

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 Program MPALG]

JONAS HUGO

KRISTOFER HANSSON

CHALMERS UNIVERSITY OF TECHNOLOGY
Department of Computer Science and Engineering
Göteborg, Sweden, June 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, June 2014.
© KRISTOFER HANSSON, June 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 troll image, and will not be used in the final report.

Department of Computer Science and Engineering
Göteborg, Sweden June 2014

Abstract

This report aims to present a way to do lexical analysis incrementally instead of the present norm: sequential analysis. In a text editor, where updates are common, an incremental lexer together with an incremental parser could be used to give users real time parsing feedback. Previous work has proven that regular expressions can be implemented incrementally [12], we make use of these findings in order to show that it can be expanded to a lexical analyzer. The results in this report shows that an incremental lexer is efficient, it can do an update in $\Theta \log n$ time which makes it suitable when updates are common. In order for an incremental lexer to be applicable it has to be precise, only correctly lexed tokens are relevant. It is required that an incremental lexer is robust, a lexical error for a partial result must be handled gracefully since it may not propagate to the final result. To achieve incrementality a divide and conquer tree structure, fingertrees, is used that stores the intermediate lexical results of all the partial trees. When an update to the tree is made only the effected node and its parents are updated. The state machine in the implementation is generated by Alex since it is efficient and enables support for lexical analysis of different languages. The report finishes with giving suggestions for improvements to the drawbacks found during the work, The suggestions given are mainly for improving space complexity. This report shows that an implementation of an incremental lexer can be precise, efficient and robust.

Acknowledgments

We would like to take the chance of thanking our supervisor at department of computer science, Jean-Philippe Bernardy, with which we have had a lot of constructive discussions. We would also like to thank our examiner at the department of computer science Bengt Nordström. Lastly we would like to thank everyone that has helped proof read this report.

Jonas Hugo & Kristofer Hansson, Göteborg June 2014

Contents

1	Introduction	1
1.1	Scope of work	1
1.2	Related 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	5
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 Three Steps	12
3.1.2	Associative Function	13
3.1.3	Time Complexity	13
3.1.4	Hands on Example	13
3.1.5	Incremental Computing	15
3.2	Fingertree	15
3.2.1	Fundamental Concepts	16
3.2.2	Simple Sequence	16
3.2.3	Double-ended Queue Operations	19
3.2.4	Concatenation Operations	20
3.2.5	Measurements	22
3.2.6	Sequences	23

3.3	Divide and Conquer Lexing in General	24
3.3.1	Tree structure	24
3.3.2	Transition map	24
3.3.3	Expected Time Complexity	26
3.4	Lexical Errors	26
4	Implementation	28
4.1	The DFA Design	28
4.2	Token data structure	29
4.2.1	Tokens	29
4.2.2	Suffix	29
4.3	Transition Map	30
4.4	Fingertree	31
4.5	Lexical routines	32
4.5.1	Combination of Tokens	32
4.5.2	Combine Tokens With Right Hand Side	32
4.5.3	Merge Two Tokens	33
4.5.4	Merging Suffixes	34
4.5.5	Append to Sequence of Tokens	35
5	Result	36
5.1	Preciseness	36
5.2	Performance	37
6	Discussion	41
6.1	Used Programming Language and Data Structure	41
6.2	Trials and Errors	42
6.3	Implementation Suggestions	43
7	Conclusion and Future Work	44
7.1	Conclusions	44
7.2	Future Work	45
	Bibliography	47
A	Java Lette Light	48
B	Incremental Lexer Source Code	49
C	Benchmarks	55

Acronyms

BNFC Backus-Naur-Form Converter.

DFA Deterministic Finite Automata.

IDE Integrated Development Environment.

NFA Non-deterministic Finite Automata.

Glossary

Alex Is a tool for generating lexical analyzers in Haskell.

Fingertree Is a tree data structure for general sequence representation with arbitrary annotations.

Incremental (computation) Is to save time by only recomputing those outputs which directly depend on the changed data.

Lazy Evaluation Is an evaluation strategy where the evaluation of an expression is delayed until its value is needed.

Lexeme Is the string representation of a token.

Lexer Is a program for converting a sequence of characters into a sequence of tokens.

Monoid Is an algebraic structure with a single associative binary operation and an identity element.

Normal Form An expression is evaluated to normal form when no further computations can be made.

Regular Expression Is a sequence of characters that forms a search pattern.

Regular Language Is a formal language that can be expressed using a regular expression.

Sequence Is an abstract data type that implements a finite ordered collection of elements.

Token Is a string of one or more characters that is significant as a group.

Transition Is a path in a DFA from an in state to an out state given a string.

Transition Map Is a collection of all possible transitions for a string.

1

Introduction

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 compiler strength 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.

In order to implement a parser with the characteristics described above; robust, precise and efficient; a lexer that has the same properties is needed. This project aims to implement such a lexer.

1.1 Scope of work

Existing lexical analyzers are sequential. When the text is updated the lexer must start the lexical analysis from the beginning. The goal of this project is to create an algorithm that, after an update to the text, only needs to recalculate the update and the part of the result effected by the update. The recalculation should have time complexity $\theta \log(n)$ in order to be run in real time, for example in a text editor with immediate update.

The report gives a general understanding of how lexical analyzers work and what tools are needed for a divide and conquer implementation to work. An overview of the ideas behind a divide and conquer lexer is explained in order to give enough understanding for the algorithm this report purposes. A robust implementation of the algorithm is presented with explanations on how different cases are handled. Tests for preciseness,

time performance and space performance are explained and their corresponding result are presented. The report will present a discussion on where a divide and conquer lexer is useful.

1.2 Related Work

This project revolves around the idea of using incremental regular expressions. Piponi wrote a blog post about how to implement incremental regular expressions using finger trees [12]. The solution to matching regular expressions incrementally in the blog post gives a good starting point to this project, however a lexer do not match a string against one expression. A lexer matches a string against a set of regular expressions and returns which expressions where matched and in what order rather then answering if the string matched the expressions [12].

Bernardy and Claessen [3] wrote a paper titled “Efficient Divide-and-Conquer Parsing of Practical Context-Free Languages” that describes an efficient parallel parser built on Valiants algorithm [17]. The paper proves that for a defined set of input the complexity for the parser will be $\theta \log^3(n)$. Since the implementation of the parser is done in BNFC [6] it uses a sequential lexer generated by alex [5]. The lexical analyzer this project purposes is a divide and conquer solution which could be integrated to the parser proposed by Bernardy and Claessen to get a divide and conquer solution from the programming code to the result of the parser.

2

Lexer

A lexer, lexical analyzer, is a pattern matcher. Its job is to divide a text into a sequence of tokens (such as words, punctuation and symbols). A Lexer is often used as a front end to a syntax analyzer [14]. The syntax analyzer in turn takes the tokens generated by the lexer and returns a set of expressions and statements. Lexing 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 analyzer. There are several reasons why a compiler should be separated into a lexical analyzer and a parser (syntax analyzer).

Firstly, simplicity of design is the most important reason. When dividing the task into these two subtasks, it allows the programmer to simplify each of these subtasks. 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 a lexer. When the two tasks have been separated into subtasks it can lead to cleaner overall design when designing a new language. The only thing the syntax analyzer will see is the output from the lexer, tokens (and lexemes). The lexer usually skips comments and white-spaces, since these are often not relevant for the syntax analyzer.

Secondly, overall efficiency of the compiler can be improved. When separating the lexical analyzer it allows for use of specialized techniques that can be used only in the lexical task.

Thirdly and last, compiler portability can be enhanced. This means that 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 [14]. Finding these errors allows the compiler to break the compilation before running the syntax analyzer, thereby saving computing time.

2.2 Token Specification

The job of the lexical analyzer 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

Regular expressions is a key part in describing patterns of a text. However, they cannot express all possible patterns on which a text can follow, but they can describing those type of pattern which are used in a lexer, the pattern of tokens.

Definition 2.2.1 (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 [1].

Example 2.2.2 (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: $letter(letter|digit)^*$.

2.2.2 Languages

An alphabet is a finite set of symbols. For example Unicode includes approximately 100,000 characters. A language is any countable set of strings of some fixed alphabet [2]. The term formal language refers to languages which can be described by a body of systematic rules. There is a subset of languages to formal languages called regular languages, these regular languages refers to those languages that can be defined by regular expressions [13].

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 the rules and the result from lexing the text. This section describe the structures which the lexical analyzer use for representing the abstract data; what these structures are used for and what is forwarded to the syntactical analyzer.

A lexical analyzer uses three different concepts. The concepts are 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 variable or identifier is commonly formed by a letter followed by a sequence of letters, digits and `_`, an integer is a sequence consisting of digits from 0 to 9. This can be described by a regular expression.

Lexemes is a sequence of characters in the code being analyzed which matches the pattern for a token and is identified by the lexical analyzer 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 [14]. In some cases it is not relevant to return a token for a pattern, in these cases the token and lexeme is discarded and the lexer continues, typical examples are whitespaces and comments which have no impact on the code [2].

$$\begin{aligned} \langle letter \rangle &\in \{ 'a' - 'z' \} \cup \{ 'A' - 'Z' \} \cup \{ '_' \} \\ \langle digit \rangle &\in \{ 0 - 9 \} \\ \langle identifier \rangle &::= \langle letter \rangle (\langle letter \rangle \mid \langle digit \rangle)^* \\ \langle integer \rangle &::= \langle digit \rangle^+ \\ \langle multi-line\ comment \rangle &::= '/*' ([\wedge '*'] \mid '*' [\wedge '/'])^* '*/' \\ \langle reserved-words \rangle &::= '(' \mid ')' \mid '{' \mid '}' \mid ';' \mid '=' \mid '++' \mid '<' \mid '+' \mid '-' \mid '*' \end{aligned}$$

Figure 2.1: Grammar rules for example 2.3.1 & example 2.4.1

An example follows how a small piece of code would be divided given the regular language described in appendix A.

Example 2.3.1 (Logical grouping [14]).

Consider the following text; to be lexed:

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

Given the regular language defined in appendix A, the lexical analyzer would use the rules defined in fig. 2.1 In order to produce the resulting tokens shown in fig. 2.2.

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

Figure 2.2: Result of lexing the code in example 2.3.1

2.4 Recognition of Tokens

The topic in the previous section covers how to represent a pattern using regular expressions and how these expressions relate to tokens. A pattern is used to determine if a string matches a token. This section is highlighting how to transform a sequence of characters into a sequence of abstract tokens using patterns.

2.4.1 Transition Diagrams

A state transition diagram, or transition diagram is a directed graph, where the nodes are labeled 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 labeled 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 [14].

