



# Cutting Languages Down to Size

## Student Project

at the Cooperative State University Baden-Württemberg Stuttgart

by

**Nahku Saidy and Hanna Siegfried**

06/08/2020

**Time of Project**  
**Student ID; Course**  
**Advisors**

07/10/2019 - 06/08/2020  
8540946, 6430174; TINF17ITA  
Prof. Dr. Stephan Schulz and  
Prof. Geoffrey Sutcliffe, Ph.D.

# **Erklärung**

Wir versichern hiermit, dass wir die Studienarbeit mit dem Thema: *Cutting Languages Down to Size* selbstständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt haben.

Wir sind uns bewusst, dass eine falsche Erklärung rechtliche Folgen haben wird.

Stuttgart, 06/08/2020

---

Nahku Saidy und Hanna Siegfried

Die Software wurde mittels Paarprogrammierung entwickelt.

Die schriftliche Ausarbeitung der Arbeit kann wie folgt aufgeteilt werden:

Kapitel	Nahku Saidy	Hanna Siegfried	Zusammen
1	•	•	gesamt
2	•	•	gesamt
3	3.7 - 3.9, 3.12 - 3.13	3.4 - 3.6, 3.10 - 3.11	3.1 - 3.3
4	4.4 - 4.6, 4.7, 4.9	4.1 - 4.3, 4.6, 4.8	•
5	•	•	gesamt
6	•	•	gesamt

## **Abstract**

Computer languages are likely to grow over time as they are getting more complex when their functionality is extended. An example for that is the TPTP language for automated theorem proving. Over time various forms of classical logics ranging from Clause Normal Form (CNF) to Typed First-order Form (TFF) have been included in and extended the TPTP language. This research describes a tool called Synplifier (Syntax simplifier) that automatically extracts sub-languages from the TPTP language. Automatic extraction instead of manually maintaining sub-languages has the advantage of avoiding maintenance overhead as well as unnoticed divergences from the full language. Sub-languages of interest are for example CNF or First-order Form (FOF) and are extracted based on the user's selection which part of the language to maintain. Synplifier has been successfully tested by extracting CNF from the TPTP language.

# Contents

<b>Acronyms</b>	<b>I</b>
<b>List of Figures</b>	<b>II</b>
<b>List of Tables</b>	<b>IV</b>
<b>Listings</b>	<b>V</b>
<b>List of Algorithms</b>	<b>VI</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 Formal languages . . . . .	3
2.2 Backus-Naur Form ( <a href="#">BNF</a> ) . . . . .	7
2.3 TPTP language . . . . .	8
2.4 Lexing . . . . .	8
2.5 Parsing . . . . .	9
2.6 Lex and Yacc . . . . .	11
2.7 Python . . . . .	12
<b>3 Concept</b>	<b>15</b>
3.1 Approach . . . . .	15
3.2 Requirements . . . . .	15
3.3 Overview . . . . .	18
3.4 Lexer . . . . .	20
3.5 Parser . . . . .	22
3.6 Graph generation . . . . .	27
3.7 Control file . . . . .	29
3.8 Maintaining comments . . . . .	30
3.9 Extraction of a sub-syntax . . . . .	33
3.10 Output generation . . . . .	35
3.11 GUI . . . . .	38
3.12 Command-line interface . . . . .	41
3.13 Counting the number of rules in the TPTP syntax . . . . .	42

<b>4 Implementation</b>	<b>44</b>
4.1 Lexer . . . . .	44
4.2 Parser . . . . .	47
4.3 Graph generation . . . . .	53
4.4 Maintaining comments . . . . .	59
4.5 Extraction of a sub-syntax . . . . .	66
4.6 Input . . . . .	73
4.7 Output generation . . . . .	74
4.8 GUI . . . . .	82
4.9 Command-line interface . . . . .	96
<b>5 Validation</b>	<b>97</b>
5.1 Automated parser generation . . . . .	97
5.2 Syntax size comparison . . . . .	99
<b>6 Conclusion</b>	<b>101</b>
6.1 Future Work . . . . .	101

## Bibliography

i

# Acronyms

<b>ASCII</b>	American Standard Code for Information Interchange
<b>ATP</b>	Automated Theorem Proving
<b>BNF</b>	Backus-Naur Form
<b>CFG</b>	Context-free grammar
<b>CNF</b>	Clause Normal Form
<b>EBNF</b>	Extended Backus-Naur Form
<b>FOF</b>	First-order Form
<b>PLY</b>	Python Lex-Yacc
<b>Synplifier</b>	Syntax simplifier
<b>TFF</b>	Typed First-order Form
<b>THF</b>	Typed Higher-order Form
<b>TPTP</b>	Thousands of Problems for Theorem Provers
<b>URL</b>	Uniform Resource Locator

# List of Figures

2.1	Finite automaton . . . . .	5
3.1	Procedure of extracting a sublanguage . . . . .	18
3.2	UML diagram of the architecture of Syntax simplifier ( <a href="#">Synplifier</a> ) . .	20
3.3	Parsing procedure . . . . .	23
3.4	Maintaining comments flow chart . . . . .	33
3.5	Procedure of generating a string representation of the grammar graph	36
3.6	Import menu . . . . .	38
3.7	Main GUI view . . . . .	39
3.8	View menu . . . . .	39
3.9	Reduce menu . . . . .	40
3.10	Save menu . . . . .	41
3.11	Control file menu . . . . .	41
4.1	Rule data types UML class diagram . . . . .	49
4.2	Comment block data type UML class diagram . . . . .	49
4.3	Production element data type UML class diagram . . . . .	50
4.4	Production property data type UML class diagram . . . . .	50
4.5	Production data type UML class diagram . . . . .	50
4.6	Productions list data type UML class diagram . . . . .	51
4.7	Grammar list data type UML class diagram . . . . .	51
4.8	Parsing example . . . . .	52
4.9	Graphbuilder UML diagram . . . . .	54
4.10	NTNode UML diagram . . . . .	55
4.11	RuleType UML diagram . . . . .	55
4.12	Split comment block by top of page line flow chart . . . . .	65
4.13	Check if node value is in terminating set . . . . .	70
4.14	Save ordered rules from grammar graph . . . . .	76
4.15	View UML class diagram . . . . .	82
4.16	Reduce and save Thousands of Problems for Theorem Provers ( <a href="#">TPTP</a> ) syntax . . . . .	89



# List of Tables

3.1	TPTP language rule types [12]	21
3.2	Command-line interface parameters	42
5.1	TPTP syntax size comparison	100

# Listings

3.1	TPTP rule example . . . . .	22
3.2	Productions of the nonterminal symbol <i>&lt;formula_role&gt;</i> . . . . .	28
3.3	Control file example . . . . .	29
3.4	Example of a comment in the TPTP syntax . . . . .	30
3.5	Comment lines split by a <i>Top of Page</i> line in the TPTP syntax . .	31
3.6	Comment syntax definition in the TPTP syntax . . . . .	37
3.7	Making the comment syntax reachable . . . . .	37
4.1	Lexer tokens . . . . .	44
4.2	Lexer regular expressions . . . . .	45
4.3	Example of rules defined over multiple lines . . . . .	46
4.4	Lexer ignore function . . . . .	46
4.5	Lexer error function . . . . .	47
4.6	Productions list parser function . . . . .	47
4.7	Grammar expression . . . . .	49
4.8	Production element . . . . .	52
4.9	Multiple rules with the same left-hand side of the rule . . . . .	56
4.10	Import TPTP syntax from the TPTP website . . . . .	74
4.11	Read text from file . . . . .	74
4.12	Print list entry . . . . .	74
4.13	Production alternatives in the TPTP syntax . . . . .	79
4.14	Implementation of menu bar . . . . .	83
4.15	Implementation of menu actions . . . . .	83
4.16	Init Tree View . . . . .	84
4.17	Reduce TPTP syntax with selection . . . . .	87
4.18	Argparse command-line parser configuration . . . . .	96
5.1	Control file to extract Clause Normal Form (CNF) . . . . .	98
5.2	<i>&lt;annotated_formula&gt;</i> production rule . . . . .	98
5.3	CNF parser main-function . . . . .	99

# List of Algorithms

1	Build nodes dictionary . . . . .	56
2	Build graph recursively . . . . .	57
3	Extract productions from productions list . . . . .	58
4	Append children to node . . . . .	58
5	Assign comments to rules . . . . .	61
6	Append comments to node . . . . .	62
7	Extend node comment block . . . . .	63
8	Removing blocked productions . . . . .	67
9	Removing nonterminating symbols . . . . .	68
10	Find nonterminating symbols . . . . .	69
11	Delete productions containing nonterminating symbols . . . . .	72
12	Delete nonterminating nodes . . . . .	72
13	Create print list . . . . .	77
14	<i>NTNode</i> string creation . . . . .	78
15	<i>PRODUCTION_LIST</i> string creation . . . . .	79
16	<i>PRODUCTION</i> string creation . . . . .	80
17	<i>PRODUCTION_ELEMENT</i> string creation . . . . .	81
18	Init tree view . . . . .	84
19	Toggle comments . . . . .	86
20	Load controlfile . . . . .	93
21	Produce controlfile . . . . .	95

# 1 Introduction

## 1.1 Problem statement and goals

Computer languages are likely to grow over time as they are getting more complex when their functionality is extended and more use cases are covered. On the one hand that leads to a more powerful language capable of handling a wide range of use cases. On the other hand increased complexity makes a language harder to learn and to use. Especially new users are discouraged to implement tools in that language.

One example of a language that has been expanding is the [TPTP](#) language for automated theorem proving. Over time various forms of classical logics ranging from [CNF](#) to Typed First-order Form ([TFF](#)) have been included in and extended the [TPTP](#) language.

This report describes a tool called [Synplifier](#) that is able to automatically extract sub-languages from the [TPTP](#) language. Sub-languages of interest are for example [CNF](#) or First-order Form ([FOF](#)) and are specified by the user using the application. The goal is maintaining the expressiveness of the whole [TPTP](#) language but allowing users to extract a sub-syntax to simplify the language for their particular use case. [Synplifier](#) processes a given grammar of a language in multiple steps. First it parses the formal grammar into a structured internal representation using Python Lex-Yacc ([PLY](#)). The processed grammar is presented to the user via a GUI. The user can select a start symbol and disable productions that should not be included in the desired sub-syntax. Using the users input, the developed application extracts the sub-syntax from the [TPTP](#) syntax and presents the sub-syntax in the same format as the original [TPTP](#) syntax. Also, comments present in the [TPTP](#) syntax are maintained and associated with the corresponding rules in the reduced syntax.

The concept of sub-syntax extraction is not only limited to the [TPTP](#) language and can also be applied to any other computer language, e.g., the SMT-LIB language [1]. To our knowledge, no other tool that extracts sub-languages from a formal language exists.

## **1.2 Structure of the Report**

The first chapter introduces the problem of complex computer languages and the goal of this report that is extracting smaller sub-languages. The second chapter provides necessary background information including the [TPTP](#) language, formal grammars, lexing and parsing. The third chapter outlines the concept of [Synplifier](#) utilising the information from chapter 2. Based on this, the implementation of [Synplifier](#) is featured in the fourth chapter. The fifth chapter presents an evaluation of the effectiveness of [Synplifier](#). Considering the evaluation, the sixth chapter sums up the results, compares them to the defined goals in the first chapter and offers an outlook for possible future research.

# 2 Background and Theory

This chapter introduces the technologies and background that will be utilized in the following chapters. First, an introduction into formal grammars is given and the Backus-Naur Form ([BNF](#)) is described. Following that, the [TPTP](#) language is introduced. Then the foundations of lexing and parsing are outlined. After that, Python and relevant Python packages that are used in the implementation of [Synplifier](#) are presented.

## 2.1 Formal languages

This section introduces terms and concepts in the area of formal languages that are necessary to understand [Synplifier](#).

### 2.1.1 Alphabet

An alphabet is a finite, non-empty set of symbols usually represented by the uppercase sigma  $\Sigma$ . An example is the binary alphabet  $\Sigma = \{0, 1\}$ . [2]

### 2.1.2 String

A string, also called word, is a finite sequence of symbols from some alphabet. For example, the string *101* is a string from the binary alphabet  $\Sigma = \{0, 1\}$ . The set of all strings over an alphabet  $\Sigma$  is denoted as  $\Sigma^*$ . [2]

### 2.1.3 Language

If  $\Sigma$  is an alphabet and  $L$  is a subset of  $\Sigma^*$ , then  $L$  is a language over  $\Sigma$  [2]. If sigma is the American Standard Code for Information Interchange ([ASCII](#)) alphabet  $\Sigma = \{a - z, A - Z, 0 - 9, \ , !, =, \#, \%, \$, \&, ', (,), *, +, , -, ., /, :, ;, <, >, =\}$

, ?, @, [ , ], \, ^, \_\_, ‘{, }, |, ~}, then for example the language where words can consist of  $a - zA - Z$  is one language over the **ASCII** alphabet.

In the example the abbreviations  $a - zA - Z$  and 0–9 are used. The first abbreviation represents all alphabetic letters in lower and upper case. The second abbreviation represents all one-digit numbers.

## 2.1.4 Finite automata

A finite automaton can recognize words of a language. The language represented by an automaton is the set of all accepted words. An automaton consists of a set of states, a set of input symbols and a transition function. The transition function takes a state and an input symbol as input and returns a state based on these input arguments. States are usually represented graphically by circles and transitions by labelled arrows. There are two special states: The start state and final state(s). The start state is the state in which the automaton starts processing the input symbols. It is represented graphically by an arrow pointing to the start state. Final states are represented graphically by a second circle embedded in the circle of a state. Starting from the start state the automaton processes the first input symbol by evaluating the transition function that takes the start state and the first input symbol as arguments. The automaton continues to process the input data by evaluating transition functions until the automaton has read every input symbol. The input word is part of the language the automaton represents if the automaton has reached a final state. If the automaton has not reached a final state, the input words is not part of the language the automaton represents. [2]

Deterministic and non-deterministic automata can be distinguished. Non-deterministic automata can have multiple transitions for the same input. Deterministic automata have unambiguous transitions for a given input. The concept of finite automata be applied to a lexical analysers for recognizing tokens (see section 2.4). [2]

Figure 2.1 shows an automaton consisting of the set of states  $\{z_0, z_1, z_2\}$  and several transitions. For convenience not all transitions are printed. If the first input symbol would be a zero, the automaton would transit in a so called error state which cannot be left regardless the input symbols.

The automaton accepts words consisting of a one, an arbitrary number of zeros and a one after that.

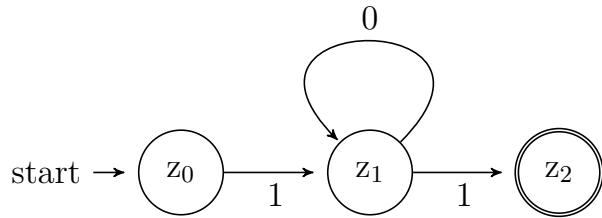


Figure 2.1: Finite automaton

## 2.1.5 Regular expression

A regular expression is an algebraic description of a regular language. Regular expressions declare strings that are part of the language and can describe the same words that can also be represented by a finite automaton. In comparison to finite automata, the words can be described in an algebraic way. Due to the similarities between regular expressions and finite automata, it is possible to convert regular expressions to finite automata and backwards. Regular expressions are often used to describe tokens that should be recognized by a lexer (see section 2.4). [2]

For example, the regular expression  $10^*$  denotes the language consisting of words that are made up a single 1 followed by an arbitrary number of 0's.

Given an alphabet  $\Sigma$ , regular expressions are formally defined as followed [2]:

1. The constants  $\epsilon$  and  $\emptyset$  are regular expressions denoting the empty string and the empty set.
2. If  $a$  is a symbol of the alphabet  $\Sigma$ , then  $a$  is a regular expression.
3. If  $a$  and  $b$  are regular expressions, then the alternation  $a|b$  (either  $a$  or  $b$ ), the concatenation  $ab$  ( $a$  followed by  $b$ ) and the Kleene star  $a^*$  (an arbitrary number of  $a$ 's) are regular expressions.

## 2.1.6 Formal grammars

Formal grammars describe how words of a language can be generated. In comparison to regular expressions, more classes of languages can be described. While regular expressions only describe regular grammars, formal grammars are able to describe all types of grammars of the Chomsky hierarchy.

The description of a grammar consists of four components:

- A set of symbols that defines the alphabet of the language. These symbols are called terminal symbols. The words of a language only consist of terminal symbols.
- A set of variables called nonterminal symbols. Each nonterminal symbol represents a set of words.
- One nonterminal symbol represents the language that is being defined. This nonterminal symbol is called the start symbol. Other nonterminal symbols help to define the language of the start symbol.
- A set of productions/rules. They recursively describe the language. Each production consists of a nonterminal symbol on the left-hand side that is being substituted by a sequence of terminal symbols and/or other nonterminal symbols on the right-hand side. There can also be empty productions, meaning a nonterminal symbol is substituted by an empty string.

[2]

### Context-free grammars

A context-free grammar is a grammar that only allows production rules that are in the form of  $V \rightarrow \beta$  with  $V$  being a nonterminal symbol and  $\beta$  being a sequence of nonterminal and terminal symbol with an arbitrary length. Context-free grammars are often the basis of a parser (see section 2.5). [2]

### Reduced grammars

Grammars are called reduced if each nonterminal symbol is terminating and reachable [3].

Given the set of terminal symbols  $\Sigma$ , a nonterminal symbol  $A$  is called terminating if there are productions  $A \xrightarrow{*} z$  (meaning some sequence of productions can be applied) so that  $z$  can be derived from  $A$  in one or multiple steps and  $z \in \Sigma^*$ . [3]

In other words, a nonterminal symbol  $A$  is terminating if there exist production rules so that  $A$  can be replaced by a string of terminal symbols.

Given the set of terminal symbols  $\Sigma$  and the start symbol  $S$ , a nonterminal symbol  $A$  is called reachable if there are production rules  $S \xrightarrow{*} uAv$  so that  $S$  can be derived to  $uAv$  and  $u, v \in \Sigma^*$ . [3]

In other words, a nonterminal symbol  $A$  is reachable if there exist production rules so that the start symbol can be replaced by a word containing  $A$ .

## 2.2 Backus-Naur Form (BNF)

The Backus-Naur Form (BNF) is a language to describe context-free grammars. In the Backus-Naur Form (BNF) nonterminal symbols are distinguished from terminal symbols by being enclosed by angle brackets, e. g.  $\langle TPTP\_File \rangle$  denotes the nonterminal symbol  $TPTP\_File$ . Productions are described using the “ $::=$ ” symbol and alternatives are specified using the “|” symbol. [4] Alternative symbols are used to conveniently describe several productions with the same left-hand side nonterminal symbol. An example for a BNF production using an alternative symbol would be

```
<TPTP_File> ::= <TPTP_Input> | <comment>
```

Without using alternative symbols it would be

```
<TPTP_File> ::= <TPTP_Input>
<TPTP_File> ::= <comment>
```

which is longer than the previous notation. Using the BNF pattern of notation whole grammars can be specified.

The Extended Backus-Naur Form (EBNF) extends the BNF with following rules:

- Optional expressions are surrounded by square brackets.
- Repetition is denoted by curly brackets.
- Parentheses are used for grouping.
- Terminals are enclosed in quotation marks.

[5]

## 2.3 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 syntax is specified in an [EBNF](#). [6]

Originally the [TPTP](#) used only [CNF](#) [7], but over the years expanded to various other types of logics including include [FOF](#) [7], Typed [FOF](#) [8, 9] and Typed Higher-order Form ([THF](#)) [10, 11].

The [TPTP](#) syntax uses an extended [BNF](#), with distinct rules for syntax, semantic constraints, tokens, and character classes which will be described in more detail in section 3.4. A small tool chain based on Lex/Yacc for building a basic parser from the syntax is available [12]. While the [TPTP](#) is a major success story, the resulting syntax is large and complex. Understanding the syntax is non-trivial. For new developers who want to implement “only” a first theorem prover for [CNF](#), identifying the relevant parts of the syntax and implementing a parser is a significant barrier to entry.

## 2.4 Lexing

Lexing or a so-called lexical analysis is the division of input into units called tokens [13].

The input is a string containing a sequence of characters. The lexer groups characters of the input string into sequences that are meaningful. These sequences are called lexemes. From each lexeme, the lexer produces a token, which consists of a token name and the lexeme string. The token name is an abstract symbol which is used during parsing. [14]

The tokens are then passed to the parser for syntax analysis. 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 [15].

A simple approach to build a lexer would be building an automaton for each token definition and then test to which automata the input corresponds. However, this would be inefficient because in the worst case the input needs to pass all automata before the belonging automata is identified. Building a single automaton that

tests each token simultaneously is more suitable. This automaton 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.

Potentially final states overlap as a consequence of one token being a substring of another token. For solving such conflicts, a lexer separates the input in order to divide it into tokens. Per convention the lexer chooses the longest input that matches any token. [15]

Furthermore, a precedence of tokens can be declared. Usually the token that is being defined first has a higher precedence and thus will be chosen if possible token matches have the same length. [15]

Besides of writing a lexer manually it can also be generated by a lexer generator. A lexer generator takes a specification of tokens as input and generates the lexer automatically. The specification of the tokens is usually written using regular expressions.

## 2.5 Parsing

Parsing is the process of analysing the grammatical structure of a sequence of tokens and can be used to verify that the input of the parser is part of the language that the parser accepts [14]. To do so, a parser builds a parse tree out of the tokens [15]. Parsers can also be generated automatically. A parser generator takes a description of the relationship among tokens in form of a formal grammar (see section 2.1.6) as input. The output is the generated parser. [13]

During the syntax analysis a parser takes a string of tokens and forms a syntax tree by finding the matching derivations. The matching derivation can be found using various approaches. The approaches that are used in this report will be introduced in the following.

### 2.5.1 Bottom-up parsing

In bottom-up parsing a parse tree for an input string is constructed beginning at the leaves working up to the root. The idea is to reduce a string to the start symbol of a grammar, constructing a derivation in reverse. [14]

At each reduction step, a specific substring matching the body of a production is

replaced by the nonterminal at the head of that production (reduction is reverse step of deviation). The key decisions are when to reduce and what production to apply. These decisions depend on the parser implementation. An example for bottom-up are shift-reduce parsers that are described in the following. [14]

### 2.5.2 Shift-reduce parsing

Shift-reduce parsers are a class of bottom-up parsers. In shift-reduce parsing a stack is used to hold grammar symbols and an input buffer holds the rest of the input string. During parsing the input is scanned from left to right. The parser shifts zero or more input symbols on the stack until it can reduce the elements at the top of the stack. The elements at the top of the stack are reduced to the head of the corresponding production. This procedure is repeated until the stack only contains the start symbol and the input buffer is empty or an error is detected. [14] There are four actions a shift-reduce parser can perform:

- **Shift:** Shift next input symbol on top of stack.
- **Reduce:** With the right end of a string to be reduced being at top of the stack and the left end of the string located within the stack, the string is replaced with a chosen nonterminal.
- **Accept:** Declare that input was successfully parsed.
- **Error:** Detect a syntax error.

