

FELADATKIÍRÁS

Modular Interpreter Written in JavaScript

To learn programming, there are some basic tools everyone needs - a compiler and some text editor or an IDE. To acquire and use these can sometimes be non-trivial, especially for the younger generations. Most of the learning material is on the Internet, but to solve the exercises one needs to use the aforementioned programs installed on their computer. By eliminating this extra step and putting the programming environment closer to the materials the life of new pupils can be made easier. The aim of this thesis is to create a modular interpreter in JavaScript, using which anyone can run any kind of code inside the browser (and since it is modular, it can be easily extended with new languages). The role of the interpreter is to build an Abstract Syntax Tree (AST) from the given source code, and to process the tree to run the equivalent JavaScript code. Using this interpreter anyone could easily build educational websites that can focus on teaching programming, since the learning materials and the environment will be in the same place - in the browser. The requirements for this thesis are the following:

- Create the core part of the interpreter
- Create an example language module



Budapest University of Technology and Economics

Faculty of Electrical Engineering and Informatics

Department of Automation and Applied Informatics

Modular Interpreter Written in JavaScript

THESIS

By

Heilig Benedek

Supervisor

Kundra László János

December 8, 2016

Contents

Kivonat	7
Abstract	9
Introduction	11
1 Coding conventions	15
1.1 About JavaScript	15
1.2 Problems with ECMAScript 6 classes	16
1.3 Node.js modules	16
2 Lexical analysis	19
2.1 Using regular expressions	19
2.2 The Lexer module	20
2.2.1 Lexer subclasses	21
2.2.2 Stateless mode	22
2.2.3 Stateful mode	22
2.2.4 Skipping whitespaces	22
2.2.5 The Lexer class	22
2.3 Possible enhancements	23
3 Syntactical analysis (Parsing)	25
3.1 Describing the grammar	25
3.2 The parsing algorithm	26
3.2.1 Naive implementation	27

3.2.2	LALR - LookAheadLR	28
3.3	The Parser module	31
3.3.1	The Parser subclasses	32
3.3.2	The Parser class	33
3.4	Possible enhancements	33
4	The Language module and two examples	35
4.1	The Language module	35
4.2	Examples	35
4.2.1	The math parser	35
4.2.2	The Logo subset parser	37
5	Summary, conclusions	39
6	Problems, possible improvements	41
	Acknowledgement	43
	Bibliography	45

HALLGATÓI NYILATKOZAT

Alulírott *Heilig Benedek*, szigorló hallgató kijelentem, hogy ezt a szakdolgozatot meg nem engedett segítség nélkül, saját magam készítettem, csak a megadott forrásokat (szakirodalom, eszközök stb.) használtam fel. Minden olyan részt, melyet szó szerint, vagy azonos értelemben, de átfogalmazva más forrásból átvettem, egyértelműen, a forrás megadásával megjelöltem.

Hozzájárulok, hogy a jelen munkám alapadatait (szerző(k), cím, angol és magyar nyelvű tartalmi kivonat, készítés éve, konzulens(ek) neve) a BME VIK nyilvánosan hozzáférhető elektronikus formában, a munka teljes szövegét pedig az egyetem belső hálóján keresztül (vagy autentikált felhasználók számára) közzétegye. Kijelentem, hogy a benyújtott munka és annak elektronikus verziója megegyezik. Dékáni engedéllyel titkosított diplomatervek esetén a dolgozat szövege csak 3 év eltelte után válik hozzáférhetővé.

Budapest, December 8, 2016

Heilig Benedek
hallgató

Kivonat

A programozás megtanulása akadályokkal teli folyamat tud lenni a kezdetekben. A szakdolgozatomban erre a problémára kívánok megoldást javasolni egy JavaScript alapú moduláris értelmező elkészítésével, ami bármilyen programozási nyelvet képes értelmezni — amennyiben elkészül a megfelelő modul hozzá — és átfordítani JavaScript kóddá.

A megoldás egy reguláris kifejezés és állapotgép alapú lexikális analízátor modulból — neve Lexer —, és egy szintaktikai elemző modulból áll, aminek neve Parser és az LALR elemző algoritmust implementálja, majd ennek segítségével dolgozza fel a bemenetét a BNF-ben megadott nyelvtan szerint.

A szakdolgozatomban leírom a megoldás menetét és elmagyarázom a használt algoritmusok működését — és hogy miért az adott algoritmus mellett döntöttem. Emellett kitérek a használt technológia sajátosságaira hogy a JavaScript mélyebb ismerete nélkül is könnyen értelmezhető legyen a dolgozatom.

Abstract

Learning how to program can be a problematic experience for new students. In my thesis I propose a solution to the problematic first experiences by creating a modular interpreter that can interpret any kind of programming language — if a module is created for it — and translate that code to JavaScript.

The solution consists of a regular expressions and state machine based lexical analyser module called Lexer, and a Parser module that implements the LALR parsing algorithm to parse grammars defined in BNF.

In my thesis I describe the solution and explain the used algorithms, and the choice behind them. I also give a brief explanation of the used technologies to make the solution easy to understand even for those who do not know much about JavaScript.

Introduction

The main goal of this thesis is to create a modular interpreter in the native language of the modern Internet, JavaScript. The idea behind an interpreter such as this is that usually for someone new to programming, it is a hassle to download and install the various compilers and IDEs needed to even begin the actual programming phase, and this could be made easier with web technologies.

Nowadays, everyone has a laptop or a smartphone (most of the times both) capable of browsing the World Wide Web, and all these devices have a browser to achieve this task. Most dynamic websites use JavaScript as a language to manipulate the content of the page (to make it dynamic), so all these browsers have a JavaScript engine built into them. I intend to use this capability to bring a more seamless first experience to those who want to learn programming and are at the beginning of their journey.

This is why I set out to create an interpreter that is:

- Written in JavaScript, thus available to be used on every kind of device
- Modular, thus making the job of creators of educational tools easier by abstracting the common things, and providing an easy way to allow execution of code inside the browser in any custom language that has a module created for it

To execute code written in a specific language, either a compiler (that translates directly into machine code, which is then executed), or an interpreter is used (which is a program that actively interprets the given source code line-by-line). JavaScript is an interpreted language — partly that is why it can be found in so many places — and thus the interpreter described in this thesis is a cross-interpreter — that is, it translates from an interpreted language to another interpreted language. But before any interpretation can begin, the code has to go through some transformations, mainly lexical analysis and syntactical analysis.

