



# Cutting Languages Down to Size

## Student Project

at the Cooperative State University Baden-Württemberg Stuttgart

by

**Nahku Saidy and Hanna Siegfried**

08.06.2020

**Time of Project**

**Student ID; Course**

**Advisors**

nothing

8540946, 6430174; TINF17ITA

Prof. Dr. Stephan Schulz and Geoff Sutcliffe

# Contents

<b>Acronyms</b>	<b>I</b>
<b>List of Figures</b>	<b>II</b>
<b>List of Tables</b>	<b>III</b>
<b>Listings</b>	<b>IV</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Problem Statement and Goals . . . . .	1
1.2 Structure of the Report . . . . .	2
<b>2 Background and Theory</b>	<b>3</b>
2.1 TPTP Language . . . . .	3
2.2 Compiler . . . . .	3
2.3 Lexer . . . . .	3
2.4 Parser . . . . .	5
2.5 Python? . . . . .	5
<b>3 Concept</b>	<b>6</b>
3.1 Overview . . . . .	6
3.2 Lexer . . . . .	6
3.3 Parser . . . . .	6
3.4 Generation of the Reduced Grammar . . . . .	7
3.5 Selection of blocked Productions . . . . .	7
3.6 Determination of the remaining reachable Productions . . . . .	7
3.7 Determination of the remaining terminating Productions . . . . .	7
<b>4 Implementation</b>	<b>8</b>
4.1 Lexer . . . . .	8
4.2 Parser . . . . .	9
4.3 GUI . . . . .	9
<b>5 Validation</b>	<b>10</b>
<b>6 Conclusion</b>	<b>11</b>
<b>Bibliography</b>	<b>i</b>

# Acronyms

<b>ATP</b>	automated theorem proving
<b>BNF</b>	Backus-Naur form
<b>CFG</b>	context-free grammar
<b>CNF</b>	first-order clause normal form
<b>EBNF</b>	extended Backus-Naur form
<b>FOF</b>	full first-order logic
<b>PLY</b>	Python Lex-Yacc
<b>TFF</b>	typed first-order logic
<b>THF</b>	typed higher-order logic
<b>TPTP</b>	Thousands of Problems for Theorem Provers

# List of Figures

# List of Tables

# Listings

# 1 Introduction

## 1.1 Problem Statement and Goals

Formal languages are likely to grow over time as they are getting more complex when their functionality is extended and more application cases are covered. On the one hand that leads to a more powerful language. However, on the other hand it becomes harder to understand the language and to implement it. Thus, it becomes harder for new users to use the language.

This problem can be addressed by dividing languages into smaller sub-languages that cover everything relevant to the specific use case. This could be done manually, but using this method is likely to raise errors or divergences from the original grammar.

Therefore, the approach considered in this report is to develop an application that is able to automatically extract sub-languages from a language. A sub-language should be specified by the user using the application.

This report focusses on the Thousands of Problems for Theorem Provers ([TPTP](#)) language for automated theorem proving. Sub-languages of interest are for example a grammar just for first-order clause normal form ([CNF](#)) or full first-order logic ([FOF](#)). The grammar of the language is provided in an extended Backus-Naur form ([EBNF](#)). The first step to divide the [TPTP](#) language in smaller sub-languages is to build a parser that parses the grammar of the [TPTP](#) language. The parser should build a parse tree that represents the grammar rules of the [TPTP](#) language. This parse tree should be visually presented to the user and the user can then choose which grammar rules should not be included in the desired sub-language. After the user specified the sub-language, the developed application should extract the sub-language from the [TPTP](#) language and present the sub-language in the same format as the original [TPTP](#) syntax. Also, comments present in the [TPTP](#) syntax should be maintained and associated with the corresponding rules in the reduced syntax.

## 1.2 Structure of the Report

This report is structured into six chapters. In the first chapter the research problem and goals of this research are stated. Then in chapter 2 the necessary background information for the following chapters 3 and 4 is provided. In chapter 3 the concept for the TODO APLICATION is developed. Based on this, the implementation of the application is featured in chapter 4. In chapter 5 the results of the TODO reduced grammar is tested on a problem which is presented in a form corresponding to the reduced grammar. Chapter 6 sums up the results achieved in this research and offers an outlook for possible future research.



## 2 Background and Theory

[Mogensen.2017] Compiler: Translate (high-level) programming language into machine language

Different phases for writing a compiler, phases are processed sequentially

### 2.1 TPTP Language

The Thousands of Problems for Theorem Provers ([TPTP](#)) is a library of problems for automated theorem proving ([ATP](#)). Problems within the library are described in the [TPTP](#) language. The [TPTP](#) language is a formal language and its grammar is specified in an [EBNF](#). [Sut17]

The [EBNF](#) [[EBNF](#)] is a standard and is used to describe context-free grammars (CFGs).

