

Just-In-Time compilation of Haskell using PyPy and GHC

Knut Halvor Skrede

December 19, 2011

Abstract

The paper describes a system for using GHC as frontend and PyPy as backend for a Haskell JIT compiler. An intermediate language in JSON format based on GHC's Core language is described. The implemented parts are a serializer written in Haskell and a deserializer written in Python, in addition to some Haskell library functions. The project is meant to serve as a base for further development.

Contents

1	Introduction	4
1.1	Project description	4
1.2	Motivation	4
1.3	System description	4
1.4	Structure of the paper	5
2	Background	5
2.1	PyPy	5
2.1.1	A quick overview	5
2.1.2	RPython	6
2.1.3	Just-in-time compilation	6
2.1.4	Haskell-Python	6
2.2	GHC	6
2.2.1	Extcore	7
3	Haskell-Python	7
4	External-core	8
4.1	Informal semantics	9
4.2	Evaluation of a program	10
5	JSON representation of Core	10
5.1	JSON	10
5.2	JSCore	11
6	Implementation	11
6.1	Tools and versions	11
6.2	Pipeline	11
6.3	Organization	12
6.4	Serializer	13
6.5	Deserializer	14
6.6	Haskell libraries	15
6.7	Evaluation	15
6.8	Tests	16
6.9	Issues	16
7	Example	18
7.1	Haskell code	18
7.2	Converted to Core	18
7.3	Converted to JSCore	18
7.4	JSCore graph	19
7.5	Result	19
8	Future work	19
8.1	Rewrite the deserializer to proper RPython	19

8.2	Implement basic Haskell libraries	19
8.3	Linking functionality	19
8.4	Test suite and benchmarking	19
8.5	Optimize JIT compiler for Haskell-Python	21
A	Formal definition of External-core	22
B	Z-Encoding	25
C	Formal definition of JSON	26
D	Formal definition of JSCore	27
E	Test programs	30
E.1	helloworld	30
E.2	helloworld2	30
E.3	factorial	30
E.4	fibonacci	30

List of Figures

1	Project overview	5
2	Code tree: organization	13
3	Pipeline with core2js	14
4	Result output from running tests	16
5	Pipeline without haskell2js	17
6	Example program translated to JSCore	20

List of Tables

4	Grammar for External-core	24
5	GHC z-encoding	25
6	Grammar for JSON	26
8	Grammar for JSCore	29

1 Introduction

1.1 Project description

The project focuses on implementing a serializer from Haskell to an intermediate format using GHC, and a deserializer from the intermediate format to an interpreter using PyPy. And then to test the result on some simple Haskell programs. The hope is that this can serve as a base for future development into a full Haskell JIT compiler.

1.2 Motivation

The motivation behind the project is to see if a *strongly-typed non-strict purely-functional* programming language, Haskell, can benefit from just-in-time (JIT) compilation.

Some people seem to think that there is little to gain from using JIT compilation of strongly typed languages since they can be so heavily optimized at compile time. However, a JIT compiler has a lot more information to work with. Speedup in JIT compilers is achieved by feeding values observed at runtime into compilation and exploit them. For example, by locating variables that change very slowly, the JIT compiler can optimize multiple instances of the code, one for each variation of the variable. [2]

PyPy is a project that implements a meta-tracing JIT. The project defines a proper subset of Python called RPython. This language has the characteristics that it is possible to perform type-inference on it. Because of this, interpreters written in RPython can be translated to efficient C code. The idea is that interpreters can be written rapidly in RPython, and the interpreter implemented in RPython will benefit from PyPy's JIT compiler. To optimize the interpreter, the RPython toolchain accepts some compiler hints to determine what parts of the code to trace, and to define the static and dynamic parts of the interpreter "memory". [2]

GHC (Glasgow Haskell Compiler) is the most advanced compiler for the Haskell programming language. The GHC team has defined an external format for Core (Core is an intermediate format used by GHC) called *External Core*. GHC's front-end translates Haskell 98 (plus some extensions) into Core. This way it is possible to reuse GHC for parsing, desugaring and typechecking, while implementing the back-end separately. [7]

1.3 System description

The system will use GHC to generate the representation of the Core language and a simple Haskell program to serialize this representation into JSON. This

entails using GHC for parsing, typechecking, desugaring and simplification. Simplification is an optimization step in GHC, performing small and simple transformations.

The resulting code will then be deserialized by the PyPy Haskell interpreter and executed. See figure 1 for a description of the implemented and to-be-implemented parts of the project.

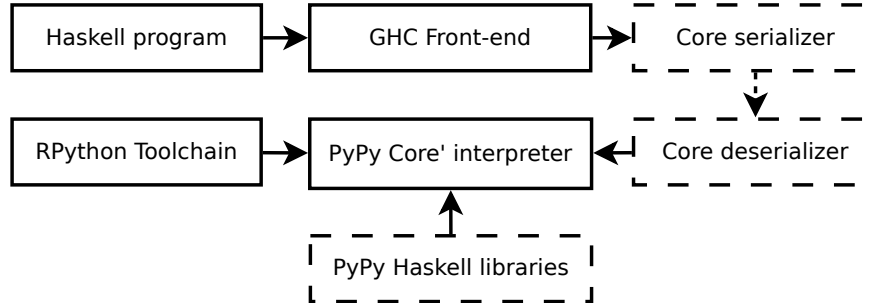


Figure 1: Project overview. Full border: implemented, Dotted border: unimplemented.

1.4 Structure of the paper

Following this introduction is a section discussing the background of the project. The background section contains some information about the tools being used, and why. Section 3 contains a description of the external language representation of the intermediate language used by GHC. Section 4 follows with the description of the intermediate format used in this project, which is a combination of JSON and External-core. Section 5 discusses the implementation of the system, the issues faced during development and some of the remaining issues. In section 6, example programs are passed through the pipeline with intermediate representations presented at each step. Section 7 discusses future work.

2 Background

2.1 PyPy

2.1.1 A quick overview

The PyPy project is basically two things:

1. the RPython toolchain, written in Python, is a set of compiler tools for programs written in RPython.

2. and an implementation of Python using these tools.

In this paper PyPy refers to former.