Lexical analysis is basically turning the source code into tokens using basic pattern matching rules, which makes the parsing stage easier. A token belongs to a class of tokens, which represent the meaning of the value of the token. The value of the token is a part of the input code that is paired with the class by pattern matching. A good and widely used tool for pattern matching is regular expressions, which is available in JavaScript as part of the language — which provides a convenient way to implement this feature of the analyser. In

case the writer of a module needs more complex rules, lexical analysers usually provide a way to use simple states and define transitions between these states to make expressing of these rules easier.

In the next stage called parsing (or syntactical analysis), the tokens are turned into a parse tree using a formal grammar. A parse tree represents the hierarchy of the matched rules. One common way to define the grammar is BNF (Backus-Naur form), which is a simple notation to describe grammars and their rules. The given grammar describes every type of statement in the language. In the definitions of the rules, the writer of the grammar can reference the token classes from the previous stage (rules are automatically added to the grammar to represent the classes). The parse tree then can be used to create an AST (Abstract Syntax Tree), which is a hierarchical representation of the structure of the input. The AST then can be used to interpret the input itself by calculating the values of the nodes of the tree, which basically represent a statement and the parameters of the statement. Some of these parameters have fixed values while others need some computation to produce a value first (which means repeating the same algorithm that calculates the value for a node for the node of the parameter itself).

To create a language module, the necessary steps are:

1. Creating the rules for the lexical analyser (defining token classes and the state machine if necessary)
2. Providing the grammar in BNF
3. Parsing the code, and iterating over the nodes of the produced AST

The last step is the actual interpretation, which can be a simple evaluation of an expression, or even a small virtual machine that keeps track of functions, variables, scopes and etc. This is the part that requires most attention when writing a language module.

Of course, something like this has been attempted many times. There are numerous solutions to creating parsers, although most of them are not written using JavaScript, they are native tools (written in compiled languages) to generate parsers. The list of some of the more popular tools include GNU Bison ¹, Jison ² (a clone of Bison in JavaScript), YACC ³ and ANTLR ⁴. But out of all of the available works, the most relevant to this thesis is Syntax ⁵, which has the same initial goals (to create a parser generator in JavaScript) and the same things in mind (educational purposes) as the idea behind this thesis. An blog post about Syntax was published just a week after the proposal for this thesis was handed in. The author of Syntax describes his motives in the post and the main differences from my thesis are the following:

¹GNU Bison <https://www.gnu.org/software/bison/>

²Jison - Your friendly JavaScript parser generator! <http://zaa.ch/jison/>

³Yet Another Compiler-Compiler <http://dinosaur.compilertools.net/>

⁴Another Tool For Language Recognition <http://www.antlr.org/>

⁵Syntax - language agnostic parser generator <https://medium.com/@DmitrySoshnikov/syntax-language-agnostic-parser-generator-bd24468d7cfc>

- He only set out to create a modular (language agnostic) parser generator
- He provides multiple parsing algorithms
- He provides multiple target languages, including JavaScript, but that only means that the parser that is generated is implemented in JavaScript (other targets include Python, Ruby and PHP)
- He provides a way to set operator precedence
- He does not provide a state based lexical analyser

In every other sense, the two projects are very similar, although I have only learned about Syntax during later stages of development.

Another similar project is an older project of mine — from which the idea came for this thesis —, which is a Logo interpreter that translates the code to JavaScript, and which was based on absolutely no research about parsers, and thus is very complicated and sometimes buggy — but generally works quite well. This project had two version, the first of which I built 4 years ago as a hobby — to help a friend of mine with her Logo homework (since the interpreter used in her class was a 16 bit program, and she had a 64 bit computer which made it impossible to run the interpreter), and the second version that built on the experience that came from building the first — so it was much more mature, and could handle Logo source reasonably well. The choice of language — other than the friend I mentioned — is because it is a good step towards understanding basic ideas of programming (in both functional and imperative styles), and also provides the user with something materialistic as a result (it is easy to create beautiful and complex drawings), which is great for young children to get a reward for their work, which is a great way to motivate them. One of the most beautiful things someone can draw with Logo easily is fractals (up to a certain level of depth), which also teach them a lot about recursion, which can be hard to understand for new pupils (fractals are a good way to visualize recursion).

Other solutions to running code on the web include cloud based solutions, which means that the code that the user inputs is actually compiled and ran on a remote virtual machine. While this is great in the sense that it is exactly as the "real thing", it has a great drawback - it requires an active connection to the Internet to use. My solution can be bundled into a desktop application (or the website can be downloaded) using tools like the electron ⁶ framework, or it can be integrated into a website. This allows the students to practice while on the go (and without an active Internet connection), which cloud based solutions can not allow.

During the implementation speed was not the main focus, since this is not a tool built for performance but rather one that is built to be available everywhere — but I did try to not be wasteful with the resources. To achieve modern, well readable source code and

⁶Electron <http://electron.atom.io/>

portability, I used the ECMAScript 6 standard, which is essentially a JavaScript standard that most modern browsers can interpret – but only the newer versions. I developed my solution in Node.js modules, which can then be converted into code that works on older browsers using Babel ⁷ and Browserify ⁸, which are Node.js modules themselves to convert to older standards and to bundle modules into one source file for browsers respectively.

Chapter 1 will outline the basic coding conventions I followed during development. In Chapter 2, I will describe the implemented lexical analyser, and it's inner workings. In Chapter 3, I will write about the parser module, and the implemented parsing algorithm. In Chapter 4, I will describe the two example modules I have created - a math expression parser, and a simple Logo subset interpreter. In chapter 5 I will give a summary of the chapters before and some conclusions while in chapter 6 I will write about problems and possible improvements for the future.

⁷Babel <https://github.com/babel/babel>

⁸Browserify <http://browserify.org/>

Chapter 1

Coding conventions

While JavaScript is a modern language, there are some parts of it that evolved quite slowly compared to other languages. The most prominent part is classes (although object-orientation is a core part of the language), which have been properly introduced into the standard with ECMAScript version 6. There are shortcomings of this implementation, although it improved on the previous version greatly. Before ECMAScript 6, a programmer who wanted to create an object or a class in JavaScript had to use a weird syntax that relied on anonymous classes and factory functions, since the language did not really have a class keyword. This allowed object oriented development, but only to a certain extent. With ECMAScript 6, a proper class keyword was introduced, alongside with keywords like `extend` and `constructor`, which made proper object-oriented code easier to write and read.