[14]

There are two kinds of conflicts that can occur during shift-reduce parsing. A shift-reduce conflict occurs if the parser cannot decide whether to shift or reduce. A reduce-reduce conflict occurs if the parser cannot decide which of multiple reductions to make. [14]

### 2.5.3 LR( $k$ ) parsing

LR( $k$ ) parsing is the most common type of bottom-up parsing. The L stands for left-to-right, meaning the input is read from left to right, and R stands for constructing a rightmost derivation in reverse. In the context of LR parsing, the rightmost

derivation refers to reducing the input to the start symbol. The parameter  $k$  refers to the number of inputs symbols that are used as a lookahead. When  $k$  is omitted,  $k = 1$  is assumed. An LR parser makes shift-reduce decisions based on transitioning to different parsing states. An LR parser consists of a stack and a parsing table defining parsing-action functions and goto functions. An action function takes a state and an input symbol as input and defines which action to perform (Shift, Reduce, Accept, Error). A goto function takes a state and a nonterminal symbol as input and is used to find the next state after a reduction. States represent set of items. An item is a production of the grammar that the parser is based on with a dot at some position on the right-hand side of the production. An item indicates how much of a right-hand side has already been seen and what the parser is expecting next to apply the production (lookahead). For example  $A \rightarrow B \cdot C$  indicates that the parser has already read the input  $B$  and expects  $C$  in order to reduce the right-hand side to  $A$ . [14]

LR parsers are complex to implement, but LR parser generators can be used, which significantly reduces the effort for creating an LR parser. [14]

### LALR parsing

An LALR parser (Look-Ahead LR parser) is a simplified LR parser. The simplification consists of merging items that have identical cores. A core is the part of the right-hand side of a production that is left of the dot. For example, an LR parser would separate the items  $A \rightarrow B \cdot C$  and  $A \rightarrow B \cdot D$  as their core is identical but their lookahead is different. An LALR parser merges the two items to one item  $A \rightarrow B \cdot \{C, D\}$ . Not considering the lookahead reduces the power of the parser because it might lead to reduce/reduce conflicts. LALR parsers are nevertheless often used because the parsing table is much smaller comparing to other LR parsing techniques due to merging items. [14]

## 2.6 Lex and Yacc

Lex and Yacc are tools for writing lexers and parsers. They often work together, meaning the parser uses the generated tokens from the lexer as input.

### 2.6.1 Lex

Lex is a tool written in C that generates a lexer by specifying regular expressions and corresponding code fragments. The specified regular expressions are translated into a program that reads an input stream. The program partitions the input stream into strings matching the regular expressions. Once an expression is recognized the corresponding code is executed. For recognizing expressions, Lex builds a deterministic finite automaton. The automaton chooses the longest possible match if the input fits multiple specified regular expressions. [16]

Flex is an open-source alternative to Lex and also compatible to Lex. [17]

### 2.6.2 Yacc

Yacc is a tool written in C that can be used for parsing the input of a computer program. For generating a parser using Yacc, a specification that specifies the structure of the input has to be given. In addition to the specification of the structure, the input to the parser generator also consists of code that is invoked once a structure has been recognized. Yacc turns the input specification into a routine that handles the input. This routine calls Lex to get tokens from the input stream and arranges them based on the specification. [18]

Bison is an open-source general-purpose parser generator that can be used as an alternative to Yacc and also accepts Yacc grammars. [19]

## 2.7 Python

Since the use cases of [Synplifier](#) are all hand-created grammars of relatively small size (dozens to thousands of rules), and the algorithms do not have high complexity, we decided to implement the tool in Python 3. Python is a simple but powerful modern multi-paradigm language with good support and excellent libraries. Python is easy to learn and its power allows it to create complex applications using Python. Libraries that are used in the implementation [Synplifier](#) are introduced in the following sections.

## 2.7.1 PLY

Python Lex-Yacc (**PLY**) is an implementation of Lex and Yacc in Python. It is compatible with Python 2 and 3 and is fully implemented in Python. The goal of **PLY** is to offer a pure-Python implementation of Lex and Yacc. The implementation aims to rebuild the functionalities of Lex and Yacc and thus includes:

- Support for LALR parsing (see chapter 2.5.3)
- Support for empty productions
- Support for ambiguous grammars
- Precedence rules
- Error checking and recovery

**PLY** consists of the two modules *lex.py* and *yacc.py* that are meant to execute together. I. e. *yacc.py* retrieves tokens from the input stream provided by *lex.py*. *Yacc.py* returns an abstract syntax tree by default. However, it is possible for the user to change the output of *yacc.py* according to his demands. [20] **PLY** is used to generate the lexer and parser for the **TPTP** syntax files.

## 2.7.2 PyQt

PyQt is a Python binding for the cross-platform GUI framework Qt. It is licensed under the GNU GPL version 3. Qt offers a set of C++ libraries and development tools for various things including tools for building GUIs, defining regular expressions, setting up SQL databases, using NFC or interacting via Bluetooth. PyQt implements over 1000 of these Qt tools in Python. A simple GUI can be created by subclassing the QMainWindow, which provides a main window for an application. The QTreeView class provides an implementation of a tree view, which can be used to display data in a hierarchically structured way. [21]

### **2.7.3 argparse**

The Python module argparse is a module for creating command-line argument parsers. It provides the means to specify input arguments and automatically creates help and usage messages. It also checks if the given arguments are valid. From the specified input arguments, the module will automatically create a parser for the specified arguments. [22]

This module is used in the implementation of the command-line interface of Synplifier.

### **2.7.4 Beautiful Soup**

Beautiful Soup is a Python library for extracting data out of HTML and XML files. That makes it especially useful for getting information from web pages. [23, 24] In Synplifier Beautiful soup is used to automatically extract the latest TPTP syntax from the HTML version on the TPTP website.

# 3 Concept

This chapter outlines the concept and the architecture of [Synplifier](#). In section 3.1 a general approach to build [Synplifier](#) is outlined and in section 3.2 requirements [Synplifier](#) needs to meet are described. Section 3.3 introduces the architecture of [Synplifier](#) and its components. The following sections describe the concept of each component that is developed.

## 3.1 Approach

The goal is to decompose the full [TPTP](#) syntax into smaller sub-syntaxes extracting desired sub-syntaxes. We propose the following approach: The user selects a start symbol and can disable undesired productions. A recursive algorithm is introduced that checks which symbols are still terminating and reachable and extracts those (and the productions that contain only them) into a sub-syntax.

## 3.2 Requirements

The following section outlines requirements of [Synplifier](#). Requirements are numbered and split into three categories:

- Requirements concerning the GUI are marked by a *G*
- Requirements concerning the Command-line interface are marked by a *C*
- General requirements concerning [Synplifier](#) are marked by an *R*

The requirements are described using capitalized keywords like MUST or SHOULD following [25].

**[G1]**

Synplifier MUST provide a GUI.

**[G2]**

Synplifier MUST provide the possibility for the user to import a TPTP syntax file into the GUI by specifying a TPTP syntax file path and a desired start symbol.

**[G3]**

The imported syntax MUST be displayed in the GUI. This includes displaying productions that are defined by the syntax as well as comments that are associated with these productions.

**[G4]**

The user MUST be able to select a new start symbol to generate a sub-syntax.

**[G5]**

The user MUST be able to select which productions should be blocked to generate a sub-syntax.

**[G6]**

Synplifier MUST provide the functionality for the user to import a control file into the GUI by specifying a control file path.

**[G7]**

After importing a control file, selections are set according to the control file content in the GUI.

**[G8]**

Synplifier MUST provide the functionality to display a generated sub-syntax in the GUI.

**[G9]**

The user MUST be able to export a control file based on the users start symbol selection (G2) and blocked productions selection (G3) in a .txt-format.

**[G10]**

Synplifier MUST provide the functionality to export a generated sub-syntax to .txt format by specifying a file path in the GUI.

**[C1]**

Synplifier MUST provide a command-line interface.

**[C2]**

The command-line interface MUST provide the functionality to extract a sub-syntax with the user specifying the **TPTP** syntax file and control file path.

**[C3]**

The extracted sub-syntax is saved in .txt format at a file path that has been specified by the user.

**[C4]**

The command-line interface MUST provide help information which describe all command-line parameters.

### [R1]

When generating a sub-syntax comments referring to rules that are present in the reduced syntax SHOULD be maintained in the sub-syntax.

### [R2]

Output [TPTP](#) sub-syntaxes SHOULD be compatible to the automated parser generator for the [TPTP](#) language [12].

### [R3]

A generated sub-syntax MUST only contain terminating and reachable symbols.

## 3.3 Overview

Figure 3.1 outlines the procedure of extracting a sublanguage of the [TPTP](#) language. The first task is to import the [TPTP](#) syntax file and extract the tokens inside that file using the lexer. The next phase is for the parser to create a data structure from the tokens, also checking if the syntax in the syntax file was correct. Then, a graph representing the imported [TPTP](#) syntax is built.

This graph is subject to manipulation by disabling certain transitions or selecting a new start symbol in the following phase. This includes computation of the remaining reachable and terminating syntax. The modified graph represents the syntax of the extracted sub-language. To make this syntax usable, lastly the syntax must be output in the same format as the original syntax.

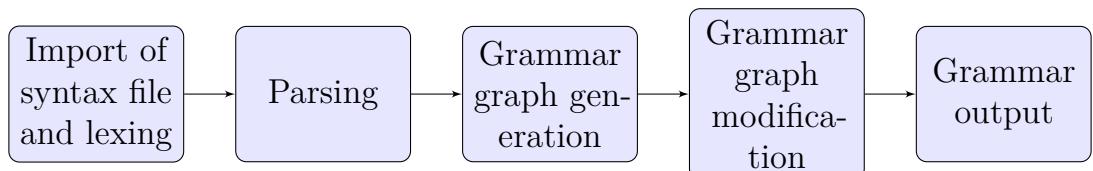


Figure 3.1: Procedure of extracting a sublanguage

### 3.3.1 Proposed architecture

The proposed architecture for [Synplifier](#) consists of seven components:

- A lexer for extracting tokens from the [TPTP](#) syntax
- A parser for creating a data structure from the tokens
- A graph builder for building a graph out of the parsers data structure and manipulating this graph
- A graphical user interface
- A console interface
- An input module that imports the [TPTP](#) syntax file
- An output module that exports extracted sub-syntaxes to local storage

Figure 3.2 contains a high-level UML diagram describing the architecture of [Synplifier](#). The user interacts either with the *Console* or *View* class. The *Console* class provides the command-line interface and the *View* class provides the GUI. Both have a reference on *Input* and *Output* for reading from and writing to files. They also have a reference on the *TPTPGraphBuilder* class. This class is responsible for building a grammar graph and extracting sub-syntaxes by graph manipulation. For that, lexing and parsing are necessary. The *TPTPGraphBuilder* uses the *Parser* class for getting a [TPTP](#) syntax representation and the *Parser* uses the *Lexer* to extract the tokens from a [TPTP](#) syntax file.

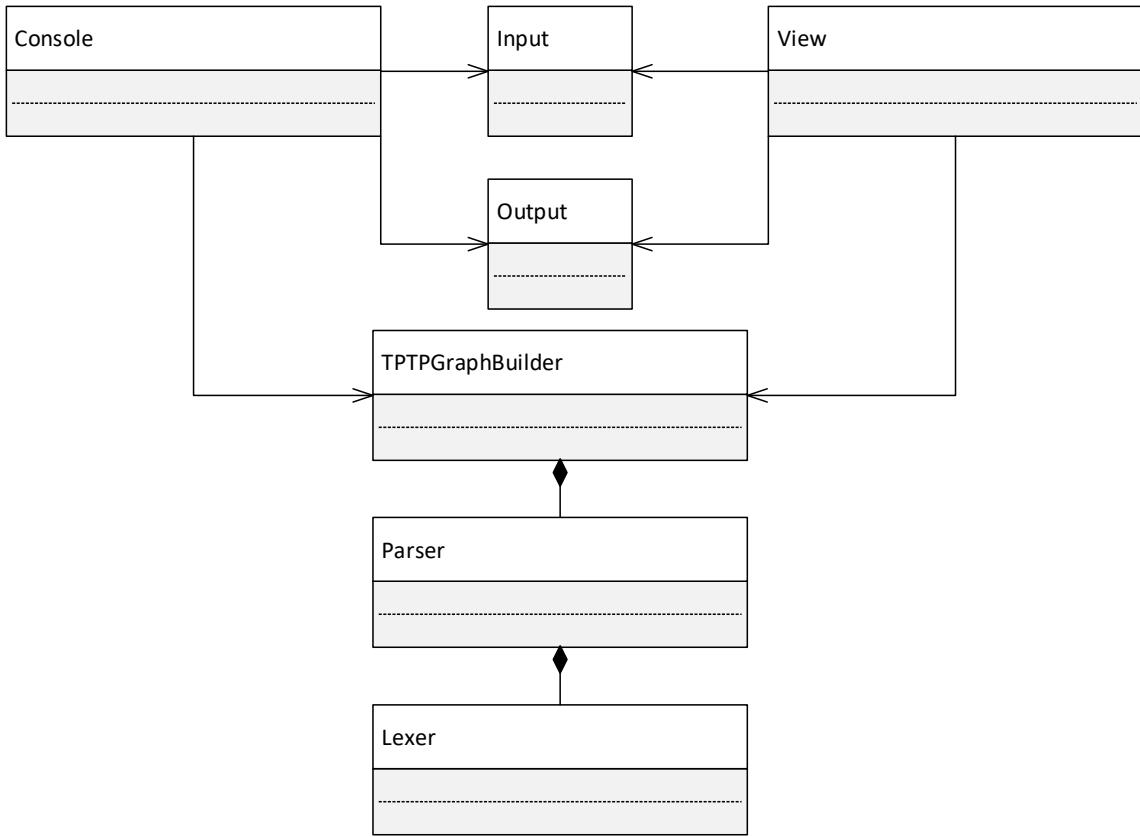


Figure 3.2: UML diagram of the architecture of [Synplifier](#)

## 3.4 Lexer

The lexer is responsible for extracting tokens from the [TPTP](#) syntax file. Using [PLY](#), a lexer can be generated by specifying tokens using regular expressions. [PLY](#) is used because it is the most popular tool for Lex and Yacc in Python. In order to find elementary tokens, the [TPTP](#) syntax has been analyzed and regular expressions that precisely describe these tokens have been defined.

The standard [EBNF](#) only uses one production symbol (" ::= "). In the [TPTP](#) syntax the standard production symbol is used for syntactic rules. Additional symbols for semantic, lexical and character-macro rules have been added. Semantic rules describe rules whose left-hand side is already part of a syntactic rules but from which only a subset makes semantic sense. Any further rules that are reached from the nonterminal symbols on the right-hand side of a semantic rule are semantic

rules as well. Lexical rules are rules that produce tokens. Character-macro rules are rules that define regular expressions for the tokens. The following table contains the production symbols for grammar (syntactic rules), strict (semantic rules), token (lexical rules) and macro (character-macro rules) rule types used in the TPTP syntax.

Table 3.1: TPTP language rule types [12]

Symbol	Rule Type
::=	Grammar
:=:	Strict
::-	Token
:::	Macro

The following paragraph introduces the tokens that are recognized by the lexer. Tokens are implemented as Python classes (see chapter 4.1) and written in bold in the following paragraphs.

Following standard BNF, **nonterminal** symbols are enclosed by the < and > symbol. In between there can be any arbitrary sequence of alphanumerical characters and underscores.

A **terminal symbol** does not have any special notation and can be any sequence of characters, excluding whitespace characters, that is not matched by any other token.

There are four **expression** token types (one for each rule type). **Expressions** are defined as a nonterminal symbol followed by a production symbol ( ::= for grammar, ::- for token, :=: for strict, :::: for macro rule type). The nonterminal symbol and the following production symbol are selected to be a single token and are not identified as two separate tokens to clearly identify the start of a new rule and therefore avoid ambiguity while parsing. The listing below features two tokens: A grammar expression and a nonterminal symbol.

```
1 <formula_role>      ::= <lower_word>
```

Listing 3.1: TPTP rule example

TPTP **comments** always occupy a full line. They start with a percentage sign and continue to the end of the line. The percentage sign can appear as a regular character in **EBNF** productions but also indicates the start of a comment. But when it is used as a terminal symbol it is always embedded in square brackets and can therefore never be the first character of a new line. Therefore, a percentage sign used as terminal symbol cannot be mistaken as the start of a comment.

Additional tokens are the meta-symbols including **open** "(" and **close parentheses** ")", **open** "[" and **close square brackets** "]", asterisks "\*" called **repetition symbols** and vertical bars "|" called **alternative symbols**.

## 3.5 Parser

The parser takes the tokens from the lexer as input and creates a data structure that represents the structure of the **TPTP** syntax. The parser is based on a context-free grammar that is described in section 3.5.1.

Figure 3.3 outlines the responsibilities of the parser component and the sequence of its sub-functions. First, the tokens generated by the lexer need to be parsed and based on that the data structure representing the **TPTP** syntax is to be created. Secondly, the rules in the data structure have to be numbered to maintain the correct order for output. The third step is the so-called disambiguation of square brackets. In the **TPTP** syntax, square brackets not necessarily denote that an expression is optional, which is the case in traditional **EBNF**. In token and macro rules they denote that an expression is optional and in grammar and strict rules square brackets are terminals. Therefore, disambiguation of square brackets is necessary. Finally, the output of the parser is returned.

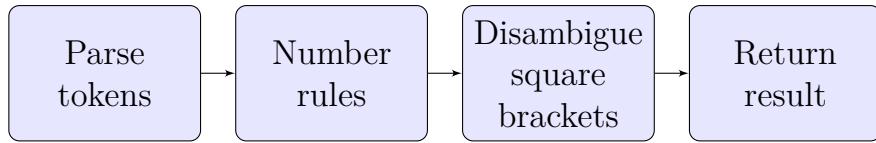


Figure 3.3: Parsing procedure

### 3.5.1 Defined grammar

The context-free grammar that the parser is based on is formally specified by  $G = (N, \Sigma, P, S)$ .

$N$  is the set of nonterminal symbols. The set includes grammar list, comment block, grammar expression, token expression, strict expression, macro expression, productions list, production, xor productions list, terminal symbol production and production element. The nonterminal symbols are described in detail in chapter 4.2.

$\Sigma$  is the set of terminal symbols. Terminal symbols of the specified grammar are the tokens generated by the lexer (see 3.4).

$P$  is the set of production rules that is presented in the following.

A grammar list implies a comment block, the four expressions (grammar expression, token expression, macro expression, strict expression) or a grammar list followed by one of the expressions.

$$\begin{aligned}
 \text{grammar list} \rightarrow & \text{comment block} \mid \text{grammar expression} \\
 & \mid \text{token expression} \mid \text{strict expression} \\
 & \mid \text{macro expression} \\
 & \mid \text{grammar list } \text{grammar expression} \\
 & \mid \text{grammar list } \text{strict expression} \\
 & \mid \text{grammar list } \text{macro expression} \\
 & \mid \text{grammar list } \text{comment block}
 \end{aligned}$$

A comment block is either a comment or a comment block with a comment.

*comment block* → **comment** | *comment block* **comment**

The four expressions are the expression token followed by an productions list. Grammar expression has a special case without a productions list because the production rule

<null> :::=

of the TPTP syntax has no production on the right-hand side.

*grammar expression* → 1 **grammar expression** *productions list*  
| 1 **grammar expression**

*token expression* → 1 **token expression** *productions list*

*strict expression* → 1 **strict expression** *productions list*

*macro expression* → 1 **macro expression** *productions list*

Productions list and xor productions list imply either a production or a productions list alternative symbol production.

*productions list* → *production*  
| *productions list* **alternative symbol production**

*xor productions list* → *production*  
| *xor productions list* **alternative symbol production**

Terminal symbol production is either a terminal symbol or a terminal symbol/repetition symbol/ alternative symbol embedded in square brackets.

*terminal symbol production* → **open square bracket** **terminal symbol**  
**close square bracket**  
| **open square bracket** **repetition symbol**  
**close square bracket**  
| **open square bracket** **alternative symbol**  
**close square bracket**  
| **terminal symbol**

A production element can be replaced by a nonterminal symbol or by a nonterminal symbol in square brackets or nonterminal symbol repetition. In the case of repetition or square brackets the production element is categorized as optional when in square brackets or as repetition when followed by the repetition symbol. The same applies to terminal symbol production. A production element can also only be square brackets only.

*production element* → **open square bracket** **nonterminal symbol**  
**close square bracket**  
| **nonterminal symbol** **repetition symbol**  
| *terminal symbol production* **repetition symbol**  
| **open square bracket** **close square bracket**  
| **nonterminal symbol**  
| *terminal symbol production*

A production is either a production element, an already existing production followed by a production element or a composition of production(s), a xor production list, parenthesis and a repetition symbol.

```
production → production element | production production element
| open parenthesis xor productions list
close parenthesis
| open parenthesis production close parenthesis
| production open parenthesis production
close parenthesis
| production open parenthesis xor productions list
close parenthesis
| open parenthesis production close parenthesis
production
| open parenthesis xor productions list
close parenthesis production
| open parenthesis production close parenthesis
repetition symbol
| production open parenthesis production
close parenthesis repetition symbol
```

The start symbol is not specifically mentioned because it is per convention the left element of the first rule.

### 3.5.2 Disambiguation of square brackets

As mentioned at the beginning of the parsing section, square brackets have different meanings depending on the rule type. The idea to solve this problem is to treat all rules the same in the first processing step. Square brackets are then interpreted as denoting the optional production property. This production property is then selected for productions that are enclosed by square brackets for all types of rules. In an additional processing step, after creating the grammar list each grammar and strict rule is iterated, exchanging the production property optional by the square bracket terminal symbols. The output of the parser is a list of the rules and the comments from the [TPTP](#) syntax file.

## 3.6 Graph generation

After parsing, the parsed TPTP syntax is stored in a grammar list. The grammar list data structure does not allow easy traversing which is needed for extracting a sub-language. That is why a new data structure is introduced that allows easy modification and traversing. The data structure that is used is a graph representing the grammar.

