



UNIVERSITY OF GOTHENBURG

Implementing incremental and parallel parsing

A subtitle that can be rather long

*Master of Science Thesis in Computer Science*

TOBIAS OLAUSSON

University of Gothenburg  
Chalmers University of Technology  
Department of Computer Science and Engineering  
Göteborg, Sweden, April 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.

## **Implementing incremental and parallel parsing**

A subtitle here that can be quite long

TOBIAS OLAUSSON

© TOBIAS OLAUSSON, April 2014

Examiner: PATRIK JANSSON

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

Cover: an image that is used as a cover image

Department of Computer Science and Engineering  
Göteborg, Sweden, April 2014

## **Abstract**

This is an abstract

# Contents

<b>1</b>	<b>Background</b>	<b>3</b>
1.1	Introduction . . . . .	3
1.1.1	Divide-and-conquer . . . . .	3
1.1.2	Incrementality . . . . .	3
1.1.3	Parallelism . . . . .	3
1.1.4	Parsing . . . . .	3
1.1.5	Motivation . . . . .	4
1.2	Lexing . . . . .	4
1.3	Context-free grammars . . . . .	4
1.3.1	Backus-Naur Form . . . . .	4
1.3.2	Chomsky Normal Form . . . . .	5
1.4	Parsing . . . . .	5
1.4.1	CYK algorithm . . . . .	5
1.4.2	Valiant . . . . .	5
1.4.3	Improvement by Bernardy & Claessen . . . . .	5
1.5	Dependently typed programming . . . . .	5
1.5.1	Kinds, Types and Values . . . . .	5
<b>2</b>	<b>Implementation</b>	<b>6</b>
2.1	Finger trees . . . . .	6
2.2	Measuring and Monoids . . . . .	6
2.2.1	Pipeline of measures . . . . .	6
2.3	Lexing . . . . .	6
2.3.1	LexGen – Alex discrepancy . . . . .	6
2.3.2	Position information . . . . .	6
2.4	Parsing . . . . .	6
2.4.1	BNFC . . . . .	6
2.4.2	Dependently typed programming with charts . . . . .	6
2.4.3	Oracle and unsafePerformIO . . . . .	6

<b>3</b>	<b>Results</b>	<b>7</b>
3.1	Final product . . . . .	7
3.1.1	Testing . . . . .	7
3.2	Measurements . . . . .	7
<b>4</b>	<b>Discussion</b>	<b>8</b>
4.1	Pitfalls . . . . .	8
4.1.1	Too many result branches . . . . .	8
4.1.2	Position information . . . . .	8
4.2	Future work . . . . .	8

# Chapter 1

## Background

### 1.1 Introduction

The topic of this thesis is to do **parsing** in an **incremental** fashion that can easily be **parallelizable**, using a **divide-and-conquer approach**. In this section, I will give a brief explanation of the topics covered, and end with a motivation for why this is interesting to do in the first place.

#### 1.1.1 Divide-and-conquer

Add stuff here from section 2 in parparse paper.

#### 1.1.2 Incrementality

Doing something incrementally means that one does it step by step, and not longer than necessary.

#### 1.1.3 Parallelism

Why is parallelism interesting and how does it apply in this case?

#### 1.1.4 Parsing

To parse is a to check if some given input corresponds to a certain language's grammar, and in this thesis it will use **context-free grammars** for programming languages.

### 1.1.5 Motivation

In compilers, lexing and parsing are the two first phases. The output of these is an abstract syntax tree (AST) which is fed to the next phase of the compiler. But an AST could also provide useful feedback for programmers, already in their editor, if the code could be lexed and parsed fast enough. With a lexer and parser that is incremental and that can also be parallelized could real-time feedback in the form of an AST easily be provided to the programmer. Most current text editors give syntax feedback based on regular expressions, which does not yield any information about depth or the surrounding AST.

Something something about connecting to a type-checker.

## 1.2 Lexing

\* Describe what lexing is \* Shortly describe LexGen and its relevance.

## 1.3 Context-free grammars

Give a light-weight description here.

Formally, a context-free grammar is a 4-tuple:  $G = (V, \Sigma, P, S)$ .  $V$  is a set of non-terminals, or variables.  $\Sigma$  is the set of terminal symbols, describing the content that can be written in the language.  $P$  is a

### 1.3.1 Backus-Naur Form

Context-free grammars are often used to describe the syntax of programming languages. Such descriptions are often given in a **labelled Backus-Naur form**, where each rule is written on the following form:

$$\textit{Label.Variable} ::= \textit{Production}$$

Grammars written in Labelled Backus-Naur Form are also used in the **BNF Converter (BNFC)**, a lexer and parser generator tool developed at Chalmers. Given such a grammar, BNFC generates, among other things, a lexer and a parser, implemented in one of several languages, for the language described in that grammar. According to documentation, usage of BNFC saves approx 90% of source code in writing a compiler front-end.

### 1.3.2 Chomsky Normal Form

Chomsky Normal Form (CNF) is a subset of context-free grammars that was first described by linguist Noam Chomsky. Productions in CNF are restricted to the following forms:

$A \rightarrow BC$ ,     $A$  variable,  $B$  and  $C$  productions  
 $A \rightarrow a$ ,     $A$  variable,  $a$  terminal symbol

Figure 1.1: Rules allowed in Chomsky Normal Form

Since grammars in CNF are restricted to branches or single terminal symbols, they are well suited for usage in divide-and-conquer algorithms. There are several existing algorithms to convert context-free grammars into CNF.

## 1.4 Parsing

### 1.4.1 CYK algorithm

### 1.4.2 Valiant

### 1.4.3 Improvement by Bernardy & Claessen

Running time analysis. Oracle for list.

## 1.5 Dependently typed programming

What is dependently typed programming, and how can it be used in Haskell. How about an example using vectors (standard example, sort of).

### 1.5.1 Kinds, Types and Values

Describe how kinds relate to types as types relate to values.



# Chapter 2

## Implementation

### 2.1 Finger trees

What are they?

### 2.2 Measuring and Monoids

#### 2.2.1 Pipeline of measures

An illustration would be good here

### 2.3 Lexing

#### 2.3.1 LexGen – Alex discrepancy

#### 2.3.2 Position information

### 2.4 Parsing

#### 2.4.1 BNFC

#### 2.4.2 Dependently typed programming with charts

#### 2.4.3 Oracle and unsafePerformIO

# Chapter 3

## Results

### 3.1 Final product

#### 3.1.1 Testing

### 3.2 Measurements

How fast is it? What is the complexity?

# Chapter 4

## Discussion

### 4.1 Pitfalls

Look through the LOG to remember whatever happened. Describe sort of chronologically?

#### 4.1.1 Too many result branches

Describe the reason for this, with the merge stuff.

#### 4.1.2 Position information

\* Tuple of monoids? \* RingP instance? \* Newtype for tuples?

### 4.2 Future work