbefehlen haben. Somit wird eine Seperation von Subbäumen, die keine Entsprechung in Maschinenbefehlen haben und denen, die eine haben bewerkstelligt wird. Eine Reine Anweisung ist Maschinenbefehlen ähnlicher als eine Unreine Anweisung. Somit sparrt man sich in der Implementierung Fallunterscheidungen, indem Reine Ausdrücke und Reine Anweisungen direkt in Maschinenbefehle übersetzt werden können und nicht unterschieden werden muss, ob darin Unreine Ausdrücke vorkommen.

#### 1.5.2 A-Normalform

Zum Verständnis dieses Kapitels sind die Begriffe Ausdruck (Definition 1.43) und Anweisung (Definition 1.44) wichtig.

Eine Programmiersprache  $L_1$  soll in eine Maschinensprache  $L_2$  kompiliert werden. Im Falle dessen, dass es sich bei einer Sprache  $L_1$  um eine höhere Programmiersprache und bei  $L_2$  um eine Maschinensprache handelt, ist es fast unerlässlich einen Pass einzuführen, der Komplexe Ausdrücke (Definition 1.50) in

<sup>&</sup>lt;sup>15</sup>Für deren Kompilierung die Abstrakte Syntax überhaupt definiert wurde.

<sup>&</sup>lt;sup>16</sup>Obwohl hinter int var ein ; steht, ist es immer noch ein Ausdruck. Allerdings gibt es keine einheitliche Festlegung, was eine Anweisung ist und was nicht, es wurde für diese Schrifftliche Ausarbeitung der Bachelorarbeit nur so definiert.

Anweisungen und Ausdrücken verhindert. Das wird erreicht, indem man aus den Komplexen Ausdrücken vorangestellte Ausdrücke macht, in denen die Komplexen Ausdrücke temporären Locations (Definiton 1.48) zugewiesen werden und dann anstelle des Komplexen Ausdrucks auf die jeweilige temporäre Location zugegriffen wird.

#### Definition 1.48: Location

Z

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).

Sollte der Komplexe Ausdruck, welcher einer temporären Location zugewiesen wird, Teilausdrücke enthalten, die komplex sind, muss das gleiche Vorgehen erneut für die Teilausdrücke angewandt werden, bis alle Komplexen Ausdrücke nur noch Atomare Ausdrücke (Definiton 1.49) enthalten, falls sie sich überhaupt in weitere Teilausdrücke aufteilen lassen.

Sollte es sich bei dem Komplexen Ausdruck um einen Unreinen Ausdruck handeln, welcher nur einen Nebeneffekt ausführt und sich nicht in Maschinenbefehle übersetzen lässt, so wird aus diesem ein vorangestellter Ausdruck gemacht, welcher einfach nur den Nebeneffekt dieses Unreinen Ausdrucks ausführt und keiner temporären Location zugewiesen wird.

#### Definition 1.49: Atomarer Ausdruck



Ein Atomarer Ausdruck ist ein Reiner Ausdruck (Definition 1.45), der keinem kompletten Maschinenbefehl entspricht, sondern nur ein Argument, wie z.B. einen Immediate in einer Folge von Maschinenbefehlen festlegt.<sup>a</sup>

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

Bei einem üblichen Compiler, bei dem z.B. Register für temporäre Zwischenergebnisse genutzt werden und der Maschinenbefehlssatz es erlaubt zwei Register miteinander zu verechnen<sup>17</sup> sind Atomare Ausdrücke z.B. eine Variable (z.B. var), eine Zahl (z.B. 12) oder ein ASCII-Zeichen (z.B. 'c'), da diese häufig direkt über Register zugreifbar sind, die direkt mit einem Maschinenbefehl verechnet werden können<sup>18</sup> 19.

Im Fall des PicoC-Compilers ist ein Zugriff auf eine Location (z.B. stack(i)) der einzige Atomare Ausdruck, da der PicoC-Compiler so umgesetzt ist, dass er alle Zwischenergebnisse auf dem Stack speichert und dort dann auf diese zugreift, um sie in Register zu laden und miteinander zu verechnen<sup>20</sup>. Aus diesem Grund braucht es mindestens einen Maschinenbefehl<sup>21</sup>, um z.B. eine Zahl überhaupt für einen Maschinenbefehl zugreifbar zu machen, was der Definition 1.49 widerspricht. Daher sind z.B. Zahlen beim PicoC-Compiler keine Atomaren Ausdrücke.

 $<sup>^{17}{\</sup>rm Z.B.}$  Addieren oder Subtraktion von zwei Registerinhalten.

