



Politechnika Wrocławska

Faculty of Computer Science and Management

Field of study: Computer Science

Bachelor Thesis

FORMAL GRAMMAR PRODUCTION RULE PARSING TOOL

Karol Belina

keywords:

combinatory parsing, context-free grammars,
Extended Backus-Naur form

short summary:

The thesis documents the process of designing and implementing a tool for parsing the production rules of context-free grammars in a textual form. It discusses the choice of Extended Backus-Naur form notation over the alternatives and provides a model for parsing such a notation. The implemented parser can turn a high-level specification of a grammar into a parser itself, which in turn is capable of constructing a parse tree from arbitrary input provided to the program with the use of parser combinators.

Supervisor	dr inż. Zdzisław Szałowski
	Title/degree/name and surname

The final evaluation of the thesis

Chairman of the Diploma Examination Committee
	Title/degree/name and surname	grade	signature

For the purposes of archival thesis qualified to:*

- a) category A (perpetual files)
- b) category BE 50 (subject to expertise after 50 years)

* Delete as appropriate

stamp of the faculty

Wrocław 2020

Abstract

The thesis presents the design and implementation of a context-free grammar parsing tool with real-time explanations and error detection. It discusses the choice of Extended Backus-Naur form notation over the alternatives and provides a model for parsing such a notation. For this purpose, the official specification of the EBNF from the ISO/IEC 14977 standard has been examined and transformed into an unambiguous form. A definition of an abstract syntax tree is proposed to act as a result of the syntactic analysis phase formed with a technique called *combinatory parsing*. A method of testing an arbitrary input against the language generated by the constructed grammar is described. The thesis shows the process of creating a simple command-line program to act as a basic tool for interfacing with the grammar parser and checker, but in order to efficiently use the library, a web-based application is designed on top of that to serve as a more visual, user-friendly and easily accessible tool. It describes the deployment of the application on a static site hosting service. The designed and implemented system gives the opportunity to extend it with other grammar specifications.

Streszczenie

Praca przedstawia proces projektowania i implementacji narzędzia służącego do analizy syntaktycznej gramatyk bezkontekstowych z naciskiem na obsługę błędów i wyjaśnień w czasie rzeczywistym. Omawia wybór rozszerzonej notacji Backusa-Naura i przedstawia model do analizy takiej notacji. W tym celu przeprowadzono analizę i przekształcenie w jednoznaczną formę oficjalnej jej specyfikacji zdefiniowanej w standardzie ISO/IEC 14977. Zaproponowano definicję abstrakcyjnego drzewa składniowego dla tej notacji, która tworzona jest w wyniku analizy syntaktycznej za pomocą techniki zwanej *kombinacją parserów*. Opisano metodę sprawdzania dowolnego ciągu znaków pod kątem języka generowanego przez analizowaną gramatykę. Praca przedstawia stworzenie prostego programu działającego z poziomu wiersza poleceń, który jest podstawowym narzędziem do analizy gramatyk, jednak by móc efektywnie korzystać ze stworzonej biblioteki, zaprojektowano aplikację webową, która służy za bardziej wizualne, przyjazne i łatwo dostępne dla użytkownika narzędzie. Praca opisuje wdrażanie aplikacji na usługę hostingową dla statycznych stron. Zaprojektowany i wdrożony system daje możliwość rozszerzenia go o inne specyfikacje gramatyk.

Contents

1	Problem analysis	1
1.1	Description and motivation	1
1.2	Goal of the thesis	1
1.3	Scope of the project	2
1.4	Glossary	3
2	Theoretical preliminaries	5
2.1	Formal grammars	5
2.1.1	Introduction to formal grammars	5
2.1.2	The Chomsky hierarchy	6
2.1.3	Parsing expression grammars	7
2.2	Why EBNF?	8
2.3	Modifying the ISO/IEC 14977 specification	9
2.4	Lexical analysis	9
2.5	Syntactic analysis	10
2.5.1	Combinatory parsing	11
2.5.2	Abstract syntax tree	12
2.5.3	The parsing process	13
3	Analysis of similar solutions	15
4	Design of the project	19
4.1	Requirements	19
4.1.1	Functional requirements	19
4.1.2	Non-functional requirements	20
4.2	User stories	20
4.3	Use case specification	22
4.3.1	Use cases	22
4.3.2	Requirements traceability graph	23
4.3.3	Use case scenarios	23
4.3.4	Activity diagrams	25
4.3.5	Sequence diagram	26
4.4	System architecture	27
4.4.1	Logical architecture	27
4.4.2	Physical architecture	27
4.5	Interface prototype	28

5	Implementation of the project	29
5.1	Software environment	29
5.1.1	Used technologies	29
5.1.2	Project structure	37
5.2	Business logic	37
5.2.1	Domain modelling	37
5.2.2	Lexical analyser	39
5.2.3	Syntactic analyser	41
5.2.4	Semantic analyser	43
5.2.5	Grammar processing	43
5.3	Command-line application	46
5.4	Web-based application	47
5.4.1	Linking the business logic	47
5.4.2	Text editor	49
5.4.3	Parse tree visualizations	50
6	Project quality study	51
6.1	Business logic testing	51
6.1.1	Unit testing	51
6.1.2	Property-based testing	53
6.2	Integration testing	54
6.3	Benchmarking	55
6.4	Auditing	57
6.5	Complexity analysis	59
7	Deployment	61
7.1	Application building	61
7.2	Production environment	61
7.3	Continuous integration and continuous deployment	62
8	Software artifacts	65
9	User manual	67
9.1	System requirements	67
9.2	Installation guide	67
9.3	Usage guide	68
10	Summary	71
	Bibliography	73
	List of Figures	77
	List of Tables	79
	List of Listings	81
A	Modified specification of the ISO/IEC 14977	83
B	CD contents	85

Thesis structure

Chapter 1 describes the aim of the thesis and the problem it is solving. It states the scope of the project and provides a glossary of common terms used throughout the thesis.

Chapter 2 introduces the reader to theoretical concepts regarding the project.

Chapter 3 analyses various solutions available on the market, which are similar to the project stated in the thesis.

Chapter 4 introduces a specification of the project in terms of requirements, use cases, user stories, and various UML diagrams representing the architecture, use case scenarios, or the flow of the business logic.

Chapter 5 talks about the implementation phase of the project. It describes the technologies used in the implementation, as well as the project structure. Finally, it presents various intricacies encountered during the implementation process.

Chapter 6 studies the quality of the implemented project by proposing unit and integration tests of the business logic, defining benchmarks, as well as describing various auditing techniques.

Chapter 7 describes the process of deploying the final product into production environment, as well as techniques related to continuous integration.

Chapter 8 presents the results of the project in form of software artifacts.

Chapter 9 describes the process of installing, building, and running the web application on a local machine, as well as lists its system and software requirements. It presents screenshots of the web application along with a user guide and an example use case scenario of the product.

Chapter 10 sums up the thesis and proposes next steps in the evolution of the project.

1. Problem analysis

1.1. Description and motivation

Programming language theory has become a well-recognized branch of computer science that deals with the study of programming languages and their characteristics. It is an active research field, with findings published in various journals, as well as general publications in computer science and engineering. But besides the formal nature of Programming language theory, many amateur programming language creators try their hand at the challenge of creating a programming language of their own as a personal project. It is certainly relevant for a person to write their own language for educational purposes, and to learn about programming language and compiler design. However, the language creator must first of all make some fundamental decisions about the paradigms to be used, as well as the syntax of the language.

The tools for aiding the design and implementation of the syntax of a language are generally called *compiler-compilers*. These programs generate parsers, interpreters, or compilers from some formal description of a programming language (usually a grammar). The most commonly used types of compiler-compilers are *parser generators*, which handle only the syntactic analysis of the language — they do not handle the semantic analysis, nor the code generation aspect. The parser generators most generally transform a grammar of the syntax of a given programming language into a source code of a parser for that language. The language of the source code for such a parser is dependent on the parser generator.

Most such tools, however, suffer from too much complexity and generally have a steep learning curve for people inexperienced with the topic. Limited availability makes them less fitted for prototyping a syntax of a language — they often require a complex setup for simple tasks, which is not welcoming for new users. The lack of visualization capabilities shipped with these tools makes them less desirable for teachers in the theory of formal languages, who often require such features for educative purposes in order to present the formulations of context-free grammars in a more visual format.

1.2. Goal of the thesis

The main goal of this thesis is to design and implement a specialized tool, that serves teachers, programmers and other kinds of enthusiasts of the theory of formal languages in the field of discrete mathematics and computer science, in order to formulate and visualize context-free grammars in the Extended Backus-Naur form. To make it more approachable to users, the tool must provide a graphical user interface. Additionally, to ensure the highest degree of approachability, the tool must be available in the form of an easily accessible web-based application that is deployed on a web page and can run in a browser without the need of installation on the user's device. The thesis itself will document the entire process of creating such a project.

The final product will be much more welcoming to new users than other similar tools. Users will be able to access the website and use the tool right from their browser without

needing to install any programs on their device or carry out a complex setup process. The tool will be able to visualize the parse tree in an interactive manner, which will provide additional benefits to the entire process of working with this tool.

In order to achieve the general goal, several sub-goals have been distinguished, all of which contribute to the main objective as a whole

- analysis of existing solutions and applications,
- presentation of the theoretical preliminaries of the project,
- definition of the outline of the project, including a description of the functional and non-functional requirements, the use case diagram, use case scenarios, and user stories,
- description of technologies used in the implementation,
- implementation of the project,
- description of the testing and deployment environments.

1.3. Scope of the project

The thesis will propose a definition of a grammar in the form of an abstract syntax tree of the Extended Backus-Naur form. It will describe the process of implementing the business logic of the application in the Rust programming language [33] compiled to WebAssembly [43]. The compiled code is then ran inside the web-based application made with the Svelte framework, which incorporates the markup, CSS styles, and JavaScript scripts in the superset of the HyperText Markup Language (HTML).

The implementation phase will include the process of tokenization — the act of dividing the grammar in a textual form into a sequence of tokens — while taking into account proper interpretation of Unicode graphemes. The whitespace-agnostic tokens will be then combined together to form an instance of a previously-defined abstract syntax tree with a technique called *combinatory parsing*. Several smaller helper parsers will be defined, all of which then will be combined into more sophisticated parsers capable of parsing entire terms, productions, and grammars. The implementation phase will also include the definition of an algorithm for detecting left recursion in the resulting grammar, as well as catching any undefined production rules referenced in non-terminals. Up to this stage, any errors encountered in the textual form of a grammar are going to be reported to the user in a friendly format with exact locations of the errors in the input. After providing a valid definition of a grammar, the user will be able to supply an arbitrary input string to check if it belongs to the language generated by that grammar. This process will produce a parse tree, which will be displayed to the user as an interactive component. The web application will provide a basic code editor for inputting the grammar, with autocompletions, syntax highlighting, and other user-friendly features. The scope of the thesis includes the implementation of a simple command-line program for basic interfacing with the grammar parser and checker.

The web application will be deployed on the GitHub Pages hosting service for static sites.

1.4. Glossary

Please note that all terms below are explained in detail further in the thesis in their respective sections. This section provides simplified definitions of these terms.

AST — Abstract syntax tree

A tree representation of the syntactic structure of computer source code written in some language.

BNF — Backus-Naur form

A notation for context-free grammars used to describe the syntax of various languages.

CFG — Context-free grammar

A type of a formal grammar, in which the production rules can be applied regardless of the context (the surrounding symbols) of a non-terminal.

DFA — Deterministic finite automaton

A finite-state machine that accepts or rejects a given string of symbols, by running through a state sequence uniquely determined by the string.

EBNF — Extended Backus-Naur form

A family of languages used to describe other languages. An extension of Backus-Naur Form; allows to more easily define optionals, repetitions, and other constructs.

Formal grammar

Formally defines the construction of strings of symbols from the language's alphabet according to the language's syntax.

LBA — Linear bounded automaton

A restricted form of Turing machine that is able to recognize context-sensitive languages described by type-1 grammars in the Chomsky hierarchy.

Parser

A software component that takes input in the form of a sequence of characters or tokens and produces a structured output.

Parser-parser

The codename of the project.

PDA — Pushdown automaton

A type of automaton that employs a stack and is able to recognize context-free languages.

PEG — Parsing expression grammar

A type of analytic formal grammar introduced by Bryan Ford in 2004, which is an unambiguous version of the context-free grammar more suitable for machine parsing.

RE — Regular expression

The most restrictive type of a formal grammar often used for lexical analysis of programming languages and for defining search patterns for text.

REPL — Read-Eval-Print loop

A simple interactive computer programming environment that takes single user inputs, executes them, and returns the result to the user.

2. Theoretical preliminaries

2.1. Formal grammars

2.1.1. Introduction to formal grammars

Formal grammar of a language defines the construction of strings of symbols from the language's *alphabet* according to the language's *syntax*. It is a set of so-called *production rules* for rewriting certain strings of symbols with other strings of symbols — it can therefore generate any string belonging to that language by repeatedly applying these rules to a given starting symbol [26]. Furthermore, a grammar can also be applied in reverse: it can be determined if a string of symbols belongs to a given language by breaking it down into its constituents and analyzing them in the process known as *parsing*.

A formal grammar is a 4-tuple (N, Σ, P, S) , where N is the set of *non-terminal* symbols, which must be rewritten into other, possibly non-terminal symbols; Σ is the set of *terminal* symbols, which cannot be rewritten further; P is the set of production rules, and S is the start symbol. Let's consider a simple example of a formal grammar, where $N = \{ S, B \}$, $\Sigma = \{ a, b, c \}$, and P is the set of the following production rules:

1. $S \rightarrow aBSc$
2. $S \rightarrow abc$
3. $Ba \rightarrow aB$
4. $Bb \rightarrow bb$

To generate a string in this language, one must apply these rules (starting with the start symbol) until a string consisting only of terminal symbols is produced. A production rule is applied to a string by replacing an occurrence of the production rule's left-hand side in the string by that production rule's right-hand side. The simplest example of generating such a string would be

$$S \Rightarrow_2 \underline{abc}$$

where $P \Rightarrow_i Q$ means that string P generates the string Q according to the production rule i , and the generated part of the string is underlined.

By choosing a different sequence of production rules we can generate a different string in that language

$$\begin{aligned} S &\Rightarrow_1 \underline{aBSc} \\ &\Rightarrow_2 aB\underline{abcc} \\ &\Rightarrow_3 \underline{aa}Bbcc \\ &\Rightarrow_4 \underline{aabb}cc \end{aligned}$$

After examining further examples of strings generated by these production rules we may come into a conclusion that this grammar generates the language $\{ a^n b^n c^n \mid n \geq 1 \}$, where x^n is a string of n consecutive x 's. It means that the language is the set of strings consisting of one or more a 's, followed by the exact same number of b 's, then followed by the exact same number of c 's.

Such a system provides us with a notation for describing a given language formally. Such a language is a usually infinite set of finite-length sequences of terminal symbols from that language.

2.1.2. The Chomsky hierarchy

In [5], Chomsky divides formal grammars into four classes and classifies them in the now called *Chomsky hierarchy*. Each class is a subset of another, distinguished by the complexity. The hierarchy is visualized in Figure 2.1.

Type-3 grammars generate the so-called *regular languages*. As described in [3], regular languages can be matched by *regular expressions* and decided by a *finite state automaton*. They are the most restricting kinds of grammars, with its production rules consisting of a single non-terminal on the left-hand side and a single terminal, possibly followed by a single non-terminal on the right-hand side. Because of their simplicity, regular languages are used for lexical analysis of programming languages [23].

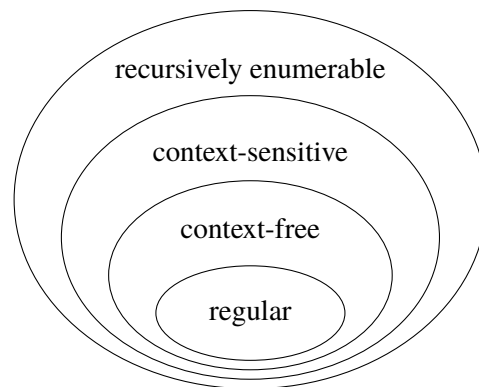


Figure 2.1. The Chomsky hierarchy visualized.

Type-2 grammars produce *context-free languages* and can be represented as a *pushdown automaton* which is an automaton that can maintain its state with the use of a stack. A pushdown automaton (PDA) varies from a finite state machine, as it is able to utilize the top of the stack in order to decide on the next transition, as well as manipulate the stack as part of performing a transition. A pushdown automaton reads the input sequentially from left to right. For every step, it decides which transition to take by taking into account the input symbol, the current state, and the current symbol on the top of the stack. A pushdown automaton can also mutate the stack by, for example, pushing a particular symbol onto the stack, or popping off the element from the top of the stack. A pushdown automaton can also just ignore the stack, and leave it in the same state as before. Context-free languages are defined by rules of the form of a single non-terminal on the left-hand side with a string of terminals or non-terminals on the right-hand side. They are the theoretical basis for the syntax and structure of most programming languages of today, even though their syntax usually comes with the problem of context-sensitive name resolution [21].

Type-1 grammars generate context-sensitive languages. The languages described by these grammars are exactly all languages that can be recognized by a *linear bounded automaton* (LBA) — a Turing machine whose tape size is finite and can be predicted based on the length of the input. This feature makes an LBA a more realistic representation of a real-world computer than a Turing machine, which by definition has a tape of infinite length. Because of the ability to perform random access on its memory, LBAs can generate languages based on an arbitrary context, hence the name of languages these grammars generate. Context-sensitive languages have less restrictive rules than type-2 and type-3 grammars — these come in the

form of $\alpha A \beta \rightarrow \alpha \gamma \beta$, where A is a non-terminal, α , β , and γ are strings of terminals, and α and β may be empty.

Type-0 grammars include all formal grammars. They generate all languages that can be recognized by a Turing machine. These languages are also known as the *Turing-recognizable* or *recursively enumerable* languages. They have no constraints on the structure of the production rules.

2.1.3. Parsing expression grammars

The syntax of various programming languages and protocols has been expressed with the use of Chomsky hierarchy of grammars — context-free grammars and regular expressions, in particular. Their original purpose is modeling natural languages, and the ability to express ambiguous constructs reflects that. This, however, makes it difficult to express deterministic, machine-oriented languages using context-free grammars. *Parsing expression grammars* (PEGs), introduced by Ford in [14] are an alternative to context-free grammars for specifying syntax formally. From the syntactic point of view, PEGs look similar to context-free grammars. They, however, have a different interpretation: most notably, the choice operator in PEGs chooses the first match, while in CFGs it happens ambiguously. As a result, in PEGs, if a string parses successfully, it must have exactly one valid parse tree.

PEGs address limitations of expressiveness of context-free grammars and regular expressions by simplifying their syntax definitions. Each parsing rule in a PEG has the form $A \leftarrow e$, where A is a non-terminal symbol, and e is a *parsing expression*. An *atomic* parsing expression is either a terminal symbol, a non-terminal symbol, or the empty string ϵ . A parsing expression which consists of a single terminal yields a *success* only if the beginning of the input string matches that terminal, and then consumes an appropriate number of input characters. In any other case, the expression yields a *failure* result. A parsing expression consisting of the empty string always succeeds without consuming any input. A parsing expression consisting of a non-terminal A depicts a recursive call to the production rule A .

Terminals, non-terminals, and empty strings may be combined to construct more complex parsing expression using the following operators:

- Sequence $e_1 e_2$, which first invokes e_1 , and in case e_1 succeeds, subsequently invokes e_2 on the unconsumed part of the input string, and returns the result. In case either e_1 or e_2 fails, the entire sequence expression $e_1 e_2$ fails consuming no input.
- Ordered choice e_1 / e_2 , which first invokes e_1 , and in case e_1 succeeds, returns the result. Otherwise, in case e_1 fails, the choice operator calls e_2 with the original input it invoked e_1 with, returning e_2 's result.
- Zero-or-more e^* , which consumes zero or more consecutive repetitions of the sub-expression e .
- One-or-more e^+ , which consumes one or more consecutive repetitions of the sub-expression e .
- Optional $e?$, which consumes zero or one occurrence of the sub-expression e .
- And-predicate $\&e$, which invokes the sub-expression e , and then succeeds if e succeeds and fails if e fails, but never consumes any input.
- Not-predicate $!e$, which succeeds if e fails and fails if e succeeds, consuming no input.

where e , e_1 , and e_2 are parsing expressions.

2.2. Why EBNF?

Backus-Naur form or *Backus normal form* (BNF) is a metasyntax notation for context-free grammars, often used similarly to Parsing expression grammars (discussed in Section 2.1.3) to describe the syntax of various languages, such as computer programming languages, instruction sets, document formats and communication protocols. They are applied in official language specifications, in manuals, and in textbooks on programming language theory — anywhere the precise descriptions of languages are required.

```
<expr> ::= <term> "+" <expr> | <term>
<term> ::= <factor> "*" <term> |
  ⇨ <factor>
<factor> ::= "(" <expr> ")" | <const>
<const> ::= "0" | "1" | ... | "9"
```

Listing 2.1. Example expression grammar in BNF.

BNF was proposed by John Backus, a programming language designer at IBM. He introduced a metalanguage of “metalinguistic formulas” to define the grammar of the new programming language IAL, recognized today as ALGOL 58. His notation was used for the first time in the ALGOL 60 report. BNF is a notation for context-free grammars defined by Chomsky.

A BNF specification is a set of production rules, written as $A ::= e$, where A is a non-terminal, and e consists of one or more sequences of symbols, where each sequence is separated by the vertical bar “|”, which represents a choice. The entire expression e is a potential substitution for the symbol on the left-hand side. Symbols that do not appear on a left side of any production rule are terminals. On the other hand, symbols that do appear on a left side are non-terminals, which must be always surrounded by angle brackets “<>”, as seen in Listing 2.1.

Extended Backus-Naur form (EBNF) is a family of extensions of the original Backus-Naur form. These are not, by any means, a superset of BNF, and there’s not a single official EBNF — each author defines their own variant of EBNF that is different from others. Any grammar described with EBNF can likewise be expressed in BNF, however, the representations of BNF are usually longer and more verbose. Constructs, such as options or repetitions are impossible to express directly in BNF. These require an intermediate production rule or, for options, an alternative operator of either nothing or the optional production, and, for repetitions, an alternative operator of nothing or the repeated production of itself, recursively. The BNF uses the symbols “<”, “>”, “|”, and “:=” for its grammar, but does not originally include quotes around terminals. As a result, these characters cannot be used in languages defined in BNF, and require a unique symbol for the empty string. In EBNF, terminal symbols must be surrounded by quotation marks (“...” or ‘...’). The angle brackets (“<...>”) used normally for non-terminals in BNF can be ignored in EBNF. BNF syntax can only represent a single rule in a single line, whereas in EBNF the end of a production rule is marked with the semicolon character “;”. Moreover, EBNF provides mechanisms for extensions, specifying the number of repetitions, excluding alternatives, comments, and so on.