The basic concept of PyPy is to use a high-level language to allow for rapid development of interpreters for a variety of platforms. By implementing a compiler for RPython, interpreters for other languages can be written in RPython and compiled to any platform supported by the PyPy toolchain. Supported platforms include CLI and JVM. [1]

PyPy uses the meta-programming argument; if a VM (virtual machine) can be written at a level of abstraction high enough, then it should be possible to automatically translate this VM to other lower-level platforms. This is what PyPy does. [6]

2.1.2 RPython

RPython (Restricted Python) is a restricted proper subset of Python that it is possible to perform type inference on. This means that it can be translated to efficient C code, and it enables easy analysis as well as efficient compilation. This also means that RPython code can be run and debugged by Python interpreters, like CPython. [1]

2.1.3 Just-in-time compilation

The JIT compiler is the reason why PyPy is able to compete with other language implementations on speed. Or rather, it's meta-tracing JIT. The JIT is implemented for the RPython compiler, but through a set of compiler hints, it is able to trace the execution of the application interpreted by the RPython program.

2.1.4 Haskell-Python

Haskell-Python is an interpreter for a subset of the Haskell language, called Core'. We call it Core' (Core marked) here because it does not directly correspond to the Core language used by GHC (Glasgow Haskell Compiler). Haskell-Python is written in RPython, and is compiled into a JIT compiler using the RPython toolchain. Our goal for this project is to extend Haskell-Python to use GHC as a frontend for compilation of Haskell programs.

2.2 *GHC*

GHC (Glasgow Haskell Compiler) is a compiler for the Haskell programming language. *Haskell* is a *strongly-typed non-strict purely-functional* programming

language, it will not be described in any detail here, since Haskell is not the language we focus on. See [3] for an introduction to Haskell.

Core is an intermediate language used by the Glasgow Haskell Compiler [4], and it is this language we wish to interpret. Core is a desugared version of Haskell, things like pattern matching and list comprehensions are transformed out to simpler constructs. [5]

2.2.1 Extcore

Extcore is a Haskell package for working with GHC's Core language. Among other things, it implements a parser for External-core, this is the part used from extcore in this project.

3 Haskell-Python

The Haskell-Python interpreter contains the following classes:

Class name	Docstring
Symbol	A cached symbol that can be compared by identity (which is not true for strings).
HaskellObject	Base class for all objects that the interpreter handles.
Value	Base class for evaluated values (i.e. already in head-normal form).
Constructor	A constructor. This is an abstract base class, there are subclasses generated below for various numbers of arguments.
ConstructorN	
Integer	
AbstractFunction	
Function	A user-defined function, i.e. written in Haskell
Rule	One rule of a user-defined function.
Substitution	The body of a function with numbered variables substituted by values.
PrimFunction	A primitive function, i.e. one not implemented in Haskell but at the machine level.
Var	
NumberedVar	
Application	A function application. This is an abstract base class, there are subclasses generated below for various numbers of arguments.
ApplicationN	

Thunk	An unevaluated function application.
StackElement	Base class of the stack elements of the evaluation stack.
CopyStackElement	Need to copy the top of the stack.
UpdateStackElement	Need to update the thunk stored in this after its content has been evaluated.

3.1 Our extensions

The extensions to Haskell-Python implemented in this project, is the deserializer.

4 External-core

The Core language is an intermediate language used by GHC. It is the internal program representation in the compilers simplification phase. It appears to be a subset of Haskell, but with explicit type annotations in the style of the polymorphic lambda calculus (F_w). External-core is an external representation of Core generated from Haskell, by using GHC with a compiler flag (*-fext-core*). By using External-core, one may implement just parts of a Haskell compiler, using the remaining parts from GHC. For this project, GHC is used to generate External-core. This way, desugaring, type checking, pattern matching and overloading is performed. The remaining task is then to interpret the External-core representation. [7]

The GHC-generated External-core uses z-encoding for special characters in names (variables, constructors, ...), for a definition of the z-encoding, see appendix B. This is a good thing, as much of the functionality represented by these names has to be implemented in RPython, and special characters (like #, or +) is not allowed in Python names (function identifiers, class identifiers, ...).

Without using GHC to produce External-core, linking code into GHC would be an optional way of achieving this, which would be a difficult and large task. Or, the GHC API could be used to do the same task more cleanly. [7]

The initial starting point of this project was to use the GHC API, the reasoning was that External-core is not fully parenthesized and more tricky to parse. Thus, generating a fully parenthesized and more machine-readable format would make sense. However, it turned out that this too was a complicated task. As the internal datatypes of GHC representing Core does not match the description of External-core. The choice was then made to use GHC-generated External-core as the base for creating a new intermediate format that could be easily generated using available packages for manipulating External-core (*extcore* is used for this).

For a formal definition of External-core, see appendix A.

4.1 Informal semantics

The first construct in the External-core representation of a program, is the *module*. This construct is represented by a *module identifier*, followed by a list of *type definitions* and *value-definitions* respectively.

The *module* directly corresponds to Haskell source modules. The *module identifier* contains information of what package it belongs to, followed by its name. *Identifiers* that are defined at the top level of the module can be internal or external. External identifiers can be referenced from other modules in the program, internal identifiers can not.

A *type definition* can be an *algebraic-type*, or a *newtype construct*. An *algebraic-type* is represented by a *type constructor* and a *type binder*, followed by a list of *constructor definitions*.

Each new *algebraic type* introduces a new *type constructor* and a set of one or more *constructor definitions*.

A *constructor definition* is represented by a *data constructor*, followed by a *type binder* and one or more *atomic types*.

A *value definition* can be *recursive* or *non-recursive*. A *recursive value definition* is represented by one or more *non-recursive value definitions*. A *non-recursive value definition* is represented by a *variable identifier* followed by a *type* and an *expression*.

An *Atomic expression* can be either of the following; a *variable*, *data constructor*, *literal* or a *nested expression*.

An *expression* can be either of the followin; an *atomic expression*, *application*, *abstraction*, *local definition*, *case expression*, *type coercion*, *expression note*, *external reference*, *dynamic external reference*, or an *external label*.

An *argument* can be either a *type argument* or a *value argument*. A *type argument* is simply an *atomic type*, and a *value argument* is an *atomic expression*.

A *case alternative* can be either of the following; a *constructor alternative*, a *literal alternative* or a *default alternative*.

A *binder* can be either a *type binder* or a *value binder*.

A *type binder* can either be a *type variable* (implicitly of kind ***), or a *type variable* followed by a *kind* (explicitly kinded).

A *value binder* is a *variable* followed by a *type*.

A *literal* can be either an *integer*, *rational*, *character* or *string*.

An *atomic type* can be either a *type variable*, *type constructor* or a *nested type*.

A *basic type* can be either of the following; an *atomic type*, *type application*, *transitive coercion*, *symmetric coercion*, *unsafe coercion*, *left coercion*, *right coercion*, or an *instantiation coercion*.

