

# ECE491 Advanced compilers final report

## Implementing a lazy functional language

Jonathan Lam

2022/05/12

### Contents

<b>1</b>	<b>Project overview</b>	<b>3</b>
1.1	Motivation and overview . . . . .	3
1.2	Commentary on the tutorial . . . . .	3
1.3	Implementation setup . . . . .	4
<b>2</b>	<b>Background</b>	<b>5</b>
2.1	Definition of a programming language . . . . .	5
2.2	Implementations of programming language . . . . .	5
2.3	The untyped $\lambda$ -calculus . . . . .	5
2.4	Functional programming . . . . .	5
<b>3</b>	<b>The Core language</b>	<b>6</b>
3.1	Terminology . . . . .	6
3.2	Syntax . . . . .	6
3.3	Dynamic semantics . . . . .	6
3.4	Sample programs . . . . .	6
3.5	Lexer . . . . .	6
3.6	Parser . . . . .	6
3.7	Pretty-print utility . . . . .	6
<b>4</b>	<b>Template instantiation (TI) evaluator</b>	<b>7</b>
4.1	The TI abstract machine . . . . .	7
4.1.1	Instantiation . . . . .	7
4.1.2	Unwinding . . . . .	7
4.1.3	Memory model . . . . .	7
4.2	Sample evaluations . . . . .	7
<b>5</b>	<b>G-Machine (GM) compiler</b>	<b>8</b>
5.1	The GM abstract machine . . . . .	8
5.1.1	List of opcodes . . . . .	8
5.2	Compilation schemes . . . . .	9

5.3	Evaluator . . . . .	9
5.4	Sample compilations and evaluations . . . . .	9
<b>6</b>	<b>Future work</b>	<b>10</b>
<b>7</b>	<b>Conclusions</b>	<b>11</b>
<b>8</b>	<b>References</b>	<b>12</b>

	Purely functional	Non-purely functional
Untyped	Core	Scheme <sup>1</sup>
Typed	Hazel <sup>2</sup>	C <sup>3</sup>

Table 1: Summary of programming language implementation projects

## 1 Project overview

### 1.1 Motivation and overview

This project was the semester-long project for the ECE491 independent study on advanced compilers. It is the implementation of a simple lazy (normal-order) evaluation functional language similar to Haskell. The work follows the tutorial, “Implementing functional languages: a tutorial” by Simon Peyton Jones at Microsoft Research [3]. This was published two years after Haskell 1.0 was defined in 1990 [1], to which SPJ was a major contributor.

This project is the culmination of my studies in programming languages and functional programming over the past two years. These studies include writing a C compiler in C (Spring 2021), a Scheme (LISP) interpreter in Scheme (Summer 2021), and this project, a Core interpreter and compiler in Haskell (Spring 2022). Additionally, my Master’s thesis (2021-2022 school year) was about updating the evaluation model of Hazel, a programming language implemented in OCaml. We summarize these languages in Table 1.

### 1.2 Commentary on the tutorial

The tutorial used for this project assumes a basic understanding of functional programming and a non-strict language such as Miranda or Haskell. For this report, we do not assume this and give a brief introduction to these ideas.

The text is largely in tutorial format, in that it provides much of the code for the reader to follow along with, but it is also in large part similar to a textbook. While much of the base code is given, much of the implementation is left in the form of non-trivial exercises to be completed using the provided theory. As the book progresses, less code is provided directly; rather, the semantics are given using the state transition notation, and the reader is asked to implement this in code.

The structure of the tutorial goes as follows: Chapter 1 introduces the Core language. Chapter 2 provides an evaluator implementation called the Template Instantiation (TI) evaluator, which we will describe in Section 4. Chapter 3 provides a compiled version of the TI evaluator, called the G-Machine<sup>4</sup> (GM), which we will describe in Section 5. We note that Haskell officially uses a Spineless Tagless G-Machine (STG) [2].

---

<sup>1</sup>Mostly functional

<sup>2</sup>Gradually-typed

<sup>3</sup>Weakly-typed

<sup>4</sup>“G” for “graph,” most likely.

Chapter 4 describes the Three-Address Machine compiled implementation, and Chapter 5 describes the Parallel G-Machine. The G-Machine (a stack-based machine) may be translated into a TAM representation. Chapter 6 introduces the  $\lambda$ -abstraction syntax using the  $\lambda$ -lifting program transformation. Chapters 4-6 were not covered for this independent study.