The first EBNF notation was developed by Niklaus Wirth, who combined concepts of *Wirth syntax notation* with a different syntax. Many of the variants of EBNF are currently still in use. The International Organization for Standardization adopted an EBNF standard (ISO/IEC 14977) in 1996 [22]. Parser-parser uses EBNF as specified by that ISO. Other EBNF variants use different syntactic conventions.

2.3. Modifying the ISO/IEC 14977 specification

There are many concerns regarding the modern use of the ISO/IEC 14977 standard [44]. The use of this standard in Parser-parser comes mostly from its popularity, however, some considerations regarding the official specification of the notation should be taken into account. Although Parser-parser conforms to the majority of the EBNF specification, there are some deviations introduced to facilitate the use of the specification in a modern environment.

First of all, ISO/IEC 14977 does not support International/Unicode characters, code points, or byte values. Instead, it is defined through ISO/IEC 646 characters. Even though it provides a special notation for describing a character in an informal way, it does not have official support for it, consequently, having no ability to directly represent the full range of characters allowed by ISO/IEC 10646 or Unicode when processing text. Additionally to text, ISO/IEC 14977 cannot describe binary formats or provide characters by their value. It cannot include Unicode characters as terminal symbols in the grammar by surrounding them with single or double-quotes. Parser-parser extends the ISO/IEC 14977 specification by taking into account Unicode characters in place of `letter` and `decimal digit` definitions in the ISO/IEC 14977 EBNF specification [22].

Secondly, the ISO/IEC 14977 specification may be difficult to understand, and many of the key concepts are not clearly defined. The specification is abstract, which may be somewhat inevitable given the subject matter. The syntax production rule in the specification is defined ambiguously in multiple sections of the specification.

This led to an eventual modification of the ISO/IEC 14977 specification for the final implementation of the EBNF parser. The modified specification can be seen in Appendix A. It addresses the problem of Unicode characters and unambiguously defines the grammar.

2.4. Lexical analysis

Prior to parsing, the input string provided by the user should be divided into individual *tokens* in the process of *tokenization*, or the *lexical analysis*. Thanks to tokenization, all symbols provided in a textual form get separated and become independent from whitespace characters or other meaningless constructs of the language. This makes the parsing phase easier to reason about — the designers do not need to concern themselves with whitespace on every occasion.

A program that performs lexical analysis may be termed a *lexer*, *tokenizer*, or *scanner*. A lexer is generally combined with a parser, which are used together to analyze the syntax of programming languages, web pages, and other similar concepts. A lexer outlines the first phase of a parser front-end in language processing. The lexical analysis commonly happens in a single pass. Lexers tend to be uncomplicated, with most of the complexity deferred to the parsing phase. However, lexers can sometimes involve more complex behavior to make the input easier and simplify the parser. Lexers may be written by hand, to either support more features or simply for performance reasons [35].

A *lexical token* or simply *token* is a string with an assigned meaning. It is defined as a pair of a *token name* and an optional *token value*. The token name is a category of the token. The EBNF specification consists of several, clearly identifiable token types, which are listed in Table 2.1.

Table 2.1. Tokens types from the modified EBNF notation.

Token name	Normal representation
Non-terminal	a string of alphanumeric characters or whitespace characters beginning with an alphabetic character
Terminal	a string of at least one character surrounded by either “'”s or “””s
Special	a string of characters surrounded by question marks (“?”)
Integer	a string of at least one decimal digit with optional whitespace in-between
Concatenation	“,”
Definition	“=”
Definition separator	“ ”, “/”, or “!”
End group	“)”
End option	“]” or “/)”
End repeat	“}” or “:)”
Exception	“_”
Repetition	“*”
Start group	“(”
Start option	“[” or “(/”
Start repeat	“{” or “(:”
Terminator	“;”

Several of these tokens must carry an additional token value to distinguish it from other tokens of the same name. These include: non-terminals, terminals, specials, and integers.

2.5. Syntactic analysis

Parsing, syntax analysis, or syntactic analysis is the process of analyzing a string of symbols, either in natural languages, computer languages, or data structures, adhering to the rules of a formal grammar. A *parser* is a program that takes input data and produces a data structure — usually, a parse tree, abstract syntax tree, or other hierarchical structure, providing a structural description of the input while simultaneously ensuring the correct syntax [2, 35]. This is usually accomplished with a context-free grammar which defines elements that make up an expression and the order in which they appear, in a recursive manner. Determining if and how the input string can be produced from the initial symbol of the grammar is the main responsibility of the parser. This usually can be accomplished in two ways:

Top-down parsing — This method of parsing can be portrayed as finding left-most derivations of an input string by searching for parse trees using a top-down expansion of the given production rules of the formal grammar. Tokens of the input string are consumed sequentially from left to right. The inclusive choice is used to support ambiguity by expanding all alternative right-hand-sides of production rules.

Bottom-up parsing — This method of parsing starts with the input string and makes an attempt to rewrite it back to the initial symbol. The parser tries to locate the most basic components, then the elements containing these, etc.

2.5.1. Combinatory parsing

Combinatory parsing is a top-down parsing technique that involves *parser combinators*. A parser combinator is a higher-order function that receives several parsers as inputs and produces an entirely new parser as its output. In this context, a parser is a function accepting a string of tokens as input and returning some structured output, usually a parse tree. Parser combinators enable a strategy of recursive descent parsing that simultaneously facilitates modular construction of parsers [40].

Parsers formulated with the use of parser combinators are straightforward to construct, readable, modular, well-structured, and easily maintainable [40]. They have been widely used in situations where complex and diverse semantic actions are closely integrated with syntactic analysis, such as compiler prototyping and processors for domain-specific languages.

Following examples of parsers will be written in the Haskell [19] programming language. Haskell is purely functional, and one of its features (heavily utilized in examples below) is higher-order functions, which are functions that can be passed to other functions as values. Every expression in Haskell has a type which is determined at compile time. Types become not only a form of guarantee, but a language for expressing the construction of programs. At the same time, the programmer does not have to state every type of every value, as types will be inferred by the compiler automatically. This makes the syntax of Haskell very concise and ideal for presenting various mathematical concepts.

The basic idea of a parser is a function from strings to some *things*. It can be easily defined in a programming language that supports higher-order functions. The definition of a parser in Haskell would look like

```
type Parser t = String -> Maybe (t, String)
```

A parser that returns *things* of an arbitrary type *t* is a function that takes a string and returns *Maybe* a pair of *thing* *t* and the unconsumed part of the input string, which can in turn be passed on to some other parser.

The *Parser* type utilizes the built-in *Maybe* type — a type that represents an optional value. *Maybe*, as an algebraic data type is defined as

```
data Maybe t = Nothing | Just t
```

Maybe is a type constructor of a single argument (here, the *t* type), where the value of type *Maybe t* either contains a value of type *t* (represented as *Just t*, where *Just* is a single-argument value constructor of *Maybe*), or it is empty (represented as *Nothing*, which takes no arguments). Using *Maybe* is a good way to deal with errors or exceptional cases, as it is in the *Parser* type. Further in the thesis, Listing 5.1 in Chapter 5 will illustrate a concrete use of algebraic data types. Please note that the *Parser* type is not an algebraic data type — it is an alias to an existing function type from *Strings* to *Maybe (t, String)*s. A type alias does not introduce any value constructors for that type.

Such a *Parser* type can be found in every library for combinatory parsing, and all of these libraries work in the same way: the user can utilize *primitives* or “building blocks” for parsers, and then combine them to build bigger, more sophisticated parsers.

One of the simple primitive parsers is a parser for parsing a single digit at the beginning of the string, here called *digit*.

```
parse digit "123"  
-- Just ('1', "23")
```

It takes a string of characters and returns a pair of the parsed digit and the remaining characters, which the user may want to subsequently parse with another parser.

Another parser might be called `char` and parse a specified character at the beginning of the input string. Here it is specialized to parse a character “a”, but since there’s no character “a” at the beginning of “bcd” it fails and returns `Nothing`.

```
parse (char 'a') "bcd"
-- Nothing
```

The code snippet shown below presents a way to combine these primitive parsers into more complex ones. For example, the `<|>` operator tries to parse either a `digit`, or if it fails it then tries to parse a `letter`. Then, the `multiple` combinator takes a parser and applies it as many times as it can to the input string.

```
parse (multiple (digit <|> letter)) "abc123"
-- Just ("abc123", "")
```

Consider the following grammar defined in EBNF:

```
expression = term, '+', expr | term;
term        = factor, '*', term | factor;
factor      = '(', expression, ')' | digit;
digit       = '0' | '1' | '2' | '3' | '4' | '5' | '6' | '7' | '8' | '9';
```

The expression production rule could be defined with parser combinators like so

```
expression = do x <- term
              char '+'
              y <- expression
              return (Addition x y)
              <|> term
```

where `term`, `expression`, and `char` are different parsers.

The `do` keyword in the example above signifies a “do” block, where several parsers are applied to the input sequentially. As all values returned by `term`, `expression`, and `char` parsers may result in an error (the `Nothing` value), it would be reasonable to not check subsequent parsers after encountering a `Nothing` value, as it is guaranteed the whole parser already failed. “Do” blocks are a concise way to represent just that — a sequential operation which might result in an error, without having to check if the `Maybe` values are `Just` values or `Nothing` values at every step. If any of the values are `Nothing`, the whole `do` expression will result in a `Nothing`.

The key observation here is that the production rule looks basically the same as the parser. With some more of these smaller parsers, a complete parser for arithmetic expressions could be implemented. This is the idea of parser combinators — parsers are, in principle, functions. The user defines a library with some basic “building blocks” or primitives and some combining forms that let them put these things together. Parsers wrote in this fashion look very similar to the grammars that are written to describe languages [13, 25].

2.5.2. Abstract syntax tree

Abstract syntax trees (AST) are data structures that represent the structure of program code, and are the result of the syntax analysis phase of a program. The design of an AST is usually tightly linked with the design of a parser and its expected features [2].

The definition of an abstract syntax tree for EBNF is closely related to the definition of Parsing expression grammars discussed in Section 2.1.3. Like PEGs, EBNF provides such constructs as sequences, alternatives, optionals, and repeated expressions. Additionally, EBNF can represent *factors*, which are repeated expressions with a number of repetitions specified, as well as *exceptions*, which are excluding alternatives that impose certain restrictions on expressions. All EBNF constructs can be represented as the nodes of the AST, which are listed in Table 2.2.

Table 2.2. Node types for the abstract syntax tree of EBNF.

Node name	Normal representation
Alternative	" $e_1 \mid e_2 \mid \dots \mid e_n$ ", where e_i is an AST node
Sequence	" e_1, e_2, \dots, e_n ", where e_i is an AST node
Optional	"[e]", where e is an AST node
Repeated	"{ e }", where e is an AST node
Factor	" $n * e$ ", where n is an integer and e is an AST node
Exception	" $e_1 * e_2$ ", where e_1 and e_2 are AST nodes
Non-terminal	A non-terminal token
Terminal	A terminal token
Special	A special token
Empty	The empty string ε

2.5.3. The parsing process

Figure 2.2 shows the entire parsing process, starting with the textual representation of the EBNF grammar. The lexical analyser produces a string of tokens, which are then provided to the syntactic analyser, which produces an abstract syntax tree, or some sort of an intermediate representation of the abstract syntax tree. The semantic analyser analyses the intermediate representation and produces a final AST, which represents the grammar the user provided. At any point, the lexer, parser, or the processor may produce an error in the form it's been given to.

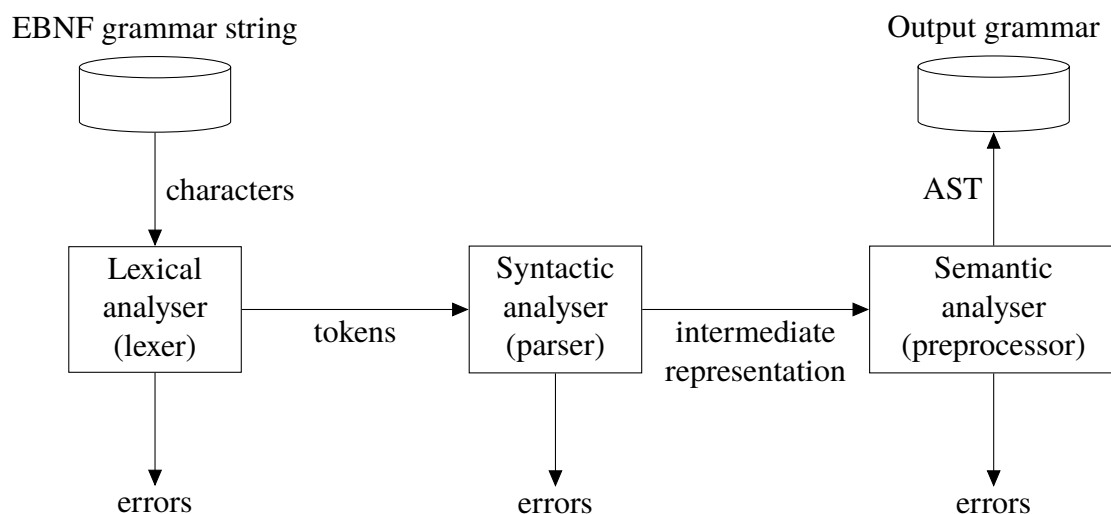


Figure 2.2. The visualization of the EBNF parsing process.

3. Analysis of similar solutions

Regex101

Regex101 [10] is a web-based tool that lets the user build regular expressions and debug them in real-time. Users are able to construct their expressions and see how it affects a data set all in a single screen at the same time. The tool was made by Firas Dib, with contributions from many other developers. It is said to be the largest regex testing service in the world.

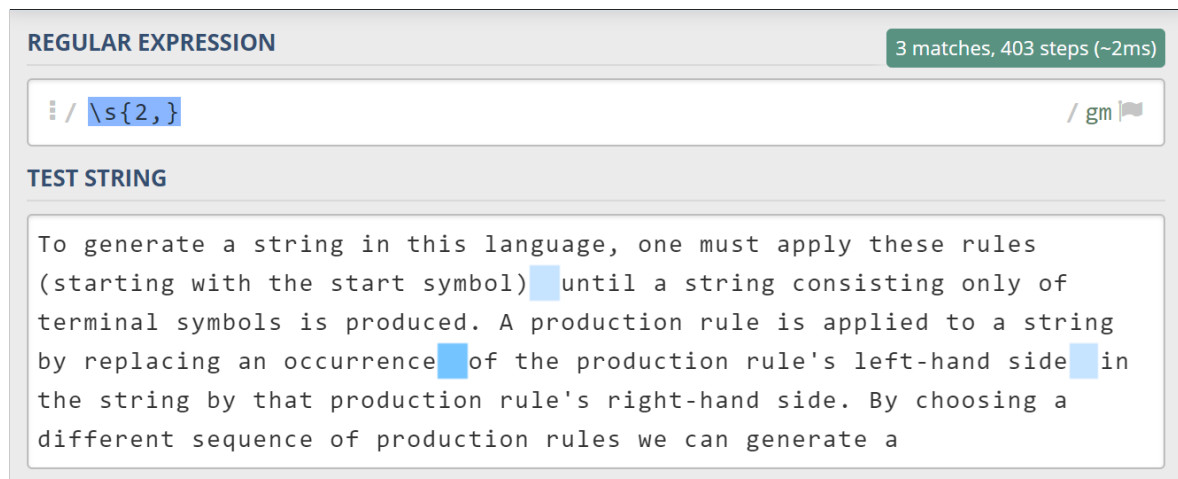


Figure 3.1. Screenshot of the Regex101’s matching functionality. The user provided the “\s\s+” regular expression, which matched every occurrence of two or more consecutive space characters in the test string.

The tool is available to users in the form of a web application and can be accessed from <https://regex101.com/>. It lets users build expressions fast and debug them along the way, for example by pasting in a set of data and then, through trial and error, building an expression with desired behavior. Figure 3.1 shows a typical usage of Regex101 — matching a pasted test string to a regular expression. The tool makes it clear if data is matching the expression or not, it even notifies users when the expression is broken, and gives some explanation of why it is not working, as seen in Figure 3.2.

These two feedback mechanisms are essential if the user is not accustomed to the regular expression language, or just does not know how to construct the appropriate expression yet. Being able to follow each step of the expression is convenient when users are not able to figure out why something is not working, or even if they are simply interested in discovering more about regular expressions. Getting this immediate feedback without Regex101 would have required users to write their expressions in a text editor and then run the code separately, without getting much feedback about why it is or isn’t working — Regex101.com eliminates this step.

Not only does Regex101.com make it easy to construct expressions, discover errors, and



Figure 3.2. Screenshot of a basic error in the regular expression reported by Regex101.

even get accustomed to the syntax, it makes looking up a metacharacter in regular expressions very easy. Always present, unless users minimize it, the *Quick Reference* tool lets them search for any token or character they need. Finally, Regex101 lets users switch which *flavor* or version of regular expressions they want to use, as they may need to integrate a regular expression into any number of other programming languages such as Python, JavaScript, Golang, etc. Regex101 has the ability to switch the version of the testing environment and will generate the code in that language for the user to use in other projects.

Parser-parser takes a lot of inspiration from Regex101 when it comes to availability — it's a web application, where all the work is done client-side. The user does not have to install any additional software except the web browser, the web application is accessed through a web page. In spite of its similar nature, Regex101 cannot be a replacement of Parser-parser — it focuses on various dialects of regular expressions rather than parsing EBNF and generating parse trees — it does, however, influence it with its accessibility and functionalities.

Pest

Pest [29] is a general purpose parser for the Rust [33] programming language. It uses its own dialect of *parsing expression grammars* as input, similarly to Parser-parser. Pest addresses the problem of hand-written parsers in Rust, which in some circumstances can become hard to maintain by their developers. Writing a specialized, domain-specific parser for a language can become tedious, so developers usually gravitate towards using a grammar-generated parser. This allows the developers to focus on the definition of the language, rather than on the implementation of the parser. Grammars that define the language offer better correctness guarantees, and any issues with the language can be fixed declaratively in the grammar itself. Memory safety of Rust additionally limits the damage bugs could cause. Static analysis and a precise low-level implementation form a stable foundation on which major performance tuning is feasible.

Developers of Pest, in spite of focusing mainly on the functionalities in the Rust programming language, also provide an online editor available from the browser on the Pest's homepage (<https://pest.rs/#editor>). The online editor allows potential future users of Pest to experience the syntactic characteristics of the Pest's dialect of PEGs and its error reporting capabilities. The editor will inform the user about any syntactic errors, as well as errors of semantic nature, such as undefined or left-recursive production rules (seen in Figure 3.3) and highlight them in their exact locations.

After parsing the grammar, Pest provides a window, which acts as an input console, where users can type string that may or may not be parsed by the parser generated by Pest. Additionally, users can choose the initial production rule from a dropdown menu, which



Figure 3.3. Screenshot of Pest's online editor example error report.

is an interesting choice, as opposed to automatically detecting the initial rule based on the dependency graph of production rules. The output window presents the parse tree, or the errors encountered in the input string in case there are any. These features can be seen in Figure 3.4.

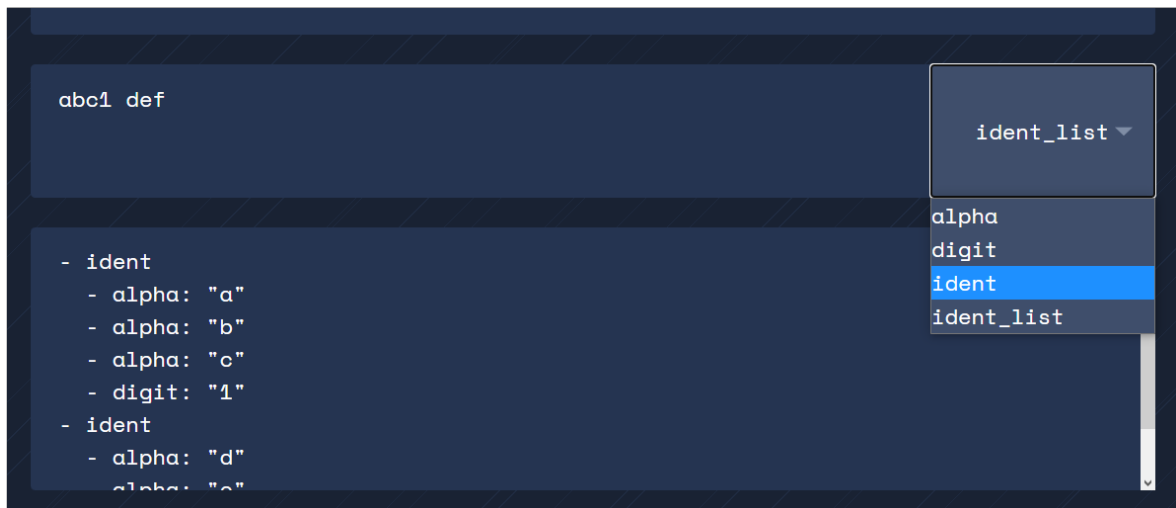


Figure 3.4. Screenshot of the input and output windows in Pest's online editor.

While Pest focuses mainly on its integration with the Rust programming language, this is not the case for Parser-parser, which aims to provide all of its functionality inside the web application. The online editor of Pest serves largely as a “try me” feature for new users, rather than a reliable tool. The editor lacks the standard editor features, such as autocompletion, code folding, search and replace interface, as well as bracket and tag matching, all of which Parser-parser does provide. The parse tree in the output window is shown in a basic textual form, without any interactive capabilities, which the user may value. Finally, while Pest's grammar is based on PEGs, and is similar in nature to EBNF, it is, in fact, not EBNF. The whole point of using EBNF and other notations discussed in Section 2.2 is that they're standardized, well-known, and accepted by the community; Pest's syntax is known only to users of Pest and requires them to learn a new, non-standard language just for the purpose of parsing grammars, where other, already established languages may have sufficed.

4. Design of the project

This chapter introduces a specification for the application described in Chapter 1. The specification is presented in forms of a list of functional and non-functional requirements in Section 4.1, and user stories in Section 4.2. Section 4.3 describes use cases and their descriptions structured in the form of a use case diagram in the Unified Modeling Language, as well as their example scenarios, also presented as activity and sequence diagrams. The chapter describes the architecture of the system from the logical and physical perspective as component and deployment diagrams in Section 4.4. The chapter does not cover any class or database diagrams, as the implementation of this project and its functionally-oriented nature, as opposed to being object-oriented, does not require them. Finally, the chapter concludes with the prototype and sketches of the user interface for the web application in Section 4.5.

4.1. Requirements

4.1.1. Functional requirements

Functional requirements shown in Table 4.1 define functionalities and features of the system. Each requirement is associated with a certain priority.

Table 4.1. The functional requirements of the project, their features, and priorities.

Id	Requirement	Features	Priority
<i>FR1</i>	Specifying the grammar	The user can specify the grammar of a given language in the EBNF notation by providing it in a textual form in a designated editor window.	high
<i>FR2</i>	Error reporting	The editor provides feedback about any syntactic or semantic ¹ errors encountered during the parsing by highlighting the exact location of the error in the provided grammar. The user can then hover the mouse pointer over the highlighted area to read the error message.	high
<i>FR3</i>	Specifying the input string	The user can specify the input string in a designated editor window to check if it belongs to the language generated by the previously-defined grammar.	high
<i>FR4</i>	Visualizing the parse tree	The application visualizes the parse tree resulting from parsing the specified input string with the parser generated by the grammar defined by the user.	high

¹Such as production rule duplication or left recursion.