1.1 About JavaScript

JavaScript is a fairly modern interpreted language with weak types — which makes development in JavaScript easy, but it also makes writing bugs easy since a variable can hold any type of data at any point in the execution process.

The language was developed in 1995 Brendan Eich ¹ in only ten days while he worked for Netscape — one of most popular early web browsers. The fact that it was developed in such short time is remarkable, but it means that there are some oddities in the language. A great presentation ² about these parts of the language was given by Gary Bernhardt in 2012. Some of the most interesting odd behaviours originate from the language being weakly typed — during comparison of values and while using operators, the interpreter tries to cast the values into other types to run the operation it was asked successfully. This ends up creating weird things, such as adding an empty array and an empty object together, which is a valid operation in JavaScript. In the case of the empty array being the first operand (`[] + {}`), the empty array is interpreted as an empty string, and so the empty object is

¹Brendan Eich's website <https://brendaneich.com/>

²Talk on JavaScript oddities <https://www.destroyallsoftware.com/talks/wat>

cast into a string, which yields `[object Object]` as the result. Surprisingly flipping the operands `{ } + []` yields `0` as a result, making the plus operand not commutative in this case.

While weak types can be problematic, they can also be a helpful tool during development, since a lot of operations are allowed this way that would require explicit casting in other languages — but the developer has to keep these weird rules in mind when doing comparison and such operations.

Nowadays JavaScript is nicely standardized, with the standard being continually improved upon. This standard is called ECMAScript ³ on which work began in 1997. Since then many versions were released, of which this thesis uses the 6th version, ECMAScript 2015.

1.2 Problems with ECMAScript 6 classes

While they are a great improvement on the previous situation, ECMAScript 6 classes have some shortcomings - a programmer can not really define a private member for the class, since that keyword is not present in the language. Getters and setters can be used, as well as static variables, but simulating private members is not easy.

Example 1.1: *"A typical class"*

```
class ClassName {
  constructor() {
    this.member = 'x'
    this._private = 1
  }

  get private() {
    return this._private
  }

  doSomething() {
    // do something here
  }
}
```

One of the ways of dealing with the absence of private members is to follow a naming convention — name everything that should be private with an underscore prefix. In fact this is the convention I followed during development.

1.3 Node.js modules

I developed the solution by splitting parts into their own Node.js modules, which allowed me to easily check my solution without even starting a browser. Node.js is a relatively new technology that provides an JavaScript interpreter — the same interpreter that Google Chrome uses, the V8 JavaScript engine ⁴ — in a binary form with interfaces for lower level

³ECMAScript 6 specification <http://www.ecma-international.org/ecma-262/6.0/ECMA-262.pdf>

⁴V8 JavaScript engine <https://chromium.googlesource.com/v8/v8.git>

tasks — such as networking — readily available. This allows programmers to create server applications in JavaScript, which means a website can be developed with the frontend and the backend written in the same language — sometimes even sharing code. Sadly the Node.js module syntax is not something that browsers can deal with by themselves, but this problem can be eliminated by using the tools described in the end of the introduction chapter.

Chapter 2

Lexical analysis

Lexical analysis is the process in which an input string is turned into tokens. This is done to make the next step — parsing — easier. For example, the parser can deal with words instead of having to try to match the grammar's rules character-by-character. This can greatly reduce the complexity of the grammar, which means it makes it easier to write a grammar for a specific language.

2.1 Using regular expressions

Regular expressions are a common tool used in lexical analysis (e.g. the lexical analyser of the popular parser generator Yacc uses this technique [8]), since during the analysis, we need to determine what class of tokens does the beginning of the given source string translates to — and regular expressions provide a simple way to write patterns that we can match against our input.

Regular expressions is a well established way of pattern matching. The concept was created in the 1950s by Stephen Cole Kleene [7]. A regular expression is a string of characters that describe a pattern, which then can be matched against a string to do various operations like finding the first matching position, finding all matching positions, replacing said matches with other strings, etc. The language of regular expressions have numerous rules which allows us to create very complicated patterns. The regular expression engine that JavaScript implements is well documented ¹, with most of the features of regular expressions available to the programmer. As an example, here are some patterns and their meaning:

- `/[0-9]+\.[0-9]+/` match strings looking like floating point numbers (a dot indicates the decimal point)
- `/[a-z][a-z0-9]*/i` match a string that must start with a character, has to be at least 1 character long and otherwise can contain alphanumeric characters (lower and upper case too)

¹JavaScript RegExp reference https://developer.mozilla.org/en/docs/Web/JavaScript/Guide/Regular_Expressions

- `/[0-9]+(?![0-9]*\.[0-9]+)/` match string that represent an integer (a sequence of numbers that is not followed by any more numbers and a dot)

The first pattern makes use of ranges. Instead of having to write out all the numbers, we can just say `[0-9]`, which will match any number. Anything in a pair of square brackets is a group of possible characters in that position, for example `[abc]` would only match `a`, `b` or `c`. The `+` sign is a quantifier, expressing that at least one character from the group before it must be present. We can see another quantifier in the second rule, the `*` (which is called a Kleene star) means zero or more matches of the group it is attached to. Of course these quantifiers can be attached to single symbols too, in which case the square brackets are not needed. In the first and the last group, we can observe a `\.` combination. This is required because `.` has a special meaning in regular expressions (it matches any character), so if we want to only match a dot, we need to escape it, which is done by adding a `\` before it. In the second rule, we can see a modifier at the end of the expression. That modifier `i` means that the matching should be done in a case-insensitive way (the default mode is case-sensitive). In the last rule, we can see one of the more complicated kinds of expressions - at the end of the expression, there is a pattern `(?![0-9]*\.[0-9]+)`, which is a negative lookahead pattern. This means that the whole expression will only match if the inside of this pattern (some numbers, a dot and some more numbers in this case) can not be matched after the examined string. The negative lookahead will not be included in the result. Several other tools are available in regular expressions to write complicated patterns (such as positive lookahead, capture groups, matching exactly `n` times, etc.), which can make lexical analysis easier, since the regular expression engine is doing some of the work for us. The JavaScript regular expression engine does not however implement lookbehinds, which is the same concept as lookaheads, but instead peeking back to before a given pattern.

Because of the expressiveness of regular expressions, my implementation of the lexical analyser will be fairly simple, without the need of explicitly working with lookahead symbols when transitioning between states and matching strings to their classes, since these special rules can be handled by carefully writing the regular expressions.