There are different type of states in the the diagram where one state state is said to be initial. 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 do not have any legal edge out of that state, the lexer can not 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 was not 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;
```

All though this text 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. The rules relevant to this example are defined in fig. 2.1 the rest of the rules can be found in appendix A.

The result can be found in fig. 2.3.

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

Figure 2.3: Result of lexing the code in example 2.4.1

2.4.3 Finite Automata

Transition diagrams of the form used in lexers are representations of a class of mathematical machines called finite automata. Finite automata can be designed to recognize members of a class of languages called regular languages, mentioned above [14]. A finite automaton is essentially a graph, like transitions diagrams, with some differences:

- Finite automata are recognizers; they 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 labeled by the same symbol out from the same state. Further, the empty string ϵ 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 recognizing 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 [15]). 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 accepting states.

The transition function does not 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 into 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 & Transition Table [2]). Given this regular expression:

$$(a|b) * abb$$

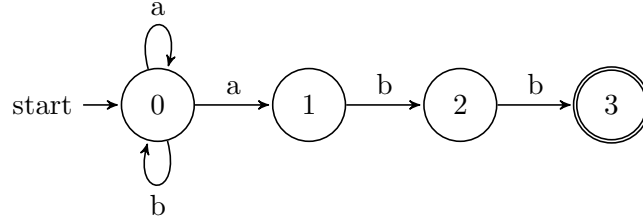


Figure 2.4: Transition Diagram, accepting the pattern $(a|b)^*abb$

The transition diagram in fig. 2.4 is representing this regular expression.

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

Figure 2.5: Transition Table, accepting the pattern $(a|b)^*abb$

It could also be converted into the transition table shown in fig. 2.5

Transition tables have the advantage that they have an 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 labeled with a .

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

Definition 2.4.4 (Deterministic Finite Automata [15]). 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 accepting states.

Example 2.4.5 (DFA representation of RegExp [2]). A DFA representation of same regular expression from example 2.4.3 is shown in fig. 2.6

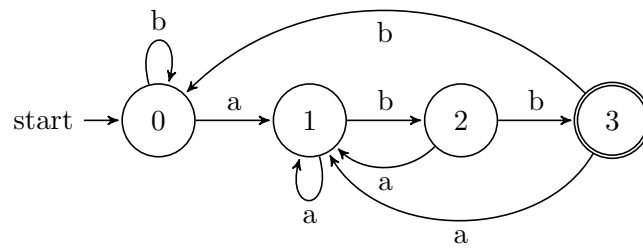


Figure 2.6: DFA, accepting the regular expression: $(a|b)^*abb$

3

Divide-and-Conquer Lexer

An incremental divide and conquer lexer works by dividing the sequence, to be lexically analyzed, into small parts and analyzes them; and then combining them. In the base case the lexical analysis is done on a single character. The conquer step then combines the smaller tokens into as large tokens as possible. The end result is a sequence of tokens that represent the code. How this is done is described in this chapter.

3.1 Divide and Conquer in General

This section gives an idea of how the Divide and Conquer algorithm works in general, before addressing in detail how to apply it to lexing. It describes the power of divide and conquer in terms of executing time and how laziness can be applied to these algorithms.

3.1.1 The Three Steps

The general idea of a divide and conquer algorithm is to divide a problem into smaller parts, solve them independently and then combine the results. A Divide and Conquer algorithm always consists of a pattern with these three steps [7].

Divide: If the input size is bigger than the base case 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 subproblems, combine the results to solve the original problem.

3.1.2 Associative Function

An associative function, or operator, is a function that does not care in what order it is applied. An example of such a function is addition (+) of numbers, which is associative since it has the property in example 3.1.1.

In divide and conquer algorithms this is essential. In the divide step of the divide and conquer algorithm, there is no certain order of how the subproblems are going to be divided. This means that the order the subproblems are being conquered can not have an impact on the algorithm, hence the conquer step must be associative.

Example 3.1.1 (Associativity of the conquer step). Let $f(x,y)$ be the conquer function, where x and y are of the same type as the result of f , then:

$$f(x, f(y, z)) = f(f(x, y), z)$$

Otherwise the algorithm can give different results for the same data.

3.1.3 Time Complexity

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

Theorem 3.1.2 (Master Theorem [4]).

Assume a function T_n constrained by the recurrence

$$T_n = \alpha T_{\frac{n}{\beta}} + f(n)$$

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

If we let $e = \log_{\beta} \alpha$, 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 $\alpha \cdot f(n/\beta) \leq c \cdot f(n)$
where $c < 1$ and all sufficiently large n

■

3.1.4 Hands on Example

The divide and conquer pattern can be performed on algorithm that solves different problems. A general problem is sorting, or more precisely sorting a sequence of integers. This example shows merge-sort.

Divide: The algorithm starts with the divide step. Given the input S the algorithm will check if the length of S is less than or equal to 1.

- If this is true, the sequence is returned. A sequence of one or zero elements is always sorted.
- If this is false, the sequence is split into two equally big sequences, S_1 and S_2 . S_1 will be the first half of S while S_2 will be the second half.

Recur: The next step is to sort the subsequences S_1 and S_2 . The sorting function sorts the subsequences by recursively calling itself twice with S_1 and S_2 as arguments respectively.

Conquer: Since S_1 and S_2 are sorted combining them into one sorted sequence is trivial. This process is what is referred to as merge in merge-sort. The resulting sequence of the merge is returned.

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

```

Given the merge-sort algorithm, time complexity can be calculated as follows using the master method. There are 2 recursive calls and the subproblems are $1/2$ of the original problem size, so $\alpha = 2$ and $\beta = 2$. To merge the two sorted subproblems the worst case is to check every element in the two list, $f(n) = 2 \cdot n/2 = n$.

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

$$e = \log_{\beta} \alpha = \log_2 2 = 1$$

Case 2 of the theorem 3.1.2 applies, since

$$f(n) = \Theta(n)$$

So the solution will be:

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

3.1.5 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 [16]. This is illustrated in fig. 3.1.

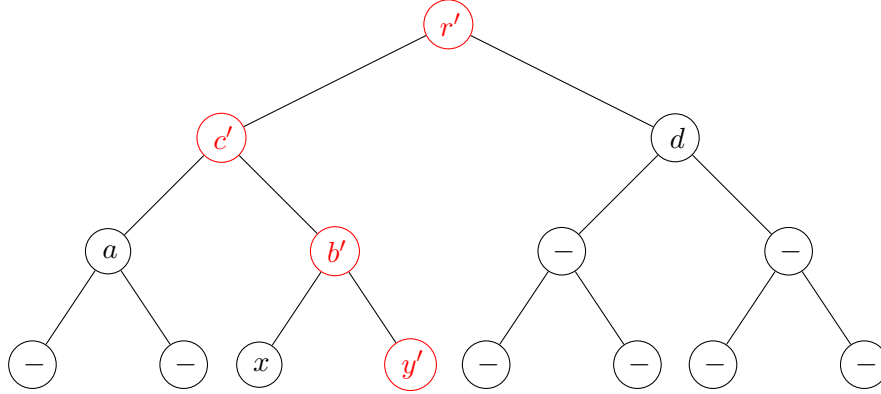


Figure 3.1: When the node y changed recomputed nodes are marked with a '.

For a divide and conquer lexer this means to 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 recalculate all the 1024 tokens the lexer only needs to do 10 recalculations. Since, $\log_2 1024 = 10$. This can be explained by the theorem 3.1.2.

Only one branch in the tree will be followed at every level and the problem is already divided. Therefore the parameters will be set to:

$$\alpha = 1, \beta = 2 \text{ and } f(n) = 1. \quad e = \log_\beta \alpha = \log_2 1 = 0$$

Case 2 of the theorem 3.1.2 applies, since

$$f(n) = \Theta(n^e)$$

The complexity is therefore:

$$T(n) = \Theta(n^e \cdot \log n) = \Theta(\log n)$$

3.2 Fingertree

Fingertree is a tree structure which is incremental in its nature and has good performance. To ensure that an incremental divide and conquer algorithm can access the intermediate states a data structure like fingertrees can be used. Before describing

how the fingertree is defined, an introduction to the fingertrees building blocks is given [8].

3.2.1 Fundamental Concepts

Fingertrees uses monoids which in abstract algebra is a set, S , and a binary operation \bullet which fulfills 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 $\exists e \in S : \forall a \in S : e \bullet a = a \bullet e = a$

Fingertrees uses Right and Left Reductions. This is a function which collapses a structure of $f\ a$ into a single value of type a . The base case for when the tree is empty is replaced with a constant value, such as \emptyset . Intermediate results are combined using a binary operation, like \bullet . Reduction with a monoid always returns the same value, independently of the argument nesting. For a reduction with an arbitrary constant and binary operation there must be a specified nesting rule. If combining operations 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 described in fig. 3.2 [8].

```
class Reduce f where
  reducer :: (a -> b -> b) -> (f a -> b -> b)
  reducel :: (b -> a -> b) -> (b -> f a -> b)
```

Figure 3.2: Reduction function in Haskell

3.2.2 Simple Sequence

The fingertrees can be described by comparing it to an already known data structure and how that data structure represent data. Lets take a look at the definition on a 2-3 fingertree and how they can implement a sequence. Lets start by looking at an ordinary 2-3 tree representing the string "thisisnotaree".

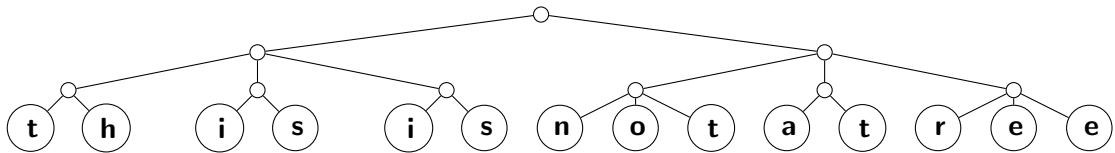


Figure 3.3: Ordinary 2-3 tree

The tree shown in the fig. 3.3 stores all it is data in the leafs. This can be expressed by defining an non-regular or nested type, as shown in fig. 3.4.

```

data Tree a = Zero a | Succ (Tree (Node a))
data Node a = Node2 a a | Node3 a a a
    
```

Figure 3.4: Definition of a 2-3 Fingertree

Operations on these types of trees usually takes logarithmic time in the size of the tree. However in a sequence representation, constant time complexity is preferable for adding or removing element from the start or end of the sequence.

A finger is a structure which provides efficient access to nodes near the distinguished location. To obtain efficient access to the starting and ending elements of the sequence represented by the tree, there should be fingers placed at these positions of the tree. In the example tree, taking hold of the end and start nodes of and lifting them up together. The result should look like in fig. 3.5

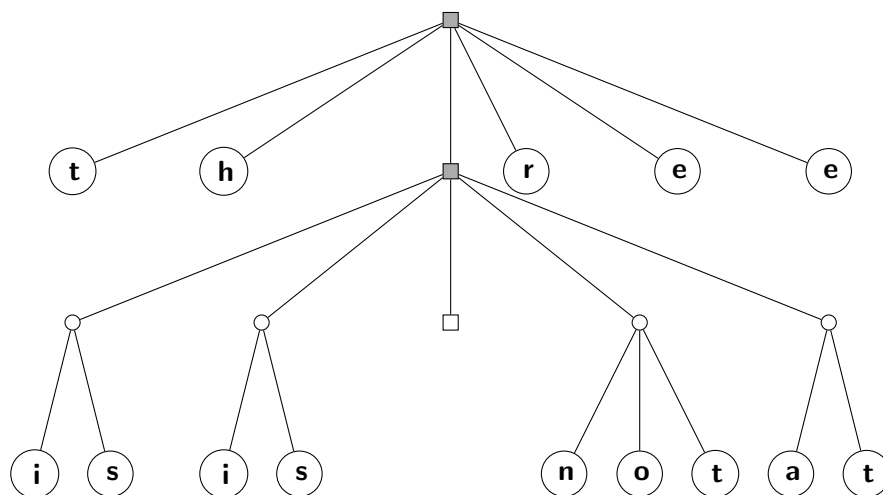


Figure 3.5: 2-3 Fingertree

Since all leafs in the 2-3 tree are at the same level, the left and right spine has the same length. Therefore the left and right spines can be paired up to create a single central spine. Branching out from the spine are 2-3 trees. At the top level there are two to three elements on each side, while the other levels have one or two subtrees, whose depth increases down the spine. Depending on if the root node had 2 or 3 branches in the original 2-3 tree, The bottom node will have either a single 2-3 tree or an empty tree. This structure can be described as shown in fig. 3.6.

Where Digit is a buffer of elements stored left to right, here represented as a list for simplicity

The non-regular definition of the *FingerTree* type determines the unusual shape of these trees, which is the key to their performance. The top level of the tree contains elements

```

data FingerTree a = Empty
                  | Single a
                  | Deep (Digit a) (FingerTree (Node a)) (Digit a)

type Digit a = [a]

```

Figure 3.6: Definition of the Fingertree data type

of type a . Next level contains elements of type $Node\ a$. At the n th level, elements are of type $Node^n\ a$, which are 2-3 trees with a depth of n . This gives that a sequence of n elements is represented by a *FingerTree* of depth $\Theta(\log n)$. An element at position d from the nearest end is stored at a depth of $\Theta(\log d)$ in the *FingerTree*.

In fingertrees and nodes the reduce function mentioned in fundamental concepts is generically defined to the following types. Reduction for the node which is shown in fig. 3.7.

```

instance Reduce Node where
  reducer (-<) (Node2 a b) z = a (-<) (b (-<) z)
  reducer (-<) (Node3 a b c) z = a (-<) (b (-<) (c (-<) z))

  reducel (>-) z (Node2 a b) = (z (>-) b) (>-) a
  reducel (>-) z (Node3 c b a) = ((z (>-) c) (>-) b) (>-) a

```

Figure 3.7: Reduction of a fingertrees node

Reduction of fingertrees single and double lifting of the binary operation as shown in fig. 3.8 [8].

```

instance Reduce FingerTree where
  reducer (-<) Empty z = z
  reducer (-<) (Single x) z = x (-<) z
  reducer (-<) (Deep pr m sf) z = pr (-<') ( m (-<'') ( sf (-<')
    z ))
    where (-<') = reducer (-<)
          (-<'') = reducer (reducer (-<))

  reducel (>-) z Empty = z
  reducel (>-) z (Single x) = z (>-) x
  reducel (>-) z (Deep pr m sf) = ((z (>-') pr ) (>-'') m ) (>-')
    sf
    where (>-') = reducel (>-)
          (>-'') = reducel (reducel (>-))

```

Figure 3.8: Reduction of a Fingertree

3.2.3 Double-ended Queue Operations

After showing how the Fingertrees basic structure is defined, lets take a look on how fingertrees makes efficient Double-ended Queue, a queue which can be accessed from both ends, where both the operations having the time complexity $\Theta(1)$.