In the graph, the nonterminal symbols of the TPTP syntax are represented by a node class that has the following attributes (see also UML class diagram 4.10):

- Value: name of nonterminal symbol
- Productions list: productions list of nonterminal symbols that has been created by the parser
- Rule type: rule type of nonterminal symbol
- Comment block: list of comments belonging to the nonterminal symbol
- Position: position of the production in the input file
- Children: nested list that contains the node of every nonterminal symbol that is part of the productions list grouped by productions in the productions list

Before generating the graph, a dictionary is created that provides an efficient way accessing the nodes in order to build the grammar graph and also during the next steps of the sub-syntax generation. This dictionary contains nodes constructed from the grammar list output by the parser. The combination of the nodes value and rule type form the unique key for each node. While constructing the dictionary, comments in the grammar list are associated to nodes using a heuristic that is described in section 3.8.

Also, a new temporary start symbol is introduced. This is necessary because one nonterminal symbol in the TPTP syntax can be mapped to multiple nodes which the following example 3.2 shows.

The example below shows the productions for the nonterminal symbol  $< formula\_role >$ . Since this nonterminal symbol has multiple types of rules one node will be created for each rule type.

```
1 <formula_role> ::= <lower_word>
2 <formula_role> ::= axiom | hypothesis | definition | assumption |
3   lemma | theorem | corollary | conjecture |
4   negated_conjecture | plain | type | fi_domain |
5   fi_functors | fi_predicates | unknown
```

Listing 3.2: Productions of the nonterminal symbol  $\langle formula\_role \rangle$

If a nonterminal symbol that has multiple rule types is selected as the desired start symbol, multiple nodes would represent that start symbol and therefore it would not be possible to select one node as the starting point of the graph generation. To solve this problem, a temporary start symbol is introduced before generating the graph. The temporary start symbol produces the start symbol that the user specified and is used as a starting point for the graph generation. If the nonterminal symbol  $\langle formula\_role \rangle$ , that was mentioned before, would be selected as start symbol by the user, a temporary start symbol representing the rule

```
<start_symbol>      ::= <formula_role>
```

would be introduced. This ensures that only one node is representing the start symbol that is used for graph generation.

Starting with the temporary start symbol, the graph is generated recursively (described in chapter 4.3).

The following example shows a production rule and an abstract representation of the resulting list of children belonging to the node. Each production alternative is embedded in the children list of the node.

```
<disjunction>  ::=  <literal> | <disjunction><vline><literal>
```

Output:

```
node.value:    <disjunction>
node.ruleType: grammar
node.children: [[<literal>], [<disjunction>, <vline>, <literal>]]
```

## 3.7 Control file

A format for specifying the desired start symbol and blocked productions is required. Using a file-based configuration enables the user to store desired configurations so that a manual selection in the GUI is not necessary. It is also essential for using the command line interface, because there manual selection is not possible. The file should be human-readable and -editable. The format should be easy to parse and allow to specify all necessary information. This includes the desired start symbol and all production rules that should be blocked. The proposed way to describe this information is to:

- Define the desired start symbol in the first line.
- Define blocked productions grouped by the nonterminal symbol and rule type symbol, separating each group by a new line. First defining the nonterminal symbol, then the production symbol and after that the indexes of the alternatives that should be blocked (indexing starts at zero).

Identifying the production symbol is necessary because there are nonterminal symbols that have productions with more than one production symbol.

The example below contains a sample control file. In this file, *<TPTP\_file>* is specified as start symbol. The second line indicates that the second grammar production alternative of the nonterminal symbol *<TPTP\_input>* should be disabled. Furthermore, the first, second, third and fifth grammar production alternative of the nonterminal symbol *<annotated\_formula>* are said to be disabled in line 3.

```
1 <TPTP_file>
2 <TPTP_input>, ::=, 1
3 <annotated_formula>, ::=, 0, 1, 2, 5
```

Listing 3.3: Control file example

This format is easy to parse and also enables users to specify their desired start symbols and blocked productions without having to use the GUI. The control file can also be generated from selections made by the user in the GUI. It can be exported as a text file. When generating a sub-syntax from the user GUI selection a control file representing that selection is created internally as input for the next processing steps.

Specifying which production should be blocked, and not the ones should be kept, typically results in a significantly smaller file. Storing the indexes of the productions that should be blocked offers that in case productions are renamed the control file would still be valid. On the other hand, if productions are added or deleted from the original [TPTP](#) syntax, the control file may have to be updated.

## 3.8 Maintaining comments

In the [TPTP](#) syntax there are comments providing supplemental information about the language and its symbols and rules. When generating a reduced grammar maintaining comments is desired. This means that comments from the original language specification should be associated with the rule they belong to and if the rule is still present in the reduced grammar, the comment should be as well. Therefore, a mechanism has to be designed for the association of comments to grammar rules.

Listing 3.4 features an example of a comment in a [TPTP](#) syntax file. This comment begins with a *Top of Page* line which, in the HTML version of the [TPTP](#) syntax, contains a hyperlink which leads to the beginning of the syntax file. The next line contains a relevant comment.

```
1 %——Top of Page——  
2 %——TFF formulae.  
3 <tff_formula> ::= <tff_logic_formula> | <tff_atom_typing> |  
4           <tff_subtype> | <tfx_sequent>
```

Listing 3.4: Example of a comment in the [TPTP](#) syntax

The heuristic matching comments to rules takes these *Top of Page* lines into account. When there is a *Top of Page* line in between comment lines it generally also splits comments semantically. In listing 3.4 can be seen that the comment in line 2 refers to the rule after. Therefore, it would be correct to associate the comment line after the *Top of Page* line to the rule after. Also, if there is one *Top of Page* line in between multiple comment lines it is highly probable that the first part of the comment lines before the *Top of Page* line refer to the rule before the comments and that the lines after the *Top of Page* line refer to the rule after the comment lines. This scenario can be seen in listing 3.5. The *Top of Page* line is in line 28 and

the comment lines before refer to the rule before. The comment line after refers to the rule after that line.

```

1 <formula_role>      ::= axiom | hypothesis | definition | assumption |
2                                lemma | theorem | corollary | conjecture |
3                                negated_conjecture | plain | type |
4                                fi_domain | fi_functors | fi_predicates | unknown
5 %----"axiom"s are accepted, without proof. There is no guarantee that the
6 %----axioms of a problem are consistent.
7 %----"hypothesis"s are assumed to be true for a particular problem, and are
8 %----used like "axiom"s.
9 %----"definition"s are intended to define symbols. They are either universally
10 %----quantified equations, or universally quantified equivalences with an
11 %----atomic lefthand side. They can be treated like "axiom"s.
12 %----"assumption"s can be used like axioms, but must be discharged before a
13 %----derivation is complete.
14 %----"lemma"s and "theorem"s have been proven from the "axiom"s. They can be
15 %----used like "axiom"s in problems, and a problem containing a non-redundant
16 %----"lemma" or theorem" is ill-formed. They can also appear in derivations.
17 %----"theorem"s are more important than "lemma"s from the user perspective.
18 %----"conjecture"s are to be proven from the "axiom"(-like) formulae. A problem
19 %----is solved only when all "conjecture"s are proven.
20 %----"negated_conjecture"s are formed from negation of a "conjecture" (usually
21 %----in a FOF to CNF conversion).
22 %----"plain"s have no specified user semantics.
23 %----"fi_domain", "fi_functors", and "fi_predicates" are used to record the
24 %----domain, interpretation of functors, and interpretation of predicates, for
25 %----a finite interpretation.
26 %----"type" defines the type globally for one symbol; treat as $true.
27 %----"unknown"s have unknown role, and this is an error situation.
28 %----Top of Page
29 %----THF formulae.
30 <thf_formula>      ::= <thf_logic_formula> | <thf_atom_typing> |
31                                <thf_subtype> | <thf_sequent>

```

Listing 3.5: Comment lines split by a *Top of Page* line in the TPTP syntax

The flow chart in figure 3.4 shows the process of matching comment blocks, that are consecutive comment lines, to rules. First, the comment block is split into multiple separate comment blocks by using *Top of Page* lines as separators.

- If this results in no comment blocks the comment block consisted only of one line which was a *Top of Page* line. Then no comment block has to be associated to a rule because *Top of Page* lines are not relevant.
- If this results in one comment block, that means that no *Top of Page* line was present in the comment block or that there was a *Top of Page* line at the beginning or end of the comment block. The comment block is associated with the rule after if the comment block is not at the end of the file. If it is at

the end of the file, it is associated with the rule before, because no rule after exists.

- If this results in two comment blocks, one *Top of Page* line was present in between comment lines. Then the comment block before the *Top of Page* line is associated with the rule before when possible and the comment block after the *Top of Page* line is associated with the rule after. If the original comment block was at the beginning of the file, both comment blocks resulting from splitting the original comment block are associated with the rule after. If the original comment block was at the end of the file, the comment blocks resulting from the split are associated with the rule before.

The case of three or more comment blocks after splitting the original comment block is not featured in the flow chart. This case does not occur in the [TPTP](#) syntax version 7.3.0. Since it might occur in a future version of the [TPTP](#) syntax it is handled by merging all comment blocks starting from the second and then following the procedure of two comment blocks in the flow chart in figure 3.4.

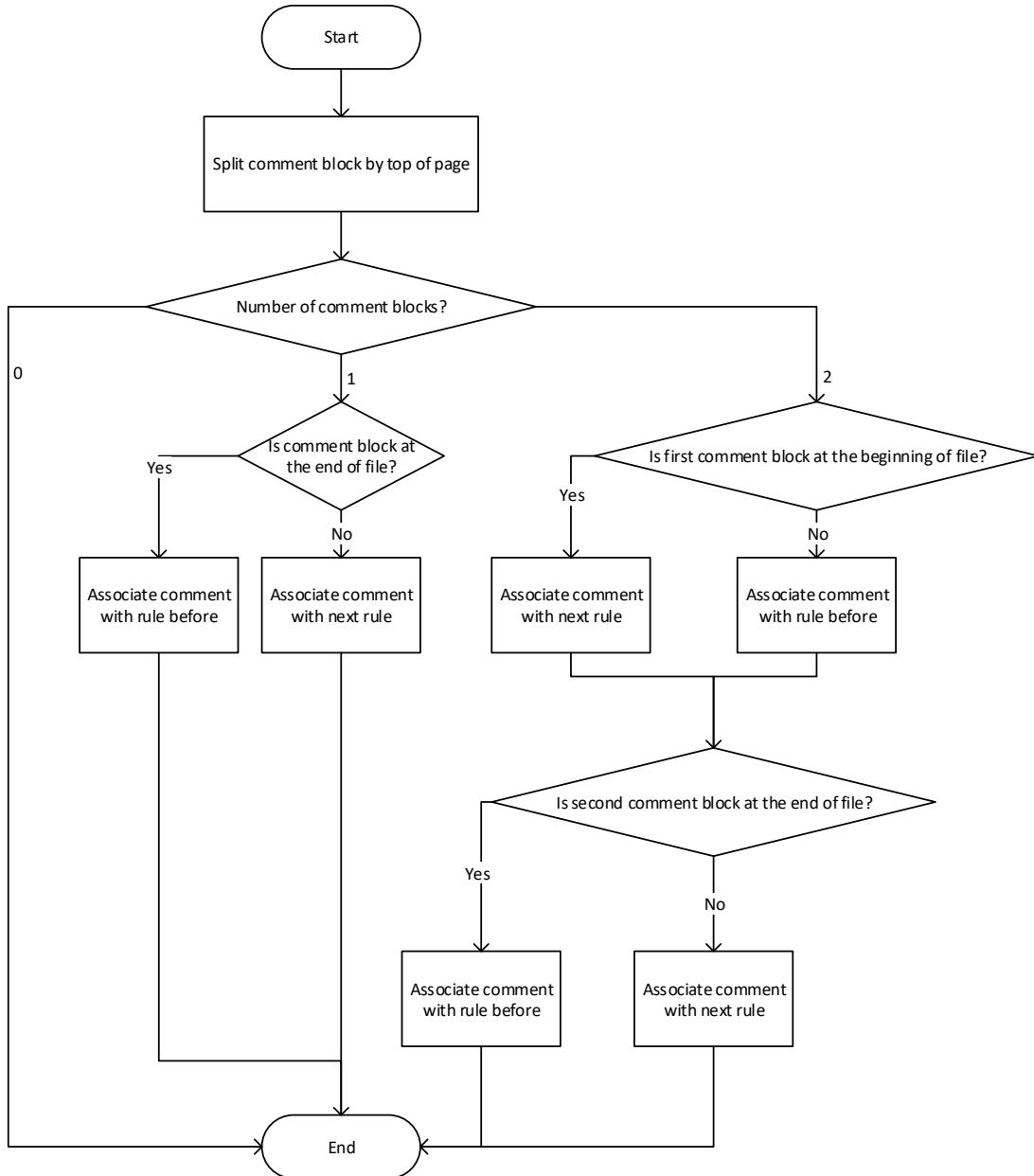


Figure 3.4: Maintaining comments flow chart

## 3.9 Extraction of a sub-syntax

This section covers the concept of how a sub-syntax is computed from the original syntax. The original syntax is represented by a grammar graph (see section 3.6) and the information on what part of the syntax to extract is described in a control

file. To extract a sub-syntax from the original grammar graph, four steps must be performed:

1. The control file has to be parsed, in order to get information on the desired start symbol and which productions should be blocked.
2. The blocked productions specified in the control file must be disabled and therefore the corresponding transitions must be removed from the grammar graph.
3. The remaining part of the grammar that is terminating must be computed.
4. From the part of the grammar that is terminating, non-reachable symbols must be removed.

### 3.9.1 Parsing a control file

The control file provides the necessary information for the extraction of a sub-syntax. The format of the control file is described in section 3.7. The start symbol is listed in the first line. Every consecutive line contains a nonterminal symbol, the corresponding rule type symbol, and the indexes of the productions that should be blocked separated by comma. The start symbol will be relevant in determining the remaining reachable part of the grammar (section 3.9.4) whereas the information on the productions that should be blocked will be needed in the next section.

### 3.9.2 Removal of blocked productions

In the control file, for each nonterminal symbol, that will be modified, its name, rule type and the indexes of the productions that should be blocked are listed. From all nodes that are addressed (by nonterminal symbol name and rule type), the indexed productions are removed. This includes deleting the corresponding element from the productions list and from the children list.

### 3.9.3 Determination of the remaining terminating symbols

After the desired productions have been deleted from the grammar graph, the next step is to remove the nonterminating symbols from the grammar graph. First it must be determined which symbols are terminating and which are nonterminating. Terminating nodes are found by iterating the nodes in the *children\_list* of a node. If the *children\_list* is empty, the production only consists of terminal symbols. If that is not the case, the node is terminating if all nodes in the *children\_list* represent a nonterminal symbol which is terminating. If a nonterminal symbol has multiple rule types, it is represented by multiple nodes. Also, if a nonterminal symbol with multiple rule types is featured in a production, all nodes representing that nonterminal symbol are in the children list corresponding to this production. However, only one of the nodes needs to be terminating for the nonterminal symbol to be terminating.

After the terminating symbols have been determined, productions that contain a nonterminating symbol must be removed from the nodes of the grammar graph. Because some nonterminal symbols might have productions that contain nonterminating symbols, but also other productions that only contain terminating symbols, only the productions containing nonterminating symbols must be removed.

### 3.9.4 Determination of the remaining reachable symbols

After removing the productions specified in the control file it has to be determined which part of the grammar still remains reachable. This is done by generating a new grammar graph starting from the desired start symbol. When generating the new grammar graph, all reachable parts of the grammar from the start symbol will be added to the new grammar graph (see section 3.6). All nodes that are not part of the new grammar graph are not reachable and can therefore be removed. They are also removed from the nodes dictionary.

## 3.10 Output generation

Within Synplifier a syntax is represented by a grammar graph. To export subsyntaxes that have been extracted from a grammar graph they have to be converted

to the original form of a **TPTP** syntax file. This process is described in the following subsection. After that, measures to provide compatibility with the automated parser generator are discussed in section 3.10.2.

### 3.10.1 Create output from grammar graph

There are three steps necessary in order to convert a syntax represented by a grammar graph to the original **TPTP** syntax file form. These steps are displayed in figure 3.5.

The first step is to traverse the grammar graph by iterating all nodes and getting a string representation of the nodes.

In order to maintain the same order that rules have had in the original **TPTP** syntax, the position of the nodes is saved together with the node string in a tuple. After generating all node string representations, in the second step all tuples are ordered based on their position.

The ordered list of nodes strings is then, in the third step, written to a file and exported. The node that represents the temporary start symbol is not printed as it does not belong to the **TPTP** syntax.

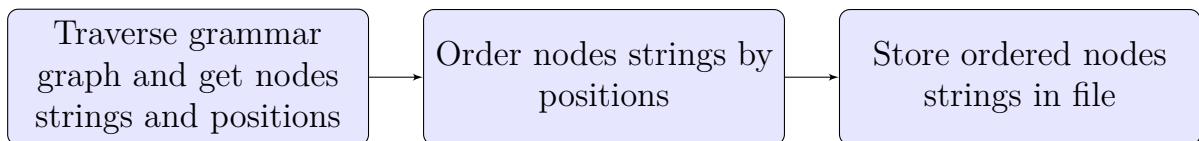


Figure 3.5: Procedure of generating a string representation of the grammar graph

### 3.10.2 Automated parser generator compatibility

The automated parser generator for the **TPTP** syntax [12] takes a **TPTP** syntax file as input and creates a Lex and Yacc file corresponding to the specification in the input syntax. It is used for validating that an extracted sub-syntax is compatible with the original **TPTP** syntax.

In the **TPTP** syntax version 7.3.0 there is a definition of the syntax of comments which is shown in listing 3.6.

```

1 <comment>           ::= <comment_line>|<comment_block>
2 <comment_line>       ::= [%]<printable_char>*
3 <comment_block>      ::= [/][*]<not_star_slash>[*][*]*[/]
4 <not_star_slash>     ::= ([^*]*[*][*]*[^/*])*[^*]*
5 <printable_char>    ::= .

```

Listing 3.6: Comment syntax definition in the **TPTP** syntax

However, this part of the syntax is not reachable from the start symbol *TPTP\_file*. When constructing the grammar graph this part of the syntax is therefore removed. Using a syntax file, in which the comment syntax is not present, with the automated parser generator results in an error because the comment symbol is not present in the lex specification but is referred to in a Yacc rule.

There are two ways to include the comment syntax in the output of **Synplifier** in order to be accepted by the automated parser generator. Either making the *<comment>* nonterminal symbol reachable by for example adding it as an alternative to *<TPTP\_input>* which can be seen in listing 3.7 or maintaining the comment syntax separately and adding it to the output syntax even though it is not reachable.

```

1 <TPTP_file>          ::= <TPTP_input>*
2 <TPTP_input>          ::= <annotated_formula> | <include> | <comment>

```

Listing 3.7: Making the comment syntax reachable

Both ways are supported by **Synplifier**.

If the *<comment>* nonterminal symbol is made reachable in the **TPTP** syntax, the grammar graph generated by **Synplifier** contains the comment syntax and no further action is necessary.

To add the comment syntax to the output if the *<comment>* nonterminal symbol is not reachable is done by using an external configuration file that contains the syntax from listing 3.6. From the syntax in this configuration file a separate grammar graph is generated. The output string that can be generated from this grammar graph is added to the output of **Synplifier** when outputting a generated sub-syntax.

The GUI provides a menu option for appending *<comment>* to the output file (see chapter 3.11.4).

## 3.11 GUI

The graphical user interface is built using PyQt. Advantages of PyQt are that it offers a treeview that is used for displaying the grammar. In comparison to the standard Python tool kit Tkinter [27] PyQt offers checkboxes in treeviews that are used for selecting blocked productions.

The menu of the GUI consists of five submenus.

### 3.11.1 Import Menu

The *Import* menu that can be seen in figure 3.6 provides the option to import a TPTP syntax file in .txt file format from local storage or to import the latest TPTP syntax version from the TPTP website. The TPTP syntax on the TPTP project website is stored in a HTML format. This file is downloaded and converted to plain text using the BeautifulSoup Python library [23].



Figure 3.6: Import menu

After selecting a file for import, a start symbol needs to be selected. Starting with the start symbol a grammar graph is generated and the corresponding text displayed. Each rule in the TPTP syntax is a top-level element of the treeview and can be expanded to show the productions alternatives. Right of the nonterminals name the rule type is displayed. Comments that have been assigned to a rule are also displayed. An example of the imported TPTP syntax is shown in figure 3.7.

### 3 Concept

The screenshot shows the TPTP GUI interface. At the top, there is a menu bar with options: Import syntax, Save syntax, Reduce, Control file, and View. Below the menu is a toolbar with icons for opening files, saving, reducing, and controlling. The main area displays a grammar tree. The root node is <TPTP\_file>, which has two children: <TPTP\_input> and <comment>. The <comment> node is expanded, showing two children: <annotated\_formula> and <include>. The <TPTP\_input> node is collapsed. To the right of the tree, there is a column labeled 'GRAMMAR' next to each node. A status bar at the bottom shows the version information: %----v7.3.0.0 (TPTP version.internal development number).

```
%----v7.3.0.0 (TPTP version.internal development number)
%
%----README ... this header provides important meta- and usage information
%
%----Intended uses of the various parts of the TPTP syntax are explained
%----in the TPTP technical manual, linked from www.tptp.org.
%
%----Four kinds of separators are used, to indicate different types of rules:
%---- ::= is used for regular grammar rules, for syntactic parsing.
%---- := is used for semantic grammar rules. These define specific values
%---- that make semantic sense when more general syntactic rules apply.
%---- :: is used for rules that produce tokens.
%---- :: is used for rules that define character classes used in the
%---- construction of tokens.
%
%----White space may occur between any two tokens. White space is not specified
%----in the grammar, but there are some restrictions to ensure that the grammar
%----is compatible with standard Prolog: a <TPTP_file> should be readable with
%----read/1.
%
%----The syntax of comments is defined by the <comment> rule. Comments may
%----occur between any two tokens, but do not act as white space. Comments
%----will normally be discarded at the lexical level, but may be processed
%----by systems that understand them (e.g., if the system comment convention
%----is followed).
%
%----Multiple languages are defined. Depending on your need, you can implement
%----just the one(s) you need. The common rules for atoms, terms, etc, come
%----after the definitions of the languages, and mostly all needed for all the
%----Files. Empty file is OK.
```

	GRAMMAR
<TPTP_file>	<TPTP_input>*
<TPTP_input>	<annotated_formula> <include>
<comment>	
<annotated_formula>	GRAMMAR
<tpi_annotated>	GRAMMAR
<tpi_formula>	GRAMMAR
<thf_annotated>	GRAMMAR
<tff_annotated>	GRAMMAR
<tcf_annotated>	GRAMMAR

Figure 3.7: Main GUI view

#### 3.11.2 View Menu

The *View* menu, shown in figure 3.8, provides the possibility to toggle the display of comments to improve readability.



