

CHALMERS



A Generator of Incremental 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, January 2014

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.

A Generator of Incremental Divide-and-Conquer Lexers
A Tool to Generate an Incremental Lexer from a Lexical Specification
JONAS HUGO,
KRISTOFER HANSSON,

© JONAS HUGO, January 2014.
© KRISTOFER HANSSON, January 2014.

Examiner: BENGT NORDSTRÖM

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:
The image is a trol image, and will not be used in the final report.

Department of Computer Science and Engineering
Göteborg, Sweden January 2014

Abstract

A text that is a crash.course of the project. What will the report talk about, what obsticals had to be conquerd. talk talk talk.

Acknowledgements

We would like to take the chance of thanking our supervisor at department of computer science, Jean-Philippe Bernardy. Also thank our parents, and last but not least. We like to thank our self!

Jonas Hugo & Kristofer Hansson, Göteborg January 2014

Contents

1	Introduction	1
1.1	Background	1
1.2	Scope of work	2
2	Lexer	3
2.1	Lexing vs Parsing	3
2.2	Token Specification	4
2.2.1	Regular Expressions	4
2.2.2	Languages	4
2.2.3	Regular Definitions	5
2.3	Tokens, Patterns and Lexemes	5
2.4	Recognition of Tokens	7
2.4.1	Transition Diagrams	7
2.4.2	Longest Match	7
2.4.3	Finite Automata	8
3	Divide-and-Conquer Lexer	12
3.1	Divide and Conquer in General	12
3.1.1	The Basics	12
3.1.2	Associative Function	13
3.1.3	Time Complexity	13
3.1.4	Hand on Example	13
3.2	Divide and Conquer Lexing in General	15
3.2.1	Treestructure	15
3.2.2	Transition map	15
3.2.3	The Base Case	15
3.2.4	Conquer Step	16
3.2.5	Longest Match	16
3.2.6	Incremental Computing	17
3.2.7	Expected Time Complexity	17

3.3	Pitfalls	18
3.3.1	Bruteforce	18
3.3.2	Dont know what to call this!	20
3.4	Lexical Errors	22
4	Implementation	23
4.1	Alex	23
4.1.1	The DFA design	23
4.2	Monoid	23
4.2.1	The Basecase	24
4.2.2	The Conquer Step	24
4.3	Fingertree	24
4.3.1	Fundemental Concepts	24
4.4	Sequences	24
4.5	Transition Map	25
4.5.1	Array Format	25
4.5.2	Function Composition Format	25
5	Result	26
5.1	Preciseness	26
5.2	Performance	27
6	Discussion	29
7	Conclusion and Futher Work	30
	Bibliography	32
A	Java Lette Light	33

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

Dan Piponi has written a blogpost on how to determine in an incremental way if a string fulfills a regular expression. This is done by using Monoids, Fingertrees and tabulate functions. [7]

This blogpost is the fundamental idea behind the project. To build one same general idea, but instead build a tool that generates a lexical analyser given a bnf file specification of a language. Where the core algorithm in the tool follows the blogposts idea. The project will use Alex [5] as much as possible, that is this algorithm will be used as a wrapper to the Alex lexing tool.

The goal of the project is to create an algorithm that can do a lexical analysis on an update to an already lexed code with a sufficient fast time cost. With a sufficient fast time means that the lexical analyser can be run in real time.

The report will start by give a more general knowledge about lexical analysis. Then start to give a overviewing image of what is needed of the algorithm to work correctly. Which building blocks needed to create the algorithm. This will lead up to the implementation of the algorithm and specific requirements on the algorithm for it to be fully correct. The report will also describe how the testing has been done. That is test for correctness, robustness and efficiency. Also present the result for the benchmarking on different computer systems. The last part of the report will give a more formal performance of the algorithm, discussion of the result, some conclusions and further work.

2

Lexer

A lexer, lexical analyser, is a pattern matcher. Its job is to divide a text into a sequence of tokens (such as words, punctuation and symbols). The Lexer is a front end of a syntax analyser [9]. The syntax analyser in turn takes the tokens generated by the lexer and returns a set of expressions and statements. This can be done by using regular expressions, regular sets and finite automata, which are central concepts in formal language theory [1]. The rest of this chapter describes the concepts of the lexer in detail.

2.1 Lexing vs Parsing

Lexers usually work as a pass before parser; giving their result to the syntax analyser. There are several reasons why a compiler should be separated in to a lexical analyser and a parser (syntax analyser). First, simplicity of design is the most important reason. When dividing the task into these two sub tasks, it allows the programmer to simplify each of these sub-tasks. For example, a parser that has to deal with white-spaces and comments would be more complex than one that can assume these have already been removed by an lexer. Also when the two tasks have been seperated into sub-tasks it can lead to cleaner overall design when designing a new language. The only thing the syntax analyser will see is the output from the lexer, tokens and lexemes. The lexer also usually skips comments and white-spaces, since these are not relevant for the syntax analyser. Second, overall efficiency of the compiler can be improved. When separating the lexical analyser it allows for use of specialised techniques that can be used only in the lexical task. Third and last, compiler portability can be enhanced. That is Input-device-specific peculiarities can be restricted to the lexical analysis [2]. Therefore the lexer can detect syntactical errors in tokens, such as ill-formed floating-points literals, and report these

errors to the user [9]. Finding these errors allows the compiler to break the compilation before running the syntax analyser, thereby saving computing time.

2.2 Token Specification

The job of lexical analyser is to translate a human readable text to an abstract computer-readable list of tokens. There are different techniques a lexer can use when finding the abstract tokens representing a text. This section describes the techniques used when writing rules for the tokens patterns.

2.2.1 Regular Expressions

Example 2.2.1 (Valid C Idents [2]). Using regular expressions to express a set of valid C identifiers is easy. given an element $letter \in \{a \dots z\} \cup \{A \dots Z\} \cup \{-\}$ and another element $digit \in \{0 \dots 9\}$ Then using a regular expression, the definition of all valid C identifiers could look like this: $letter(letter|digit)^*$.