Adding an element to the beginning of the sequence is trivial, except when the initial buffer (*Digit*) already is full. In this case, push all but one of the elements in the buffer as a node, leaving behind two elements in the buffer, shown in fig. 3.9.

```

(<|)                :: a -> FingerTree a -> FingerTree a
a (<|) Empty        = Single a
a (<|) Single b      = Deep [a] Empty [b]
a (<|) Deep [b,c,d,e] m sf = Deep [a, b] (Node3 c d e (<|) m) sf
a (<|) Deep pr m sf  = Deep ([a] ++ pr) m sf
    
```

Figure 3.9: Adding an element to the beginning of the sequence

Adding to the end of the sequence is a mirror image of the code in fig. 3.9 and is shown in fig. 3.10.

```

(|>)                :: FingerTree a -> a -> FingerTree a
Empty                (|>) a = Single a
Single b             (|>) a = Deep [b] Empty [a]
Deep pr m [e,d,c,b] (|>) a = Deep pr (m (|>) Node3 e d c) [b,a]
Deep pr m sf         (|>) a = Deep pr m (sf ++ [a])
    
```

Figure 3.10: Adding an element to the end of the sequence

In the basic 2-3 tree, where the data is stored in the leaves, an insertion operation is done with a time complexity of $\Theta(\log n)$. The expected time complexity of a fingertree can be expressed in this way: Digits of two or three elements (which is being of similar structure to elements of type *Node a*) are classified as safe since these can easily be mapped to a *Node* and those of one or four elements are classified as dangerous. A double-ended queue operation can only propagate to the next level from a dangerous element. By doing so making that dangerous element safe, which means that the next operation reaching that digit will not propagate. This will result in that at most 1/2 of the operations descend one level, at most 1/4 two levels, and so on. This will give that in a over time the average cost of operations is constant.

The same bounds hold in a persistent setting if subtrees are suspended using lazy evaluation. Laziness makes sure that changes deep in the spine do not take place until a subsequent operation need to go that far. From the above properties of safe and dangerous digits, by the time a change deep down in the tree is needed enough cheap shallow operations will have been performed to pay for the more expensive operation [8]. This

can be expressed more in detail by the bankers method.

The Bankers Method

The bankers method is a technique used to calculate the practical time assumption where it accounts for accumulated debt. Each debit represents a constant amount of suspended work. When a computation initially suspends, it create a number of debits proportional to it is shared cost and associate each debit with a location in the object. The choice of location for each debit depends on the nature of the computation. If the computation is monolithic (i.e., once begun, it runs to completion), then all debits are usually assigned to the root of the result, which the fingertree is not. But if the computation is lazy, like the fingertree, then the debits may be distributed among the roots of the partial results.

The amortized cost of an operation is the unshared cost of the operation plus the number of debits discharged by the operation. Note that the number of debits created by an operation is not included in its amortized cost. The order in which debits should be discharged depends on how the object will be accessed; debits on nodes likely to be accessed soon should be discharged first.

Incremental functions play an important role in the bankers method because they allow debits to be dispersed to different locations in a data structure, each corresponding to a nested suspension. Then, each location can be accessed as soon as its debits are discharged, without waiting for the debits at other locations to be discharged. This means that the initial partial results of an incremental computation can be paid for very quickly, and that subsequent partial results may be paid for as they are needed [10].

Banker Method on the Fingertree

The amortized time can be argued for by using the Banker method. This is done by assigning the suspension of the middle tree in each Deep node as many debits as the node has safe digits. (0,1 or 2) A double-ended queue operation which descends k levels turns k dangerous digits into safe digits. By doing so creates k debits to pay for the work done. Applying the bankers method of debit passing to any debits already attached to these k nodes. It can be shown that each operation must discharge at most one debit. Therefore the double-ended queue operations run in $\Theta(1)$ amortized time [8].

3.2.4 Concatenation Operations

Concatenation is a simple operation for most cases, except for the case when two *Deep* trees are being concatenated. Since *Empty* is the identity element, concatenation with

an *Empty* yields the other tree. Concatenation with a *Single* will reduce to $< |$ or $| >$. For the hard part when there are two *Deep* trees, the prefix of the first tree will be the final prefix. Suffix of the second tree will be the suffix of the final tree. The recursive function *app3* shown in fig. 3.11 combines two trees and a list of *Nodes* (basically the old prefix and suffixes down the spines of the old trees).

```
app3 :: FingerTree a -> [a] -> FingerTree a -> FingerTree a
app3 Empty ts xs      = ts (<|') xs
app3 xs ts Empty      = xs (|>') ts
app3 (Single x) ts xs = x (<|) (ts (<|') xs)
app3 xs ts (Single x) = (xs (|>') ts) (|>) x
app3 (Deep pr1 m1 sf1) ts (Deep pr2 m2 sf2)
    = Deep pr1 (app3 m1 (nodes (sf1 ++ ts ++ pr2)) m2) sf2
```

Figure 3.11: Help function for concatenating two fingertrees

Where $(<|')$ and $(|>')$ are the functions defined in fig. 3.12 and *nodes* groups a list of elements into *Nodes* as shown in fig. 3.13.

```
(<|') :: (Reduce f) => f a -> FingerTree a -> FingerTree a
(<|') = reducer (<|)

(|>') :: (Reduce f) => FingerTree a -> f a -> FingerTree a
(|>') = reducer1 (|>)
```

Figure 3.12: Help function for transforming a list of element into a list of Nodes

```
nodes :: [a] -> [Node a]
nodes [a, b]          = [Node2 a b]
nodes [a, b, c]       = [Node3 a b c]
nodes [a, b, c, d]    = [Node2 a b, Node2 c d]
nodes (a : b : c : xs) = Node3 a b c : nodes xs
```

Figure 3.13: Help function for transforming a list of element into a list of Nodes

The concatenation of the Fingertrees calls on *app3* with an empty list between the two trees, as shown in fig. 3.14.

```
(><) :: FingerTree a -> FingerTree a -> FingerTree a
xs (><) ys = app3 xs [] ys
```

Figure 3.14: Concatenation function for Fingertree

The time spent on concatenation can be reasoned in this way. Each invocation of *app3* arising from $(><)$ the argument list has a length of at most 4, which means that each

of these invocations takes $\Theta(1)$ time. The recursion terminates when the bottom of the shallower tree has been reached, with up to 4 insertions. So the total time complexity is $\Theta(\log \min\{n_1, n_2\})$ where n_1 and n_2 are the number of elements in the two trees [8].

3.2.5 Measurements

Fingertrees has been shown to work well as catenable double-ended queues. A measurement is a property describing the state of the tree. This section present Paterson and Hinze modification of the fingertree, which provides positional and set-theoretic operations. For example taking or dropping the first n elements. These operations involve searching for an element with a certain property. To implement additional operations with good performance there must be a way to steer this search. A way to measure the tree.

A measurement can be viewed as a cached reduction with some monoid. Reductions, possibly cached, are captured by the class declaration in fig. 3.15. Where a is the type of a tree and v the type of an associated measurement. v must be of a monoidal structure so measurements of subtrees easily can be combined independently of the nesting. Take the size of a tree as an example. Measure maps onto the monoid over the set of natural numbers and with the binary operator of addition [8].

```
class (Monoid v) => Measured a v where
  || · || :: a -> v
```

Figure 3.15: Definition of the Measure class

Caching measurements

It should be cheap to obtain a measurement. The fingertree should ensure that an measurement can be obtain with a bounded number of \bullet operations. Therefore fingertrees cache the measurements in the 2-3 nodes. In fig. 3.16 the Measure of Node is shown. The constructors and the instance declaration are completely generic: they work for arbitrary annotations.

Digits are measured on the fly. As the length of the buffer Digit is bounded by a constant, the number of \bullet operations is also bounded.

Fingertrees are modified in a similar manner to 2-3 nodes. The top level of a measured fingertree contains elements of type a , the second level of type $Node\ v\ a$, the third of type $Node\ v\ (Node\ v\ a)$, and so on. The Measure function is shown in fig. 3.18. The tree type a changes from level to level, whereas the measure type v remains the same. This means that Fingertree is nested in a , but regular in v [8].

```

data Node v a = Node2 v a a | Node3 v a a a

node2 :: (Measured a v) => a -> a -> Node v a
node2 a b = Node2 (||a|| • ||b||) a b

node3 :: (Measured a v) => a -> a -> a -> Node v a
node3 a b c = Node3 (||a|| • ||b|| • ||c||) a b c

instance (Monoid v) => Measured (Node v a) v where
    measure (Node2 v _ _) = v
    measure (Node3 v _ _ _) = v
    
```

Figure 3.16: Measure of the data type Node

```

instance (Measured a v) => Measured (Digit a) v where
    measure xs = reducel (\i a -> i • ||a||) • xs
    
```

Figure 3.17: Measure of the data type Digit

```

data FingerTree v a = Empty
| Single a
| Deep v (Digit a) (FingerTree v (Node v a)) (Digit a)

deep :: (Measured a v) =>
    Digit a -> FingerTree v (Node v a) -> Digit a -> FingerTree v a
deep pr m sf = Deep (||pr|| • ||m|| • ||sf||) pr m sf

instance (Measured a v) => Measured (FingerTree v a) v where
    measure Empty      = ∅
    measure (Single x) = measure x
    measure (Deep v)   = v
    
```

Figure 3.18: Fingertrees Measure function

3.2.6 Sequences

A sequence in Haskell is a special case of the fingertree that has no measure. The performance is therefore superior to that of standard lists. 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. Results in 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 an element anywhere in the sequence can be done in worst case, $\Theta(\log n)$ [8].

3.3 Divide and Conquer Lexing in General

In the last section we covered the general divide and conquer algorithm. This section covers the general data structures and algorithms for an incremental divide and conquer lexer.

3.3.1 Tree structure

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 base case is reached. This is clearly a tree structure of solutions, where a leaf is a token for a single character, and the root is a sequence of all tokens in the code.

3.3.2 Transition map

When storing a result of a lexed string it is a good idea to store more than just the tokens. In particular the in and out states are needed when combining the lexed string with another string. We will henceforth refer to this as a *transition*.

```
type Transition = (State,[Token],State)
```

Since the lexer does not know if the current string is a prefix of the entire code or not it can not make any assumptions on the in state. Because of this the lexer needs to store a transition for every possible in state, we will henceforth refer to this as a *transition map*.

```
type Transition_map = [Transition]
```

The Base Case When the lexer tries to lex one character it will create a transition map using the DFA for the language. It will for each state create a transition that has the state as in state, a list containing the character as the only token and by using the DFA, lookup what out state the transition should have. For the character '/' part of a transition map might look like the following.

In the examples below the first number refers to the in state, the middle part is the sequence of tokens and the second number is the out state, that can be accepting.

$$\left[\begin{array}{ccc} 10 & ['/'] & 10 \\ 11 & & NoState \\ 12 & ['/'] & 10 \end{array} \right]$$

NoState transition is used to tell the lexer that using that particular transition will result in a lexical error. For reasons being covered later in this section, they can not be discarded.

Conquer Step The conquer step of the algorithm is to combine two transition maps into one transition map. This is done by, for every transition in the left transition map, combining the transition with the transition in the right transition map that has the same in state as the left transitions out state. This can be described by the following logical statement where T_1 and T_2 refers to the first and second transition map.

$$\forall.t_1 \in T_1 \exists.t_2 \in T_2 \ o_1 = i_2, o_1 = \text{outState}(t_1), i_2 = \text{inState}(t_2) \vdash t_{\text{new}} = \text{merge}(t_1, t_2)$$

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

Concat: If the out state of the first transition is accepting, the sequence in the transition that starts in the starting state of the second transition map will be appended to the first.

```
concat :: Tokens -> Tokens -> Tokens
concat tokens1 tokens2 = tokens1 <> tokens2
```

Combine: If the out state of the first transition is not accepting, the transition in the second transition map with the same in state as the out state of the first transition will be used. The last token of the sequence from the first transition will be combined with the first token in the second transition into one token and put between the two sequences.

```
combine :: Tokens -> Tokens -> Tokens
combine tokens1 tokens2 = prefix1 |> newToken <> suffix2
  where prefix1 |> token1' = tokens1
        tokens2' <| suffix2 = tokens2
        newToken           = token1' 'combinedWith' tokens2'
```

For both the cases the in state of the first transition will be the new in state and the out state of the second transition will be the new out state.