Types are built from *type constructors* and *type variables* using *type application* and *universal quantification*. There are a number of primitive *type constructors* defined in the "GHC.Prim" module. *algebraic-type* and *newtype* constructs introduce new *type constructors*. The *type constructors* are distinguished by name only.

A *type* can be either a *basic type*, a *type abstraction* or a *arrow type construction*.

An *atomic kind* can be either of the following; a *lifted kind* (*), *unlifted kind* (#), *open kind* (?), *equality kind*, or a *nested kind*.

A *kind* can be either an *atomic kind* or an *arrow kind*.

4.2 Evaluation of a program

A program is evaluated by reducing the expression "main:ZCMain.main" (note that qualified names in a module named *m* must have module name *m*, this is the only exception) to *weak-head-normal-form* (WHNF), i.e. a primitive value, lambda abstraction, or fully applied data constructor. A heap is used to make sure evaluation is shared. The heap contains two types; a *thunk*, or a *WHNF*. A *thunk* is an unevaluated expression, also called a *suspension*. A *WHNF* is an evaluated expression, the result of evaluating a *thunk* is a *WHNF*. [7]

5 JSON representation of Core

5.1 JSON

JavaScript Object Notation (JSON) is a lightweight data interchange format.

Since a library for manipulating JSON is available for Haskell, this makes it a good choice for the project. In addition, it is easy to parse and the Haskell library contains a pretty-printer, making the result easier to inspect.

For a formal definition of JSON see appendix C

5.2 JSCore

In order to work with JSON and External-core, a format was defined that expresses the Core program in JSON notation. Most of the right hand side of the grammar evaluates to JSON Values. Even though the grammar is changed to support JSON, an effort was made to keep it similar to the original Core

grammar for easy referencing. The size of the resulting files was not considered to be an issue. Other than the syntax, the JSCore language is the same as External-core.

For a formal definition of JSCore see appendix D.

6 Implementation

6.1 Tools and versions

The following tools and package versions was used in the implementation:

- GHC version 7.0.3
- extcore version 1.0.1
- PyPy current head branch (Last tested 09/11/2011)
- Haskell-Python interpreter: Interpreter of Core' written in RPython
- CPython 2.7: Used to test the Core' interpreter without it having to be correct RPython. (pypy-python could also have been used)
- Git was used as the projects version control system.

6.2 Pipeline

The project implements the following pipeline (see figure 3):

1. Serialize Haskell program:
 - (a) Create External-core file from Haskell program using GHC.
 - (b) Create JSCore from External-core using the extcore and JSON packages.
2. Deserialize JSCore:
 - (a) Parse JSCore using the parsing tools available for PyPy
 - (b) Build Core AST from resulting JSON datatype
3. Evaluate program:
 - (a) Evaluation is done by the already implemented PyPy Core' interpreter, Haskell-Python. Additional functionality had to be built on top of this, mostly Haskell library functions.

6.3 Organization

The implementation is organized as represented by figure 2. The main folder (interpreter) contains the main program, and a program for generating dot¹ files, "makegraph.py" (used to create graphs of parsed JSCore files using graphviz²). The "haskell" folder contains the Haskell-Python Core' interpreter code, used by the main program for evaluation, and by the parser to generate the abstract-syntax-tree (AST). In addition to this, the subfolder "ghc" implements some simple functionality to be used by the test-programs. Among others, a very simple IO function for printing text to the terminal (putStrLn). These are categorised into packages corresponding to GHC packages, and libraries, corresponding to the libraries in the GHC packages. The RPython modules in these "packages" are loaded, and references to the attributes they contain are used during the creation of the AST. See figure 3 for a simple description of the pipeline.

The "core" folder contains a Haskell program for generating JSCore files from External-core files (core2js.hs), the JSCore parser (parser.py), and a simple datastructure representing Core modules (module.py).

The "tests" folder naturally contains tests, each test is in its own subfolders.

¹A dot file is a file written in dot-language, used by the *dot* command line tool from the graphviz software package

²Graphviz is an open source software package for graph visualization

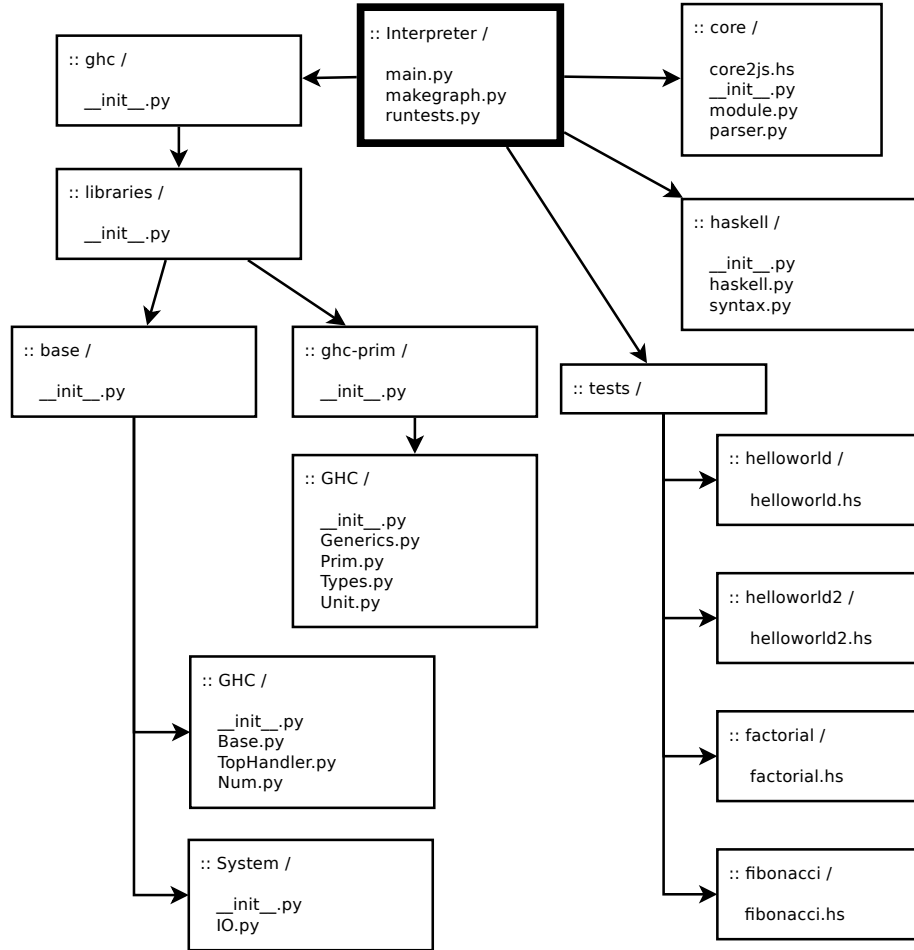


Figure 2: Code tree: The top of the boxes is the folder name, and the rest is the source files. Arrows represent subfolders.

