



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.

Abstract

Contents

Acronyms	I
List of Figures	II
List of Tables	III
Listings	IV
1 Introduction	1
1.1 Problem statement and goals	1
1.2 Structure of the Report	2
2 Background and Theory	3
2.1 TPTP language	3
2.2 Formal languages	3
2.2.1 Finite automata	4
2.2.2 Regular expression	4
2.2.3 Formal grammars	4
2.3 Backus-Naur Form (BNF)	5
2.4 Lexing	5
2.5 Parsing	6
2.6 Python	7
2.6.1 PLY	7
2.6.2 PyQt	7
2.6.3 argparse	7
3 Concept	8
3.1 Requirements	8
3.2 Overview	8
3.2.1 Proposed architecture	9
3.2.2 Implementation language	10
3.3 Lexer	10
3.4 Parser	13
3.4.1 Data structures and data types	13
3.4.2 Production rules	16
3.4.3 Disambiguation of square brackets	16
3.5 Graph generation	16

3.6	Control file	18
3.7	Maintainig comments	19
3.8	Extraction of a sub-syntax	23
3.8.1	Parsing control file	23
3.8.2	Removing of blocked productions	23
3.8.3	Determination of the remaining reachable symbols	24
3.8.4	Determination of the remaining terminating symbols	24
3.9	Output generation	25
3.10	GUI	25
3.10.1	Menu	25
3.10.2	Display rules	26
3.11	Command-line interface	26
4	Implementation	28
4.1	Lexer	28
4.2	Parser	30
4.2.1	Data types	30
4.2.2	Defined grammar	31
4.3	Graph generation	35
4.3.1	Removing of blocked productions	37
4.3.2	Determination of the remaining reachable productions	39
4.4	GUI	39
4.4.1	Display rules	40
4.4.2	Toggle comments	40
4.4.3	Import control file	40
4.4.4	Import Thousands of Problems for Theorem Provers (TPTP) syntax from the internet	40
4.5	Command-line interface	40
5	Validation	43
5.1	Comment Association	43
5.2	Automated Parser Generation	43
5.2.1	Comment handling	43
5.2.2	Building a basic parser	43
5.3	Syntax size comparison	43
6	Conclusion	44
6.1	Future Work	44
	Bibliography	i
	Appendix	iii

Acronyms

ATP	Automated Theorem Proving
BNF	Backus-Naur Form
CFG	Context-free grammar
CNF	Clause Normal Form
EBNF	Extended Backus-Naur Form
FOF	First-order Form
PLY	Python Lex-Yacc
TFF	Typed First-order Form
THF	Typed Higher-order Form
TPTP	Thousands of Problems for Theorem Provers

List of Figures

3.1	Procedure of extracting a sublanguage	9
3.2	UML diagram of the architecture of the software tool	10
3.3	Parsing procedure	13
3.4	Maintaining comments flow chart	22
4.1	UML diagram for expressions	30
4.2	Symbols UML diagram	31
4.3	Parsing example	35
4.4	View UML class diagram	40

List of Tables

3.1	TPTP language production symbols [11]	11
4.1	Command-line interface parameters	41

Listings

3.1	Grammar expression	14
3.2	Graph generation	17
3.3	Control file	18
3.4	Comment in the TPTP syntax	19
3.5	Comment lines split by a <i>Top of Page</i> line in the TPTP syntax . .	20
4.1	Multi line production rule	28
4.2	Commented out production rule	30
4.3	Production element	35
4.4	Argparse command-line parser configuration	41

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 Clause Normal Form ([CNF](#)) to Typed First-order Form ([TFF](#)) have been included in and extended the [TPTP](#) language.

This report describes a tool 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. The developed tool 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.

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. By means of the background information, the third chapter outlines the concept of the developed tool. Based on this, the implementation of the tool is featured in the fourth chapter. The fifth chapter presents an evaluation of the effectiveness? of the tool. Considering the evaluation, the sixth chapter sums up the results of the developed tool, compares the results to the defined goals in first chapter and offers an outlook for possible future research.

2 Background and Theory

This chapter introduces the technologies and background that will be utilised in the following chapters. First, an introduction into the [TPTP](#) language is given. Then, formal grammars and the [BNF](#) are described. Following that, the foundations of lexing and parsing are outlined. Finally, Python and relevant Python modules that are used in the implementation are presented.

2.1 TPTP language

The Thousands of Problems for Theorem Provers ([TPTP](#)) is a library of problems for Automated Theorem Proving ([ATP](#)). Problems within the library are described in the [TPTP](#) language. The [TPTP](#) language is a formal language and its syntax is specified in an Extended Backus-Naur Form ([EBNF](#)). [1]

TODO more detailed

2.2 Formal languages

A formal language is a set of words over an alphabet.

An alphabet is a finite, nonempty set of symbols usually represented by Σ . An example is the binary alphabet $\Sigma = \{0, 1\}$. A string is a finite sequence of symbols from some alphabet. For example the string 101 is a string over the binary alphabet $\Sigma = \{0, 1\}$. A language is set of strings. If Σ is an alphabet, then $L(\Sigma)$ is a language over Σ . [2] - Vocabulary

2.2.1 Finite automata

2.2.2 Regular expression

A regular expression is an algebraic description of a regular/formal language. Regular expressions declare strings that are part of the language. [2] For example the regular expression $10+1^*$ denotes the language consisting of a single 1 followed by a single 0 or any number of 1's.

2.2.3 Formal grammars

Unlike regular expressions, grammars not only describe a language but also define a structure of the words of a language.

A grammar is a list of rules that defines the relationships among tokens [3]. These rules are also referred to as production rules. Given a start symbol, this symbol can be replaced by other symbols using the production rules. Using a recursive notation, production rules define derivations for words. The derived symbols can then once again be replaced until the derivation is a terminal symbol. Terminal symbols describe symbols that cannot be further derived. The alphabet of the described language is build by the set of terminal symbols. Nonterminal symbols however can be further derived and build merged with the terminal symbols the vocabulary of a grammar. Nonterminal symbols and terminal symbols are disjoint.