Figure 3.8: View menu

In order to extract a sub-syntax, the user has to choose a new start symbol by checking one of the check boxes left of the nonterminal symbols name. By default, the start symbol that has been selected after importing the syntax is selected. For selecting blocked productions, the user has the choice of on the one hand expanding the rules and uncheck the checkboxes belonging to the productions the user wishes to block. On the other hand, the user can import a control file. If the user imports a control file, the checkboxes are set accordingly.

#### 3.11.3 Reduce Menu

In the *Reduce* menu that can be seen in figure 3.9 the user can reduce the syntax according to his selections. The syntax is reduced and the reduced syntax is displayed afterwards. The rules maintain the same order as in the original TPTP syntax.

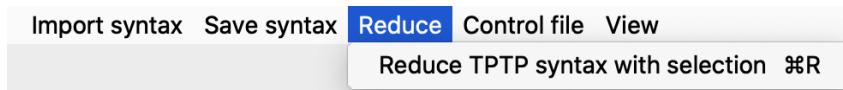


Figure 3.9: Reduce menu

#### 3.11.4 Save syntax Menu

In the *Save syntax* menu that is shown in figure 3.10 the user can reduce the syntax based on his GUI selection or an imported control file. The difference between the reduce option and the save menu is that the save menu saves the generated reduced syntax to local storage and does not display it. If the user wishes that the *< comments >*, which is necessary for using the automated parser generator (see chapter 3.10.2), production is including in the reduced syntax the user has the option to save the reduced syntax with the *< comment >* production.

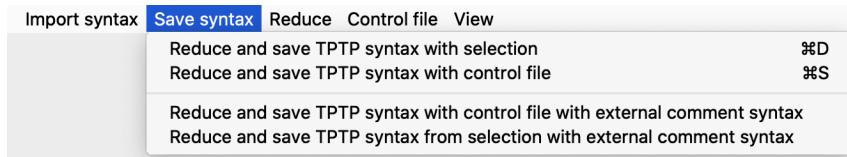


Figure 3.10: Save menu

### 3.11.5 Control file Menu

The *Control file* menu shown in figure 3.11 gives the user the option to import a control file instead of blocking productions manually. If the user imports a control file, the checkboxes are set accordingly. The user also has the option of generating a control file based on his GUI selection. This is particularly helpful if the user uses the same blocked productions multiple times. This way the user can use the import control file option later instead of selecting the blocked productions in the GUI.



Figure 3.11: Control file menu

## 3.12 Command-line interface

In addition to the GUI, [Synplifier](#) offers a command-line interface to provide the means for convenient automation of sub-syntax extraction. It takes a [TPTP](#) syntax file and a control file as input and outputs the resulting sub-syntax. It is implemented using the Python module argparse [22], which provides a framework for creating command-line argument parsers.

The table below provides an overview about the implemented command-line arguments. The syntax file location and the control file location have to be specified. Specifying an output path and filename is optional, by default the output filename will be *output.txt*. The *-ex* flag enables additional output of the comment syntax (see section 3.10.2). Additionally, the help description can be opened by using the *-h* option.

Table 3.2: Command-line interface parameters

Name	Short form	Default	Description
--grammar	-g	None	TPTP syntax file path and filename
--control	-c	None	Control file path and filename
--output	-o	output.txt	Output file path and filename (optional)
--external_comment	-ex	False	Flag to include external comment syntax (optional)

The implementation of the command-line interface is described in section 4.9. More complex actions like control file generation can be done more comfortably by using GUI which is why this is not a feature of the command-line interface.

The design decision to separate the GUI from the command-line interface (see section 3.3.1) has the advantage that if a user chooses to solely use the command-line interface, it is not necessary to install the library used in GUI implementation (PyQt 5).

## 3.13 Counting the number of rules in the TPTP syntax

Rules and productions are counted to compare syntax sizes. The results are described in chapter 5.

### 3.13.1 Counting rules in a grammar graph

Rules are counted by counting nodes in the grammar graph. Productions are counted by traversing the nodes in the grammar graph and counting the number of elements in each productions list.

In addition to the total number of rules and productions, the number of grammar rules and productions are filtered and counted separately. The total number of rules and productions and the number of grammar rules and productions is printed to the command-line after the import of a syntax and after a sub-syntax generation. Since the grammar graph only contains reachable and terminating symbols, the

total number of rules and productions in the original [TPTP](#) syntax version 7.3.0 cannot be counted, since some symbols are not reachable. To compare generated sub-syntaxes to the original [TPTP](#) a script has been developed, which is described in the following.

### 3.13.2 Counting rules in the original [TPTP](#) syntax

To count rules of the original [TPTP](#) syntax (including non-reachable and nonterminating symbols) a script has been developed. It uses the parser to parse the original [TPTP](#) syntax which results in a grammar list. This grammar list is iterated and rules and productions are counted in order to get the total numbers. Also, grammar rules and production are counted separate too.

Counting rules based on the full [TPTP](#) syntax gives a better overview when comparing syntax sizes.

The script can be used independently of [Synplifier](#).

# 4 Implementation

Based on the previous introduction of the concept of Synplifier, this chapter describes the implementation of Synplifier in detail. The chapter covers lexing and parsing the [TPTP](#) syntax as well as building the grammar graph, extracting a sub-syntax and maintaining comments. Furthermore, the implementation of the GUI and the command-line interface is introduced including needed input and output functionalities.

## 4.1 Lexer

As mentioned in chapter [3.4](#) the [PLY](#) package uses the definition of tokens in form of regular expressions to generate the lexer. The following section presents defined tokens that have been introduced in chapter [3.4](#) and the regular expressions that describe these tokens in detail.

### 4.1.1 Tokens

The names of the tokens indicate what they represent. All tokens that have been defined are combined to a token list. A token list is always required when writing a lexer using [PLY](#) [20]. The list of tokens is also used by the parser to identify terminal symbols [20]. Listing [4.1](#) shows the token list of the implemented lexer.

```
1 # List of token names
2 tokens = (
3     'LGRAMMAR_EXPRESSION',
4     'LTOKEN_EXPRESSION',
5     'LSTRICT_EXPRESSION',
6     'LMACRO_EXPRESSION',
7     'NT_SYMBOL',
8     'COMMENT',
9     'OPEN_PARENTHESIS',
10    'CLOSE_PARENTHESIS',
11    'OPEN_SQUARE_BRACKET',
12    'CLOSE_SQUARE_BRACKET',
13    'ALTERNATIVE_SYMBOL',
```

```

14     'REPETITION_SYMBOL' ,
15     'T_SYMBOL'
16 )

```

Listing 4.1: Lexer tokens

Tokens are specified by regular expressions. Regular expressions, that define tokens, have to be defined in a Python function in order for PLY to recognize these tokens. The name of the function is the name of the corresponding token in the token list with the prefix `t_`. Inside a function, it is possible to modify the value of the generated LexToken object and then return it. If not modified, the value of the object is not modified and remains to be the actual text that has been matched. The regular expression is specified as a string in the beginning of a function.

An extract of the functions can be seen in listing 4.2. The function for the token *COMMENT* consists of the regular expression that matches a comment which has already been described in chapter 3.4 and returns the LexToken object unmodified. The functions of *OPEN\_SQUARE\_BRACKET*, *NT\_SYMBOL* and *T\_SYMBOL* (and the other tokens in the tokens list) are implemented similar to *COMMENT*. The function of *LGRAMMAR\_EXPRESSION* and the functions of the other expressions modify the value of the LexToken object. Instead of the value being the actual text matched, first the last three characters are removed. These characters represent the production symbol which indicates the rule type (for an overview of the present rule types see table 3.1). The rule type is already known because depending on the rule type of a production different tokens are matched. Then the Python method `rstrip()` removes all trailing white space characters from the string value. Thus, the function returns a LexToken object with the value of the text belonging to the nonterminal symbol that the expression represents.

```

1 def t_COMMENT(self , t):
2     r'^%.*$'
3     return t
4
5 def t_OPEN_SQUARE_BRACKET(self , t):
6     r'\['
7     return t
8
9 def t_LGRAMMAR_EXPRESSION(self , t):
10    r'<\w+>[\s]*::='
11    t.value = t.value[:-3]
12    t.value = t.value.rstrip()
13
14

```

```

15 def t_NT_SYMBOL(self, t):
16     r'<\w+>'
17     return t
18
19 def t_T_SYMBOL(self, t):
20     r'[$\backslash\\.,a-zA-Z\_0-9\0-9-\>&@!:{}~?^+=/^!^~/@/+-/%;][a-zA-Z\_0-9\0-9\-/"!?/@
21     /+-/*%;\->&+=\$\\\\.,]*'
22     return t
22 )

```

Listing 4.2: Lexer regular expressions

### 4.1.2 Building the lexer

In order to build the lexer, the function `lex.lex()` is used. After the lexer has been created, the function `lex.token()` returns matched tokens. The function `lex.token()` is called until it returns a token of the type `None` that marks the end of the input stream.

In PLY it is possible to declare characters that should be ignored. This means that the characters are ignored in the input stream of the lexer. However, if one of those characters is part of a regular expression of a token it is not ignored and will be used for token matching. In this project tabs, white spaces and newline characters are ignored as they do not have any special meaning other than providing better readability. With exception of the comment token, the information about newline characters is not relevant due to rules being defined over multiple lines, which can be seen in the listing below.

```

1 <annotated_formula> ::=      <thf_annotated> | <tff_annotated> |
2                                <tcf_annotated> | <fof_annotated> |
3                                <cnf_annotated> | <tpi_annotated>

```

Listing 4.3: Example of rules defined over multiple lines

The variable `t_ignore` shown in listing 4.4 is used for declaring characters that should be ignored in the input stream.

```

1 t_ignore = '\t\n '

```

Listing 4.4: Lexer ignore function

The function `t_error` is called when an error occurs during lexing. The defined error function shown in listing 4.5 prints the value of the illegal character and skips the character.

```
1 def t_error(self, t):
2     print("Illegal character '%s'" % t.value[0])
3     t.lexer.skip(1)
4 )
```

Listing 4.5: Lexer error function

## 4.2 Parser

The parser takes the matched tokens from the lexer and constructs a grammar list based on a specified context free grammar (see chapter 3.5). The rules of this grammar are defined by Python functions that are described in the following paragraph. The left-hand side nonterminal symbols of the rules are implemented using data structures that are introduced in section 4.2.1. After that, in section 4.2.2 it is explained how to build the parser with PLY.

The docstrings of the Python functions contain a specification of the grammar's rules. Functions accept one argument which is a list of grammar symbols and the bodies of the functions contain semantic actions. An example for a defined function can be seen in listing 4.6. This function defines the rule that has the nonterminal symbol *productions\_list* on the left side. The docstring of the function contains the specification of the rule and the code inside the function determines the actions that are taken once the function is called. If the input has the length 2, meaning the input list consists of two entries, a new *PRODUCTION\_LIST* is created with the second input element and is saved in the first input element. If the length of the input equals four, a *productions\_list* already exists and is expanded by a new production alternative. This is done by appending the fourth input element to the existing list that is stored in the second input element. The result is again being stored in the first input element.

```
1 def p_productions_list(self, p):
2     """
3         productions_list : production
4                 | productions_list ALTERNATIVE_SYMBOL production
5         """
6     pass
```

```

6  if len(p) == 2:
7      p[0] = PRODUCTIONS_LIST([p[1]])
8  elif len(p) == 4:
9      p[1].list.append(p[3])
10     p[0] = p[1]

```

Listing 4.6: Productions list parser function

### 4.2.1 Nonterminal symbols

The nonterminal symbols of the defined context-free grammar, that the parser is based on, are implemented using multiple data structures. These data structures are described in the following.

A *T\_SYMBOL\_PRODUCTION* consists of the token *T\_SYMBOL* (see chapter 4.1.1) and additionally a production property. The production property can take two values and denotes whether a *T\_SYMBOL* is embedded in square brackets or not. In the production rules a *T\_SYMBOL* token, either embedded in square brackets or without brackets, can be derived from a *T\_SYMBOL\_PRODUCTION*. Asterisk and vertical bar tokens embedded in square brackets can be derived from a *T\_SYMBOL\_PRODUCTION* production as well. The *T\_SYMBOL\_PRODUCTION* is necessary because usually square brackets denote that the term within the brackets is optional which is not the case with terminal symbols. Square brackets embedding terminal symbols denote that the symbol within the brackets is terminal. For example, an asterisk embedded in square brackets refers to the terminal symbol. An asterisk not embedded in square bracket corresponds to the production property *optional* (see figure 4.4).

A *RULE* consists of an expression token (see chapter 4.1.1), a production list and a position. The position denotes at which position in the TPTP syntax the rule was listed. This information is needed to maintain the original order of the rules when printing the reduced syntax. For each rule type there is a data type. This means that *GRAMMAR*, *TOKEN*, *STRICT* and *MACRO\_RULE* nonterminal symbols are introduced. They inherit from the base class *RULE* which can be seen in figure 4.1. Listing 4.7 contains an example of a line in a TPTP syntax file that is represented by the *GRAMMAR\_RULE* data type. The nonterminal symbol name which is produced is *<tff\_formula>*. The production list consists of two productions, as can be seen in the listing.

```
1 <tff_formula> ::= <tff_logic_formula> | <tff_atom_typing>
```

Listing 4.7: Grammar expression

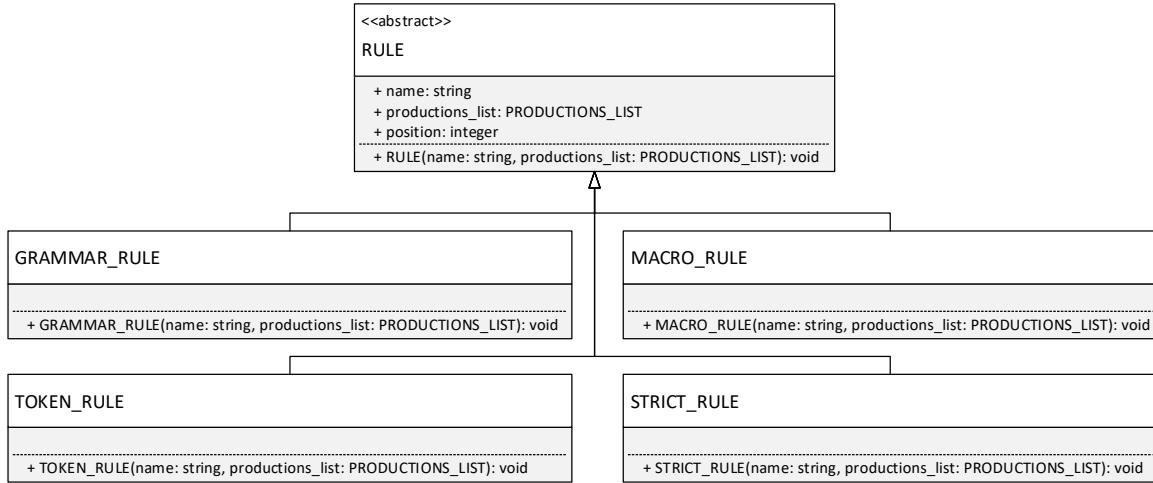


Figure 4.1: Rule data types UML class diagram

A *COMMENT\_BLOCK* is a list of consecutive comment lines. Comment lines can also be extended after initialising a comment block. Figure 4.2 shows the UML denotation of a comment block.

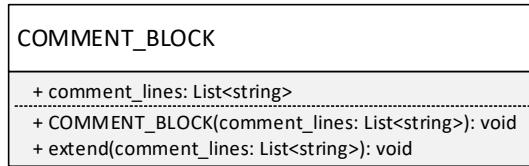


Figure 4.2: Comment block data type UML class diagram

A *T\_SYMBOL\_PRODUCTION* or a *NT\_SYMBOL* is a *PRODUCTION\_ELEMENT*. Therefore, it has an instance of *SYMBOL* as an attribute which can be seen in figure 4.3. Additionally, a production element has a production property that is

shown in figure 4.4. The production property can take one of three values and denotes whether a production is optional, can be repeated multiple times or does not have any special property. In the original TPTP syntax file, the property is represented by square brackets or the repetition symbol.

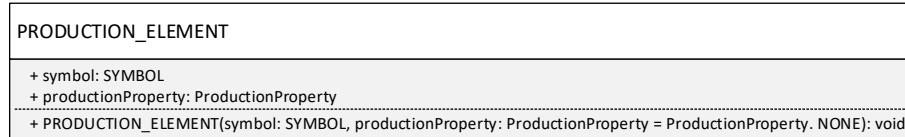


Figure 4.3: Production element data type UML class diagram

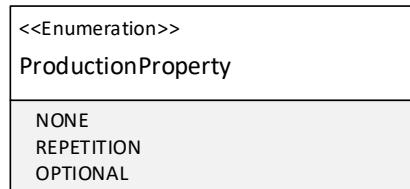


Figure 4.4: Production property data type UML class diagram

A *PRODUCTION* is one of the production alternatives of a rule. As can be seen in figure 4.5 it consists of a list of *PRODUCTION\_ELEMENTS* and/or nested *PRODUCTIONs*. It also has a *ProductionProperty* as attribute.

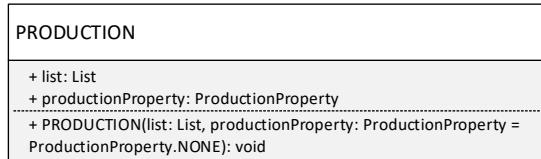


Figure 4.5: Production data type UML class diagram

A *PRODUCTIONS\_LIST* contains a list of productions where each production is one alternative in the description of a rule.

The *XOR\_PRODUCTIONS\_LIST* represents a *PRODUCTIONS\_LIST* with multiple alternatives enclosed by parentheses. It contains a list of the alternate productions. Figure 4.6 shows the UML denotation of a *PRODUCTIONS\_LIST* and a *XOR\_PRODUCTIONS\_LIST* which inherits from the *PRODUCTION\_LIST* class.

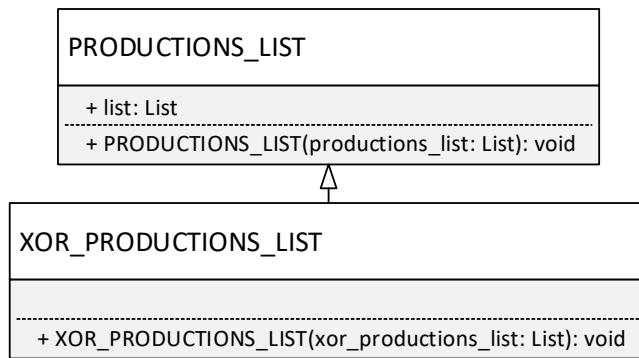


Figure 4.6: Productions list data type UML class diagram

The *GRAMMAR\_LIST* class shown in figure 4.7 is the top-level data structure. It contains a list of all rules that are in the TPTP syntax file. This includes any type of rules (grammar, token, strict and macro) and comment blocks.

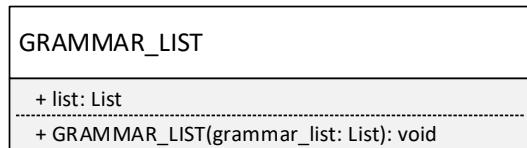


Figure 4.7: Grammar list data type UML class diagram

## Example

Listing 4.8 shows the production rule of the nonterminal symbol `<tfx tuple>` as well as the tokens that have been generated by the lexer. Using this single production, the functionality of the parser should be exemplified.

```

1 <tfx_tuple>      ::= [] | [<tff_arguments>]
2 consists of tokens:
3 l grammar expression open square bracket close square bracket alternative symbol open
    square bracket nt symbol close square bracket

```

Listing 4.8: Production element

The resulting output of the parser can be seen in figure 4.8. The tokens have been reduced to the start symbol grammar list and are displayed in a form of a tree. Internally, the grammar list is stored as a nested list.

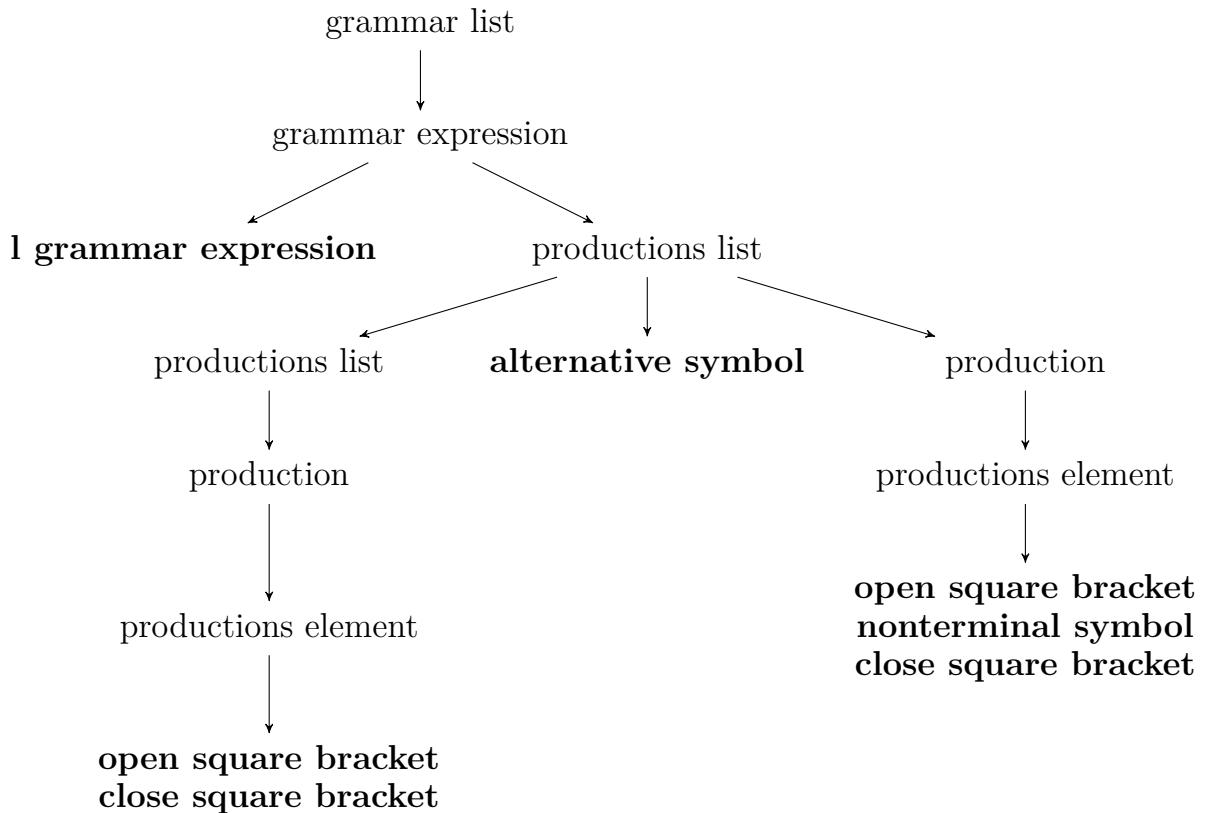


Figure 4.8: Parsing example

### 4.2.2 Building the parser

To build the parser, the **PLY** function `yacc.yacc()` is called that constructs a parser from a specified grammar. To run the generated parser `parse()` must be called. The function that is defined first in the Python file determines the start symbol to which the grammar should be reduced to. The parser stops when no more input exists and the start symbol has been reached or if an error has occurred.