6.4 Serializer

The serializer consists of two parts; GHC generating External-core, and a Haskell program to generate JSCore (core2js.hs).

External-core is easily generated by using a compiler flag:

```
ghc -fext-core {path-to-program}
```

The Haskell program generating JSCore uses the extcore packages. This package implements functionality for working with External-core. The result is a datastructure mapping directly to the External-core format defined in [7]. By

traversing this structure, the program builds up a JSON tree using the Haskell JSON package. The result is a tree of JSON constructs (corresponding to the grammar defined in table 8), this is then pretty-printed and dumped to a file.

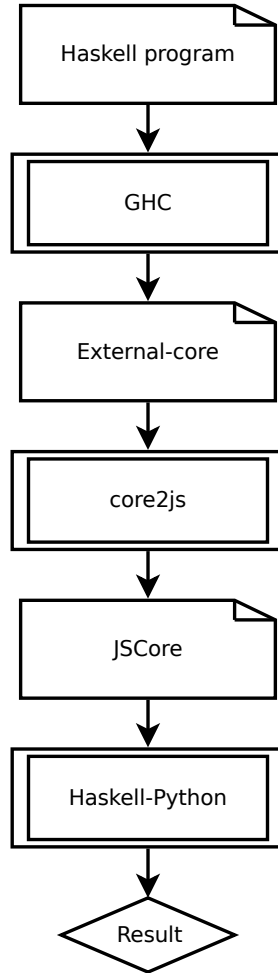


Figure 3: Pipeline implementation with core2js using extcore

6.5 Deserializer

The deserializer implements a JSON parser. The resulting datastructure is then traversed, building up an AST using the constructs defined in the Haskell-Python interpreter.

PyPy implements a parser generator, this simply takes a grammar defined as a string, written in extended-backus-naur-form (EBNF), and generates a parser.

This parser is then used to create a JSON datastructure, as represented by table 6. The resulting datastructure is then traversed. By checking the contents of the JSON constructs with the actual External-core format, the Core AST is built. External functionality is imported from the Python implementations of the Haskell libraries as it is encountered.

After this is done, we are left with a "module" object, corresponding to the initial Haskell module.

6.6 Haskell libraries

To make some simple test-cases work, some basic Haskell functionality had to be implemented. Some of this functionality was implemented already in the Haskell-Python Core' interpreter. The work done here was mostly to organize the functionality into modules corresponding to Haskell modules. The functionality implemented in these modules does however, not correspond to the Haskell implementations. This is left for future work, as this is a large task.

From figure 6 (representing a "hello world" program in JSCore), the expression "base:SystemzIO.putStrLn" (See appendix B for a definition of the z-encoding, used in this expression identifier) corresponds to the Haskell function "putStrLn", which is located in the Haskell module "System.IO". This is translated into a reference to the function "putStrLn" defined in the python module "ghc.libraries.base.System.IO". See figure 2.

6.7 Evaluation

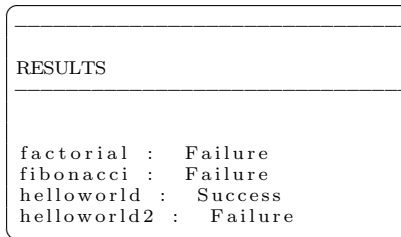
In order to evaluate the Haskell programs correctly, the expression "main:ZCMain.main" is reduced to *WHNF*. However, to do this correctly would require a lot of the functionality used by GHC to be implemented. Specifically, "GHC.TopHandler.runMainIO". In order to implement this function a lot of other functionality would have to be implemented. The function is a wrapper around "main:Main.main", it catches uncaught exceptions and flushes stdout/stderr before exiting.

Currently, a simple hack is implemented; the function "runMainIO" simply returns the function "main:Main.main". The "main:Main.main" function is the entry-point to our program. This function is then evaluated when reducing "main:ZCMain.main" to *WHNF*.

Implementing this functionality in a manner similar to GHC is a goal for future development. Currently, the hacks suffice for testing. Simple Haskell functionality is implemented at a high level, such as the "putStrLn" function which is implemented as a simple "print" Python function.

6.8 Tests

Testing is done by a simple Python script. The test directory has a subdirectory for each test. This directory contains a Haskell program. The program is converted to External-core, then to JSCore, and then executed. If the program returns "0" it passes the test. The results of running these tests are given in figure 4.



```
RESULTS

factorial : Failure
fibonacci : Failure
helloworld : Success
helloworld2 : Failure
```

Figure 4: Result output from running tests (runtests.py)

The tests are supposed to be in the following increasing order of difficulty: helloworld, helloworld2, factorial and fibonacci. All the failures are due to the same problem; unimplemented primitives and library functions. See appendix E for test programs.

6.9 Issues

Very few of the initial plans regarding this project was actually realized in the implementation. The GHC API was intended to be used to generate JSCore, but this failed do to a lack of experience with Haskell and the GHC API (see figure 5 for a description of the initially intended pipeline). It was then discovered that a lot of the necessary code was already written in GHC (the code that generates the External-core files), however, this code was not exported in any way. It was also deeply embedded in the GHC code. After several unfruitful attempts at creating the JSCore format by using the GHC API and by reimplementing the GHC functions generating External-core. And some emails back and forth with the GHC team, these methods where abandoned.

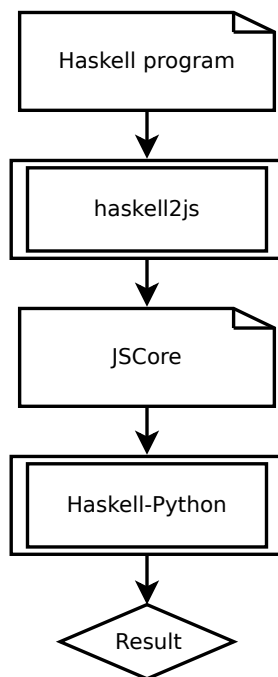


Figure 5: Pipeline with implementation of `haskell2js` using the GHC API (Not implemented)

The `extcore` package was then introduced, as it was thought to serve nicely for our purpose. This does however add an extra step to the generation of JSCore, having to generate External-core first. There is no way of working with External-core without parsing it from a file. See figure 3.