- Beispiel

Context-free grammar

Reduced grammars

Grammars are called reduced if each nonterminal symbol is terminating and reachable [4].

Given the set of terminal symbols Σ , a nonterminal symbol ξ is called terminating if there are productions $\xi \xrightarrow{*} z$ so that z can be derived from ξ and $z \in \Sigma^*$.

In other words, a nonterminal symbol ξ is terminating if there exist production rules so that ξ can be replaced by a string of terminal symbols. [4]

Given the set of terminal symbols Σ and the start symbol S , a nonterminal symbol

ξ is called reachable if there are production rules $S \xrightarrow{*} u\xi v$ so that S can be derivated to $u\xi v$ and $u, v \in \Sigma^*$.

In other words, a nonterminal symbol ξ is reachable if there exist production rules so that the start symbol can be replaced by a word containing ξ . [4]

todo beispiel

2.3 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. [5] An example for a BNF production would be $TPTP_File ::= \langle TPTP_Input \rangle \mid \langle comment \rangle$. Using this pattern of notation whole grammars can be specified.

The EBNF extends the BNF by 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.

[6]

2.4 Lexing

Lexing or a so-called lexical analysis is the division of input into units called tokens [3]. Tokens are for example variable names or keywords. The input is a string containing a sequence of characters, the output is a sequence of tokens. Afterwards, the output can be used for further processing like parsing. 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. [7]

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. More suitable is building a single automata that tests each token simultaneously. This automata can be build 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 is separating the input in order to divide it into tokens. Per convention the lexer chooses the longest input that matches any token. [7]

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. [7]

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 is usually written using regular expressions.

2.5 Parsing

The aim of parsing is to establish a relationship among tokens generated by a lexer [3]. For doing so, a parser builds a syntax tree out of the generated tokens [7].

Similar to lexers, parsers can be generated automatically. A parser generator takes as input a description of the relationship among tokens in form of a formal grammar (see). The output is the generated parser. [3]

During the syntax analysis a parser takes a string of tokens and forms a syntax tree with this construct by finding the matching derivations. The matching derivation can be found by using different approaches for example random guessing (predictive parsing) or LR parsing. Input: description of grammar [3] Output: parser [3]

-bottom up (LR parsing): parser takes inputs and searches for production where input is on the right side of a production rule and then replaces it by the left side

-top down (predictive parsing): parser takes input and searches for production where input is on the left side of a production rule

2.6 Python

2.6.1 PLY

Python Lex-Yacc ([PLY](#)) [8] is an implementation of lex and yacc in python. [LALR-parsing] consists of `lex.py` and `yacc.py`

`lex.py` tokenizes an input string

2.6.2 PyQt

PyQt is a Python binding for the cross-platform GUI framework Qt [9]. It is licensed under the GNU GPL version 3.

`tkinter`

2.6.3 argparse

The python module `argparse` is a module for creating command line interfaces. It provides the means to specify input arguments and [10] automatically creates help and usage messages. It also checks if the given arguments are valid. After specifying input arguments, the module will automatically create a parser for the specified arguments.

3 Concept

This chapter outlines the concept and the architecture of the software tool. First, in section 3.1, the requirements the software tool needs to meet are described. Then, in section 3.2, the components needed are introduced. Then the proposed software architecture is described. After that the concept of each component is developed.

3.1 Requirements

The tool should meet the following requirements:

The tool has a GUI that is the interface between the tool and a user. Hence, the user communicates with the tool via the GUI. The user is able to import a syntax file. After the syntax file is imported, it should be displayed. This includes displaying by the syntax defined productions as well as comments that are associated with these productions. The user can select a new start symbol and can select which productions should be blocked. After the user made his choice, the new sub-syntax is generated and displayed. The tool can also generate a control file listing blocked productions and the start symbol. Furthermore, the tool is able to import a control file and extract a sub-syntax based on this control file instead of extracting a sub-syntax based on a users selection of blocked productions. The new sub-syntax can be exported to .txt format. Also, comments referring to the remaining productions are kept and comments referring to productions that were discarded are not be included in the sub-syntax. The tool also provides a console interface. This interface accepts a [TPTP](#) syntax file and a control file and output the sub-syntax described in the control file. It is possible to specify the output path and filename.

3.2 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 should be 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 grammar. That new graph represents the syntax of the extracted sub-language. To make this grammar usable, lastly the syntax has to be output, based on the new graph, in the same format as the original syntax.

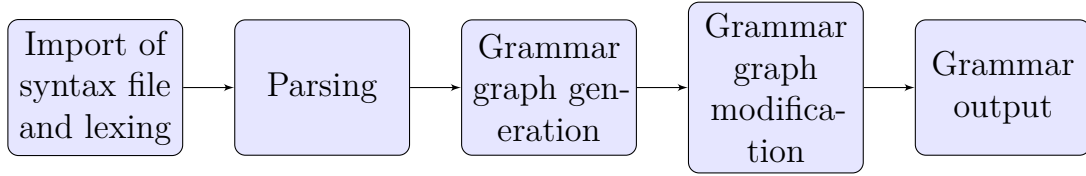


Figure 3.1: Procedure of extracting a sublanguage

3.2.1 Proposed architecture

The architecture of the software tool should take the procedure of extracting a sublanguage (section 3.2) into consideration. From that, five main components can be identified: An import module responsible for importing the **TPTP** syntax from a file; A lexer for extracting tokens from the language specification; A parser for creating a data structure from the tokens; A graph builder and manipulator; An export module for exporting the graph in a text representation corresponding to the original language specification.

In addition to the components that provide the main functionality a graphical user interface and a console interface for user convenience is desired.

Figure 3.2 contains a high-level UML diagram describing the architecture of the software tool. 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.