Definition 2.2.2 (Regular Expressions [1]).

1. The following characters are meta characters $meta = \{ '|', '(', ')', '*', '\}$.
2. A character $a \notin meta$ 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 that matches the string xy iff x matches r_1 and y matches r_2 .
5. If r is a regular expression 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 omitted by adopting the convention that the *Kleene closure* operator $*$ has the highest precedence, the *concat* operator $(r_1)(r_2)$ the second highest and last the *or* operator $|$. The two binary operators, *concat* and *or* are left-associative.

2.2.2 Languages

An alphabet is a finite set of symbols, for example Unicode, which includes approximately 100,000 characters. A language is any countable set of some fixed alphabet [2]. The term

”formal languages” refers to languages which can be described by a body of systematic rules. There is a subset of languages to formal languages called regular language, these regular languages refers to those languages that can be defined by regular expressions [8].

2.2.3 Regular Definitions

When defining a language it is useful to give the regular expressions names, so they can for example be used in other regular expressions. These names for the regular expressions are themselves symbols. If Σ is an alphabet of basic symbols, then a regular definition is a sequence of definitions of the form:

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

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

When rules have been defined for a language, the lexer needs structures to represent rules and the result from lexing the code-string. This section describe the structures which the lexical analyser use for representing the abstract data; what these structures are for and what is forwarded to the syntactical analyser.

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

Token is a pair consisting of a token name and an optional attribute value. The token name is an abstract symbol corresponding to a lexical unit [2]. For example, a particular keyword, data-type or identifier.

Pattern is a description of what form a lexeme 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. Can be described by a regular expression.

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

As mentioned before a token consists of a token name and an optional attribute value. This attribute is used when one pattern can match more than one lexeme. 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. Therefore the lexer often returns not just the token but also an attribute value that describes the lexeme found in the source program corresponding to this token [2].

A lexer collects chars into logical groups and assigns internal codes to these groups according to their structure, where the groups of chars are lexemes and the internal codes are tokens [9]. In some cases it is not relevant to return a token for a pattern, in these cases the token and lexeme is simply discarded and the lexer continues, typical examples of this is whitespaces and comments which have no impact on the code [2]. An example follows how a small piece of code would be divided.

Example 2.3.1 (Logical grouping [9]).

Consider the following text; to be lexed:

```
result = oldsum - value /100;
```

Given the regular language defined in appendix A, the lexical analyser would use the following rules:

Identifier \rightarrow *letter ident**

Integer \rightarrow *digit+*

Reserved \rightarrow $\backslash (| \backslash) | \backslash \{ | \backslash \} | ; | = | \backslash + \backslash + | < | \backslash + | - | \backslash * |$ Reserved characters

In order to produce the following tokens.

<u>Token</u>	<u>Lexeme</u>
Identifier	result
Reserved	=
Identifier	oldsum
Reserved	—
Identifier	value
Reserved	/
Integer	100
Reserved	;

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 giving some basic understanding with transition diagrams.

2.4.1 Transition Diagrams

A state transition diagram, or just transition diagram is a directed graph, where the nodes are labelled with state names. Each node represents a state which could occur during the process of scanning the input, looking for a lexeme that matches one of several patterns [2]. The edges are labelled with the input characters that causes transitions among the states. An edge may also contain actions that the lexer must perform when the transition is a token [9].

Some properties of transition diagrams follow. 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. If the found token is the longest match (see section 2.4.2) then the token will be returned with any additional optional values, mentioned in previous section, and the transition is reset to the initial state [2].

2.4.2 Longest Match

If there are multiple feasible solutions when performing the lexical analysis, the lexer will return the token that is the longest. To manage this the lexer will continue in the transition diagram if there are any legal edges leading out of the current state, even if it is an accepting state.[2].

The above rule is not always enough since the lexer has to explore all legal edges, even if the current state is accepting. If the lexer is in a state that is not accepting and don't have any legal edge out of that state, the lexer can't return a token. To solve this the lexer has to keep track of what the latest accepting state was. When the lexer reaches a state with no legal edge out of it, the lexer returns the token corresponding to the last accepting state. The tail of the string, the part that wasn't in the returned token, is then lexed from the initial state as part of a new token.[2]

Example 2.4.1 (Longest Match). Consider the following text; to be lexed.

```
/* result = oldsum - value /100;
```

Although this tex is not legal code, there is no lexical errors in it. Since the text starts with a multi line comment sign the lexer will try to lex it as a comment. When the

lexer encounters the end of the text it will return the token corresponding to the last accepting state and begin lexing the rest from the initial state. Below follows the rules relevant to this example followed by the result to the lexer, the rest of the rules can be found in appendix A.

The Rules:

	$/\backslash^*) ([\backslash^*] (\backslash^*) [\backslash^*])^* (\backslash^*)+ /$	Multi line comment
Identifier	\rightarrow <i>letter ident*</i>	
Integer	\rightarrow <i>digit+</i>	
Reserved	\rightarrow $\backslash (\backslash) \backslash \{ \backslash \} ; = \backslash + \backslash + < \backslash + - \backslash ^* $	Reserved characters

The result:

<u>Token</u>	<u>Lexeme</u>
Reserved	/
Reserved	*
Identifier	result
Reserved	=
Identifier	oldsum
Reserved	—
Identifier	value
Reserved	/
Integer	100
Reserved	;

2.4.3 Finite Automata

Transition diagrams of the formed used in lexers are representations of a class of mathematical machines called finite automata. Finite automata can be designed to recognise members of a class of languages called regular languages, mentioned above [9]. A finite automaton is essentially a graph, 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 two different forms:

Non-deterministic Finite Automata (NFA) which have no restriction of the edges, several edges can be labelled by the same symbol out from the same state. Further 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. The empty string ϵ is not a valid label.