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

This will not work as a lexer for most languages since the longest match rule is not implemented. For example, 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.

Longest Match Instead of taking the naive approach, 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 is not *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.

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

The second case is when the out state for the right token list is *NoState*. This means that the two lists of tokens can not 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.

3.3.3 Expected Time Complexity

Incremental computing states 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 recalculate every node on this path. As shown by fig. 3.1. Only one subproblem is updated in every level of the tree. Using the master method to calculate the expected time complexity gives: $e = \log_b a$ 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, the number of needed update calls at every level of the tree is 1, therefore $a = 1$. The constant b is still 2. This will give $e = \log_2 1 = 0$. Thus the update function of the incremental algorithm will have an expected time complexity of $\Theta(n^0 \cdot \log n) = \Theta(\log n)$

Since the fingertree is lazy, when an element is added to the root level of the tree, root elements might be pushed down in the tree. The measure of the lower levels does not need to be immediately recalculated. Instead they are recalculated when they are used. Paying for this expensive operation like described in the section about bankers method [8].

3.4 Lexical Errors

Since the lexer has to be able to handle any kind of possible not "complete" tokens, error handling has to take this into account. for instance if an illegal character would

be found somewhere in the code the lexer can not return a lexical error right away since said character could be part of a comment. One way to handle this is to tie the lexical errors to transitions. in doing so the transition map can have one transition that is valid and another that produces a lexical error. To include this the transition data type can be modified as follows:

```
type Transition = (State, Maybe ([Token], State))
```

If the transition is $(State_i, Nothing)$ there is a lexical error for the string lexed if the incoming state is $State_i$.

Example 3.4.1 (A lexer that only lexes letters). #This needs to be rewritten!

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

In this chapter the tools, data structure and implementation of the incremental divide and conquer lexer is explained. The implementation of the incremental divide and conquer lexer uses fingertrees for storing the intermediate tokens and the lexed text. It has an internal representation of the tokens to keep track of the data needed when two fingertrees are combined. The lexical routines for combining the internal token data type take advantage of functional composition in order to get lazy updating of the tokens when two fingertrees are combined. The complete implementation can be found in appendix B.

4.1 The DFA Design

The DFA used in the incremental lexer was created using Alex. Alex is a Haskell tool for generating lexical analyzers given a description of the language in the form of regular expressions, it is similar to lex and flex in C and C++. The resulting lexer is Haskell 98 compatible and can easily be used with the parser Happy, a parser generator for Haskell [5]. Alex is notably used in BNFC which is a program to generate among other things a lexer, parser and abstract syntax from Backus-Naur Form [6].

The reason for using Alex to generate the DFA is that it optimizes the number of elements in the transition table. Instead of having an array for every possible character and state combination, 5 arrays are generated that takes advantage of the fact that for most characters the same state will be used the majority of time. This saves a lot of elements that would otherwise be the same in the array.

The trade off for using the Alex generated DFA is that some minor arithmetic operations are used and some extra lookups are needed. These operations are far less time

consuming then the rest of the lexical operations.

4.2 Token data structure

To keep all the information that might be needed when combining two texts, a data structure for the tokens was created. This data type contains more information about the last token then what a sequential lexer would save, exactly what is explained in section 4.2.2.

Since this project is about creating a real-time lexing tool, performance is important. Therefore there are advantages of using sequences instead of lists, since they have better time complexity. The most notable place where this is used is in the measure of the fingertree, where the tokens are stored in a sequence rather than a list. Sequences are also used elsewhere in the project but the measure is the most notable place since it is frequently updated.

4.2.1 Tokens

The internal structure used to store lexed tokens is called *Tokens*. There are three constructors in the *Tokens* data type, see fig. 4.1.

```
data Tokens    = NoTokens
                | InvalidTokens (Seq Char)
                | Tokens { currentSeq :: (Seq Token)
                        , lastToken  :: Suffix
                        , outState   :: State }
```

Figure 4.1: Tokens Data Type

NoTokens is a representation of when an empty string has been lexed. *InvalidTokens* represents a lexical error somewhere in the text that was lexed, the sequence of characters is the lexical error or last token lexed. The *Tokens* constructor is the case when legal tokens have been found. *currentSeq* are all the currently lexed tokens save for the last, *lastToken* are all the possible ways that the last token can be lexed, in this implementation this is referred to as the suffix and what it is and why it is needed will be explained next.

4.2.2 Suffix

When a text is lexed it is uncertain that the last token is the actual end of the file since it may be combined with something else. To ensure that all possible outcomes will be

handled the last token can be on one of three different forms. The part of the text lexed can end in:

- a legal state that is not accepting.
- an accepting state.
- a legal state that is not accepting, but the text can also be a sequence of multiple tokens.

To keep track of these cases a data structure that captures this was implemented, see fig. 4.2.

```
data Suffix    = Str (Seq Char)
               | One Token
               | Multi Tokens
```

Figure 4.2: Suffix Data Type

The *Str* constructor is used to keep track of partially complete tokens, an example of this is when a string is started but the end quotation character have not yet been found.

The *One* constructor is used one exactly one token have been found, it may or may not be the token that is used in the final result of the lexing. Since this constructor is a special case of the *Multi* constructor it can be omitted. However the *One* constructor makes certain cases redundant since the lexer makes assumptions that can not be made for the *Multi* constructor.

The *Multi* constructor is used when at least one token have been found but the lexeme for the suffix does not match exactly one token. The entire suffix still need to have a legal out state. This type of suffix can typically be found when the beginning of a comment are lexed. for example the text */*hello world* would be lexed to a sequence of complete tokens, */*, ***, *hello* and *world*, but the lexer still needs to keep track of the fact that it may be in the middle of a multi-line comment. Note that in this case the *Tokens* data structure would have one out state, the state for the middle of a comment, and the suffix would have another, the end of an ident.

4.3 Transition Map

The transition map is a function from an in state to *Tokens*. As shown in fig. 4.1 the *Tokens* data type contains the out state.

This data type is used in the lexical routines. The reason for using transition maps is that the lexer does not know what the in state for a lexed text is, hence the tokens for all possible in states must be stored. The transition map can be implemented in two ways, a table format and a function composition format.

```
type State = Int

type Transition = State -> Tokens

getTokens :: Transition -> State -> Tokens
getTokens trans state = trans state
```

Figure 4.3: Transition Data Type

The table format uses an array to store the currently lexed tokens where the index of the array represents the in state for that sequence of tokens. This is useful when the tokens needs to be stored since it ensures that the tokens are computed.

When combining lexed tokens it is useful to use functional composition since it ensures that no unnecessary states will be computed. The drawback is that it does not guarantee that the actual tokens are computed which may result in slow performance at a later stage in the lexing.

Both these representations are used in the incremental divide and conquer lexer. The table format is used when storing the tokens in the fingertree to allow for fast access. The function composition is used when combining tokens to ensure that only needed data is computed.

4.4 Fingertree

The fingertree is constructed with the characters of a text being the leaves and with the table format transition map as it is measure. The *Table* data type has to be a monoid in order to be a legal measure of the fingertree.

```
type LexTree = FingerTree Table Char

type Table = Array State Tokens
```

Figure 4.4: The data type for storing the tokens and text

The monoid class in Haskell have two different functions, *empty* which is the identity element and *mappend* which is an associative operator that describes how two elements are combined. As can be seen in fig. 4.5, *empty* creates an array filled of empty *Tokens*. *mappend* extracts the functions from the old tables, combines them using *combineTokens* then creates a new table filled with the combination.

There are two helper functions that convert between the table format that is stored as the measure and the function composition format that is used in the lexical routines.

```

tabulate :: (State, State) -> Transition -> Table
access  :: Table -> Transition

tabulate range f = listArray range [f i | i <- [fst range..snd
    range]]
access a x = a ! x

instance Monoid Table where
    mempty = tabulate stateRange (\_ -> emptyTokens)
    f `mappend` g = tabulate stateRange $ combineTokens (access f)
        (access g)

```

Figure 4.5: The *tabulate* functions and monoid implementation

4.5 Lexical routines

The lexical routines are divided into five functions. They each handle different parts of the lexical steps that are needed in an incremental divide and conquer lexer.

4.5.1 Combination of Tokens

```

combineTokens :: Transition -> Transition -> Transition
combineTokens trans1 trans2 in_state
    | isInvalid toks1 = toks1
    | isEmpty toks1   = trans2 in_state
    | otherwise      = combineWithRHS toks1 trans2
    where toks1 = getTokens trans1 in_state

```

Figure 4.6: The *combineTokens* function

combineTokens is the function called when two fingertrees are combined. The function starts by checking if the tokens generated from *in_state* from the first transition is empty or invalid in which case the output is trivial. If the tokens generated are valid, the tokens are passed on to *combineWithRHS* together with the second transition.

4.5.2 Combine Tokens With Right Hand Side

combineWithRHS checks how the tokens from the first transition are to be combined with the second transition.

combineWithRHS starts by creating tokens from the second transition, *toks2*, using the out state from the first tokens, this can result in three different cases, the definition of the variable names can be found in fig. 4.7.

```

combineWithRHS :: Tokens -> Transition -> Tokens
combineWithRHS toks1 trans2 | isEmpty toks2 = toks1
                           | isValid toks2 =
    let toks2' = mergeTokens (lastToken toks1) toks2 trans2
    in appendTokens seq1 toks2'
                           | otherwise = case lastToken toks1 of
Multi suffToks ->
    let toks2' = combineWithRHS suffToks trans2
    in appendTokens seq1 toks2'
One tok -> appendTokens (seq1 |> tok) (getTokens trans2
startState)
Str s -> invalidTokens s
where toks2 = getTokens trans2 (outState toks1)
      seq1 = currentSeq toks1

```

Figure 4.7: *CombineWithRHS* function

isEmpty If *toks2* is empty *toks1* is returned.

isValid If *toks2* is valid it means that the last token from the *toks1* can be combined into one token with the first token in *toks2*.

otherwise If *toks2* is not valid the lexer checks the suffix of *toks1* to see if it ends in an accepting state or a valid state.

- if the *One* constructor is found the suffix ends in an accepting state which means that tokens created from the start state can be appended to *toks1*.
- If the *Multi* constructor is found the tokens from the suffix, *suffToks*, is extracted and a recursive call to *combineWithRHS* is made with *suffToks* as argument instead.
- If the *Str* constructor is found the suffix does not end in a valid state and *InvalidTokens* will be returned.

4.5.3 Merge Two Tokens

mergeTokens combines the last token from the first tokens with the first token of the second tokens, for the code see fig. 4.8.