The reason for not using `extcore` from the beginning was that it was thought to not be well supported. As well as the seeming lack of support for the External-core format.

A large amount of time was spent trying to generate this intermediate format. And the greatest obstacle was that the tools being used was either incompatible, or lacking in documentation.

The method described here worked well for very trivial Haskell programs, however, it turned out that the `extcore` package was not able to parse any nontrivial External-core files generated by the GHC version used. It was thought that this was due to the fact that the `extcore` package was written for earlier versions of GHC. After unsuccessfully attempting to make the tools work with earlier versions (GHC 6.10 and 6.12 on linux and 6.10 on windows), it turned out that the `extcore` package implemented two parsers. One of which was outdated. This was not documented anywhere, but realized after e-mailing the package maintainer, who turned out to be very helpful.

7 Example

The example is a simple "hello world" program, as this was practically the only program that was able to pass through the entire pipeline. Following is the "hello world" program written in Haskell.

7.1 Haskell code

```
main = putStrLn "Hello , world!"
```

7.2 Converted to Core

After the program has passed through GHC and the External-core file has been generated, the program looks like this:

```
%module main:Main
main:Main.main :: (ghczmpirim:GHCziTypes.IO
                  ghczmpirim:GHCziUnit.Z0T) =
  base:SystemziIO.putStrLn
  (base:GHCziBase.unpackCStringzh
   ("Hello , world!" :: ghczmpirim:GHCziPrim.Addrzh));
main:ZCMain.main :: (ghczmpirim:GHCziTypes.IO
                   ghczmpirim:GHCziUnit.Z0T) =
  base:GHCziTopHandler.runMainIO @ ghczmpirim:GHCziUnit.Z0T
  main:Main.main;
```

7.3 Converted to JSCore

In the next step it is parsed and dumped to JSCore:

```
{"%module": "main:Main", "tdefg": [],
 "vdefg": [{"qvar": "main:Main.main",
  "ty": {"bty": {"qtycon": "ghczmpirim:GHCziTypes.IO"},
  "aty": {"qtycon": "ghczmpirim:GHCziUnit.Z0T"}},
  "exp": {"aexp": {"qvar": "base:SystemziIO.putStrLn"},
  "args": {"aexp": {"aexp": {"qvar": "base:GHCziBase
    .unpackCStringzh"},
  "args": {"aexp": {"lit": "Hello ,
    world!"}}}}}],
  {"qvar": "main:ZCMain.main",
  "ty": {"bty": {"qtycon": "ghczmpirim:GHCziTypes.IO"},
  "aty": {"qtycon": "ghczmpirim:GHCziUnit.Z0T"}},
  "exp": {"aexp": {"aexp": {"qvar": "base:GHCziTopHandler.
    runMainIO"},
  "args": {"aty": {"qtycon": "ghczmpirim:
    GHCziUnit.Z0T"}},
  "args": {"aexp": {"qvar": "main:Main.main"}}}]}
```

7.4 JSCore graph

Using some libraries from PyPy we can generate a nice graph, directly corresponding to the datastructure generated from parsing. See figure 6. By simply traversing this datastructure we can generate the AST for the Core' interpreter (Haskell-Python).

7.5 Result

The program results as expected, outputting the string "Hello, world!".

8 Future work

8.1 Rewrite the deserializer to proper RPython

The deserializer is currently not written in proper RPython. Converting the code to RPython will allow it to be compiled to a JIT interpreter by the RPython toolchain. This should not be a very big task, but it requires understanding of the RPython coding style. There are also restrictions on how one may use the PyPy parser tools.

8.2 Implement basic Haskell libraries

Implementing Haskell libraries is necessary to run any Haskell program passed through GHC. One option may be to implement (or automatically generate from GHC code) Haskell primitive types, and to convert the Haskell base libraries to JSCore. Compiling GHC with the *-fext-core* flag should be possible, and should generate External-core files for all library modules.

8.3 Linking functionality

Currently no thought has been given to the linking between multiple modules. This must be implemented for any non-trivial Haskell programs to function.

8.4 Test suite and benchmarking

A framework for testing and benchmarking should be set up. Allowing us to compare speed to other implementations, and progress of development. GHC uses a test framework relying on Python and GNU Make that should be looked into.