<i>FR5</i>	Syntax highlighting	The editor highlights parts of the specified grammar with a different syntactic meaning in a different manner with the use of multi-colored fonts.	medium
<i>FR6</i>	Autocompletion of non-terminals	The editor predicts the identifier of a non-terminal a user is typing by providing a list of possible non-terminals, which then can be chosen by the user.	low
<i>FR7</i>	Production rule folding	The editor provides the ability to hide and reveal a production rule of the grammar inside the editor window.	low
<i>FR8</i>	Search and replace interface	The user can search for any occurrences of a phrase in the editor window and possibly replace them with a different phrase. The search and replace functionality should also support regular expressions.	low

4.1.2. Non-functional requirements

Table 4.2 describes requirements of the non-functional nature of the system, which focus on aspects of usability, availability, and compatibility of the system.

Table 4.2. The non-functional requirements of the project and their priorities.

Requirement	Priority
The web application should be available 24 hours a day, 7 days a week.	medium
Page loading time should be less than 1 second with internet download speed of 80 Mbps. Parsing and checking times should both be less than 50 milliseconds.	high
The application must work and display correctly in <ul style="list-style-type: none"> • Chrome version 86 or later, • Safari version 14 or later, • Edge version 86 or later, • Firefox version 82 or later, • Opera version 71 or later. 	high
	medium
The source code of the product should be open source and freely available for possible modification and redistribution.	high
The project should include the documentation necessary for extension and maintenance of the system.	high
The system should provide high degree of integrability with future components which extend the functionalities of the system.	high

4.2. User stories

Stories in Table 4.3 are short descriptions of a feature told from the perspective of the person who desires a new functionality in the system.

Table 4.3. The user stories.

Id	User story
US1	As the user, I want to be able to paste the contents of my clipboard into the editor window in the application.
US2	As the user, I want to be able to type in the editor window with my keyboard.
US3	As the user, I want to be able to appreciate the multi-colored appearance of the text that represents the syntax that I provided.
US4	As the user, I want to be able to select a portion of the text in the editor window and copy it to the clipboard using a keyboard shortcut.
US5	As the user, I want to be able to hold the <i>Alt</i> key on my keyboard to define multiple cursors in the editor window.
US6	As the user, I want to have the ability to autocomplete the non-terminal I am typing that has already been declared elsewhere in the code.
US7	As the user, I want to be able to hide any existing production rules that might appear too long, to increase the degree of clarity and readability of the grammar I'm working on.
US8	As the user, I want to be able to show any previously hidden production rules of the grammar.
US9	As the user, I want to have the ability to press a certain key combination on my keyboard that would allow me to type a specific phrase in the popup window, which would then find all the occurrences of that phrase in the editor window.
US10	As the user, I want to be able to provide a regular expression for the <i>find</i> functionality that would allow me to find all occurrences of phrases that pattern match that specific regular expression.
US11	As the user, I want to be able to replace some of the occurrences of phrases found with the <i>find</i> functionality with another phrase provided in a popup window.
US12	As the user, I want to be able to specify the initial production rule in the process of checking the input string against the grammar I provided.
US13	As the user, I want to be able to see errors in the syntax of the provided grammar in the form of underlined text in the location of where the errors actually occur.
US14	As the user, I want to have the ability to hover the mouse pointer over the underlined text to read the error message at that location. Alternatively, I want to be able to hover over the error indicator, which appears next to the line number.
US15	As the user, I want to be able to see the parse tree of the recognized input string that I provided.
US16	As the user, I want to have the ability to collapse any nodes in the visualized parse tree that might appear too long.

4.3. Use case specification

4.3.1. Use cases

Figure 4.1 shows the use case diagram of the system. Each use case also presented in Table 4.4 along with a short description.

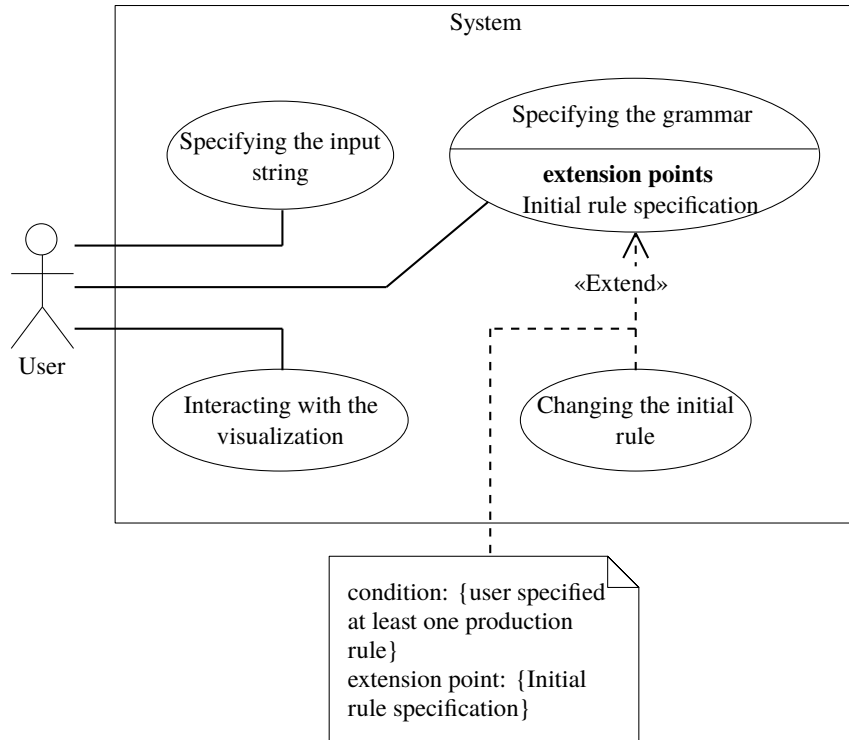


Figure 4.1. The use case diagram.

Table 4.4. Descriptions of the use cases.

Id	Name	Description
UC1	Specifying the grammar	Allows the user to specify the grammar of a given language in the EBNF notation by providing it in a textual form in a designated editor window.
UC2	Specifying the input string	Allows the user to specify the input string in a designated editor window to check if it belongs to the language generated by the previously-defined grammar.
UC3	Interacting with the visualization	Allows the user to observe the visualized parse tree of the provided input string and interact with it by expanding and collapsing the tree nodes.
UC4	Changing the initial rule	Allows the user to specify the initial production rule used in the process of checking the provided input string against the defined grammar.

4.3.2. Requirements traceability graph

Figure 4.2 presents the relationship between functional requirements, user stories and use cases in the form of a requirements traceability graph. It shows that every user story is connected with at least one functional requirement and vice versa, and that every use case is associated with at least one user story and vice versa.

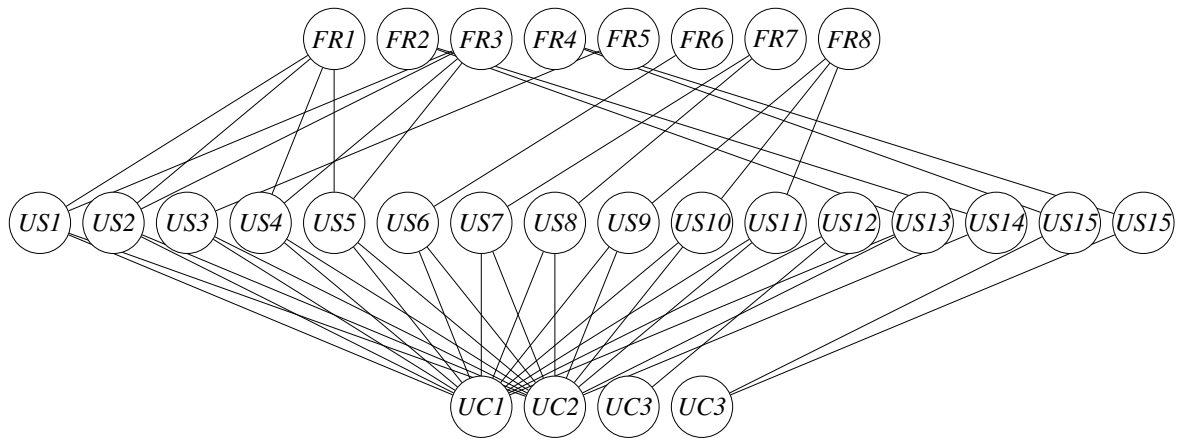


Figure 4.2. The requirements traceability graph.

4.3.3. Use case scenarios

Tables 4.5, 4.6, 4.7, and 4.8 describe the scenarios of each use case in the system. Every scenario is defined by its pre-conditions, its post-conditions, and a list of steps made by the system or the user required to complete it.

Table 4.5. Use case scenario of *UC1* Specifying the grammar.

Identifier	<i>UC1</i>
Name	Specifying the grammar
Summary	Allows the user to specify the grammar of a given language in the EBNF notation by providing it in a textual form in a designated editor window.
Pre-conditions	None.
Post-conditions	The grammar has been correctly defined by the user with no syntactic errors.
Main scenario	<ol style="list-style-type: none"> 1. The system shows a grammar editor window to the user. 2. The user provides a syntactically and semantically correct definition of a grammar. 3. The system shows an icon indicating no errors detected in the grammar. <p>End of scenario.</p>
Alternative scenario	<ol style="list-style-type: none"> 2a.1. The user provides an invalid definition of a grammar. 2a.2. The system highlights the text in the grammar editor window at the error location. <p>Return to step 2.</p>

Table 4.6. Use case scenario of *UC2* Specifying the input string.

Identifier	<i>UC2</i>
Name	Specifying the input string
Summary	Allows the user to specify the input string in a designated editor window to check if it belongs to the language generated by the previously-defined grammar.
Pre-conditions	None.
Post-conditions	The input string has been correctly entered by the user.
Main scenario	<ol style="list-style-type: none"> 1. The system shows a input string editor window to the user. 2. The user provides a desired input string. 3. A valid grammar has been provided by the user in the grammar editor window. 4. The system shows the result of the checker in the result window. <p>End of scenario.</p>
Alternative scenario	<p>3a.1. The user did not provide a valid grammar in the grammar editor window.</p> <p>3a.2. The system does not show a result of the checker.</p> <p>End of scenario.</p>

Table 4.7. Use case scenario of *UC3* Interacting with the visualization.

Identifier	<i>UC3</i>
Name	Interacting with the visualization
Summary	Allows the user to observe the visualized parse tree of the provided input string and interact with it by expanding and collapsing the tree nodes.
Pre-conditions	The user has provided a valid definition of a grammar, as well as an input string, that belongs to the language generated by that grammar.
Post-conditions	None.
Main scenario	<ol style="list-style-type: none"> 1. The system shows the result of a checker in the form of a visualized parse tree. 2. The user clicks the nodes of the parse tree to collapse or expand them. <p>End of scenario.</p>

Table 4.8. Use case scenario of *UC2* Specifying the input string.

Identifier	<i>UC4</i>
Name	Changing the initial rule
Summary	Allows the user to specify the initial production rule used in the process of checking the provided input string against the defined grammar.

Pre-conditions	The user has provided a valid definition of a grammar.
Post-conditions	The initial production rule has been successfully changed to the desired one.
Main scenario	<ol style="list-style-type: none"> 1. The system shows a button the current initial production rule written on top. 2. The user clicks on the button. 3. The system shows a dropdown menu with a list of all production rules defined in the provided grammar. 4. The user clicks on an item of the list corresponding to the desired initial production rule. 5. The system changes the identifier of the initial production rule on the button. <p>End of scenario.</p>

4.3.4. Activity diagrams

Figures 4.3, 4.4, 4.5, and 4.6 are the graphical representations of use case scenarios defined in Section 4.3.3, represented in the form of UML activity diagrams.

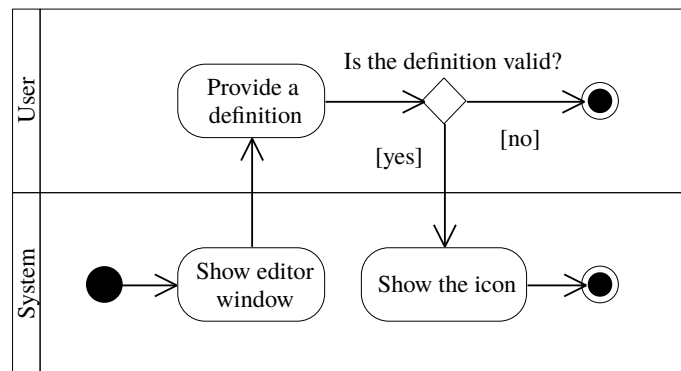


Figure 4.3. The activity diagram of *UC1* Specifying the grammar.

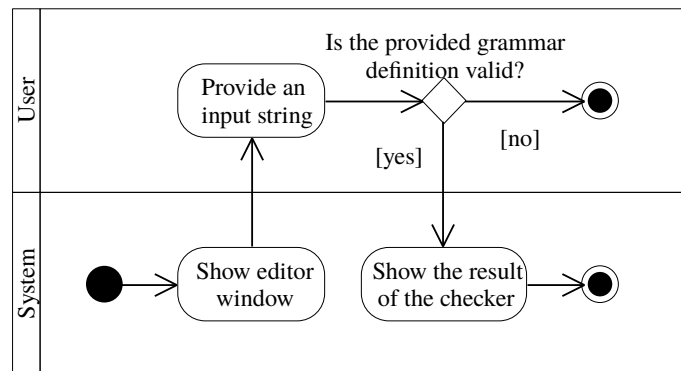


Figure 4.4. The activity diagram of *UC2* Specifying the input string.

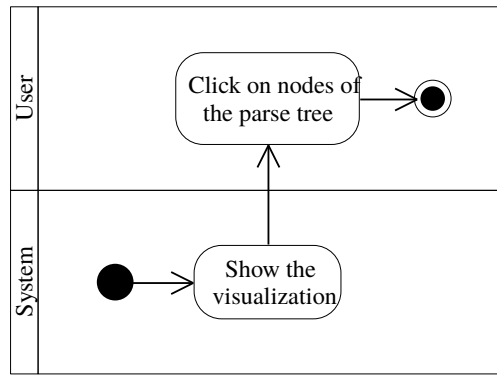


Figure 4.5. The activity diagram of *UC3* Interacting with the visualization.

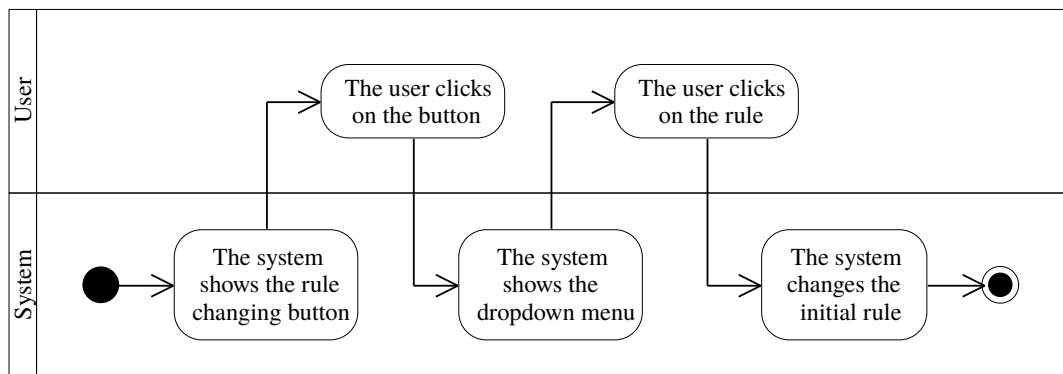


Figure 4.6. The activity diagram of *UC4* Changing the initial rule.

4.3.5. Sequence diagram

Figure 4.7 shows a sequence diagram, that is in essence an interaction diagram that details how operations in the system are carried out and visualizes interactions between objects and components. It captures interactions from every use case, all of which were defined in Section 4.3.1.

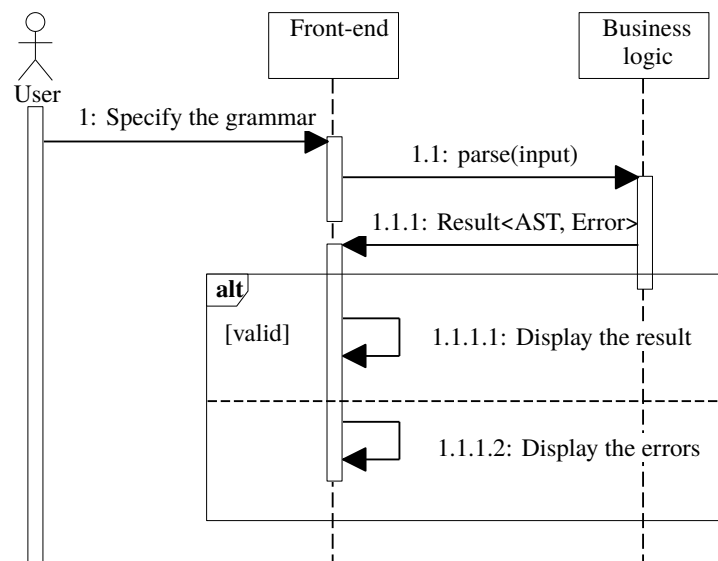


Figure 4.7. The sequence diagram representing the specification of the grammar.

4.4. System architecture

4.4.1. Logical architecture

Logical architecture of a system can be represented by UML component diagrams, which focus on a system's components that are often used to model the static implementation view of a system. A component diagram divides the system into various high levels of functionality. Each component is responsible for one distinct aim within the entire system and only communicates with other necessary elements on the so-called need-to-know basis. In a system with a functional-oriented approach it is more suitable for modelling interactions between components. The logical architecture of Parser-parser is modelled with such a diagram and can be seen in Figure 4.8.

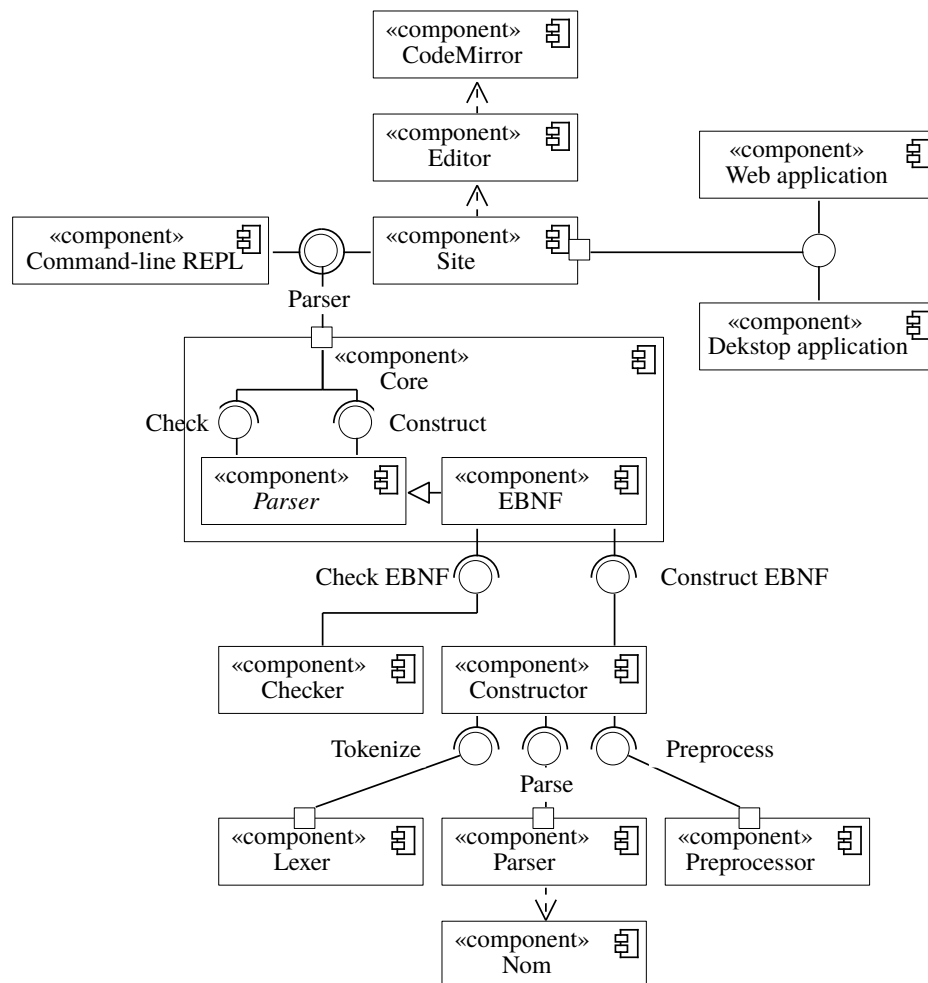


Figure 4.8. The logical architecture of the system represented with a UML component diagram.

4.4.2. Physical architecture

A deployment diagram in the Unified Modeling Language models the physical deployment of artifacts on nodes and can represent a physical architecture of a system. Diagram shown on Figure 4.9 visualizes the architecture for Parser-parser.

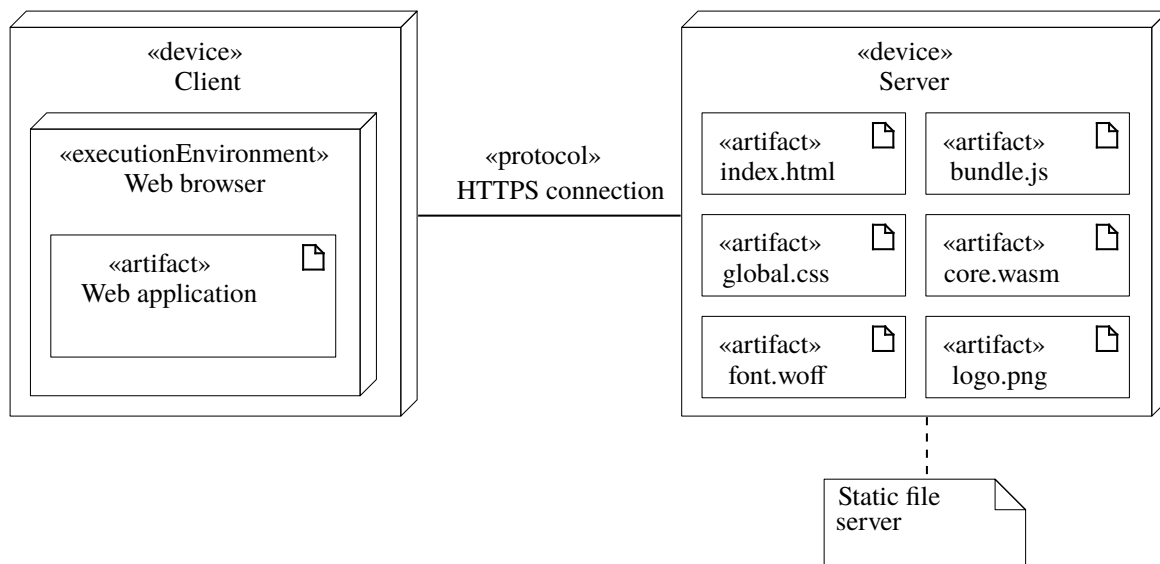


Figure 4.9. The physical architecture of the system represented with a UML deployment diagram.