- If there are more then one token in *toks2*, *suff1* is combined into one token with the first token in *toks2* and the rest of the tokens in *toks2* is appended and returned.
- If there is exactly one token in *toks2*, the suffix from *toks1* is combined with *suff1*. When two suffixes are combined some extra checks are needed. If *toks2* has an accepting out state, the two suffixes can be combined into one token. If *toks2* does not have an accepting out state the work is passed on to *mergeSuff*.

```

mergeTokens :: Suffix -> Tokens -> Transition -> Tokens
mergeTokens suff1 toks2 trans2 = case view1 (currentSeq toks2) of
  token2 :< seq2' -> let newToken = mergeToken suff1 token2
                    in toks2 {currentSeq = newToken <| seq2'}
EmptyL -> case alex_accept ! out_state of
  [] -> let newSuff = mergeSuff suff1 (lastToken toks2) trans2
        in toks2 {lastToken = newSuff}
  acc -> let lex = suffToStr suff1 <
        suffToStr (lastToken toks2)
        in toks2 {lastToken = One $ createToken lex acc}
where out_state = outState toks2

```

Figure 4.8: *MergeTokens* function

4.5.4 Merging Suffixes

mergeSuff checks which pairs of suffixes it has and takes the appropriate actions.

```

mergeSuff :: Suffix -> Suffix -> Transition -> Suffix
mergeSuff (Multi toks1) suff2 trans2 = Multi $
  let newToks = combineWithRHS toks1 trans2
  in if isValid $ newToks
    then newToks
    else let newSuff = mergeSuff (lastToken toks1) suff2 trans2
        in toks1 {lastToken = newSuff}
mergeSuff (Str s1) suff2 _ = Str $ s1 < suffToStr suff2
mergeSuff (One token1) (Str s) trans2 =
  let toks2 = getTokens trans2 startState
  in if isValid toks2
    then Multi $ toks2 {currentSeq = token1 <| currentSeq toks2}
    else Multi $ createTokens (singleton token1) (Str s) (-1)
mergeSuff suff1 (One token2) _ = One $ mergeToken suff1 token2
mergeSuff suff1 (Multi toks2) trans2 =
  Multi $ mergeTokens suff1 toks2 trans2

```

Figure 4.9: *MergeSuff* function

- If the first suffix is of type *Multi* the function calls *combineWithRHS*. If the resulting tokens is invalid a recursive call is made with the suffix from the new tokens as first suffix.
- If the first suffix is of type *Str* the result will always be another *Str* no matter what is in the second suffix so the string is extracted and appended.
- if the first suffix is of type *One* and the second *Str* a new *Multi* suffix is created. A new second tokens is created using the start state on the second suffix, if this

results in a valid *Tokens*, the token from the first suffix is prepended. If it is not valid the *Str* is just added to the end of the new suffix.

- When both suffix are *One* they can be combined into a single token.
- When the first suffix is *One* and the second is *Multi* it is passed onto *mergeTokens*.

4.5.5 Append to Sequence of Tokens

appendTokens checks if there is a lexical error in *toks2*. if there is an error, that error is returned, otherwise *toks2* is appended to *toks1*.

```
appendTokens :: Seq IntToken -> Tokens -> Tokens
appendTokens seq1 toks2 | isValid toks2 =
    toks2 {currentSeq = seq1 <> currentSeq toks2}
    | otherwise = toks2
```

5

Result

The incremental lexer has three requirements, it should be robust, efficient and precise. Robustness means that the lexer does not 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 does not 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 text. The implementation this report propose is robust since it stores errors in the data structure rather than returning an error.

For it to be efficient the feedback to the user must be instant, or more formally the combination of two strings should be handled in $\Theta(\log(n))$ time.

Finally to be precise the lexer must give a correct result. This chapter is describing how these requirements are tested and what the results are.

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 was 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 part of a text and be able to combine two partial texts. If the lexical analysis of one partial text does not result in any legal tokens it must be able to be combined with other partial texts that makes it legal tokens. The lexical analysis of a text might not always result in the same tokens that the combination of the text with another text would give.

To test these cases a test was constructed which did a lexical analysis on two partial texts 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 entire text using the incremental lexer and the result using a sequential lexer.

It is not enough to test if the combination of two partial texts yields the same sequence of tokens as the text. To test that the result of the incremental lexer was the correct sequence of tokens, it was compared to what a sequential lexer generated. This comparison was an equality test that compared token for token that they were the same kind of token and had 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`.

Tests were performed on different files that were cut in different places when the update was tested. In all the tests the incremental lexer produced the same tokens as the sequential lexer. When these tests were done no text that would produce a lexical error was used. Some partial texts did produce lexical errors but the texts passed through the sequential lexer and compared to the result of the incremental lexer did not produce lexical errors.

5.2 Performance

All tests, where time performance was measured, were done with the help of `Criterion`, a Haskell library. `Criterion` has tools to ensure that the functions being tested is evaluated to normal form. `Criterion` runs the tests 100 times with a warm up run by default. The warm up run makes sure that all inputs to the functions being tested are evaluated before the actual testing begins, this ensures that nothing but the function being tested is measured [11].

To make sure that the time performance tests were not skewed under a certain system they were tested on different hardware and operating systems. The results were similar on all the systems tested. The results presented below were done on an intel i5 quad core at 2900MHz with 8GB memory under the linux Red Hat operating system.

To measure the performance of the incremental step two fingertrees were created, each representing one half of a text. By creating the two fingertrees the transition map for the code in those trees are created as well. The benchmarking was then done on the combination of the two trees. The results of the incremental lexer benchmarking suggested a running time of $\Theta(\log(n))$. To get a reference point the same text was lexed using a sequential lexer. The benchmarks can be found in fig. 5.3 and fig. 5.2. More details on the performance benchmarks can be found in appendix C.

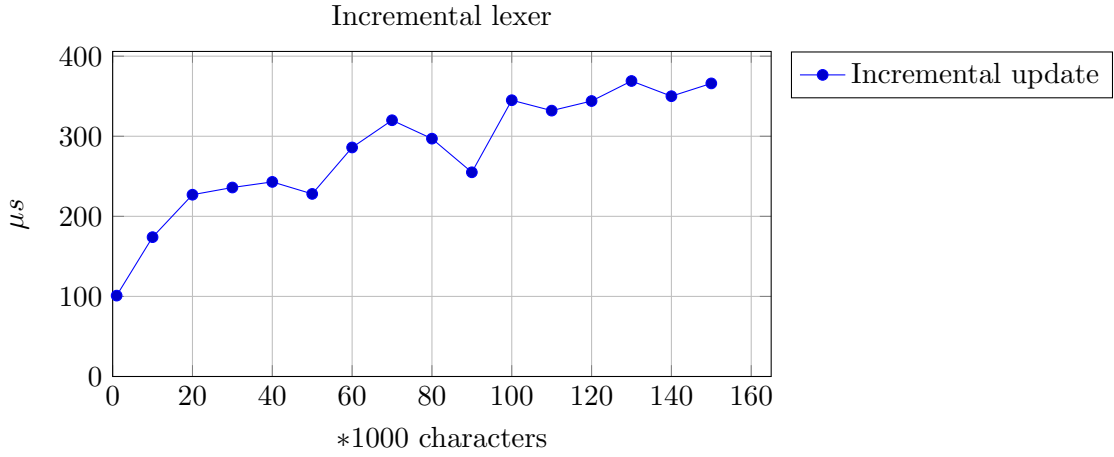


Figure 5.2: Benchmarking times of an incremental update

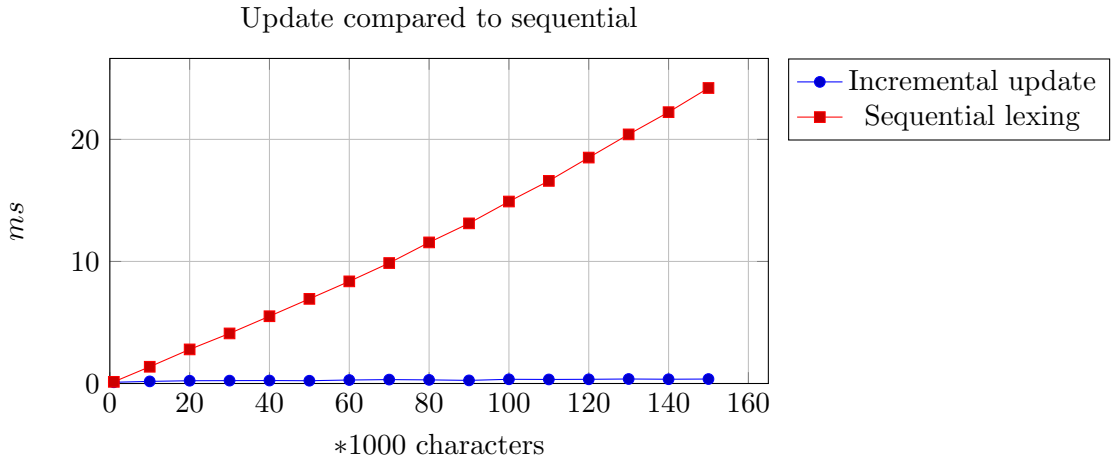


Figure 5.3: Comparison between an incremental update and sequential lexer

The running time when constructing the tree was not as fast as either the update or the sequential lexer. The tests for the running time, when constructing a new tree, suggested $\Theta(n \log(n))$, this was expected since there are $\log(n)$ updates for every character. The result can be found in fig. 5.4.

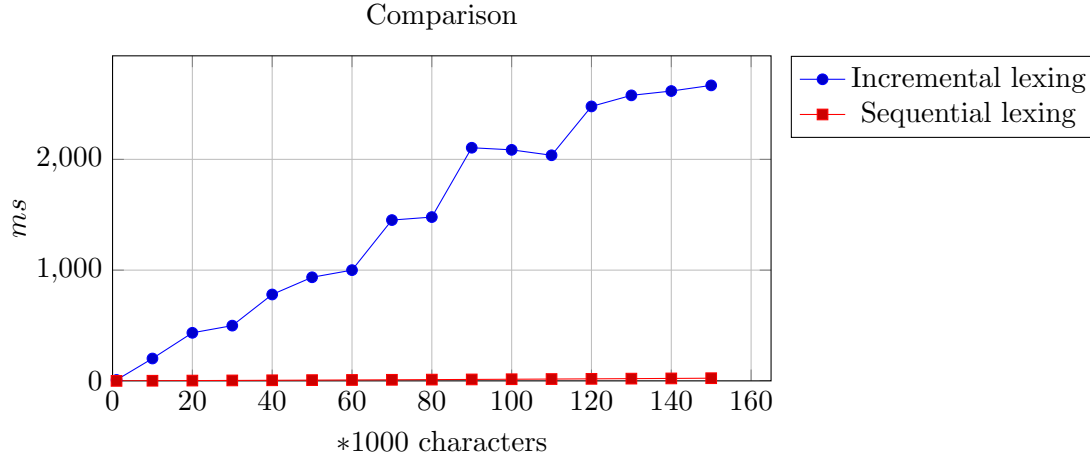


Figure 5.4: Comparison between the sequential lexer and the incremental lexer when lexing an entire file

Since the incremental lexer uses a tree structure there will be an increase in data-space used compared to a sequential lexer. Since the root of every partial tree stores the lexed tokens for that tree, every level of the tree will store approximately the size of the complete tokens for the entire text. A lexer that uses a DFA with m states will produce token data structure that grows with $\Theta(mn)$ since each in state will have the tokens for that state. The depth of the tree is $\log(n)$, this suggest a space complexity of $\Theta(mn \log(n))$.

To test the space complexity a DFA for an early Java version that has 90 states was used. To see how much the data structure grew the transition map was serialized and stored to the disc using the Haskell library *Data.Binary*. The test suggested that the size of the transition map grows linearly with the size of text being lexed, the results can be found in fig. 5.5.

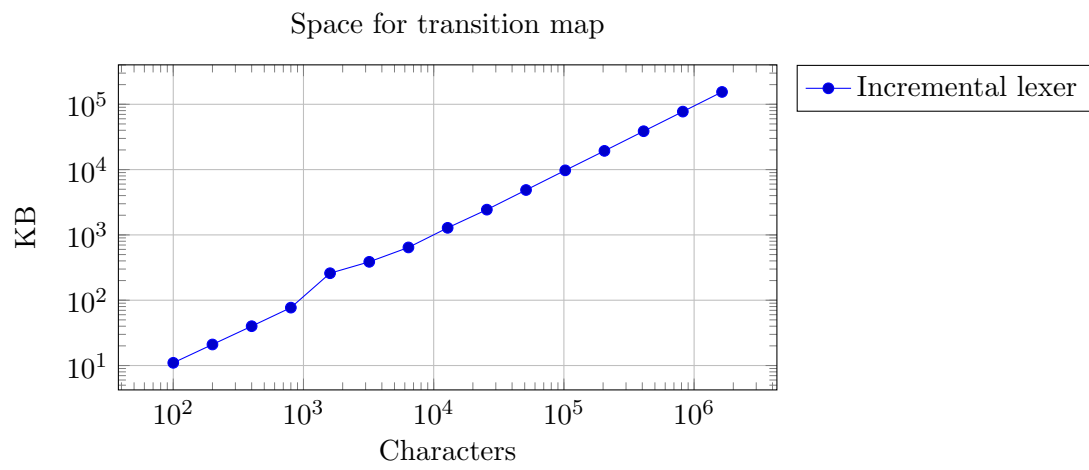


Figure 5.5: Space usage of the transition map using a DFA with 90 states

6

Discussion

During the course of the project there were some setbacks. The first setback was that our initial solution had bad running time, hard to understand and did not handle longest matching correctly.