Both these forms of finite automata are capable of recognising the same subset of languages, all regular languages [2].

Non-deterministic Finite Automata

An NFA accepts the input; x if and only if there is a path in the transition diagram from the start state to one of the accepting states, such that the symbols along the way spells out x [2]. The formal definition of a non-deterministic finite automaton follows:

Definition 2.4.2 (Non-deterministic Finite Automata [10]). A finite automata is a 5-tuple $(Q, \Sigma, \delta, q_0, F)$, where

1. Q is a finite set called the states,
2. Σ is a finite set called alphabet,
3. $\delta : Q \times \Sigma \rightarrow P(Q)$ is a transition function,
4. $q_0 \in Q$ is the start state, and
5. $F \subseteq Q$ is the accept state.

The transition function doesn't map to one particular state from a state and element tuple. This is because one state may have more than one edge per element. An example of this can be seen in example 2.4.3.

There are two different ways of representing an NFA which this report will describe. One is by transition diagrams, where the regular expression will be represented by a graph structure. Another is by transitions table, where the regular expression will be converted in to a table of states and the transitions for these states given the input. The following examples shows how the transition diagram and transition table representation will look like for a given regular expression.

Example 2.4.3 (RegExp to Transition Diagram [2]). Given this regular expression: $(a|b)^*abb$ the transition diagram in fig. 2.1 representing this regular expression.

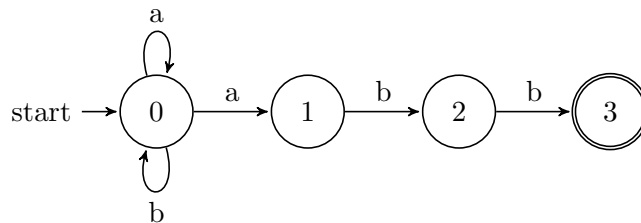


Figure 2.1: Transition Diagram, accepting the pattern $(a|b)^*abb$

Example 2.4.4 (RegExp to Transition Table [2]). Given the regular expression from example 2.4.3 it can be converted into the transition table shown in fig. 2.2

State	a	b	ϵ
0	$\{0, 1\}$	$\{0\}$	\emptyset
1	\emptyset	$\{2\}$	\emptyset
2	\emptyset	$\{3\}$	\emptyset
3	\emptyset	\emptyset	\emptyset

Figure 2.2: Transition Table representation of regular expression in example 2.4.3

Transition tables have the advantage that they have a quick lookup time. But instead it will take a lot of data space, when the alphabet is large. Most states do not have any move on most of the input symbols [2].

Deterministic Finite Automata

DFA is a special case of an NFA where,

1. there are no moves on input ϵ and
2. for each state s and input symbol a , there is exactly one edge out of s labelled with a .

While a NFA is an abstract representation of an algorithm to recognise the string of a language, the DFA is a simple concrete algorithm for recognising strings. Every regular expression can be converted in to a NFA and every NFA can be converted in to a DFA and then converted back to a regular expression [2]. It is the DFA that is implemented and used when building lexical analysers. The formal definition of a deterministic finite automaton follows:

Definition 2.4.5 (Deterministic Finite Automata [10]). A finite automata is a 5-tuple $(Q, \Sigma, \delta, q_0, F)$, where

1. Q is a finite set called the states,
2. Σ is a finite set called alphabet,
3. $\delta : Q \times \Sigma \rightarrow Q$ is a transition function,
4. $q_0 \in Q$ is the start state, and
5. $F \subseteq Q$ is the set of accept states.

Example 2.4.6 (DFA representation of RegExp [2]). A DFA representation of same regular expression from example 2.4.3 is shown in fig. 2.3

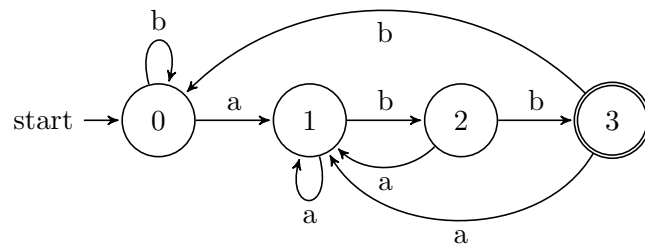


Figure 2.3: DFA, accepting the regular expression: $(a|b)^*abb$

3

Divide-and-Conquer Lexer

An incremental divide and conquer lexer works by dividing the sequence, to be lexically analysed, into smallest part and analyse them; and then combining them. In the base case the lexical analysis is done on a single character. The conquer step is then to combine the smaller tokens into as large tokens as possible. The end result should be a sequence of token that represent the code. How this is done will be described below.

3.1 Divide and Conquer in General

This section will try to give a idea of how the Divide and Conquer algorithm works in a general case. So when talking about the incremental lexer there is an understanding about the underlying techniques.

3.1.1 The Basics

The general idea of a divide and conquer is to divide a problem into to smaler parts and solve them separatly and then combine the results. A Divide and Conquer algorithm always consist of a pattern with these three steps [6].

Divide: If the input size is bigger than the basecase then divide the input into subproblems. Otherwise solve the problem using a straightforward method.

Recur: Recursively solve the subproblems associated with the subset.

Conquer: Given the solutions to the subproblem, merge them into to the original problem.

3.1.2 Associative Function

The associative property stated for when an expression containing two or more of the same operator, or function, in a row. The order of which of the operations that are executed first has no impact on the result.

In divide and conquer this is essential. On the conquer step in the divide and conquer algorithm, there is no certain order of which how the subproblems are going to be conquered. So the result of running any part with another part must give the same result as doing it in any other way.

3.1.3 Time Complexity

To calculate the running time of an divide and conquer algorithm the master method can be applied [4]. This method is based on the following theorem.