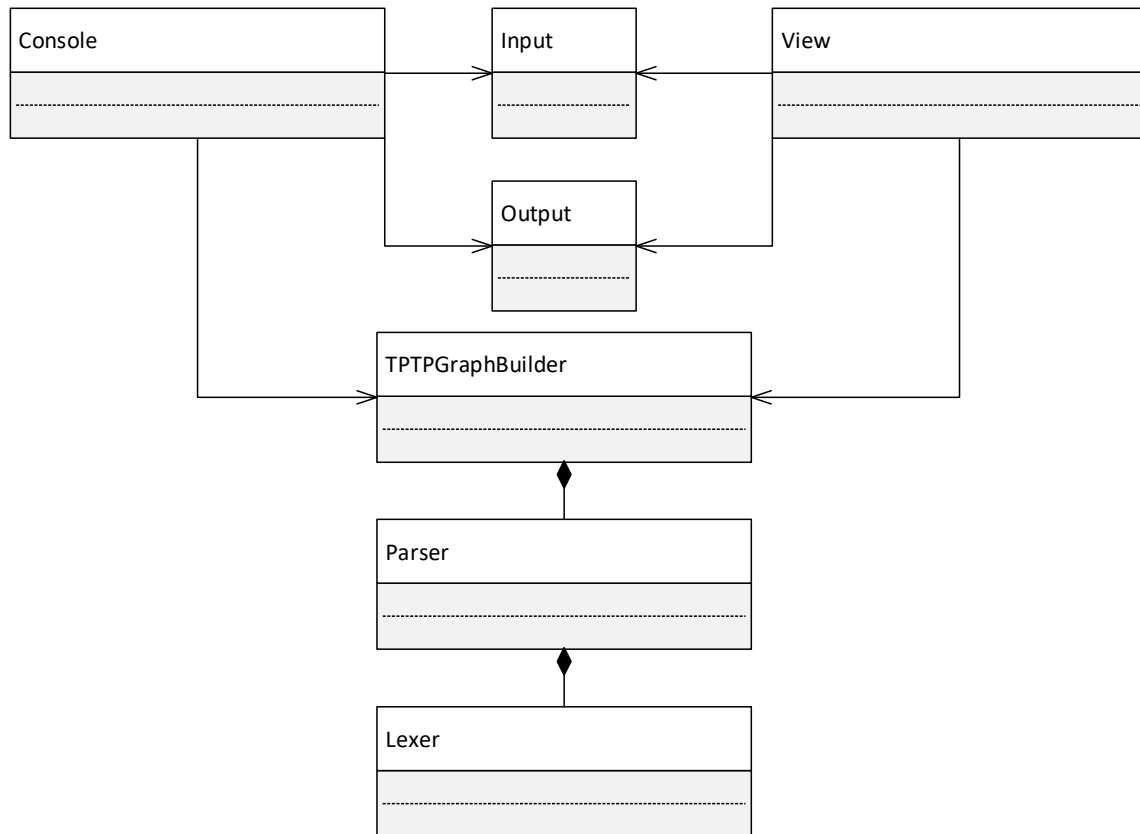


Figure 3.2: UML diagram of the architecture of the software tool

3.2.2 Implementation language

todo why python

3.3 Lexer

The lexer is responsible for extracting tokens from the TPTP language grammar specification file. Using **PLY** a lexer can be built by specifying tokens as regular expressions.

Therefore the TPTP language grammar specification needs to be analysed in order to find elementary tokens and regular expressions, that precisely describe these tokens.

Token identification

The syntax of the TPTP language is specified in a modified EBNF [??](#). Therefore there are deviations from standard EBNF (see [2.3](#)) that need to be analysed to specify elementary tokens. The standard EBNF only uses one production symbol (" ::= "). In the TPTP syntax additional production symbols have been added. The following table [3.1](#) contains the production symbols used in the TPTP syntax, that also have to be recognized by the lexer.

Table 3.1: TPTP language production symbols [\[11\]](#)

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

Another deviation from EBNF is that repetition is not denoted by surrounding curly brackets, but with a trailing * symbol.

Curly brackets have no special meaning in the TPTP syntax and can be treated as terminal symbols.

The meaning of the alternative symbol | is unchanged and also parentheses and square brackets can appear as meta symbols.

Also, there are line comments in the TPTP syntax. A comment starts with the % symbol at the beginning of a line and ends at the end of that line.

Following standard BNF, nonterminal symbols are enclosed by the < and > symbol and terminal symbols are written without any special marking.

The outcome of this are 13 defined tokens that are explained in the following.

Expressions can either be of the type grammar, token, strict or macro. They are defined as a nonterminal symbol followed by the production symbol itself

(::=, :=, ::, ...). The nonterminal symbol and the production are merged to a single token and are not identified as two tokens to avoid ambiguity while parsing. If not it would be difficult for the parser to determine whether the non terminal symbol that describes the rule is the start of a new rule or does still belong to the previous rule because as mentioned rules can cover multiple lines.

Nonterminal symbol

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

Terminal symbol

terminal symbol

Comment

A *comment* is identified by the lexer as a start of a new line followed by a percentage sign followed by an arbitrary character and ends with a newline. Because the percentage sign is also part of the terminal symbols, it is necessary to check whether the percentage sign is the first character of a new line. The percentage symbol when used as terminal symbol is embedded in square brackets and can never be the first character of a new line.

Meta symbol

Meta symbols include open and close parentheses "()", open and close square brackets "[]", asterisks "*" and vertical bars "|". They are recognized by the symbol itself and have a special meaning for the parser as they impact the to be build data structures.

3.4 Parser

The parser takes the tokens from the lexer as input and creates a data structure that represents the structure of the **TPTP** syntax.

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. The rules in the data structure have to be numbered, to maintain the correct order for output, after creating the grammar tree in the next step (see section ??).

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.

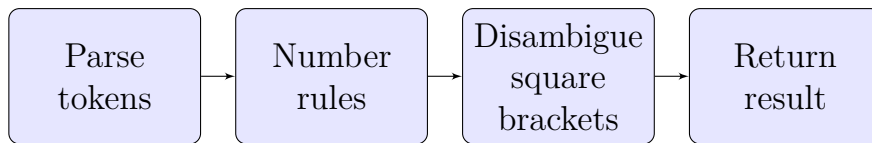


Figure 3.3: Parsing procedure

3.4.1 Data structures and data types

To build the representative data structure, data types that represent the data stored in the **TPTP** syntax have to be defined. The following section describes the data structure and data types that are used and created by the parser in the parsing process.