After the first solution came a solution that was easier to understand but still had problem with the longest match. The running time was greatly improved from the first version but was still not faster than sequential lexers.

To solve this our last implementation which is described in chapter 4 made use of arrays and a DFA from Alex. The longest match problems that existed in the earlier versions of the lexer was mainly due to difficulty finding the correct solution.

6.1 Used Programming Language and Data Structure

The project is written in Haskell. One of the reasons is that similar research and projects have been done in Haskell, for instance [12] and [8]. Haskell also has the tools and the data structures used in the project. For instance Alex was used in two parts, first the DFA generated was used and second a lexer generated by Alex was used to get a comparison of the time for lexing.

There are other advantages of using Haskell as implementation language, namely higher order functions, lazy evaluation and function composition. Functional composition is useful in the lexical routines since the transitions are implemented as functions and are more or less just evaluated by composing. The lexical routines can be implemented as lazy, however this makes the incremental step in the lexer ineffective. Because of this the lexical routines are implemented strict, that is all values are always calculated.

The project could have been done in other languages. There are for example lexer generators in other languages, `lex`, `Flex` and `Jlex`, which can be used to create an efficient DFA. There have also been earlier articles that handle problems similar to this project written in Java [9].

The advantage of using finger trees in an incremental lexer is that it is easy to keep track of which tokens correspond to which part of the text, since the tokens are in the measure of the tree. Finger trees also has the advantage of keeping track of earlier result. When a tree is split up the tokens that match the partial texts are already calculated. The time complexity for combining two trees is low, $O(\log n)$, where n is the size of the smallest tree. The lexical routines revolve around combining two lexical results. Because of this it is advantageous to use finger trees since combination of two finger trees are fast. The main reason to not use any other type of tree is that another tree would not keep track of measure for the partial trees. In the case of splitting a tree the tokens would have to be recalculated.

6.2 Trials and Errors

The first version of the lexical routines only had the goal to get an idea of how a divide and conquer lexer could be implemented. As a result a lot of unnecessary information was stored and computed and some necessary information was stored and computed more than once. The solution was to calculate uncertain tokens until a satisfactory result was found or all possibilities was exhausted. This meant that for each combination the worst case was $O(2^n)$.

The next step was to make sure that the information was stored in the right places and that no unnecessary information was stored. To solve this an overhaul of the projects data structure was made. The result is close to the structure used in the final result. This solution still had some problem with the running time. The main reason was that the lexical result was not explicitly stored in the finger trees. That is, the functional composition was stored in the finger trees and since Haskell is lazy in nature the result wasn't computed resulting in slow running time when combining trees.

The final version of the project ensured that the lexer ran fast and made sure that fringe cases were computed correctly. The solution to this was to use arrays in the measure of the finger trees since arrays are always explicitly evaluated to normal form and that they have quick lookup time $O(1)$. Since the lexical routines can take advantage of the transitions in the form of functions that representation was used in the lexical routines and functions for converting between the arrays and function representation was implemented.

6.3 Implementation Suggestions

The updating step of the lexer is fast since the only computation needed is the combination of the last token of the first tree with the first token of the second tree and the combination of the actual trees. If an incremental lexer is used the fingertrees should be stored instead of the raw text. If the fingertrees are not stored the fingertree would need to be recalculated each time the file was opened.

7

Conclusion and Future Work

As mentioned in the result chapter the incremental lexer was both robust and precise. This means that without considering the time and space efficiency an incremental lexer will produce the same result as a sequential lexer with the difference being how lexical errors are handled. The incremental lexer is efficient in the sense that updates are done in $\Theta \log(n)$ time. However when a tree is built up from scratch the incremental lexer takes $\Theta n \log(n)$ time compared to the sequential lexer that takes Θn time.

7.1 Conclusions

Incremental lexers is not suited to be used in a stand alone lexer since a sequential lexer is more efficient then an incremental lexer when an entire text is being lexed. If a development environment that uses an incremental lexer was used, the stand alone lexer can be omitted since the tokens are already generated, saving one step in the compilation process.

It is however suited in an environment where updates are likely to happen, for example to give lexical feedback in a text editor where each key stroke would be an update. Insertion of a character in the text will be faster with an incremental lexer compared to a sequential lexer since the lexical analysis does not need to be done on the entire text. This means that the lexer could be run in real time without a user noticing it. The result from an incremental lexer can be passed to an incremental parser, giving parsing feedback to the user instead. However, loading times when opening files will be longer if the tree containing the tokens are not stored.

The space requirements for the incremental lexer grows with the tree. There is information for the entire text in all levels of the tree and each level has information for all

possible in states. The space of a tree grows with $\Theta mn \log(n)$, where m is the number of states in the DFA. This means that the memory usage will be big for large files and complex languages.

7.2 Future Work

To solve the problem with the space requirements for this implementation an implementation using sequence of characters could be used instead. That is, instead of using a character as the base case, a sequence of characters is used which is sequentially lexed, an example of a sequence is one line of code. This would shrink the tree from $\log(n)$ to $\log(n/x)$ where x is the mean length of a line. Since lines in the code is roughly of the same length there will be no impact on the worst case scenario time, for instance lexing 10 characters always takes the same amount of time. Since the lines in general are short updating a line will not take long time.

Another solution which could be used is to limit how big a tree can be. When that text is bigger then what fits in a tree, the tree is split into two trees. This will result in smaller trees at the expense of run time since the combination of the trees needs to be calculated on the fly.

In general a lot of in states will have the same sequence of tokens. The implementation suggested by this report will store all such sequences separately. An improvement would be if somehow the sequences of tokens that are identical could be stored in a separate table and the in state in the transition map points to the corresponding sequence for that in state. This would not improve the space complexity, but the practical space needed would shrink.

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] Jean-Philippe Bernardy and Koen Claessen. Efficient divide-and-conquer parsing of practical context-free languages. In *Proceeding of the 18th ACM SIGPLAN international conference on Functional Programming*, pages 111–122, 2013.
- [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. <http://www.haskell.org/alex/doc/html/index.html>, May 2014.
- [6] Markus Forsberg and Aarne Ranta. Bnf converter. In *Proceedings of the 2004 ACM SIGPLAN Workshop on Haskell*, Haskell '04, pages 94–95, New York, NY, USA, 2004. ACM.
- [7] Michael T Goodrich and Roberto Tamassia. *Data Structures and Algorithms in Java, 4th Edition*. John Wiley & Sons, 2005.
- [8] RALF HINZE and ROSS PATERSON. Finger trees: a simple general-purpose data structure. *Journal of Functional Programming*, 16:197–217, 3 2006.
- [9] Eugene Kirpichov. Incremental regular expressions. <http://jkff.info/articles/ire>, May 2014. English version of published russian article.
- [10] Chris Okasaki. *Purely Functional Data Structures*. Cambridge University Press, New York, NY, USA, 1999.

BIBLIOGRAPHY

- [11] Bryan O’Sullivan. The criterion package. <https://hackage.haskell.org/package/criterion>, May 2014.
- [12] Dan Piponi. Fast incremental regular expression matching with monoids. <http://blog.sigfpe.com/2009/01/fast-incremental-regular-expression.html>, January 2009.
- [13] Aarne Ranta and Markus Forsberg. *Implementing Programming Languages*, chapter Lexing and Parsing, pages 38–47. College Publications, London, 2012.
- [14] R.W. Sebesta. *Concepts of Programming Languages [With Access Code]*. Always learning. Pearson Education, Limited, 2012.
- [15] M. Sipser. *Introduction To The Theory Of Computation*. Advanced Topics Series. Thomson Course Technology, 2006.
- [16] 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.
- [17] L.G. Valiant. General context-free recognition in less than cubic time. *Journal of computer and system sciences*, 10(2):308–314, 1975.

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 analyzer will read “while” as an identifier, The syntactical analyzer will later determine if it is a loop. The expressions that matches rules without a name are discarded since they are not 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]

White space characters

Rules

white+

// [*'*]*

Single line comment

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

Multi line comment

Identifier → *letter ident**

Integer → *digit*+

Double → *digit*+ \. *digit*+

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

B

Incremental Lexer Source Code

Below follows the main part of the lexical routines and the data structures used in this project.

```
type State = Int
type Transition = State -> Tokens — Transition from in state to Tokens
data Tokens = NoTokens
             | InvalidTokens !(Seq Char)
             | Tokens { currentSeq :: !(Seq IntToken)
                       , lastToken  :: !Suffix
                       , outState   :: !State}
— The suffix is the sequence of as long as possible accepting tokens.
— It can itself contain a suffix for the last token.
               deriving Show
— This is either a Sequence of tokens or one token if it hits an accepting
— state with later characters
data Suffix  = Str !(Seq Char)
             | One !IntToken
             | Multi !Tokens
               deriving Show
type Size    = Sum Int
type LexTree = FingerTree (Table State Tokens,Size) Char
data IntToken = Token { lexeme    :: !(Seq Char)
                       , token_id :: Accepts}
type Accepts  = [AlexAcc (Posn -> Seq Char -> Token) ()]

tabulate :: (State,State) -> (State -> b) -> Table State b
access  :: Table State b -> (State -> b)
```

```

{— Functional Table variant
newtype Table a b = Tab {getFun :: a -> b}
tabulate - f = Tab f
access a x = (getFun a) x
—}

type Table a b = Array State b
tabulate range f = listArray range [f i | i <- [fst range..snd
    range]]
access a x = a ! x

instance Monoid (Table State Tokens) where
    mempty = tabulate stateRange (\_ -> emptyTokens)
    f 'mappend' g = tabulate stateRange $ combineTokens (access f)
        (access g)

— The base case for when one character is lexed.
instance Measured (Table State Tokens, Size) Char where
    measure c =
        let bytes = encode c
            cSeq = singleton c
            baseCase in_state | in_state == -1 = InvalidTokens cSeq
                               | otherwise = case foldl automata
                                   in_state bytes of
                -1 -> InvalidTokens cSeq
                os -> case alex_accept ! os of
                    [] -> Tokens empty (Str cSeq) os
                    acc -> Tokens empty (One (createToken cSeq acc)) os
            in (tabulate stateRange $ baseCase, Sum 1)

createToken :: (Seq Char) -> Accepts -> IntToken
createToken lex acc = Token lex acc

createTokens :: Seq IntToken -> Suffix -> State -> Tokens
createTokens seq suf state = if null seq
    then NoTokens
    else Tokens seq suf state

invalidTokens :: (Seq Char) -> Tokens
invalidTokens s = InvalidTokens s

emptyTokens :: Tokens
emptyTokens = NoTokens

————— Combination functions, the conquer step

— Combines two transition maps
combineTokens :: Transition -> Transition -> Transition

```

```

combineTokens trans1 trans2 in_state | isInvalid toks1 = toks1
                                     | isEmpty toks1   = trans2
                                     | in_state
                                     | otherwise =
                                     combineWithRHS toks1
                                     trans2

where toks1 = trans1 in_state

— Tries to merge tokens first, if it can't it either appends
  the token or calls
— itself if the suffix contains Tokens instead of a single
  token.
combineWithRHS :: Tokens -> Transition -> Tokens
combineWithRHS toks1 trans2 | isEmpty toks2 = toks1
                             | isValid toks2 =
    let toks2' = mergeTokens (lastToken toks1) toks2 trans2
    in appendTokens seq1 toks2'
                             | otherwise      = case lastToken
                             toks1 of

Multi suffToks ->
    let toks2' = combineWithRHS suffToks trans2 — try to
    combine suffix with transition
    in appendTokens seq1 toks2'
One tok -> appendTokens (seq1 |> tok) (trans2 startState)
Str s -> invalidTokens s
where toks2 = trans2 $ outState toks1
      seq1 = currentSeq toks1

— Creates one token from the last token of the first sequence
  and and the first
— token of the second sequence and inserts it between the init
  of the first
— sequence and the tail of the second sequence
mergeTokens :: Suffix -> Tokens -> Transition -> Tokens
mergeTokens suff1 toks2 trans2 = case view1 (currentSeq toks2) of
  token2 :< seq2' -> let newToken = mergeToken suff1 token2
                    in toks2 {currentSeq = newToken <| seq2'}
EmptyL -> case alex_accept ! out_state of
  [] -> toks2 {lastToken = mergeSuff suff1 (lastToken toks2)
              trans2}
  acc -> let lex = suffToStr suff1 <> suffToStr (lastToken
              toks2)
        in toks2 {lastToken = One $ createToken lex acc}
where out_state = outState toks2

— Creates on token from a suffix and a token
mergeToken :: Suffix -> IntToken -> IntToken
mergeToken suff1 token2 = token2 {lexeme = suffToStr suff1 <>
  lexeme token2}