The output of the parser is further modified to include a numbering of rules and to handle square brackets correctly (see chapter 3.5.2). Rules are numbered by assigning the position of a rule within the grammar list as its position. Comment blocks are skipped.

The disambiguation of square brackets replaces the *ProductionProperty Optional* by the terminal symbols open and closing square bracket for *GRAMMAR* and *STRICT* expressions. By doing so, the productions in the productions list of a rule are iterated. If the production has the production property *Optional* the production property is changed to *None* and square brackets are inserted in the production. This is done for productions as well as production elements that have the production property *Optional* inside a production.

## 4.3 Graph generation

The class *TPTPGraphBuilder* that can be seen in figure 4.9 has an instance of the parser that has been introduced in chapter 4.2. With the output of the parser, which is a grammar list, the graph builder builds a nodes dictionary and a grammar graph. This is done by calling the `run` method that calls `init_tree` for building the graph. The methods assigned to a) in figure 4.9 are used for building a graph and are introduced in the following. The graph builder is also responsible for reducing the **TPTP** syntax, the methods used for that are assigned to b) and introduced in chapter 4.5. Moreover, with the graph builder, comments are assigned to rules and nodes. Belonging methods are assigned to c) and described in chapter 4.4.

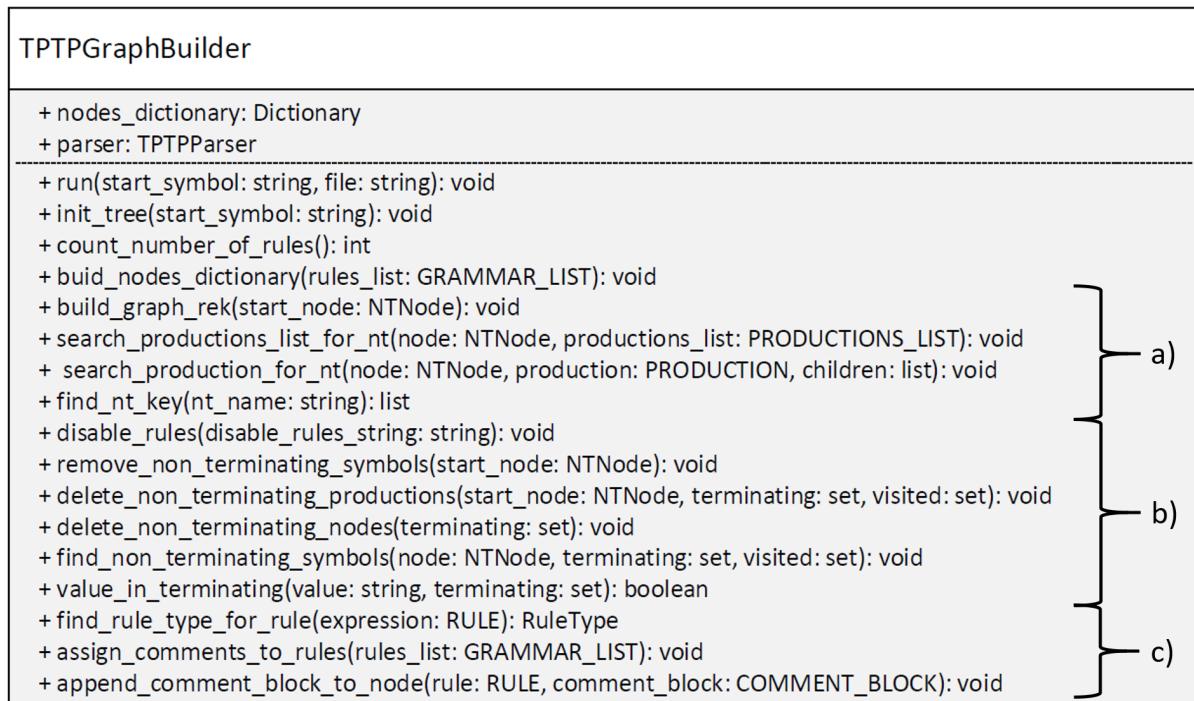


Figure 4.9: Graphbuilder UML diagram

Before generating a graph, a node dictionary consisting of nodes is built.

Figure 4.10 shows the class *NTNode* that has already been introduced in chapter 3.6. The class has three class methods: *Add\_children* adds a list of children to the existing nodes children list. *Extend\_comment\_block* appends a comment block to the existing comment block of a node. The *\_\_str\_\_()* method is used for exporting the grammar.

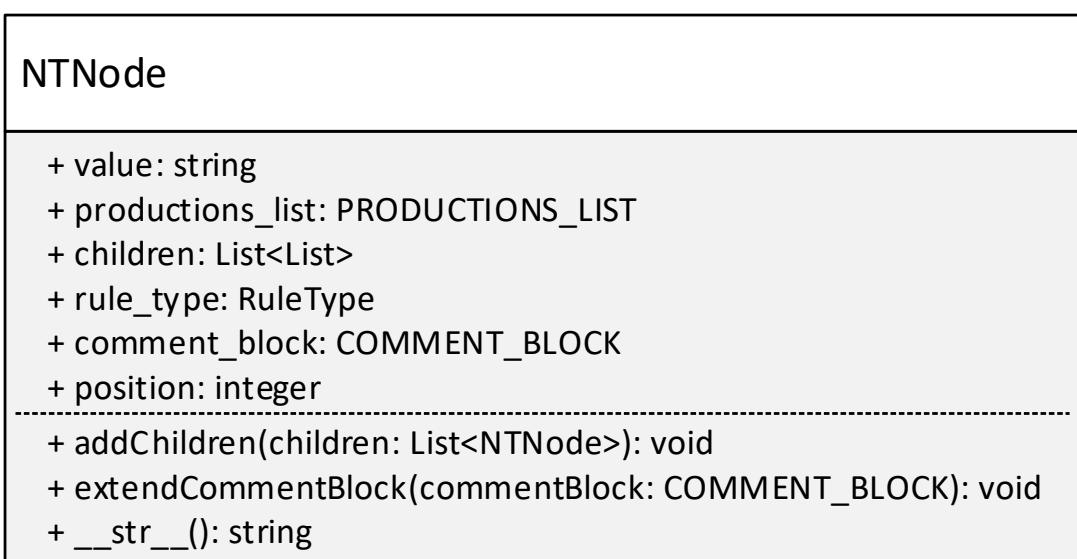


Figure 4.10: NTNode UML diagram

The rule type of a node is either grammar, token, strict or macro and is implemented as an enumeration which can be seen in figure 4.11.

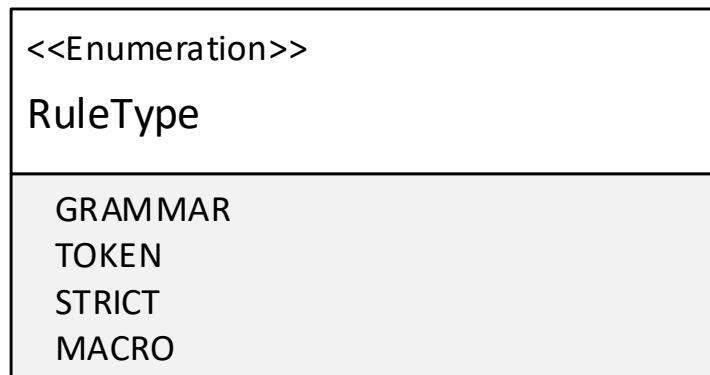


Figure 4.11: RuleType UML diagram

Algorithm 1 shows how the nodes dictionary is built. Iterating the rules list, a new dictionary item is created that consists of the key rule name and rule type and the

*NTNode* as value. This is only the case if the rule is not a comment block. Listing 4.9 shows that multiple rules can exist that correspond to the same node.

```

1 <defined_proposition> ::= <defined_predicate>
2 <defined_proposition> ::= $true | $false
3 <defined_predicate> ::= <atomic_defined_word>
4 <defined_predicate> ::= $distinct |
5                               $less | $lesseq | $greater | $greatereq |
6                               $is_int | $is_rat |
7                               $box_P | $box_i | $box_int | $box |
8                               $dia_P | $dia_i | $dia_int | $dia

```

Listing 4.9: Multiple rules with the same left-hand side of the rule

If that would not have been considered, the second rule would create a new node and replace the existing one in the dictionary because their dictionary keys are the same. To not overwrite the node, the algorithm checks whether the key already exists in the dictionary. If the key already exists, the productions lists of the two rules are merged into one productions list that is appended to the node. After iterating, comments are assigned (see section 4.4).

---

### Algorithm 1 Build nodes dictionary

---

**Input:** rules\_list

```

1: for all rule in rules_list do
2:   if rule is not a comment block then
3:     rule_type = this.find_rule_type_for_rule(rule)
4:     node_key = Node_Key(rule.name, rule_type)
5:     if node_key is not in this.nodes_dictionary then
6:       node = NTNode(rule.name, rule.productions_list, rule_type, None,
    rule.position)
7:       this.nodes_dictionary.update(node_key:node)
8:     else
9:       node = this.nodes_dictionary[node_key]
10:      node.productions_list.list.extend(rule.productions_list.list)
11:    end if
12:  end if
13: end for
14: Assign comments to rules

```

---

Starting with the temporary start symbol, the graph is generated recursively. Iterating over each nonterminal symbol in the productions list of the start symbol, the corresponding nodes are identified. These nodes are then appended to the list of children of the start symbol. If a nonterminal symbol corresponds to two nodes with different rule types, the two nodes are stored in a list and are appended to the list of children resulting in a nested list. The identified children may again have children. This process is repeated until a node has no children because there are only terminal symbols in the productions list of a nonterminal symbol.

Since it is possible for a nonterminal symbol to be on the right side as well as on the left side of the same production rule, a node can also be its own child. To avoid revisiting the same node infinitely, it is checked whether a node already has children so that it will not be visited again. This also improves the performance of the tool as a nonterminal symbol that has already been visited will not be visited again independent of circular dependencies.

To generate the graph of a given grammar three algorithms are needed that will be explained in the following.

Algorithm 2 calls function 3 that appends children of a node to the nodes list of children. The algorithm is first called with the start symbol. After the children of a node have been appended to the node, the algorithm is called for every child resulting in appending their children to their children's list.

---

**Algorithm 2** Build graph recursively
 

---

**Input:** node

```

1: searchProductionsListForNT(node, node.productionsList)
2: if node has children then
3:   for all children do
4:     buildGraphRek(child)
5:   end for
6: end if
```

---

The right side of a production rule is stored in a productions list. For identifying the nonterminal or terminal symbols in the productions lists, a loop iterates all elements of the productions list. Each element is a production and calls the function 3. This

function identifies the children of the given element which are then appended to the node.

---

**Algorithm 3** Extract productions from productions list

---

**Input:** node, productionsList

```

1: for all elements in productionsList do
2:   children = new empty list
3:   searchProductionForNT(node, element in productionsList, children)
4:   append children to node
5: end for
```

---

For identifying the nonterminal symbols, it is checked if the production is a nested production and if so, the same function is called again. If the production is a XOR production list the function 3 is called to break down the productions list. If the production element is a nonterminal symbol, the element is searched in the node dictionary to get the node where the element is on the left side. This element is then appended to a list of children. It is possible that an element appears multiple times on the left side if it is presented by multiple expressions. In this case each element is appended to the list of children.

---

**Algorithm 4** Append children to node

---

**Input:** node, productionsList, children

```

1: for all elements in productionsList do
2:   if element is a production then
3:     searchProductionForNT(node, element, children)
4:   else if element is a XOR productions list then
5:     searchProductionsListForNT
6:   else if element is a nonterminal symbol then
7:     find element(s) in node dictionary
8:     append element(s) to children
9:   end if
10: end for
```

---

## 4.4 Maintaining comments

After building the nodes dictionary and the grammar graph, comments of the original **TPTP** syntax are identified and assigned to the corresponding nodes. This section describes the implementation of associating comments with nodes in order to maintain comments in a generated sub-syntax. The implementation follows the concept outlined in section 3.8. Comments are associated after the grammar graph has been generated. They are associated to the node, that represents the rule they belong to. For that, the *NTNode* class has a *COMMENT\_BLOCK* attribute (see *NTNode* UML diagram in figure 4.10).

Associating comments with nodes has the advantage, that when the grammar graph is reduced, only the comments associated with the rules that are part of the sub-syntax are maintained and can also be printed when outputting the sub-syntax.

### 4.4.1 Associating comments with grammar graph nodes

Algorithm 5 describes the procedure of matching comments to nodes. The algorithm follows the concept outlined in section 3.8 and is used within the *TPTPGraphBuilder* after generating the grammar graph.

The input is the *grammar\_list*, that the parser created from parsing the **TPTP** syntax file. The *grammar\_list* contains rules and comment blocks in the order they appeared in the original **TPTP** syntax file. The *grammar\_list* is iterated (line 1). If the current element of the *rules\_list* is a *COMMENT\_BLOCK*, it is split by top of page lines (line 3). The procedure of splitting a comment block by top of page lines is described in section 4.4.2. It returns a list of comment block objects.

In the following if-clause it is checked, how many comment block objects are returned.

In the first case of one comment block, this comment block should be associated with the following rule by default. If that is not possible, because the comment block is at the end of the syntax file, it should be associated with the rule before. In line 5 it is checked whether the comment block is not at the end of the syntax file. To associate the comment block with the following rule, the next rule (rule after the comment block) is obtained by indexing the next element of the *rules\_list*. Then, the comment block is appended to the node, specified by the rule, using the

*append\_comments\_to\_node* method. This method is described in the following subsection.

If the comment block is at the end of the syntax file, it has to be associated with the rule before. This case is handled in the else branch (line 8). The procedure of adding the comment block to the desired node is essentially the same as described before. The difference is that the previous rule is addressed (line 9).

If splitting the original comment block by top of page lines results in two comment blocks the first is associated with the rule before and the second with the rule after on condition that the original comment block was neither at the end nor the beginning of the syntax file (lines 14, 15, 21, 22).

If the original comment block was at the beginning of the syntax file, the first comment block is appended to the rule after too (lines 17, 18).

If the original comment block was at the end of the syntax file, the second comment block is associated with the rule before.

The case of two or more top of page lines does not occur in the [TPTP](#) syntax version 7.3.0. and is not included in the pseudo-code.

This hypothetical case is handled by using the same procedure as in the case of two comment blocks but treating all comment blocks after the second comment block like the second. This means that all comment blocks from the second on would be associated with the rule after, except the case that the original comment block was at the end of the syntax file. In that case all comment blocks would be appended to the rule before.

The algorithm relies on the *rules\_list* never having two or more comment blocks in a row (meaning they are always separated by at least one rule). This is ensured by the parser, because all consecutive comment lines in the syntax file are included into one comment block. Therefore, consecutive comment block objects do not occur in the *rules\_list*.

---

**Algorithm 5** Assign comments to rules

---

**Input:** grammar\_list

```

1: for all index, expression in rules_list.list do
2:   if expression is COMMENT_BLOCK then
3:     comment_block_list = expression.split_comment_block_by_top_of_page()
4:     if length(comment_block_list) == 1 then
5:       if index < length(rules_list.list) - 1 then
6:         next_rule = rules_list.list[index + 1]
7:         this.append_comments_to_node(next_rule, comment_block_list[0])
8:       else
9:         previous_rule = rules_list.list[index - 1]
10:        this.append_comments_to_node(previous_rule, comment_block_list[0])
11:      end if
12:    else if length(comment_block_list) == 2 then
13:      if index != 0 then
14:        previous_rule = rules_list.list[index - 1]
15:        this.append_comments_to_node(previous_rule, comment_block_list[0])
16:      else
17:        next_rule = rules_list.list[index + 1]
18:        this.append_comments_to_node(next_rule, comment_block_list[0])
19:      end if
20:      if index < length(rules_list.list) - 1 then
21:        next_rule = rules_list.list[index + 1]
22:        this.append_comments_to_node(next_rule, comment_block_list[1])
23:      else
24:        previous_rule = rules_list.list[index - 1]
25:        this.append_comments_to_node(previous_rule, comment_block_list[1])
26:      end if
27:    else if length(comment_block_list) > 2 then
28:      %excluded
29:    end if
30:  end if
31: end for

```

---

### Appending a comment block to a node

In Algorithm 6 it is described, how a comment block is appended to a node from the grammar graph, using information from the original rule created by the parser to address the node.

To associate the *COMMENT\_BLOCK* with the node representing the rule, the node needs to be accessed. It is more efficient to use the created *nodes\_dictionary* to access the desired nodes than to traverse the grammar graph for each wanted node.

In the *nodes\_dictionary*, the key for each node consists of the nonterminal symbol name and the *RuleType* is stored (see section 4.3).

The rule type of the input rule is indicated by the type of the *rule* object (see class diagram in figure 4.11). To index the node in the *nodes\_dictionary* the *RuleType* corresponding to the type of the rule object is needed. The *find\_rule\_type\_for\_rule* method determines and returns the *RuleType* corresponding to the *rule* object type. Using the determined *RuleType* and the nonterminal name, which is stored in the *rule*, the key for the *NodesDictionary* can be constructed (line 8). Then, the node, to which the *COMMENT\_BLOCK* should be added is addressed in the *NodesDictionary* and the *extend\_comment\_block* method of the node called, with the *COMMENT\_BLOCK* as argument. The *extend\_comment\_block* method is described in algorithm 7. It appends the *COMMENT\_BLOCK* to the *COMMENT\_BLOCK* attribute of the node object.

---

**Algorithm 6** Append comments to node
 

---

**Input:** rule, comment\_block

- 1: rule\_type = find\_rule\_type\_for\_rule(rule)
  - 2: node\_key = Node\_Key(rule.name,rule\_type)
  - 3: this.nodes\_dictionary[node\_key].extend\_comment\_block(comment\_block)
- 

### Extending a comment block of a node

Algorithm 7 shows the procedure that is used in the *extend\_comment\_block* method of the *NTNode* class. If no *comment\_block* has not been associated with the node before, the input *comment\_block* is associated with the node. If the node already has a *comment\_block*, this *comment\_block* is extended, with the comments lines from the input *comment\_block*.

This procedure is needed because it is possible that multiple comment blocks, that have been generated on splitting comment blocks, have to be associated with the same node.

**Algorithm 7** Extend node comment block

---

**Input:** comment\_block

```
1: if this.comment_block is None then,  
2:   this.comment_block = comment_block  
3: else  
4:   this.comment_block.extend(comment_block.comment_lines)  
5: end if
```

---

#### 4.4.2 Splitting comment blocks

The flow chart in figure 4.12 shows how a comment block is split into multiple comment blocks separating it by top of page lines.

First the indexes of the top of page lines are collected in a list (using method of *COMMENT\_BLOCK*). After that, the *comment\_block\_list*, which will contain the generated *COMMENT\_BLOCKs*, is initialized. If the list, storing top of page line indexes, is empty, the original comment block is added to the *comment\_block\_list*, because the comment block does not need to be split into multiple comment blocks because no top of page line is present. Then the *comment\_block\_list* is returned. If the list of top of page line indexes is not empty, the first top of page line index is retrieved from the top of page indexes list. It is also checked whether the first top of page index is unequal zero.

If that is the case, the first top of page index is not the first line of the original comment block, and a comment block can be constructed from the line(s) before. In this case the first comment block can be constructed from the comment lines of the original comment block from index zero to the first top of page index. This comment block is added to the *comment\_block\_list*.

Then the *index* variable, which represents an index on the *top\_of\_page\_indexes* list, is initialized. The *index* is used to access the corresponding element of the *top\_of\_page\_indexes* list.

Then it is checked if top of page line that corresponds to the top of page line index is the last line of the original comment block. In that case, no further comment block needs to be added to the *comment\_block\_list* and the *comment\_block\_list* is returned.

Otherwise the start index of the next comment block is computed as the current *top\_of\_page\_index* incremented by one (line after the top of page line). There are two cases that have to be considered when determining the end index of the

comment block. Either there is another top of page line, in this case the comment block ends before that, or the currently indexed top of page line was the last one, in that case the end index is the end of the original comment block. Then the new command block is constructed with the before determined *start* and *end*. This comment block is added to the *comment\_block\_list*. If the *index* is at the end of the *top\_of\_page\_index* list, the *comment\_block\_list* is returned. If not, the *index* variable is incremented and the process is repeated with the next top of page index.