Terminal symbol

The terminal symbol data type has one attribute, which is the name of the terminal symbol it represents.

-todo Production Property

Nonterminal symbol

Analogue to the terminal symbol data type, the nonterminal symbol also has its name as an attribute.

Rules

A rule consists of the nonterminal symbol name which is produced, 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 (see table 3.1) there is a data type. This means that grammar, token, strict and macro rule data types are introduced.

Listing 3.1 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 3.1: Grammar expression

Comment block

A comment block is a list of consecutive comment lines.

Production element

A production element is either a terminal or nonterminal symbol. Additionally a production symbol has a production property.

Production property

The production property can take one of three values and denotes whether a production is optional, can be repeated any number of times or does not have any special property. In the original [TPTP](#) syntax file this was represented by square brackets or the repetition symbol.

Production

A production is one production alternative specified in any expression. It consists of a list of production elements and has a production property. Productions can also be nested. Therefore the list can also contain further productions

-show example

Productions list

A productions list contains a list of productions where each production is one alternative in the description of an expression.

XOR Productions list

The XOR productions list represents multiple alternatives enclosed by parentheses. It contains a list of the alternate productions.

Grammar list

The grammar list is the top level data structure. It contains a list of all elements that were in the [TPTP](#) syntax file. This includes any type of rules (grammar, token, strict and macro) and comment blocks.

3.4.2 Production rules

When using the [PLY](#) parser generator, production rules have to be defined. The rules describe how the tokens are to be processed.

todo describe production rules

3.4.3 Disambiguation of square brackets

As mentioned before, 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 would then be interpreted as denoting the optional production property. This production property would then be 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 can be iterated, exchanging the production property optional by the square bracket terminal symbols.

todo vor- nachteile

The output of the parser is a list of the rules and the comments from the [TPTP](#) syntax file.

3.5 Graph generation

todo nodes dictionary

todo überarbeiten und genauer schreiben Furthermore, the start symbol can have multiple rule types resulting in two start symbols. As it is not possible in the application to have two start symbols, a new start symbol is generated that implies the two original start symbols.

The [TPTP](#) syntax, extracted from the [TPTP](#) syntax file, needs to be stored in a data structure that allows for modification and traversing. The data structure that is used is a graph consisting of nodes. Every nonterminal symbol that is on the left side of a production in the [TPTP](#) syntax is represented by a node and instantiates a defined class "Node" that has the following attributes:

- value: name of nonterminal symbol
- productions list: productions list of nonterminal symbol (REFERENCE)
- 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: list containing all children of a node TODO A child is ..

todo mention nodes Starting with the start symbol, the graph is generated recursively. Iterating over each nonterminal symbol that is on the right side of a production rule, the corresponding node is identified. These nodes are then appended to the list of children of the nonterminal on the left side of the rule. The identified children may again have children. This process is repeated until a node has no children because there are only terminal symbols on the right side of the production rule.

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 wont be visited again independent of circular dependencies.

The following example 3.2 shows a production rule and the resulting list of children belonging to the node. Each production alternative(REFERENCE) has its own list of children.

```
1 Production rule:
2 <disjunction> ::= <literal> | <disjunction><vline><literal>
3 Output:
4 node.value: <disjunction>
5 node.ruleType: grammar
6 node.children: [[<literal>],[<disjunction>,<vline>,<literal>]]
```

Listing 3.2: Graph generation

3.6 Control file

In the following section a format for specifying the desired start symbol and blocked productions is described. Using a file-based configuration enables the user to store desired configurations and for example a manual selection in the graphical user interface is not necessary. It also helps with 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 nonterminal symbol and production symbol separating each group by a new line. First defining the nonterminal symbol, then the production symbol and after that the index of the alternatives that should be blocked (indexing starts at zero).

Identifying the production symbol is necessary because there may be a nonterminal symbol that has productions with more than one production symbol.

Listing 3.3 contains a sample control file. In this file in the first line `<TPTP_file>` is specified as start symbol. The second line means, that the second grammar production alternative of the nonterminal symbol `<TPTP_input>` should be disabled. Analogue to that, 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

This format is relatively easy to parse and also enables users to specify their desired start symbols and blocked productions without having to use the GUI.

pro: 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.

3.7 Maintainig 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, also the comment should be.

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: Comment in the **TPTP** syntax

todo check if listing is handled correctly

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 sematically. todo maybe proof 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 (see section 3.4.1), 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 and 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. todo why

- If this results in two comment blocks, one *Top of Page* line was present. Then the comment block before the *Top of Page* line is associated with the rule before when possible. If this comment block is at the beginning of the file it is associated with the rule after. The comment block after the *Top of Page* line is associated with the rule after. If it is at the end of the file it is 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. Therefore it is not particularly relevant. 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](#).