```

```

— Creates the apropiet new suffix from two suffixes
mergeSuff :: Suffix -> Suffix -> Transition -> Suffix
mergeSuff (Multi toks1) suff2 trans2 = Multi $
  let newToks = combineWithRHS toks1 trans2
  in if isValid $ newToks
    then newToks
    else toks1 {lastToken = mergeSuff (lastToken toks1) suff2
      trans2}
mergeSuff (Str s1) suff2 _ = Str $ s1 <> suffToStr suff2
mergeSuff (One token1) (Str s) trans2 =
  let toks2 = trans2 startState
  in if isValid toks2
    then Multi $ toks2 {currentSeq = token1 <| currentSeq toks2}
    else Multi $ createTokens (singleton token1) (Str s) (-1)
mergeSuff suff1 (One token2) _ = One $ mergeToken suff1 token2
mergeSuff suff1 (Multi toks2) trans2 = Multi $ mergeTokens suff1
  toks2 trans2

— Prepends a sequence of tokens on the sequence in Tokens
appendTokens :: Seq IntToken -> Tokens -> Tokens
appendTokens seq1 toks2 | isValid toks2 =
  toks2 {currentSeq = seq1 <> currentSeq toks2}
  | otherwise = toks2

————— Constructors

makeTree :: String -> LexTree
makeTree = fromList

measureToTokens :: (Table State Tokens, Size) -> Seq Token
measureToTokens m = case access (fst $ m) startState of
  InvalidTokens s -> error $ "Unacceptable_token:_" ++ toList s
  NoTokens -> empty
  Tokens seq suff out_state ->
    snd $ foldlWithIndex showToken (Pn 0 1 1, empty) $ intToks
    seq suff
where showToken (pos, toks) _ (Token lex accs) =
  let pos' = foldl alexMove pos lex
  in case accs of
    [] -> (pos', toks)
    AlexAcc f:_ -> (pos', toks |> f pos lex)
    AlexAccSkip:_ -> (pos', toks)
  intToks seq (Str str) = error $ "Unacceptable_token:_"
    ++ toList str
  intToks seq (One token) = seq |> token
  intToks seq (Multi (Tokens seq' suff' _)) = intToks (seq
    <> seq') suff'

```

```

treeToTokens :: LexTree -> Seq Token
treeToTokens = measureToTokens . measure

----- Util funs

isValid :: Tokens -> Bool
isValid (Tokens _ _ _) = True
isValid _ = False

isEmpty :: Tokens -> Bool
isEmpty NoTokens = True
isEmpty _ = False

isInvalid :: Tokens -> Bool
isInvalid (InvalidTokens _) = True
isInvalid _ = False

suffToStr :: Suffix -> Seq Char
suffToStr (Str s) = s
suffToStr (One token) = lexeme token
suffToStr (Multi toks) =
  concatLexemes (currentSeq toks) <> suffToStr (lastToken toks)

isAccepting :: Tokens -> Bool
isAccepting (Tokens _ suff _) = case suff of
  Str _ -> False
  One _ -> True
  Multi toks -> isAccepting toks
isAccepting NoTokens = True
isAccepting (InvalidTokens _) = False

concatLexemes :: Seq IntToken -> Seq Char
concatLexemes = foldr ((<>) . lexeme) mempty

insertAtIndex :: String -> Int -> LexTree -> LexTree
insertAtIndex str i tree =
  if i < 0
  then error "index_must_be_>=0"
  else l <> (makeTree str) <> r
    where (l,r) = splitTreeAt i tree

splitTreeAt :: Int -> LexTree -> (LexTree, LexTree)
splitTreeAt i tree = split (\(_,s) -> getSum s>i) tree

size :: LexTree -> Int
size tree = getSum . snd $ measure tree

— Starting state
startState = 0

```



```

— A tuple that says how many states there are
stateRange = let (start,end) = bounds alex_accept
               in (start-1,end)

— Takes an in state and a byte and returns the corresponding
  out state using
— the DFA generated by Alex
automata :: Int -> Word8 -> Int
automata (-1) _ = -1
automata s c = let base    = alex_base ! s
                  ord_c    = fromEnum c
                  offset   = base + ord_c
                  check    = alex_check ! offset
               in if (offset >= (0)) && (check == ord_c)
                  then alex_table ! offset
                  else alex_deflt ! s

```

C

Benchmarks

Below follows three benchmarks of different code sizes. All three benchmarks originate from the same code base, but the first is bench is on a code 10 times the origin code, the second is 100 times the origin and the third is 1000 times the origin. The code given as a string kept in the memory, that is the code is not read from a file.

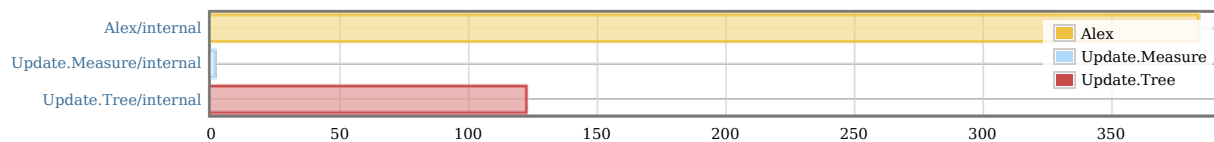
These test are done a system with the following specification:

Thing	Spec
-------	------

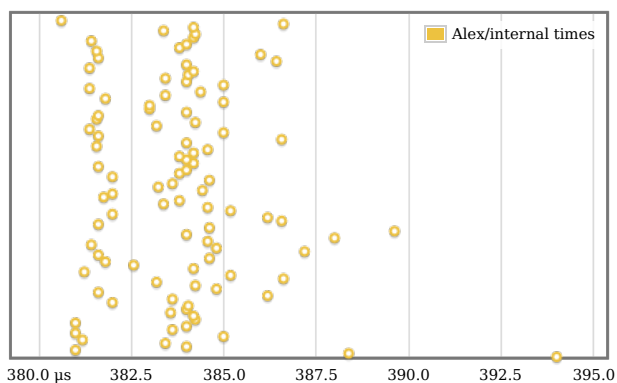
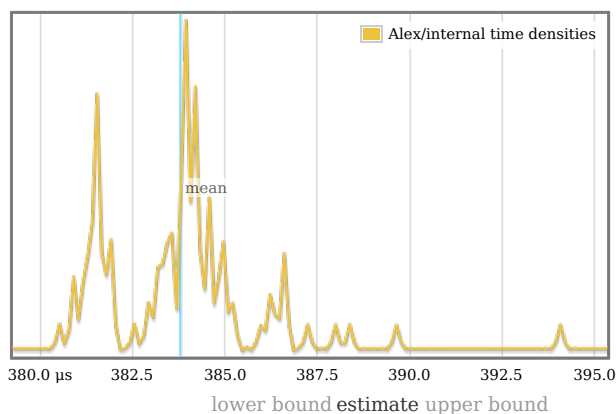
There are benchmarks done on other systems as well and they follow the same trend.

Performance Measurements (initial code * 10)

Overview

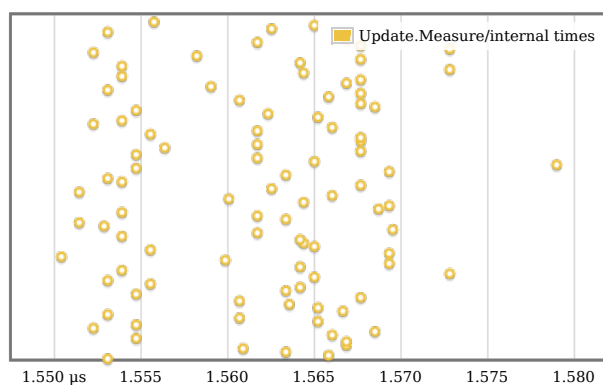
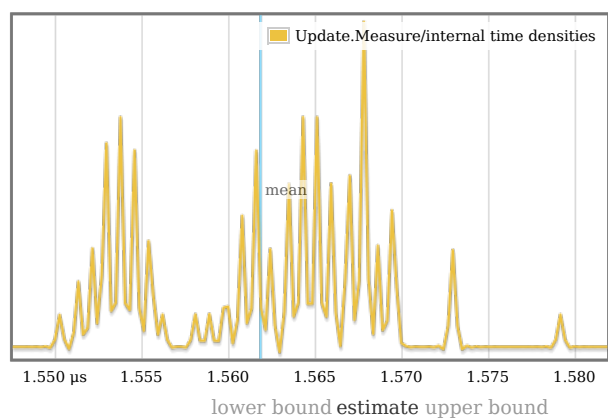


Alex - internal code



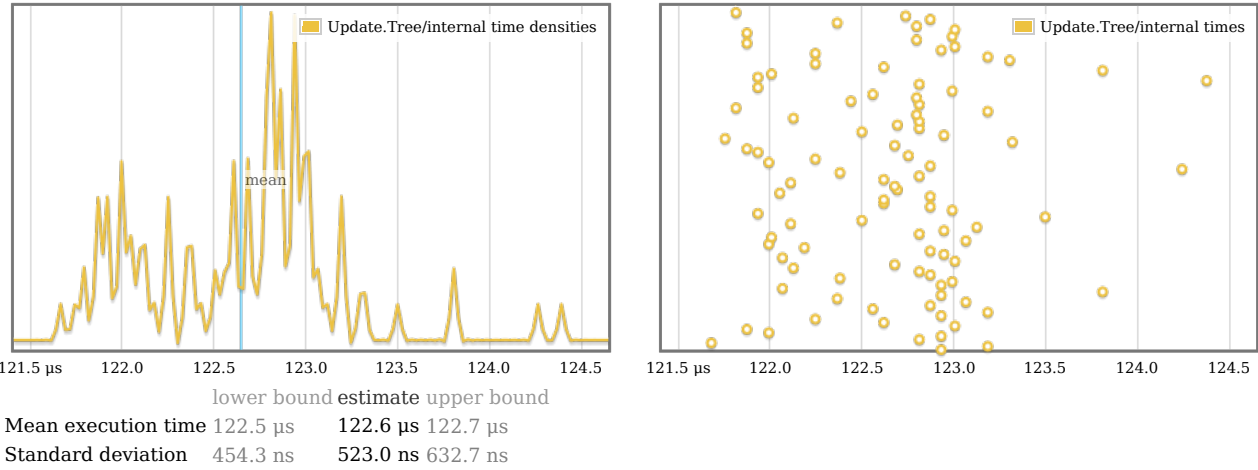
Outlying measurements have no (1.0%) effect on estimated standard deviation.

Update Measure - internal code



Outlying measurements have no (1.0%) effect on estimated standard deviation.

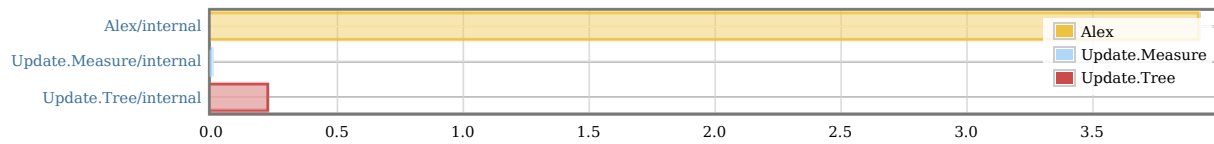
Update FingerTree - internal code



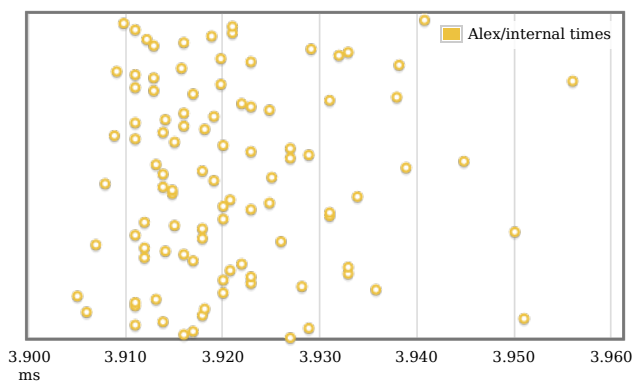
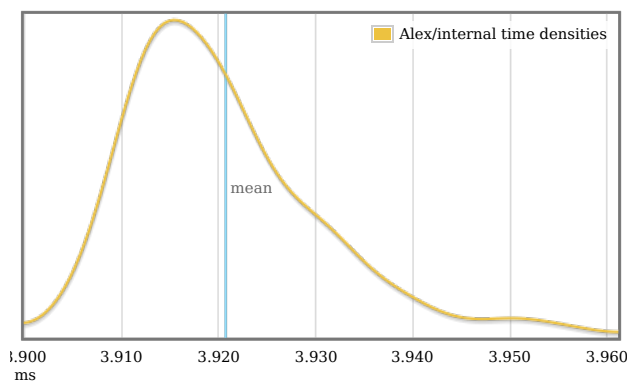
Outlying measurements have no (1.0%) effect on estimated standard deviation.