2.2 The Lexer module

The lexical analyser module is named `Lexer` — which is a JavaScript class with several sub-classes. `Lexer` implements a regular expression based matching system with an additional ability to specify states, and transitions between these states. This is optional however, and `Lexer` can work without keeping track of states.

2.2.1 Lexer subclasses

Lexer.Token

The Token class represents a token — it has a value, which is the parsed string that belongs to the token, and a class, which is a reference to the TokenClass the token belongs to.

Lexer.TokenClass

The TokenClass represents a class of tokens — it has a name, an associated regular expression, and a state, which defaults to the DefaultState. It also has a match method that tries to match against a given string, and returns the consumed part of the string alongside with the remaining part. If the regular expression can not be successfully matched against the string, the consumed part will have a value of false, indicating no successful match. Two special subclasses are always added to the list of classes, which are the EOLTokenClass — it matches the end-of-line characters — and the EOFTokenClass — which matches the end-of-file indicator. Another special subclass is the StatelessTokenClass, which associates the AnyState with itself, which means it will match in any non-strict state.

Lexer.TokenClassGroup

The TokenClassGroup is a container for TokenClasses. It can be used to define transitions if any of the classes in it is matched. It has an interface that consists of a push and pop method, a mapping function to transform all elements and a find function that checks if a TokenClass is part of the group or not.

Lexer.State

The State class represents a state in stateful mode — it has a name, and a strict parameter, which is false by default. Two special subclasses exist, one is the AnyState, which is used for tokens that shall match in any state that is non-strict, and other is the DefaultState, which is the one the state machine starts in, and which all classes belong to if no state parameter is given to them.

Lexer.StateTransitionBase

This class is the common core of the state transition classes — it has a from and to member, which represent the states to transition from and to, and it also has a method that takes the Lexer, and checks and retrieves the states for the two members (initially, the to and from members contain the name of the states only).

Lexer.StateTransition

The `StateTransition` class is the traditional state transition, which has a token class member alongside the inherited `to` and `from`. The instances can be created with two or three parameters — in the case of two parameters, the `from` state will be derived from the state associated with the token class, while if the class and the two states are provided, the `from` state will be explicitly set. This allows state transitions based on stateless token classes.

Lexer.GroupStateTransition

The `GroupStateTransition` is similar to the `StateTransition` class, but it defines a transition that happens when any of the group's token classes match in a certain state. The three parameters are the `TokenClassGroup`, the `from` and the `to` states.

2.2.2 Stateless mode

In stateless mode, `Lexer` only uses the default state, and does not transition into any other state — and thus the token classes do not need to be initialised with data about which state they belong to.

2.2.3 Stateful mode

In stateful mode, `Lexer` takes advantage of the rules defined by the `States`, `StateTransitions` and `GroupStateTransitions` given by the user. `Lexer` keeps track of the current state, and prioritizes group based transitions when deciding what state to transition to.

2.2.4 Skipping whitespaces

To make the lexical analysis easier, `Lexer` can be told to ignore whitespace characters. This is a very useful option, and thus the default value is `true` — but some languages require us to keep track of the indentation (e.g. Python), in which case the rules have to deal with the whitespace characters.

2.2.5 The Lexer class

The `Lexer` class is the core of this module. It provides facilities to add states, transitions, token classes, and has a method that runs the analysis on the given string. This method is called `tokenize`, and it follows an iterative algorithm — it looks for the first rule that matches the beginning of the string, creates a token based on that rule, does the necessary state transitions and starts the process all over again. The return value is a dictionary with three values:

- success - indicates whether the whole string was consumed or not
- tokens - the list of tokens the process generated
- rest - the rest of the string that was not parsed, if any

2.3 Possible enhancements

In the current state, the Lexer module is easily usable — however, the user needs to define first the states, then the token classes and groups and finally the transitions. This is a quite verbose way of defining the rules for the analysis — which could be done by using a custom data format that describes the rules. This data could be read from a file and preprocessed, which could make the process of defining the rules more robust. It could also be a more compact way of describing the rules.

Currently, the Lexer module can not transform the value of the tokens it creates. This could also make parsing easier, since e.g. integers and floats could be transformed into their numerical values during this stage, rather than having to deal with this process during the traversal of the AST.

Chapter 3

Syntactical analysis (Parsing)

Parsing (or syntactical analysis) is the process of turning the tokens that the lexical analysis produces into a parse tree by using a grammar. The lexical analysis phase only determines the types of things in the source code, while during parsing the higher level structures are revealed. In any human language, the lexical analysis would only reveal whether the word is a noun, verb, adverb etc., while syntactical analysis would turn these words into sentences, paragraphs, chapters and so on.

The input of the parser is the grammar and the tokens, while the output is a hierarchical structure that represents the structure of the source code (or the tokens, in this case).

3.1 Describing the grammar

There are well established ways to describe grammars — one of them is a notation called Backus-Naur form, or BNF (named after John Backus — who proposed the notation — and Peter Naur — who edited the first document that used it, a 1963 report on Algol60 [1]). This notation is widely used, and provides an easy way to describe context-free grammars. Context-free grammars are a subset of formal grammars which do not use context to determine which rule to be used in a given situation.

A grammar consists of basic production rules which match one non-terminal to zero or several terminals and non-terminals, for example:

$$\begin{aligned}\langle Sum \rangle &::= \langle Sum \rangle '+' \langle One \rangle \\ &| \langle One \rangle \\ \langle One \rangle &::= '1'\end{aligned}$$

This grammar describes a language of only two tokens (or terminals), '+' and '1', and describes two rules:

- A Sum can be a Sum plus ('+') a One, or a One on it's own

- A One is a '1' character

Or in short, an infinite amount of ones added to other ones is valid in a language using this grammar, and the sum of the ones can be easily calculated recursively after parsing. The '|' symbol means 'or', while other rules are referenced with their name written in '<>'s, and string literals are represented in quotes.

BNF has been extended over its existence, the two notable extensions is EBNF and ABNF, where the first characters stand for extended and augmented respectively. The parser implemented in this thesis uses the original BNF syntax, and thus the two extended versions are not described here.