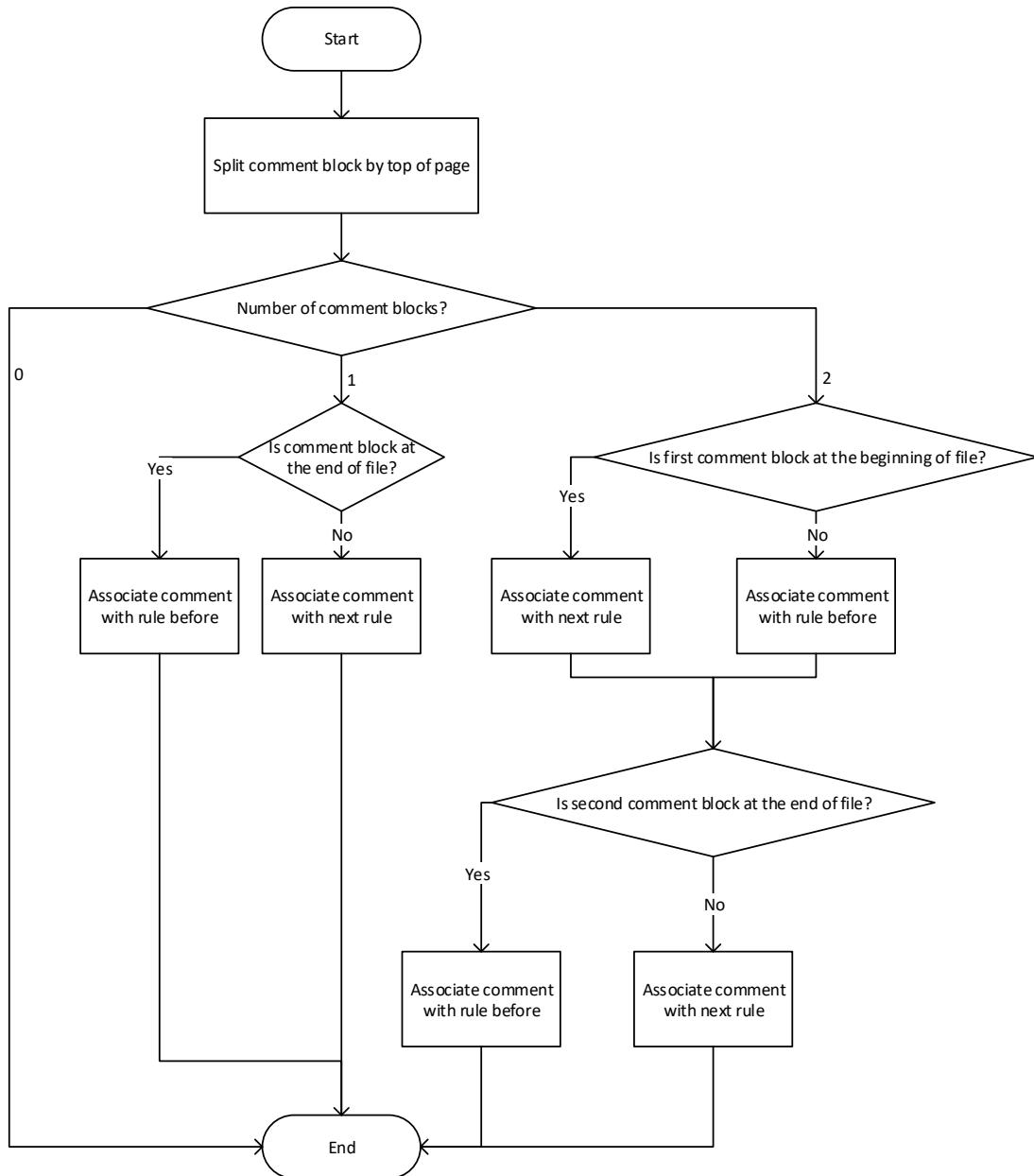


Figure 3.4: Maintaining comments flow chart

todo describe split comment block by top of page

todo explain could use content of comment

3.8 Extraction of a sub-syntax

This section covers the concept of how a sub-syntax can be computed from the original syntax. The original syntax is represented by a grammar graph (see section 3.5) and the information on what part of the grammar should be extracted is specified in the control file. For that, four steps must be performed: todo what

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 reachable part of the grammar must be computed.
4. Starting from the still reachable part of the grammar, non terminating productions must be removed. todo check order

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

3.8.1 Parsing 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 3.6. The start symbol is list 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 3.8.4) whereas the information on the productions that should be blocked will be needed in the next section.

3.8.2 Removing of blocked productions

The productions that are to be blocked have to be removed from the grammar graph before determining which symbols remain reachable.

-delete the productions that are specified in the control file -go through all nodes of which the nonterminal name is in the control file -delete the productions that are indexed in the control file -(see section before) and delete corresponding children

3.8.3 Determination of the remaining reachable symbols

After blocking the desired productions it has to be determined which part of the grammar still remains reachable. This can be done, by starting at the start symbol, going through all possible productions from there and at the end, removing all symbols, that have not been reached. These symbols are useless since they can not be reached.

The todo ?easiest way? to compute the remaining reachable grammar is to generate a new grammar graph (see section 3.5). The resulting graph will contain all reachable productions.

-determine which nonterminal symbols are still reachable, after removing productions specified in control file -determine remaining reachable by generating new graph -generate new graph again starting with temporary start symbol
-resulting graph only includes reachable part of the grammar

dynamic programming

3.8.4 Determination of the remaining terminating symbols

After all reachable symbols have been determined, the last step of extracting a sub-syntax is to determine and remove the non-terminating productions and symbols. Non-terminating productions are productions in which at least one non-terminal symbol is non-terminating. todo check reference -after reachable grammar has been determined -last step is to remove non terminating symbols dynamic programming -reachable means reachable from start symbol -remove non reachable productions -determine terminating symbols -start with start symbol -traverse through graph and compute terminating symbols

todo sources

3.9 Output generation

3.10 GUI

In this section ... The graphical user interface should display the grammar similar to the original language grammar specification file. It should also be possible to make selections in the GUI instead of having to use a control file. -show rules similar to file Selection of a new start symbol and productions that should be possible in the GUI and also with the import of a control file.

-show extracted grammar -export exported grammar -include + toggle comments
-> algorithm -import/export control file extra: -web import

3.10.1 Menu

The menu of the GUI has to provide multiple operations ...:

Import operations

An operation to import a [TPTP](#) syntax file has to exist. The latest [TPTP](#) syntax is provided on the [TPTP](#) project website [[TPTP](#)].

Providing the option of importing the [TPTP](#) syntax from the project website in addition to importing a [TPTP](#) syntax file from storage ensures the possibility of getting the latest version of the [TPTP](#) syntax. On the [TPTP](#) project website the [TPTP](#) syntax is stored in a HTML format. This file needs to be downloaded and converted to plain text.

-on import start symbol has to be selected

It would also be convenient if it was possible to import a control file, which would set selections in the GUI according to that control file.

