

Maturaarbeit

Compiler Construction

Fabio Stalder

Betreut durch
Thomas Jampen

5. August 2024



Gymnasium Kirchenfeld
Abteilung MN

Inhaltsverzeichnis

| | | |
|----------|-----------------------------|----------|
| 1 | Was ist ein Compiler | 3 |
| 1.1 | Scanning | 3 |
| 1.2 | Parsing | 4 |
| 1.3 | Semantic Analysis | 5 |
| 1.4 | Code Generation | 5 |
| 1.5 | Optimization | 5 |
| 2 | Meine Idee | 7 |

1 Was ist ein Compiler

In der Informatik beschreibt Compiler ein Programm, das Code aus einer Programmiersprache in eine andere übersetzt. In dieser Hinsicht gleichen Compiler Übersetzern für Menschensprache. Jedoch unterscheidet sich ein Compiler grundsätzlich von Übersetzern in der Erwartungshaltung, die an sie gestellt wird. Menschensprache ist sehr komplex und [...]

Ein Compiler ist traditionell nach folgendem Schema aufgebaut.

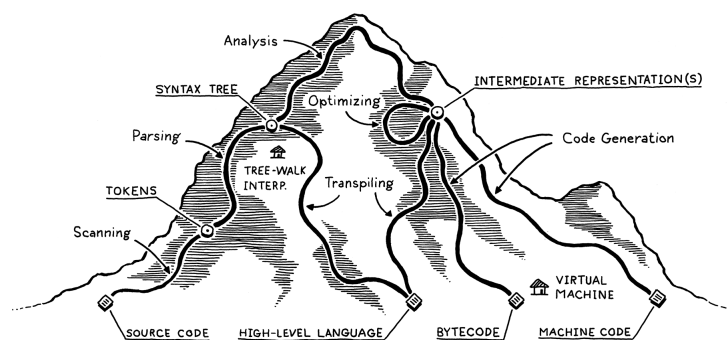


Abbildung 1.1: Schritte, die ein Compiler durchläuft [1]

In diesem **TEXT** werde ich mich nur auf die im unteren Schema dargestellten Schritte fokussieren.

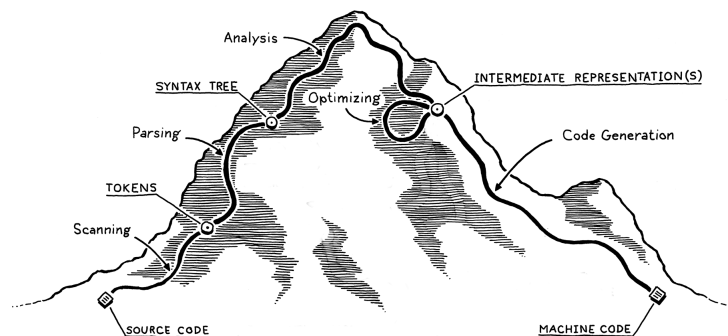


Abbildung 1.2: Schritte, die in diesem **TEXT** behandelt werden (Basierend auf Figure 1.1)

1.1 Scanning

Meist werden Programme so geschrieben, dass wir Menschen es lesen und verstehen können. Dafür verwendet man Buchstaben und Zahlen, Zeichen, wie +, *, oder Klammern, und Whitespaces, wie Leerzeichen oder Absätze. Diese Zeichen sind jedoch für den Computer unverständlich. Der erste Schritt beim compilieren ist daher das Scannen. Dies wird von einem Teil des Compilers, dem Lexer, durchgeführt. Die

Aufgabe dieses Lexers ist es den Input File zu scannen und die gescannten Zeichen in sogenannte Tokens zu verwandeln. Diese Tokens sind Datenstrukturen, die der Compiler kennt und mit denen er weiterarbeiten kann.

Als Beispiel:

Listing 1.1: C code vor dem Scannen

```
int foo()
{
    if (bar == 0)
    {
        return 0;
    }

    return 1;
}
```

Würde hierbei zu einem Array von Token Objekten umgewandelt werden:

Listing 1.2: Tokens nach dem Scannen

```
KeywordToken (keyword="int ")
IdentifierToken (id="foo ")
LParenthesisToken
RParenthesisToken
KeywordToken (keyword="if ")
LParenthesisToken
IdentifierToken (id="bar ")
OperatorToken (operator=ComparisonEqual)
LiteralIntToken (value=0)
[...]
```

Der Lexer legt hierbei fest welche Zeichen die Input-Programmiersprache enthalten darf und welche Bedeutung ihnen zugesprochen wird. So ist zum Beispiel im Lexer festgelegt, dass ein `+` Zeichen als Addition interpretiert wird. Genauso wie im Listing 1.2 `'if'` als `KeywordToken` gescannt wird, lässt sich im Lexer auch bestimmen, dass ein Wort wie `'print'` als `Keyword` angesehen werden soll.