Rather than developing each implementation vertically (e.g., finishing the lexer, then the parser, then the intermediate representation, then exporting opcodes), the tutorial uses a horizontally incremental development method. In this style of development, we first complete a base implementation of only the core language features, and then incrementally add support for new language features. The tutorial calls these horizontal versions “marks,” i.e., “TI Mark 3” or “GM Mark 7.”

I greatly enjoyed the tutorial and found the exercises delightfully challenging and enlightening. My only complaints with the tutorial are that sometimes the exercises seemed to require knowledge from future marks, and that the names of certain concepts seem to be arbitrarily named<sup>5</sup>.

### 1.3 Implementation setup

My implementation is called `flc`, short for “fun(ctional) lazy compiler.” `flc` is written in Haskell, similar to the tutorial<sup>6</sup>. Instructions for use may be found on the GitHub’s README.

The implementation may be found at `jlam55555/fun-lazy-compiler`. A transpiler for the GM Mark 1 to x86-64 assembly code may be found at `jlam55555/flc-transpiler`.

---

<sup>5</sup>For example, the compilation schemes are named with arbitrary letters, and the “G” in “G-Machine” and the “I” in “Iseq” are never explained.

<sup>6</sup>The tutorial says it was written in Miranda, but I am not sure if this is correct. All of the code samples run successfully in Haskell, and I believe that Miranda’s syntax is somewhat different.

$$e ::= \lambda x. e \mid x \mid e e$$

Figure 1: Grammar of  $\Lambda$

## 2 Background

This section provides brief background on programming language theory from a functional perspective.

### 2.1 Definition of a programming language

*Interfaces* are necessary for efficient and effective communication. A *programming language* serves as an interface between humans and computers. To rigorously work with a programming language, we define its *syntax* and *semantics*.

The syntax of a programming language describes the way valid expressions and programs are formed. It is specified using a *grammar*. The grammar for the untyped  $\lambda$ -calculus  $\Lambda$  is shown in Figure 1.

The *semantics* of a programming language describes its behavior. The *static semantics* denotes the behavior of processes that happen prior to evaluation, such as type checking. The *dynamic semantics* describes the behavior of *evaluation*. Evaluation is the process of reducing an expression down to a *value*, or irreducible expression.

For this project, we actually define two languages: the Core language, and the abstract

### 2.2 Implementations of programming language

TODO: compilers, interpreters

### 2.3 The untyped $\lambda$ -calculus

### 2.4 Functional programming

TODO: notation with spaces for function application and curried-by-default application

TODO: strict languages like the ML family

TODO: haskell and miranda

TODO: pure functional programming and lack of side effects

TODO: algebraic datatypes

## 3 The Core language

### 3.1 Terminology

TODO: lazy evaluation (call-by-value, call-by-name, call-by-need, applicative-order, normal-order), supercombinators, currying, lambda-abstraction, function application, ADTs (product and sum types) and structured data, whnf, standard representations (booleans, list constructors), standard definitions/prelude (e.g., if, simple combinators, standard representations), scrutinee, definiens

### 3.2 Syntax

### 3.3 Dynamic semantics

TODO: function application and basic forms have the same dynamic semantics as the lambda calculus, except using an environment to store variable bindings rather than using substitution

### 3.4 Sample programs

TODO: show basic forms

    TODO: twice twice twice

    TODO: show algebraic datatypes

### 3.5 Lexer

TODO: very simple tokenization

### 3.6 Parser

TODO: build up LL(1) parser from useful subparsers

### 3.7 Pretty-print utility

TODO: constant-time append with indentation, only linear-time rendering

## 4 Template instantiation (TI) evaluator

TODO: only cover this briefly, relation to the GM machine

### 4.1 The TI abstract machine

TODO: process: instantiate and unwind

#### 4.1.1 Instantiation

TODO: conditionals

#### 4.1.2 Unwinding

#### 4.1.3 Memory model

### 4.2 Sample evaluations

TODO: sample updating

TODO: sample lazy evaluation

TODO: tco

TODO: arithmetic

## 5 G-Machine (GM) compiler

### 5.1 The GM abstract machine

#### 5.1.1 List of opcodes

**TODO:** list each opcode, describe it, show its behavior using the state machine

**Pushglobal  $f$**  Push the address of the global (supercombinator)  $f$  onto the stack<sup>7</sup>.

**Pushint  $n$**  Push the address of an integer node representing the integer  $n$  onto the stack. If such an integer node representing the integer  $n$  already exists, this may use that existing address; otherwise, this will allocate a new integer node  $n$ .