<sup>&</sup>lt;sup>18</sup>Mit dem RETI-Befehlssatz wäre das durchaus möglich, durch z.B. MULT ACC IN2.

<sup>&</sup>lt;sup>19</sup>Werden allerdings keine Register für Zwischenergebnisse genutzt werden, braucht man mehrere Maschinenbefehle, um die Zwischenergebnisse auf den Stack zu speichern und ein Stackpointer Register anzupassen.

<sup>&</sup>lt;sup>20</sup>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. Beim PicoC-Compiler sollte sich an die Paradigmen aus der Vorlesung Scholl, "Betriebssysteme" gehalten werden.

 $<sup>^{21}{\</sup>rm Genauerge sagt}$ 4.

#### Definition 1.50: Komplexer Ausdruck

/

Ein Komplexer Ausdruck ist ein Ausdruck, der nicht atomar ist.<sup>a</sup>

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

Im Fall des PicoC-Compilers sind Komplexe Ausdrücke z.B. 5 % 4, -1, fun(12) oder int var, da diese zur Berechnung auf jeden Fall mehrere Maschinenbefehle benötigen, was Definition 1.50 widerspricht oder unrein sind. Die Teilausdrücke 4, 5, 1 müssen erst auf den Stack geschrieben werden, um dann in Register geladen zu werden, damit dann der gesamte Komplexe Ausdruck berechnet werden kann. Die Ausdrücke fun(12) und int var sind unrein und daher Komplexe Ausdrücke.

In Abbildung 1.8 ist zur besseren Vorstellung die Einteilung von Komplexen, Atomaren, Unreinen und Reinen Ausdrücken veranschaulicht. Des Weiteren sind in der Abbildung alle Ausdrücke ausgegraut, welche die Monadische Normalform nicht erfüllen. Hierbei wird vom PicoC-Compiler ausgegangen, bei dem nur ein Zugriff auf eine Location (z.B. stack(i)) einen Atomaren Ausdruck darstellt.

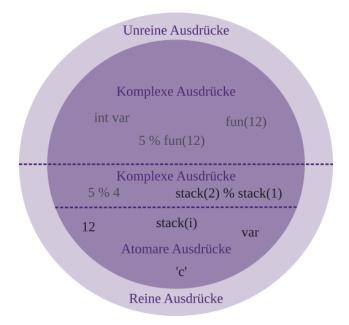


Abbildung 1.8: Übersicht über Komplexe, Atomare, Unreine und Reine Ausdrücke.

Es wird in dem gerade beschriebenen Pass dafür gesorgt, dass alle Anweisungen und Ausdrücke in A-Normalform<sup>22</sup> (Definition 1.51) sind. Wenn eine Konkrete Grammatik in A-Normalform ist, ist diese per Definition 1.51 auch automatisch in Monadischer Normalform, genauso, wie ein Atomarer Ausdruck nach Definition 1.49 auch ein Reiner Ausdruck ist.

#### Definition 1.51: A-Normalform (ANF)

Code ist in A-Normalform, wenn dieser nach einer Konkreten Grammatik in A-Normalform abgeleitet ist.

Eine Konkrete Grammatik ist in A-Normalform, wenn sie in Monadischer Normalform (Definition 1.47) ist und wenn alle Komplexen Ausdrücke nur Atomare Ausdrücke enthalten,

<sup>&</sup>lt;sup>22</sup>Das 'A' kommt vermutlich von "atomar" bzw. engl. "atomic", weil alle Komplexen Ausdrücke nur noch Atomare Ausdrücke enthalten dürfen.

falls sie sich überhaupt in weitere Teilausdrücke einteilen lassen.

Eine Abstrakte Grammatik ist in A-Normalform, wenn die Konkrete Grammatik für welche sie definiert wurde in A-Normalform ist. abc

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

Ein Beispiel für dieses Vorgehen, Code in die A-Normalform zu bringen, ist in Abbildung 1.9 zu sehen. Der Einfachheit halber wurde auf die Darstellung in Abstrakter Syntax verzichtet, welche allerdings zum großen Teil in dieser Schrifftlichen Ausarbeitung der Bachelorarbeit verwendet wird. Die Codebeispiele in Abbildung 1.7 sind daher in Konkreten Syntax<sup>23</sup> aufgeschrieben.

Um konsistent mit der Implementierung zu sein und später keine Verwirrung zu erzeugen, wird beim Beispiel in Abbildung 1.9 vom PicoC-Compiler ausgegangen, bei dem Variablen (z.B. var), Zahlen (z.B. 12) oder ASCII-Zeichen (z.B. 'c') Komplexe Ausdrücke darstellen.