3.2 The parsing algorithm

One of the big choices someone who sets out to write a parser shall make is which parsing algorithm to use. The algorithms can be divided into two groups, namely top-down parsing algorithms and bottom up parsing algorithms. These differ in the following ways:

- Top-down parsers try to fit a parsing tree to the input string by looking at the tokens left to right, and looking for any matching grammar rules. This means that the parse tree is built from the higher level structures, and the most basic rules get parsed the last — hence the name top-down parsing. Since two rules can have similar prefixes, in ambiguous cases all the possible choices are explored. The list of top-down parsing algorithms include:
 - Recursive-descent parser [9] – these parsers use recursive functions to parse the input. These parsers closely resemble the grammars they are generated from, but they tend to have performance problems because of the presence of recursion. They can be applied to any grammar, but the algorithm is not guaranteed to terminate.
 - LL(k) parser [10] – an LL(k) parser is a state machine based parser which parses the tokens from left to right and performs leftmost derivation (the first non-terminal from left to right is rewritten), with k lookahead. What this means is that during parsing, the decisions are based on not just the current token, but also the next k tokens. These parsers are generally easy to construct, so many computer languages are built to be LL(1) — which means the language's grammar can be parsed by an LL(1) parser. LL(1) grammars can also be parsed by a recursive-descent algorithm, and in that case the algorithm is guaranteed to terminate.
- Bottom-up parsers do parsing the other way around — they first find the most basic structures, and build the higher level structures from the already recognized lower ones. A different name for bottom-up parsers is shift-reduce parsers, since these

parsers uses these main actions during parsing. Some bottom-up parsing algorithms are:

- LR(k) parser [10] – similar to the LL(k) parser in the sense that it reads from left to right, it is state machine based and uses k lookahead, but these parsers perform rightmost derivation (usually in reverse). These parsers can handle the grammars they accept (unambiguous, context-free) very efficiently (in linear time).
- LALR parser – a variation of the LR(1) parser which is more memory efficient, since a smaller table is used by the state machine. The downside is that an LALR parser only accept a subset of LR(1) grammars. Nevertheless this algorithm is widely used by parser generator tools — e.g. in Yacc and GNU Bison.

During development I have tried two of these algorithms. One was a naive implementation that was essentially a recursive-descent-like algorithm, while the final version implements an LALR parser construction algorithm.

3.2.1 Naive implementation

My first idea to parse a source string was to build a parser following an intuitive, greedy algorithm. The basic idea was to read the grammar rules, and try to match the beginning of the string with one of the rules by recursively matching the right hand side of each rule. Surprisingly this approach worked quite well in some cases, but had some major flaws which led me to abandon this idea and look for an algorithm that was proven to be effective. The problems were the following:

- The parser was greedy, it accepted the rule that consumed the most of the input string, so rules that had the same prefix were not handled well by the algorithm.
- Left-recursive rules caused an infinitely deep recursion, since the parser always tried to match the same rule over and over again. This is something that is a common problem in the family of top-down parser, although an algorithm was constructed [4] that is top-down and can handle left-recursive grammars.

After realizing that the naive algorithm has severe limitations, I looked for a better alternative — one that was not too difficult to implement, but was also efficient. My first idea was to use an LL parser, but once I realized that this would not solve the left-recursion problem, I turned to the slightly more complicated LR parsers. After a bit of research, I have found a very easy-to-follow description of an LALR parser construction algorithm, and so I have based my work on this algorithm.

3.2.2 LALR - LookAheadLR

The invention of LALR algorithm is attributed to Frank DeRemer [3], but he did not actually create an algorithm to construct these kinds of parsers. Since then multiple algorithms have been created to construct LALR parsers. The algorithm I have used - and which is described by Stephen Jackson in his "A Tutorial Explaining LALR(1) Parsing" [5] is an algorithm that first builds a full LR(1) table, which is then reduced to an LALR table. This is not the most efficient way to construct an LALR parser, since merging of some rules only happens during a late stage, while merging could be done in one of the earlier stages too ¹, with the downside of the merging being slightly more complicated to do. I have decided to closely follow the tutorial, since memory during table construction should not be a huge problem (nowadays there is plenty enough memory in the use-case scenarios of this parser generator, we are not aiming to create a parser for an embedded system with limited memory), but some steps of the algorithm — namely the calculation of First and Follow sets — required additional research to be implemented effectively. The construction of an LALR parser by reducing an LR(1) is done in the following steps:

Adding a START rule to the grammar

This is more of a preparation step for the real algorithm - we need to have a specific starting rule for our grammar, which must only have one item on the right hand side. We could let the writer of the grammar deal with this convention, but it is easier to just create a starting rule for ourselves that has the first rule of the grammar on the right hand side. In my case, this rule gets the special name "#S" which is otherwise not a valid name, so this rule can not be changed later during the grammar definition by adding alternatives. This rule will be the new first rule of the grammar.

Finding the item sets

The first step is finding the item sets — or to be more precise the canonical collection of LR(0) items. These items represent partially parsed rules, with a \bullet after the last parsed item on the right hand side (or in the case of a rule that has no parsed elements, at the beginning). The following item represents a *Sum* rule that expects a '+' as the next parsed token:

$$\langle Sum \rangle ::= \langle Sum \rangle \bullet '+' \langle One \rangle$$

The item sets are constructed using the following algorithm:

1. The first set starts with the starting rule, with the \bullet at the beginning of the right hand side.

¹LALR parsing — Stanford University handout <http://dragonbook.stanford.edu/lecture-notes/Stanford-CS143/11-LALR-Parsing.pdf>

2. We expand the set by adding every rule (with the \bullet at the beginning) that maps from anything on the right hand side of the rules in the set. We do this until the set does not change any more.
3. We calculate what inputs the set is expecting (what items are after the \bullet in each rule) and we begin a new set using each of these inputs. We begin a new set by adding each item from the current set that expects the input we chose, but with the \bullet moved one position forward (after the input, meaning that we have moved on to the next input).
4. We expand each of the newly created sets by repeating step 2 on them, until we eventually expand all sets with all possible inputs.

Finding the translation table elements

Using the sets we have built in the previous step, we can find the translation table elements. We can treat the sets from now on as the states of the parser, and the translation table will tell us which state is next if we get a certain input. This can be easily derived from the previous step, since any time we create a new set, we know which input lead to that state.

Finding the extended grammar rules

The extended grammar consists of extended grammar rules - which are basically the original grammar rules with some added notation. Every terminal and non-terminal in the left and right hand side is augmented with some additional information about state transitions (a *from* and a *to set*), resulting in an extended grammar item following, and following this algorithm:

1. The left hand side's *from set* is the item set the rule belongs to, while the *to set* is read from the translation table of the set — specifically the set that we would transition to if the left hand side was given as input. In the case of the starting rule, this set does not exist, but this special case is avoided because the starting rule can not appear on the right hand side of any rule, since that is a special rule that we ourselves added.
2. The right hand side is calculated iteratively. The first *from* and *to set* of the first item is set similarly to the left hand side, while for the following items the *from set* will be the *to set* of the previous item — basically forming a chain of transitions. This will allow us later to determine which is the final set that the rule transitions to.