Theorem 3.1.1 (Master Theorem [4]).

Assume a function T_n constrained by the recurrence

$$T_n = aT_{\frac{n}{b}} + f(n)$$

This is typically the equation of a divide-and-conquer algorithm, where a is the number of sub-problems at each recursive step, n/b is the size of each sub-problem, and $f(n)$ is the running time of dividing up the problem space into a parts, and combining the sub-results together.

If we let $e = \log_b a$, then

1. $T_n = \Theta(n^e)$ if $f(n) = O(n^{e-\epsilon})$ and $\epsilon > 0$
2. $T_n = \Theta(n^e \log n)$ if $f(n) = \Theta(n^e)$
3. $T_n = \Theta(f(n))$ if $f(n) = \Omega(n^{e+\epsilon})$ and $\epsilon > 0$ and $a \cdot f(n/b) \leq c \cdot f(n)$ where $c < 1$ and all sufficiently large n

■

3.1.4 Hand on Example

The divide and conquer pattern can be preformed on different sorts of algorithm that solves different problems. A general problem is sorting, or more precies sorting a list of integers. This example will show merge-sort.

1. The algorithm starts with the divide step. Given the input S the algorithm will check if S consist of a number of elements less or equal to 1.
 - If this is true, return the sequence. A sequence of one or zero elements is always sorted.

- If this is false, split the sequence into two (as close as possible) equally big sequences. S_1 and S_2 , S_1 is equal $S.\text{getRange}(1, n/2)$. S_2 is equal $S.\text{getRange}((n/2)+1, n)$. The indexing of the sequence starts at 1.
- 2. The next step is to sort the subsequences S_1 and S_2 . This is done by recursively call it self with the subsequences as input. There will be 2 calls, one with S_1 as input and one with with S_2 as input.
- 3. The to sequences S_1 and S_2 is has now be recursively sorted and can be merged. Assign S to the result of the merging.

Algoritihm 1 shows a more formal definition of merge-sort.

Algorithm 1: MergeSort

Data: Sequence of integers S containing n integers

Result: Sorted sequence S

```

1 if  $\text{length}(S) \leq 1$  then
2   return  $S$ 
3 else
4    $(S_1, S_2) \leftarrow \text{splitAt}(S, n/2)$ 
5    $S_1 \leftarrow \text{MergeSort}(S_1)$ 
6    $S_2 \leftarrow \text{MergeSort}(S_2)$ 
7    $S \leftarrow \text{Merge}(S_1, S_2)$ 
8   return  $S$ 

```

Given the mergesort algorithm, time complexity is calculated this way. There are 2 recursive calls, the subproblems are 1/2 of the original problem size. To merge the two sorted subproblems the worst case is to check every element in the two list, $2 \cdot n/2 = n$.

$$T(n) = 2T(n/2) + n.$$

Results in:

$$a = 2, b = 2 \text{ and } f(n) = n \Rightarrow n^{\log_b a} = n^{\log_2 2} = n$$

Case 2 of the master theorem applies, since

$$f(n) = O(n)$$

So the solution will be:

$$T(n) = \Theta(n^{\log_2 2} \cdot \log n) = \Theta(n \cdot \log n)$$

3.2 Divide and Conquer Lexing in General

This section will give a general view of the techniques for an incremental divide and conquer lexer. Also an overview why it is a good idea using these techniques, before going into deep details about the incremental lexer this report relates to.

3.2.1 Treestructure

The incremental divide and conquer lexer should use a structure where the code-lexemes can be related to its tokens, current result can be saved and easily recalculated. A divide and conquer lexer should therefore use a tree structure to save the lexed result in. Since every problem can be divided into several subproblems, until the basecase is reached. This is clearly a tree structure of solutions, where the leaves are the tokens for a single character. and the root is a sequence of all tokens in the code.

3.2.2 Transition map

To be able to make a divide-and-conquer lexer work you can use a transition map for the subexpression that has currently been lexed. That is you store for each subexpression a list of tuples which are built up of an in state, a list of tokens and an out state. This could in Haskell look something like this.

```
type Transition = (State,[Tokens],Suffix,State)
type Transitions = [Transition]
```

The States

List of Tokens

The Suffix

3.2.3 The Base Case

When the lexer tries to lex one character it will create a transition table using the DFA for the language. It will for each state in the DFA lookup what out state the character would yield and create a Tokens type of Single. For the character '/' part of a transition map might look like the following.

$$\begin{bmatrix} 10 & \text{Single}'/' & 10 \\ 11 & \text{Single}'/' & \text{NoState} \\ 12 & \text{Single}'/' & 10 \end{bmatrix}$$

The reason we keep track of *NoState* will be evident once we show how the conquer step work.

3.2.4 Conquer Step

The most general case is a naive lexer that takes the first accepting state it can find. When two lists of tokens are combined there are two different outcomes:

Concat: create two tokens if the out state of the first list is accepting then the second list will be appended to the first list.

Combine: The other case is when the first list of tokens does not end in an accepting state. In this case the lexer will try to find an in state in the second list that is the same as the out state of the first transition.

$$\begin{bmatrix} 0 & \text{Single}'/' & 1 \\ 1 & \text{Single}'/' & \text{Accepting5} \end{bmatrix} \text{'combineTokens'} \begin{bmatrix} 0 & \text{Single}'/' & 1 \\ 1 & \text{Single}'/' & \text{Accepting5} \end{bmatrix} = \begin{bmatrix} 0 & \text{Single}'/' & \text{Accepting5} \\ 1 & \text{Multiple}'/'[]'/' & \text{Accepting1} \end{bmatrix}$$

This won't work as a lexer for most languages since it will lex a variable to variables where the length is a single character, for example "os" will be lexed as two tokens, "o" and "s". To solve this some more work is needed to be done.