4.5. Interface prototype

Because of the visual nature of the web application, a prototype of the user interface design should be established to allow the developer to plan out the implementation of the front-end aspect of the application. Parser-parser, being a rather simple application, will consist of a single view (as seen in Figure 4.10), which is made out of several components.

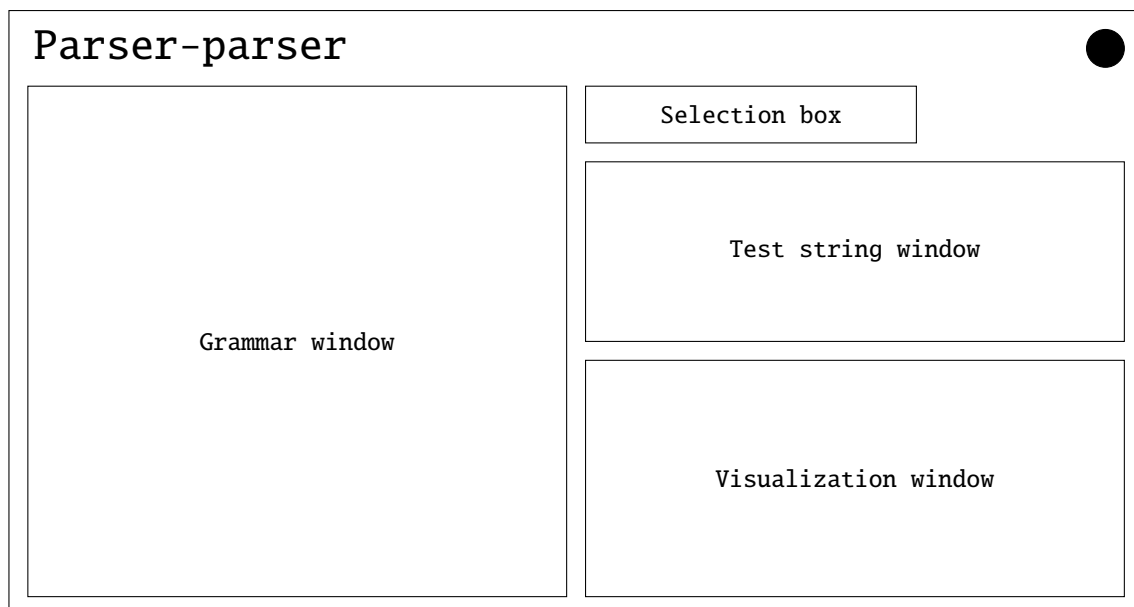


Figure 4.10. The user interface sketch.

5. Implementation of the project

5.1. Software environment

5.1.1. Used technologies

Visual Studio Code

Visual Studio Code [27] is a free, open-source text editor made by Microsoft for Windows, Linux and macOS. It is designed to write code and features syntax highlighting, code completion, snippets, code refactoring, and code debugging. The editor can be used with various programming languages, and supports extensions, which can be installed through a central repository called VS Code Marketplace available in the editor itself. The extensions may provide feature additions to the editor, as well as the support for various programming languages in the form of code linters, static code analysers, and debuggers. The editor is integrated with various version control systems, including Git and Subversion

According to the 2019 Developers Survey of Stack Overflow, Visual Studio Code ranked #1 among the top popular developer tools, with 50.7 % of the 87317 respondents using it [37].

The extensions for the editor are made by the members of Visual Studio Code community. Two main extensions used by the author to develop the project were:

rust-analyzer [12] — An implementation of the Language Server Protocol for the Rust programming language, which provides features such as code completion, messages for syntax and semantic errors, code actions, diagnostics, “go to definition” and other editor actions.

Svelte for VS Code [39] — An implementation of the Language Server Protocol for the Svelte framework. The extension provides diagnostic messages for warnings and errors, support for Svelte pre-processors that provide source maps, as well as the support for Svelte-specific formatting (via prettier-plugin-svelte). Besides the Svelte language, the extension supports features such as hover info, messages for syntax and lint errors, and autocompletions for HTML, CSS/SCSS/LESS, as well as TypeScript and JavaScript.

The extensions have not proven to be crucial for the development of the project, but were an excellent addition to the workflow.

Besides the editor extensions, the terminal integrated with Visual Studio Code editor has been a valuable feature throughout the development process. The command line is a substantial factor in the development of modern applications, so a built-in terminal window allows the user to swiftly switch between the code editor and the command line.

The support for the Git version control system has also been advantageous when it comes to code editing. Every added, modified, or removed line of code is highlighted with an appropriate color in the code editor. This greatly improves the readability of the code, and

allows the users to revert the code to its previous state right from the editor without any external tools.

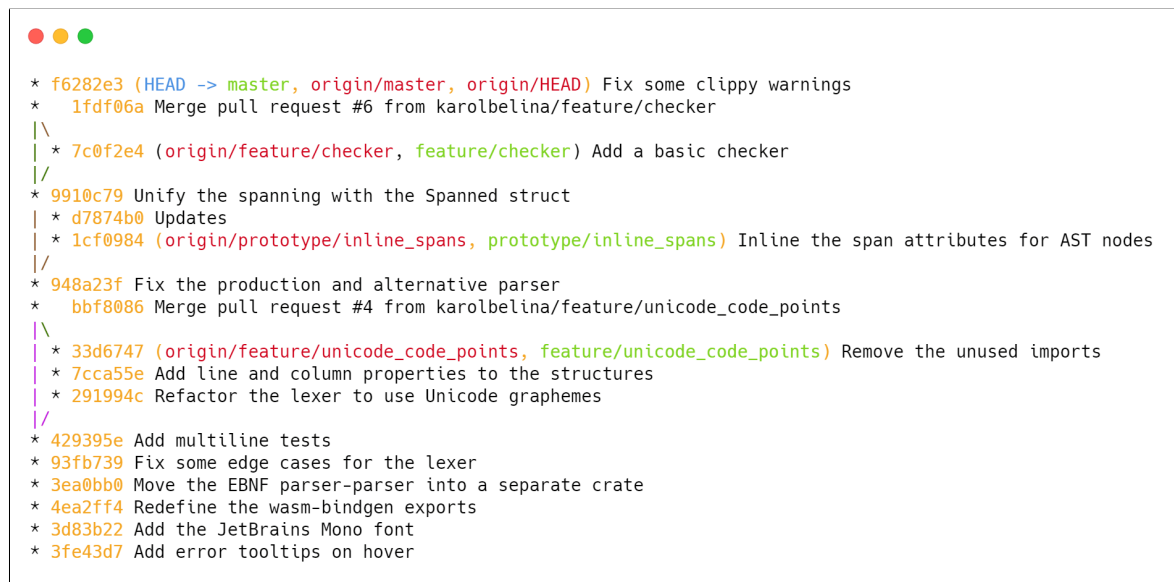
Git

Git [16] is a free and open source distributed version control system. It has been a major part of the development process for the project, and has been used mainly as a tool for keeping track of the changes made to the source code and for integrating features in a smooth, non-disruptive manner.

Git supports branching and merging, which means that several project features may be implemented simultaneously and independently on separate *branches* and then *merged* into the main project. Every major code change has been implemented on a designated branch and was merged into the main branch only after a thorough testing process — this has made parallel development very easy, by isolating new development from finished work. This style of a workflow is known as GitFlow, made popular by Vincent Driessen [11], it has shown itself to be very effective for projects of any scale. Efficient switching between different versions of project files enables developers to work effectively on the project. Git includes specific tools for visualizing and navigating a non-linear development history. The author used [4] as a reference for using the tool.

Git is now the most widely used source-code management tool, with 87.2 % of the 74298 respondents of the 2018 Developers Survey of Stack Overflow reporting that they use Git as their primary source-control system [36].

The main client of Git used in the project was the command-line tool on the Ubuntu operating system running on Windows Subsystem for Linux. Figure 5.1 shows an example of GitFlow's *feature branches* and changes in the project repository in the Git version control system.



```
* f6282e3 (HEAD -> master, origin/master, origin/HEAD) Fix some clippy warnings
* 1fdf06a Merge pull request #6 from karolbelina/feature/checker
| \
| * 7c0f2e4 (origin/feature/checker, feature/checker) Add a basic checker
| /
* 9910c79 Unify the spanning with the Spanned struct
| * d7874b0 Updates
| * 1cf0984 (origin/prototype/inline_spans, prototype/inline_spans) Inline the span attributes for AST nodes
| /
* 948a23f Fix the production and alternative parser
* bbf8086 Merge pull request #4 from karolbelina/feature/unicode_code_points
| \
| * 33d6747 (origin/feature/unicode_code_points, feature/unicode_code_points) Remove the unused imports
| * 7cca55e Add line and column properties to the structures
| * 291994c Refactor the lexer to use Unicode graphemes
| /
* 429395e Add multiline tests
* 93fb739 Fix some edge cases for the lexer
* 3ea0bb0 Move the EBNF parser-parser into a separate crate
* 4ea2ff4 Redefine the wasm-bindgen exports
* 3d83b22 Add the JetBrains Mono font
* 3fe43d7 Add error tooltips on hover
```

Figure 5.1. Screenshot of the command-line interface of the Git version control system.

GitHub

GitHub [17] is a for-profit company owned by Microsoft that offers a cloud-based Git repository hosting service. As a company, GitHub makes money by selling hosted private

code repositories, as well as other business-focused plans that make it easier for organizations to manage team members and security. The author used the free GitHub plan as the service for hosting the project's Git repository. Having the source code on an external server protected the project against data loss and allowed the developer to work on the project from any device at any convenient time.

In addition to using GitHub as a hosting service, one can also exploit its project management features. Developers can define project boards related to the project's code repository, which are simple kanban [1] boards that can help organize and prioritize the work. With projects, the developers have the flexibility to manage boards for an entire project, or just for specific features. Figure 5.2 shows an example project board from Parser-parser.

Project boards contain *issues* and *pull requests*, which can be moved from one kanban column to another, indicating that some work is currently “to do”, work in progress, or complete. These work “cards” contain information about the author, assignees, the status, as well as simple textual notes. The *issues* are a way of reporting ideas, bugs, enhancements, or tasks natively on GitHub. After completing the work on an issue, a developer might make a *pull request* to allow other developers on the project to review and discuss the changes made to the code, and then deploy the changes by “pulling” the code to the central code repository.

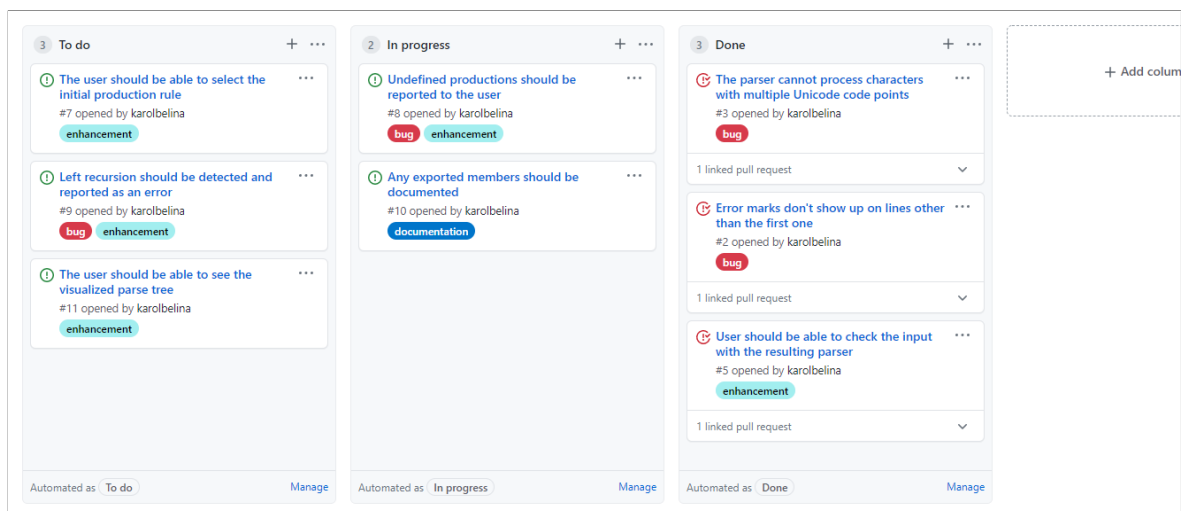


Figure 5.2. Screenshot of one of the project's kanban boards on GitHub.

GitHub supports Continuous integration and Continuous Delivery functionalities in form of *Actions* and *Pages*. GitHub Actions are a way to automate and execute any software development workflow after any change to the code in the repository. The user may set up many various actions for testing the changes on many development environments and operating systems at the same time, as well as building and deploying the code as a package or an arbitrary artifact. An action consists of jobs, which are defined by a list of steps required to execute them.

The GitHub Actions are used by the author to automate the testing and build process on every change made to the code repository. The built application is then deployed to a static site hosting service called GitHub Pages, which integrates itself seamlessly with Actions and GitHub repositories. GitHub Pages allows the user to host a website directly from a GitHub repository by combining static HTML, CSS, JavaScript, and other files straight from a repository into a website and publishing it on a `github.io` domain or a custom one.

Rust

Rust [33] is the main programming language used in Parser-parser — it powers the business logic part of the project. The language has been the most loved language for four years in a row in the Stack Overflow’s survey [37]. The core idea of the language is memory safety — the language enforces certain rules checked at compile time, which guarantee that the program is safe from bugs like dereferencing null or dangling pointers, as well as making it difficult for the programmer to leak memory. Rust does this through a system of ownership and borrowing. The language, besides the safety, focuses on speed — its design lets the developer design programs that have the performance and control of a low-level language, but with the powerful abstractions of a high-level language.

Rust’s design borrows heavily from the one of Haskell [19] — both languages feature a rich type system, both are immutable-by-default, avoid mutation of shared references et cetera. Many developers tend to write Rust code in a functional style and adhere to the principles of functional programming, even though the language is multi-paradigm.

Without the need of a garbage collector, Rust projects are well-suited to be used as libraries by other programming languages. The language over the last few years has manifested itself in several distinct domains, including command-line tools, networking, and embedded systems. Rust is supported on multiple operating systems and targets multiple platforms, has notable documentation, a user-friendly compiler with convenient error messages, and excellent tooling and ecosystem. For referencing the language, the author used [24], which covers many features and concepts of Rust.

WebAssembly

WebAssembly [43] (abbreviated *Wasm*) is a safe, portable, low-level code format designed for efficient execution and compact representation. Its main goal is to enable high performance applications on the Web, working alongside JavaScript, but not to be a replacement of it. Its design makes it portable, compact, and execute at native or near-native speeds. Although it has currently attracted attention in the JavaScript community and other Web communities in general, Wasm does not make any assumptions about its host environment. WebAssembly is supported as a target for many programming languages, including C# via Blazor, C++ via EmScripten, and the main language used in Parser-parser — Rust. The author compiles Rust code to WebAssembly to be then used in a web environment for several reasons:

- Code size is an important aspect, as the `.wasm` file must be downloaded over the network. Rust does not require a runtime and does not come with a garbage collector, allowing for small Wasm sizes, as there is no extra code included,
- Rust and WebAssembly integrate with existing JavaScript tooling — the developer can continue using the tooling they are already familiar with, like npm and Webpack,
- JavaScript Web applications struggle to achieve and maintain stable performance. The code is required to be ran frequently, so Wasm can solve this kind of problem with greater memory and CPU performance at a lower level compared to the interpreter of JavaScript.
- The Rust language itself, with a strong package manager, high performance, memory safety, and zero-cost abstractions.

Cargo

Cargo [32] is the Rust’s package manager. It downloads the dependencies of a Rust package and compiles them, making sure that the developer will always be presented with a repeatable build. To achieve this, Cargo includes two metadata files with various information about the packages, it downloads and builds the dependencies, runs the Rust compiler with appropriate parameters to build the package, and introduces various conventions to make working with Rust packages more comfortable.

Rust provides first-class support for unit and integration testing, and Cargo allows the developer to execute all tests with a single command. Additionally, Cargo allows the developer to install extensions, which enhance the workflow and the development process. One of extensions useful for the author was Clippy — a collection of lints to catch common mistakes and improve the Rust code.

`crates.io` is the primary package registry of the Rust community that is used to discover and download packages. It is used by Cargo as a default way to find requested packages. The project uses several dependencies, the most important of which include:

`nom` [7, 8] — A combinatory parsing library for the Rust programming language. It provides tools to build memory-safe parsers without compromising memory consumption or speed. To achieve this, Nom uses Rust’s strong typing and memory safety very extensively to deliver performant parsers and provides various functions, macros, and traits to hide most of the error-prone details.

Language parsers are often written manually for extra flexibility and performance. Nom can be (and has been) successfully used as a parser for prototyping a language. The resulting code of the parser is concise and resembles the grammar the developer would have written with other parser approaches. Parsers written with Nom are small and easy to maintain, as well as easy to test separately.

`unicode-segmentation` [41] — A library with a set of iterators which split strings on *grapheme clusters*, *words* or *sentence boundaries*, according to the Unicode Standard Annex #29 [9] rules.

`wasm-bindgen` [31] — A Rust library and CLI tool that makes high-level interactions between Wasm modules and JavaScript easier. Besides integers and floats, this library allows JavaScript and Wasm to communicate with strings, JS objects, classes, and other structures. Notable features of this project include:

- Importing JS functionality into Rust such as DOM manipulation, console logging, or performance monitoring.
- Working with rich types like strings, numbers, classes, closures, and objects.
- Automatically generating TypeScript bindings for Rust code being consumed by JS.

Wasm-bindgen only generates bindings for the JavaScript imports that are actually being used, and Rust functionality that is being exposed.

`quickcheck` [15] — A property-based testing framework inspired by the QuickCheck framework for Haskell. The crate comes with the ability to randomly generate and shrink integers, floats, tuples, booleans, lists, strings, options and results. All QuickCheck needs is a property function — it will then randomly generate inputs to that function

and call the property for each set of inputs. If the property fails (whether by a runtime error like index out-of-bounds or by not satisfying the property), the inputs are "shrunk" to find a smaller counter-example.

criterion [20] — Provides a powerful but simple way to measure software performance. It provides both a framework for executing and analyzing benchmarks and a set of driver functions that makes it easy to build and run benchmarks, and to analyse their results.

structopt [18] — Parses command-line arguments by defining a struct. StructOpt combines the Clap library with a custom derive for marking a struct containing values that will be translated into specific command-line arguments based on their types and markers. Besides the desired arguments, StructOpt automatically generates the *help* string, which can be invoked with the `--help` flag, as well as a `--version` flag for checking the version of the program.

All of the above dependencies are available under the MIT license.

Wasm-pack

Wasm-pack [42] is a tool for building and working with Rust-generated WebAssembly that the developer would like to interop with JavaScript, in the browser, or with Node.js. The generated WebAssembly packages then could be published to the npm registry, or otherwise used alongside any JavaScript packages in workflows that the developer already uses, such as Webpack or Rollup. The tool interoperates and utilizes the `wasm-bindgen`, another tool, to provide a bridge between the types of JavaScript and Rust. It allows JavaScript to call a Rust API with a string, or a Rust function to catch a JavaScript exception. `wasm-pack` wraps the CLI portion of the `wasm-bindgen` tool. This results in wrapping the WebAssembly module in JS wrappers which make it easier to interact with the module. `Wasm-bindgen` supports both ES6 modules and CommonJS and the developer can use `wasm-pack` to produce either type of package.

Svelte

Svelte [38] is a free and open-source front end JavaScript framework. Svelte has its own compiler for converting app code into client-side JavaScript at build time. The developer writes the components using HTML, CSS, and JavaScript, which Svelte compiles during the build process into small standalone JavaScript modules. Unlike frameworks like React and Vue, which do most of their computations on the client-side while the app is running, Svelte turns that process into a compile-time step that happens only when the developer builds the app, which produces highly-optimized JavaScript. The compiler can ensure the user's web browser does as little work as possible by statically analyzing the templates of the components. Thanks to this approach, Svelte does not only provide smaller application bundles and better performance, but also an experience that is more friendly for people that are not familiar with the modern tooling ecosystem. Svelte is particularly appropriate to tackle the following situations:

- Welcoming people with basic web development knowledge: Web developers with basic HTML, CSS, and JavaScript knowledge can easily learn Svelte in a short time and start building web applications.

- Interactive pages or complex state-based visualizations: If the user is building data-visualizations that require displaying a large number of elements, the performance gains thanks to no runtime overhead will ensure that user interactions are responsive.
- Web applications intended for low power devices: Applications built with Svelte have smaller bundle sizes, which is ideal for devices with slow network connections and limited processing power.

Svelte, being a compiler, can extend features of HTML, CSS, and JavaScript, and generate optimal JavaScript code with no runtime overhead. To accomplish this, Svelte extends core web technologies and only modifies them in very specific circumstances and only in the context of Svelte components.

Rollup

Rollup [30] is a module bundler for JavaScript which is able to compile smaller code modules into a complex library or application. It uses the standardized ES module format for code, which lets the developer smoothly and seamlessly connect individual functions and external modules. Rollup is able to optimize modules for faster loading in modern browsers, or output a legacy format of modules.

By splitting the project into smaller distinct pieces, the development process is oftentimes more straightforward, since that generally removes unexpected interactions and reduces the complexity of the issues the developer needs to work on — simply writing smaller projects isn't necessarily the answer to that. Historically, JavaScript has not included the feature of modules as a core feature in the language. This changed with the ES6 revision of JavaScript, which includes syntax for importing and exporting functions and data so they can be shared between separate scripts. The specification is now fixed, however, it is implemented only in modern browsers. Rollup allows the user to write code using the new module system, and will then compile it back down to existing supported formats such as CommonJS modules, AMD modules, and IIFE-style scripts. This means that the developer gets to write future-proof code.

Besides enabling the use of ES modules, Rollup also statically analyzes the imported code and will exclude anything that isn't used in the code. This allows the user to build on top of existing tools and modules without adding extra dependencies or increasing the size of the project. Rollup only includes the minimum, which results in lighter, faster, and less complicated libraries and applications. Since this approach can utilise explicit `import` and `export` statements, it is more powerful than running an automated *minifier* to discover unused variables in the compiled JavaScript code.

npm

Node Package Manager [28] is a package manager for the JavaScript programming language. It includes a command-line client, and an online database of public and paid-for private packages called the npm registry. The registry can be accessed through the client application, and the available packages can also be browsed and searched via the npm website.

Npm provides several built-in scripts and allows users to define their own. An npm script is a convenient way to bundle common shell commands for the project. They are usually commands, or strings of commands, which would normally be entered at the command line in order to do something with the application. Scripts are stored in a project's configuration file, which means they're shared amongst everyone using the codebase. They aid in automating

repetitive tasks and allow the user to have to learn fewer tools. Scripts also make sure that every programmer is using the same command with the same parameters. Common use cases for npm scripts include building the project, starting a development server, compiling CSS, linting, or minifying.