Die Ausdrücke 4;, x;, usw. für sich sind in diesem Fall Komplexe Ausdrücke, deren Wert einer Location, in diesem Fall einer Speicherzelle des Stack zugewiesen werden. Auf das Ergebnis dieser Komplexen Ausdrücke wird mittels stack(2) und stack(1) zugegriffen, um diese in Register zu schreiben und dann z.B. im Komplexen Ausdruck stack(2) % stack(1) miteinander zu verrechnen und das Ergebnis wiederum einer Location, in diesem Fall ebenfalls einer Speicherzelle des Stack zuzuweisen. Dieses Ergebnis wird dann vom Stack stack(1) in die Globalen Statischen Daten global(0) gespeichert, wo die Variable x allokiert ist. Die Zahlen 0, 1, 2 sind dabei hierbei relative Adressen auf dem Stack und in den Globalen Statischen Daten.

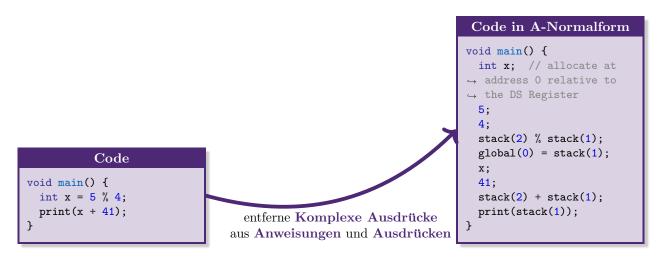


Abbildung 1.9: Codebeispiel für das Entfernen Komplexer Ausdrücke aus Operationen.

Ein Pass, wie er gerade beschrieben wurde hat vor allem in erster Linie die Aufgabe, den Abstrakten Syntaxbaum der Syntax von Maschinenbefehlen besonders dadurch anzunähren, dass er die Anweisungen weniger komplex macht und diese dadurch den ziemlich simplen Maschinenbefehlen syntaktisch ähnlicher sind. Des Weiteren vereinfacht dieser Pass die Implementierung der nachfolgenden Passes enorm, indem weniger Fallunterscheidungen nötig sind, da Anweisungen wie z.B. Zuweisungen nur noch die Form global(rel\_addr) = stack(1) haben, welche zudem viel einfacher verarbeitet werden kann.

<sup>&</sup>lt;sup>b</sup>Bolingbroke und Peyton Jones, "Types are calling conventions".

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

<sup>&</sup>lt;sup>23</sup>Für deren Kompilierung die Abstrakte Syntax überhaupt definiert wurde.

Alle weiteren denkbaren Passes sind zu spezifisch auf bestimmte Anweisungen 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 ?? erklärt.

#### 1.5.3 Ausgabe des Maschinencodes

Nachdem alle Passes durchgearbeitet wurden, ist es notwendig aus dem finalen Abstrakten Syntaxbaum den eigentlichen Maschinencode in Konkreter Syntax zu generieren. In üblichen Compilern wird hier für den Maschinencode eine binäre Repräsentation gewählt<sup>24</sup>. Der Weg von der Abstrakten Syntax zur Konkreten Syntax ist allerdings wesentlich einfacher, als der Weg von der Konkreten Syntax zur Abstrakten Syntax, für die eine gesamte Syntaktische Analyse, die eine Lexikalische Analyse beinhaltet durchlaufen werden muss.

Jeder Knoten des Abstrakten Syntaxbaumes 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; öffnenden- und schließenden runden Klammern (), öffnenden- und schließenden geschweiften Klammern {} usw. ausgibt. Dabei wird der gesamte Abstrakte Syntaxbaum durchlaufen und die Methode to\_string zur Ausgabe der Textrepräsentation der verschiedenen Knoten aufgerufen, die wiederum die Methode to\_string ihrer Kinder aufrufen und die zurückgegebene Textrepräsentation passend zusammenfügen und selbst zurückgebeben.