If we take every item from every item set with the \bullet at the beginning of the right hand side, and apply this procedure to those items, we get the full extended grammar.

Finding the First and Follow sets

The first sets are calculated for the extended grammar items with the following algorithm:

1. If x is a terminal, $\text{First}(x) = \{ x \}$
2. If V is not a terminal, then examine the rules that have V as the left hand side, and iterate over the right hand side using the following rules (ϵ represents the empty string):
 - (a) If the current item is a terminal, stop the iteration
 - (b) If the current item is not a terminal, add the First set of that item to the set of V (minus ϵ). If the current item's First set does not contain ϵ , stop the iteration, otherwise continue with the next item.
 - (c) If we are currently at the last item of the right hand side and the item's first set contains ϵ , add ϵ to the First set of V

This algorithm can be calculated iteratively — using these rules we mutate the First sets until we get to a pass over the sets that does not mutate the sets any more.

The final sets are calculated for the same items, but using the Firsts sets in the algorithm, which is the following:

1. Add the end-of-file token to the Follow set of the starting rule
2. If x is a terminal, $\text{Follow}(x) = \{ \}$
3. If V is not a terminal, then examine the rules that have V appear on the right hand side, and do one of the following (ϵ represents the empty string again):
 - (a) If V is not the last item on the right hand side, add everything from the First set of the next item (minus ϵ) to the Follow set of V . If the next item's First set contained ϵ , add everything from the left hand side's Follow set to the Follow set of V too
 - (b) If V is the last item on the right hand side, add everything from the left hand side's Follow set to the Follow set of V

The Follow sets can be calculated iteratively the same way we calculated the First sets.

Merging the rules of the extended grammar

We now can merge the extended grammar rules. This is the step that reduces the size of the Action/Goto table, and thus the required memory for the parser.

We can merge two rules if their original rule is the same (that is, they are created from the same original grammar rule, so stripping away the state transition informations would leave us with identical rules) and their final sets are also the same. When we merge two rules, the merged rule's Follow set will be the union of the Follow sets of the merged rules.

Creating the Action/Goto table

Now the core of the parser, the Action/Goto table can be constructed. Using this table, we can then parse any input by applying simple rules — the Actions. We do this by keeping track of the following: a stack of states, the input, and an array of parse tree nodes which is the output of the parser.

The Goto part of the table can be read from the translation table directly — every non-terminal transition is copied as a Goto. The Action part of the table can contain 4 types of actions:

- Accept – this means that the input has been parsed, and the algorithm is finished
- Error – this represents a syntax error — a token was given as input that was not expected. This action is the default if none specified directly
- Shift – pushes the state associated with it to the parser stack and removes the first token from the input
- Reduce – Removes as many items from the state stack as many items the right hand side of the associated rule contains. The same amount of nodes are removed from the nodes list, and a new node is added with the removed nodes as children, and the action's rule as the rule for the node. The next state is deduced from the Goto table — we use the state on top of the stack as a temporary state, we use the left hand side of the rule as an input to the Goto table, and we push the item in the Goto table to the stack.

We start out in state 0, which contains the starting rule. Then we can start parsing by examining the beginning of the input and reading the action we need to perform from the table. Doing this until we hit an Accept or Error action will result in either a syntax error, or a valid parse tree. This tree then can be transformed into an abstract syntax tree by removing unnecessary data — we merge nodes that only have one children with their child, we remove children which have terminal values originally — these contain no information —, and we remove empty strings from the children after reducing them too.

3.3 The Parser module

The Parser module is the one that implements the algorithms and the classes required. While constructing this module, I tried to keep things separated — I've ended up with one private method for the main class for each step of the parsing algorithm, and a parse method to parse code and return the AST.

3.3.1 The Parser subclasses

Parser.RuleTerminalBase

The base class for the following two classes. It defines the `isEpsilonRule` and `isTerminalRule` methods.

Parser.BNFRule

Represents a rule read from the BNF definition. Has a name, and a list of subrules (which corresponds to the right hand side of the rule in a two dimensional array, the alternatives are on the first level and the actual rules on the second). The token classes are represented using this class, in that case a `tokenClass` member is set, and the subrules list is empty.

Parser.BNFTerminal

This class represents a terminal value read from the BNF definition.

Parser.LR0Item

This class represents an item of the canonical collection of $LR(0)$ items. It has a rule, an index to indicate which subrule alternative to use, and a dot value to indicate which item the \bullet is before.

Parser.LR0ItemSet

It represents a set of LR0Items. It stores the items and the translation table, and has methods used during calculation of the item sets (getting a list of items after the \bullet symbol, creating new items for a new set, expanding the set).

Parser.ExtendedGrammarItem

It represents a terminal or non-terminal of the original grammar, but with the additional information about *from* and *to sets*, and also the First and Follow sets. A set of these items is maintained during parser generation.

Parser.ExtendedGrammarRule

This class is the wrapper for an extended grammar rule, and it keeps track of the left hand side and the right hand side of the original rule. It has methods to check if it is mergeable with another rule, and to get the final set of the rule.

Parser.Action

This is the base class for the actions in the Action/Goto table. For convenience reasons, the Goto class is also a subclass of this. All subclasses implement an execute method that modify the parser's state (the parser is passed as a parameter). It also has a member to store the input in (the token that triggers the action).

Parser.Accept, Parser.Shift, Parser.Reduce, Parser.Goto

These classes represent the items in the Action/Goto table, each implementing the algorithm associated with the action in their execute method.

Parser.Node

This is the class that represents a node in the parse tree and AST. Each node has a reduce method which does the conversion for the node and the node's subtree from a parse tree node into an AST node (removing unnecessary data and merging).

3.3.2 The Parser class

The Parser class is the class that encapsulates all the algorithms and does the actual parsing of the source tokens in the parse method. This is the class that is exported from the module.

3.4 Possible enhancements

The current implementation runs quite well, can handle LALR grammars with left-recursive rules and generate ASTs from those. A good way to improve the parser would be to implement different algorithms for parsing, which would ease the limitation of only accepting LALR grammars.

Chapter 4

The Language module and two examples

4.1 The Language module

The Language module consists of only one class — the Language class — which encapsulates the Lexer and Parser classes. It has two getters — one for the Lexer and one for the Parser — and a method to build and return the AST from the raw source code string. The module exports the Language class, which has the Lexer and Parser classes as subclasses too.