Performance Measurements (initial code * 100)

Overview



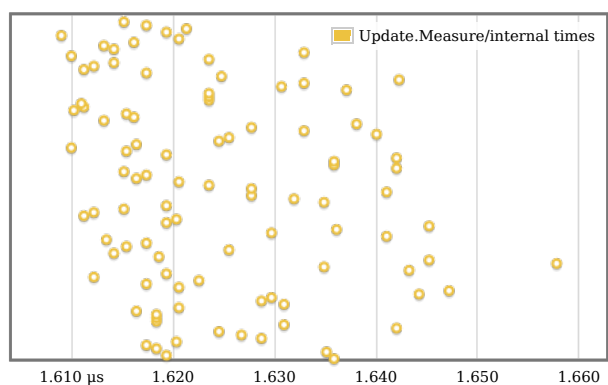
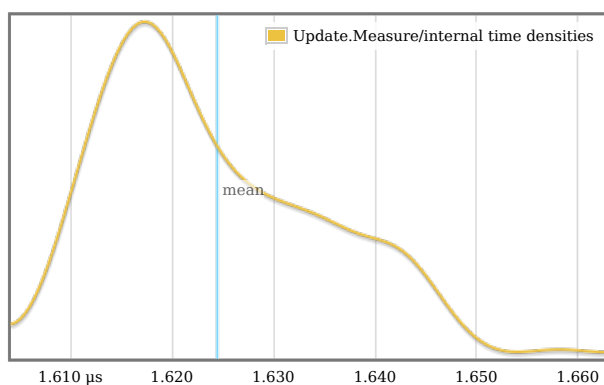
Alex - internal code



Mean execution time 3.918 ms 3.920 ms 3.922 ms
Standard deviation 8.489 μ s 10.07 μ s 12.27 μ s

Outlying measurements have no (1.0%) effect on estimated standard deviation.

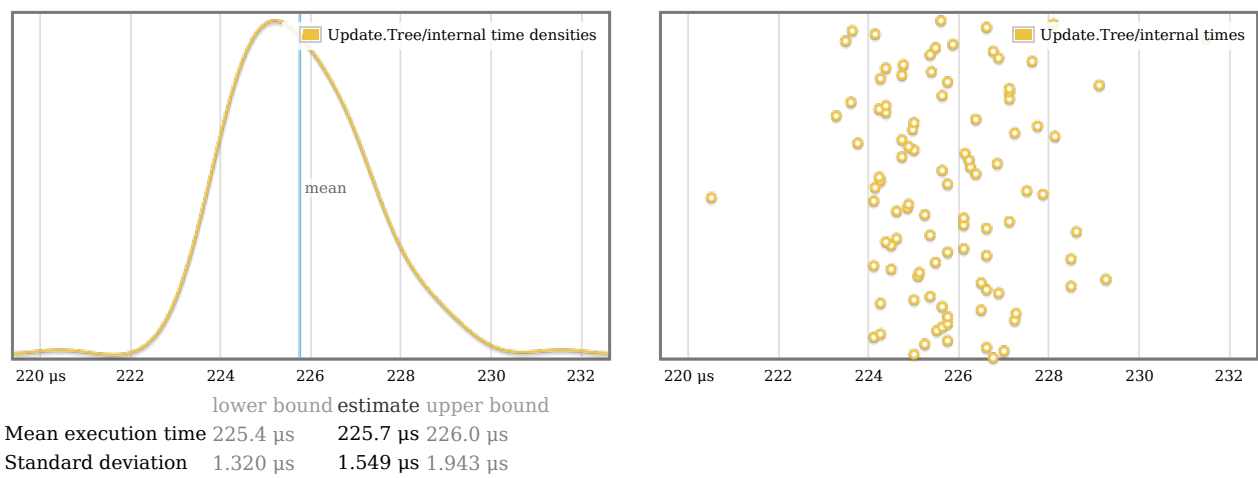
Update Measure - internal code



Mean execution time 1.622 μ s 1.624 μ s 1.626 μ s
Standard deviation 9.517 ns 10.68 ns 12.35 ns

Outlying measurements have no (1.0%) effect on estimated standard deviation.

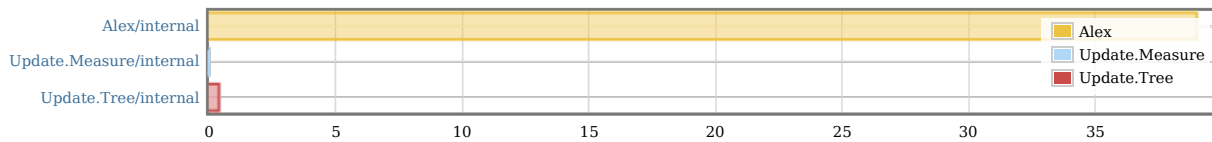
Update FingerTree - internal code



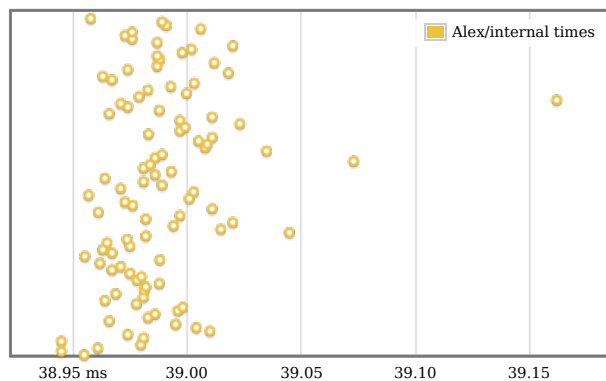
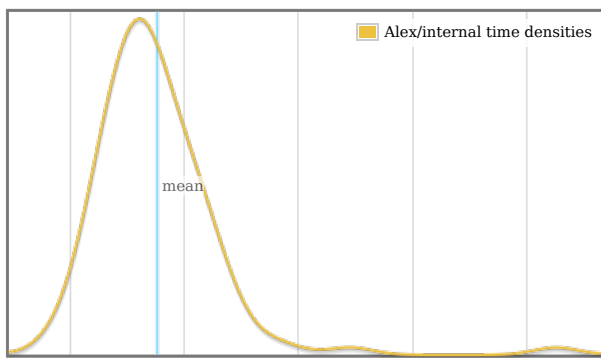
Outlying measurements have no (1.0%) effect on estimated standard deviation.

Performance Measurements (initial code * 1000)

Overview



Alex - internal code



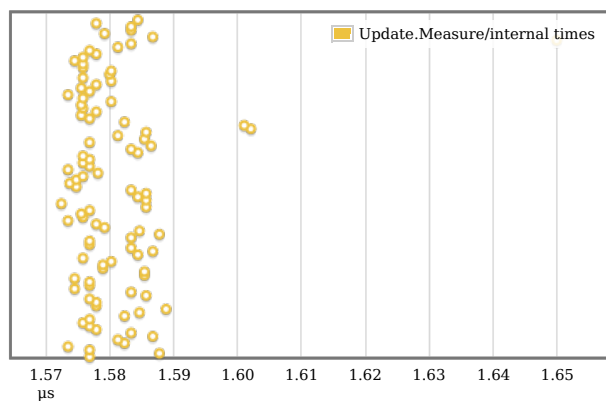
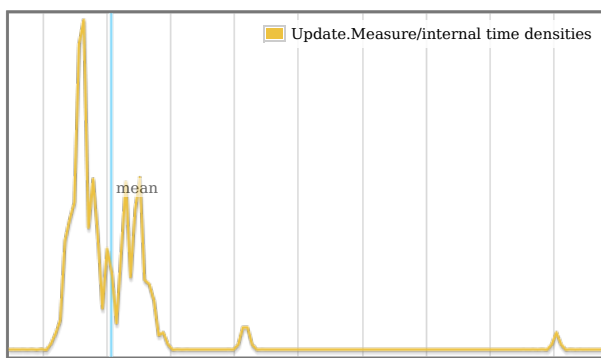
lower bound estimate upper bound

Mean execution time 38.98 ms 38.98 ms 38.99 ms

Standard deviation 19.40 μ s 27.17 μ s 45.78 μ s

Outlying measurements have no (1.0%) effect on estimated standard deviation.

Update Measure - internal code



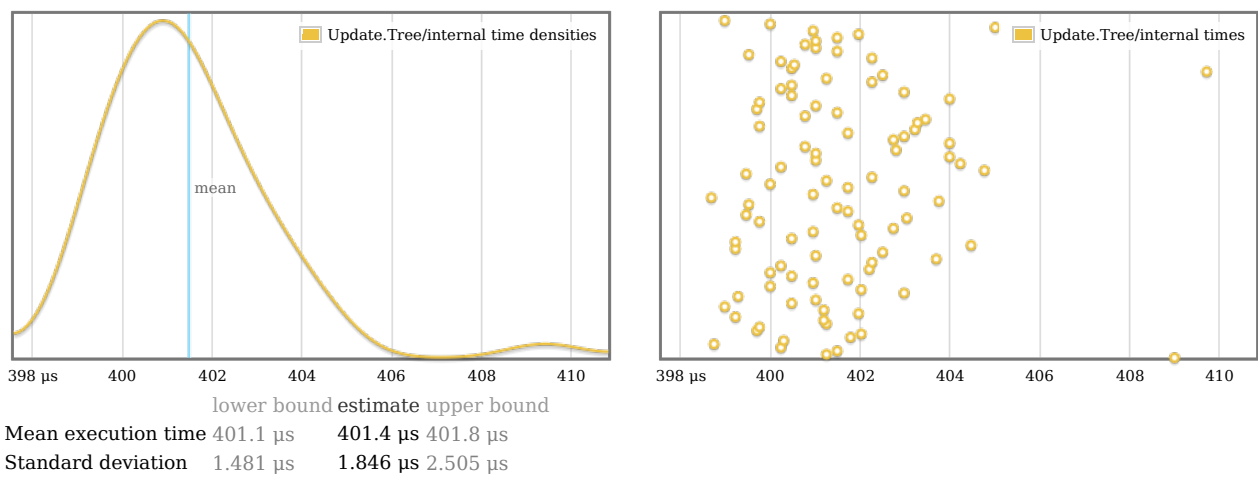
lower bound estimate upper bound

Mean execution time 1.579 μ s 1.580 μ s 1.583 μ s

Standard deviation 4.833 ns 8.727 ns 17.44 ns

Outlying measurements have no (1.0%) effect on estimated standard deviation.

Update FingerTree - internal code



Outlying measurements have no (1.0%) effect on estimated standard deviation.

Understanding the Graphs

In this report, each function benchmarked by criterion is assigned a section of its own. In each section, we display two charts, each with an x axis that represents measured execution time. These charts are active; if you hover your mouse over data points and annotations, you will see more details.

- The chart on the left is a [kernel density estimate](#) (also known as a KDE) of time measurements. This graphs the probability of any given time measurement occurring. A spike indicates that a measurement of a particular time occurred; its height indicates how often that measurement was repeated.
- The chart on the right is the raw data from which the kernel density estimate is built. Measurements are displayed on the y axis in the order in which they occurred.

Under the charts is a small table displaying the mean and standard deviation of the measurements. We use a statistical technique called the [bootstrap](#) to provide confidence intervals on our estimates of these values. The bootstrap-derived upper and lower bounds on the mean and standard deviation let you see how accurate we believe those estimates to be. (Hover the mouse over the table headers to see the confidence levels.)

A noisy benchmarking environment can cause some or many measurements to fall far from the mean. These outlying measurements can have a significant inflationary effect on the estimate of the standard deviation. We calculate and display an estimate of the extent to which the standard deviation has been inflated by outliers.

These reports were created using the [criterion](#) benchmark execution and performance analysis tool. Criterion is developed and maintained by [Bryan O'Sullivan](#)