



# 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, XXXXX; 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 Backus-Naur form (BNF) . . . . .	3
2.3 Parser . . . . .	3
<b>3 Concept</b>	<b>5</b>
<b>4 Implementation</b>	<b>6</b>
<b>5 Validation</b>	<b>7</b>
<b>6 Conclusion</b>	<b>8</b>
<b>Bibliography</b>	<b>i</b>

# Acronyms

<b>BNF</b>	Backus-Naur form
<b>CNF</b>	First order clause normal form
<b>FOF</b>	Full first order logic
<b>TFF</b>	Typed first order logic
<b>THF</b>	Typed high 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

Programming languages and other formal languages (e.g. data representation languages like JSON) are often described via a context-free formal grammar, e.g. in the Extended Backus–Naur form that can be automatically processed by tools like Yacc and Bison. These languages tend to grow over time, as more features are added and the language covers more application cases. On the one hand, this makes the language more powerful. But on the other hand, this often makes the language harder to understand and harder to implement, increasing the barrier to entry for new developers.

This problem can, of course, be addressed by manually maintaining simplified or partial grammars. However, this creates not only a maintenance overhead, but may also lead to undesired and even unnoticed divergences of the full and the simplified grammar. Instead, we propose to use automatic extraction of self-contained parts of a formal grammar to create simpler sub-languages.

Languages are likely to grow over time as the language is getting more complex to extend its functionality. On

- languages grow + more functionality - harder to understand the whole language and to implement it - hard for new users to apply - therefore, divide language into smaller sub-language that cover everything that is relevant to the specific use case
- could be made manually, but that is likely to raise errors or divergences from the original grammar - our approach is to develop an application that is able to automatically subtract a sub-language from a language

Our particular use case is the TPTP syntax for automated theorem proving tools. TPTP defines a family of languages, from pure first-order clause normal form (CNF) and full first order (FOF) via typed first-order (TFF), eventually to typed higher-order logic (THF). The language can be used to write both input problems (i.e. lists of logical formulae), but also derivations (where some of the formulas are

input formulas, and the others are justified by reference to existing formulas and some mechanism for deriving the former from the latter. Most of the extensions are modular and even conservative (i.e. FOF is mostly a superset of CNF, and normally FOF is used as shorthand for CNF+FOF), and we are interested in e.g. extracting just a grammar for CNF or CNF+FOF. We are also interested extracting languages that do not allow certain features, e.g. a language in which only the pure input format is specified, not logical derivations. The task of this thesis is the design of a tool that performs this extraction and presents the result in a format that is both machine readable, but also friendly to use for human users.

- in specific, the considered language is the tptp language for automated theorem proving - language in bnf

One approach to this would be to a) select a given non-terminal node as the start symbol for the new grammar and b) block the undesired productions starting at certain non-terminal symbols. The tool can then start at the selected new start symbol and follow all legal transition to collect the remaining reachable grammar. It can also clean-up the resulting grammar by inlining rules for non-terminals with only a single production. An outstanding solution would also heuristically associate comments with grammar productions, and keep the comments referring to the remaining productions in place.

The approach is to first build a parser that parses the grammar of the Thousands of Problems for Theorem Provers (**TPTP**) language that is provided in an extended **BNF**. 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 subtract 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



## 2 Background and Theory

### 2.1 TPTP Language

[1]

### 2.2 Backus-Naur form (BNF)

### 2.3 Parser

#### 2.3.1 Lex

Lexing/lexical analysis: Division of input into units so called tokens [2]

Input: description of tokens - lex specification, regular expressions][2] Output: routine that identifies those tokens [2]

#### 2.3.2 Yacc

Parsing: establish relationship among tokens [2] Grammar: list of rules that defines the relationships [2]

Input: description of grammar [2] Output: parser [2]

### 2.3.3 PLY

Python implementation of lex and yacc [LALR-parsing] consists of lex.py and yacc.py

lex.py tokenizes an input string

<http://www.dabeaz.com/ply/ply.html>

## 3 Concept

## 4 Implementation

## 5 Validation

# 6 Conclusion

Outlook

# Bibliography

## Publikationen

- [1] G. Sutcliffe. “The TPTP Problem Library and Associated Infrastructure. From CNF to TH0, TPTP v6.4.0”. In: *Journal of Automated Reasoning* 59.4 (2017), pp. 483–502.
- [2] John Levine, Tony Mason, and Doug Brown. *Lex & Yacc*. O’Reilly Media Inc., 1992. ISBN: 9781565920002.