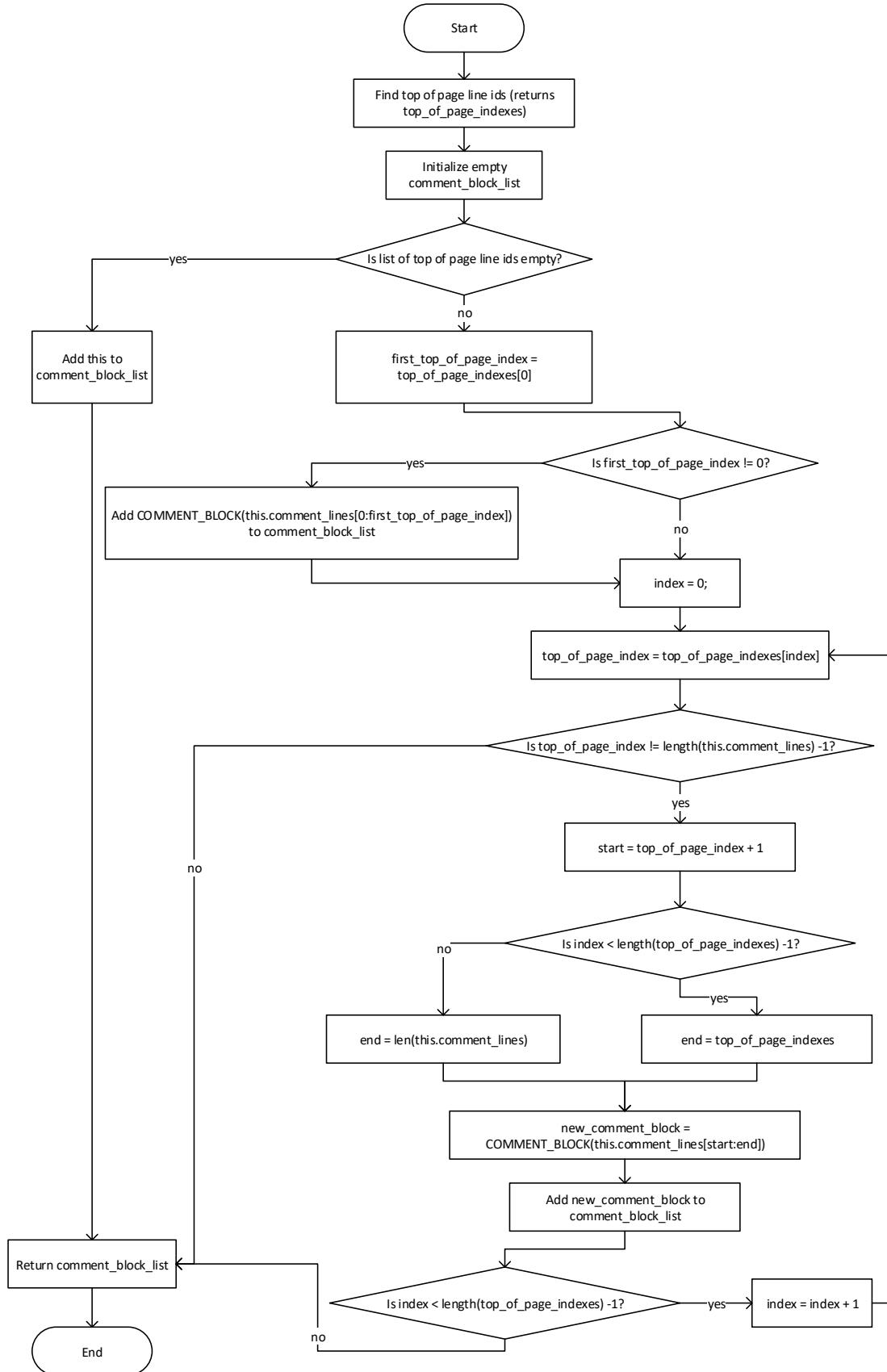


Figure 4.12: Split comment block by top of page line flow chart

## 4.5 Extraction of a sub-syntax

The following sections describe the algorithms used for the extraction of a sub-syntax based on the concept outlined in section 3.9.

### 4.5.1 Removal of blocked productions

Algorithm 8 takes the text of the control file as an input argument (control file format description can be found in section 3.7) and removes the productions that should be blocked from the nodes. The control file string is split into a list of the lines (line 1). The first line of the control file string, which describes the start symbol, is deleted from the list (line 2) because this information is not needed at this step.

Then the following lines are iterated:

For each line the content is split by the comma symbol, which produces a list of strings. The first element of that list is the nonterminal symbol name of the symbol from which productions should be disabled. The second element is the corresponding rule type symbol. The nonterminal symbol name and the rule type, that is determined from the rules symbol (line 7), are stored in separate variables. Then these elements are removed from the list and the remaining elements are converted to integers. They are the indexes of the productions that should be removed. The elements are ordered descending (line 10). Ordering the indexes in that order has the advantage that, when the list of indexes is iterated to delete productions, the indexes of other productions that should be deleted are not changed by the deletion of productions before. To delete the productions from the specified node the node object has to be accessed. This is done by addressing the node in the *nodes\_dictionary* with its key consisting of the nonterminal symbol name and the rule type. From the node, in a loop, the elements of the *productions\_list* and the *children* list corresponding to the indexes in the control file are deleted.

---

**Algorithm 8** Removing blocked productions

---

**Input:** control\_string

```

1: lines = control_string.splitlines()
2: delete lines[0]
3: for all line in lines do
4:   data = line.splitBy(",")
5:   nonterminal_name = data[0]
6:   rule_symbol = data[1]
7:   rule_type = determineRuleType(rule_symbol)
8:   delete data[0:2]
9:   data = parseInteger(data)
10:  data.sortReverse()
11:  node    = this.nodes_dictionary.get(Node_Key(nonterminal_name,
      rule_type))
12:  for all index in data do
13:    delete node.productions_list.list[index]
14:    delete node.children[index]
15:  end for
16: end for
```

---

#### 4.5.2 Determination of the remaining terminating symbols

After the desired productions have been deleted from the grammar graph, the next step is to remove the nonterminating symbols from the grammar graph. Algorithm 9 shows the procedure of removing nonterminating symbols. First it has to be determined which symbols are terminating and which are nonterminating.

On the one hand, starting from the terminal symbols all nonterminal symbols that derive terminal symbols could be found. However, the graph data structure can only be traversed top down not bottom up. Therefore, starting at the start node, the graph is traversed and terminating symbols are found. A set of terminating nodes and a temporary set of terminating nodes are initialized. In a while loop the recursive algorithm *find\_non\_terminating\_symbols(start\_node, temp\_terminating, visited)* (algorithm 10) finds terminating nodes with the start\_node, an initialized set of known terminating symbols and a set of already visited nodes. This algorithm is called with the updated set of terminating symbols until the set does not differ from the run before. When that is the case, all terminating nodes have been found. After all terminating nodes have been found, the productions that contain nonter-

minating symbols are deleted from the nodes and after that nodes representing nonterminating symbols are removed.

In the following parts of the section, first the *find\_non\_terminating\_symbols* algorithm is described. Then the *delete\_non\_terminating\_productions* algorithm and after that the *delete\_non\_terminating\_nodes* algorithm are outlined.

---

**Algorithm 9** Removing nonterminating symbols
 

---

**Input:** start\_node

```

1: terminating = new set()
2: temp_terminating = new set()
3: while True do
4:   visited = new set()
5:   this.find_non_terminating_symbols(start_node, temp_terminating, visited)
6:   if terminating == temp_terminating then
7:     break
8:   else
9:     terminating = temp_terminating
10:  end if
11: end while
12: visited = new set()
13: delete_non_terminating_productions(start_node, terminating, visited)
14: delete_non_terminating_nodes(terminating)
```

---

Algorithm 10 recursively determines whether a node can be identified as terminating. The input is a set of already known terminating nodes and nodes that have already been visited. If the node has been not been visited yet, it is added to the visited set. If it has been visited, the method ends to prevent infinite recursion. If it has not been visited yet, the *children* list of the node is iterated for each children list in *node.children*. If a *children\_list* is empty, this means that in a production there are no nonterminal symbols and the production only consists of terminal symbols. If that is the case the node is terminating.

If it is not the case, the node is terminating if all nodes in the *children\_list* represent a terminating nonterminal symbol. To check that, the algorithm is called recursively for all nodes in the *children\_list*. If a nonterminal symbol has multiple rule types, it is represented by multiple nodes. Also, if a nonterminal symbol with multiple rule

types is featured in a production, all nodes representing that nonterminal symbol are in the children list corresponding to this production. Only one of the nodes needs to be terminating in order for the nonterminal symbol to be terminating. This is considered in the procedure of checking if a node represents a terminating symbol (line 11), which is described in figure 4.13. If every node in the *children\_list* represents a terminating symbol, the node is added to the known terminating nodes.

---

**Algorithm 10** Find nonterminating symbols
 

---

**Input:** node, terminating\_set, visited\_set

```

1: node_key = Node_Key(node.value, node.rule_type)
2: if node_key not in visited_set then
3:   visited_set.add(node_key)
4:   for all children_list in node.children do
5:     if len(children_list) == 0 then
6:       terminating_set.add(node_key)
7:     else
8:       terminating_flag = True
9:       for all child in children_list do
10:        find_non_terminating_symbols(child, terminating_set, visited_set)
11:        if not value_in_terminating(child.value, terminating_set) then
12:          terminating_flag = False
13:        end if
14:      end for
15:      if terminating_flag then
16:        terminating_set.add(node_key)
17:      end if
18:    end if
19:  end for
20: end if
```

---

The flow diagram in figure 4.13 shows the procedure of checking whether a non-terminal symbol, described by its name, is in the set that contains information on which nodes are known to be terminating. Theoretically there can be up to four nodes representing a nonterminal symbol, because it can have up to four rule type (and a node represents the combination of nonterminal symbol name and rule type). However, according to the TPTP syntax only the rule types grammar and strict

can belong to one nonterminal symbol.

In the *terminating\_set* the *NodeKeys* of known terminating nodes are stored. A nonterminal symbol is known to be terminating when any node representing that nonterminal symbol is known to be terminating. If a *NodeKey* of any rule type and with the nonterminal name is in the *terminating\_set* then the symbol is terminating and True is returned, otherwise False is returned.

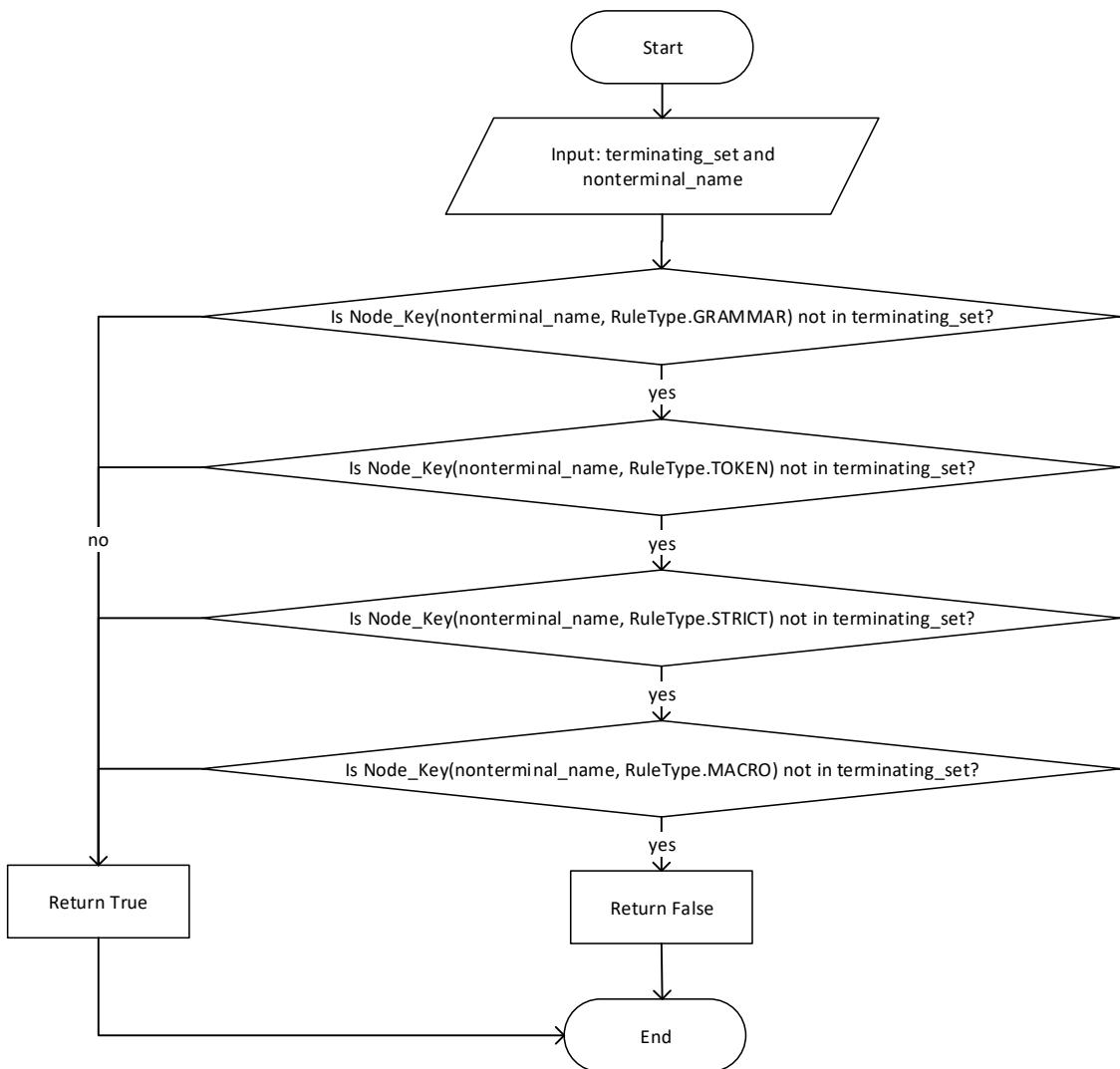


Figure 4.13: Check if node value is in terminating set

After the terminating symbols have been determined, productions that contain a nonterminating symbol have to be removed from the nodes of the grammar graph.

Only productions and not the node itself is deleted because a node might also have productions that only contain terminating symbol. Therefore, only the productions containing nonterminating symbols have to be removed.

Algorithm 11 describes how productions that include nonterminating symbols are removed from the grammar graph. The start node of the graph, the set of known terminating nodes and an empty set of already visited nodes are input at the first call of the algorithm. If the input node has been visited before, it is not visited again and the algorithm ends. If it has not been visited before the corresponding *Node\_Key* is added to the set of visited nodes. Then the *children* list of the node is iterated (line 5). The index created in line 4 corresponds to the current element of the *children* list, that is iterated. It is needed to delete the *children* list and the corresponding *productions\_list*, if a production contains nonterminating symbols. The *children* list is iterated in reverse, because, if a *children* list and *productions\_list* is deleted this will not affect the indexing of the next *children* list to analyse. The *children* list contains lists of nodes for each production (see section 3.6). Each list, nested in the *children* list, is iterated and for each node the algorithm is called recursively (line 7 and 8).

If a node represents a nonterminating symbol the *not\_terminating* flag is set to true. If this flag is true after iterating one element of the *children* list, it means that the production corresponding to the *children* list element contains at least one nonterminating symbol. If that is the case, the corresponding element of the *children* list and of the *productions\_list* is removed from the node. This is repeated for all elements of the *children* list.

As a result of this algorithm, all productions that contained nonterminal symbols have been removed from the grammar graph.

---

**Algorithm 11** Delete productions containing nonterminating symbols

---

**Input:** node, terminating\_set, visited\_set

```

1: node_key = Node_Key(node.value, node.rule_type)
2: if node_key not in visited_set then
3:   visited_set.add(node_key)
4:   index = len(node.children) - 1
5:   for all children_list in reversed(node.children) do
6:     not_terminating = False
7:     for all child in children_list do
8:       delete_non_terminating_productions(child, terminating_set, vis-
    ited_set)
9:       if not value_in_terminating(child.value, terminating_set) then
10:        not_terminating = true
11:      end if
12:    end for
13:    if not_terminating then
14:      delete node.children[index]
15:      delete node.productions_list.list[index]
16:    end if
17:    index = index - 1
18:  end for
19: end if

```

---

After the productions containing nonterminating symbols have been removed from the grammar graph, the nodes that represent nonterminating symbols can also be removed from the *nodes\_dictionary*. Algorithm 12 is responsible for this. It takes the set of *Node\_Keys* as input, that specifies all terminating nodes. To remove the nonterminating nodes, in a for-loop, all terminating nodes that are addressed by the *Node\_Keys* are stored in a temporary dictionary. This dictionary replaces the original *nodes\_dictionary* of the *TPTPGraphBuilder*.

---

**Algorithm 12** Delete nonterminating nodes

---

**Input:** terminating\_set

```

1: temporary_dictionary = {}
2: for all node_key in terminating_set do
3:   value = this.nodes_dictionary.get(node_key, None)
4:   temporary_dictionary.update(node_key: value)
5: end for
6: this.nodes_dictionary = temporary_dictionary

```

---

### 4.5.3 Determination of the remaining reachable symbols

Now, that nonterminating symbols have been removed from the grammar graph, the remaining reachable symbols have to be determined. As described in section 3.9.2 the determination of the remaining reachable is done by creating a new grammar graph, with the desired start symbol from the control file. The new grammar graph will only contain the reachable symbols (as described in section 3.9.4).

## 4.6 Input

There are two input operations that are needed for [Synplifier](#). One is reading the content of text files which is needed for importing [TPTP](#) syntax files and also control files. The other input operation is the import of the current full [TPTP](#) syntax from the [TPTP](#) website. This increases the convenience for users because no local copy of the original [TPTP](#) syntax is needed for sub-syntax extraction. The two methods representing these operations are described in the following.

### Importing the [TPTP](#) syntax from the [TPTP](#) website

The method *import\_tptp\_syntax\_from\_web* shown in listing 4.10 is used to get a text file of the [TPTP](#) syntax from the [TPTP](#) website.

For opening a Uniform Resource Locator ([URL](#)), the library `urllib` is used.

`Urllib.request` opens a specified [URL](#) and sends an HTTP request. `Urlopen` returns a response object for the initiated request. The response is a file-like object that can be read using `read()`. [28]

The file can then be converted from html to text format using BeautifulSoup's html parser. The header from the file is deleted because it does not belong to the [TPTP](#) syntax but is appended by `urllib`. The header is deleted by splitting the text file of the grammar at every new line and joining all lines excluding the first one.

```

1 with urllib.request.urlopen("http://www.tptp.org/TPTP/SyntaxBNF.html") as url:
2     html_doc = url.read()
3 soup = BeautifulSoup(html_doc, 'html.parser')
4 tptp_grammar = soup.get_text()
5 # Delete Header
6 tptp_grammar = '\n'.join(tptp_grammar.split('\n')[1:])

```

Listing 4.10: Import TPTP syntax from the TPTP website

### Reading TPTP syntax file and control file

*Read\_text\_from\_file* in listing 4.11 uses the build-in function *open* to read a text file. The filename of the text file has to be specified before.

```

1 with open(filename, 'r') as text_file:
2     text = text_file.read()

```

Listing 4.11: Read text from file

## 4.7 Output generation

For saving the grammar graph in a text file, the graph needs to be traversed and converted to a string format that can be saved in a text document. The grammar graph is converted to a string format by converting every rule to a string and cohering them. Every rule is represented by a node including the nodes children list. For converting the graph to a string format, the graph is traversed and string representation of the nodes are created. The final output is stored in a *print\_list*. A *print\_list* is a list of *print\_list\_entries* that can be seen in listing 4.12. A *print\_list\_entry* is a tuple consisting of the position of the rule that in the original TPTP syntax is represented by the node and a string version of the rule or rather the node and its children list.

```

1 PrintListEntry = namedtuple("PrintListEntry", ["position", "string"])

```

Listing 4.12: Print list entry

The detailed process of saving a string version of a rule is described in the following sections.

### 4.7.1 Saving string representation of graph in text file

In order to create the *print\_list*, an empty set *visited* is created that is used for memorising which nodes have already been visited. Moreover, an empty *print\_list* and *print\_string* are created. The set *visited* and the empty *print\_list* are used for calling the function *create\_print\_list* (see following section) which returns the final *print\_list*. Because the rules in the *print\_list* are not in the same order as the rules in the original **TPTP** syntax, the *print\_list* is sorted based on the rule's position that is stored in the first element of the *PrintListEntry* tuple. For converting the *print\_list* to one string, the *print\_string*, the *print\_list* is iterated. Except for the first element of the list, every element of the list is added to the *print\_string*. This is done by adding the strings that are stored in the second element of the *PrintListEntry* tuple and separating them by a new line. The first element is not added because it represents the general start symbol that has been added to the grammar (see chapter 3.6) but is not a part of the **TPTP** syntax. A text document with a filename that has been selected by the user is then opened and the *print\_string* is written to the file. The simplified process can be seen in the flow diagram in figure 4.14.

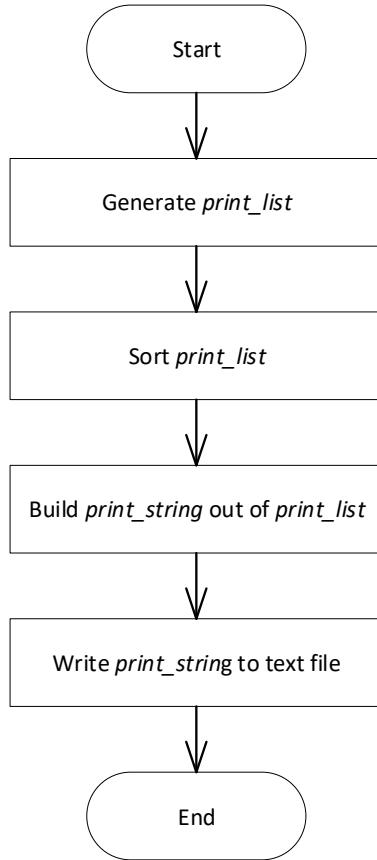


Figure 4.14: Save ordered rules from grammar graph

#### 4.7.2 Traverse graph to create string representations of objects

The graph is traversed using function 13. For every node a new *print\_list\_entry* is created. The str function that is called to create the *print\_list\_entry* is described in section 4.7.3. The *print\_list\_entry* is then appended to the *print\_list*. Once visited a node is added to the set *visited* so that it will not be visited again. Recursively all children in the node's children list are traversed. Because a child can once again have a list, every element of the child's list is considered and calls *create\_print\_list* to complete the final *print\_list*. The visited elements are then also added to visited set. The function ends with returning his version of the *print\_list* that will be complete list once the first node who called the function return its *print\_list*.

**Algorithm 13** Create print list

---

**Input:** node, visited, print\_list

```
1: Create PrintListEntry with node.position and str(node)
2: Append PrintListEntry to print_list
3: Add node to visited
4: for all children in node.children do
5:   for all elements in children do
6:     if element is not in visited then
7:       print_list = create_print_list(element, visited, print_list)
8:       Add element to visited
9:     end if
10:   end for
11: end for
12: return print_list
```

---

### 4.7.3 Create string representations of objects

This section covers the string representation creation for the data types that are used within the grammar graph. For creating the string representation, the Python string class is used [26]. The string class creates a string version of the object passed to the constructor. To do that it calls the `__str__()` of the passed object, if the object provides that method.

The algorithms outlined in the following are embedded in the `__str__()` method of each class.

The goal is for the generated string representations to be in the same format as the original TPTP syntax.

#### **NTNode string representation creation**

An *NTNode* object represents a rule statement. This may include a comment block associated with that rule statement. The UML diagram of the *NTNode* class can be seen in figure 4.10. The algorithm that describes how a string representation of an *NTNode* object is generated is featured in algorithm 14. If a comment block is associated with an *NTNode*, its string representation is added to the *NTNode* string representation at the beginning. Then the value of the nonterminal symbol (left hand side of the production rule the *NTNode* represents) is added

to the string representation. After that the corresponding production symbol corresponding to the *RuleType* (see table 3.1) of the *NTNode* is added to the string representation separated to the nonterminal symbol by a space character. To align the productions in the output file, a whitespace padding is added. Next, the string representation of the *productions\_list* is added to the string representation of the *NTNode*. The procedure of how the string representation of the *productions\_list* is created is outlined in the following after the description of the comment block string generation.