Export operations

- export reduced grammar -export selection as control file -> to not always have to select the same selection or to use in command-line interface -
- export <comment> part of grammar additionally to make parser generator work
- other operations -reduce grammar -show/hide comments

3.10.2 Display rules

The rules from a [TPTP](#) syntax need to be displayed in order for the user to select a start symbol and productions that should be blocked when generating a sub-syntax.

- display rules from syntax -display enabled disabled rule -display comment -similar to file for intuitive understanding -select start symbol and blocked productions in rule
- todo mention check

3.11 Command-line interface

The goal of the command-line interface is to provide means for convenient automation of sub-syntax extraction. It should take a [TPTP](#) syntax file and a control file as input and output the resulting sub-syntax. Also basic help information should be accessible over the command-line interface.

This allows to use the software tool in scripts and enables automation for repeated tasks.

- command-line interface allows for automation with scripts for repeated tasks -basic functionality extract sub-syntax by providing base syntax file and control file, -output sub-syntax with input control file -provide basic information with help menu
- more complex actions like control file generation can more comfortably be done by using gui -gui package not needed

todo count rules

cd /

4 Implementation

4.1 Lexer

As mentioned in chapter 3.3 the implementation of the lexer consists of the definition of tokens in form of regular expression. The following paragraph presents defined tokens and their regular expressions.

Ignored symbols

It is possible to declare symbols that should be ignored. However, if a symbol is declared as ignored but is specially mentioned in another token, then if the sequence of characters represent that token, the ignored symbol is not ignored. In this project, tabs and white spaces are ignored as they do not have any special meaning other than providing clarity. Also, newlines are generally ignored because as can be seen in listing 4.1 there are rules that cover multiple lines.

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

Listing 4.1: Multi line production rule

Apart from the ignored symbols, there are 13 defined tokens:

Expressions

Expressions can either be of the type grammar, token, strict or macro. It is defined as a nonterminal symbol followed by the production symbol itself ($::=, :=, ::, \dots$). The nonterminal symbol and the production are merged to a single token and are not identified as two tokens to avoid ambiguity while parsing. If not it would be difficult for the parser to determine whether the nonterminal symbol that describes the rule is the start of a new rule or does still belong to the previous rule because as mentioned rules can cover multiple lines.

Regular expression of grammar expression: ' $\langle w+ \rangle [\backslash s]^ ::=$ '*

$\backslash w+$ matches any alphanumeric and underscore character that can occur more than one time. $[\backslash s]^*$ matches an arbitrary amount of white spaces.

Nonterminal symbol

A nonterminal symbol starts with " \langle " and ends with " \rangle ". In between there is any arbitrary sequence of numbers, underscores and small or capital letters.

Terminal symbol

Comment

A comment is identified by the lexer as a start of a new line followed by a percentage sign followed by an arbitrary character and ends with a newline. Because the percentage sign is also part of the terminal symbols, it is necessary to check whether the percentage sign is in a newline because the terminal symbol is not because the percentage symbol when used as terminal symbol is embedded in square brackets.

Meta-Symbols

Meta-Symbols include open and close parentheses " $()$ ", open and close square brackets " $[]$ ", asterisks " $*$ " and vertical bars " $|$ ".

They are recognized by the symbol itself and have a greater meaning for the parser as they impact the to be build data structures.

Ambiguity

The following example could either be matched as one comment token or as comment, grammar expression, non terminal symbol, terminal symbol, non terminal symbol. This ambiguity is solved because by convention the lexer matches the longest possible token, the sequence of characters is matched as one comment.

```
1 %—— <formula_role> ::= <user_role><source>
```

Listing 4.2: Commented out production rule

4.2 Parser

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

comment block reimplement equal operator

4.2.1 Data types

Figure 4.1 contains the UML modelling of the data types described in section 3.4.1.

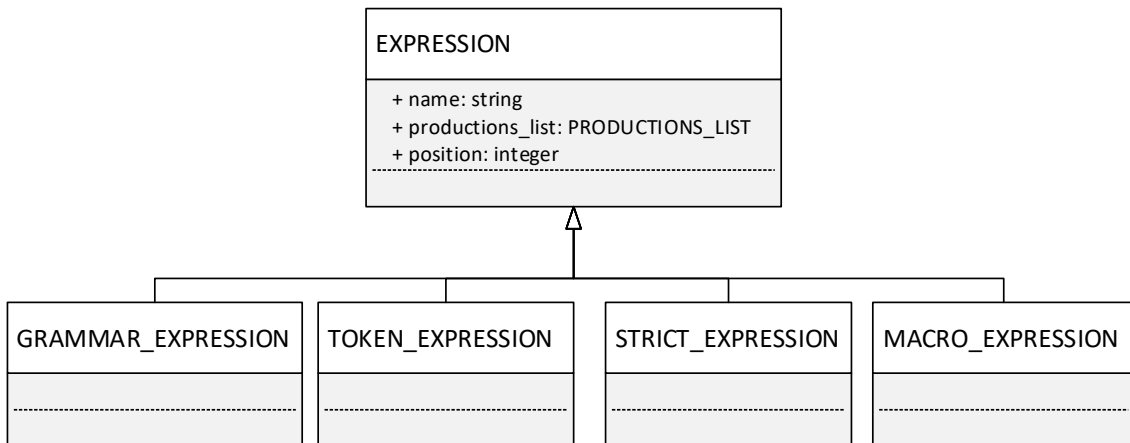


Figure 4.1: UML diagram for expressions

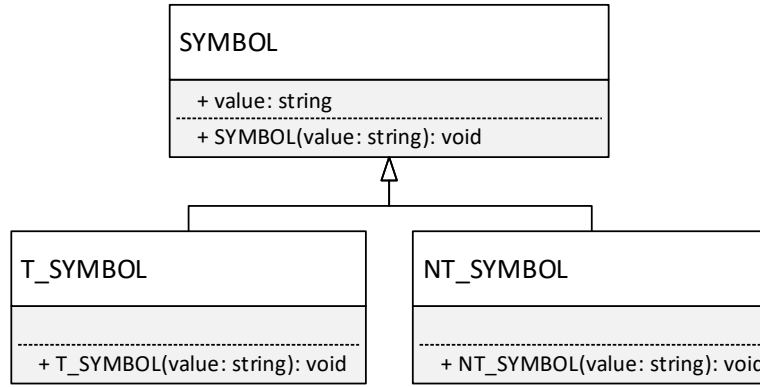


Figure 4.2: Symbols UML diagram

4.2.2 Defined grammar

- Input of parser is formal grammar

The grammar is formally specified $G = (N, \Sigma, P, S)$.