4.2 Examples

I have written two example language modules to test the core modules — one is a math expression parser, and the other is a very small subset of the Logo programming language — a for cycle, and the drawing statements are accepted, and of course the parameters can be math expressions.

4.2.1 The math parser

The way I imagined people could create language modules is by extending the Language class — and this is what I have done in this case. In the constructor I initialize the Lexer token classes, and the parser with the BNF grammar. I added two more methods — an executeOne and an execute method. The execute method builds the AST using the parser's parse method, and calls executeOne on the main node. The executeOne method then recursively executes the subtrees of the nodes — it checks for known types based on the names of the token classes, and parses the integers and floats into their numerical values. If the current node is not a terminal, then it first calls itself on that node to determine its value. Once all the numerical parameters are ready, it does the calculation, and returns

the resulting number. This is done using switch statements to handle each rule differently if necessary.

```
const Language = require('../lib/language')
const { Lexer, Parser } = Language

class BasicMath extends Language {
  constructor() {
    super()

    this.lexer.addTokenClasses([
      new Lexer.TokenClass('int', /[0-9]+(?![0-9]*\[0-9]+\+)/),
      new Lexer.TokenClass('float', /[0-9]+\.[0-9]+\+/),
      new Lexer.TokenClass('char', /\S/)
    ])

    this.parser.fromBNF(
      '<S> ::= <S> <SEP> <E> | ""'
      '<SEP> ::= ",", " " | <Token-EOL>'
      '<E> ::= <E> <PM> <T> | <T>'
      '<T> ::= <T> <MD> <H> | <H>'
      '<H> ::= <H> "^" <F> | <F>'
      '<PM> ::= "+" | "-"'
      '<MD> ::= "*" | "/"'
      '<F> ::= "(" <E> ")" | <Token-int> | <Token-float>'
    )
  }

  executeOne(node) {
    if (node.rule.class !== undefined) {
      switch (node.rule.class.name) {
        case 'int':
          return parseInt(node.rule.value)
          break;
        case 'float':
          return parseFloat(node.rule.value)
          break;
        default:
          return node.rule.value
      }
    }
  }

  let parameters = node.children.map(child => {
    return this.executeOne(child)
  })

  switch (node.rule.name) {
    case 'E': case 'T':
      switch (parameters[1]) {
        case '-':
          return parameters[0] - parameters[2]
          break
        case '+':
          return parameters[0] + parameters[2]
          break
        case '*':
          return parameters[0] * parameters[2]
          break
        case '/':
          return parameters[0] / parameters[2]
      }
    }
  }
```

```

        break
      default:
        throw Error('Unknown rule')
    }
    break
  case 'H':
    return Math.pow(parameters[0], parameters[1])
    break
  case 'F':
    return parameters[0]
    break
  case 'S':
    if (parameters.length === 1)
      return parameters[0]
    return parameters[0] + ', ' + parameters[parameters.length - 1]
    break
  default:
    throw Error('Unknown rule')
}
}

execute(code) {
  let ast = this.buildAST(code)
  return this.executeOne(ast[0])
}
}

module.exports = new BasicMath()

```

I have added some extra rules to verify that the parser can handle optional rules (rules with ϵ as an alternative) and left-recursion. By carefully constructing the grammar I have ensured that the operator precedence is as it should be (*nthpower* > ***/*>* *+/-*).

```

const BasicMath = require('./examples/math')
global.BasicMath = BasicMath
console.log(
  BasicMath.execute('10 ^ ((10 - 5 + 4) / (6 - 3)), 10 22/10
  432.432/10')
)

```

Running this piece of code results in the following: '1000, 10, 2.2, 43.2432', which is indeed the expected output.

4.2.2 The Logo subset parser

The way the Logo subset parser is constructed is similar to the math expression parser. I uses math expressions as the parameters for the drawing expressions, and it has a for cycle that runs the code given to it in a list 'n' times.

It uses the following grammar:

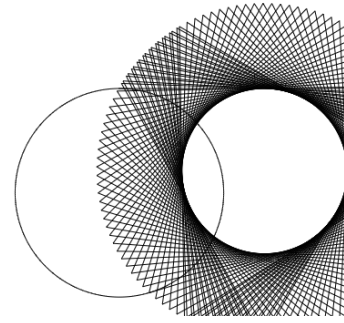
$$\begin{aligned}
\langle Program \rangle &::= \langle Program \rangle \langle Expression \rangle \\
&| \text{ ' ' } \\
\langle Expression \rangle &::= \langle Command \rangle \\
&| \langle For \rangle \\
\langle Command \rangle &::= \langle Command-name \rangle \langle Math \rangle \\
\langle Command-name \rangle &::= \text{ 'f' } \\
&| \text{ 'b' } \\
&| \text{ 'l' } \\
&| \text{ 'r' } \\
\langle For \rangle &::= \text{ 'for' } \langle Math \rangle \text{ '[' } \langle Program \rangle \text{ ']' } \\
\langle Math \rangle &::= \langle Math \rangle \langle PM \rangle \langle T \rangle \\
&| \langle T \rangle \\
\langle T \rangle &::= \langle T \rangle \langle MD \rangle \langle H \rangle \\
&| \langle H \rangle \\
\langle H \rangle &::= \langle H \rangle ^ \langle F \rangle \\
&| \langle F \rangle \\
\langle PM \rangle &::= \text{ '+' } \\
&| \text{ '-' } \\
\langle MD \rangle &::= \text{ '*' } \\
&| \text{ '/' } \\
\langle F \rangle &::= \text{ '(' } \langle Math \rangle \text{ ')' } \\
&| \langle Token-int \rangle \\
&| \langle Token-float \rangle
\end{aligned}$$


Figure 4.1: *Result of the Logo example*

The resulting module is then similarly exported into a global variable, which means that in the browser, it will be available as a property of the window object, while in Node.js it will become a global variable.

By providing a callback parameter, and implementing simple drawing on a HTML5 canvas in JavaScript, if we bundle the interpreter script and handle the necessary things in the browser, we can actually create drawings fairly easily.

```
logo.execute('for 360 [f 2 r 1] for 10^2+50 [f 321 r 121 for 2 [r 1 l 1]]')
```

This code generates the drawing on Figure 4.1

Chapter 5

Summary, conclusions

The original goal of this thesis was to create a usable parser module in JavaScript that is as easy to use as to understand. Since my background in formal languages was limited, I aimed to use an algorithm that I could completely understand. This is the reason behind my first solution — I tried to solve the problem using only intuition, to see what can be done without previous knowledge of the subject.

