

CHALMERS



A Generator of Divide-and-Conquer Lexers

A Tool to Generate an Incremental Lexer from a
Lexical Specification

Master of Science Thesis [in the Programme MPALG]

JONAS HUGO

KRISTOFER HANSSON

CHALMERS UNIVERSITY OF TECHNOLOGY
Department of Computer Science and Engineering
Göteborg, Sweden, July 2013

The Author grants to Chalmers University of Technology and University of Gothenburg the non-exclusive right to publish the Work electronically and in a non-commercial purpose make it accessible on the Internet. The Author warrants that he/she is the author to the Work, and warrants that the Work does not contain text, pictures or other material that violates copyright law.

The Author shall, when transferring the rights of the Work to a third party (for example a publisher or a company), acknowledge the third party about this agreement. If the Author has signed a copyright agreement with a third party regarding the Work, the Author warrants hereby that he/she has obtained any necessary permission from this third party to let Chalmers University of Technology and University of Gothenburg store the Work electronically and make it accessible on the Internet.

An not to long headline describing the content of the report
A Subtitle that can be Very Much Longer if Necessary
JONAS HUGO,
KRISTOFER HANSSON,

© JONAS HUGO, July 2013.
© KRISTOFER HANSSON, July 2013.

Examiner: NAME A. FAMILYNAME

Chalmers University of Technology
University of Gothenburg
Department of Computer Science and Engineering
SE-412 96 Göteborg
Sweden
Telephone + 46 (0)31-772 1000

Cover:
an explanatory caption for the (possible) cover picture
with page reference to detailed information in this essay.

Department of Computer Science and Engineering
Göteborg, Sweden July 2013

Abstract

Lorem ipsum dolor sit amet, consectetur adipisicing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

Acknowledgements

Lorem ipsum dolor sit amet, consectetur adipisicing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

The Authors, Location 11/9/11

Contents

1	Introduction	1
1.1	Background	1
1.2	Scope of work	1
2	Lexer	2
2.1	Lexing vs Parsing	2
2.2	Token Specification	2
2.2.1	Languages	2
2.2.2	Regular Expressions	3
2.2.3	Regular Definitions	3
2.3	Tokens, Patterns and Lexemes	4
2.4	Recognition of Tokens	4
2.4.1	Transition Diagrams	5
2.4.2	Finite Automata	5
3	Incremental Lexer	6
3.1	FingerTree	6
3.2	Devide and Conquer	6
3.3	Branch and Bound	6
4	Parallelism	7
5	Used Structures	8
6	Testing	9
7	Result	10
8	Performance Analysis	11
9	Discussion	12

10 Conclusion and Futher Work	13
Bibliography	14

1

Introduction

This master-thesis is carried out at Chalmers, on the department of computer science.

1.1 Background

Editors normally have regular-expression based parsers, which are efficient and robust, but lack in precision: they are unable to recognize complex structures. Parsers used in compilers are precise, but typically not robust: they fail to recover after an error. They are also not efficient for editing purposes, because they have to parse files from the beginning, even if the user makes incremental changes to the input. More modern IDEs use compilerstrength parsers, but they give delayed feedback to the user. Building a parser with good characteristics is challenging: no system offers such a combination of properties

1.2 Scope of work

*Usage of BNFC *With help of regexp build a finit state machine that will lex a code string. *Give finite states with corresponding Monoid data type. *Flag for errors from the Lexer, give meningfull info to the user, and stop the worklow after lexer, until new updated text. *If no errors, handel layout *Parse the Monoid data type tree, AKA integrate the result with an existing parser. *Smile and be happy!

2

Lexer

A Lexer, lexical analyser, is a program which job is to convert a string of a formal language into a sequence of tokens. #Hitta REF. This can be done by using regular expressions, regular sets and finite automata. Which are central concepts in formal language theory. [1]

2.1 Lexing vs Parsing

There are several reasons why a compiler should be separated in to lexical analyser and a parser (syntax analyser) phases. Simplicity of design is the most important reason. When dividing the task in to these to sub task, it allows the system to simplify one of these sub-tasks. For example, a parser that has to deal with white-spaces and comments as syntactical units would be more complex then one that can assume white-spaces and comments have already been removed by an lexer. Also when the two tasks are divided into sub-tasks it can lead to cleaner overall design when designing a new language [2] Also overall efficiency of the compiler can be improved. When separating the lexical analyser it allows for appliance of specialised techniques that serve only the lexical task. [2] Last compiler portability can be enhanced. That is Input-device-specific peculiarities can be restricted to the lexical analysis. [2]

2.2 Token Specification

This section will describe how to write rules for the tokens patterns.

2.2.1 Languages

An alphabet is an finite set of symbols, such an alphabet is for example the unicode, which includes approximately 100,000 characters. A language is any countable set of

strings over some fixed alphabet. [2]

#More text on languages!!!!

Like any formal language, a regular language is a set of strings. In other words a sequence of symbols, from a finite set of symbols. Only some formal languages are regular; in fact, regular languages are exactly those that can be defined by regular expressions. [3]

2.2.2 Regular Expressions

Say that we want to express the set of valid C identifiers. Use of regular expressions make it very easy, as shown in example 2.2.1.

Example 2.2.1 (Valid C Idents). Say we have a element $letter \in \{a \dots z\} \cup \{A \dots Z\} \cup \{-\}$ and another element $digit \in \{0 \dots 9\}$ Then with help of regular expressions the definition of all valid C identifiers would look like this: $letter(letter|digit)^*$. [2]

In definition 2.2.2 is the formal definition for regular expressions.