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, t symbol production and production element.

Σ is the set of terminal symbols. Terminal symbols of the specified grammar are the tokens generated by the lexer (see 3.3).

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

A grammar list implies a comment block, the four expressions or a grammar list followed by the expressions.

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

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

$$\textit{comment block} \rightarrow \mathbf{comment} \mid \textit{comment block} \mathbf{comment}$$

The four expressions are the expression token followed by their productions list. Grammar expression has a special case without a productions list because the production rule $\langle \textit{null} \rangle ::=$ has nothing on the right side.

$$\begin{aligned}
 \textit{grammar expression} \rightarrow & \mathbf{l} \mathbf{grammar expression} \textit{ productions list} \\
 & \mid \mathbf{l} \mathbf{grammar expression}
 \end{aligned}$$

$$\textit{token expression} \rightarrow \mathbf{l} \mathbf{token expression} \textit{ productions list}$$

$$\textit{strict expression} \rightarrow \mathbf{l} \mathbf{strict expression} \textit{ productions list}$$

$$\textit{macro expression} \rightarrow \mathbf{l} \mathbf{macro expression} \textit{ productions list}$$

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

productions list \rightarrow *production*
 | *productions list* **alternative symbol** *production*

xor productions list \rightarrow *production*
 | *xor productions list* **alternative symbol** *production*

T symbol production is either a t symbol or a t symbol/repetition symbol/ alternative symbol embedded in square brackets.

t symbol production \rightarrow **open square bracket** t symbol
 close square bracket
 | **open square bracket** repetition symbol
 close square bracket
 | **open square bracket** alternative symbol
 close square bracket
 | t symbol

Production element can be replaced by a nt symbol or by a nt symbol in square brackets or nt 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 t symbol production. A production element can also only be square brackets only.

production element → **open square bracket** **nt symbol**
close square bracket
| **nt symbol** **repetition symbol**
| *t symbol production* **repetition symbol**
| **open square bracket** **close square bracket**
| **nt symbol**
| *t symbol production*

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**

- *S*: grammar list, start symbol is not mentioned, per convention by [PLY](#) first rule found (top level rule)

Listing 4.3 shows the production rule of the nonterminal $\langle tfx_tuple \rangle$ as well as the tokens that have been generated by the lexer.

```

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

```

Listing 4.3: Production element

The resulting parse tree can be seen in figure 4.3.

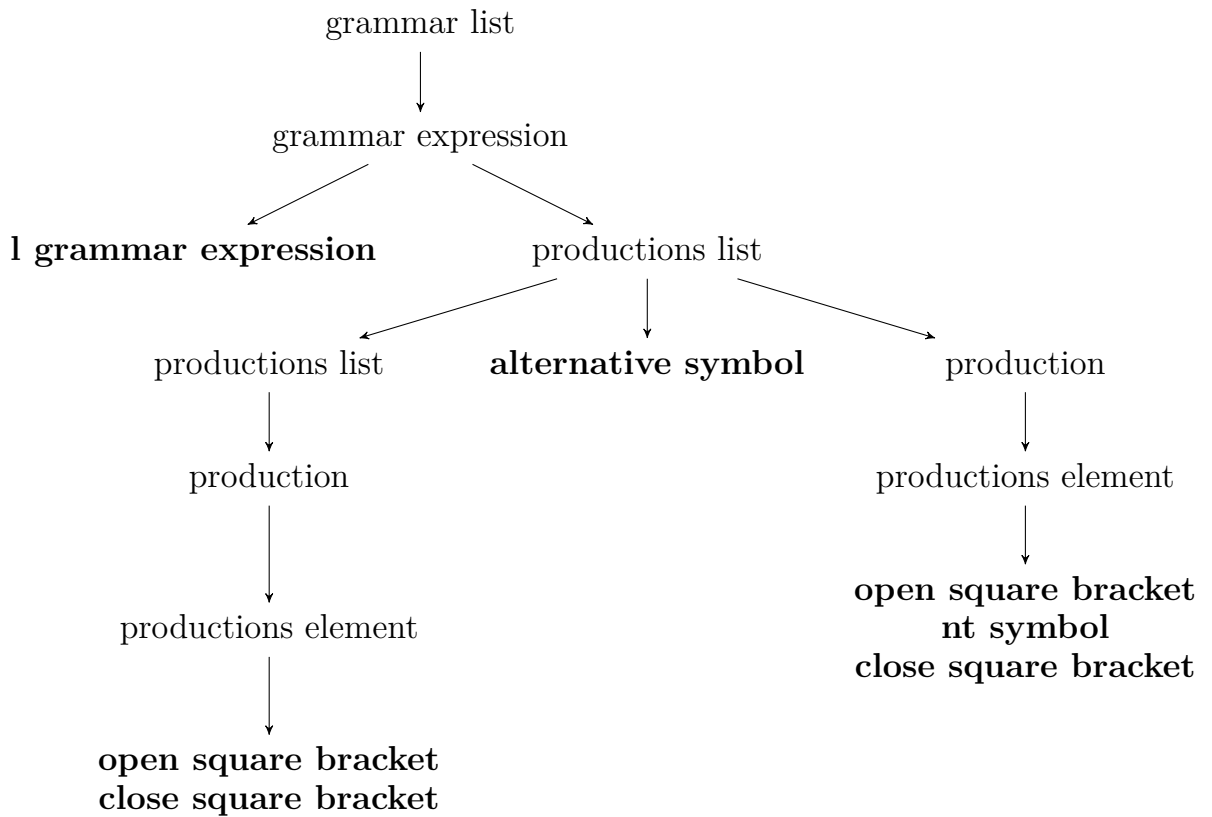


Figure 4.3: Parsing example

4.3 Graph generation

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

The algorithm *buildGraphRek* calls the function *searchProductionsListForNT* 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, every child calls the algorithm resulting in appending their own children to their children's list.