**Push  $n$**  Push the  $n$ -th address of the stack onto the stack.

**Update  $n$**  Set the node pointed to by the  $n$ -th address of the stack to an indirection node pointing to the top element of the stack, then pop one element from the stack.

**Pop  $n$**  Pop  $n$  elements from the stack.

**Alloc  $n$**  Push  $n$  invalid addresses onto the stack, or decrease the stack pointer by  $n$ <sup>8</sup>.

**Slide  $n$**  “Slide” the top element of the stack up  $n$  slots. In other words, the  $(n + 1)$ -th element of the stack will be replaced with the top of the stack, and then  $n$  elements will be popped.

**Unwind** Same as in the TI implementation Section 4.1.2.

**Mkap** Pop two elements from the stack, allocate a new application node (with the first popped element in function position, and the second popped element in argument position) and push the application node onto the stack.

**Eval** Strictly evaluate the element on the top of the stack to WHNF on a new stack (frame)<sup>9</sup>.

**Add, Sub, Mul, Div, Neg** Arithmetic binary and unary operators. Assumes the two elements on the top of the stack are evaluated to integers (otherwise will crash the program). Pops the top two elements and leaves a number node on the top of the stack representing the arithmetic result.

---

<sup>7</sup>This is overloaded in Mark 6 for a technicality regarding unsaturated data constructors. This is Exercise 3.38.

<sup>8</sup>This is used for **letrec** expressions. These addresses initially point to an invalid node, and will be updated by the evaluation of the definiens.

<sup>9</sup>This pushes a new stack on the dump and evaluates the node on the new stack. This is more or less equivalent to a function call in the conventional procedural ABIs.



**Eq, Ne, Lt, Le, Gt, Ge** Relational binary operators. Assumes the two elements on the top of the stack are evaluated to booleans in the standard representation (otherwise will crash the program). Pops the top two elements and leaves a node on the top of the stack representing the boolean result in the standard representation.

**Pack  $t\ n$**  Pops  $n$  addresses from the stack. Allocate and push a structured data node representing structured data with tag  $t$  and the  $n$  popped arguments.

**Casejump  $rules$**   $rules$  is of the form  $\{t_j \rightarrow [i_1, i_2, \dots]\}$ : a mapping of tags  $t_j$  to code sequences  $[i_1, i_2, \dots]$ . Assumes the top of the stack is the scrutinee of a **case** expression, and thus a data node with tag  $t$ . Appends the instruction list of the appropriate rule onto the current instruction queue. Throws an error if the scrutinee is not a structured data node or if there is no match.

**Split  $n$**  Assumes the top of the stack is a structured data node. Pops it from the stack, and then pushes its  $n$  arguments onto the stack.

**Print** Output the top element on the stack. Assumes that it is a data node (integer or structured data node). Recursively generates additional **Print** and **Eval** opcodes for structured data.

## 5.2 Compilation schemes

**TODO:** list each compilation scheme, describe it, show its behavior on expression types

## 5.3 Evaluator

**TODO:** the evaluator performs the opcode's descriptions; maybe show its behavior using state machine here

## 5.4 Sample compilations and evaluations

## 6 Future work

TODO: garbage collection

TODO: last marks of TI and GM

TODO: 3-address machine and parallel g-machine

TODO: lambda lifting

TODO: study implementation of STG

## 7 Conclusions

TODO: got things to work as expected

TODO: learned a lot about lazy evaluation

TODO: lots of future work to do still

## 8 References

- [1] Paul Hudak et al. “A history of Haskell: being lazy with class”. In: *Proceedings of the third ACM SIGPLAN conference on History of programming languages*. 2007, pp. 12–1.
- [2] Simon L Peyton Jones. “Implementing lazy functional languages on stock hardware: the Spineless Tagless G-machine”. In: *Journal of functional programming* 2.2 (1992), pp. 127–202.
- [3] Simon L Peyton Jones and David R Lester. “Implementing functional languages: a tutorial”. In: *Department of Computer Science, University of Glasgow* (2000).