While this turned out to be somewhat of a waste of time, I did learn a lot from this failure. I had a better grasp at what has to be done, and I was eager to find out how a proper solution worked. This lead me to choose one of the state-machine based algorithms, and to eliminate the left-recursion problem, I ended up looking at LR parsers. I decided to use the LALR parser algorithm because of it being a bit more memory efficient, and not the easiest to implement — but not the hardest either.

While creating an implementation for the steps of LALR parser construction, I tried to use algorithms that are easy to understand, so that others can build upon and understand my solution later.

The two examples show how easy it is to run code in any custom language in the browser with this modular interpreter — and how easy it is to write modules for languages. This satisfies the initial goal of making the creation of educational websites easier. For example using this parser, HTML based presentations could be built for programming classes that have live demo codes on the slides that can be modified even during presentation. In the more traditional sense, simple programming education websites could use this interpreter to provide a simple programming environment to the students, keeping them close to the learning material.

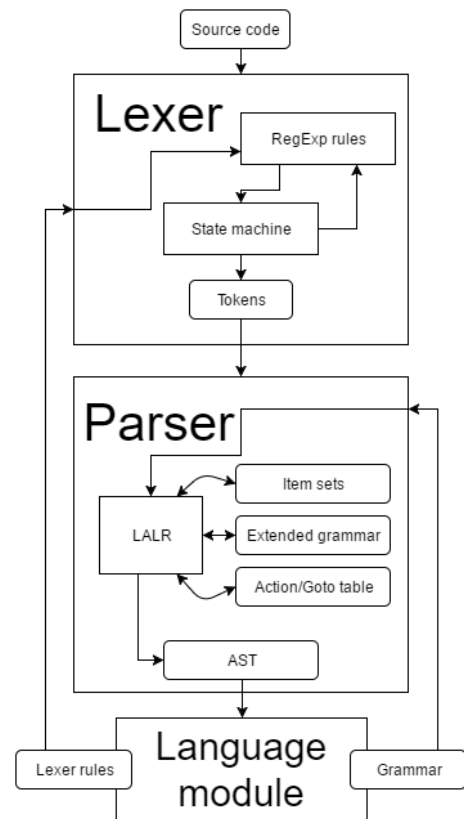


Figure 5.1: *High level overview of the interpreter*

The source code of the solution is available on GitHub in a public repository (<https://github.com/brenca/szakdolgozat>). This thesis is also available in the repository, since it serves as a detailed documentation for the solution.

Chapter 6

Problems, possible improvements

The implementation that is described in this paper works quite well, but the goal of development was not to create an optimized, efficient algorithm but to create something that is easy to understand and to use.

This means that while the examples show that the basic functions work properly, there might be ways to improve on the speed of the algorithms.

- The First and Follow sets are calculated iteratively. While this is a simple-to-understand algorithm, it is not the most efficient — in fact a lot of set unions are recalculated a lot of times, which makes this algorithm rather wasteful. In the future, better, more efficient algorithms [2] could be used to calculate these sets.
- The LALR parser generator algorithm first generates the items sets and the extended grammar for the LR(1) parser table, and only reduces the table size while calculating the Action/Goto table. This could be improved upon by merging during the calculation of the item sets, as is mentioned in Chapter 3.2.2.

Another problem is that the current implementation only accepts LALR grammars. While this is not a severe limitation as simple languages, or subsets of proper programming languages (or even full programming languages like Java [6]) can be implemented using LALR grammars, there are programming languages that have a non-LALR grammar, like C++ [11]. This problem could be solved by implementing other parsing algorithms, and providing the possibility for the module writer of choosing the algorithm they want to use.

Some improvements could be made in the area of error handling. Currently if an unexpected symbol is encountered by the parser, it throws a syntax error, that has information about which line and which token in the line caused the error. This however does not mean that an exact character position for the error can be determined, since tokens can have length larger than one character, and in the default case whitespaces are skipped, which means information about their length is lost. Improving this would make writing an web based programming editor for the interpreter easier.

Acknowledgements

I would like to thank my friend Blumenschein Volfram for helping me by checking my work for grammatical errors. I would also like to thank Stephen Jackson for his tutorial, which helped me to understand the basic idea of LALR parser construction algorithm - he tried to make the subject easy to understand, and I think he succeeded.

Bibliography

- [1] John W Backus, FL Bauer, J Green, C Katz, J McCarthy, P Naur, AJ Perlis, H Rutishauser, K Samelson, B Vauquois, et al. Revised report on the algorithmic language algol 60. *The Computer Journal*, 5(4):349–367, 1963.
- [2] Frank DeRemer and Thomas Pennello. Efficient computation of lalr (1) look-ahead sets. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 4(4):615–649, 1982.
- [3] Franklin L DeRemer. *Practical translators for LR (k) languages*. PhD thesis, MIT Cambridge, Mass., 1969.
- [4] Richard A Frost, Rahmatullah Hafiz, and Paul C Callaghan. Modular and efficient top-down parsing for ambiguous left-recursive grammars. In *Proceedings of the 10th International Conference on Parsing Technologies*, pages 109–120. Association for Computational Linguistics, 2007.
- [5] Stephen Jackson. *A Tutorial Explaining LALR(1) Parsing*, 2009. URL: <http://web.cs.dal.ca/~sjackson/lalr1.html>.
- [6] Bill Joy, Guy Steele, James Gosling, and Gilad Bracha. Java (tm) language specification. *Addison-Wesley*, June, 2000.
- [7] Stephen Cole Kleene. Representation of events in nerve nets and finite automata. Technical report, DTIC Document, 1951.
- [8] M. E. Lesk and E. Schmidt. *Lex - A Lexical Analyzer Generator*. Bell Laboratories, July 1975. URL: <http://hpd.c.syr.edu/~chapin/cis657/lex.pdf>.
- [9] Roman R Redziejowski. Parsing expression grammar as a primitive recursive-descent parser with backtracking. *Fundamenta Informaticae*, 79(3-4):513–524, 2007.
- [10] Seppo Sippu and Eljas Soisalon-Soininen. *Parsing Theory: Volume II LR (k) and LL (k) Parsing*, volume 20. Springer Science & Business Media, 2013.
- [11] Edward D Willink. Meta-compilation for c++. 2001.