Algorithm 1 Graph Generation Algorithm: buildGraphRek

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 through all elements of the productions list. Each element is a production and calls the function *searchProductionForNT*. This function identifies the children of the given element who are then appended to the node.

Algorithm 2 Algorithm for extracting productions from productions list: searchProductionsListForNT

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

```

The goal is to identify the nonterminal symbols. Therefore 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 *searchProductionsListForNT* 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 3 Algorithm for appending children to node: searchProductionForNT

Input: node, productionsElement, children

```

1: for all elements in productionsElementList 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

```

4.3.1 Removing of blocked productions

The productions that should be blocked are specified in the control file. Algorithm ??

Algorithm 4 Removing blocked productions

Input: control_string

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

4.3.2 Determination of the remaining reachable productions

Removing non-terminating symbols

Algorithm 5 Removing blocked productions

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: delete_non_terminating Productions_(start_node, terminating, visited)
13: delete_non_terminating_nodes(terminating)
```

4.4 GUI

Tkinter and PyQt have been evaluated as a basis for the GUI. -pyqt offers checkboxes in treeviews

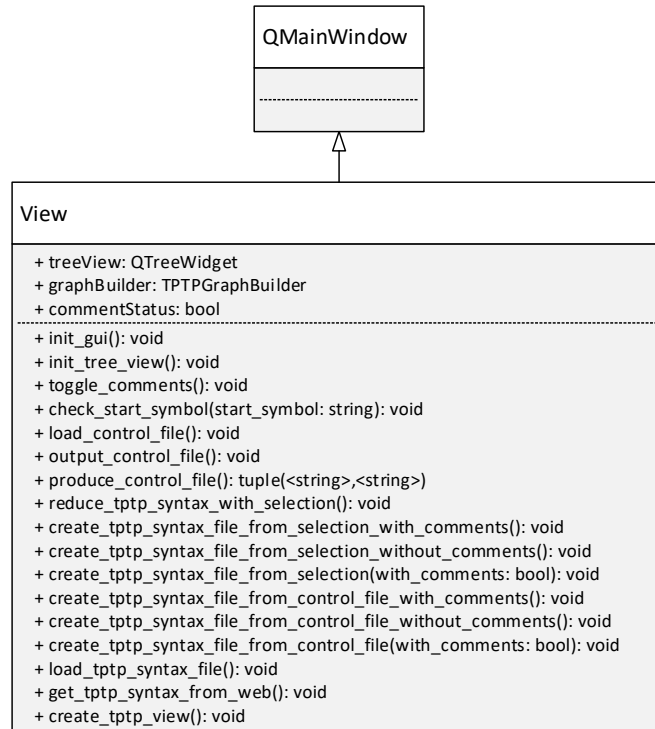


Figure 4.4: View UML class diagram

-todo uml diagram View

4.4.1 Display rules

4.4.2 Toggle comments

4.4.3 Import control file

4.4.4 Import **TPTP** syntax from the internet

4.5 Command-line interface

For the command-line interface (see section 3.11) the python module argparse is used. todo more, explain and explain why

Table 4.1 provides an overview about the needed command-line arguments. The syntax file location and the control file location has to be specified. Specifying an output path and file name is optional, by default the output filename will be *output.txt*.

Additionally, the help description can be opened by using the *-h* option

Table 4.1: 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)

With the `argparse` module a command-line parser object can be created and parameters can be added to that object. In the first line of 4.4 the command-line parser object is created with a description. Lines two to four 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*. `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
  TPTP syntax file and a control file')
2 self.argument_parser.add_argument('-g', '--grammar', metavar='', type=str, required=
  True, help='path of the TPTP syntax file')
3 self.argument_parser.add_argument('-c', '--control', metavar='', type=str, required=
  True, help='path of the control file')
4 self.argument_parser.add_argument('-o', '--output', metavar='', type=str, required=
  False, help='optional output file name (default output.txt)', default= "output.
  txt")

```

Listing 4.4: Argparse command-line parser configuration

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

-argparse Python offers a library for command-line interfaces. -options for path to grammar and control file -option for output path

-todo describe separation of gui and why

5 Validation

back to back testing show advantages and useful for tptp users... show size before after

5.1 Comment Association

comment association

5.2 Automated Parser Generation

5.2.1 Comment handling

5.2.2 Building a basic parser

5.3 Syntax size comparison

6 Conclusion

6.1 Future Work

comment association

Bibliography

Publications

- [1] G. Sutcliffe. “The TPTP Problem Library and Associated Infrastructure. From CNF to TH0, TPTP v6.4.0”. In: *Journal of Automated Reasoning* 59.4 (2017), pp. 483–502.
- [2] John E. Hopcroft, Rajeev Motwani, and Jeffrey D. Ullman. *Introduction to Automata Theory, Languages, and Computation 3rd Edition*. Pearson Education, Inc., 2007. ISBN: 0321455363.
- [3] John Levine, Tony Mason, and Doug Brown. *Lex & Yacc*. O’Reilly Media Inc., 1992. ISBN: 9781565920002.
- [4] Armin Cremers and Seymour Ginsburg. “Context-free grammar forms”. In: *Journal of Computer and System Sciences* 11 (1975), pp. 86–117.
- [5] 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>.
- [6] 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>.
- [7] Torben Aegidius Mogensen. *Introduction to Compiler Design*. Springer, 2017. ISBN: 9783319669656.

- [11] 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.

Online sources

- [8] David Beazley. *PLY (Python Lex-Yacc)*. URL: <https://www.dabeaz.com/ply/> (visited on 01/26/2020).
- [9] 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).
- [10] 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).

Appendix