The project is dependent on several npm packages:

CodeMirror [6] — A versatile text editor implemented in JavaScript for the browser. It is specialized for editing code and includes several built-in language modes and addons that implement more advanced functionalities. A rich programming API and a CSS theming system are available for customizing CodeMirror to fit the needs of the user, as well as making extensions with new functionalities. CodeMirror is used in the developer tools for Firefox, Chrome, and Safari, in Light Table, Adobe Brackets, Bitbucket, and many other projects.

CodeMirror supports a wide variety of configurations — the basic version of the editor without any addons provides the support for over 100 languages, autocompletion, code folding, configurable keybindings, search and replace interface, bracket and tag matching, support for split view, linter integration, various themes, and many more.

svelte-tree [34] — A tree-like view component for Svelte. The component has the ability to display a collapsable tree structure based on a provided tree of JavaScript objects with custom name and children properties. The component provides a slot space to display custom nodes, which will give the tree node DOM/components the access to the nodes being rendered.

The configuration file also lists dependencies for the development process, which can be divided into several categories:

- Rollup and its plugins
 - rollup,
 - @rollup/plugin-alias,
 - @rollup/plugin-commonjs,
 - @rollup/plugin-node-resolve,
 - @wasm-tool/rollup-plugin-rust,
 - rollup-plugin-copy,
 - rollup-plugin-livereload,
 - rollup-plugin-svelte,
 - rollup-plugin-terser
- Svelte and its plugins
 - svelte,
 - svelte-check,
 - svelte-loader,
 - svelte-preprocess
- Miscellaneous dependencies
 - gh-pages,
 - rimraf,
 - sirv-cli

5.1.2. Project structure

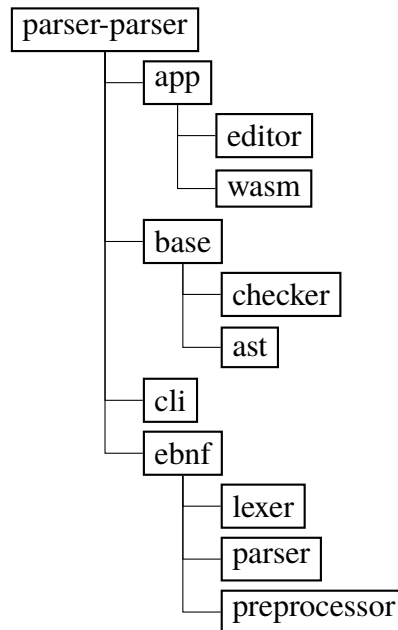


Figure 5.3. A simplified directory tree of the Parser-parser project.

The project structure of Parser-parser, visualized in Figure 5.3, consists of several components:

- `app` — An npm package containing the web application written primarily with Svelte. The component also defines a simple Rust crate, which acts as a link between the *base* and *ebnf* crates and compiles them to WebAssembly. The structure of the web application is described in detail in Section 5.4,
- `base` — Rust crate containing the definition of the AST along with the *checker* module, which will be described in detail in Section 5.2.5,
- `cli` — The auxillary command-line application written in Rust, explored further in Section 5.3,
- `ebnf` — The core business logic of the application in the form of a Rust crate. This crate is further divided into modules: the lexical, syntactic, and semantic analysers, all of which will be discussed in Sections 5.2.2, 5.2.3, and 5.2.4.

5.2. Business logic

5.2.1. Domain modelling

Domain Modeling is a way to describe and model entities and the relationships between them, which as a whole describe the problem domain space. Types can be utilized to describe the domain in a detailed way. Oftentimes, types can even be used to encode business rules so that the developer cannot write invalid code. Static type checking can be used as an immediate unit test — making sure that the code is valid at compile time. Types are the rules that describe what can happen in the domain and could be utilized to prevent anyone else from putting the system in a state invalid to the domain. Forcing illegal states to be

unrepresentable is the idea of statically proving that all values in the runtime correspond to valid objects in the business domain. That makes the code much easier to reason about and gives the developer confidence that the business rules are being respected.

If the logic of the system is represented by types, it is automatically self-documenting, and any changes to the business rules will immediately produce breaking changes, which is generally a welcome feature. This way the programmer can encode business requirements by writing code and, in the development process, produce documentation enforced by the compiler.

Using algebraic data types is a powerful technique for designing with types and making illegal states unrepresentable. Constructs such as sum types and product types provide us with an expressive method of modelling the business rules. This method also allows the developer to utilize property-based testing — letting the computer generate test cases.

For modelling the domain, the author will use the Haskell [19] programming language, with its expressive data types and highly reusable abstractions, as well as a concise syntax. However, modelling based on algebraic data types is also practical for other languages with complex enough type systems — Rust, used as the main language in Parser-parser, is one example of such a language. Please note that the author uses Haskell only in the thesis document — the project and its implementation uses the Rust programming language.

Token type definition

```

1 data Token
2   = Nonterminal String
3     | Terminal String
4     | Special String
5     | Integer Integer
6     | Concatenation
7     | Definition
8     | DefinitionSeparator
9     | EndGroup
10    | EndOption
11    | EndRepeat
12    | Exception
13    | Repetition
14    | StartGroup
15    | StartOption
16    | StartRepeat
17    | Terminator

```

Listing 5.1. Definition of the Token type in Haskell.

The tokenization process, described in Section 2.4, converts a stream of characters into a stream of tokens. A set of valid tokens can be represented as a sum type of all individual token types, shown in Listing 5.1. Several type constructors carry additional information about the token:

`Non-terminal` is specified by the textual form of the meta-identifier represented by the `String` type,

`Terminal` is specified by the contents of the terminal in the form of a `String`,

`Special` carries with it exact contents of the special sequence specified in the grammar, to be processed further,

`Integer` is specified by an actual numeric value encoded as Haskell's `Integer` type.

An equivalent of such a type in Rust could be defined as an *enumeration*, also referred to as an *enum*. Enums allow the developer to define a type by enumerating its possible variants. Enums variants in Rust, besides the kind of a variant, can also carry additional data, either in the form of a tuple, or a record. This way, the `Nonterminal`, `Terminal`, `Special`, and `Integer` variants can have an associated value built right into the type system.

Grammar type definition

The definition of the abstract Syntax Tree for EBNF in Section 2.5.2 gives an intuition for the definition of types related to grammars. The definition of the `Expression` type, shown is Listing 5.2, is recursive, meaning it has another instance of the enumeration as the associated value for one or more of the enumeration variants. In this case, any instance of `Alternative`, `Sequence`, `Optional`, `Repeated`, `Factor`, or `Exception` is considered a `branch node` in the AST, where, on the other hand, any instance of `Nonterminal`, `Terminal`, `Special`, or `Empty` is a `leaf node`.

```
1 type Grammar = HashMap String Expression
2 data Expression
3   = Alternative Expression Expression [Expression]
4   | Sequence Expression Expression [Expression]
5   | Optional Expression
6   | Repeated Expression
7   | Factor Integer Expression
8   | Exception Expression Expression
9   | Nonterminal String
10  | Terminal String
11  | Special String
12  | Empty
```

Listing 5.2. Definition of the types related to the AST in Haskell.

The `Expression Expression [Expression]` construct indicates a list of `Expressions` that contains at least two elements at all times. This way, the information about the minimum size of a list is built right into the type system. Such a construct could be abstracted away into its own definition, but in the process the `Expression` type would lose its pattern matching capabilities.

5.2.2. Lexical analyser

The lexical analyser, also known as *the lexer* or *the tokenizer*, performs the tokenization described in Section 2.4. A simplified version of the EBNF tokenizer could be modelled as a Deterministic Finite-state Automaton (DFA) (see Figure 5.4). This version, however, would not support nested comments defined in the specification — any nested structure cannot be tokenized with the use of regular languages, and those are, in fact, equivalent to DFAs.

The implementation does not follow the pure DFA approach. Instead, the lexer stores the current state of the tokenization process in various forms. For the purpose of tokenizing comments, the lexer remembers the *nest level*, which essentially counts the number of recursively-nested comments, and ends the comment only when the nest level reaches zero.

The lexer’s main control flow is a simple infinite loop followed by a pattern match, whose cases are the initial characters of each token type, and each case may be a loop to consume the rest of the token and return its type. Every token is preceded by the “whitespace-comment-whitespace” search, which skips any whitespace characters and (possibly nested) comments. Any whitespace inside integers or meta-identifiers is handled in the corresponding token loop.

which encodes every possible error the tokenizer can report. The `InvalidSymbol` case contains the additional information about the actual invalid symbol.

The whole tokenization process is preceded by the procedure of splitting the input string into individual Unicode graphemes according to the Unicode Standard Annex #29 [9] rules with the use of the `unicode-segmentation` crate. As each grapheme may consist of several characters, it is encoded a string, which means the lexer cannot use native functions for checking if a character is whitespace, alphabetic, alphanumeric, or a digit — these have to be defined separately with graphemes in mind.

The code related to the lexer is contained in the `lexer` module and is split into several files:

`mod.rs` is the main module file with the core business logic of the lexer and contains the `scan` and `lex` functions, as well as the utility functions related to graphemes,

`token.rs` contains the definition of the `Token` type,

`error.rs` contains the definition of the `lexer::Error` type,

`tests.rs` holds unit tests related to lexical analysis, which will be later discussed in Section 6.1.

5.2.3. Syntactic analyser

The syntactic analysis phase, also known as the parsing phase (described in Section 2.5) is conducted by the parser module. The goal of this thesis is not implementing a parser combinator library from scratch, so instead of writing every parser by hand, `Parser-parser` uses `Nom` (see Section 5.1.1) as the parser combinator library of choice, which ships with a large number of utility parsers and combinators already defined. These are in turn combined into more sophisticated parsers that are capable of parsing certain EBNF-like structures.

```
1 fn sequence(i: Tokens) -> IResult<Tokens, Spanned<Expression>,
  ↳ Spanned<Error>> {
2     map(
3         separated_list1(concatenation_symbol, term),
4         |nodes| match nodes.len() {
5             1 => nodes[0].clone(),
6             _ => Expression::Sequence {
7                 first: Box::new(nodes[0].clone()),
8                 second: Box::new(nodes[1].clone()),
9                 rest: nodes[2..].to_vec(),
10            }
11         .spanning(Span::combine(&nodes[0].span, &nodes[nodes.len() -
  ↳ 1].span)),
12     },
13     )(i)
14 }
```

Listing 5.3. The sequence parser.

For instance, the parser seen in Listing 5.3 parses the sequence of terms separated by commas with the use of the `separated_list_1` native to `Nom`. It also utilizes the `map`

combinator to transform the result into an appropriate type — it returns the `Sequence` of terms in case it parsed two or more terms, or just the single term in case it was the only one in the sequence. Note that this parser uses the `term` and `concatenation_symbol` parsers defined in a similar fashion.

All the more complicated parsers eventually use the so-called *literals* — the most simple parsers capable of parsing single tokens. Every token from the `Token` type (defined in Section 5.2.1) has an equivalent literal parser.

All parsers in `Nom` are basically functions, which most generally take a string as the input, and return an `IResult` — either the parsed *thing* along with the rest of the unconsumed input, or an error of some sort indicating that the parser failed. In the case of `Parser`, however, the input is not a string of characters, but a string of tokens. Fortunately, `Nom`, thanks to its high extensibility, can parse any type as long as that type implements certain traits. This resulted in the definition of the `Tokens` data structure, which encapsulates the `&[Spanned<Token>]` type — a slice of a sequence of (spanned) tokens. `Tokens` implements such traits as

- `InputLength`
- `InputIter`
- `InputTake`
- `UnspecializedInput`
- `Compare<Tokens>`
- `Slice<Range<usize>>`
- `FindSubstring<&[Spanned<Token>]>`

The `Tokens` type can then be used by `Nom` in a very efficient manner — these traits are one of the reasons why `Nom` is so performant.

The parsers generally return a single AST node as their output. These nodes are then combined into other nodes, which are then combined into `Productions`, and finally into a `Grammar`. In fact, the main parse function, the header of which is

```
fn parse(tokens: &[Spanned<Token>]) -> Result<Spanned<Grammar>,
↳ Spanned<Error>>
```

returns a `Grammar`, which contains the information about the whole syntax tree parsed from the provided sequence of tokens.

The code related to syntactic analysis can be found in the `parser` module. It is divided among several files:

`mod.rs` defines all major parsers along with the main `parse` function exported from the module,

`ast.rs` contains the definition of the AST,

`tokens.rs` holds the definition of the `Tokens` type and its trait implementations,

`error.rs` contains the definition of the `parser::Error` type,

`utils.rs` defines the utility functions and *literal* parsers,

`tests.rs` contains unit tests related to syntactic analysis, which are going to be discussed in Section 6.1.

5.2.4. Semantic analyser

Grammar preprocessing is a phase that tries to detect semantic errors in the AST. This phase does not *produce* anything from the AST, and serves only as a *guard* for any semantic inaccuracies, which include undefined production rules, or direct or indirect left recursion.

Before checking for left recursion, the algorithm checks for any non-terminals in the AST that do not have a definition. These non-terminals are invalid and any occurrence of such a non-terminal should be reported as an error. The basic principle of detecting undefined non-terminals, and so — production rules — is a recursive walk of the AST in search of any undefined non-terminals. A non-terminal is undefined if there is no `Production` in the `Grammar` with an appropriate identifier in the lhs (left-hand side). When such a non-terminal is found, the algorithm reports an `UndefinedRule` error at the position of the non-terminal.

In the formal language theory, a production rule is left-recursive if the leftmost symbol of it is itself (in the case of direct left recursion) or can be made itself by some sequence of substitutions (in the case of indirect left recursion). A PEG is called *well-formed* if it contains no left-recursive rules, i.e., rules that allow a non-terminal to expand to an expression in which the same non-terminal occurs as the leftmost symbol.

Consider the following example:

```
integer = ? [0-9.]+ ?;  
value   = integer | '(', expr, ')';  
product = expr, {'*' | '/'}, expr;  
sum      = expr, {'+' | '-'}, expr;  
expr     = product | sum | value;
```

In this grammar, matching an `expr` first requires testing if a `product` matches while, at the same time, matching a `product` requires testing if an `expr` matches. As the term appears in the leftmost position, these rules make up a circular definition and cannot be resolved.

Left recursion regularly poses problems for parsers as it leads them into the state of infinite recursion. As left-recursive rules can always be rewritten, grammars are often preprocessed to eliminate them. Parser-parser, however, does not attempt to do any rewrites. Instead, it only detects the presence of direct or indirect left recursion in the provided grammar. The process of detecting left recursion can be seen in Algorithm 5.1.

The code related to preprocessing, or the syntactic analysis, can be found in the `preprocessor` module of the crate. The module is divided into several separate files:

`mod.rs` contains the algorithms for undefined rule and left recursion detection, as well as the main `preprocess` function,

`error.rs` defines the `preprocessor::Error` type,

`tests.rs` contains unit tests related to preprocessing, which are further mentioned in Section 6.1.

5.2.5. Grammar processing

After creating the AST of a grammar, it can be analysed along with an input string to check if that input string belongs to the language generated by the grammar. For this purpose, the program recursively checks nodes of the AST to see if they match the currently scanned part of the input string. First, the user must provide the initial production rule, from which the process will begin. The recursive function must either return *success* or *failure* to indicate

Algorithm 5.1: Detecting left recursion in the set of production rules of a grammar.

```
input: Dictionary of production rules  $R$ 
1 function CheckExpr( $e, t$ ) begin
  inputs: The current expression  $e$ ; stack of identifiers  $t$ 
2  switch  $e$  do
3    case Alternative( $S$ ) do
4      foreach  $s \in S$  do
5        | CheckExpr( $s, t$ )
6    case Sequence(( $s_1, s_2, \dots, s_n$ )) do
7      /* Skip the last expression */
8      foreach  $s \in (s_1, s_2, \dots, s_{n-1})$  do
9        /* check if  $s$  can be an empty expression (e.g. Optional) */
10       | if  $s \neq \varepsilon$  then
11         | CheckExpr( $s, t$ )
12       | CheckExpr( $s_n, t$ ) // Check the last expression
13    case Optional( $s$ ) do
14      | CheckExpr( $s, t$ )
15    case Repeated( $s$ ) do
16      | CheckExpr( $s, t$ )
17    case Factor( $n, s$ ) do
18      if  $n > 0$  then
19        | CheckExpr( $s, t$ )
20    case Exception( $s, r$ ) do
21      | CheckExpr( $s, t$ ) CheckExpr( $r, t$ )
22    case Nonterminal( $i$ ) do
23      if  $i = t[0]$  then
24        | t.push( $i$ )
25        | exit() // Report the left recursion error with trace  $t$ 
26      if  $i \notin t$  then
27        | t.push( $i$ )
28        | CheckExpr( $R[i], t$ )
29        | t.pop()
28 foreach ( $i, e$ )  $\in R$  do
29   | CheckExpr( $e, R$ )
```

whether the parsing has succeeded. The program processes the input differently depending on the type of the AST node:

Alternative — The program processes each case of the *Alternative* sequentially and returns the first one to return *success*. If no case returned *success*, it returns *failure*,

Sequence — The program processes each expression of the *Sequence* sequentially and returns *success* if and only if every processed case returned *success*, otherwise it returns *failure*,

Optional — The program processes the expression inside of the *Optional* and returns *success* regardless of the result,

Repeated — The program repeatedly processes the expression inside of the *Repeated* and returns *success* regardless of any results,

Factor — The program processes the expression inside of the *Factor* *N* times, where *N* is the number of repetitions of the *Factor*. It returns *success* if the processing succeeds all *N* times, otherwise it returns *failure*,

Exception — The program processes the *subject* of the *Exception* and in case of *success*, it stores the processed input string to be then processed by the *restriction* of the *Exception*. If the restriction succeeds and the resulting processed input string is the same as the input string in case of the subject, it returns *failure*. In every other case it returns *success*,

Non-terminal — The program recursively processes the production rule specified by the identifier inside of the *Non-terminal*,

Terminal — The program checks if the processed input string starts with the input specified by the *Terminal* and returns *success* if it does,

Special — Currently, the program always returns *failure* for any special sequence,

Empty — The program always returns *success* without processing the input string.

The recursive function produces a parse tree along the way, which is represented by the `Node` type defined in Listing 5.4. Each `Node` represents a tree node, where `Nonterminal` variants are *branches* (containing multiple children), and `Terminal` variants are *leaves* of the tree. Additionally, each `Nonterminal` has an associated identifier, and each `Terminal` carries information about the input parsed by that terminal. The construction of `Nodes` follows a similar behavior to the checking:

```
pub enum Node {  
    Nonterminal(String, Vec<Node>),  
    Terminal(String),  
}
```

Listing 5.4. The definition of the `Node` type.

- `Alternatives` return the node of the first successful parse,
- In the case of `Sequences`, all nodes parsed from that sequence are appended to the `Vec` of children nodes,
- `Optionals` may not return a node at all,
- `Repeateds` append the repeated node to the `Vec` of children as many times as the parser is successful,
- `Factors` append the repeated node to the `Vec` of children as many times as the `Factor` defines,
- In the case of `Exceptions`, the node of the *subject* is returned as long as the exception is valid,
- `Nonterminals` return the `Nonterminal` node with an appropriate identifier and a recursive call,
- `Terminals` return the `Terminal` node along with the parsed input.

5.3. Command-line application

The command-line application serves as a most basic tool for interfacing with the grammar parser and checker. It is intended for people who want to, for example, automate the parsing and checking process, do it locally, or just plan to use it without the graphical user interface.

```
$ parser-parser --help
parser-parser 0.1.0

USAGE:
  parser-parser [OPTIONS] <GRAMMAR FILE>

FLAGS:
  -h, --help          Prints help information
  -V, --version        Prints version information

OPTIONS:
  -t, --test <TEST STRING FILE>  Test string
  ↪ file

ARGS:
  <GRAMMAR FILE>  Grammar file path
```

Listing 5.5. The output of the program ran with the `--help` flag.

the program respectively.

The tool reads the arguments to the program from the command line. The user must provide a path to the file containing a grammar in the EBNF format as an obligatory positional argument. The program provides an additional optional argument for passing a file containing the test input string denoted with the `-t`, or the `-test` flag. In case the test string is provided, the program parses the grammar, and then checks the provided input from a file. However, if the test string is not provided, the program, after parsing the grammar, enters the REPL mode. The program provides the `--help` and `--version` flags for displaying the help menu (Figure 5.5) and the current version of

For parsing command-line arguments, in its core `Parser-parser` uses the *Clap* crate for Rust — a Command Line Argument Parser. It is a simple-to-use, efficient, and full-featured library for parsing command-line arguments and subcommands when writing console/terminal applications. *Clap* is used to parse and validate the string of command-line arguments provided by a user at runtime. The developer provides the list of valid possibilities, and *Clap* handles the rest. This means the developer can focus on the application’s functionality, and less on the parsing and validating of arguments.

Clap provides many things with no configuration, including the traditional version and help switches (or flags) along with associated messages. If the user is using subcommands, *Clap* will also auto-generate a help subcommand and separate associated help messages. Once *Clap* parses the user provided string of arguments, it returns the matches along with any applicable values. If the user made an error or typo, *Clap* informs them with a friendly message and exits the program.

Besides *Clap*, `Parser-parser` also uses the *StructOpt* crate, which parses command line arguments by defining a struct. It combines *Clap* with a custom derive. By defining a regular struct and marking it with a specific derive, *StructOpt* can automatically generate a command-line argument parser based on the values in the struct and their types. *StructOpt* is easy to use and is a convenient method of parsing the arguments into a single structure, which then can be used in the actual program. The *StructOpt* struct for `Parser-parser` is defined in Listing 5.6, where it takes two file paths: one obligatory and positional, and one optional and marked with an appropriate flag.

```

1 #[derive(Debug, StructOpt)]
2 pub struct Config {
3     /// Grammar file path
4     #[structopt(name = "GRAMMAR FILE", parse(from_os_str))]
5     pub grammar_path: PathBuf,
6     /// Test string file path
7     #[structopt(short = "t", long = "test", name = "TEST STRING FILE",
8     ↪ parse(from_os_str))]
9     pub test_string_path: Option<PathBuf>,
10 }

```

Listing 5.6. The StructOpt struct defining the command-line arguments for the program.

In computing, a *REPL*, or the *read-eval-print loop*, is a simple interactive computer programming environment that takes single user inputs, executes them, and returns the result to the user. A program written in a REPL environment is executed piecewise. The term is usually used to refer to programming interfaces similar to the classic Lisp machine interactive environment. Common examples include command line shells and similar environments for programming languages, and the technique is very characteristic of scripting languages.

The name *read-eval-print loop* comes from the steps the program executes:

- The *read* step accepts an input from the user,
- The *eval* step takes the provided input and “evaluates” it, or in this case it checks the input against the provided grammar,
- The *print* step takes the result yielded by `eval` and prints it out to the user.

The environment then returns to the *read* state, creating a loop, which terminates when the program is closed. An example terminal output of the REPL mode can be seen in Figure 5.7.