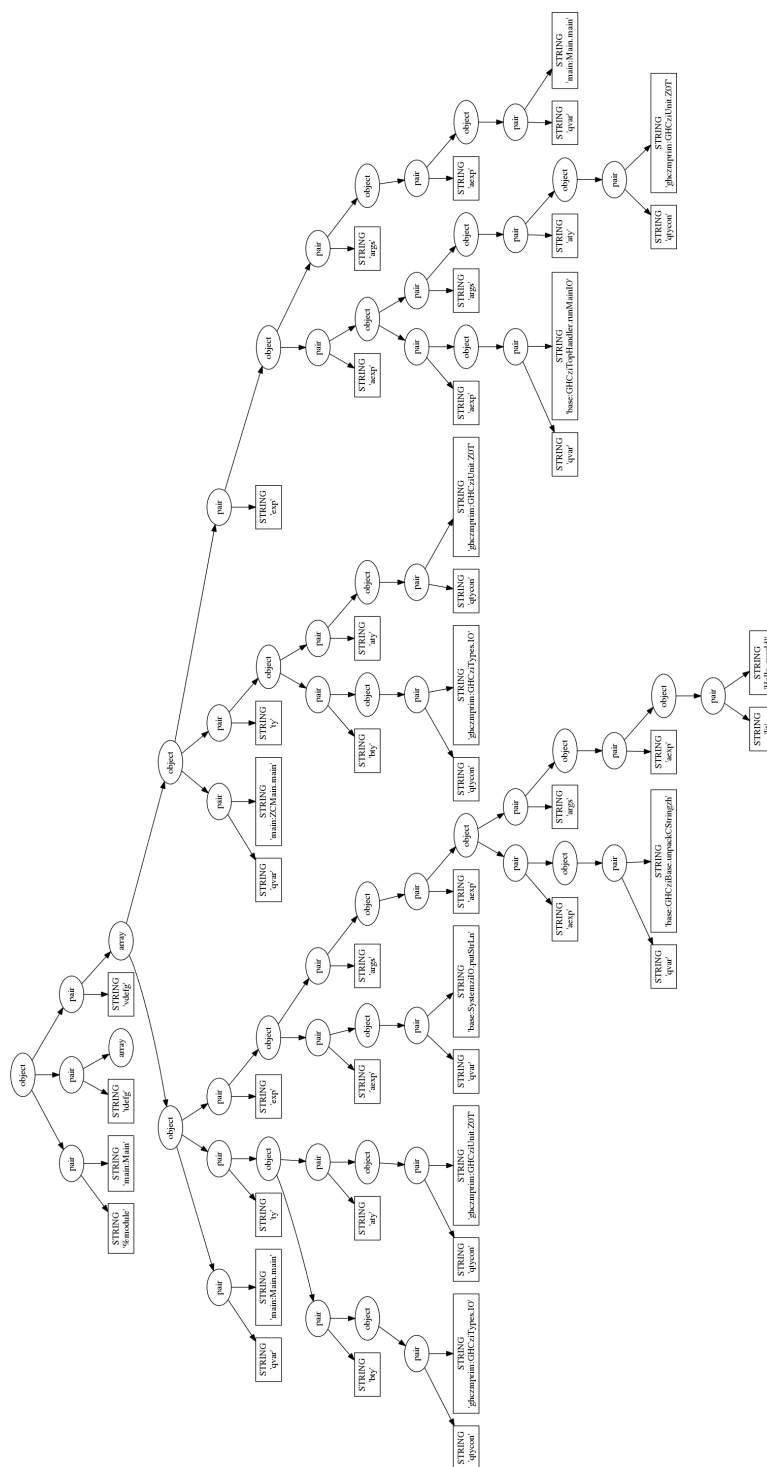


Figure 6: Example program translated to JSCore

8.5 Optimize JIT compiler for Haskell-Python

Further optimization of the JIT compiler for the Haskell-Python interpreter will be necessary to achieve good speed.

Appendices

A Formal definition of External-core

The following semantics is used to define the Core grammar, as seen in [7]:

$[\text{pat}]$: optional
 $\{ \text{pat} \}$: zero or more repetitions
 $\{ \text{pat} \}^+$: one or more repetitions
 $\text{pat}_1 | \text{pat}_2$: choice

Module

$\text{module} \rightarrow \% \text{module } mident \{ tdefg ; \} \{ vdefg ; \}$

Type defn.

$tdefg \rightarrow \% \text{data } qtycon \{ tbind \} = \{ [cdef \{ ; cdef \}] \}$ algebraic type
 $| \% \text{newtype } qtycon \{ tbind \} = ty$ newtype

Constr. defn.

$cdef \rightarrow qdcon \{ @ tbind \} \{ aty \}^+$

Value defn.

$vdefg \rightarrow \% \text{rec } \{ vdef \{ ; vdef \} \}$ recursive
 $| vdef$ non-recursive
 $vdef \rightarrow qvar :: ty = exp$

Atomic expr.

$aexp \rightarrow qvar$ variable
 $| qdcon$ data constructor
 $| lit$ literal
 $| (exp)$ nested expr.

Expression

$exp \rightarrow aexp$ atomic expr.
 $| aexp \{ arg \}^+$ application
 $| \{ binder \} \rightarrow exp$ abstraction
 $| \% \text{let } vdefg \% \text{in } exp$ local definition
 $| \% \text{case } (aty) exp \% \text{of } vbind \{ alt \{ ; alt \} \}$ case expr.
 $| \% \text{cast } exp aty$ type coercion
 $| \% \text{note } " \{ char \} " exp$ expression note
 $| \% \text{external ccal } " \{ char \} " aty$ external reference
 $| \% \text{dynexternal ccal } aty$ external reference (dynamic)
 $| \% \text{label } " \{ char \} "$ external label

Argument

$arg \rightarrow @ aty$ type argument
 $| aexp$ value argument

Case alt

A FORMAL DEFINITION OF EXTERNAL-CORE

<i>alt</i>	→	<i>qdcon</i> { @ <i>tbind</i> } { <i>vbind</i> } -> <i>exp</i>	constructor alternative
		<i>lit</i> -> <i>exp</i>	literal alternative
		%_ -> <i>exp</i>	default alternative
Binder			
<i>binder</i>	→	@ <i>tbind</i>	type binder
		<i>vbind</i>	value binder
Type binder			
<i>tbind</i>	→	<i>tyvar</i>	implicit of kind *
		(<i>tyvar</i> :: <i>kind</i>)	explicitly kinded
Value binder			
<i>vbind</i>	→	(<i>var</i> :: <i>ty</i>)	
Literal			
<i>lit</i>	→	([] { <i>digit</i> } ⁺ :: <i>ty</i>)	integer
		([] { <i>digit</i> } ⁺ % { <i>digit</i> } ⁺ :: <i>ty</i>)	rational
		(' <i>char</i> ' :: <i>ty</i>)	character
		(" { <i>char</i> } " :: <i>ty</i>)	string
Character			
<i>char</i>	→	Any ASCII character in range 0x20-0x7E except 0x22, 0x27, 0x5c	
		\x <i>hex hex</i>	ASCII code escape sequence
<i>hex</i>	→	0 ... 9 a ... f	
Atomic type			
<i>aty</i>	→	<i>tyvar</i>	type variable
		<i>qtycon</i>	type constructor
		(<i>ty</i>)	nested type
Basic type			
<i>bty</i>	→	<i>aty</i>	atomic type
		<i>bty aty</i>	type application
		%trans <i>aty aty</i>	transitive coercion
		%sym <i>aty</i>	symmetric coercion
		%unsafe <i>aty aty</i>	unsafe coercion
		%left <i>aty</i>	left coercion
		%right <i>aty</i>	right coercion
		%inst <i>aty aty</i>	instantiation coercion
Type			
<i>ty</i>	→	<i>bty</i>	basic type
		%forall { <i>tbind</i> } ⁺ . <i>ty</i>	type abstraction
		<i>bty</i> -> <i>ty</i>	arrow type construction
Atomic kind			
<i>akind</i>	→	*	lifted kind
		#	unlifted kind
		?	open kind
		<i>bty</i> :=: <i>bty</i>	equality kind
		(<i>kind</i>)	nested kind
Kind			
<i>kind</i>	→	<i>akind</i>	atomic kind
		<i>akind</i> -> <i>kind</i>	arrow kind