1.2 Parsing

Nun versteht der Compiler was mit den Zeichen im Input File gemeint ist, jedoch fehlt noch etwas bis tatsächlich in eine andere Programmiersprache übersetzt werden kann. Und das ist Verständnis für Syntax. Die meisten High-Level Programmiersprachen weisen Syntaxregeln auf. Diese beinhalten, wie Funktionen und Variablen definiert werden oder mit welchen Punktvorstrich-Regeln Expressions evaluiert werden. Die beim Scanning gefundenen Tokens werden nun ineinander verschachtelt und in einen sogenannten Abstract Syntax Tree (AST) überführt.

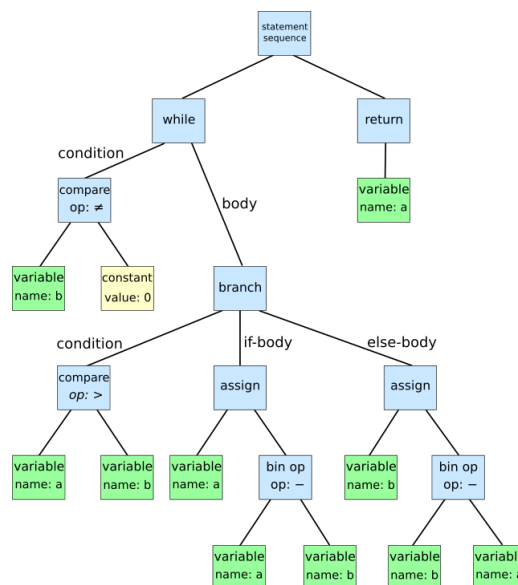


Abbildung 1.3: Abstract Syntax Tree zum Euklidischen Algorithmus

Ein AST enthält somit nicht nur Informationen über die Tokens, sondern über die gesamte Struktur die sich aus den Tokens ergibt. Variabel- und Funktionsdefinitionen oder komplexe Statements wie 'if' oder 'for' sind hierbei im AST enthalten. [...]

1.3 Semantic Analysis

Wie auch bei den meisten Menschensprachen gibt es auch für Programmiersprachen eine Semantik und diese muss natürlich vom Compiler verstanden werden. [...]

1.4 Code Generation

Code Generation ist der finale und oft auch komplexeste Schritt, der ein Compiler ausführen muss. Nun da unser Input-Code nicht mehr nur als Textfile, sondern als Intermediate Representation vorliegt, kann endlich Output-Code generiert werden. Jedoch lässt sich über diesen Schritt fast am wenigsten sagen, da er je nach Output-Sprache sehr unterschiedlich aussehen kann.

1.5 Optimization

Code Generation ist zwar der letzte Schritt beim Compilieren, trotzdem wurde eine wichtige Aufgabe des Compilers noch nicht betrachtet. Optimization ist ein Sprache die zwischen jedem der genannten Schritte geschieht. Dabei geht es darum den Output-Code so effizient wie möglich zu machen. Effizient kann hierbei jedoch viel Verschiedenes bedeuten. Der Output-Code muss so schnell wie möglich ausgeführt werden können, Memory sparsam verwenden und am besten auch noch ein kleiner File sein. Optimization

reicht vom Entfernen der Kommentare beim Scannen oder umstellen von mathematischen Operationen bis zu entfernen von ungebrauchten Variablen und Deadstores. Es muss von CPU Registern profitiert, mit Heap-Memory umgegangen und von inline Funktionen Gebrauch gemacht werden. Compiler Optimization ist somit ein sehr vielseitiges Problem, dass hierbei nicht weiter thematisiert werden sollte.

2 Meine Idee

Wie im vorherigen Abschnitt gezeigt, ist ein Compiler ein äusserts komplexes Programm, mit vielen verschiedenen Schritten. Jedoch ist die zugrundeliegende Aufgabe gar nicht so kompliziert. Man braucht ja nur, ein Dokument mit Text der bestimmten Regeln folgt, in Text mit anderen Regeln verwandeln. Natürlich ist dies etwas salopp ausgedrückt, trotzdem fragte ich mich, ob es nicht möglich sei einen viel einfacheren Compiler zu schreiben

Abbildungsverzeichnis

| | | |
|-----|---|---|
| 1.1 | Schritte, die ein Compiler durchläuft [1] | 3 |
| 1.2 | Schritte, die in diesem TEXT behandelt werden (Basierend auf Figure 1.1) | 3 |
| 1.3 | Abstract Syntax Tree zum Euklidischen Algorithmus | 5 |

Literaturverzeichnis

[1] Bob Nystrom. A map of territory (mountain.png), 2021. [Online; accessed 2024-08-05].