5.4. Web-based application

5.4.1. Linking the business logic

To use the business logic (described in Section 5.2) in a web environment, the Rust code must be compiled to a WebAssembly module. This is done through Wasm-pack — a tool for building and working with WebAssembly generated from Rust code. The generated WebAssembly module can then be used like any other JavaScript module. Wasm-pack tool interoperates with `wasm-bindgen`, which has the ability to generate bindings, that allows JavaScript and Rust interact with each other through numbers, strings, or even JavaScript objects and arrays.

Wasm-pack automatically ensures that Rust 1.30 or newer and the `wasm32-unknown-unknown` target is installed via *rustup*, compiles the Rust sources into a WebAssembly `.wasm` binary

```

$ parser-parser ../grammar.ebnf
parser-parser 0.1.0

Successfully parsed the provided grammar

Entering REPL mode...
> 4+12*6
true
> x*3f
false
>

```

Listing 5.7. Sample of the output of the program ran in REPL mode.

via Cargo, and finally uses `wasm-bindgen` to generate the JavaScript API for using the Rust-generated WebAssembly. To do all of that, the developer needs to run the `wasm-pack build` command inside the project directory. When the build completes, its artifacts can be found in the `pkg` directory. The directory contains several files:

- The `.wasm` file is the WebAssembly binary that is generated by the Rust compiler from the Rust sources. It contains the compiled-to-wasm versions of all of the Rust functions and data.
- The `.js` file is generated by `wasm-bindgen` and contains JavaScript “glue” for importing DOM and JavaScript functions into Rust and exposing an API to the WebAssembly functions to JavaScript. The code will verify parameters passed across the boundary of WebAssembly and JavaScript and invoke appropriate functions.
- The `.d.ts` file contains TypeScript type declarations for the JavaScript “glue”. If the developer is using TypeScript, they can have their calls into WebAssembly functions type checked, and the IDE can provide autocompletions and suggestions. Parser-parser does not use TypeScript, so the file can be ignored.
- The `package.json` file contains metadata about the generated JavaScript and WebAssembly package. This is used by npm and JavaScript bundlers to determine dependencies across packages, package names, and their versions. It helps us integrate with JavaScript tooling and allows the developer to publish the package to npm.

```

1  #[wasm_bindgen]
2  pub struct EbnfParserParser {
3      grammar: base::Grammar,
4  }
5
6  #[wasm_bindgen]
7  impl EbnfParserParser {
8      #[wasm_bindgen(constructor)]
9      pub fn new(input: &str) ->
10         ↪ Result<EbnfParserParser, JsValue> {
11         // ...
12     }
13
14     #[wasm_bindgen(getter = productionRules)]
15     pub fn get_production_rules(&self) -> Array {
16         // ...
17     }
18
19     pub fn check(&self, input: &str, initial_rule:
20         ↪ &str) -> Option<Object> {
21         // ...
22     }
23 }

```

Listing 5.8. The definition of the EBNF parser struct that encapsulates the grammar.

Since the `Grammar` struct does not to be exported directly to JavaScript, `Wasm-bindgen` exposes the `EbnfParserParser` struct encapsulating the `Grammar` (Listing 5.8), which JavaScript treats as a regular class. `Wasm-bindgen` allows to mark the `new` function as a class constructor, so the construction of the `EbnfParserParser` object can be done via the `new` keyword. Besides the constructor, `Wasm-bindgen` marks the `get_production_rules` method as a regular getter, which will be exposed as a JavaScript property, instead of a method. Finally, after generation, the class will also contain the `check` method, which will check the provided `input` against the encapsulated grammar, starting with the `initial_rule`.

The return types of these functions get exposed to JavaScript in the following way:

- `Object` and `Array` types are translated directly into JavaScript objects and arrays respectively. These types come from the `js_sys` crate. The `Array` type may be produced by coercing a `Vec` of `JsValues`. The `Object` type, on the other hand, can be produced with `Object::new()` and have its properties set using the `Reflect` module and the `Reflect::set` function.

- The `Option<T>` type translates into a nullable value, so the `Some(T)` variant becomes `T`, and the `None` variant becomes `null`.
- A function returning `Result<T, E>`, after exposing the function to JavaScript, will return `T` in case of the `Ok(T)` variant of the result, and throw the value of `E` in case of the `Err(E)` variant.

5.4.2. Text editor

Parser-parser utilizes CodeMirror as its text and code editor. CodeMirror is a code editor component that can be embedded in web pages. The core library provides only the editor component, no accompanying buttons, auto-completion, or other IDE functionality. It does provide a rich API on top of which such functionality can be straightforwardly implemented. The distribution includes addons, which are reusable implementations of extra features. CodeMirror works with language-specific modes, which are JavaScript programs that help color (and optionally indent) text written in a given language. The distribution comes with a number of modes, including the EBNF mode, which is used in Parser-parser.

To integrate CodeMirror with Svelte, the author implemented a wrapper Svelte component on top of the native CodeMirror component. The implemented component utilizes the state-driven approach of Svelte, and exposes callbacks for associating functionality with events of the editor. CodeMirror works on top of a `<textarea>` HTML element of the web document:

```
<textarea bind:this="{textAreaRef}" readonly></textarea>
```

The `<textarea>` is bound to a `textAreaRef` variable with the use of the `bind:this` construct. Through JavaScript, the editor is constructed in the `onMount` event of the component by invoking the `fromTextArea` function with an appropriate binding. This will, among other things, ensure that the `textarea`'s value is updated with the editor's contents when the form is submitted.

```
editor = CodeMirror.fromTextArea(textAreaRef, opts);
```

where `opts` is a dictionary of configuration options provided by Svelte. Any option not supplied like this will be taken from `CodeMirror.defaults`, an object containing the default options. Options are not checked in any way, so setting undefined option values is bound to lead to odd errors. Some more notable options include:

mode — The mode to use. When not given, this will default to the first mode that was loaded. It may be a string, which either simply names the mode or is a MIME type associated with the mode. Alternatively, it may be an object containing configuration options for the mode, with a `name` property that names the mode (for example `name: "javascript", json: true`). The value `null` indicates no highlighting should be applied.

theme — The theme to style the editor with. The developer must make sure the CSS file defining the corresponding `.cm-s-[name]` styles is loaded from the theme directory in the distribution. The default is “default”, for which colors are included in `codemirror.css`.

tabSize — The width of a tab character. Defaults to 4.

extraKeys — Can be used to specify extra key bindings for the editor, alongside the ones defined by **keyMap**. Should be either null, or a dictionary of key bindings.

lint — Available after loading the Lint addon, which defines an interface component for showing linting warnings, with pluggable warning sources. The **lint** option can be set to a warning source, which is a validator function that returns a list of errors from the provided string.

5.4.3. Parse tree visualizations

```

  ✓ expression
    ✓ term
      > factor
        "+"
    ✓ term
      ✓ factor
        ✓ constant
          ✓ digit
            "3"
```

Section 5.2.5 described the **Node** structure, which represents the parse tree returned from the checker. This structure is utilized in the visualization of the parse tree on the front-end. Each **Node** converted to JavaScript objects is composed of a name parameter, and, in the case of branch nodes, a **children** parameter with an array of all children nodes. This structure is compatible with *Svelte-tree* package — a tree-like view component for Svelte. The component is used on the front-end to display a collapsible tree based on the structure described above. An example of the visualized parse tree is presented in Figure 5.5.

Figure 5.5. A screenshot of the visualized parse tree.

Svelte-tree allows the user to associate a node with an identifier via the **let:node** construct, and this way access the names of nodes being rendered at appropriate DOM elements.

The *Svelte-tree* is abstracted away into its own Svelte component responsible for visualizing the parse tree, which can be seen in Listing 5.9. In the **<script>** section it imports the **Tree** component from the *Svelte-tree* package, and exposes the **tree** parameter as its input. Each node in the tree is marked with a **node** class, which, in the **<style>** section, is getting an appropriate styling with CSS.

```

1 <script>
2   import Tree from "svelte-tree";
3   export let tree;
4 </script>
5
6 <main>
7   <Tree tree={[tree]} let:node>
8     <div class="node">{node.name}</div>
9   </Tree>
10 </main>
11
12 <style>
13   .node {
14     float: left;
15     font-family: "JetBrains Mono", Consolas, monospace;
16     color: #928374;
17     padding-left: 24px;
18   }
19 </style>
```

Listing 5.9. Contents of **ParseTree.svelte**, the parse tree visualization component.

6. Project quality study

6.1. Business logic testing

6.1.1. Unit testing

Unit testing is a type of software testing where individual units or components of the software are being tested. The purpose of unit testing is to validate that each unit of the software code performs in an expected way. Unit tests assist in fixing bugs early in the development cycle, help to understand the code base and make changes to the code quicker, and can also serve as project documentation. Unit testing allows the programmer to refactor code, and ensure the refactored module still works correctly. The process is to write test cases for all functionalities so that whenever a change causes an error, it can be quickly identified and be dealt with. Unit testing, however, can't be expected to catch every error in a program. It is not possible to evaluate all execution paths even in the most trivial programs. Unit testing by its definition focuses on a unit of code — it can't catch integration errors or broad system-level errors.

A programmer usually writes a section of code in the application just to test a function. Isolating the tested code helps in disconnecting unnecessary dependencies from the code being tested and other units in the project — testing should be focused on only one piece of code at a time, and individual unit test cases should be independent. In case of any changes in requirements, unit test cases should not be affected. If units are made less interdependent to make unit testing possible, the unintended impact of changes to any code is less noticeable.

Unit tests in Parser-parser focus on testing the individual components of the EBNF parser, that is: the lexer, the parser, and the preprocessor. If the developer is sure that all of these components work correctly in isolation, the same thing can be said about the whole system, because there are no hidden dependencies between these components.

Tests related to the lexer test the functionality of the lexer, that is, test the ability of the lexer to turn the textual representation of EBNF into individual tokens. An typical unit test related to the lexer can be seen in Listing 6.1, where it makes an assertion that the result of the `lex` function on input “`abc\n = 'def'` ” is a success (the `Ok` case), and it's a vector of three certain tokens.

```
1 #[test]
2 fn test_multiline() {
3     assert_eq!(
4         lex(" abc \n = 'def' "),
5         Ok(vec![
6             Token::Nonterminal("abc".to_owned()).spanning(Span::from(((1, 0), (4, 0)))),
7             Token::Definition.spanning(Span::from(((1, 1), (2, 1)))),
8             Token::Terminal("def".to_owned()).spanning(Span::from(((3, 1), (8, 1))))
9         ])
10    );
11 }
```

Listing 6.1. A unit test related to the lexer testing the proper tokenization of the input string.

All test cases related to the lexer are listed in Table 6.1 in a simplified format, where the result can either be a success case or a failure case, but does not provide information about the lexed tokens for simplicity. The actual test cases for the lexer are located in the file `ebnf/src/lexer/tests.rs`.

Table 6.1. Tests cases related to the lexer along with the expected results.

Test case	Result	Test case	Result	Test case	Result
,	success	""	failure	abc12_3_	success
//!	success	' '	failure	_x_	success
abc_=_b;	success	_?_test_?	success	_+_	failure
(//_/)_	success	?a\nbc?_	success	_,_\n,	success
(/)	failure	_?bbb_	failure	_(*_test*_)_	success
/	success	??	success	_(*_test*_	failure
(:*)__{_}	success	_123_	success	_(*_(failure
(:)	failure	_1_2_3_	success	_,_(*,_)_,_	success
_ "ab_c_"	success	_01234_5"	success	_,_(*,_(*)__,_*)_,_	success
_ _' "aba ' _	success	_0_	success	_(*_(*_*_)_	failure
_ _'a_"	failure	_abc_	success	_(*_)_	failure
"bbb'_ _	failure	a_ _bc_	success	_abc_\n_=_ 'def' _	success

Tests related to the parser test the functionality of the parser, that is, test the ability to transform a list of tokens into an AST. To make the unit tests of the parser independent from the implementation of the lexer, the input to the `parse` function is a vector of tokens, rather than a result of the `lex` function. An example test case related to the parser is seen in Listing 6.2, where the `ok_case` macro represents a test case that, after testing the `factor` parser, should result in a success, parse 3 tokens, and return a `Factor` AST node. In the same file, that is `ebnf/src/parser/tests.rs` one can also find the `error_case` macro, which represents a test case that should fail with a specific error.

```

1 ok_case!(
2     factor,
3     &vec![
4         Token::Integer(2).spanning(Span::from(((0, 0), (1, 0)))),
5         Token::Repetition.spanning(Span::from(((2, 0), (3, 0)))),
6         Token::Terminal("terminal".to_owned()).spanning(Span::from(((4, 0), (14, 0))))
7     ],
8     3,
9     Expression::Factor {
10         count: 2.spanning(Span::from(((0, 0), (1, 0)))),
11         primary: Box::new(
12             Expression::Terminal("terminal".to_owned()).spanning(Span::from(((4, 0), (14, 0))))
13         )
14     }
15     .spanning(Span::from(((0, 0), (14, 0))))
16 );

```

Listing 6.2. A unit test related to the parser, where the ability to turn a string of tokens into an AST is tested.

Tests related to the preprocessor test the functionality of the preprocessor, that is, test the ability to detect undefined rules and left recursion. Test case in Listing 6.3 tests the grammar

```

a = b;
b = a;

```

and the detection of indirect left recursion $b \rightarrow a \rightarrow b$.

```
1  #[test]
2  fn test_indirect_left_recursion() {
3      assert_eq!(
4          preprocess(
5              Grammar {
6                  productions: vec![
7                      Production {
8                          lhs: "a".to_owned().spanning(Span::from(((0, 0), (1, 0)))),
9                          rhs: Expression::Nonterminal("b".to_owned())
10                             .spanning(Span::from(((4, 0), (5, 0))))
11                      },
12                      .spanning(Span::from(((0, 0), (6, 0)))),
13                      Production {
14                          lhs: "b".to_owned().spanning(Span::from(((0, 1), (1, 1)))),
15                          rhs: Expression::Nonterminal("a".to_owned())
16                             .spanning(Span::from(((4, 1), (5, 1))))
17                      },
18                      .spanning(Span::from(((0, 1), (6, 1))))
19                  ]
20              },
21          ),
22      ),
23      Err(
24          Error::LeftRecursion(vec!["b".to_owned(), "a".to_owned(), "b".to_owned()])
25              .spanning(Span::from(((4, 0), (5, 0))))
26      )
27  );
28 }
```

Listing 6.3. A preprocessor unit test testing the left recursion detection in an AST.

All tests related to business logic are written in Rust, as the business logic is also written in Rust. All unit tests in the project can be ran with the use of Cargo with the `cargo test` command in the terminal.

6.1.2. Property-based testing

Test engineers write mostly example-based tests where only one input scenario is tested. Property-based testing is a useful addition to a test suite because it runs one statement hundreds of times with different inputs. Property-based testing frameworks use almost every conceivable input that could break the code, such as empty lists, negative numbers, and really long lists or strings. Property based testing has become quite famous in the functional world. Mainly introduced by QuickCheck framework in Haskell, it suggests another way to test software.

Property-based tests are designed to test the aspects of a property that should always be true. They allow for a range of inputs to be programmed and tested within a single test, rather than having to write a different test for every value that the programmer wants to test. These tests are useful when a range of inputs needs to be tested on a given aspect of a software property, or if the developer is concerned about finding missed edge cases.

```

1 use quickcheck_macros::quickcheck;
2
3 #[quickcheck]
4 fn test_arbitrary_input(input: String) {
5     let _ = lex(&input);
6 }

```

Listing 6.4. A QuickCheck test testing arbitrary inputs on the lexer.

Parser-parser uses the `quickcheck` crate for Rust, which is a property-based testing framework inspired by Haskell’s QuickCheck and has been described in Section 5.1.1. `quickcheck` tests the “correctness” of `lex` and `parse`, i.e. it checks if these functions do not emit a *panic* for an arbitrary input, whether it’s a random string of characters for the case of `lex`, or a random string of tokens for `parse`. Example `quickcheck` tests can be seen in Listings 6.4 and 6.5, where the `quickcheck` macro is used to streamline the usage of the library. `quickcheck` catches any panics that may occur in `lex` and `parse` functions and reports them to the user along with the input provided to these functions. In the case of `parse`, to generate an arbitrary string of Tokens, it was necessary to implement the `Arbitrary` trait provided by `quickcheck` for `Token`. `Spanned<T>`, `Span`, and `Location` types (most of this has been omitted in Listing 6.5).

```

1 // ...
2 impl<T: Arbitrary> Arbitrary for Spanned<T> {
3     fn arbitrary<G: Gen>(g: &mut G) -> Spanned<T> {
4         Spanned {
5             node: T::arbitrary(g),
6             span: Span::arbitrary(g),
7         }
8     }
9 }
10
11 #[quickcheck]
12 fn test_arbitrary_input(tokens: Vec<Spanned<Token>>) {
13     use super::parse;
14     let _ = parse(tokens.as_slice());
15 }

```

Listing 6.5. A QuickCheck test related to the parser, testing strings of arbitrary tokens.

6.2. Integration testing

Integration testing is defined as a type of testing where software modules are integrated logically and tested as a group. A typical software project consists of multiple software modules, coded by different programmers. The purpose of this level of testing is to expose defects in the interaction between these software modules when they are integrated.

Parser-parser utilizes integration testing thanks to *wasm-pack*. This tool allows the developer you build rust-generated WebAssembly packages, as well as test them in a headless web browser via the `wasm-pack test` command. A tool used to define the integration tests is the `wasm-bindgen-test` crate — an experimental test harness for Rust programs compiled

to wasm using wasm-bindgen and the wasm32-unknown-unknown target. The wasm-pack test command wraps the wasm-bindgen-test-runner CLI allowing the developer to run wasm tests in different browsers without needing to install the different webdrivers. An example integration test ran with wasm-bindgen-test can be seen in Listing 6.6.

```
1 use ebnf::parse;
2 use wasm_bindgen_test::*;
3
4 #[wasm_bindgen_test]
5 fn test_ebnf() {
6     assert!(parse(" abc = 'def'; ").is_ok());
7     assert!(parse(" (* test *) ").is_err());
8     assert!(parse(" (* test *").is_err());
9     assert!(parse("a = b;").is_err());
10    assert!(parse("a = 'a';").is_err());
11    assert!(parse("a = ;").is_ok());
12    assert!(parse("a = a;").is_err());
13 }
```

Listing 6.6. An integration test ran in a headless browser, which tests various grammars in a textual form.

6.3. Benchmarking

Benchmarks check the performance of the code. Parser-parser’s benchmark tool of choice is Criterion.rs — a statistics-driven micro-benchmarking tool. It is a Rust port of Haskell’s Criterion library. Criterion.rs benchmarks collect and store statistical information from run to run and can automatically detect performance regressions as well as measuring optimizations. Listing 6.7 shows a Criterion.rs benchmark that parses a sample input grammar multiple times.

```
1 use criterion::{black_box, criterion_group, criterion_main, Criterion};
2 use ebnf::parse;
3
4 const GRAMMAR: &str = "
5 expression = term, { ('+' | '-'), term };
6 term       = factor, { ('*' | '/'), factor };
7 factor     = constant | variable | '(', expression, ')';
8 variable   = 'x' | 'y' | 'z';
9 constant   = digit, { digit };
10 digit      = '0' | '1' | '2' | '3' | '4' | '5' | '6' | '7' | '8' | '9';
11 ";
12
13 fn criterion_benchmark(c: &mut Criterion) {
14     c.bench_function("parse", |b| b.iter(|| parse(black_box(GRAMMAR))));
15 }
16
17 criterion_group!(benches, criterion_benchmark);
18 criterion_main!(benches);
```

Listing 6.7. A benchmark testing the speed of parsing a sample grammar.

Criterion.rs can generate a number of useful charts and graphs which the developer can check to get a better understanding of the behavior of the benchmark. These charts will be generated with *gnuplot* by default, and the examples below were generated using the *gnuplot* backend.

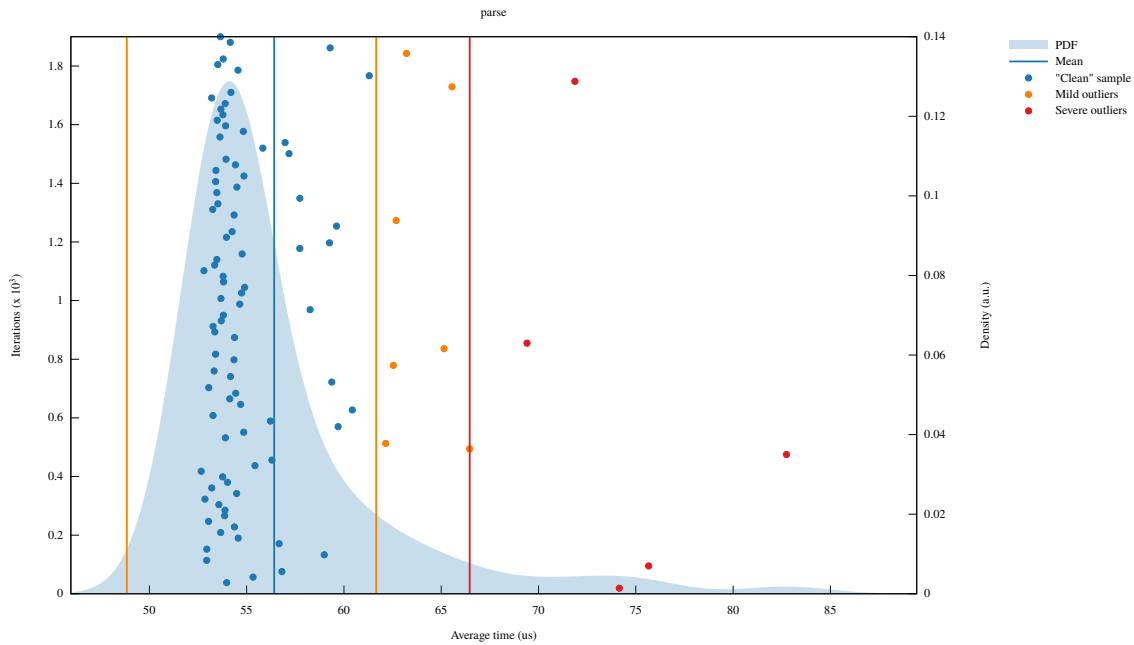


Figure 6.1. Probability Distribution Function chart generated by Criterion.rs.