Identifier

<i>mident</i>	\rightarrow	<i>pname</i> : <i>uname</i>	module
<i>tycon</i>	\rightarrow	<i>uname</i>	type constr.
<i>qtycon</i>	\rightarrow	<i>mident</i> . <i>tycon</i>	qualified type constr.
<i>tyvar</i>	\rightarrow	<i>lname</i>	type variable
<i>dcon</i>	\rightarrow	<i>uname</i>	data constr.
<i>qdcon</i>	\rightarrow	<i>mident</i> . <i>dcon</i>	qualified data constr.
<i>var</i>	\rightarrow	<i>lname</i>	variable
<i>qvar</i>	\rightarrow	[<i>mident</i> .] <i>var</i>	optionally qualified variable

Name

<i>lname</i>	\rightarrow	<i>lower</i> { <i>namechar</i> }
<i>uname</i>	\rightarrow	<i>upper</i> { <i>namechar</i> }
<i>pname</i>	\rightarrow	{ <i>namechar</i> } ⁺
<i>namechar</i>	\rightarrow	<i>lower</i> <i>upper</i> <i>digit</i>
<i>lower</i>	\rightarrow	a b ... z _
<i>upper</i>	\rightarrow	A B ... Z _
<i>digit</i>	\rightarrow	0 1 ... 9

Table 3: Grammar for External-core

B Z-Encoding

Character	Code
Tuples: (##) () (,,)	Z1H Z0T Z3T
Constructors: () [] : Z	ZL ZR ZM ZN ZC ZZ
Variables: z & — ^ \$ = i # . i - ! + , \ / * - % c	zz za zb zc zd ze zg zh zi zl zm zn zp zq zr zs zt zu zv znnnU

Table 4: GHC z-encoding

C Formal definition of JSON

Object	
<i>object</i>	→ { }
	{ <i>members</i> }
<i>members</i>	→ <i>pair</i>
	<i>pair</i> , <i>members</i>
<i>pair</i>	→ <i>string</i> : <i>value</i>
Array	
<i>array</i>	→ []
	[<i>elements</i>]
<i>elements</i>	→ <i>value</i>
	<i>value</i> , <i>elements</i>
Value	
<i>value</i>	→ <i>string</i>
	<i>number</i>
	<i>object</i>
	<i>array</i>
	true
	false
	null
String	
<i>string</i>	→ ""
	" <i>chars</i> "
<i>chars</i>	→ <i>char</i>
	<i>char</i> <i>chars</i>
<i>char</i>	→ any Unicode character except " or \ or control characters: \\ \/ \b \f \n \r \t \u four-hex digits
Number	
<i>number</i>	→ <i>int</i>
	<i>int</i> <i>frac</i>
	<i>int</i> <i>exp</i>
	<i>int</i> <i>frac</i> <i>exp</i>
<i>int</i>	→ <i>digit</i>
	<i>digit</i> 1 – 9 <i>digits</i>
	- <i>digit</i>
	- <i>digit</i> 1 – 9 <i>digits</i>
<i>frac</i>	→ . <i>digits</i>
<i>exp</i>	→ e <i>digits</i>
<i>digits</i>	→ <i>digit</i>
	<i>digit</i> <i>digits</i>
<i>e</i>	→ e
	e+
	e-
	E
	E+
	E-

Table 5: Grammar for JSON

D Formal definition of JSCore

$[pat]$:	Zero or more repetitions of pat surrounded by $[]$ and comma separated (A JSON Array).
$[pat]^+$:	One or more repetitions of pat surrounded by $[]$ and comma separated (A JSON Array).
$\{ pat \}$:	Represents a JSON Object, pat is a JSON <i>members</i> .
$pat_1 \mid pat_2$:	Choice.
$\ pat \ $:	Optional

JSCore grammar:

Module			
$module$	\rightarrow	$\{ \%module : mident , \%tdefg : [tdefg] , \%vdefg : [vdefg] \}$	
Type defn.			
$tdefg$	\rightarrow	$\{ \%data : qtycon , \%tbind : [tbind] , \%cdef : [cdef] \}$	algebraic type
	\mid	$\{ \%newtype : qtycon , \%qtycon : qtycon , \%tbind : [tbind] , \%ty : ty \}$	newtype
Constr. defn.			
$cdef$	\rightarrow	$\{ \%qdcon : qdcon , \%tbind : [tbind] , \%aty : [aty]^+ \}$	
Value defn.			
$vdefg$	\rightarrow	$\{ \%rec : [vdef]^+ \}$	recursive
	\mid	$vdef$	non-recursive
$vdef$	\rightarrow	$\{ \%qvar : qvar , \%ty : ty , \%exp : exp \}$	
Atomic expr.			
$aexp$	\rightarrow	$\{ \%qvar : qvar \}$	variable
	\mid	$\{ \%qdcon : qdcon \}$	data constructor
	\mid	$\{ \%lit : lit \}$	literal
	\mid	$\{ \%exp : exp \}$	nested expr.
Expression			
exp	\rightarrow	$aexp$	atomic expr.
	\mid	$\{ \%aexp : aexp , \%args : [arg]^+ \}$	application
	\mid	$\{ \%lambda : [binder] , \%exp : exp \}$	abstraction
	\mid	$\{ \%let : vdefg , \%in : exp \}$	local definition
	\mid	$\{ \%case : aty , \%exp : exp , \%of : vbind , \%alt : [alt]^+ \}$	case expr.
	\mid	$\{ \%cast : exp , \%aty : aty \}$	type coercion
	\mid	$\{ \%note : " \{ char \} " , \%exp : exp \}$	expression note
	\mid	$\{ \%external ccal : " \{ char \} " , \%aty : aty \}$	external reference
	\mid	$\{ \%dynexternal ccal : aty \}$	external reference (dynamic)
	\mid	$\{ \%label : " \{ char \} " \}$	external label
Argument			
arg	\rightarrow	$\{ \%aty : aty \}$	type argument
	\mid	$\{ \%aexp : aexp \}$	value argument
Case alt			
alt	\rightarrow	$\{ \%qdcon : qdcon , \%tbind : [tbind] , \%vbind : [vbind] , \%exp : exp \}$	constructor alternative
	\mid	$\{ \%lit : lit , \%exp : exp \}$	literal alternative
	\mid	$\{ \%_ : exp \}$	default alternative
Binder			