3.2.5 Longest Match

Instead of taking the naive aproach where a token is created if the lexer finds an accepting state, the rule for creating a new token will instead be when the combination of two transitions yields *NoState* the lists will be appended. That is, when there is an out state from the first transition that corresponds to an in state of the second transition and the out state of the second transition isn't *NoState*, the last token of the first transition and the first token of the second transition will become one token, otherwise append the second list to the first list.

$$\begin{bmatrix} 0 & \text{Single}'/' & \text{Accepting5} \\ 1 & \text{Multiple}'/'[]'/' & 1 \end{bmatrix} \text{'combineTokens'} \begin{bmatrix} 0 & \text{Single}'\backslash n' & \text{Accepting6} \\ 1 & \text{Single}'\backslash n' & 1 \\ 5 & \text{Single}'\backslash n' & \text{NoState} \end{bmatrix} =$$

$$\begin{bmatrix} 0 & \text{Multiple}'//'\backslash n' & \text{Accepting6} \\ 1 & \text{Multiple}'//'\backslash n' & 1 \end{bmatrix}$$

The second case is when the out state for the right token list is *NoState*. This means that the two lists of tokens can't be combined. In this case the first token in the second list will be viewed as the start of a token and the last token in the first list will be viewed as the end of a token.

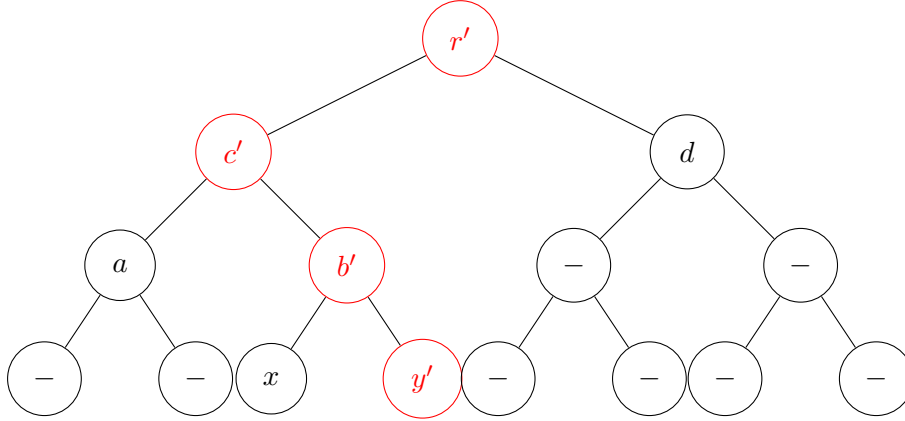


Figure 3.1: Incremental Computing, the updated nodes when a leaf changes

3.2.6 Incremental Computing

To be incremental means that, whenever some part of the data to the algorithm changes the algorithm tries to save time by only recomputing the changed data and the parts that depend on this changed data. [11]

For a divide and conquer lexer this would mean only recompute the changed token and the token to the right of the changed token. This is done recursively until the root of the tree is reached. The expected result of this would be that when a character is added to the code of 1024 tokens, instead of relex all 1024 tokens the lexer will only do 10 recalculations for new tokens. Since, $\log_2 1024 = 10$. This can be explained by the theorem 3.1.1. How this is calculated for an incremental divide and conquer lexer is described more in detail in the next sub-section.

3.2.7 Expected Time Complexity

Since incremental computing stated that only content which depends on the new data will be recalculated. That is, follow the branch of the tree from the new leaf to the root and recalculated every node on this path. As shown by fig. 3.1. Only one subproblem is updated in every level of the tree. Back to the master theorem. Let put this in to

numbers, $e = \log_b a$ where a is number of recursive calls and n/b is size of the subproblem where n is the size of the original problem. As shown by the fig. 3.1 number of needed update calls is 1, therefor $a = 1$. The constant b is still 2. This will give $e = \log_2 1 = 0$. Thus the update function of the incremental algorithm will have a time complexity of $\Theta(n^0 \cdot \log n) = \Theta(\log n)$

3.3 Pitfalls

This section will describe techniques that were tried under the constuction of the incremental divide and conquer lexer but was shown to give bad results.

3.3.1 Bruteforce

The first naive solution was to "bruteforce" to find the lex. This was showned to be to resource-eating. But it describes the general idea of how the problem could be solved. Why it was a bad solution will be described futher on in the text.

The above rules will work for very simple languages. When comments are introduced you will get the problem that the whole code can be one long partial comment token. To remedy this you can add two rules:

- Every time you combine two tokens you only do so if the combination has a transition from the starting state.
- If two tokens can be combined completly, check if the next token can be combined aswell.

This ensures that every token starts in the starting state and that each token is as long as it can be.

This also has some problems though. When keywords like "else if" are introduced the lexer will start to lex like in example 3.3.1. To solve this the lexer checks when two tokens are completely uncombinable if the first of these have an accepting state as outgoing state. If the token don't have an accepting out state, the lexer tries to break up the token until it does. The exception to this rule is single characters which are permitted to not have no accepting out states.

Example 3.3.1 (else if lexing). Somewhere in the middle of the code "... 1 else 0 ..."

String	Type
1	<i>Number</i>
-	<i>Space</i>
else_	<i>Nothing</i>
0	<i>Number</i>

Example 3.3.2 (Devide and Append). The lexer will always try to build as lage tokens as possible. When it realizes that this cant be done it has to backup and try to combine the parts in a different way. This example will show how this is done in theory.

The code segment for this example is:

"else return".