The *PDF* chart in Figure 6.1 shows the *probability distribution function* for the samples. It also shows the ranges used to classify samples as outliers. In this example we can see that the performance trend does not change noticeably across the whole benchmark and stays around 55 μ s

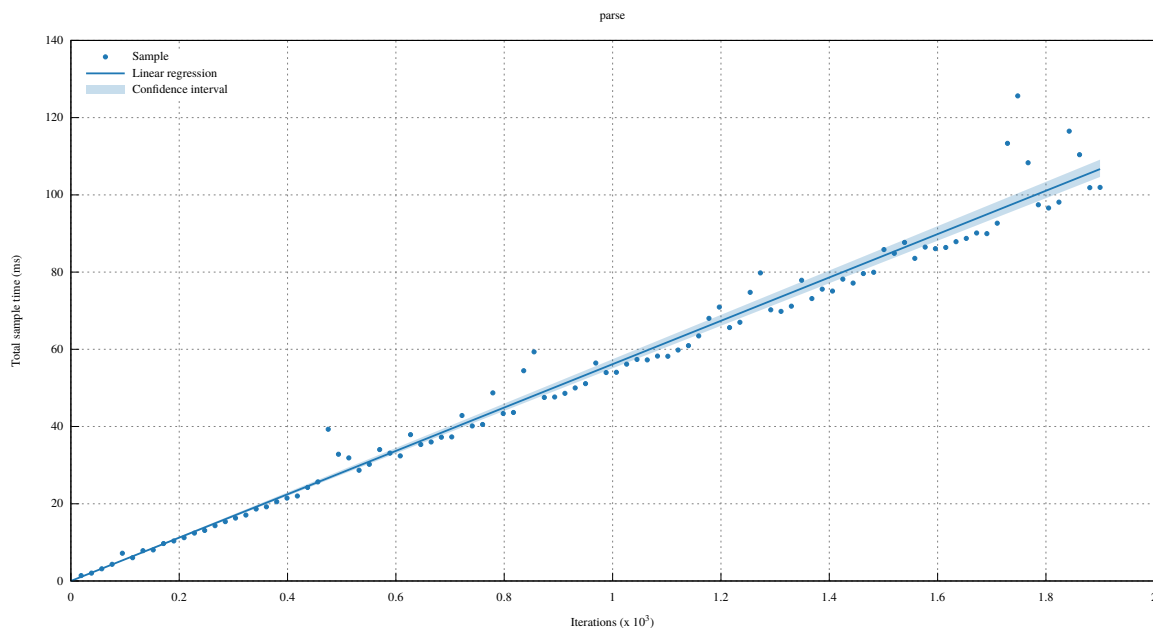


Figure 6.2. Regression plot generated by Criterion.rs.

The *regression plot*, shown in Figure 6.2, shows each data point plotted on an X-Y plane showing the number of iterations versus the time taken. It also shows the line representing Criterion.rs’ best guess at the time per iteration. A good benchmark will show the data points all closely following the line. If the data points are scattered widely, this indicates that there is a lot of noise in the data and that the benchmark may not be reliable. If the data points follow a consistent trend but don’t match the line (eg. if they follow a curved pattern or show several discrete line segments) this indicates that the benchmark is doing different amounts of work depending on the number of iterations, which prevents Criterion.rs from generating accurate statistics and means that the benchmark may need to be reworked.

The graphics that Criterion.rs generates are perfect for contributors of the project as there is no dearth of info. Criterion generates graphics that break down mean, median, standard deviation, MAD, etc., which are invaluable when trying to pinpoint areas of improvement.

6.4. Auditing

A software audit is an internal or external review of a software program to check its quality, progress or adherence to plans, standards and regulations. Parser-parser utilizes Google Lighthouse — an open-source, automated tool for improving the quality of web pages. Developers can run it against any web page, public or requiring authentication. It has audits for performance, accessibility, progressive web apps, SEO and more. When the tool finishes analyzing a web page, it returns a report with the calculated scores for each metric, a list of problems with the page, and general, or sometimes specific, recommendations regarding solving those problems. Figure 6.3 shows a short summary of Lighthouse ran against the deployed web application. For most metrics, Lighthouse calculates a score by comparing the page with the FCP data present in HTTP Archive. The tool uses a color-coding system to display how well a page performs according to a particular metric.

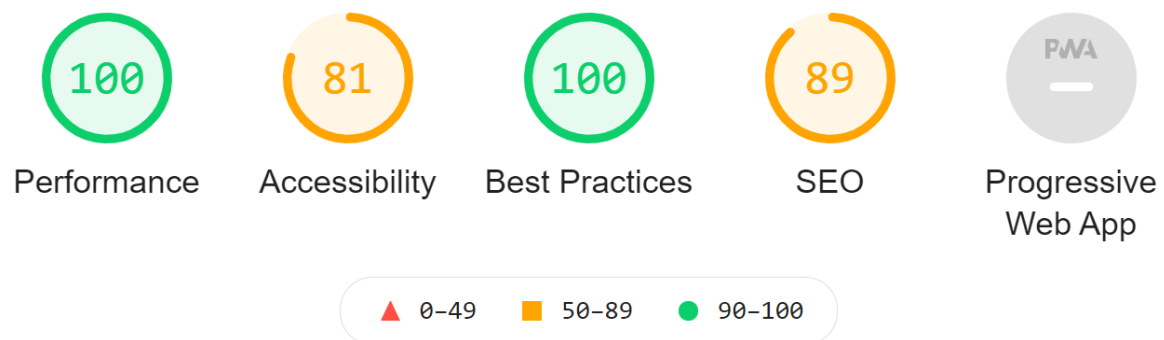


Figure 6.3. Google Lighthouse’s scores of the deployed web application.

Parser-parser scores 100 points in the performance category. This category has metrics that together reflect how fast the page is perceptually:

First Contentful Paint — Shows how long it takes for a browser to render DOM content,

Speed Index — Shows how quickly the contents of a page load visually. To do this, Lighthouse records a video of your page loading and then computes a visual progression between frames,

Largest Contentful Paint — Reports the render time of the largest image or text block within the viewport (i.e., the visible part of the page),

Time To Interactive — Measures how long it takes a page to become *fully* interactive, i.e. useful content (measured by FCP) is displayed, JavaScript event handlers are bound to visible elements' events, and the page responds to user interaction within 50 milliseconds,

Total Blocking Time — The sum of all time records between FCP and TTI when a page is blocked from user interaction for more than 50 milliseconds,

Cumulative Layout Shift — Metric for showing how aggressively elements shift each other during the loading.

All values for performance metrics generated by Lighthouse can be seen in Table 6.2. Besides the metrics, the website passed additional 29 audits, such as *Minify CSS*, *Minify JavaScript*, *Enable text compression*, *Initial server response time was short*, *Remove duplicate modules in JavaScript bundles*, *Avoids an excessive DOM size* and more.

Table 6.2. Performance metrics generated by Lighthouse.

Metric	Value
First Contentful Paint	0.2 s
Speed Index	0.4 s
Largest Contentful Paint	0.4 s
Time to Interactive	0.4 s
Total Blocking Time	0 ms
Cumulative Layout Shift	0

In terms of best practices, the web application scores 100 points, thanks to HTTPS usage provided by GitHub Pages, safe links to cross-origin destinations, avoiding requesting geolocation and notification permissions on page load, and avoiding front-end JavaScript libraries with known security vulnerabilities.

Besides performance audits, Lighthouse also provides audits regarding accessibility, SEO (Search Engine Optimization), and PWA (Progressive Web App). These audits were not the main focus of the project, however they provide useful information about possible future improvements. Lighthouse informs that background and foreground colors do not have a sufficient contrast ratio and that low-contrast text is difficult or impossible for many users to read, which is a valid concern that should be taken into consideration. The accessibility audit also suggests providing labels for associated form elements, which most likely relates to `textarea` grammar and input string fields. In terms of Progressive Web App audits, Lighthouse suggests using a service worker, which enables the web app to be reliable in unpredictable network conditions and use many Progressive Web App features, such as *offline*, *add to homescreen*, and *push notifications*.

6.5. Complexity analysis

Software complexity is a way to describe a specific set of characteristics of the code. These characteristics all focus on how the code interacts with other pieces of code.

A popular metric for measuring code complexity is *cyclomatic complexity*. Cyclomatic complexity uses a mathematical model to assess methods, producing accurate measurements of the effort required to test them, but inaccurate measurements of the effort required to understand them.

Cognitive complexity, on the other hand, breaks from the practice of using mathematical models to assess software maintainability. It starts from the precedents set by cyclomatic complexity, but uses human judgement to assess how structures should be counted, and to decide what should be added to the model as a whole. As a result, it yields method complexity scores which strike programmers as fairer relative assessments of maintainability than have been available with previous models.

Calculating cognitive complexity in Rust can be achieved with the *rust-clippy* tool. Clippy is a collection of lints to catch common mistakes and improve Rust code. There are over 400 lints included in the tool, and are divided into categories, each with a default lint level. One way to use Clippy is by installing Clippy through *rustup* as a Cargo subcommand. Clippy can then be invoked with the `cargo clippy` command.

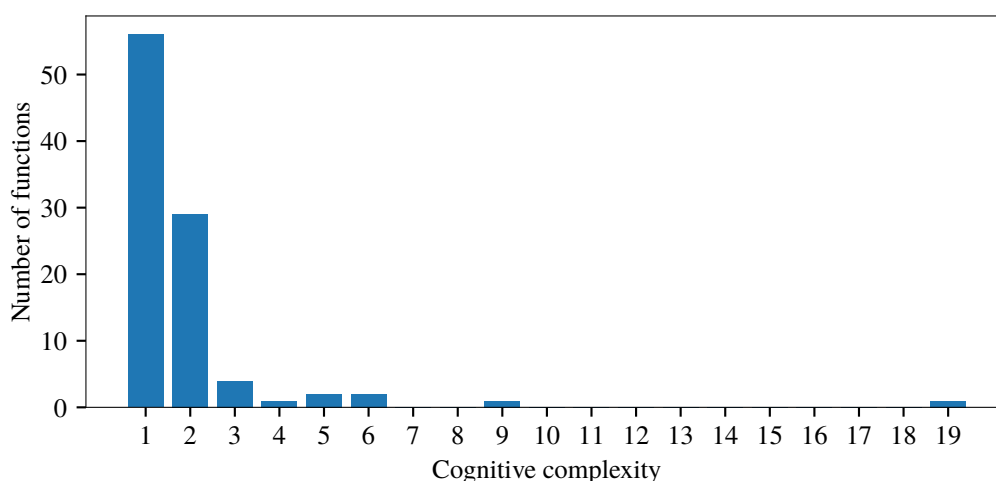


Figure 6.4. Number of functions with various degrees of cognitive complexity.

Figure 6.4 shows the bar plot of a number of functions of a particular cognitive complexity value. Most of functions in the project have a low cognitive complexity of around 1 or 2. Acceptable values for cognitive complexity of functions generally are in the range between 0 and 10. Based on this, there is one function with value above that range — the `lex` function. However, the complexity of the `lex` function is deliberate — the lexer was written with complex behavior in a single procedure in mind, for performance reasons. The written code is simple and low-risk when it comes to maintenance.

7. Deployment

7.1. Application building

Deployment involves packaging up the web application and putting it in a production environment that can run the app. To package, or *bundle* the web application, Parser-parser uses Rollup — a simple JavaScript module bundler, described in Section 5.1.1, which can also serve as a *build tool* for JavaScript. A bundler is a tool that recursively follows all imports from the entry point of the application and bundles them up into a single JavaScript file. What makes Rollup great, is its ability to keep files small — compared to the other tools for creating JavaScript bundles, Rollup will almost always generate a smaller, faster bundle. This happens, because Rollup is based on ES2015 modules, which are more efficient than CommonJS modules, which are what other bundlers use. It's also much easier for Rollup to remove unused code from modules using a technique called *tree-shaking*, thanks to which only the code the application actually needs is included in the final bundle. Tree-shaking becomes important when the application includes third-party tools or frameworks that have dozens of functions and methods available.

Rollup uses the `rollup.config.js` configuration file, which includes such configuration options as the entry point, the destination folder, the output format, but most importantly it allows to include plugins, particularly for preprocessing Svelte files, including CSS files, aliasing, static file copying, and wasm-pack support for building Rust project to WebAssembly.

To bundle the web application, the only command the developer has to invoke is `rollup -c`, which is aliased in the `package.json` file to the `build` script for accessibility. The command will generate a new folder called `dist` in the project, that contains all generated files ready for deployment.

7.2. Production environment

GitHub Pages are public web pages for users, organizations, and repositories, that are freely hosted on GitHub's `github.io` domain or on a custom domain name of choice. GitHub users can design and host both personal websites and websites related to specific GitHub projects. GitHub Pages provides each user with one user page (`username.github.io`) and an unlimited number of project pages (`username.github.io/projectname`). Each organization can also have one organization site and an unlimited number of project pages. Parser-parser's code is hosted on the `karolbelina/parser-parser` repository, which means that the GitHub Pages page, hosted in the same repository, is available on

`https://karolbelina.github.io/parser-parser/`

The developer has to decide how to organize the code so that files with ordinary application code are separated from files necessary for deployment on GitHub Pages. User and organiza-

tion pages must be loaded from the `main/master` branch, and project pages will work from any branch, but by convention it's usually called `gh-pages`.

To deploy the bundled web application to GitHub pages, Parser-parser uses the `gh-pages` npm package, which with a single command allows to push desired files to any git branch of the current repository. By default, the `gh-pages` package we will push the file tree to the `gh-pages` branch and make one if necessary. This makes the main branch and all other branches stay untouched without the need to explicitly *checkout* the `gh-pages` branch.

GitHub Pages and its `github.io` domains provide HTTPS, that adds a layer of encryption that prevents others from tampering with traffic to the site. All GitHub Pages sites, including sites that are correctly configured with a custom domain, support HTTPS and HTTPS enforcement. HTTPS enforcement is required for GitHub Pages sites using a `github.io` domain that were made after June 15, 2016. The user can customize the domain name of a GitHub Pages site, however this was not necessary in the final product.

7.3. Continuous integration and continuous deployment

Continuous integration allows different developers to upload and merge code changes in the same repository branch on a frequent basis. Once the code has been uploaded, it's validated automatically by means of unit and integration tests (both the uploaded code and the rest of the components of the application). In case of an error, this can be corrected more simply.

In continuous delivery, the new code that has been introduced and that has passed the CI process is automatically published in a production environment. What's intended is that the repository code is always in a state that allows its implementation in a production environment. In continuous deployment, all changes in the code that have gone through the previous two stages are automatically implemented in production.

Parser-parser uses GitHub Actions — a new tool to perform the CI/CD process, introduced in October 2018, launched in beta in August 2019, and finally distributed in November 2019. GitHub Actions is also paid, although it has a free version with some limitations.

To set up a GitHub Action workflow, the developer generates an appropriate file — in the case of CI/CD (Continuous integration and continuous deployment), which is the only workflow in the project, the file is named `ci.yml` and sits in the `.github` directory at the root of the repository. The format of the workflow file is YAML — a human-readable data-serialization language that is commonly used for configuration files. A part of the file can be seen in Listing 7.1. With the `on` keyword it lists the events that invoke the workflow — in this case, any push or pull request to the *master* branch.

Each workflow is made up of one or more jobs. The `runs-on` parameter contains the type of virtual machine in which the code will be executed. In this case, the workflow uses `ubuntu-latest` — the latest supported version of the Ubuntu operating system. The `steps` parameter is a sequence of tasks — in the example of CI/CD, the first step is to *checkout* the repository so the workflow can access it. Then, the workflow sets up the toolchain: it installs Node.js, Rust, and `wasm-pack`, to finally install the npm dependencies via `npm ci` (which is usually faster than `npm install`), build the project with `npm run build`, and deploy it to GitHub Pages using the `npm run deploy` command. The deployment is preceded with an appropriate Git setup, where the workflow configures Git username and e-mail to “fake” credentials of a GitHub Actions bot, as well as a GitHub token to authorize the push to the `gh-pages` branch. The token is stored in the repository via *the secrets*, which allow the developers to store sensitive information in the repository or organization, and are available

to use in GitHub Actions workflows. GitHub uses a *libsodium sealed box* to help ensure that secrets are encrypted before they reach GitHub, and remain encrypted until the developer uses them in a workflow. Only people with access to the repository can view the secrets.

The workflow uses third party GitHub Actions to set up the toolchain. The `setup-node@v1` action installs an appropriate version of Node.js along with npm. The `actions-rs/toolchain@v1` sets up the nightly version of the Rust compiler and the Cargo dependency manager. Finally, `jetli/wasm-pack-action@v0.3.0` installs the latest version of `wasm-pack` for the Web-Assembly compilation.

```
1 name: CI/CD
2
3 on:
4   push:
5     branches: [ master ]
6   pull_request:
7     branches: [ master ]
8
9 jobs:
10  deploy:
11    runs-on: ubuntu-latest
12
13    steps:
14      - name: Checkout sources
15        uses: actions/checkout@v2
16      - name: Set up Node.js
17        uses: actions/setup-node@v1
18        with:
19          node-version: '12'
20      - name: Install nightly Rust toolchain
21        uses: actions-rs/toolchain@v1
22        with:
23          profile: minimal
24          toolchain: nightly
25          override: true
26      - name: Install wasm-pack
27        uses: jetli/wasm-pack-action@v0.3.0
28        with:
29          version: 'latest'
30      - run: cd web
31      - run: npm ci
32      - run: npm run build
33      - name: Deploy to gh-pages
34        run: |
35          git config --global user.name $user_name
36          git config --global user.email $user_email
37          git remote set-url origin https://${github_token}@github.com/${repository}
38          npm run deploy
39      env:
40        user_name: 'github-actions[bot]'
41        user_email: 'github-actions[bot]@users.noreply.github.com'
42        github_token: ${ secrets.ACTIONS_DEPLOY_ACCESS_TOKEN }
43        repository: ${ github.repository }
```

Listing 7.1. The deployment part of the GitHub Action workflow file related to CI/CD.

8. Software artifacts

Artifacts described in this chapter are by-products of the design, implementation, and testing phases described in Chapters 4, 5, and 6, which are the core of the software development part of the project.

The software artifact provides a template to build on to allow continuity instead of starting from scratch or terminating the software — every program requires a constant update for continuity, lest they become redundant. The need to continually do software updates can become a significant source of work overload if the developer has to begin the developmental process all over again. However, using artifacts allow developers to have a stored template to fall back on when they need to update or upgrade their programs.

Design document

This thesis serves as the design document of the project. It describes the use cases — descriptions of how users are meant to perform tasks on the website. These relate directly to the function of the site, making them important artifacts. UML diagrams are a way of visualizing and plotting out the way a piece of software works. It works to map out links, processes, etc. Like use cases, UML doesn't directly help software run, but it's a key step in designing and programming a piece of software. Sequence diagrams are a way to map out the structure of the application — they are used to map out links and processes that happen between clicks in a more visual way. Diagrams are a great way to visualize the internal processes of a given program. These will be designed throughout the coding process, particularly in the preliminary stages.

Source code

The source code of the application was produced throughout the whole development period. It encapsulates the functionalities of the application defined in the design document, as well as test suites — coded test to run against the program in order to make sure a certain process is working. The project is open-source and is available on a public GitHub repository at

<https://github.com/karolbelina/parser-parser>

The source code consists of three Rust crates and one npm package, as seen in Section 5.1.2. Most of these work together to compose the single web application, with the npm package providing the tool chain necessary to build the project. The Rollup configuration provides the build tool necessary to bundle the source code in the form of Rust projects and Svelte files into a compiled web application.

The source code provides the testing infrastructure for the `ebnf` crate, with unit tests defined for major components of the system, along with integration tests which focus on the whole functionality compiled to a WebAssembly module. With a code artifact, a software

programmer can test the program in detail and perfect things before launching the software. It defines benchmarks for the parsing and checking functionality of the `ebnf` crate with the use of `Criterion.rs`.

The definition of the GitHub Actions workflow in the `ci.yml` file along with the `Dockerfile` show the build process for the web application that may be utilized by deploying the final product in production, as well as to ascertain the compatibility of each developed program with a specific machine.

Compiled applications

The web application is the source code compiled into a usable application. This is the final artifact, and one of the only ones a typical user will care about. The compiled application will let the user access it on their machine, and use it as its meant to be used. The command-line application is an auxillary compiled artifact of the project, that may also be used by the user by downloading the binary of the program and using it on their own machine.

Documentation

The standard Rust distribution ships with a tool called `rustdoc`. Its job is to generate documentation for Rust projects. On a fundamental level, `Rustdoc` takes as an argument either a crate root or a Markdown file, and produces HTML, CSS, and JavaScript. The documentation can be built for all crates that are specified in the `Cargo.toml` by using the `cargo doc` command. This compiles into documentation the *doc comments* in functions or modules, which are very useful for big projects that require documentation. They are denoted by a `“///”`, and support Markdown. The compiled documentation is put into the generated HTML documentation in the `target/doc` directory. For convenience, running `cargo doc -open` will build the HTML for the current crate’s documentation (as well as the documentation for all of your crate’s dependencies) and open the result in a web browser.

9. User manual

9.1. System requirements

The development process gives the user much flexibility when it comes to the choice of an operating system. Any tool in the project's toolchain is available on all major operating systems, whether it's Windows, GNU/Linux or macOS.

To install Rust along with Cargo please follow the official installation guide for your operating system available at

<https://www.rust-lang.org/tools/install>

Node.js and npm can be downloaded through

<https://www.npmjs.com/get-npm>

Wasm-pack's installer and installation guide can be found at

<https://rustwasm.github.io/wasm-pack/installer/>

All three tools must be installed on the machine of the user in order to run the application locally.

9.2. Installation guide

The source code of the application can be obtained either from the CD appended to the thesis, or from the official repository hosted on GitHub at

<https://github.com/karolbelina/parser-parser>

which can be downloaded as a ZIP file, or cloned from the command-line with the `git clone` command (in this case, make sure that Git is installed).

After downloading, please make sure all necessary components are in place. To list all files in the project's directory, first change the current directory to `parser-parser` with the `cd` command and then list all files using the `ls` command on Unix-like systems and `dir` on Windows. Listing 9.1 shows the process of downloading the repository from GitHub and listing all files in that repository from the command-line.

```
$ git clone https://github.com/karolbelina/parser-parser.git
$ cd parser-parser
$ ls
Dockerfile  LICENSE  README.md  app  base  cli  ebnf
```

Listing 9.1. Command-line output after downloading the source code on Ubuntu 20.04.

To run the web application, the user must first install npm dependencies of the `app` component of the project. To do this, one must change the current directory to `app`, and then, by using `npm`, install the necessary dependencies with the `npm install` command. This requires a connection with the internet and will take several minutes to complete. To finally build and run the application, the user must run the `npm run start` command, which will start the build process of the web application. This will build all Rust components of the application, as well as the Svelte application itself, and can take several additional minutes. After completion, the user should see the “Your application is ready” notification, and below, the local IP address where the application is currently being served, as seen in Listing 9.2. To see the application, the user should copy the address and paste it into the address bar in any supported web browser.

```
$ cd app
$ npm install
$ npm run build
```

Your application is ready~!

- Local: `http://localhost:5000`
- Network: Add `--host` to expose`

----- LOGS -----

Listing 9.2. Installation of the dependencies and building the web application.