$binder \rightarrow \begin{cases} \text{"tbind"} : tbind \\ \text{"vbind"} : vbind \end{cases}$

type binder
value binder

Type binder

$tbind \rightarrow \begin{cases} \text{"tyvar"} : tyvar \\ \text{"tyvar"} : tyvar, \text{"kind"} : kind \end{cases}$

implicit of kind *
explicitly kinded

Value binder

$vbind \rightarrow \{ \text{"var"} : var, \text{"ty"} : ty \}$

Literal

$lit \rightarrow \begin{cases} jsstring \\ jsnumber \end{cases}$

string
number

JSON String

$jsstring \rightarrow \begin{cases} "" \\ \text{" } jschars \text{ " } \end{cases}$
 $jschars \rightarrow \begin{cases} jschar \\ jschar jschars \end{cases}$
 $jschar \rightarrow \begin{cases} \text{any Unicode character except " or \ or control characters:} \\ \backslash \\ \backslash / \\ \backslash b \\ \backslash f \\ \backslash n \\ \backslash r \\ \backslash t \\ \backslash u \text{ four-hex digits} \end{cases}$

JSON Number

$jsnumber \rightarrow \begin{cases} jsint \\ jsint jsfrac \\ jsint jsexp \\ jsint jsfrac jsexp \end{cases}$
 $jsint \rightarrow \begin{cases} jsdigit \\ jsdigit1 - 9 jsdigits \\ - jsdigit \\ - jsdigit1 - 9 jsdigits \end{cases}$
 $jsfrac \rightarrow . jsdigits$
 $jsexp \rightarrow jse jsdigits$
 $jsdigits \rightarrow \begin{cases} jsdigit \\ jsdigit jsdigits \end{cases}$
 $jse \rightarrow \begin{cases} e \\ e+ \\ e- \\ E \\ E+ \\ E- \end{cases}$

Atomic type

$aty \rightarrow \begin{cases} \text{"tyvar"} : tyvar \\ \text{"qtycon"} : qtycon \\ \text{"ty"} : ty \end{cases}$

type variable
type constructor
nested type

Basic type

$bty \rightarrow \begin{cases} aty \\ \text{"bty"} : bty, \text{"aty"} : aty \\ \text{"%trans"} : aty, \text{"aty"} : aty \\ \text{"%sym"} : aty \end{cases}$

atomic type
type application
transitive coercion
symmetric coercion

D FORMAL DEFINITION OF JSCORE

		{ "%unsafe" : <i>aty</i> , "aty" : <i>aty</i> }	unsafe coercion
		{ "%left" : <i>aty</i> }	left coercion
		{ "%right" : <i>aty</i> }	right coercion
		{ "%inst" : <i>aty</i> , "aty" : <i>aty</i> }	instantiation coercion
Type			
<i>ty</i>	→	<i>bty</i>	basic type
		{ "%forall" : [<i>tbind</i>] ⁺ , "ty" : <i>ty</i> }	type abstraction
		{ "bty" <i>bty</i> , "ty" : <i>ty</i> }	arrow type construction
Atomic kind			
<i>akind</i>	→	*	lifted kind
		#	unlifted kind
		?	open kind
		{ "bty" : <i>bty</i> , "bty" : <i>bty</i> }	equality kind
		{ "kind" : <i>kind</i> }	nested kind
Kind			
<i>kind</i>	→	{ "akind" : <i>akind</i> }	atomic kind
		{ "akind" : <i>akind</i> , "kind" : <i>kind</i> }	arrow kind
Identifier			
<i>mident</i>	→	" <i>pname</i> : <i>uname</i> "	module
<i>tycon</i>	→	" <i>uname</i> "	type constr.
<i>qtycon</i>	→	" <i>mident</i> . <i>tycon</i> "	qualified type constr.
<i>tyvar</i>	→	" <i>lname</i> "	type variable
<i>dcon</i>	→	" <i>uname</i> "	data constr.
<i>qdcon</i>	→	" <i>mident</i> . <i>dcon</i> "	qualified data constr.
<i>var</i>	→	" <i>lname</i> "	variable
<i>qvar</i>	→	" <i>mident</i> . <i>var</i> "	optionally qualified variable
Name			
<i>lname</i>	→	<i>lower</i> { <i>namechar</i> }	
<i>uname</i>	→	<i>upper</i> { <i>namechar</i> }	
<i>pname</i>	→	{ <i>namechar</i> } ⁺	
<i>namechar</i>	→	<i>lower</i> <i>upper</i> <i>digit</i>	
<i>lower</i>	→	a b ... z _	
<i>upper</i>	→	A B ... Z _	
<i>digit</i>	→	0 1 ... 9	

Table 7: Grammar for JSCore

E Test programs

E.1 helloworld

```
main = putStrLn "Hello , world!"
```

E.2 helloworld2

```
main = putStrLn ("Hello , " ++ "world!")
```

E.3 factorial

```
module Main where

main = do
  let n = 10
  putStrLn $ "Factorial 10 is: " ++ (show $ fac n)

fac :: Int -> Int
fac 0 = 1
fac n = n*fac (n-1)
```

E.4 fibonacci

```
module Main where

main = do
  let n = 10
  putStrLn $ "The 10. fibonacci number is: " ++ (show $ fib n)

fib :: Integer -> Integer
fib 0 = 0
fib 1 = 1
fib n = fib (n-1) + fib (n-2)
```

References

- [1] D. Ancona, M. Ancona, A. Cuni, and N.D. Matsakis. RPython: a step towards reconciling dynamically and statically typed OO languages. In *Proceedings of the 2007 symposium on Dynamic languages*, pages 53–64. ACM, 2007.
- [2] C.F. Bolz, A. Cuni, M. Fijałkowski, M. Leuschel, S. Pedroni, and A. Rigo. Runtime feedback in a meta-tracing JIT for efficient dynamic languages. In *Proceedings of the 6th Workshop on Implementation, Compilation, Optimization of Object-Oriented Languages, Programs and Systems*, page 9. ACM, 2011.
- [3] P. Hudak, S. Peyton Jones, P. Wadler, B. Boutel, J. Fairbairn, J. Fasel, M.M. Guzmán, K. Hammond, J. Hughes, T. Johnsson, et al. Report on the programming language Haskell: a non-strict, purely functional language version 1.2. *ACM SigPlan notices*, 27(5):1–164, 1992.
- [4] S.L.P. Jones, C. Hall, K. Hammond, J. Cordy, H. Kevin, W. Partain, and P. Wadler. The Glasgow Haskell compiler: a technical overview. 1992.
- [5] S.P. Jones and A. Santos. Compilation by transformation in the Glasgow Haskell Compiler. *Functional Programming, Glasgow*, pages 184–204, 1994.
- [6] A. Rigo and S. Pedroni. PyPy’s approach to virtual machine construction. In *Companion to the 21st ACM SIGPLAN symposium on Object-oriented programming systems, languages, and applications*, pages 944–953. ACM, 2006.
- [7] A. Tolmach and T. Chevalier. the GHC Team. An external representation for the GHC Core language, 2010.