---

**Algorithm 14** *NTNode* string creation
 

---

```

1: if this.comment_block is not None then
2:   node_string += str(this.comment_block)
3: end if
4: node_string += this.value
5: add whitespace padding
6: if this.rule_type == RuleType.GRAMMAR then
7:   node_string += " ::= "
8: else if this.rule_type == RuleType.TOKEN then
9:   node_string += " :- "
10: else if this.rule_type == RuleType.STRICT then
11:   node_string += " === "
12: else if this.rule_type == RuleType.MACRO then
13:   node_string += " :: "
14: end if
15: node_string += str(this.productions_list) return node_string
  
```

---

***COMMENT\_BLOCK* string representation creation**

The *COMMENT\_BLOCK* class has a list of comment lines as an attribute (see UML diagram in figure 4.2). The string representation of the *COMMENT\_BLOCK* objects is created by concatenating the comment lines, separating the individual comment lines by newline characters.

### ***PRODUCTION\_LIST* string representation creation**

The *PRODUCTION\_LIST* class contains a list of production alternatives of a rule (the right-hand side of a rule). The production alternatives are separated by " | " in the original TPTP syntax, which can be seen in the following listing.

```
1 <annotated_formula> ::= <thf_annotated> | <tff_annotated> | <tcf_annotated> |
2                                <fof_annotated> | <cnf_annotated> | <tpi_annotated>
```

Listing 4.13: Production alternatives in the TPTP syntax

Algorithm 15 shows the process of creating the string representation of a *PRODUCTION\_LIST*. The productions in *PRODUCTION\_LIST* are iterated and the string representation of every production is created. The strings of the productions are then separated by the " | " character. The counter in line 3 of algorithm 15 is used to identify the last production and thus not to print alternative symbol after that production.

---

#### **Algorithm 15 *PRODUCTION\_LIST* string creation**

---

```
1: productions_list_string = ""
2: length = length(this.list)
3: counter = 1
4: for all production in this.list do
5:     productions_list_string += str(production)
6:     if counter < length then
7:         productions_list_string += " | "
8:     end if
9:     counter += 1
10: end forreturn productions_list_string
```

---

### ***PRODUCTION* string representation creation**

The *PRODUCTION* class has a list of elements and a *ProductionProperty* as an attribute (see UML diagram in figure 4.5). The elements of the list can be of type *PRODUCTION*, *PRODUCTION\_ELEMENT* or *XOR\_PRODUCTION\_LIST*. Algorithm 16 describes the procedure of creating the string representation of a *PRODUCTION*. First, the *production\_string* is initialised. Then, depending on

the *ProductionProperty* characters that represent these productions properties are added to the *production\_string*. If the *ProductionProperty* is *NONE* no special characters need to be added. In case of the *ProductionProperty REPETITION*, the open parenthesis character "(" is added to the *production\_string* in order to enclose the *list* elements. In case of the *ProductionProperty OPTIONAL*, the elements have to be enclosed in square brackets. Therefore the "[" character is added to the *production\_string* in that case.

After that, the *list* is iterated and the string representation of each element is created and added to the *production\_string*.

Again depending on the *ProductionProperty* additional characters might have to be added. In case of the *ProductionProperty REPETITION* the open parenthesis is closed with the close parenthesis and the repetition symbol is added to denote the repetition property. In case of the *ProductionProperty OPTIONAL* the open square bracket is closed. Then the created *production\_string* is returned.

---

**Algorithm 16** *PRODUCTION* string creation
 

---

```

1: production_string = ""
2: if self.productionProperty == ProductionProperty.REPETITION then
3:   production_string += "("
4: else if element.productionProperty == ProductionProperty.OPTIONAL then
5:   production_string += "["
6: end if
7: for all element in this.list do
8:   production_string += str(element)
9: end for
10: if self.productionProperty == ProductionProperty.REPETITION then
11:   production_string += ")**"
12: else if element.productionProperty == ProductionProperty.OPTIONAL then
13:   production_string += "]"
14: end if
15: Return production_string
  
```

---

***XOR\_PRODUCTIONS\_LIST* string representation creation**

The procedure to create the string representation of the *XOR\_PRODUCTIONS\_LIST* class is almost equivalent to the *PRODUCTIONS\_LIST* process described in algo-

rithm 15. The only difference is, that generated strings are enclosed by parentheses.

### ***PRODUCTION\_ELEMENT* string representation creation**

The *PRODUCTION\_ELEMENT* data type consists of a *SYMBOL* and a *ProductionProperty*, which can be seen in the UML diagram in section 4.3. The procedure of creating the string representation of a *PRODUCTION\_ELEMENT* is outlined in the following algorithm 17. The string representation is created similar to the creation of the string representation of a *PRODUCTION*.

---

#### **Algorithm 17 *PRODUCTION\_ELEMENT* string creation**

---

```

1: production_element_string = ""
2: if element.productionProperty == ProductionProperty.NONE then
3:   production_string += str(this.symbol)
4: else if element.productionProperty == ProductionProperty.REPETITION
      then
5:   production_string += "("
6:   production_string += str(this.symbol)
7:   production_string += ")"
8:   production_string += "*"
9: else if element.productionProperty == ProductionProperty.OPTIONAL then
10:  production_string += "["
11:  production_string += str(this.symbol)
12:  production_string += "]"
13: end if
14: production_element_string

```

---

### **Symbol string representation creation**

To create the string representation of a nonterminal symbol, its value can be returned as-is. The terminal symbol data type additionally has a production property as a class attribute.

To create the string representation of a terminal symbol, in the case of the property *OPTIONAL*, square brackets have to be printed around the terminal symbol value. In case of the production property *REPETITION* the asterisk symbol has to be

printed after the terminal symbol value. In case of the production property *NONE*, the string representation consists only of the terminal symbol value.

## 4.8 GUI

The GUI consists of the class *View* that can be seen in Figure 4.15. The class has the attributes *treeView* which is an instance of the PyQt *QTreeWidget*, *graphBuilder* which is an instance of the TPTPGraphBuilder and the boolean value *commentStatus*. The class functions and their implementation are described in the following sections based on the five sub-menus that have been introduced in chapter 3.11.

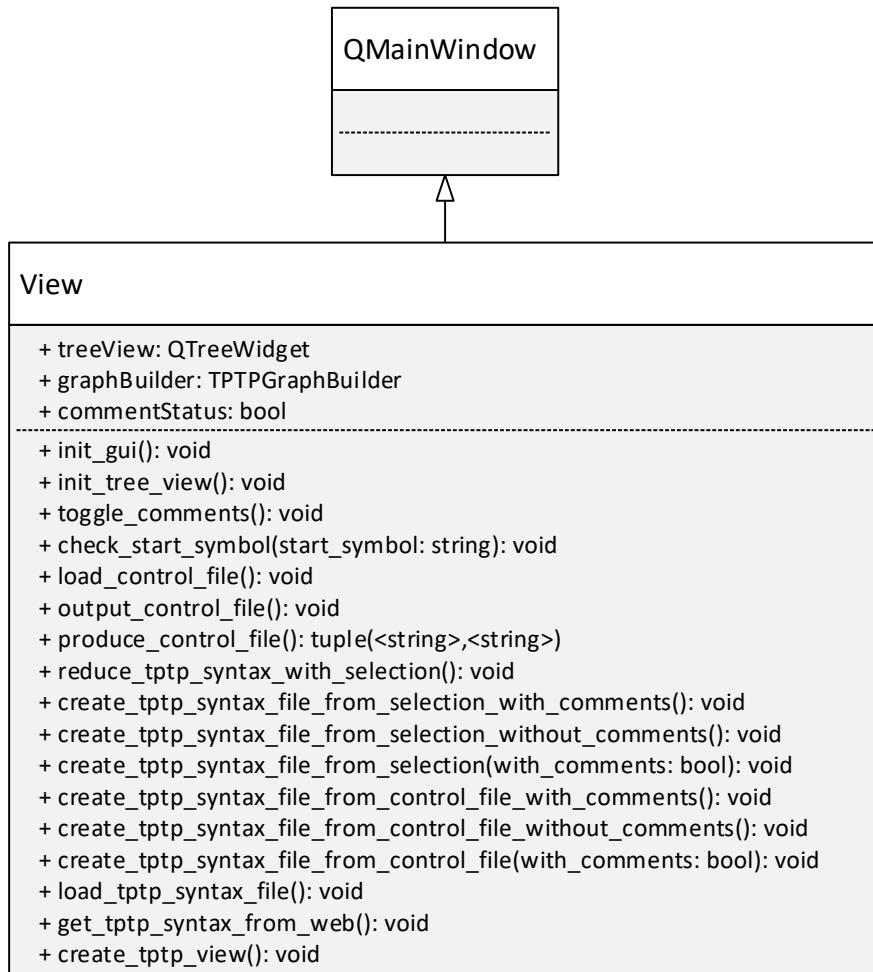


Figure 4.15: View UML class diagram

### 4.8.1 Main application window

#### Init GUI

The function *init\_gui* sets up the menu bar in the main application window. The menu bar consists of the sub-menus "Import syntax", "Save syntax", "Reduce", "Control file" and "View". To each sub-menu menu actions that have been described in chapter 3.11 are added and implemented using QActions. The following listing 4.14 shows the Python code for adding a menu bar and menu actions for example for the "Import syntax" menu bar.

```
1 import_menu = menubar.addMenu("Import syntax")
2 import_menu.addAction(import_tptp_file_action)
3 import_menu.addAction(import_tptp_file_from_web_action)
```

Listing 4.14: Implementation of menu bar

Besides a name, a short cut for executing the menu action is defined. Also, triggers have been defined that trigger functions once the menu button is pressed. Listing 4.15 shows how a short cut and a trigger have been defined for the action *ImportTPTPSyntaxfile*.

```
1 import_tptp_file_action = QAction('&Import TPTP syntax file', self)
2 import_tptp_file_action.setShortcut('Ctrl+O')
3 import_tptp_file_action.triggered.connect(self.load_tptp_syntax_file)
```

Listing 4.15: Implementation of menu actions

#### Init TreeView

Init TreeView is used for displaying a grammar in the main application window. The method is called after importing or reducing a grammar. The nodes of the grammar are stored in a nodes dictionary whose values (the nodes) are then converted to a list. This list is displayed by transforming it into a PyQt TreeView. A tree view can be used for displaying nested lists and allowing a user to navigate through them. The tree view consists of all nodes whose leaves are the productions in their (nested) productions list. The nodes can be expanded to see their productions which are then displayed indented and underneath the node.

A new tree view is defined using QTreeWidget with the header labels "Nonterminal", "Production Type" and "Production" (line 1 and 2 in listing 4.16). All nodes of

the nodes dictionary of the GraphBuilder instance are extracted and sorted based on their position in the original TPTP syntax file (line 3 and 4 in listing 4.16). The pseudo code in 18 shows the code used for filling the created treeView. For every node a new QTreeWidgetItem is generated. The item consists of the nodes (node.value) and its rule type. Also a new checkbox is defined that is displayed next to the item whose initial value is unchecked. The background colour is set to gray. For every production in the nodes production list, a new QTreeWidgetItem and checkbox are created as well. However, the default value of the checkbox is set to unchecked. The generated item is then added as a child to the node item. A potential comment block is added to the treeView and the item is added to the treeView as well.

```

1 self.treeView = QTreeWidget()
2 self.treeView.setHeaderLabels(['Non Terminal', 'Production Type', 'Production'])
3 nodes_list = list(self.graphBuilder.nodes_dictionary.values())
4 nodes_list.sort(key=lambda x: x.position)

```

Listing 4.16: Init Tree View

---

**Algorithm 18** Init tree view

```

1: for all nodes in nodes_list do
2:   rule_type = node.rule_type
3:   item = new QTreeWidgetItem
4:   Uncheck checkbox of item
5:   Set background to grey
6:   for all productions in node.productions_list do
7:     child_item = new QTreeWidgetItem
8:     Check checkbox of child item
9:     Add child_item to item
10:  end for
11:  if node.comment_block exists then
12:    comment_item = new QTreeWidgetItem
13:    Add comment_item to treeView
14:  end if
15:  Add item to treeView
16: end for

```

---

## 4.8.2 Import menu

### Import TPTP syntax file

The TPTP syntax file is imported using a QFileDialog which is returning the filename of the imported file selected by the user. If the filename is not empty, the application opens a dialog window where the user has to specify a desired start symbol. The method *read\_text\_from\_file* returns the text of the file belonging to the filename. With the text file and the start symbol, the method *create\_tptp\_view* is called.

### Import TPTP syntax from the TPTP website

The TPTP syntax file is extracted from the TPTP website using the Input method *import\_tptp\_syntax\_from\_web* (see chapter REF). This method returns a text file. The application opens a dialog window where the user has to specify a desired start symbol. With the text file and the start symbol, the method *create\_tptp\_view* is called.

*create\_tptp\_view*

The method creates a new instance of the *TPTPGraphBuilder*, calls *graphBuilder.run*, *init\_tree\_view* and *check\_start\_symbol*.

### Check start symbol

This function checks the checkbox of the defined start symbol. The QTreeView method *findItem(string itemName)* is used to find the start symbol in the treeview and return all items that match the item name. The found items are then checked.

### 4.8.3 View menu

#### Toggle comments

Algorithm 19 shows the algorithm used for toggling comments. The algorithm should either show or hide the comments based on the current status. The View class attribute *commentStatus* indicates whether comments are displayed (true) or not (false). For toggling comments, a new status is set that is the opposite of *commentStatus*. While looping through every item in *treeView*, comment items can be found by checking whether the item is checkable because by default a comment is not checkable by the user. A QTreeWidgetItem has flags including a flag "ItemIsUserCheckable" that is used for determining whether the item is checkable. The QTreeWidgetItem function *setHidden(bool hide)* is applied on every comment with the value of the new status. CommentStatus is set to the new status.

---

#### Algorithm 19 Toggle comments

---

```
1: newStatus = not commentStatus
2: for all items in treeView do
3:   if item is not checkable then
4:     item.setHidden(newStatus)
5:   end if
6: end for
7: commentStatus = newStatus
```

---

### 4.8.4 Reduce menu

#### Reduce TPTP syntax with selection

The method reduce TPTP syntax with selection reduces the syntax and displays the resulting syntax in the GUI. For reducing the TPTP syntax based on the user's selection in the GUI, a control file is produced. If no grammar has been imported, no start symbol or multiple start symbols have been selected an error is raised. The error is displayed using a QMessageBox. With the produced control file, the

Graphbuilder disables blocked productions, treeView is initialized and the start symbol is checked (see listing 4.17).

```
1 control_string, start_symbol = self.produce_control_file()
2 self.graphBuilder.disable_rules(control_string)
3 self.init_tree_view()
4 self.check_start_symbol(start_symbol)
```

Listing 4.17: Reduce TPTP syntax with selection

## 4.8.5 Save syntax menu

### Create TPTP syntax file from selection

Figure 4.16 shows a flow diagram of the process of reducing and saving a TPTP syntax file from a user's selection. The input of this method is a boolean value *with\_comments* that describes if the syntax should be saved with or without comments for the automatic parser generator (see chapter 5.1.2). First, the filename is import using a *QFileDialog*. If the import is successful, meaning a filename has been selected, a control file based on the user's selection is produced. If the import is not successful, the method terminates.

In comparison to *Reduce TPTP syntax with selection* the reduced syntax is not displayed in the GUI but only saved to local storage and the GUI displays still the full syntax. If the user wishes to display the reduced syntax in the GUI and save it, he can use the button *Reduce TPTP syntax with selection* and then use the button *Reduce and save TPTP syntax with selection*. Because the full syntax is still displayed, the graphBuilder object storing the syntax cannot be modified. Therefore, the old graphBuilder is copied to a new graphBuilder object. Because the old object should not change, the new object should have copies of all objects that are part of the graphBuilder rather than references. To copy the graphBuilder object which is a compound object it is necessary to use *copy.deepcopy* instead of *copy.copy*. Deepcopy creates a new object and inserts copies of all objects that are part of to be copied object. Copy however, creates a new object and inserts references to the to be copied objects.

The new graphBuilder is initialized with the start symbol that has been provided by the generated control file and rules that include blocked productions are disabled. During the initialization the new start symbol *<start\_symbol>* has been added (see chapter 3.6 for explanation). This start symbol is searched in the reduced syntax.

If it can be found, the syntax is saved with or without comments based on the input parameter. However, if the start symbol cannot be found in the syntax no part of the grammar is terminating. In this case the saved file is empty. To save the grammar the method *Output.save\_ordered\_rules\_from\_graph(\_with\_comments)* is used.

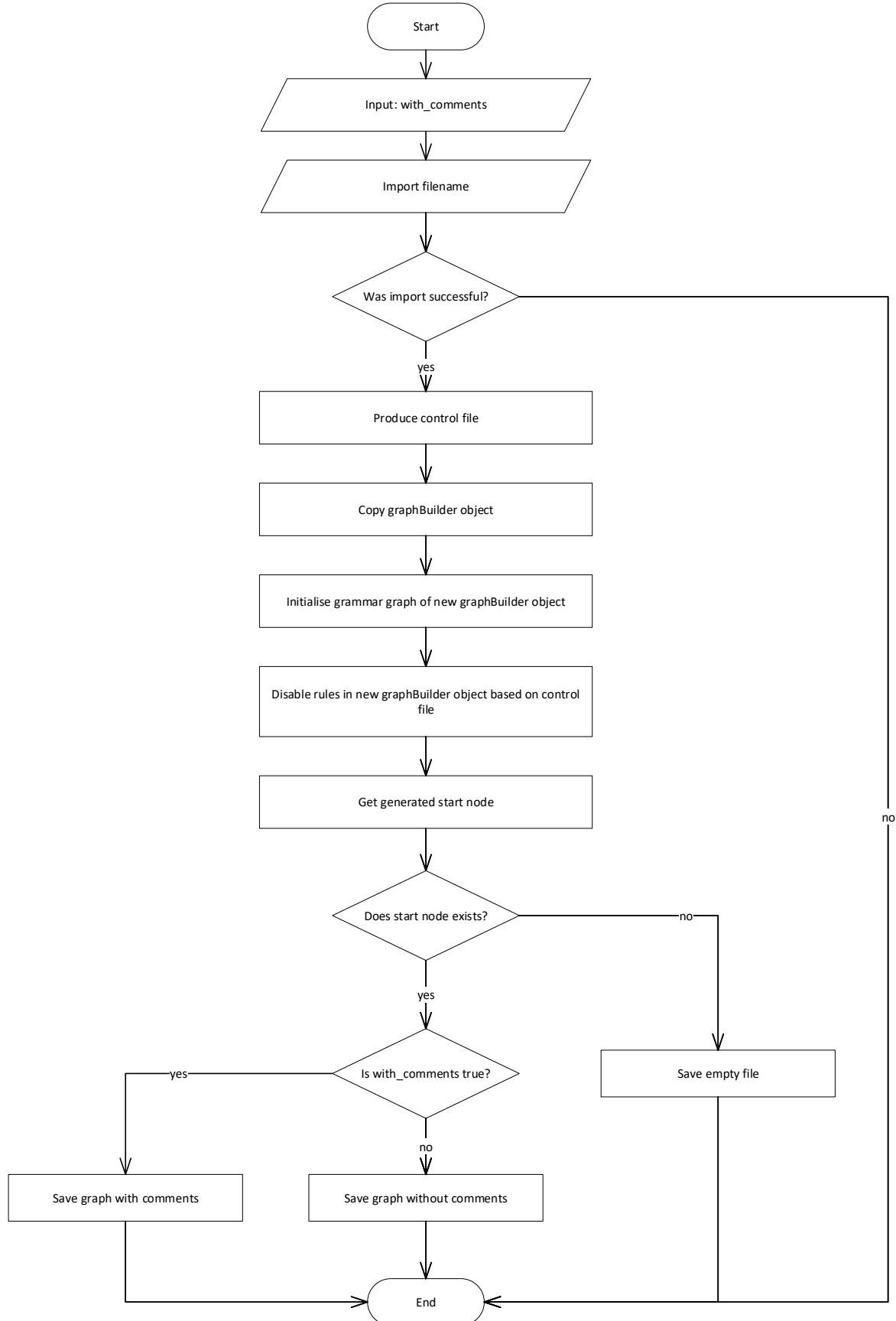


Figure 4.16: Reduce and save TPTP syntax

In the GUI there are two options of reducing and saving the **TPTP** syntax namely *Reduce and save TPTP syntax with selection* and *Reduce and save TPTP syntax from selection with external comments*. These two options both call the function described above and in figure 4.16. However, one calls it and passes *with\_comments* as true and the other as false. Because each button in the GUI needs its own designated function to call, the functions *create\_tptp\_syntaxfile\_from\_selection\_with\_comments* and *create\_tptp\_syntaxfile\_from\_selection\_without\_comments* are introduced. They simply set the parameter *with\_comments* accordingly and call the function *Create TPTP syntax file from selection*.

### **Create TPTP syntax file from control file**

The previous section explained how a **TPTP** syntax is reduced and saved based on a user's selection in the GUI. The described function in this section also reduces and saves a **TPTP** syntax but this time based on an imported control file. The function is similar to the previous function *Create TPTP syntax file from selection*. Figure 4.17 shows the same flow diagram as figure 4.16 but highlights applied changes in red. Changes are that besides importing a filename, a name for the control file has to be imported as well. Instead of producing a control file based on the user's selection, the specified control file is read and the grammar graph is initialised with the start symbol that is specified in the first line of the control file.

Also similar to the previous section, the functions

*create\_tptp\_syntax\_from\_control\_file\_with\_comments* and

*create\_tptp\_syntax\_from\_control\_file\_without\_comments* are introduced that set the parameter *with\_comments* accordingly and call the function *Create TPTP syntax file from control file*.