The tree in fig. 3.2 shows the first step of the token combine rutine. Clearly this returns a nonexsisting token. From here when the lexer has found that there are no tokens for

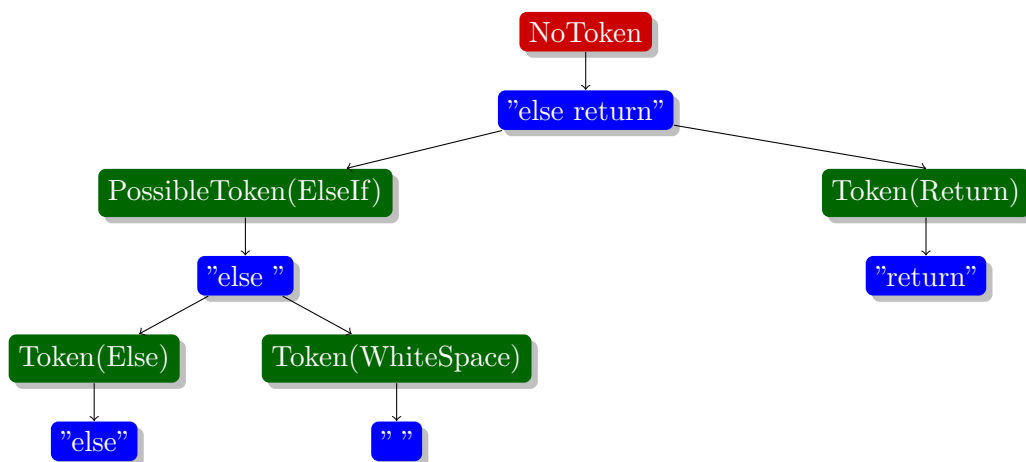


Figure 3.2: Lexer thinks "else " is an "else if" pattern.

this lexeme it will try to split the left child token.

split("else ") => ["else"," "]

Now the lexer has a pair of two lexems that represent valid tokens. The lexer knows that combining these two lexems in the pair returns in a NoToken result. So The only thing to do is to try to combine the right token in the pair with the right child token and let the token to the left in the pair stand alone. This also return a NoToken. So the

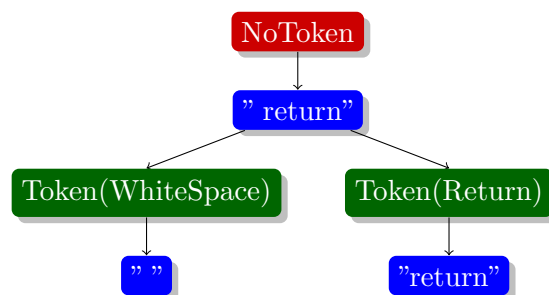


Figure 3.3: Lexer tries to combine an white space with a return statement

same thing will be done again. The lexer tries to split the left child before NoToken was

given. In this case the whitespace.

`split(" ") => []`

But because the whitespace is of the lowest form and is not build up by smaller tokens the resulting list from the split function will be empty. Now the lexer knows that this token must be by it self. The "return" is the last lexeme in this example code so the lexer can't combine it futher. Thus the lexer has found the resulting sequence of tokens:

`[(Token(Else), "else"), (Token(WhiteSpace), " "), (Token(Return), "return")]`

3.3.2 Dont know what to call this!

When the code is divided the lexer doesn't know if the string (or character) it lexes is the first, last or is somewhere in the middle of a token. Instead of checking what type of token the string will be (if it were to begin from the starting state) it saves all the possible state transitions for that string.

In the examples that follow below state 0 is considered the starting state and state 1 – 6 are considered accepting.

Example 3.3.3 (Transition map for a token). A hypothetical transition map for the char 'i'.

'i'	
<i>in</i>	<i>out</i>
0	1
1	1
8	7

In the base case the lexer will map all the transitions for all individual characters in the code and construct partial tokens of them. The conquer step will then combine two of these at a time by checking which possible outgoing states from the first token can be matched with incoming states from the second token. If there are such pairs of outgoing states with incomming states, then a new partial token is created.

Example 3.3.4 (Combining two tokens). 'if' can be an ident (state 1) or part of 'else if' (state 5).

'i'		'f'		'if'	
<i>in</i>	<i>out</i>	<i>in</i>	<i>out</i>	<i>in</i>	<i>out</i>
0	1	0	1	0	1
1	1	1	1	8	5
8	7	7	5		

[illegible]

However the remainder may not have the start state as a possible incoming state. In this case the lexer tries to find the largest possible token (that has the starting state as incoming state) and tries to construct a token of the rest of the remainder, repeating this procedure until the entire remainder has been split into acceptable tokens. All the tokens except the one that is on the very end of the sequence will have all but their accepting states stripped. This case does occur quite frequently since most languages have comments and strings which can contain anything.

		$'\text{--}'$		$'i'$	
$'i'$		<i>in</i>	<i>out</i>	<i>in</i>	<i>out</i>
<i>in</i>	<i>out</i>	= 0	2	<i>'combineToken'</i>	0 1
9	7	2	2		1 1
		9	8		8 7

$$checkRemainder \begin{pmatrix} ' _i ' \\ in & out \\ 9 & 7 \end{pmatrix} = \begin{pmatrix} ' _ ' \\ in & out \\ 0 & 2 \end{pmatrix} ++ \begin{pmatrix} ' i ' \\ in & out \\ 0 & 1 \end{pmatrix}$$

When all partial tokens has been combined in this way the resulting sequence of tokens represents the the code the lexer was run on.

3.4 Lexical Errors

Since the lexer has to be able to handle any kind of possible not "complete" tokens, error handling can be done in different ways. One approach is to simply return as many tokens as possible from the code and where there might be lexical errors the lexer returns the error in as small parts as possible.

Example 3.4.1 (A lexer that only lexes letters). When the lexer encounters the string "what @ day" it would return:

String	Type
What	<i>Word</i>
' '	<i>Space</i>
'@'	<i>No-Token</i>
' '	<i>Space</i>
day	<i>Word</i>

4

Implementation

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??

4.1 Alex

Alex is a tool for generating lexical analysers built in Haskell, given a description of the language in form of regular expressions. The result will be Haskell 98 compatible and can easily be used with the parser Happy, a parser generator for Haskell.