### 2.2 Compiler

[Mogensen.2017] Compiler: Translate (high-level) programming language into machine language

Different phases for writing a compiler, phases are processed sequentially

### 2.3 Lexer

Lexing -> Input, Output -> Final Automata -> Language, Regular Expression (Shortcuts) -> Lexer Generator

Lexing or a so called lexical analysis is the division of input into units so called tokens [LexYacc.1992]. Tokens are for example variable names or keywords. The

input is a string containing a sequence of characters that thus the output is a sequence of tokens. The output can now be used for further processing e.g. 2.4.1.

Consequently, a lexer takes a set of characters and tries to match them with a token. This can lead to reading the same sequence of characters multiple times until the matching token has been identified.

Due to the given complexity (TOO), the lexer is often generated by a lexer generator and not written manually. +shorter

A lexer generator takes a specification of tokens as input and generates the lexer automatically.

Lexer - finite automata

### Regular Expression

The specification of tokens is usually written using regular expressions. A regular expression describes a set of strings. This set of strings can be characterized as a formal language. A formal language describes a set of words belonging to the language. These words are built over the alphabet of the language, and can be described by a finite automata.

Shorthands are common to simplify a regular expression. For example all alphabetic letters in lower and upper case are combined and represented by `[a-zA-Z]`. The same principle can also be applied to represent a set of numbers. However, using not clearly defined intervals e.g. `[0-b]` is not common as it has different interpretations by different lexer generators and thus can lead to mistakes. [Mogensen.2017]

A lexer needs to distinguish different types of tokens and furthermore decide which token to use if there are multiple ones that fit the input. [Mogensen.2017]

A simple approach to build a lexer is to build an automata for each token definition and then test to which automata the input corresponds. However, this would be slow as all automatas need to be passed through in the worst case. Therefore, it is convenient to build a single automata that tests each token simultaneously. This automata can be built by combining all regular expressions by disjunction. Each final state from each regular expression is marked to know which token has been identified.

It is possible, that final states overlap as a consequence of one token being a subset of another token. For solving such conflicts a precedence of tokens can be declared. Usually the token that is being defined the earliest has a higher precedence and thus will be chosen if multiple tokens fit the input. [Mogensen.2017]

Another task of the lexer is separating the input in order to divide it into tokens. Per convention the longest input that matches any token is chosen. [Mogensen.2017]

## 2.4 Parser

### 2.4.1 Yacc

Building a syntax tree out of the generated tokens [Mogensen.2017]

Parsing: establish relationship among tokens [LexYacc.1992] Grammar: list of rules that defines the relationships [LexYacc.1992]

Input: description of grammar [LexYacc.1992] Output: parser [LexYacc.1992]

### 2.4.2 PLY

Python Lex-Yacc (PLY) [PLY] is an implementation of lex and yacc in python. [LALR-parsing] consists of lex.py and yacc.py

lex.py tokenizes an input string

### 2.4.3 Nondeterministic Finite Automata

## 2.5 Python?

# 3 Concept

This chapter outlines the

## 3.1 Overview

## 3.2 Lexer

-deviations from plain ebnf

## 3.3 Parser

### 3.3.1 Data Structure

to store the [TPTP](#) Grammar, nested

### **3.4 Generation of the Reduced Grammar**

### **3.5 Selection of blocked Productions**

### **3.6 Determination of the remaining reachable Productions**

### **3.7 Determination of the remaining terminating Productions**

bei tree building temporäres startsymbol nutzen (da mehrere Startsymbole möglich)

# 4 Implementation

## 4.1 Lexer

-Definition of tokens

-Tabs and newlines ignored

-Newline would be helpful to identify comments because a comment is a newline followed by the percentage sign, as well as new rules if each rule would be represented in one line. However, there are rules that cover multiple lines. That is the main reason newlines are ignored.

-A comment is identified by the lexer as a percentage sign followed by an arbitrary character excluding "]". This is followed by any arbitrary character. A comment can not only be identified by a percentage sign as the percentage sign is also part of the terminal symbols. However, the percentage symbol when used as terminal symbol is embedded in square brackets.

Tokens: LGRAMMAR/TOKEN/STRICT/MACRO EXPRESSION:

Any arbitrary symbol that is the name of the rule followed by the symbol itself  
(:==,:::,...)

Non terminal symbol:

A non terminal symbol starts with "<" and ends with ">". In between there is any arbitrary sequence of numbers, underscores and small or capital letters.

T SYMBOL:

COMMENT:

OPEN SQUARE BRACKET/CLOSE SQUARE BRACKET, OPEN/CLOSE PARENTHESIS, ALTERNATIVE SYMBOL, REPETITION SYMBOL:

-is recognized and represented by the symbol itself

test

## 4.2 Parser

The parser is taking the tokens from the lexer and matches them to defined production rules.

## 4.3 GUI

# 5 Validation

back to back testing show advantages and useful for tptp users...



# 6 Conclusion

Outlook

# Bibliography