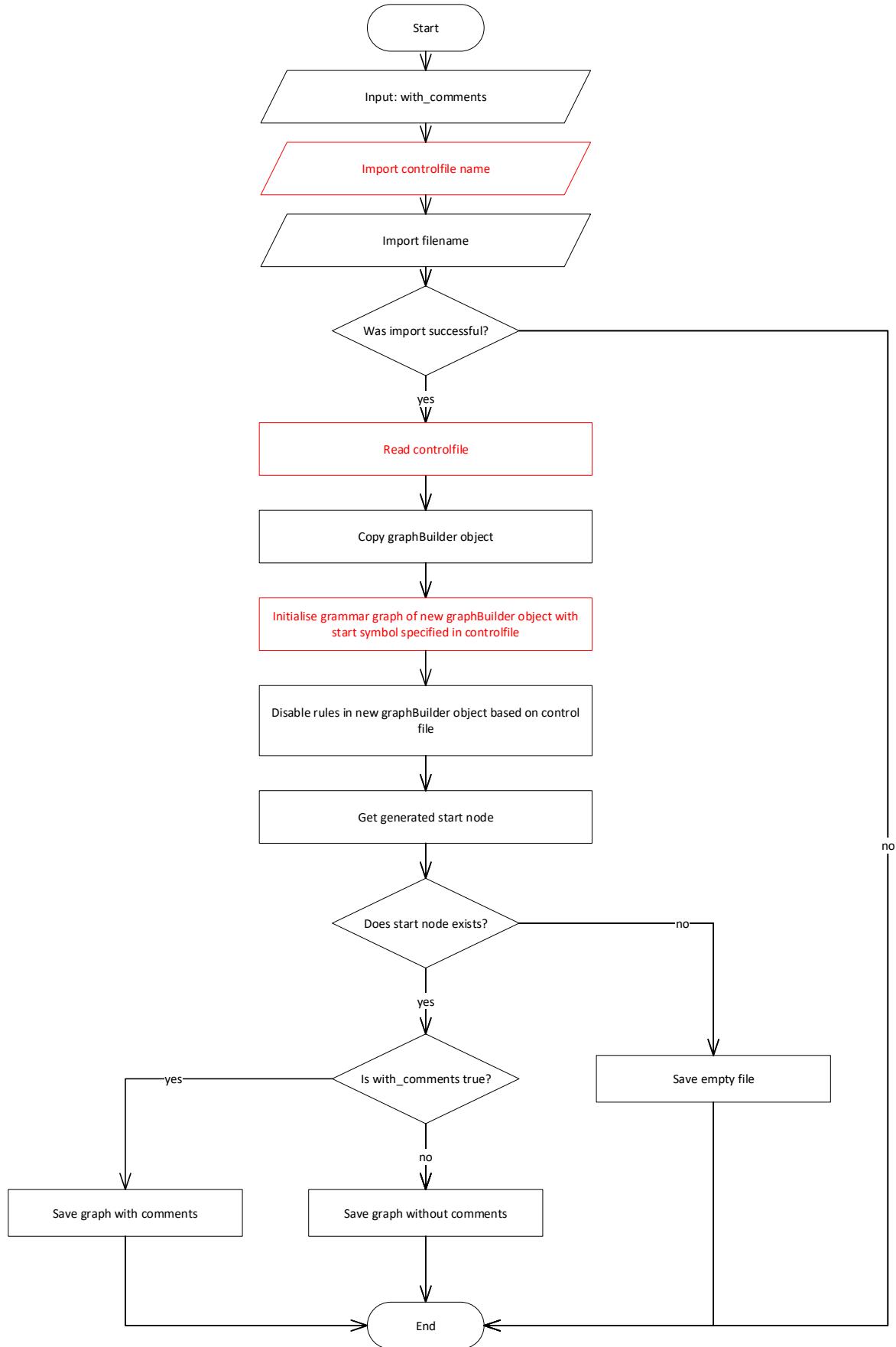


Figure 4.17: Reduce and save TPTP syntax from control file

## 4.8.6 Control file menu

### Import control file

If a control file is imported, checkboxes have to be set accordingly. The pseudo code in 20 shows the process. The name of the control file is read in using a QFileDialog. The file itself is imported using the method *Input.read\_text\_from\_file* (see chapter 4.6). If the file exists and if a tree view exists, meaning a grammar has been imported, all nodes are unchecked and all leaves/productions are checked. The different items representing nodes and leaves can be distinguished by looking for whether an item has parents and if not if it is checkable. The checkboxes are set unchecked because the default mode is that all nodes are unchecked and all productions are checked. After setting the checkboxes to the default mode, the first line of the control string is read and the start symbol is set accordingly. If there are multiple start symbols, every start symbol is checked. Iterating every line of the control string, the *nt\_name* which is the first element of a line is saved and the *rule\_type* which is the second element as well. The item that belongs to *nt\_name* is then searched and every production whose index is listed in the control file is unchecked.

**Algorithm 20** Load controlfile

---

```
1: Import Control File
2: if Control File exists then
3:   for all items in tree view do
4:     if parent of item is None then
5:       if item is not a comment then
6:         Uncheck item
7:       end if
8:     else
9:       Check item
10:    end if
11:  end for
12: Get first line of control file
13: Search item that represents start symbol
14: Check item
15: for all lines in control file do
16:   Split line by comma
17:   nt_name = first element of split
18:   rule_type = second element of split
19:   Search nt_name node in tree view
20:   if parent is None and rule_type matches rule_type then
21:     for all indices in line do
22:       if child exists then
23:         Uncheck child
24:       end if
25:     end for
26:   end if
27: end for
28: end if
```

---

**Produce control file**

Algorithm 21 shows how the control file used for reducing the TPTP syntax is produced. The index of productions that have been blocked by the user in the GUI are identified and stored in a dictionary until further processing. The key of the dictionary is a tuple made of a value and a rule type that has to be defined alongside the dictionary. Besides the blocked productions, the start symbol is featured in the control file as well, therefore an empty list for listing possibly multiple start symbols is created. Multiple start symbols can occur when the nonterminal symbol has multiple rule types or if the users makes an invalid input. The dictionary is

used for collecting all blocked productions. Looping through every node of the tree view it is checked, if the item has parents and if its checkbox is unchecked. If so, the rule type of the parent is identified and a new tuple is defined. The index of the production in the parent's productions list is identified as well. If the tuple is not already in the dictionary, a new dictionary entry is created with the tuple as key and the index of the item as value. If the tuple is already in the dictionary, the index of the item is added to the existing indices.

However, if the item is checked and has no parents, the name of the item is appended to the list of start symbols.

If there are multiple start symbols in the list that do not have the same nonterminal value or no start symbol at all, the user made an invalid input and an error is raised. The final control string is made up of first entry of the start symbol list. The first entry is sufficient because the control file only specifies the start symbols name and not the rule type. If there are multiple start symbols, they have the same nonterminal name. Each in a new line, the nonterminal symbol name, rule type and indices of blocked productions are appended to the control string. Within the line, the name, rule type and indices are separated by a comma. Moreover, an error is risen if there is no imported grammar, for example the user presses the Reduce grammar button in the GUI without importing a grammar.

**Algorithm 21** Produce controlfile

---

```

1: Define new tuple Entry("Entry", ["value", "rule_type"])
2: Create empty dictionary
3: Create empty list start_symbol
4: for all items in treeView do
5:   Get parent of item
6:   if item has parent and is unchecked then
7:     Get rule type of parent
8:     Create entry tuple with name of parent and its rule type
9:     Get index of child from parent
10:    if entry tuple is not in dictionary then
11:      Create new dictionary entry with entry tuple as key and list containing
        the index of the child as value
12:    else
13:      Append index of child to list
14:    end if
15:   else if item has no parent and is checked then
16:     Append item name to start symbol list
17:   end if
18: end for
19: control_string = first entry of start symbol list
20: Append all entries of dictionary to control_string

```

---

## Output control file

A control file can also be generated from GUI selections. For that, the information on start symbol selection and the selection of blocked productions is extracted from the GUI and a control file is generated from that information. The control file is then stored in a text file at the location the user specified.

### 4.8.7 Exception handling

Exceptions are used to handle wrong input from the user. There are three types of exceptions that are raised whenever the user makes a wrong input:

- Multiple start symbols
- No start symbol

- No imported grammar

The exception *multiple start symbols* is raised if the user selects multiple start symbols in the GUI and presses any Reduce button. On the contrary the exception *no start symbol* is raised if the user selects no start symbol in the GUI and presses any Reduce button. The exception *no imported grammar* is raised if the user starts the application and does not import a grammar but presses a button that requires a grammar for example in the Reduce, View or Control File menu.

## 4.9 Command-line interface

With the argparse module a command-line parser object can be created and parameters can be added to that object following the concept described. In the first line of listing 4.18 the command-line parser object is created including a description. Lines two to five contain the specification of the accepted arguments. In addition to the name and short form of the name, the type, a help message, and whether the parameter is optional or not can be specified. Default values for arguments can also be specified. If no default value is specified and if the argument is not passed it will have the value *None*. The *action = 'store\_true'* in line 5 means, that if this argument flag is present in the command-line call, the value of this argument will be set to *True*. If the flag is not present its value will be set to *False*. Argparse automatically checks the given conditions, for example if a required argument is not given and displays an error message if that is the case.

```
1 self.argument_parser = argparse.ArgumentParser(description='Extract sub-syntax using  
2 TPTP syntax file and a control file')  
3 self.argument_parser.add_argument('-g', '--grammar', metavar='', type=str, required=  
4 True, help='path of the TPTP syntax file')  
5 self.argument_parser.add_argument('-c', '--control', metavar='', type=str, required=  
6 True, help='path of the control file')  
7 self.argument_parser.add_argument('-o', '--output', metavar='', type=str, required=  
8 False, help='optional output file name (default output.txt)', default="output.  
9 txt")  
10 self.argument_parser.add_argument('--external_comment', action='store_true',  
11 help="flag - include external comment syntax")
```

Listing 4.18: Argparse command-line parser configuration

Argparse will also automatically create the help output by using the descriptions provided when configuring the argument parser.

# 5 Validation

This chapter reflects on the sub-syntaxes generated by [Synplifier](#). It demonstrates that sub-syntaxes are compatible with the [TPTP](#) syntax format and compares the sizes of several sub-syntaxes of interest.

## 5.1 Automated parser generation

A goal of using [Synplifier](#) is to be able to use an extracted sub-syntax with the automated parser generator for the [TPTP](#) syntax [12]. To ensure compatibility it is possible to export an extracted sub-syntax and adding the part of the syntax concerning comments, even though it is not reachable in the original syntax (see section [3.10.2](#)). If this part of the syntax would not be part of the output sub-syntax file automated parser generation would result in an error because this syntax part is expected to be present.

Also, the automated parser generator is used to check if the output sub-syntax follows the original [TPTP](#) syntax format.

### 5.1.1 Building a basic parser

To demonstrate the usability and capability of [Synplifier](#) a parser parsing only [CNF](#), that counts the number of [CNF](#) clauses, is used. The creation of the parser can be divided in the following steps:

1. Extract the [CNF](#) sub-syntax from the original [TPTP](#) syntax using [Synplifier](#).
2. Generate Lex and Yacc file based on the sub-syntax using the automated parser generator.
3. Modify the generated Yacc parser to count [CNF](#) clauses.

## CNF sub-syntax extraction

The following listing 5.1 contains the control file content, that extracts CNF from the TPTP syntax version 7.3.0.0. The start symbol is *TPTP\_file*.

```
1 <TPTP_file>
2 <annotated_formula> ::= ,0,1,2,3,5
3 <annotations> ::= ,0
```

Listing 5.1: Control file to extract CNF

All productions except the *<cnf\_annotated>* are disabled from the *<annotated\_formula>* grammar rule (line 2 of listing 5.1). The *<annotated\_formula>* grammar rule can be seen in the following listing 5.2. Annotations are also disabled (line 3 in listing 5.1).

```
1 <annotated_formula> ::= <thf_annotated> | <tff_annotated> | <
   tcf_annotated> |
2                                <fof_annotated> | <cnf_annotated> | <
                                tpi_annotated>
```

Listing 5.2: *<annotated\_formula>* production rule

To extract the CNF sub-syntax using the control file content from listing 5.1 either the command-line interface or the GUI can be used.

## Lex and Yacc file generation

The generated sub-syntax is input to the automated parser generator. The automated parser generates a Lex and Yacc file, and also the corresponding C-files from that.

## CNF clause counter implementation

To count the CNF clauses the counter *cnf\_counter* has been added to the parser. Also, a main function is added, that can be seen in listing 5.3. Using this function, either a file can be passed to the parser or the input can be provided via the command-line. After parsing is complete, the total number of CNF clauses is output to the console.

```
1 int main( int argc, char **argv ){
2     ++argv, --argc; /* skip over program name */
3     if(argc>0){
4         yyin = fopen(argv[0], "r");
5     }
6     else{
7         yyin = stdin;
8     }
9     yyparse();
10    printf("Total count of cnf clauses: %d\n", cnf_counter);
11 }
```

Listing 5.3: CNF parser main-function

The incrementation of the *cnf\_counter* is done in the Yacc rule-action, that has been created by the automated parser generator, for the *cnf\_annotated* symbol.

### 5.1.2 Testing the generated parser

The generated parser has been successfully tested on TPTP problems. It returns an error if other logics than CNF are present and correctly counts the number of CNF clauses.

The automated parser generator can be used with any sub-syntax generated by Synplifier, for example sub-syntaxes with FOF, or FOF and CNF.

## 5.2 Syntax size comparison

Based on the extracted CNF and FOF sub-syntax, the size of the sub-syntaxes in comparison to the full TPTP syntax is analyzed. Table 5.1 contains the number of rule statements and rules of the original TPTP syntax and extracted sub-syntaxes. A rule statement in this case means the combination of left-hand side nonterminal symbol, rule type and rule alternatives. Since rule alternatives are just a convenient way to display multiple rules with same nonterminal symbol on the left-hand side and the same rule type, in addition to the number of rule statements the number of rules is compared.

The reachable part of the TPTP syntax in the second row of the table contains only the reachable symbols counted from the start symbol *<TPTP\_file>*. Rule

statements are counted by counting the nodes in grammar graph and the number of rules is counted by counting the productions list entries of all nodes.

Table 5.1 implies that the **FOF** sub-syntax contains only 35 % of the rule statements of the full syntax. The **CNF** sub-syntax contains only circa 30 % of the rule statements of the full syntax. Compared with the total number of reachable rules in the **TPTP** syntax, **FOF** contains circa 33 % of the total reachable rules and **CNF** contains circa 28 % of the total reachable rules. This shows a significant syntax reduction which was one goal of the tool usage.

Table 5.1: TPTP syntax size comparison

Syntax	Rules	Productions
TPTP syntax v. 7.3.0	313	635
Reachable part of TPTP syntax v. 7.3.0	285	595
FOF	108	198
CNF	92	165

# 6 Conclusion

This paper described **Synplifier**, a software tool that automatically extracts sub-syntaxes from the **TPTP** syntax. The tool includes lexing and parsing the **TPTP** syntax, building a grammar graph that represents the syntax and extracting a sub-syntax based on the graph. The user can interact with the tool using a command-line interface or a GUI. Both the GUI and the command-line interface use a developed control-file format that specifies blocked productions and the start symbol. For the GUI a concept to display the grammar in a clear and structured way using a tree structure has been developed and implemented. The functionality of the tool has been demonstrated by extracting the **CNF** sub-syntax and creating a parser based on that syntax, which has been successfully tested with a number of **TPTP** problems.

All requirements, outlined in section 3.2, have been fulfilled. In addition, features like importing the original **TPTP** syntax from the **TPTP** website, toggling comments in the GUI and counting rules as well as rule productions have been included to increase the user-friendliness of **Synplifier**.

The tool reduces the entry barrier for new users to use the **TPTP** language because desired sub-syntaxes of the **TPTP** syntax can be extracted while maintaining consistency with the original **TPTP** syntax.

In short, the developed software tool enables users to consistently, conveniently, and automatically extract sub-languages from the **TPTP** syntax, turning this one language into a family of languages, coming from a single source, but optimized for particular use cases.

## 6.1 Future Work

The concept of sub-syntax extraction is not only limited to the **TPTP** language and can also be applied to any other computer language for example the SMT-lib [1].

In order to adapt [Synplifier](#) to extract sub-syntaxes from other languages, mainly the lexer and the parser components would have to be modified to parse the new syntax. Other components can be used with only minor modifications because the logic of the sub-syntax extraction would mostly remain the same. Furthermore, comment association would have to be adapted based on the specific comment convention (if comments are present in the syntax description).

It is also possible to create a tool that accepts plain [BNF](#) or [EBNF](#) since many computer languages are described in these forms. This would make the developed tool a standard tool since it could then be used with any arbitrary computer languages that is described in [BNF/EBNF](#).

Apart from adopting the tool to other languages, it can be investigated if comment association based on the comments content would lead to a significantly better result than using the heuristic approach that we have introduced.

# Bibliography

## Publications

- [1] Clark Barrett, Pascal Fontaine, and Cesare Tinelli. *The SMT-LIB Standard: Version 2.6*. Tech. rep. Available at [www.SMT-LIB.org](http://www.SMT-LIB.org). Department of Computer Science, The University of Iowa, 2017.
- [2] John E. Hopcroft, Rajeev Motwani, and Jeffrey D. Ullman. *Introduction to Automata Theory, Languages, and Computation 3rd Edition*. Pearson Education, Inc., 2007. ISBN: 0321455363.
- [3] Armin Cremers and Seymour Ginsburg. “Context-free grammar forms”. In: *Journal of Computer and System Sciences* 11 (1975), pp. 86–117.
- [4] Donald E. Knuth. “Backus Normal Form vs. Backus Naur Form”. In: *Commun. ACM* 7.12 (Dec. 1964), pp. 735–736. ISSN: 0001-0782. DOI: [10.1145/355588.365140](https://doi.org/10.1145/355588.365140). URL: <https://doi.org/10.1145/355588.365140>.
- [5] Niklaus Wirth. “What Can We Do about the Unnecessary Diversity of Notation for Syntactic Definitions?” In: *Commun. ACM* 20.11 (Nov. 1977), pp. 822–823. ISSN: 0001-0782. DOI: [10.1145/359863.359883](https://doi.org/10.1145/359863.359883). URL: <https://doi.org/10.1145/359863.359883>.
- [6] 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.
- [7] G. Sutcliffe. “The TPTP Problem Library and Associated Infrastructure. The FOF and CNF Parts, v3.5.0”. In: *Journal of Automated Reasoning* 43.4 (2009), pp. 337–362.

## PUBLICATIONS

---

- [8] G. Sutcliffe, S. Schulz, K. Claessen, and P. Baumgartner. “The TPTP Typed First-order Form with Arithmetic”. In: *Proceedings of the 18th International Conference on Logic for Programming, Artificial Intelligence, and Reasoning*. Ed. by N. Bjørner and A. Voronkov. Lecture Notes in Artificial Intelligence 7180. Springer-Verlag, 2012, pp. 406–419.
- [9] Jasmin Christian Blanchette and Andrei Paskevich. “TFF1: The TPTP typed first-order form with rank-1 polymorphism”. In: *Proc. of the 24th CADE, Lake Placid*. Ed. by Maria Paola Bonacina. Vol. 7898. LNAI. Springer, 2013, pp. 414–420.
- [10] G. Sutcliffe and C. Benzmüller. “Automated Reasoning in Higher-Order Logic using the TPTP THF Infrastructure”. In: *Journal of Formalized Reasoning* 3.1 (2010), pp. 1–27.
- [11] Cezary Kaliszyk, Geoff Sutcliffe, and Florian Rabe. “TH1: The TPTP Typed Higher-Order Form with Rank-1 Polymorphism”. In: *Proceedings of the 5th Workshop on Practical Aspects of Automated Reasoning*. Ed. by P. Fontaine, S. Schulz, and J. Urban. CEUR Workshop Proceedings 1635. 2016, pp. 41–55. URL: <http://ceur-ws.org/Vol-1635/>.
- [12] A. Van Gelder and G. Sutcliffe. “Extending the TPTP Language to Higher-Order Logic with Automated Parser Generation”. In: *Proceedings of the 3rd International Joint Conference on Automated Reasoning*. Ed. by U. Furbach and N. Shankar. Lecture Notes in Artificial Intelligence 4130. Springer-Verlag, 2006, pp. 156–161.
- [13] John Levine, Tony Mason, and Doug Brown. *Lex & Yacc*. O'Reilly Media Inc., 1992. ISBN: 9781565920002.
- [14] Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman. *Compilers : Principles, Techniques, & Tools, Second Edition*. Ed. by Michal Hirsch, Matt Goldstein, Katherine Harutunian, and Jefferey Holocumb. Second. Addison-Wesley, 2007.

- [15] Torben Aegidius Mogensen. *Introduction to Compiler Design*. Springer, 2017.  
ISBN: 9783319669656.
  
- [25] Scott Bradner. *A practical test for univariate and multivariate normality*.  
Tech. rep. Internet Engineering Task Force, Mar. 1997. URL: <https://tools.ietf.org/html/rfc2119>.

## Online sources

- [16] M. E. Lesk and E. Schmidt. *Lex - A Lexical Analyzer Generator*. URL: <http://dinosaur.compilertools.net/lex/index.html> (visited on 05/16/2020).
  
- [17] Vern Paxson. *Flex - a scanner generator - Description*. URL: [http://dinosaur.compilertools.net/flex/flex\\_4.html#SEC4](http://dinosaur.compilertools.net/flex/flex_4.html#SEC4) (visited on 05/16/2020).
  
- [18] Stephen C. Johnson. *Yacc: Yet Another Compiler-Compiler*. URL: <http://dinosaur.compilertools.net/yacc/index.html> (visited on 05/16/2020).
  
- [19] Charles Donnelly and Richard Stallman. *Bison 1.25 - Introduction*. URL: [http://dinosaur.compilertools.net/bison/bison\\_1.html#SEC1](http://dinosaur.compilertools.net/bison/bison_1.html#SEC1) (visited on 06/03/2020).
  
- [20] David Beazley. *PLY (Python Lex-Yacc)*. URL: <https://www.dabeaz.com/ply/> (visited on 01/26/2020).
  
- [21] Riverbank Computing Limited. *Introduction - PyQt v5.14.0 Reference Guide*. URL: <https://www.riverbankcomputing.com/static/Docs/PyQt5/introduction.html> (visited on 03/09/2020).
  
- [22] Python Software Foundation. *argparse - Parser for command-line options, arguments and sub-commands - Python 3.8.2 documentation*. URL: <https://docs.python.org/3/library/argparse.html> (visited on 03/09/2020).
  
- [23] Leonard Richardson. *beautifulsoup4 4.9.0*. URL: <https://pypi.org/project/beautifulsoup4/> (visited on 04/08/2020).

## ONLINE SOURCES

---

- [24] Leonard Richardson. *Beautiful Soup Documentation*. URL: <https://www.crummy.com/software/BeautifulSoup/bs4/doc/> (visited on 05/10/2020).
- [26] Python Software Foundation. *The Python Standard Library: Built-in Types*. URL: <https://docs.python.org/3/library/stdtypes.html#str> (visited on 05/10/2020).
- [27] Python Software Foundation. *tkinter — Python interface to Tcl/Tk*. URL: <https://docs.python.org/3/library/tkinter.html> (visited on 04/08/2020).
- [28] Michael Foord. *HOWTO Fetch Internet Resources Using The urllib Package*. URL: <https://docs.python.org/3/howto/urllib2.html> (visited on 04/18/2020).