4.1.1 The DFA design

4.2 Monoid

In abstract algebra a monoid is a set, S , and a binary operation \bullet which furfills the following three properties:

Closure $\forall a, b \in S : a \bullet b \in S$

Associativity $\forall a, b, c \in S : (a \bullet b) \bullet c = a \bullet (b \bullet c)$

Identity element $\exists e \in S : \forall a \in S : e \bullet a = a \bullet e = a$

4.2.1 The Basecase

4.2.2 The Conquer Step

4.3 Fingertree

The incremental lexer uses a tree structure to save already lexed content. This tree is of form Fingertree. What that is and the basics of it is described in this section.

4.3.1 Fundamental Concepts

Before describing the functionality, lets take a look on which buildingblocks the fingertrees uses.

First, the fingertree uses monoids which has been described earlier in this chapter.

Second, the fingertree uses Right and Left Reductions. This is a function which collapses a structure of $f\ a$ into a single value of type a . The basecase for when the tree is empty is replaced with a constant value, such as \emptyset . intermediate results are combined using a binary operation, like the monoids \bullet . Reduction with a monoid always return the same value, independent of the argument nesting. But for a reduction with an arbitrary constant and binary operation there must be a specified nesting rule. If combining operation are only nested to the right, or to the left, the obtained result will be a skewed reductions, which can be singled out as a type class. [?]

4.4 Sequences

A sequence in Haskell is a form of sophisticated list. That is a list with better performance than the basic `[]` list notation. Where a list in Haskell has $\Theta(n)$ for finding, inserting or deleting elements, that is in a list there is only known current element and the rest of the list. Which will result in for example finding the last element of a list, the computer must look at every element until the empty list has been found as the rest of the list. Where in a sequence the last element can be obtained in $\Theta(1)$ time. Adding a element anywhere in the sequence can be done in worst case, $\Theta(\log n)$ [?].

Since this project is about creating a real-time lexing tool, performance is very important. String in Haskell are just a list of characters, `[Char]`, and therefore has $\Theta(n)$ in time cost. Lexing is working with strings in a high frequency and therefore there is an idea of instead define a string as a sequence of characters instead of a list of character.

4.5 Transition Map

4.5.1 Array Format

4.5.2 Function Composition Format

5

Result

The incremental lexer has as mentioned before three requirements, it should be Robust, Efficient and Precise. Robustness means that the lexer doesn't crash when it encounter an error in the syntax. That is, if a string would yield an error when lexed from the starting state the lexer doesn't return that error but instead stores the error and lexes the rest of the possible input states since the current string might not be at the start of the code. For it to be efficient the feedback to the user must be instant, or more formally the combination of two string should be handled in $O(\log(n))$ time. Finally to be precise the lexer must give a correct result. This chapter will talk about how these requirements are tested and what results the from the tests were.

In the sections Below, any mention of a sequential lexer refers to a lexer generated by Alex using the same alex file as was used when creating the incremental lexer [5]. The reason why Alex was used is because the dfa generated by Alex was used in the incremental lexer, thus ensuring that only the lexical routines differs.

5.1 Preciseness

For an incremental lexer to work, the lexer must be able to do lexical analysis of any subtext of a text and be able to combine two subtexts. If the lexical analysis of one subtext doesn't result in any legal tokens it must be able to be combined with other subtexts that makes it legal tokens. The lexical analysis of a subtext might not always result in the same tokens that the combination of the subtext with another text would give.

To test these cases a test was constructed that does a lexical analysis on two subtexts using the incremental lexer and then combining the results into one text. The result

of the combination should be the same as the lexical analysis of the text using the incremental lexer and the result using a sequential lexer.

It is not enough to test if the combination of two subtexts yields the same sequence of tokens as the text. To test that the result of the incremental lexer is the correct sequence of tokens generated, it is compared to what a sequential lexer generates. This comparison is an equality test of the text, it checks token for token that they are the same kind of token and have the same lexeme. fig. 5.1 shows the test for equality:

```
checkCorrectTokens :: IncLex.Tokens -> Alex.Tokens -> Boolean
checkCorrectTokens itoks atoks =
  let tokTupple = zip itoks atoks
  in [] == filter (\(iToken, aToken) -> iToken `notEquals` aToken) tokTupple
```

notEqual function is a function which pattern-match on the two different tokens and returns true if they are not of the same type.

Figure 5.1: Code for testing tokens from IncLex is equal to tokens from Alex.

5.2 Performance

To measure the performance of the incremental step we created the fingertree for two pieces of code. By creating the two fingertrees the transition map for the code in those trees are created aswell. The benchmarking was then done on the combination of the two trees. The results of the incremental lexer benchmarking suggests a running time of $O(\log(n))$.