9.3. Usage guide

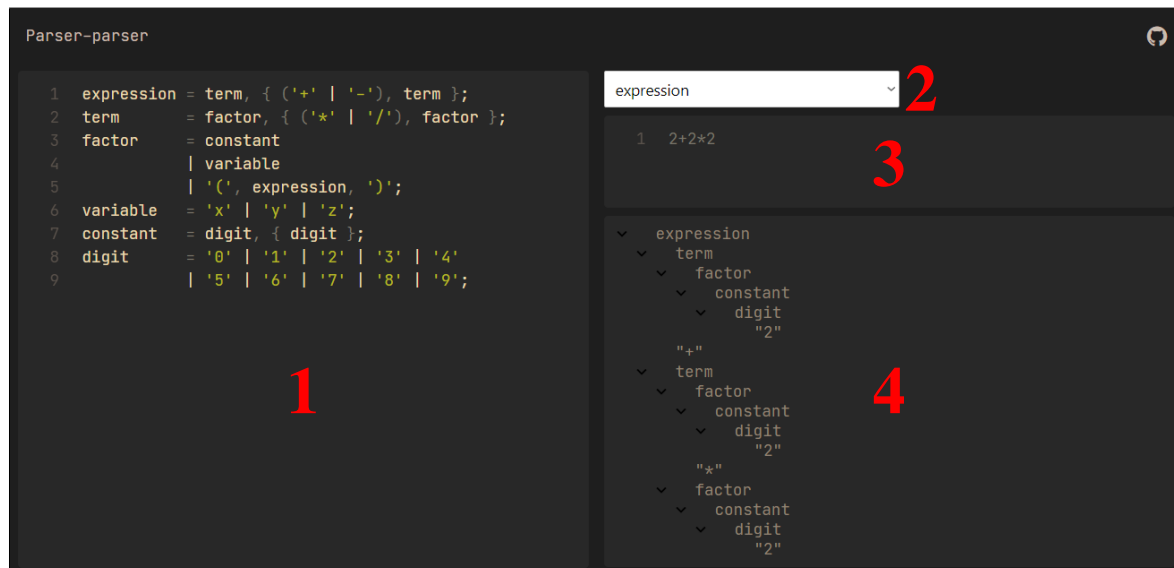


Figure 9.1. Screenshot of the Parser-parser web application with annotations in red.

Figure 9.1 shows the web application served at the local address. The application window is divided into several sections, all of which have been annotated on Figure 9.1:

1. The grammar window, where the user can provide a grammar definition in the EBNF format.

2. The initial production rule selection. After providing a valid grammar, the user can select an initial production rule from the ones provided.
3. The test string window, where the user, after providing a valid grammar, can test whether an arbitrary input string belongs to the language generated by that grammar.
4. The visualization of the parse tree, shown only if the user provides a string that belongs to the language generated by the grammar. The visualization is interactive — the nodes of the parse tree can be collapsed and expanded.

The screenshot presents an example use case of the application, in which the user clicked on the grammar window (1) with the left mouse button and began to type the definition of a desired grammar with the use of the keyboard. The grammar, defined as

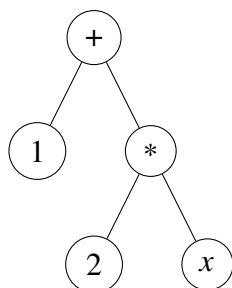
```

1 expression = term, { ('+' | '-'), term };
2 term       = factor, { ('*' | '/'), factor };
3 factor     = constant
4             | variable
5             | '(', expression, ')';
6 variable   = 'x' | 'y' | 'z';
7 constant   = digit, { digit };
8 digit      = '0' | '1' | '2' | '3' | '4'
9             | '5' | '6' | '7' | '8' | '9';

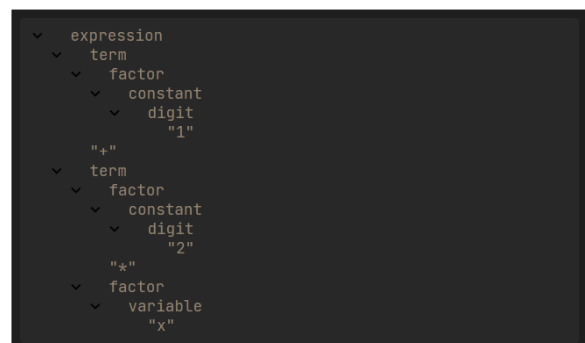
```

is able to generate simple arithmetic expressions made out of natural numbers and variables x , y , and z . The supported operators are: addition, subtraction, multiplication, and division, all of which keep an appropriate operator precedence order. The expressions can be parenthesized to alter the natural order of precedence. The precedence support for operators is built into the definition of the grammar — it is not a special behavior of Parser-parser.

After providing the grammar, the user clicked on the initial rule selection (2) and selected the initial production rule to be the “expression” rule. The user had the option to select any of the six rules defined in the grammar, which are updated in real-time, but has chosen the rule related to expressions. In the test string window (3) the user provided a test string “2+2*2” to check if it belongs to the language generated by the previously provided grammar.



(a) Expected parse tree structure.



(b) Parse tree in Parser-parser.

Figure 9.2. A parse tree of a simple arithmetic expression and its operator precedence.

As the user typed the input string, the visualization of the parse tree (4) was being updated. Initially, the output showed “failure”, as the initial empty string did not belong to

the language, but as the user typed more and more characters into the test string window, the output window updated accordingly. By the time the user typed the entire desired test string, the output window showed a visualization of the parse tree of the “2+2*2” expression. Figure 9.2 shows an expected, mathematically correct parse tree of a slightly altered $1 + 2 \cdot x$ expression, as well as the parse tree of the same expression “1+2*x” in Parser-parser. The expected parse tree and the one generated in Parser-parser have the same structure, where the multiplication operation would be evaluated first, as it is one level below the addition operator.

The user had the option to collapse or expand individual nodes of the parse tree dynamically in case the resulting parse tree was too large to visualize as a whole. At any point, the user had the option to change the initial production rule of the grammar, as well as the grammar itself.

```

1 digit excluding zero
2   = "1" | "2" | "3" | "4" | "5"
3   | "6" | "7" | "8" | "9";
4 digit
5   = "0" | digit excluding zero;
6 natural number = digit excluding zero, { digit };
7 integer
8   = "0" | [ "-" ], natural number;

```

integer

rule digitexcludingzer is undefined

(a) Undefined rule detection.

```

1 digit = '0' | '1' | '2' | '3' | '4'
2       | '5' | '6' | '7' | '8' | '9';
3 integer = digit, { digit };
4 value = integer | '(', expr, ')';
5 product = expr, {'*' | '/'}, expr;
6 sum = expr, rule product is left recursive (product -> expr -> product)
7 expr = product | sum | value;

```

expression

1 2+3

(b) Left recursion detection and its cycle.

```

1
2
3 a = 5 * "a", 5 | "b";

```

repetition symbol expected

(c) Simple syntax error.

Figure 9.3. Three examples of errors detected by Parser-parser.

As the user provides a definition of a grammar into the grammar window, the syntax and semantics of that grammar are being automatically validated. In case the definition is not valid, the reason of that is automatically reported to the user. Figure 9.3 shows three examples of invalid grammars, where the exact reason and its place in the definition is highlighted by underlining it with the red color. The user can then hover the cursor over the highlighted area to read the error, and then correct it accordingly.

As long as the definition of the grammar is invalid, the user cannot choose an initial production rule or test an input string.

10. Summary

The thesis described the process of designing and implementing a tool for formulating and visualizing context-free grammars in the Extended Backus-Naur form. First, it provided a basic analysis of the advantages of the Extended Backus-Naur form notation over its predecessor — Backus-Naur Form. It transformed the official specification of the ISO/IEC 14977 standard by removing inaccuracies and ambiguities, as well as to make it more suitable for modern software environments by providing support for Unicode as the base character set of the specification. A definition of tokens, as well as the grammar in the form of an abstract syntax tree, was proposed to specify the structure of the Extended Backus-Naur form in the implementation phase of the project.

The design process carried out by the author was a key step in designing and implementing the project. Planned use cases allowed the developer to map out key features of the final product and defined the “correct” way of performing tasks by the user, which related directly to the functionalities of the product. The thesis defined various UML diagrams to facilitate the design phase — diagrams such as activity and sequence diagrams have been used to map out links and processes that happen between user interactions in a more visual manner. Sequence diagrams provided a way of visualizing the internal processes of the application.

The implementation phase of the thesis described the process of implementing the business logic of the application. Various components of the EBNF parser: the lexer, the parser, and the preprocessor, were implemented in the Rust programming language, which were then compiled to WebAssembly to allow them to run in a web environment. The checker component, which checks if a given input string provided by the user belongs to the language generated by a given grammar, was produced and compiled in a similar way. The business logic was put together in the form of a command-line application as a most basic method of interfacing with the implemented library. To fully address the issue of accessibility of the application, a web-based application with a graphical user interface was implemented and deployed on a web page and can run in a browser without the need for installation on the user’s device. The web application was implemented in the Svelte framework, which incorporated the markup of the site, the CSS styling, as well as JavaScript scripts used for user interaction.

The phase described the process of tokenization and proposed it in the form of a Deterministic Finite Automaton. The lexer takes into account the proper interpretation of Unicode graphemes. A technique of combinatory parsing was then used to combine the tokens into an abstract syntax tree. The implementation phase included the definition of an algorithm for detecting left recursion and undefined production rules in a grammar. Occurrences of these are reported to the user in a readable way that highlights the exact positions of errors in the input string. The web application provides a basic code editor for inputting the grammar, with such features as autocompletions and syntax highlighting. The tool is able to visualize the parse tree to the user in an interactive way. The final product is welcoming to new users than other similar tools. Users are able to access the website and use the tool right from their browser without needing to carry out a complex setup process, thanks to the deployment on the GitHub Pages hosting service for static sites.

The thesis also analyzed the quality of the project by defining various unit and integration tests to ensure the correctness of the business logic of the application. The project utilized property-based testing as a useful addition to a test suite that generates hundreds of test cases automatically. Integration tests helped to combine implemented components into a single library and compile it to the WebAssembly target. With the help of Google Lighthouse, the web application underwent an external review to check its quality, progress, or adherence to plans, modern web standards, and regulations. The thesis provided a way to analyze the business logic of the application in terms of the so-called cognitive complexity.

The designed and implemented system gives the opportunity to extend it with various functionalities. First of all, the tool is designed from the start with other grammar specifications in mind. Specifications like BNF, ABNF, Wirth syntax notation, or other EBNF variations like the EBNF notation from the W3C Extensible Markup Language (XML) are a great fit for further expansion and a way to extend the list of supported languages. Each language would have an associated parser that compiles it down to a common notation. Secondly, the functionalities of the application could be expanded by providing a code-generation aspect to the grammars given by users. Users would be able to “kickstart” their parsers by generating code in various popular programming languages that represent the provided grammar (again, not necessarily in the EBNF notation).

Bibliography

- [1] AHMAD, M. O., MARKKULA, J., AND OIVO, M. Kanban in software development: A systematic literature review. In *2013 39th Euromicro Conference on Software Engineering and Advanced Applications* (2013), pp. 9–16.
- [2] AHO, A., ULLMAN, J. D., LAM, M. S., AND RAVI, S. *Kompilatory: reguły, metody, narzędzia*. Wydawnictwo Naukowe PWN, Warsaw, 2019.
- [3] AHO, A. V. Algorithms for finding patterns in strings. In *Algorithms and Complexity*. Elsevier, 1990, pp. 255–300.
- [4] CHACON, S., AND STRAUB, B. *Pro Git*. Apress, Berkeley, CA New York, NY, 2014.
- [5] CHOMSKY, N. Three models for the description of language. *IEEE Transactions on Information Theory* 2, 3 (Sept. 1956), 113–124.
- [6] CODEMIRROR TEAM. CodeMirror homepage. <https://codemirror.net/>. Last accessed 09.11.2020.
- [7] COUPRIE, G. Nom GitHub page. <https://github.com/Geal/nom>. Last accessed 08.11.2020.
- [8] COUPRIE, G. Nom, a byte oriented, streaming, zero copy, parser combinators library in Rust. In *2015 IEEE Security and Privacy Workshops* (May 2015), IEEE.
- [9] DAVIS, M., AND CHAPMAN, C. Unicode standard annex #29. <http://www.unicode.org/reports/tr29/>. Last accessed 08.11.2020.
- [10] DIB, F. Regex101 homepage. <https://regex101.com/>. Last accessed 24.10.2020.
- [11] DRIESSEN, V. A successful Git branching model. <https://nvie.com/posts/a-successful-git-branching-model/>, 2010. Last accessed 07.11.2020.
- [12] FERROUS SYSTEMS. rust-analyzer homepage. <https://rust-analyzer.github.io/>. Last accessed 07.11.2020.
- [13] FOKKER, J. Functional parsers. In *Advanced Functional Programming*. Springer Berlin Heidelberg, 1995, pp. 1–23.
- [14] FORD, B. Parsing expression grammars. In *Proceedings of the 31st ACM SIGPLAN-SIGACT symposium on Principles of programming languages - POPL '04* (2004), ACM Press.
- [15] GALLANT, A., ET AL. Quickcheck GitHub page. <https://github.com/BurntSushi/quickcheck>. Last accessed 13.12.2020.

- [16] GIT COMMUNITY. Git homepage. <https://git-scm.com/>. Last accessed 07.11.2020.
- [17] GITHUB INC. GitHub homepage. <https://github.com/>. Last accessed 08.11.2020.
- [18] GUILLAUME, P., ET AL. StructOpt GitHub page. <https://github.com/TeXitoi/structopt>. Last accessed 13.12.2020.
- [19] HASKELL TEAM. Haskell homepage. <https://www.haskell.org/>. Last accessed 13.12.2020.
- [20] HEISLER, B., ET AL. Criterion.rs GitHub page. <https://github.com/bheisler/criterion.rs>. Last accessed 13.12.2020.
- [21] HOPCROFT, J., MOTWANI, R., AND ULLMAN, J. D. *Wprowadzenie do teorii automatów, języków i obliczeń*, second ed. Wydawnictwo Naukowe PWN, Warsaw, 2012.
- [22] ISO/IEC. *ISO/IEC 14977:1996(E) — Information technology, syntactic metalanguage, Extended BNF*. Geneva, 1996.
- [23] JOHNSON, W. L., PORTER, J. H., ACKLEY, S. I., AND ROSS, D. T. Automatic generation of efficient lexical processors using finite state techniques. *Communications of the ACM* 11, 12 (Dec. 1968), 805–813.
- [24] KLABNIK, S., AND NICHOLS, C. *The Rust programming language*. No Starch Press, Inc, San Francisco, 2018.
- [25] LEIJEN, D., AND MEIJER, E. Parsec: Direct style monadic parser combinators for the real world. <https://www.microsoft.com/en-us/research/wp-content/uploads/2016/02/parsec-paper-letter.pdf>, 2001. Last accessed 13.11.2020.
- [26] MEDUNA, A. *Formal Languages and Computation: Models and Their Applications*. CRC Press, Taylor & Francis Group, Boca Raton, 2014.
- [27] MICROSOFT. Visual Studio Code homepage. <https://code.visualstudio.com/>. Last accessed 07.11.2020.
- [28] NPM TEAM. Npm homepage. <https://www.npmjs.com/>. Last accessed 13.12.2020.
- [29] PEST TEAM. Pest homepage. <https://pest.rs/>. Last accessed 14.11.2020.
- [30] ROLLUP TEAM. Rollup homepage. <https://rollupjs.org/guide/en/>. Last accessed 09.11.2020.
- [31] RUST AND WEBASSEMBLY TEAM. Wasm-bindgen GitHub page. <https://github.com/rustwasm/wasm-bindgen>. Last accessed 08.11.2020.
- [32] RUST TEAM. Crates.io homepage. <https://crates.io/>. Last accessed 13.12.2020.
- [33] RUST TEAM. Rust homepage. <https://www.rust-lang.org/>. Last accessed 08.11.2020.
- [34] SHIN, E. Svelte-tree GitHub page. <https://github.com/esinx/svelte-tree>. Last accessed 13.11.2020.

- [35] SIPSER, M. *Wprowadzenie do teorii obliczeń*. Wydawnictwa Naukowo-Techniczne, Warsaw, 2020.
- [36] STACK OVERFLOW. Developer Survey Results 2018. <https://insights.stackoverflow.com/survey/2018>. Last accessed 07.11.2020.
- [37] STACK OVERFLOW. Developer Survey Results 2019. <https://insights.stackoverflow.com/survey/2019>. Last accessed 07.11.2020.
- [38] SVELTE TEAM. Svelte homepage. <https://svelte.dev>. Last accessed 08.11.2020.
- [39] SVELTE TEAM. Svelte Language Tools GitHub page. <https://github.com/sveltejs/language-tools>. Last accessed 07.11.2020.
- [40] SWIERSTRA, S. D. Combinator parsing: A short tutorial. In *Language Engineering and Rigorous Software Development*. Springer Berlin Heidelberg, 2009, pp. 252–300.
- [41] UNICODE-RS TEAM. Unicode-segmentation GitHub page. <https://github.com/unicode-rs/unicode-segmentation>. Last accessed 08.11.2020.
- [42] WASM-PACK TEAM. Wasm-pack homepage. <https://rustwasm.github.io/wasm-pack/>. Last accessed 13.12.2020.
- [43] WEBASSEMBLY TEAM. WebAssembly homepage. <https://webassembly.org/>. Last accessed 08.11.2020.
- [44] WHEELER, D. A. Don't Use ISO/IEC 14977 Extended Backus-Naur Form (EBNF). <https://dwheeler.com/essays/dont-use-iso-14977-ebnf.html>. Last accessed 09.12.2020.

List of Figures

2.1	The Chomsky hierarchy visualized.	6
2.2	The visualization of the EBNF parsing process.	13
3.1	Screenshot of the Regex101's matching functionality. The user provided the " <code>\s\s+</code> " regular expression, which matched every occurrence of two or more consecutive space characters in the test string.	15
3.2	Screenshot of a basic error in the regular expression reported by Regex101.	16
3.3	Screenshot of Pest's online editor example error report.	17
3.4	Screenshot of the input and output windows in Pest's online editor.	17
4.1	The use case diagram.	22
4.2	The requirements traceability graph.	23
4.3	The activity diagram of <i>UC1</i> Specifying the grammar.	25
4.4	The activity diagram of <i>UC2</i> Specifying the input string.	25
4.5	The activity diagram of <i>UC3</i> Interacting with the visualization.	26
4.6	The activity diagram of <i>UC4</i> Changing the initial rule.	26
4.7	The sequence diagram representing the specification of the grammar.	26
4.8	The logical architecture of the system represented with a UML component diagram.	27
4.9	The physical architecture of the system represented with a UML deployment diagram.	28
4.10	The user interface sketch.	28
5.1	Screenshot of the command-line interface of the Git version control system.	30
5.2	Screenshot of one of the project's kanban boards on GitHub.	31
5.3	A simplified directory tree of the Parser-parser project.	37
5.4	The DFA representation of the lexer. Note that this DFA does not support nested comments. Σ is the set of all characters, A is the set of all alphabetic characters, N is the set of all numeric characters, D is the set of all ten decimal digits, and W is the set off all whitespace characters. Set notation on the transitions is omitted.	40
5.5	A screenshot of the visualized parse tree.	50
6.1	Probability Distribution Function chart generated by Criterion.rs.	56
6.2	Regression plot generated by Criterion.rs.	56
6.3	Google Lighthouse's scores of the deployed web application.	57
6.4	Number of functions with various degrees of cognitive complexity.	59
9.1	Screenshot of the Parser-parser web application with annotations in red.	68
9.2	A parse tree of a simple arithmetic expression and its operator precedence.	69
9.3	Three examples of errors detected by Parser-parser.	70

List of Tables

2.1	Tokens types from the modified EBNF notation.	10
2.2	Node types for the abstract syntax tree of EBNF.	13
4.1	The functional requirements of the project, their features, and priorities. . .	19
4.2	The non-functional requirements of the project and their priorities.	20
4.3	The user stories.	21
4.4	Descriptions of the use cases.	22
4.5	Use case scenario of <i>UC1</i> Specifying the grammar.	23
4.6	Use case scenario of <i>UC2</i> Specifying the input string.	24
4.7	Use case scenario of <i>UC3</i> Interacting with the visualization.	24
4.8	Use case scenario of <i>UC2</i> Specifying the input string.	24
6.1	Tests cases related to the lexer along with the expected results.	52
6.2	Performance metrics generated by Lighthouse.	58

List of Listings

2.1	Example expression grammar in BNF.	8
5.1	Definition of the <code>Token</code> type in Haskell.	38
5.2	Definition of the types related to the AST in Haskell.	39
5.3	The <code>sequence</code> parser.	41
5.4	The definition of the <code>Node</code> type.	45
5.5	The output of the program ran with the <code>--help</code> flag.	46
5.6	The <code>StructOpt</code> struct defining the command-line arguments for the program.	47
5.7	Sample of the output of the program ran in REPL mode.	47
5.8	The definition of the EBNF parser struct that encapsulates the grammar.	48
5.9	Contents of <code>ParseTree.svelte</code> , the parse tree visualization component.	50
6.1	A unit test related to the lexer testing the proper tokenization of the input string.	51
6.2	A unit test related to the parser, where the ability to turn a string of tokens into an AST is tested.	52
6.3	A preprocessor unit test testing the left recursion detection in an AST.	53
6.4	A QuickCheck test testing arbitrary inputs on the lexer.	54
6.5	A QuickCheck test related to the parser, testing strings of arbitrary tokens.	54
6.6	An integration test ran in a headless browser, which tests various grammars in a textual form.	55
6.7	A benchmark testing the speed of parsing a sample grammar.	55
7.1	The deployment part of the GitHub Action workflow file related to CI/CD.	63
9.1	Command-line output after downloading the source code on Ubuntu 20.04.	67
9.2	Installation of the dependencies and building the web application.	68
A.1	Modified version of the EBNF language specification defined in [22].	83

A. Modified specification of the ISO/IEC 14977

Listing A.1 shows the modified version of the specification of the ISO/IEC 14977 standard. The modification is described in detail in Section 2.3 of the thesis.

```
1 character
2   = ? any Unicode non-control character ?;
3 letter
4   = ? any Unicode alphabetic character ?;
5 digit
6   = ? any Unicode numeric character ?;
7 whitespace
8   = ? any Unicode whitespace character ?;
9 comment
10  = '(*', {comment | character}, '*)';
11 gap
12  = (whitespace | comment), {whitespace}, {{comment}, {whitespace}};
13 identifier
14  = letter, {{whitespace}, letter | digit};
15 factor
16  = [[gap], digit, {{whitespace}, digit}, [gap], '*'],
17    [gap], [(identifier
18      | ('[' | '(/', alternative, (']' | '/)')
19      | ('{' | '(:', alternative, ('}' | ':)')
20      | '(', alternative, ')')
21      | '"', character - '"', {character - '"', '"'}
22      | "'", character - "'", {character - "'", "'}
23      | '?', {{whitespace}, character - '?', '?'}, [gap]]];
24 term
25   = factor,
26     ['- ', ? a factor that could be replaced
27       by a factor containing no identifiers ?];
28 sequence
29   = term, {' ', term};
30 alternative
31   = sequence, {'|' | '/' | '!'}, sequence;
32 production
33   = [gap], identifier, [gap], '=', alternative, ( ';' | '.' ), [gap];
34 grammar
35   = production, {production};
```

Listing A.1. Modified version of the EBNF language specification defined in [22].

B. CD contents

The disc included with the thesis contains:

- the source code of the project presented in Section 5.1.2 — 320 KB,
- the source code of the thesis in \LaTeX — 1.28 MB,
- the thesis PDF document — 980 KB,