Definition 2.2.2 (Regular Expressions). [1]

1. The following characters are meta characters $\{'|', '(', ')', '*', '\}$.
2. A none meta character a is a regular expression that matches the string a .
3. If r_1 and r_2 are regular expressions then $(r_1|r_2)$ is a regular expression that matches any string that matches r_1 or r_2 .
4. If r_1 and r_2 are regular expressions. $(r_1)(r_2)$ is a regular expression of the form that matches the string xy iff x matches r_1 and y matches r_2 .
5. If r is a regular expression the r^* is a regular expression that matches any string of the form $x_1, x_2, \dots, x_n, n \geq 0$. Where r matches x_i for $1 \leq i \leq n$, in particular $(r)^*$ matches the empty string, ϵ .
6. If r is a regular expression, then (r) is a regular expression that matches the same string as r .

Many parentheses can be reduced by adopting the convention that the Kleene closure operator $*$ has the highest precedence, then concat and then or operator $|$. The two binary operators, concat and $|$ are left left-associative. [1]

2.2.3 Regular Definitions

In a definition of a language it is useful to give regular expressions names, so they can for example be used in other regular expressions, as these names where themself symbols. If Σ is an alphabet of basic symbols, then a regular definition is a sequence of definitions of the form:

$$\begin{array}{ccc}
d_1 & \rightarrow & r_1 \\
d_2 & \rightarrow & r_2 \\
\vdots & \rightarrow & \vdots \\
d_n & \rightarrow & r_n
\end{array}$$

where:

1. Each d_i is a new symbol, not in Σ and not the same as any other of the d 's.
2. Each r_i is a regular expression over the alphabet $\Sigma \cup \{d_1, d_2 \dots d_{i-1}\}$

By restricting r_i to Σ and previously defined d 's the regular definitions avoid recursive definitions. [2]

2.3 Tokens, Patterns and Lexemes

A lexical analyser uses three different terms. All which is described here below.

Token is a pair consisting of a token name and an optional attribute value. The token name is a abstract symbol corresponding to a lexical unit [2]. For example, a particular keyword, datatype or identifier. The token names is what is given to the parser.

Pattern is a description of what form a lexemes of a token may take. [2] For example, a keyword is just the sequence of characters that forms the keyword, an int is just a sequence consisting of just numbers.

Lexemes is a sequence of characters in the code that is being analysed which matches the pattern for a token and is identified by the lexical analyser as an instance of a token. [2]

As mention before a token consist of token name and a optional attribute value. This attribute is used when one lexeme can match more then one pattern. [2] For example the pattern for a digit token matches both 0 and 1, but it is important for the code generator to know which lexeme was found. Therefor the lexer often return not just the token but also an attribute value that describes the lexeme found in the source program corresponding to this token. [2]

2.4 Recognition of Tokens

In previous section the topic have been, how to represent a pattern using regular expressions and how these expressions relates to tokens. This section will highlight how to transform a sequence of characters into a sequence of abstract tokens. First some basic understanding with transition diagrams.

2.4.1 Transition Diagrams

A transition diagram is a graph consisting of nodes and edges. Each node represent a state which could occur during the process of scanning the input looking for lexeme that matches one of several patterns. Each edge is labelled with a symbol or set symbols, which tells which input must come next to be able to advance to this next state. Here follows some properties for a transition diagram, One state is said to be initial state. The transition diagram always begins at this state, before any input symbols have been read. Some states are said to be accepting (final). They indicate that a lexeme has been found. The found token should then be returned with any additional optional values, mentioned in previous section. [2]

2.4.2 Finite Automata

A finite automata are essentially graphs, like transitions diagrams, with some differences:

- Finite automata are recognizers; they simply say "YES" or "NO" about each possible input string.
- Finite automata comes in to different forms:

Nondeterministic Finite Automata (NFA) which have no restriction of the edges, several edges can be labelled by the same symbol out from the same state. ϵ , the empty string, is a possible label.

Deterministic Finite Automata (DFA) for each state and for each symbol of its input alphabet exactly one edge with that symbol leaving that state

Both these forms of finite automate are capable of recognising the same language, called regular languages. These are languages that regular expressions can describe. [2]

Nondeterministic Finite Automata

Deterministic Finite Automata

3

Incremental Lexer

Here we should talk about how a incremental lexer work. What specife techniques it uses, how this can be usefull in the real world and so on!

3.1 FingerTree

We should here also talk about some datastructers that is needed for a incremental lexer to work.

3.2 Devide and Conquer

#how dose ot work, how can we use it, what will it do for us?

3.3 Branch and Bound

#Maybe something about this???

4

Parallelism

Our solution should be able to run on several cores. This chapter should be about why how and so on.

5

Used Structures

#Unertain of the name of this chapter. But here we should talk about bnfc and alex. What we use from the differnet programs. How this is usefull is it becouse of lazynes or are the existing solutions good??

6

Testing

#How is the testing of this project preform!?

7

Result

#What have been the result?

8

Performance Analysis

#How fast is the lexer. How have we come to this conclusion?

9

Discussion

#Discuss discuss!!

10

Conclusion and Futher Work

#what will our Minions do???

Bibliography

- [1] A. V. Aho, Handbook of theoretical computer science (vol. a), MIT Press, Cambridge, MA, USA, 1990, Ch. Algorithms for finding patterns in strings, pp. 255–300.
URL <http://dl.acm.org/citation.cfm?id=114872.114877>
- [2] A. V. Aho, M. S. Lam, R. Sethi, J. D. Ullman, Compilers: Principles, Techniques, and Tools (2nd Edition), Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 2006.
- [3] A. Ranta, M. Forsberg, Implementing Programming Languages, College Publications, London, 2012, pp. 38–47.