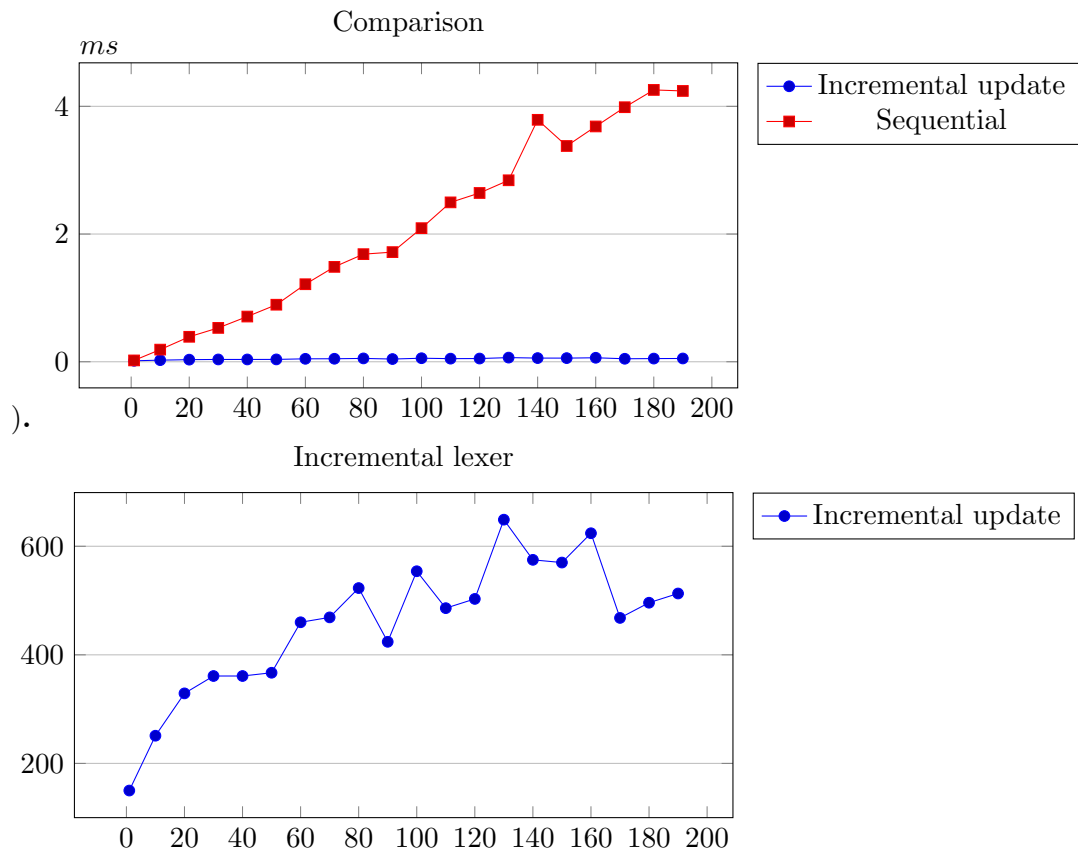
To get a reference point the same text was lexed using a sequential lexer.

#Lägga till nåt om minnes utrymmet som krävs.

#Bättre förklaring till graferna.

#Tester för nybyggnad av träd?

Example 5.2.1 (Benchmarking times of the incremental and sequential lexer



The space complexity for the fingertree structure is $O(n)$. According to our tests 100 characters of code translates to roughly 1MB of datastructure.

6

Discussion

#Discuss discuss!!

One would think that haskells quickcheck would be a good way to generate in-data for the lexer. Since quickcheck is built to generate good input for testing a function for any arguments [3]. But the problem is not to test any string representation, it is instead to test valid code segments and any substring of this code segments. Also invalid pieces of text, to see that the lexer informs the user of syntactical errors in these texts. To write a input generator in quickcheck which would generate full code with all of its components and all the different properties would have to high cost in develop time for the outcome. It would be more time efficient to test the lexer on several different code files. There for the testing of the incremental lexer has not been done with the help of quickcheck.

#possibility of using sequential lexing first time?

#When is it advantagous with incremental lexing

7

Conclusion and Futher Work

#what will our Minions do???

#implementation of ropes

Bibliography

- [1] Alfred V. Aho. *Handbook of theoretical computer science (vol. A)*, chapter Algorithms for finding patterns in strings, pages 255–300. MIT Press, Cambridge, MA, USA, 1990.
- [2] Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman. *Compilers: Principles, Techniques, and Tools (2nd Edition)*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 2006.
- [3] Koen Claessen and John Hughes. Quickcheck: a lightweight tool for random testing of haskell programs. *SIGPLAN Not.*, 35(9):268–279, September 2000.
- [4] Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. *Introduction to Algorithms, Third Edition*. The MIT Press, 3rd edition, 2009.
- [5] Chris Dornan, Isaac Jones, and Simon Marlow. *Alex User Guide*. Haskell. www.haskell.org/alex/doc/html/index.html.
- [6] Michael T Goodrich and Roberto Tamassia. *Data Structures and Algorithms in Java, 4th Edition*. John Wiley & Sons, 2005.
- [7] Dan Piponi. Fast incremental regular expression matching with monoids. blog.sigfpe.com/2009/01/fast-incremental-regular-expression.html, 2009.
- [8] Aarne Ranta and Markus Forsberg. *Implementing Programming Languages*, chapter Lexing and Parsing, pages 38–47. College Publications, London, 2012.
- [9] R.W. Sebesta. *Concepts of Programming Languages [With Access Code]*. Always learning. Pearson Education, Limited, 2012.
- [10] M. Sipser. *Introduction To The Theory Of Computation*. Advanced Topics Series. Thomson Course Technology, 2006.
- [11] R. S. Sundaresh and Paul Hudak. A theory of incremental computation and its application. In *Proceedings of the 18th ACM SIGPLAN-SIGACT symposium on*

Principles of programming languages, POPL '91, pages 1–13, New York, NY, USA, 1991. ACM.

A

Java Lette Light

Here is a simplified version of java that only have variables, numbers and some simple operators. The language includes *while* loops but the lexical analyser will read “while” as an identifier, The syntactical analyser will later determine if it is a loop. The expressions that matches rules without a name are discarded since they aren’t needed for the syntactical analyzer.

Character sets

capital → [A-Z]

lower → [a-z]

letter → [a-zA-Z]

digit → [0-9]

ident → *letter* | *digit* | [-’]

Identifier Characters

white → [\t\r\n\v\f]

The white space characters

cal analyzer. Rules

whie+

// [.]*

Single line comment

/\^* ([\^*] | (\^*) [\^])^* (\^*)+ /

Multi line comment

Identifier → *letter ident**

Integer → *digit*+

Double → *digit*+ \. *digit*+

Reserved → \ (| \) | \ { | \ } | ; | = | \ + \ + | < | \ + | - | \ * | Reserved characters