#### 1.6 Fehlermeldungen

Wenn bei einem Compiler ein unerwünschtes Verhalten der folgenden Kategorien<sup>25</sup> eintritt:

- 1. In der Lexikalischen oder Syntaktischen Analyse tritt ein Fall ein, der nicht in der Syntax der Sprache des Compilers abgedeckt ist, z.B.:
  - Der Lexer kann ein Lexeme nicht mit der Konkreten Grammatik für die Lexikalische Analyse  $G_{Lex}$  ableiten. Der Lexer ist genaugenommen ein Teil des Parsers und ist damit bereits durch den nachfolgenden Punkt "Parser" abgedeckt. Um die unterschiedlichen Ebenen, Lexikalische und Syntaktische Analyse gesondert zu betrachten wurde der Lexer an dieser Stelle ebenfalls kurz eingebracht.
  - Der Parser<sup>26</sup> entscheidet das Wortproblem (Definition 1.20) für ein Eingabeprogramm<sup>27</sup> mit 0, also das Eingabeprogramm lässt sich nicht durch die Konkrete Grammatik  $G_{Lex} \uplus G_{Parse}$  des Compilers ableiten.
- 2. In den Passes tritt ein Fall ein, der nicht in der Syntax der Sprache des Compilers abgedeckt ist, z.B.:
  - Eine Variable wird verwendet, obwohl sie noch nicht deklariert ist.
  - Bei einem Funktionsaufruf werden mehr oder weniger Argumente, Argumente des falschen Datentyps oder Argumente in der falschen Reihenfolge übergeben, als sie im Funktionsprototyp angegeben sind.
- 3. Während der Laufzeit des Compilers tritt ein Ereignis ein, das nicht durch die Semantik der Sprache des Compilers abgedeckt ist oder welches das Betriebssystem nicht erlaubt, z.B.:

<sup>&</sup>lt;sup>24</sup>Da der PicoC-Compiler vor allem zu Lernzwecken konzipiert ist, wird bei diesem der Maschinencode allerdings in einer menschenlesbaren Repräsentation ausgegeben

 $<sup>^{25}</sup>Errors\ in\ C/C++$  - Geeks for Geeks.

 $<sup>^{26}\</sup>mathrm{Bzw.}$  der **Erkenner** innerhalb des Parsers.

<sup>&</sup>lt;sup>27</sup>Bzw. ein **Wort**.

- 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 wenn während 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 oder mehr als einer anderen Objektdatei bereitsgestellt.

wird eine Fehlermeldung (Definition 1.52) ausgegeben.

#### Definition 1.52: Fehlermeldung



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

- über stdout oder stderr im einem Terminal Emulator oder richtigen Terminal.
- über eine Dialogbox in einer Graphischen Benutzeroberfläche<sup>b</sup> oder Zeichenorientierten Benutzerschnittstelle<sup>c</sup>.
- über ein Register oder eine spezielle Adresse des Hauptspeichers mithilfe eines Wertes.
- über eine Logdateid auf einem Speichermedium.

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

<sup>&</sup>lt;sup>b</sup>In engl. Graphical User Interface, kurz GUI.

<sup>&</sup>lt;sup>c</sup>In engl. Text-based User Interface, kurz TUI.

 $<sup>^</sup>d$ In engl. log file.

#### Literatur

#### Online

- 2.1.7 Vorrangregeln und Assoziativität. URL: https://www.tu-chemnitz.de/urz/archiv/kursunterlagen/C/kap2/vorrang.htm (besucht am 05.09.2022).
- A-Normalization: Why and How (with code). URL: https://matt.might.net/articles/a-normalization/(besucht am 23.07.2022).
- Errors in C/C++ GeeksforGeeks. URL: https://www.geeksforgeeks.org/errors-in-cc/ (besucht am 10.05.2022).
- JSON parser Tutorial Lark documentation. URL: 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).
- Transformers & Visitors Lark documentation. URL: https://lark-parser.readthedocs.io/en/latest/visitors.html (besucht am 09.07.2022).
- 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).
- Nystrom, Robert. Parsing Expressions · Crafting Interpreters. Genever Benning, 2021. URL: https://www.craftinginterpreters.com/parsing-expressions.html (besucht am 09.07.2022).

#### Artikel

• Earley, J. und Howard E. Sturgis. "A formalism for translator interactions". In: *CACM* (1970). DOI: 10.1145/355598.362740.

#### Vorlesungen

- Nebel, Bernhard. "Theoretische Informatik". Vorlesung. Vorlesung. Universität Freiburg, 2020. URL: http://gki.informatik.uni-freiburg.de/teaching/ss20/info3/index\_de.html (besucht am 09.07.2022).
- Scholl, Christoph. "Betriebssysteme". Vorlesung. Vorlesung. Universität Freiburg, 2020. URL: https://abs.informatik.uni-freiburg.de/src/teach\_main.php?id=157 (besucht am 09.07.2022).
- Thiemann, 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/swtvl (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. 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).